

# Likelihood Theory and Methods for Generalized Linear Mixed Models

by **Aishwarya Bhaskaran**

Thesis submitted in fulfilment of the requirements for  
the degree of

**Doctor of Philosophy in Mathematics**

under the supervision of **Prof. Matt P. Wand** and  
**Dr. Joanna Wang**

University of Technology Sydney  
Faculty of Science

Submitted in September 2022

# Certificate of original authorship

I, **Aishwarya Bhaskaran**, declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy in Mathematics, in the School of Mathematical and Physical Sciences at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

This research is supported by an Australian Government Research Training Program.

**Signature:**

Production Note:  
Signature removed prior to publication.

**Date:** September 23, 2022

# Acknowledgements

First and foremost, I would like to express my sincere gratitude to my supervisor Matt Wand. This journey was definitely a challenging one and I would like to thank you for all your patience, guidance and constant encouragement throughout the research process. I have indeed learnt a lot from your valuable insights and advice. Thank you so much for always looking out for my well-being as well. This milestone would not have been possible without your help and I am truly grateful.

Next, I would like to thank all my dear friends whom I have met throughout my life. Thank you all for constantly cheering me up, being such great listening ears and for always encouraging me to go the extra mile. I am really glad to have met each and every one of you.

Last but definitely not the least, I would like to give special thanks to my dear family. Thank you for always supporting me in my decisions and for offering me advice when I needed it. Thanks for also being so understanding especially these past three years and for always making sure that I am alright. I really appreciate everything that you all have done and will continue to do for me.

# List of papers/publications

The following list of paper and publications awards relate to the work presented in this thesis:

- Jiang, J., Wand, M.P. and Bhaskaran, A. (2022), Usable and precise asymptotics for generalized linear mixed model analysis and design. *Journal of the Royal Statistical Society, Series B*, 84: 55-82. DOI: 10.1111/rssb.12473.
- Bhaskaran, A. and Wand, M.P.(2023), Dispersion parameter extension of precise generalized linear mixed model asymptotics. *Statistics and Probability Letters*, 193, Article 109691.

# Notation

In this chapter, we introduce acronyms that are frequently used throughout this thesis.

## Acronyms

Table 1: Table with acronyms used in the thesis with their meanings.

<b>Acronym</b>	<b>Meaning</b>
GLMM	Generalized linear mixed model
GVA	Gaussian variational approximation
MLE	Maximum likelihood estimation
TAP	Thouless-Anderson-Palmer

# Contents

<b>1</b>	<b>Introduction and Background</b>	<b>2</b>
1.1	Introduction . . . . .	2
1.2	Thesis Aim . . . . .	3
1.3	Outline . . . . .	3
1.4	Matrix Theory . . . . .	4
1.4.1	Difference Between Two Matrix Inverses . . . . .	4
1.4.2	Other Useful Matrix Identities . . . . .	4
1.4.3	Block Matrix Inversion . . . . .	5
1.4.4	The vec and vech Operators . . . . .	6
1.4.4.1	The Commutation Matrix . . . . .	7
1.4.4.2	The Duplication Matrix . . . . .	7
1.4.5	Kronecker Products and Related Properties . . . . .	8
1.4.6	Vector and Matrix Norms . . . . .	8
1.4.6.1	Euclidean Norm . . . . .	8
1.4.6.2	Frobenius Norm . . . . .	8
1.4.6.3	Spectral Norm . . . . .	9
1.4.7	Eigenvalue Bound Results . . . . .	10
1.4.7.1	Matrix Identities from Harville (1977) . . . . .	10
1.4.8	Vector Differential Calculus . . . . .	11
1.5	Key Integral Results . . . . .	12
1.5.1	Useful Integral Results . . . . .	12
1.5.2	Integral Form of the Matrix Square Root . . . . .	13
1.6	Key Expectation Results . . . . .	13
1.6.1	Law of Total Expectation . . . . .	13
1.6.2	Jensen's Inequality . . . . .	13
1.6.3	Markov's Inequality . . . . .	14
1.6.4	Cauchy-Schwarz Inequality . . . . .	14
1.7	Exponential Families . . . . .	14
1.7.1	One-Parameter Exponential Families . . . . .	14
1.7.2	Two-Parameter Exponential Families . . . . .	17
1.8	Generalized Linear Mixed Models . . . . .	17
1.9	Maximum Likelihood for Generalized Linear Mixed Models . . . . .	19
1.9.1	The Likelihood Function . . . . .	19
1.9.2	Maximum Likelihood Estimation . . . . .	20

1.9.3	Asymptotic Properties of Maximum Likelihood Estimators for Generalized Linear Mixed Models . . . . .	21
1.10	Asymptotics . . . . .	23
1.10.1	Convergence of Random Variables . . . . .	23
1.10.1.1	Convergence in Probability . . . . .	23
1.10.1.2	Convergence in Distribution . . . . .	23
1.10.1.3	Continuous Mapping Theorem . . . . .	24
1.10.1.4	Slutsky's Theorem . . . . .	24
1.10.1.5	Cramér-Wold Device . . . . .	24
1.10.2	Stochastic Order Notation . . . . .	25
1.10.3	Other Tools for Working with Asymptotic Expansions . . . . .	26
1.10.3.1	Inversion of Asymptotic Series . . . . .	26
1.11	Frequentist Variational Approximations . . . . .	27
1.11.1	Thouless-Anderson-Palmer Variational Approach . . . . .	29
<b>2</b>	<b>Preliminary Lemmas and Their Proofs</b>	<b>31</b>
2.1	Lemma 1 . . . . .	31
2.2	Lemma 2 . . . . .	32
2.3	Lemma 3 . . . . .	33
2.4	Appendix . . . . .	34
2.4.1	Proof of Lemma 1 . . . . .	34
2.4.2	Proof of Lemma 2 . . . . .	35
2.4.2.1	A Fundamental Inequality for the Spectral Norm of a Vectorised Matrix . . . . .	35
2.4.2.2	Notational Definitions . . . . .	35
2.4.2.3	Derivation of (2.7) . . . . .	36
2.4.2.4	Expression for (2.6) with Lagrange Form of Remainder . . . . .	37
2.4.2.5	Spectral Norm Bounding of (2.9) . . . . .	38
2.4.2.6	Strategy for Proving (2.10) . . . . .	38
2.4.2.7	Proof of Result (2.16) . . . . .	39
2.4.2.8	Proof of Result (2.17) . . . . .	45
2.4.2.9	Summary of Moment Assumptions . . . . .	48
2.4.2.10	Succinct Expression for Moment Assumptions . . . . .	49
2.4.2.11	A Sufficient Condition for the Moment Assumptions . . . . .	50
2.4.3	Proof of Lemma 3 . . . . .	51
2.4.3.1	Matrix Extension of Results Concerning Integrals of Half-Cauchy Forms . . . . .	51
2.4.3.2	Derivation of Integrand Expressions . . . . .	52
2.4.3.3	Succinct Expressions for the Components in (2.53) . . . . .	59
2.4.3.4	Simplification of Integrals . . . . .	60
2.4.3.5	Explicit Expressions for (2.54) . . . . .	60
2.4.3.6	Convergence in Probability Limits of the Functions in (2.55) . . . . .	60
2.4.4	Multivariate Integral Limits for the Matrix Square Root Result . . . . .	61
2.4.4.1	Overview of this Appendix . . . . .	61
2.4.4.2	Computing Spectral Norms . . . . .	61

2.4.4.3	Verifying Convergence in Probability Limits of the Functions in (2.55) . . . . .	68
2.4.4.4	Conclusion for Multivariate Integral Limits for the Matrix Square Root Result . . . . .	71
<b>3</b>	<b>Usable Asymptotic Normality Results and Inference for Gaussian Response Linear Mixed Models</b>	<b>72</b>
3.1	Model Description . . . . .	73
3.2	Notation Required for Fisher Information Calculations . . . . .	74
3.3	Asymptotic Normality Theorem . . . . .	74
3.4	Appendix . . . . .	75
3.4.1	Linear Mixed Models with Multivariate Fixed and Random Effects	75
3.4.2	Expression for Top Left Block of Fisher Information Matrix . . .	77
3.4.2.1	Top Left Block of (3.5) . . . . .	77
3.4.2.2	Top Right Block of (3.5) . . . . .	77
3.4.2.3	Bottom Left Block of (3.5) . . . . .	78
3.4.2.4	Bottom Right Block of (3.5) . . . . .	78
3.4.3	Expression for Bottom Right Block of Fisher Information Matrix	82
3.4.3.1	Top Left Block of (3.10) . . . . .	82
3.4.3.2	Top Right Block of (3.10) . . . . .	83
3.4.3.3	Bottom Left Block of (3.10) . . . . .	84
3.4.3.4	Bottom Right Block of (3.10) . . . . .	84
3.4.4	The Inverse of the Fisher Information Matrix . . . . .	85
3.4.4.1	Expression for Top Left Block of Inverse Fisher Information Matrix . . . . .	85
3.4.4.2	Expression for Bottom Right Block of Inverse Fisher Information Matrix . . . . .	87
3.4.5	Derivation of the Final Asymptotic Normality Result for Gaussian Response Linear Mixed Models . . . . .	90
<b>4</b>	<b>Usable Asymptotic Normality Results and Inference for Generalized Linear Mixed Models</b>	<b>93</b>
4.1	Model Description . . . . .	94
4.2	Notation . . . . .	95
4.3	Asymptotic Normality Theorem . . . . .	96
4.4	Dispersion Parameter Extension . . . . .	97
4.5	Appendix . . . . .	98
4.5.1	Multivariate Extension of (2.6) of Tierney et al. (1989) . . . . .	98
4.5.1.1	Overview . . . . .	98
4.5.1.2	Multivariate Derivative Notation . . . . .	99
4.5.1.3	Check of the Miyata (2004) Appendix A Result for the Univariate Case . . . . .	99
4.5.1.4	The Multivariate Case . . . . .	100
4.5.1.5	Final Expression for the Multivariate Extension of (2.6) of Tierney et al. (1989) . . . . .	101
4.5.2	Proof of Theorem 12 . . . . .	101
4.5.2.1	Constructing the Fisher Information Matrix . . . . .	101



4.5.2.2	Expression for Conditional Density Function . . . . .	102
4.5.2.3	Introduction of Useful Notation and its Properties . . .	103
4.5.2.4	Computing an Asymptotic Approximation for the First Entry in (4.7) . . . . .	105
4.5.2.5	Computing an Asymptotic Approximation for the Sec- ond Entry in (4.7) . . . . .	109
4.5.2.6	Computing an Asymptotic Approximation for the Third Entry in (4.7) . . . . .	112
4.5.2.7	The Quadratic Conditional Expectations of the Scores .	115
4.5.2.8	Treating the Leading Term of the (2,2)-Entry of the Fisher Information Matrix . . . . .	122
4.5.2.9	The Fisher Information Matrix . . . . .	123
4.5.2.10	The Inverse of the Fisher Information Matrix . . . . .	123
4.5.2.11	Derivation of the Final Asymptotic Normality Result for Generalized Response Linear Mixed Models . . . . .	126
4.5.3	The Reciprocal Dispersion Parameter Fisher Information Block for Gamma Responses . . . . .	128
4.5.3.1	The Conditional Density Function . . . . .	128
4.5.3.2	The Score of the Reciprocal Dispersion Parameter . . .	131
4.5.3.3	Computing the Fisher Information Block for the Recip- rocal Dispersion Parameter . . . . .	132
4.5.3.4	Asymptotic Normality and Variance Results for the Max- imum Likelihood Estimator of the Reciprocal Dispersion Parameter . . . . .	145
4.5.3.5	Asymptotic Normality and Variance Results for the Max- imum Likelihood Estimator of the Dispersion Parameter	145
<b>5</b>	<b>Consequences and Applications of Asymptotic Normality Results</b>	<b>146</b>
5.1	Asymptotically Valid Inference . . . . .	147
5.1.1	Construction of Asymptotically Valid Confidence Intervals . . . .	147
5.1.2	Simulation Study . . . . .	148
5.2	Approximate Optimal Design . . . . .	155
5.2.1	Background and Model Description . . . . .	155
5.2.2	Approximate Locally D-Optimal Design Determination . . . . .	156
5.2.3	Illustration of Theorem 13 . . . . .	158
5.3	Appendix . . . . .	159
5.3.1	Model Description . . . . .	159
5.3.2	Asymptotic Assumption for Support Point Sample Sizes . . . . .	161
5.3.3	Useful Notation . . . . .	161
5.3.4	Key Moment Results . . . . .	162
5.3.5	The Fisher Information Matrix . . . . .	164
5.3.6	The Asymptotic D-Optimality Criterion . . . . .	164
5.3.7	Alternative Final Asymptotic D-optimality Criterion . . . . .	166
5.3.8	Special Distribution Cases . . . . .	167
<b>6</b>	<b>Thouless-Anderson-Palmer Enhancement of Generalized Linear Mixed Models</b>	<b>169</b>

6.1	Model Description . . . . .	170
6.2	The Gaussian Variational Approximate Log-Likelihood . . . . .	170
6.3	Overview of Thouless-Anderson-Palmer Enhancement . . . . .	172
6.4	The Thouless-Anderson-Palmer Approximate Negative Log-Likelihood . . . . .	173
6.5	Thouless-Anderson-Palmer Enhancement for Poisson Generalized Linear Mixed Models . . . . .	175
6.5.1	The Gaussian Variational Approximate Log-Likelihood for Simu- lation Set-Up . . . . .	175
6.5.2	The Thouless-Anderson-Palmer Negative Approximate Log-Likelihood for Simulation Set-Up . . . . .	176
6.5.3	Optimisation Issues . . . . .	176
6.5.3.1	A Simplified Version of the Optimisation Problem . . . . .	177
6.5.3.2	Simplified Simulation Study . . . . .	178
6.5.3.3	Results and Conclusion . . . . .	179
6.5.4	Simulation Study . . . . .	181
6.6	Appendix . . . . .	185
6.6.1	Proof of Result 2 . . . . .	185
6.6.1.1	Main Quantity in Onsager's Correction Term . . . . .	185
6.6.1.2	An Explicit Expression for the First Term in (6.10) . . . . .	185
6.6.1.3	An Explicit Expression for the Second Term in (6.10) . . . . .	189
6.6.1.4	An Explicit Expression for the Third Term in (6.10) . . . . .	190
6.6.1.5	The Resultant Expression for the Main Quantity in the Onsager's Correction Term . . . . .	190
6.6.2	Expressing the Main Quantity in the Onsager's Correction Term Using Integral Families . . . . .	192
6.6.3	Proof of Result 3 . . . . .	193
6.6.3.1	Simplifications in Poisson Case . . . . .	193
<b>7</b>	<b>Extensions to Noncanonical Link Generalized Linear Mixed Models</b>	<b>195</b>
7.1	Asymptotic Normality Results Involving Noncanonical Links . . . . .	196
7.1.1	Model Description . . . . .	196
7.1.2	Notation . . . . .	197
7.1.3	Asymptotic Normality Theorem . . . . .	198
7.2	Thouless-Anderson-Palmer Approach Involving Noncanonical Links . . . . .	199
7.2.1	Model Description . . . . .	199
7.2.2	The Gaussian Variational Approximate Log-Likelihood . . . . .	200
7.2.3	Overview of Thouless-Anderson-Palmer Enhancement . . . . .	200
7.3	Appendix . . . . .	202
7.3.1	Constructing the Fisher Information Matrix . . . . .	202
7.3.2	Expression for Conditional Density Function . . . . .	202
7.3.3	Deriving Expressions for the Expectation and Variance of the Response Variable . . . . .	203
7.3.4	Introduction of Useful Notation and Its Properties . . . . .	205
7.3.5	Key Conditional Moment Results . . . . .	205
7.3.6	Computing an Asymptotic Approximation for the First Entry in (7.9) . . . . .	207

---

7.3.6.1	The First Term of the First Score . . . . .	210
7.3.6.2	The Other Terms of the First Score . . . . .	210
7.3.6.3	Overall Leading Term Expression for the First Score . .	210
7.3.7	Computing an Asymptotic Approximation for the Second Entry in (7.9) . . . . .	210
7.3.7.1	The First Term of the Second Score . . . . .	211
7.3.7.2	The Other Terms of the Second Score . . . . .	211
7.3.7.3	Overall Leading Term Expression for the Second Score	211
7.3.8	Computing an Asymptotic Approximation for the Third Entry in (7.9) . . . . .	211
7.3.8.1	The First Term of the Third Score . . . . .	212
7.3.8.2	The Other Terms of the Third Score . . . . .	212
7.3.8.3	Overall Leading Term Expression for the Third Score .	212
7.3.9	The Quadratic Conditional Expectations of the Scores . . . . .	212
7.3.9.1	The Conditional Expectation of the Square of the First Score . . . . .	213
7.3.9.2	The Conditional Expectation of the Square of the Second Score . . . . .	213
7.3.9.3	The Conditional Expectation of the Square of the Third Score . . . . .	215
7.3.10	The Fisher Information Matrix . . . . .	216
7.3.11	The Inverse of the Fisher Information Matrix . . . . .	218
7.3.12	Final Asymptotic Normality Result . . . . .	218
<b>8</b>	<b>Discussion and Conclusion</b>	<b>219</b>
<b>9</b>	<b>References</b>	<b>223</b>

# Abstract

Generalized linear mixed models are an essential group of models for analysing many present-day complex data sets, especially those that contain non-normal and correlated response data. Despite the large volume of research concerning this group of models, there is very little theory concerning the statistical properties of maximum likelihood estimators for generalized linear mixed models. Existing theoretical results available for the asymptotic variance-covariance matrix for such estimators contain limits and expectations over the response distribution, hence such results are not in ready-to-use forms when carrying out tasks such as constructing studentized confidence intervals or optimal design determination. In this thesis, we derive precise asymptotic results for likelihood-based generalized linear mixed model analysis. The novel asymptotic normality results are derived for both cases involving either a canonical or noncanonical link function. In our approach, we derive the exact leading term behaviour of the Fisher information matrix when both the number of groups and number of observations within each group diverge. This leads to asymptotic normality results with explicit and simple studentizable forms. The implications of these results in optimal design theory is also explored, leading to simpler and more direct determination of approximate locally D-optimal designs. Towards the end of this thesis, a Thouless-Anderson-Palmer approach is introduced for modern statistical inference for generalized linear mixed models. Such methods have proven to provide accurate approximations to problems arising in machine learning contexts. However, statistical applications such as generalized linear mixed model analysis have not been investigated. Thus, we derive results for implementing the Thouless-Anderson-Palmer frequentist variational approach to generalized linear mixed models and analyse the accuracy of its variational estimates.

# Chapter 1

## Introduction and Background

This chapter serves as an introduction for the thesis and also presents the necessary background information and theory required for the remaining chapters in this thesis. Note that the results, definitions and work presented in this chapter are not novel.

### 1.1 Introduction

In recent developments in statistical analysis, generalized linear mixed models (GLMMs) have become an essential group of models for analysing many present-day complex data sets, causing it to become a rapidly growing area of research. These models have been deemed useful and practical when accounting for overdispersion is necessary, a common occurrence when working with outcomes that have underlying Poisson or binomial distributions. In addition, GLMMs can also be applied widely in areas such as longitudinal data analysis and disease mapping (Breslow and Clayton, 1993).

A popular approach for fitting GLMMs is maximum likelihood estimation (MLE). Problems however arise as while the log-likelihood of a GLMM can be expressed mathematically, it involves integration over the random effects component of the model, which cannot be evaluated as closed form integrals. In some cases, to counter this hindrance caused by integral intractability and to accurately evaluate these integrals, standard quadrature techniques such as Gauss-Hermite quadrature can be used. Many software packages are also available to fit these GLMMs. One such example is the package `lme4` (Bates et al., 2015) in the R computing environment (R Core Team, 2022) and contains the function `glmer()` which is used to fit GLMMs.

While estimation by maximum likelihood for GLMMs is well and widely established, asymptotic normality results that can be used for practical purposes such as constructing confidence intervals and Wald tests via Studentization are currently unavailable in the existing generalized linear mixed model literature.

With regards to possible developments in methodology used for GLMMs, one may consider the Thouless-Anderson-Palmer (TAP) paradigm, which was developed in statistical physics literature for spin glass models (Thouless et al., 1977). Recently, there has been a realisation that it provides accurate approximations to problems arising in machine learning contexts. The TAP approach is also able to overcome issues involving intractable integrals, a common problem that hinders frequentist inference carried out on GLMMs. However, statistical applications such as longitudinal data analysis and multilevel models analysis, which may benefit from using TAP methodology, have not been investigated.

## 1.2 Thesis Aim

This thesis aims to address some of the gaps currently present in the immense statistical literature available for GLMMs in the frequentist setting. We firstly aim to develop novel asymptotic theory for maximum quasi-likelihood estimators for GLMMs. The goal is to develop such theory for both cases where one may choose to either use a canonical link or noncanonical link as part of the GLMM, dependent on the type of data being used. Following that, we aim to assess the efficacy of the confidence intervals constructed using the newly derived asymptotic normality distribution for maximum quasi-likelihood estimators for GLMMs. We also aim to explore applications for the asymptotic results, such as its implications in optimal design theory. Last but not least, we aim to develop theory for implementing the TAP frequentist variational approach to generalized linear mixed models and analyse the accuracy of its variational estimates.

## 1.3 Outline

Following the introductory chapter, Chapter 2 introduces lemmas that serve as essential statistical tools for carrying out the asymptotic derivations present in this thesis. Chapter 3 then starts off by dealing with the Gaussian generalized linear mixed model and presents a theorem concerning the joint asymptotic normality of all of the maximum

quasi-likelihood estimators for such a model. The theorem is then extended to the generalized linear mixed model case in Chapter 4. Consequences and applications of the novel asymptotic normality results are then investigated and discussed in Chapter 5. Chapter 6 then turns to newly developed methodology for GLMMs and develops theory for usage of the TAP variational method for GLMMs and investigates if there are improvements in the statistical approximations. Since only canonical links are considered in Chapters 4 and 6, Chapter 7 further extends the novel asymptotic results and TAP variational approximation derivations to cater for noncanonical links as well. The thesis ends off with a discussion and conclusion based on the work presented so far.

## 1.4 Matrix Theory

In this section, we present some background on the matrix theory required for the derivations present in this thesis.

### 1.4.1 Difference Between Two Matrix Inverses

For any two equal-sized invertible matrices  $\mathbf{A}$  and  $\mathbf{B}$ , we have,

$$\mathbf{A}^{-1} - \mathbf{B}^{-1} = \mathbf{A}^{-1}(\mathbf{B} - \mathbf{A})\mathbf{B}^{-1}. \quad (1.1)$$

### 1.4.2 Other Useful Matrix Identities

Let  $\mathbf{M}$  and  $\mathbf{N}$  be invertible square matrices of the same size. Using an iterative application of the Sherman-Morrison-Woodbury formula, we have,

$$(\mathbf{M} - \mathbf{N})^{-1} = \sum_{k=0}^{\infty} (\mathbf{M}^{-1}\mathbf{N})^k \mathbf{M}^{-1}, \quad (1.2)$$

for matrices such that the spectral radius of  $\mathbf{M}^{-1}\mathbf{N}$  is less than 1. Let  $\mathbf{A}$  and  $\mathbf{B}$  be invertible square matrices and  $\mathbf{I}$  be an identity matrix with all the matrices having the same size. Using (1.2) and setting  $\mathbf{M} = \mathbf{I}$  and  $\mathbf{N} = -\mathbf{B}^{-1}\mathbf{A}^{-1}$  results in the following

matrix identity

$$\begin{aligned}
 (\mathbf{I} + \mathbf{AB})^{-1} \mathbf{A} &= [(\mathbf{AB})\{(\mathbf{AB})^{-1} + \mathbf{I}\}]^{-1} \mathbf{A} \\
 &= (\mathbf{I} + \mathbf{B}^{-1} \mathbf{A}^{-1})^{-1} \mathbf{B}^{-1} \mathbf{A}^{-1} \mathbf{A} \\
 &= (\mathbf{I} + \mathbf{B}^{-1} \mathbf{A}^{-1})^{-1} \mathbf{B}^{-1} \\
 &= (\mathbf{I} + \mathbf{B}^{-1} \mathbf{A}^{-1} + \dots) \mathbf{B}^{-1}
 \end{aligned} \tag{1.3}$$

Let  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  be invertible square matrices and  $\mathbf{I}$  be an identity matrix with all the matrices having the same size. Using (1.2) and setting  $\mathbf{M} = \mathbf{I}$  and  $\mathbf{N} = -\mathbf{B}^{-1} \mathbf{A}^{-1}$  results in the following matrix identity

$$\begin{aligned}
 (\mathbf{I} + \mathbf{AB})^{-1} \mathbf{C} &= [(\mathbf{AB})\{(\mathbf{AB})^{-1} + \mathbf{I}\}]^{-1} \mathbf{C} \\
 &= (\mathbf{I} + \mathbf{B}^{-1} \mathbf{A}^{-1})^{-1} \mathbf{B}^{-1} \mathbf{A}^{-1} \mathbf{C} \\
 &= (\mathbf{I} + \mathbf{B}^{-1} \mathbf{A}^{-1} + \dots) \mathbf{B}^{-1} \mathbf{A}^{-1} \mathbf{C}.
 \end{aligned} \tag{1.4}$$

Let  $\mathbf{A}$  and  $\mathbf{B}$  be invertible square matrices and  $\mathbf{I}$  be an identity matrix with all the matrices having the same size. Using (1.2) and setting  $\mathbf{M} = \mathbf{A}$  and  $\mathbf{N} = -\mathbf{B}^{-1}$  results in the following matrix identity

$$\begin{aligned}
 \mathbf{B}(\mathbf{I} + \mathbf{AB})^{-1} &= (\mathbf{B}^{-1})^{-1} (\mathbf{I} + \mathbf{AB})^{-1} \\
 &= \{(\mathbf{I} + \mathbf{AB})(\mathbf{B}^{-1})\}^{-1} \\
 &= (\mathbf{B}^{-1} + \mathbf{A})^{-1} \\
 &= \{\mathbf{A} - (-\mathbf{B}^{-1})\}^{-1} \\
 &= \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B}^{-1} \mathbf{A}^{-1} + \dots
 \end{aligned} \tag{1.5}$$

### 1.4.3 Block Matrix Inversion

The following definition instructs how a block matrix can be inverted.

**Result 1.** *Consider the following matrix which has been partitioned into four blocks*

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}$$

where  $\mathbf{A}$  and  $\mathbf{D}$  are square blocks of arbitrary size and blocks  $\mathbf{B}$  and  $\mathbf{C}$  are conformable such that the matrix can be properly partitioned. The matrix can then be inverted



blockwise as follows

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1} & -\mathbf{A}^{-1}\mathbf{B}(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1} \\ -(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1} & (\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1} \end{bmatrix}. \quad (1.6)$$

By permuting the blocks, an equivalent result is

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & -(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} \end{bmatrix}. \quad (1.7)$$

Note that matrices  $\mathbf{A}$  and  $\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B}$  must be invertible when using the block matrix inversion result presented in (1.6). Likewise, matrices  $\mathbf{D}$  and  $\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}$  must be invertible when using the block matrix inversion result presented in (1.7).

#### 1.4.4 The vec and vech Operators

Let  $\mathbf{A}$  be a  $m \times n$  matrix and let  $a_{ij}$  represent the element in the matrix located in the  $i$ -th row and  $j$ -th column. For any matrix  $\mathbf{A}$ ,  $\text{vec}(\mathbf{A})$  is defined as the  $mn \times 1$  vector which is constructed from the columns in  $\mathbf{A}$  being stacked on top each other, one column after the other, from left to right. If  $\mathbf{A}$  is a square  $d \times d$  matrix, then  $\text{vech}(\mathbf{A})$  is defined as a  $\frac{1}{2}d(d+1) \times 1$  vector, where the entries including and below the diagonal of  $\mathbf{A}$ , are stacked on top each other, one column after the other, from left to right.

For example, if  $\mathbf{A}$  is a  $3 \times 3$  square matrix, then we have the following matrix,  $\text{vec}$  and  $\text{vech}$  operators, where

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

with

$$\text{vec}(\mathbf{A}) = \begin{bmatrix} a_{11} \\ a_{22} \\ a_{31} \\ a_{12} \\ a_{22} \\ a_{32} \\ a_{13} \\ a_{23} \\ a_{33} \end{bmatrix} \quad \text{and} \quad \text{vech}(\mathbf{A}) = \begin{bmatrix} a_{11} \\ a_{21} \\ a_{31} \\ a_{22} \\ a_{32} \\ a_{33} \end{bmatrix}.$$

#### 1.4.4.1 The Commutation Matrix

The *commutation matrix* of order  $d$ , denoted by  $\mathbf{K}_d$ , allows for the conversion between the operators  $\text{vec}(\mathbf{A})$  and  $\text{vec}(\mathbf{A}^T)$ . It is the  $d^2 \times d^2$  matrix containing only zeroes and ones such that

$$\mathbf{K}_d \text{vec}(\mathbf{A}) = \text{vec}(\mathbf{A}^T)$$

for all  $d \times d$  matrices  $\mathbf{A}$ . The following useful property regarding commutation matrices also exists (Magnus and Neudecker, 1999):

$$\mathbf{K}_d^T = \mathbf{K}_d^{-1} = \mathbf{K}_d. \quad (1.8)$$

#### 1.4.4.2 The Duplication Matrix

The *duplication matrix* of order  $d$ , denoted by  $\mathbf{D}_d$ , allows for the conversion between the operators  $\text{vec}(\mathbf{A})$  and  $\text{vech}(\mathbf{A})$ .  $\mathbf{D}_d$  is the unique  $d^2 \times \frac{1}{2}d(d+1)$  matrix of zeros and ones such that

$$\text{vec}(\mathbf{A}) = \mathbf{D}_d \text{vech}(\mathbf{A})$$

for all  $d \times d$  matrices  $\mathbf{A}$ . Finally, the Moore-Penrose inverse of  $\mathbf{D}_d$  is defined as

$$\mathbf{D}_d^+ = (\mathbf{D}_d^T \mathbf{D}_d)^{-1} \mathbf{D}_d^T.$$

### 1.4.5 Kronecker Products and Related Properties

Let  $\mathbf{A}$  be a  $m \times n$  matrix and let  $\mathbf{B}$  be a  $p \times q$  matrix. The *Kronecker product* of matrices  $\mathbf{A}$  and  $\mathbf{B}$  is denoted as  $\mathbf{A} \otimes \mathbf{B}$  and it is the  $mp \times nq$  matrix defined by

$$\begin{bmatrix} a_{11}\mathbf{B} & \dots & a_{1n}\mathbf{B} \\ \vdots & & \vdots \\ a_{m1}\mathbf{B} & \dots & a_{mn}\mathbf{B} \end{bmatrix}.$$

Let  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  and  $\mathbf{D}$  be square matrices. Then some of the properties regarding Kronecker products are as follows:

$$\begin{aligned} \mathbf{A} \otimes (\mathbf{B} + \mathbf{C}) &= \mathbf{A} \otimes \mathbf{B} + \mathbf{A} \otimes \mathbf{C}, \\ (\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) &= \mathbf{AC} \otimes \mathbf{BD}, \\ (\mathbf{A} \otimes \mathbf{B})^{-1} &= \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}, \\ \text{vec}(\mathbf{ABC}) &= (\mathbf{C}^T \otimes \mathbf{A})\text{vec}(\mathbf{B}) \quad \text{and} \\ \text{tr}(\mathbf{ABCD}) &= \text{vec}(\mathbf{D})^T (\mathbf{A} \otimes \mathbf{C}^T) \text{vec}(\mathbf{B}^T). \end{aligned}$$

### 1.4.6 Vector and Matrix Norms

In this subsection, we present a few vector and matrix norms and their properties.

#### 1.4.6.1 Euclidean Norm

Let  $\mathbf{x} \in \mathbb{R}^d$ . Then, the *Euclidean Norm* of  $\mathbf{x}$  is

$$\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}}.$$

#### 1.4.6.2 Frobenius Norm

The *Frobenius norm* of a general matrix  $\mathbf{A}$  is

$$\|\mathbf{A}\|_F \equiv \sqrt{\text{trace}(\mathbf{A}^T \mathbf{A})}.$$

### 1.4.6.3 Spectral Norm

For a general symmetric matrix  $\mathbf{M}$ , let us define the following:

$\lambda_{\min}(\mathbf{M}) \equiv$  smallest eigenvalue of  $\mathbf{M}$  and  $\lambda_{\max}(\mathbf{M}) \equiv$  largest eigenvalue of  $\mathbf{M}$ .

The *spectral norm* of a general matrix  $\mathbf{A}$  is such that

$$\|\mathbf{A}\|_s \equiv \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}.$$

If  $\mathbf{A}$  is symmetric then

$$\|\mathbf{A}\|_s = \lambda_{\max}(\mathbf{A}).$$

Also, if  $\mathbf{A}$  is symmetric and positive definite, then the spectral decomposition of  $\mathbf{A}$  is

$$\mathbf{A} = \mathbf{U} \text{diag}(\boldsymbol{\lambda}) \mathbf{U}^T$$

where  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$  and  $\boldsymbol{\lambda}$  is the vector containing the eigenvalues of  $\mathbf{A}$ . We then have

$$\mathbf{A}^{-1} = \mathbf{U} \text{diag}(\mathbf{1}/\boldsymbol{\lambda}) \mathbf{U}^T$$

and therefore

$$\|\mathbf{A}^{-1}\|_s = 1/\lambda_{\min}(\mathbf{A}).$$

The spectral norm also possesses the following sub-multiplicity property

$$\|\mathbf{AB}\|_s \leq \|\mathbf{A}\|_s \|\mathbf{B}\|_s$$

for any pair of matrices  $\mathbf{A}$  and  $\mathbf{B}$  such that the matrix product  $\mathbf{AB}$  is defined.

Finally, suppose that  $\mathbf{A}$  is a  $d \times d$  matrix and  $\mathbf{1}_d$  is the vector of ones. Then, using the sub-multiplicity property of the spectral norm, we can claim that

$$\|\mathbf{1}_d^T \mathbf{A} \mathbf{1}_d\|_s \leq \|\mathbf{1}_d^T\|_s \|\mathbf{A}\|_s \|\mathbf{1}_d\|_s.$$

Then

$$\|\mathbf{1}_d^T\|_s = \sqrt{\text{largest eigenvalue of } \mathbf{1}_d \mathbf{1}_d^T} = \sqrt{d}$$

and

$$\|\mathbf{1}_d\|_s = \sqrt{\text{largest eigenvalue of } \mathbf{1}_d^T \mathbf{1}_d} = \sqrt{d}.$$

Hence

$$\|\mathbf{1}_d^T \mathbf{A} \mathbf{1}_d\|_s \leq d \|\mathbf{A}\|_s.$$

### 1.4.7 Eigenvalue Bound Results

Let the eigenvalues of a  $d \times d$  matrix  $\mathbf{M}$  be denoted by

$$\lambda_1(\mathbf{M}), \dots, \lambda_d(\mathbf{M}).$$

Theorem 8.1.5 of Golub and Van Loan (2013) states that, for any  $d \times d$  matrices  $\mathbf{A}$  and  $\mathbf{E}$  such that  $\mathbf{A}$  and  $\mathbf{A} + \mathbf{E}$  are symmetric, we have

$$\lambda_j(\mathbf{A}) + \lambda_{\min}(\mathbf{E}) \leq \lambda_j(\mathbf{A} + \mathbf{E}) \leq \lambda_j(\mathbf{A}) + \lambda_{\max}(\mathbf{E}) \quad \text{for all } 1 \leq j \leq d.$$

In particular, by choosing  $\lambda_j$  to correspond to  $\lambda_{\min}$ , we have

$$\lambda_{\min}(\mathbf{A} + \mathbf{E}) \geq \lambda_{\min}(\mathbf{A}) + \lambda_{\min}(\mathbf{E}). \quad (1.9)$$

#### 1.4.7.1 Matrix Identities from Harville (1977)

For a general linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha} + \boldsymbol{\varepsilon},$$

assume that  $\mathbf{y}$  is a  $n \times 1$  response vector,  $\mathbf{X}$  and  $\mathbf{Z}$  are  $n \times p$  and  $n \times q$  matrices respectively,  $\boldsymbol{\beta}$  is a  $p \times 1$  vector of unobservable fixed effects,  $\boldsymbol{\alpha}$  is a  $q \times 1$  vector of unobservable random effects and  $\boldsymbol{\varepsilon}$  is a  $n \times 1$  vector of unobservable random errors. In addition, the following properties apply where  $E(\boldsymbol{\alpha}) = \mathbf{0}$ ,  $E(\boldsymbol{\varepsilon}) = \mathbf{0}$  and  $E(\boldsymbol{\alpha}^T \boldsymbol{\varepsilon}) = \mathbf{0}$ . Also let  $\mathbf{D} = \text{Cov}(\boldsymbol{\alpha})$ ,  $\mathbf{R} = \text{Cov}(\boldsymbol{\varepsilon})$  and  $\mathbf{V} = \mathbf{Z}\mathbf{D}\mathbf{Z}^T + \mathbf{R}$  such that  $\text{Cov}(\mathbf{y}) = \mathbf{V}$ , where  $\text{Cov}(\cdot)$  denotes the covariance function. Harville (1977) then provides the following matrix identities.

$$\mathbf{V}^{-1} \equiv \mathbf{R}^{-1} - \mathbf{R}^{-1} \mathbf{Z} \mathbf{D} (\mathbf{I} + \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{Z} \mathbf{D})^{-1} \mathbf{Z}^T \mathbf{R}^{-1}, \quad (1.10a)$$

$$\mathbf{Z}^T \mathbf{V}^{-1} \equiv (\mathbf{I} + \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{Z} \mathbf{D})^{-1} \mathbf{Z}^T \mathbf{R}^{-1}. \quad (1.10b)$$

### 1.4.8 Vector Differential Calculus

Let  $f$  be a scalar-valued function with a vector  $\mathbf{x} \in \mathbb{R}^d$  as its argument. Then the  $1 \times d$  *derivative vector* of  $f$ , denoted by  $\mathbf{D}f(\mathbf{x})$ , has entries

$$\frac{\partial f(\mathbf{x})}{\partial x_i}$$

where  $1 \leq i \leq d$ . The  $d \times d$  *Hessian matrix* of  $f$  is then defined as follows

$$\mathbf{H}f(\mathbf{x}) = \mathbf{D}\{\mathbf{D}f(\mathbf{x})^T\}$$

with  $(i, j)$  entry equal to

$$\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j}$$

where  $1 \leq i \leq d$  and  $1 \leq j \leq d$ . When vectors and matrices are involved, it is more appropriate to use vector differential calculus rather than ordinary scalar differential calculus. In order to compute the derivative vector, we use the following definition (Magnus and Neudecker, 1999):

**Theorem 1.** *First Identification Theorem: If  $\mathbf{a}$  and  $\mathbf{x}$  are  $1 \times d$  vectors such that*

$$df(\mathbf{x}) = \mathbf{a} d\mathbf{x}$$

*then*

$$\mathbf{a} = \mathbf{D}f(\mathbf{x}).$$

The Hessian matrix can then be computed as follows:

**Theorem 2.** *Second Identification Theorem: If  $\mathbf{x}$  is a  $1 \times d$  vector and  $\mathbf{A}$  is a  $d \times d$  matrix such that*

$$d^2 f(\mathbf{x}) = (d\mathbf{x})^T \mathbf{A} d\mathbf{x}$$

*then*

$$\mathbf{A} = \mathbf{H}f(\mathbf{x}).$$

Further rules regarding vector differential calculus are provided in Magnus and Neudecker (1999) and Wand (2002).

## 1.5 Key Integral Results

In this section, we present key integral results used in this thesis.

### 1.5.1 Useful Integral Results

Let

$$a_1, a_2, a_3, b > 0$$

be strictly positive real numbers such that

$$a_1 a_2 > b.$$

Then, from Wolfram Research Inc. (2022), we have the following integral results:

$$\begin{aligned} \int_0^\infty \frac{dx}{a_1 + x^2} &= \frac{\pi}{2\sqrt{a_1}}, \\ \int_0^\infty \frac{dx}{(a_1 + x^2)^2} &= \frac{\pi}{4a_1\sqrt{a_1}}, \\ \int_0^\infty \frac{dx}{(a_1 + x^2)(a_2 + x^2)} &= \frac{\pi}{2\sqrt{a_1 a_2}(\sqrt{a_1} + \sqrt{a_2})}, \\ \int_0^\infty \frac{dx}{(a_1 + x^2)^2(a_2 + x^2)} &= \frac{\pi(2\sqrt{a_1} + \sqrt{a_2})}{4a_1\sqrt{a_1 a_2}(\sqrt{a_1} + \sqrt{a_2})^2}, \\ \int_0^\infty \frac{x^2 dx}{(a_1 + x^2)(a_2 + x^2)} &= \frac{\pi}{2(\sqrt{a_1} + \sqrt{a_2})}, \\ \int_0^\infty \frac{dx}{(a_1 + x^2)(a_2 + x^2)(a_3 + x^2)} &= \frac{\pi(\sqrt{a_1} + \sqrt{a_2} + \sqrt{a_3})}{2\sqrt{a_1 a_2 a_3}(\sqrt{a_1} + \sqrt{a_2})(\sqrt{a_1} + \sqrt{a_3})(\sqrt{a_2} + \sqrt{a_3})} \end{aligned}$$

and

$$\begin{aligned} &\int_0^\infty \frac{x^2 dx}{(a_1 + x^2)(a_2 + x^2) - b} \\ &= \frac{\pi}{\sqrt{2} \left\{ \sqrt{a_1 + a_2 + \sqrt{(a_1 - a_2)^2 + 4b}} + \sqrt{a_1 + a_2 - \sqrt{(a_1 - a_2)^2 + 4b}} \right\}}. \end{aligned}$$

Also note that,

$$\frac{\pi(2\sqrt{a_1} + \sqrt{a_2})}{4a_1\sqrt{a_1 a_2}(\sqrt{a_1} + \sqrt{a_2})^2} < \frac{2\pi(\sqrt{a_1} + \sqrt{a_2})}{4a_1\sqrt{a_1 a_2}(\sqrt{a_1} + \sqrt{a_2})^2} = \frac{\pi}{2a_1\sqrt{a_1 a_2}(\sqrt{a_1} + \sqrt{a_2})}$$

which leads to the bound

$$\int_0^\infty \frac{dx}{(a_1 + x^2)^2(a_2 + x^2)} < \frac{\pi}{2a_1\sqrt{a_1a_2}(\sqrt{a_1} + \sqrt{a_2})}. \quad (1.11)$$

### 1.5.2 Integral Form of the Matrix Square Root

The integral form of the square root of a matrix, for a matrix  $\mathbf{A}$  having no eigenvalues on  $\mathbb{R}^-$ , is given by Higham (2008) as

$$\mathbf{A}^{1/2} = \frac{2}{\pi} \int_0^\infty \mathbf{A}(\mathbf{A} + t^2\mathbf{I})^{-1} dt. \quad (1.12)$$

## 1.6 Key Expectation Results

In this section, we present some key expectation results used frequently in the derivations in this thesis.

### 1.6.1 Law of Total Expectation

**Theorem 3.** *Let  $X$  and  $Y$  be random variables defined on the same probability space and where the expected value of  $X$ ,  $E(X)$ , is defined. Then the law of total expectation is defined as follows*

$$E(X) = E\{E(X|Y)\}.$$

### 1.6.2 Jensen's Inequality

**Theorem 4.** *Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a convex function and  $X$  be a random variable such that its expected value,  $E(X)$ , is finite. Then, Jensen's inequality states the following*

$$f(E(X)) \leq E(f(X)).$$



### 1.6.3 Markov's Inequality

**Theorem 5.** *Let  $X$  be a random variable and  $a$  be a scalar. Markov's inequality states that if  $X$  is non-negative and  $a > 0$ , then the probability that  $X$  is greater than or equal to  $a$  is at most the expectation of  $X$  divided by  $a$ ,*

$$P(X \geq a) \leq \frac{E(X)}{a}.$$

### 1.6.4 Cauchy-Schwarz Inequality

**Theorem 6.** *Let  $X$  and  $Y$  be random variables. Then the Cauchy-Schwarz inequality is defined as follows*

$$|E(XY)|^2 \leq E(X^2)E(Y^2).$$

## 1.7 Exponential Families

*Exponential families* consist of a set of probability distributions that can be written in certain parametric forms. These exponential family forms can be used to provide several parametric families of distributions with alternative parametrizations, in terms of natural parameters, which possess useful algebraic properties. Several commonly used probability distributions fall under the umbrella of exponential family density functions, including but not restricted to the normal, binomial, Poisson and gamma distributions.

### 1.7.1 One-Parameter Exponential Families

In this subsection, we present the class of *one-parameter exponential family* probability distributions.

**Definition 1.** *The class of one-parameter exponential family density, or probability mass, functions have generic form*

$$p(y; \eta) = \exp \{y\eta - b(\eta) + c(y)\} h(y) \tag{1.13}$$

where  $\eta$  is the natural parameter and the functions  $b(\cdot)$ ,  $c(\cdot)$  and  $h(\cdot)$  are defined according to the desired response distribution.

Explicit examples of the functions  $b(\cdot)$ ,  $c(\cdot)$  and  $h(\cdot)$  for the binomial and Poisson family of distributions have been provided in Table 1.1.

Family	$b(\eta)$	$c(y)$	$h(y)$
Binomial	$\log(1 + e^\eta)$	0	$I(y \in \{0, 1\})$
Poisson	$e^\eta$	$-\log(y!)$	$I(y \in \{0\} \cup \mathbb{N})$

Table 1.1: Examples of one-parameter exponential families and their  $b$ ,  $c$  and  $h$  functions.

Here,  $I(\mathcal{P}) = 1$  if the condition  $\mathcal{P}$  is true and  $I(\mathcal{P}) = 0$  if  $\mathcal{P}$  is false.

If the random variable  $Y$  has density, or probability mass, function as in (1.13), then  $E(Y) = b'(\eta)$  and  $\text{Var}(Y) = b''(\eta)$ . A common modelling extension is to account for overdispersion. Overdispersion occurs when the variability present in the data is larger than what the proposed statistical model can account for. To model the variance flexibly, a dispersion parameter  $\phi > 0$  is introduced and  $\log\{p(y; \eta)\}$  is replaced by a *quasi-likelihood* function as shown in Definition 2.

**Definition 2.** *When accounting for overdispersion, the quasi-likelihood for the class of one-parameter exponential family density, or probability mass, functions have generic form*

$$\{y\eta - b(\eta) + c(y)\}/\phi + d(y, \phi) \quad (1.14)$$

where  $\eta$  is the natural parameter,  $\phi > 0$  is the dispersion parameter and the functions  $b(\cdot)$ ,  $c(\cdot)$  and  $d(\cdot, \cdot)$  are defined according to the desired response distribution.

Now, if the random variable  $Y$  has quasi-likelihood as in (1.14), then  $E(Y) = b'(\eta)$  and  $\text{Var}(Y) = \phi b''(\eta)$ .

Now, let  $\mathbf{Y}$  be a vector of independent observations,  $\mathbf{X}$  be a matrix of known covariates and  $\boldsymbol{\beta}$  be a vector of unknown regression coefficients (fixed effects). In generalized linear models, it is common to model the mean function  $\boldsymbol{\mu} = E(\mathbf{Y})$  as some non-linear function of the linear predictor or natural parameter such that

$$g(\boldsymbol{\mu}) = \boldsymbol{\eta}$$

where

$$\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}.$$

Here,  $g$  is known as the *link* function. If  $g^{-1} = b'$ , then  $g$  is called the *canonical link*

function and the following useful relationship exists

$$\boldsymbol{\mu} = g^{-1}(\boldsymbol{\eta}) = b'(\boldsymbol{\eta}).$$

Selecting the link function to be a canonical link leads to simpler likelihood expressions and provides useful sufficient statistics.

However, one may choose to use a *noncanonical link* function if, for example, using a noncanonical link leads to a better data fit. Some examples of noncanonical links are the probit link ( $\Phi^{-1}$ ) for binary regression and the log link (log) for Gamma regression. The following two definitions provide exponential family forms written in terms of the natural parameter,  $\eta$ , when noncanonical links are used (Fan et al., 1995).

**Definition 3.** *When using noncanonical links, the class of one-parameter exponential family density, or probability mass, functions have a generic form in terms of  $\eta$  as follows*

$$p(y; \eta) = \exp [y(g \circ b')^{-1}(\eta) - \{b \circ (g \circ b')^{-1}\}(\eta) + c(y)] h(y) \quad (1.15)$$

where  $\eta$  is the natural parameter,  $g$  is the link function and the functions  $b(\cdot)$ ,  $c(\cdot)$  and  $h(\cdot)$  are defined according to the desired response distribution.

If the random variable  $Y$  has density, or probability mass, function as in (1.15), then  $E(Y) = g^{-1}(\eta)$  and  $\text{Var}(Y) = \{b'' \circ (b')^{-1} \circ g^{-1}\}(\eta)$ . As mentioned earlier in this section, one can apply a modelling extension to account for overdispersion which leads to the expression for  $\log\{p(y; \eta)\}$  as in (1.15) being replaced by a quasi-likelihood function as shown in Definition 4.

**Definition 4.** *When accounting for overdispersion and using noncanonical links, the quasi-likelihood for the class of one-parameter exponential family density, or probability mass, functions have a generic form in terms of  $\eta$  as follows*

$$[y(g \circ b')^{-1}(\eta) - \{b \circ (g \circ b')^{-1}\}(\eta) + c(y)] / \phi + d(y, \phi) \quad (1.16)$$

where  $\eta$  is the natural parameter,  $\phi > 0$  is the dispersion parameter and the functions  $b(\cdot)$ ,  $c(\cdot)$  and  $d(\cdot, \cdot)$  are defined according to the desired response distribution.

Now, if the random variable  $Y$  has quasi-likelihood as in (1.16), then  $E(Y) = b'(\eta)$  and  $\text{Var}(Y) = \phi\{b'' \circ (b')^{-1} \circ g^{-1}\}(\eta)$ .

### 1.7.2 Two-Parameter Exponential Families

In this subsection, we present the class of *two-parameter exponential family* probability distributions.

**Definition 5.** *The class of two-parameter exponential family density, or probability mass, functions have generic form*

$$p(y; \eta, \phi) = \exp [\{y\eta - b(\eta) + c(y)\} / \phi + d(y, \phi)] h(y)$$

where  $\eta$  is the natural parameter and the functions  $b(\cdot)$ ,  $c(\cdot)$ ,  $d(\cdot, \cdot)$  and  $h(\cdot)$  are defined according to the desired response distribution.

Explicit examples of the functions  $b(\cdot)$ ,  $c(\cdot)$ ,  $d(\cdot, \cdot)$  and  $h(\cdot)$  for the Gaussian and gamma family of distributions have been provided in Table 1.2.

Family	$b(\eta)$	$c(y)$	$d(y, \phi)$	$h(y)$
Gaussian	$\frac{1}{2}\eta^2$	$-\frac{1}{2}y^2$	$-\log(2\pi\phi)$	1
Gamma	$-\log(-\eta)$	$\log(y)$	$-\log(\phi\Gamma(1/\phi)) - \log(y)$	$I(y > 0)$

Table 1.2: Examples of two-parameter exponential families and their  $b, c, d$  and  $h$  functions.

## 1.8 Generalized Linear Mixed Models

We start off with classical linear models, where the mean of the response, often required to be normally distributed, can be expressed as a linear combination of the unknown model parameters and the predictor variables. In other words, a linear model has the mean

$$E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$$

where  $\mathbf{Y}$  is a vector of independent observations,  $\mathbf{X}$  is a matrix of known covariates and  $\boldsymbol{\beta}$  is a vector of unknown regression coefficients (fixed effects). These models, however, fall short when the observations are correlated or when the mean of the response cannot be written as a linear function of the covariates.

*Generalized linear mixed models* serve as an extension of linear models in two distinct ways. Firstly, GLMMs allow for the modelling of correlated data through the inclusion of random effects. Secondly, the mean,  $\mu$ , is linked to the linear predictor through

a known function,  $g$ , known as the link function. With these additional properties, GLMMs have become an essential group of models for analysing many present-day complex data sets, which contain non-normal and correlated response data.

A summary of the various linear model effects structures ranging from linear models to generalized linear mixed models is further detailed in Table 1.3.

Model	(conditional) Mean Response
Linear Models (LMs)	$\mathbf{X}\boldsymbol{\beta}$
Linear Mixed Models (LMMs)	$\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{U}, \quad \mathbf{U} \sim (\mathbf{0}, \boldsymbol{\Sigma})$
Generalized Linear Models (GLMs)	$g^{-1}(\mathbf{X}\boldsymbol{\beta})$
Generalized Linear Mixed Models (GLMMs)	$g^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{U}), \quad \mathbf{U} \sim (\mathbf{0}, \boldsymbol{\Sigma})$

Table 1.3: Summary of various linear model effects structures.

To specify the structure of a generalized linear mixed model, we first define the conditional distribution of the response,  $Y_{ij}$ , given its associated random effect  $\mathbf{U}_i$ . Let there be  $m$  groups and  $n_i$  observations within each group. We also assume the random effects to be independent normally distributed variables and the response  $Y_{ij}$ , conditional on the random effects, to be from an exponential family  $f$ . Then the generalized linear mixed model has the following generic form:

$$Y_{ij}|\mathbf{U}_i \stackrel{\text{ind.}}{\sim} f_{Y_{ij}|\mathbf{U}_i}(y_{ij}|\mathbf{u}_i), \quad \mathbf{U}_i \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma}) \quad (1.17)$$

where  $\stackrel{\text{ind.}}{\sim}$  means ‘independently distributed as’ and with natural parameter

$$\eta_{ij} = \mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{U}_i,$$

for  $1 \leq i \leq m$  and  $1 \leq j \leq n_i$ . Here,  $\mathbf{X}_{ij}$  is a  $d_F \times 1$  vector of predictors having a fixed effects coefficient vector  $\boldsymbol{\beta}$  and  $\mathbf{Z}_{ij}$  is a  $d_R \times 1$  vector of predictors having a  $d_R \times 1$  random effects coefficient vector  $\mathbf{U}_i$ . For this generalized linear mixed model, the conditional mean of  $Y_{ij}$  is

$$E(Y_{ij}|\mathbf{U}_i) = \mu_{ij},$$

and there is a known link function,  $g$ , linking together the conditional mean and natural parameter such that

$$g(\mu_{ij}) = \eta_{ij} = \mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{U}_i.$$

Estimation of model parameters in generalized linear mixed models can be carried out using maximum likelihood estimation, with more details being provided in the next

section.

A detailed overview of the usefulness and difficulties of GLMM-based analysis can be found in ? and Jiang and Nguyen (2021). Inferential methods for GLMMs other than the maximum likelihood approach, such as generalized estimating equations and penalized quasi-likelihood have also been discussed in these books.

## 1.9 Maximum Likelihood for Generalized Linear Mixed Models

In this section, the *maximum likelihood* approach to estimating model parameters in a generalized linear mixed model is presented.

### 1.9.1 The Likelihood Function

We begin by describing the *likelihood function* for a general statistical model. Consider a statistical model, parametrized by a vector of model parameters  $\boldsymbol{\theta}$ , with probability density function  $f(\mathbf{y}; \boldsymbol{\theta})$ , where  $\mathbf{y}$  is a vector of random variables. Here, if  $\boldsymbol{\theta}$  is assumed to be known, then  $f(\mathbf{y}; \boldsymbol{\theta})$  is viewed to be the probability density function for  $\mathbf{y}$ . On the other hand, when  $\mathbf{y}$  represents a vector of known observations and  $\boldsymbol{\theta}$  is unknown, then  $f(\mathbf{y})$  is simply a function of  $\boldsymbol{\theta}$ . This is known as the likelihood function of  $\mathbf{y}$  and is usually represented as  $\mathcal{L}(\boldsymbol{\theta}; \mathbf{y})$  to emphasise that  $\boldsymbol{\theta}$  is unknown and  $\mathbf{y}$  is known. Note that mathematically, we have that

$$f(\mathbf{y}; \boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta}; \mathbf{y}).$$

Now let us consider the likelihood for a generic generalized linear mixed model. By letting  $\boldsymbol{\theta} = (\boldsymbol{\beta}, \boldsymbol{\Sigma})$  and using the model description in (1.17), the likelihood can be written as follows

$$\mathcal{L}(\boldsymbol{\beta}, \boldsymbol{\Sigma}; \mathbf{y}) = \prod_{i=1}^m \int_{\mathbb{R}^{d_R}} \prod_{j=1}^{n_i} f_{Y_{ij}|U_i}(y_{ij}|\mathbf{u}_i) f_{U_i}(\mathbf{u}_i) d\mathbf{u}_i. \quad (1.18)$$

These likelihood functions form the basis for maximum likelihood estimation which is explained in the next subsection.

### 1.9.2 Maximum Likelihood Estimation

The maximum likelihood estimation method estimates the values of model parameters such that under the fitted statistical model, the observed data is most probable. In order to find these maximum likelihood estimates, one would need to find the values of  $\boldsymbol{\theta}$  that maximise the likelihood  $\mathcal{L}(\boldsymbol{\theta}; \mathbf{y})$ , where the maximization is carried out within the permissible range of values for  $\boldsymbol{\theta}$ . For example, if one of the elements of  $\boldsymbol{\theta}$  represents a variance or covariance parameter, then its range of permissible values is restricted to non-negative values. This aspect of maximum likelihood estimation is critical for estimating variances and covariances of random effects variables.

Note that finding the values of  $\boldsymbol{\theta}$  that maximise the likelihood,  $\mathcal{L}(\boldsymbol{\theta}; \mathbf{y})$ , is equivalent to finding the values of  $\boldsymbol{\theta}$  that maximise the *log-likelihood*,  $\log \mathcal{L}(\boldsymbol{\theta}; \mathbf{y})$ , since the log function is a monotonic increasing function. The log-likelihood, commonly denoted as  $\ell(\boldsymbol{\theta})$ , is often a more convenient mathematical expression to work with. Hence the maximum likelihood estimator for  $\boldsymbol{\theta}^0$ , the true value of the parameter  $\boldsymbol{\theta}$ , in a general statistical model can now be expressed as follows

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}). \quad (1.19)$$

Now, we aim to define the maximum likelihood estimators for generalized linear mixed models. Taking the log function on both sides of (1.18), we obtain the log-likelihood of a generic generalized linear mixed model as follows

$$\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \sum_{i=1}^m \left( \log \int_{\mathbb{R}^{d_R}} \prod_{j=1}^{n_i} f_{Y_{ij}|\mathbf{U}_i}(y_{ij}|\mathbf{u}_i) f_{\mathbf{U}_i}(\mathbf{u}_i) d\mathbf{u}_i \right). \quad (1.20)$$

Then, for any  $\boldsymbol{\beta}$  ( $d_F \times 1$ ) and  $\boldsymbol{\Sigma}$  ( $d_R \times d_R$ ) that is symmetric and positive definite, the *maximum likelihood estimator* of  $(\boldsymbol{\beta}^0, \boldsymbol{\Sigma}^0)$ , the true values of the parameters  $\boldsymbol{\beta}$  and  $\boldsymbol{\Sigma}$ , is

$$(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\Sigma}}) = \operatorname{argmax}_{\boldsymbol{\beta}, \boldsymbol{\Sigma}} \ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}). \quad (1.21)$$

Note that *maximum quasi-likelihood estimators* for GLMMs can also be defined if one works with quasi-likelihoods such as those in (1.14) or (1.16).

### 1.9.3 Asymptotic Properties of Maximum Likelihood Estimators for Generalized Linear Mixed Models

Under certain regularity conditions, the maximum likelihood estimates of  $\boldsymbol{\theta}^0$  obtained using (1.19) are consistent and asymptotically normally distributed according to the theorem below (Knight, 2000; ?).

**Theorem 7.** *Let  $\mathbf{Y}$  be a vector of independently and identically distributed random variables and let  $\boldsymbol{\theta}$  denote the vector of model parameters used to parametrize a statistical model, such that the density function  $f(\mathbf{y}; \boldsymbol{\theta})$  satisfies the following regularity conditions:*

1. *The true value  $\boldsymbol{\theta}^0$  of  $\boldsymbol{\theta}$  is interior to the parameter space  $\Theta$ , which has finite dimension and is compact;*
2. *The set  $A = \{\mathbf{y} : f(\mathbf{y}; \boldsymbol{\theta}) > \mathbf{0}\}$  does not depend on  $\boldsymbol{\theta}$ ;*
3.  *$f(\mathbf{y}; \boldsymbol{\theta})$  is three times continuously differentiable with respect to  $\boldsymbol{\theta}$  for all  $\mathbf{y}$  in  $A$ ;*
4.  *$E[\ell'(\boldsymbol{\theta})] = \mathbf{0}$  for all  $\boldsymbol{\theta}$  and  $\text{Var}[\ell'(\boldsymbol{\theta})] = -E[\ell''(\boldsymbol{\theta})] = I(\boldsymbol{\theta})$  where  $\mathbf{0} < I(\boldsymbol{\theta}) < \infty$  for all  $\boldsymbol{\theta}$ ;*
5. *For each  $\boldsymbol{\theta}$  and  $\boldsymbol{\delta} > \mathbf{0}$ , there exists  $|\ell'''(\mathbf{t}; \mathbf{y})| < M(\mathbf{y})$  for  $|\boldsymbol{\theta} - \mathbf{t}| \leq \boldsymbol{\delta}$  where  $E[M(\mathbf{Y})] < \infty$ .*

*Then, the following asymptotic normality result for maximum likelihood estimators exists*

$$\widehat{\boldsymbol{\theta}} \stackrel{asy.}{\approx} N(\boldsymbol{\theta}^0, I(\boldsymbol{\theta}^0)^{-1}) \quad (1.22)$$

*where  $\stackrel{asy.}{\approx}$  means ‘asymptotically distributed as’ and with mean equal to the vector of true model parameters  $\boldsymbol{\theta}^0$  and asymptotic variance-covariance matrix equal to the inverse of the Fisher information matrix,  $I(\boldsymbol{\theta}^0)$ .*

As stated in the theorem above, there are two ways to derive the Fisher information matrix. The first approach involves computing the first derivative of  $\ell(\boldsymbol{\theta})$  with respect to  $\boldsymbol{\theta}$ . Then the Fisher information matrix can be defined as follows where

$$I(\boldsymbol{\theta}) = E \left\{ \frac{\partial \ell}{\partial \boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \boldsymbol{\theta}} \right)^T \right\}.$$

Alternatively, one could compute the second derivative of  $\ell(\boldsymbol{\theta})$  with respect to  $\boldsymbol{\theta}$  and



use the following definition of the Fisher information matrix where

$$I(\boldsymbol{\theta}) = -E \left\{ \frac{\partial^2 \ell}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right\}.$$

Despite the large volume of research concerning generalized linear mixed models, there is very little theory concerning the statistical properties of maximum likelihood estimators or maximum quasi-likelihood estimators for these models.

Recent related literature published includes Hall et al. (2011), who derive precise asymptotic normality results for estimators for models that fall under the generalized linear mixed models framework. In their case, these results are derived for Gaussian variational approximation (GVA) estimators for a single-predictor Poisson mixed model. Let  $m$  be the number of groups or subjects and let  $n$  be the number of observations within each group. The asymptotic results are derived for the case where both  $m$  and  $n$  in the model diverge. The final results obtained give rise to asymptotically valid statistical inference where Gaussian variational approximations are concerned. In this case, the aim of this thesis differs as we aim to derive asymptotic normality results for maximum likelihood estimators and for the general class of GLMMs.

Nie (2007) presents properties of maximum likelihood estimators in generalized linear and non-linear mixed effects models. In this article, the convergence rates of the asymptotic variances of these maximum likelihood estimators were investigated in three cases. In terms of the notation presented in Section 1.8, one of these cases caters to when both  $m$  and  $n$  tend to infinity while the other two cases concentrate on either  $m$  or  $n$  diverging towards infinity while the other quantity remains finite. By using the Fisher information matrix, the convergence rates of the MLEs were determined by finding out the orders of the leading terms and the remainder terms in the asymptotic variances for each estimator. However, the leading terms involved were not explicitly derived. This leaves a gap in terms of deriving the asymptotic distributions explicitly for such GLMMs by investigating the exact expressions of the leading terms involved in the asymptotic variances of the estimators.

The work presented in this thesis addresses this gap in the current statistical literature for generalized linear mixed models.

## 1.10 Asymptotics

In this section, we present the statistical tools used to carry out the asymptotic derivations present in this thesis.

### 1.10.1 Convergence of Random Variables

In this subsection, we look at two different types of convergence for sequences of random variables. We also consider the properties of these sequences of random variables when various algebraic operations are applied.

#### 1.10.1.1 Convergence in Probability

The notion of *convergence in probability* for a sequence of random variables deals with the convergence of the random variables themselves and is defined below,

**Definition 6.** Let  $\{X_n\}$  be a sequence of random variables. Then  $\{X_n\}$  converges in probability to the random variable  $X$  as  $n \rightarrow \infty$ , or  $X_n \xrightarrow{P} X$ , if for all  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P(|X_n - X| > \varepsilon) = 0.$$

It is common for the limiting random variable  $X$  to be a constant  $c$ , for which we then have  $X_n \xrightarrow{P} c$ .

#### 1.10.1.2 Convergence in Distribution

The notion of *convergence in distribution* for a sequence of random variables deals with the convergence of the distribution functions of the random variables and is defined below,

**Definition 7.** Let  $\{X_n\}$  be a sequence of random variables. Then  $\{X_n\}$  converges in distribution to the random variable  $X$  as  $n \rightarrow \infty$ , or  $X_n \xrightarrow{D} X$ , if

$$\lim_{n \rightarrow \infty} P(X_n \leq x) = P(X \leq x) = F(x)$$

for each point  $x \in \mathbb{R}$  at which  $F(x)$  is continuous, where  $F(x)$  is the cumulative distribution function of the random variable  $X$ .

### 1.10.1.3 Continuous Mapping Theorem

In probability theory, the *continuous mapping theorem* states that continuous functions of sequences of random variables preserve limits. A formal definition is provided below.

**Theorem 8.** *Let  $X$  be a random variable and  $\{X_n\}$  be a sequence of random variables. If  $g$  is a continuous function, then the continuous mapping theorem states the following*

$$X_n \xrightarrow{\mathcal{D}} X \text{ implies } g(X_n) \xrightarrow{\mathcal{D}} g(X) \text{ and } X_n \xrightarrow{P} X \text{ implies } g(X_n) \xrightarrow{P} g(X).$$

### 1.10.1.4 Slutsky's Theorem

*Slutsky's theorem*, which is partly derived using the continuous mapping theorem, provides useful results when dealing with algebraic operations involving two sequences of random variables, where one sequence converges in distribution to a random variable while the other sequence converges in probability to a constant.

**Theorem 9.** *Let both  $\{X_n\}$  and  $\{Y_n\}$  be sequences of random variables. If  $X_n \xrightarrow{\mathcal{D}} X$  and  $Y_n \xrightarrow{P} c$ , then the following properties of algebraic operations involving both  $\{X_n\}$  and  $\{Y_n\}$  exist*

$$X_n + Y_n \xrightarrow{\mathcal{D}} X + c, \quad X_n Y_n \xrightarrow{\mathcal{D}} Xc \quad \text{and} \quad X_n / Y_n \xrightarrow{\mathcal{D}} X/c.$$

### 1.10.1.5 Cramér-Wold Device

The *Cramér-Wold device* is a useful result that can be used to prove the joint convergence of random variables. A formal definition is provided below.

**Theorem 10.** *Let  $\mathbf{X}$  be a random variable and  $\{\mathbf{X}_n\}$  be a sequence of random variables, where  $\mathbf{X}, \mathbf{X}_n \in \mathbb{R}^d$ . If  $\mathbf{a}^T \mathbf{X}_n \xrightarrow{\mathcal{D}} \mathbf{a}^T \mathbf{X}$  for all  $\mathbf{a} \in \mathbb{R}^d$ , then  $\mathbf{X}_n \xrightarrow{\mathcal{D}} \mathbf{X}$ .*

### 1.10.2 Stochastic Order Notation

It is convenient to have notation that represent sequences of random variables that converge in probability to zero or sequences of random variables that are bounded in probability (van der Vaart, 1998). We make use of stochastic order notation such as  $o_P(1)$  and  $O_P(1)$  for such purposes. The notation  $o_P(1)$ , used for the notion of convergence in probability to zero, is formally defined below

**Definition 8.** *Let  $\{X_n\}$  be a sequence of random variables. It is convenient to write  $X_n = o_P(1)$  to represent that  $X_n$  converges in probability to zero, or  $X_n \xrightarrow{P} 0$ , if for every  $\varepsilon > 0$  we have,*

$$P(|X_n| > \varepsilon) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

Sequences that are bounded in probability can be represented using the  $O_P(1)$  notation, which is formally defined below

**Definition 9.** *Let  $\{X_n\}$  be a sequence of random variables. It is convenient to write  $X_n = O_P(1)$  to represent that  $X_n$  is bounded in probability, if for every  $\varepsilon > 0$ , there exists  $M_\varepsilon > 0$  such that,*

$$P(|X_n| > M_\varepsilon) < \varepsilon, \quad \text{for all } n.$$

Note that for a sequence of random variables  $\{X_n\}$ , if  $X_n = o_P(1)$ , then  $X_n = O_P(1)$  as well. Using Definitions 8 and 9, we also have the following general results, where for sequences of random variables  $\{X_n\}$ ,  $\{Y_n\}$  and  $\{R_n\}$ , we have,

$$X_n = o_P(R_n) \quad \text{if and only if} \quad X_n = Y_n R_n \quad \text{and} \quad Y_n = o_P(1).$$

and

$$X_n = O_P(R_n) \quad \text{if and only if} \quad X_n = Y_n R_n \quad \text{and} \quad Y_n = O_P(1).$$

Lastly, there are useful rules of calculus concerning  $o_P$  and  $O_P$  symbols. Some of

these rules are presented here:

$$\begin{aligned} o_P(1) + o_P(1) &= o_P(1), \\ o_P(1) + O_P(1) &= O_P(1), \\ o_P(1)O_P(1) &= o_P(1). \end{aligned}$$

Let  $\{a_n\}$  and  $\{b_n\}$  be sequences of positive real numbers. Then we also have the following rules:

$$\begin{aligned} o_P(a_n)o_P(b_n) &= o_P(a_nb_n), \\ o_P(a_n)O_P(b_n) &= o_P(a_nb_n), \\ o_P(a_n) + o_P(b_n) &= o_P(\max\{a_n, b_n\}). \end{aligned}$$

### 1.10.3 Other Tools for Working with Asymptotic Expansions

The *stochastic Taylor formula*, *inversion formula for an asymptotic series* and the *Laplace expansion for evaluating an integral* are useful tools for working with asymptotic approximations and expansions. We will highlight and present the inversion formula for an asymptotic series while the details regarding the other tools can be found in Pace and Salvan (1997).

#### 1.10.3.1 Inversion of Asymptotic Series

In this subsection, we present an approach for inverting a univariate asymptotic series. Details regarding the derivation of the inversion formula for both the univariate and multivariate cases can be found in Pace and Salvan (1997).

Let  $y = f(x)$ ,  $x \in \mathbb{R}$  be a real smooth function which admits the following power series expansion

$$y = x + a_1x^2 + a_2x^3 + \dots \quad (1.23)$$

Assume that the terms in (1.23) depend on an asymptotic parameter  $n$ . Specifically, let  $x = O(n^{-\alpha})$ ,  $\alpha > 0$  and let  $a_i = O(1)$ ,  $i = 1, 2, \dots$ . Suppose that we invert the function  $y = f(x)$  as  $x = g(y)$  and wished to express  $g(y)$ , in the neighbourhood of  $y = 0$ , as a power series expansion as follows

$$x = y + b_1x^2 + b_2x^3 + O(n^{-4\alpha}), \quad (1.24)$$

where  $b_1$  and  $b_2$  can be expressed in terms of constants  $a_1$  and  $a_2$ . Ignoring terms of

order  $O(n^{-4\alpha})$ , the final expressions for  $b_1$  and  $b_2$  are as follows

$$b_1 = -a_1, \quad b_2 = -(a_2 - 2a_1^2).$$

## 1.11 Frequentist Variational Approximations

*Variational approximations* have roots in variational calculus and serve as an approach for performing approximate inference on model parameters in complex statistical models. This class of methods is commonly used in Bayesian inference and in recent years, it has become a popular alternative to existing methods such as Markov chain Monte Carlo and Laplace approximation methods. It is shown that the same ideas can also be transferred to frequentist contexts (Ormerod and Wand, 2010). In this section, we will delve into how variational approximations can be used in frequentist contexts.

In frequentist inferential problems, variational approximation methods mainly benefit inference carried out on statistical models where the vector of observations  $\mathbf{y}$  is conditioned on a latent variable vector  $\mathbf{u}$ . In the context of generalized linear mixed models, the vector of latent variables essentially corresponds to the vector of random effects as shown in (1.17).

Let  $\boldsymbol{\theta}$  be a vector of model parameters. When conditioning over the vector of latent variables  $\mathbf{u}$  is present, the log-likelihood for a general statistical model parametrized by  $\boldsymbol{\theta}$  is as follows

$$\ell(\boldsymbol{\theta}) = \log p(\mathbf{y}; \boldsymbol{\theta}) = \log \int p(\mathbf{y}|\mathbf{u}; \boldsymbol{\theta})p(\mathbf{u}; \boldsymbol{\theta})d\mathbf{u}. \quad (1.25)$$

However, the integral in (1.25) may be intractable. Thus,  $\ell(\boldsymbol{\theta})$  may not have a closed form and maximum likelihood estimation is hindered. The variational approximation method works around the intractability issue to provide a variational approximation to the maximum likelihood estimation approach, explained further below.

Let us define  $q(\mathbf{u})$  to be an arbitrary density function in  $\mathbf{u}$ . Then the expression for

the log-likelihood satisfies the following mathematical steps

$$\begin{aligned}
\ell(\boldsymbol{\theta}) &= \ell(\boldsymbol{\theta}) \int q(\mathbf{u}) d\mathbf{u} \\
&= \int q(\mathbf{u}) \ell(\boldsymbol{\theta}) d\mathbf{u} \\
&= \int q(\mathbf{u}) \log \left\{ \frac{p(\mathbf{y}, \mathbf{u}; \boldsymbol{\theta}) / q(\mathbf{u})}{p(\mathbf{u} | \mathbf{y}; \boldsymbol{\theta}) / q(\mathbf{u})} \right\} d\mathbf{u} \\
&= \int q(\mathbf{u}) \log \left\{ \frac{p(\mathbf{y}, \mathbf{u}; \boldsymbol{\theta})}{q(\mathbf{u})} \right\} d\mathbf{u} + \int q(\mathbf{u}) \log \left\{ \frac{q(\mathbf{u})}{p(\mathbf{u} | \mathbf{y}; \boldsymbol{\theta})} \right\} d\mathbf{u} \\
&\geq \underline{\ell}(\boldsymbol{\theta}; q)
\end{aligned} \tag{1.26}$$

where

$$\underline{\ell}(\boldsymbol{\theta}; q) \equiv \int q(\mathbf{u}) \log \left\{ \frac{p(\mathbf{y}, \mathbf{u}; \boldsymbol{\theta})}{q(\mathbf{u})} \right\} d\mathbf{u}. \tag{1.27}$$

The inequality exists since

$$\int q(\mathbf{u}) \log \left\{ \frac{q(\mathbf{u})}{p(\mathbf{u} | \mathbf{y}; \boldsymbol{\theta})} \right\} d\mathbf{u} \geq 0 \tag{1.28}$$

across all densities  $q$ . Equality in (1.28) occurs if and only if

$$q(\mathbf{u}) = p(\mathbf{u} | \mathbf{y}; \boldsymbol{\theta})$$

almost everywhere. The integral expression in (1.28) is known as the *Kullback-Leibler divergence* between  $q$  and  $p(\cdot | \mathbf{y})$ .

One may now select a density  $q$ , where  $q(\mathbf{u})$  approximates  $p(\mathbf{u} | \mathbf{y}; \boldsymbol{\theta})$ , such that  $\underline{\ell}(\boldsymbol{\theta}; q)$  is more tractable than  $\ell(\boldsymbol{\theta})$ . One must also simultaneously aim to minimize the Kullback-Leibler divergence between  $q(\mathbf{u})$  and  $p(\mathbf{u} | \mathbf{y}; \boldsymbol{\theta})$ , which can be achieved by maximizing  $\underline{\ell}(\boldsymbol{\theta}; q)$ , as shown in (1.26).

Suppose that  $q$  is restricted to a family of parametric densities  $\{q(\mathbf{u}; \boldsymbol{\xi}) : \boldsymbol{\xi} \in \Xi\}$  parametrized by a vector of variational parameters  $\boldsymbol{\xi}$ . Then, the expression for  $\underline{\ell}(\boldsymbol{\theta}; q)$  in (1.27) becomes

$$\underline{\ell}(\boldsymbol{\theta}, \boldsymbol{\xi}; q) \equiv \int q(\mathbf{u}; \boldsymbol{\xi}) \log \left\{ \frac{p(\mathbf{y}, \mathbf{u}; \boldsymbol{\theta})}{q(\mathbf{u}; \boldsymbol{\xi})} \right\} d\mathbf{u}. \tag{1.29}$$

One now maximises over the vector of model parameters  $\boldsymbol{\theta}$  and the vector of variational parameters  $\boldsymbol{\xi}$ , in order to maximise the approximate log-likelihood,  $\underline{\ell}(\boldsymbol{\theta}, \boldsymbol{\xi}; q)$ , and to minimize the Kullback-Leibler divergence between  $q(\mathbf{u}; \boldsymbol{\xi})$  and  $p(\mathbf{u} | \mathbf{y}; \boldsymbol{\theta})$  respectively.

This leads to the following altered maximization problem where

$$(\hat{\underline{\theta}}, \hat{\underline{\xi}}) = \operatorname{argmax}_{\underline{\theta}, \underline{\xi}} \ell(\underline{\theta}, \underline{\xi}; q).$$

Then,  $\hat{\underline{\theta}}$  is the variational approximation to the maximum likelihood estimator  $\hat{\underline{\theta}}$  as defined in (1.21).

When  $q$  is chosen to be a Gaussian density function, then this particular class of variational approximation methods is known as Gaussian variational approximations. In GVA, the variational parameters are the mean and variance (or covariance matrix) parameters of the approximating normal distribution.

### 1.11.1 Thouless-Anderson-Palmer Variational Approach

The TAP variational framework, which builds on the GVA approach, has recently surfaced in statistical literature and potentially provides more accurate approximations as compared to the GVA approach.

The *Thouless-Anderson-Palmer* paradigm (Thouless et al., 1977) was first developed in statistical physics literature and gained traction as the authors provided TAP equations as an alternative approach to the solution for certain spin glass models. The work was further built on in Plefka (1982) where it is shown that the power expansion of the Gibbs potential of the infinite-ranged Ising spin glass model of Sherrington and Kirkpatrick (Sherrington and Kirkpatrick, 1975) up to the second order in the exchange couplings leads to the TAP equations.

In machine learning contexts, with the help of Plefka's expansion, it has been shown that better approximations arise by minimizing the TAP free energy, instead of the mean field free energy typically used in variational inference (Fan et al., 2021).

Recent theoretical work by Professor Song Mei from University of California, Berkeley, U.S.A and Professor Iain Johnstone from Stanford University, U.S.A, compares the estimates obtained from maximum likelihood estimation, GVA and TAP variational approaches for GLMMs. Let  $\underline{\theta}$  be a vector of model parameters. Then, under some conditions, recent work by Professors Song Mei and Iain Johnstone (private communication) shows that

$$\|\underline{\theta}_{\text{GVA}} - \underline{\theta}_{\text{MLE}}\|_2 \approx Cn^{-2} \quad \text{and} \quad \|\underline{\theta}_{\text{TAP}} - \underline{\theta}_{\text{MLE}}\|_2 \approx Cn^{-3},$$



---

where  $C$  denotes a constant independent of  $\boldsymbol{\theta}$ . Note that the approximation error between the TAP and MLE estimates is smaller than the approximation error between the GVA and MLE estimates.

Despite the potential improvement in the accuracy of approximations that the TAP variational method can provide, statistical applications such as longitudinal data analysis and multilevel models analysis, which may benefit from using TAP methodology, have not been investigated.

Therefore, towards the end of this thesis, we apply the Thouless-Anderson-Palmer methodology to GLMMs and evaluate the approach via simulation studies.

## Chapter 2

# Preliminary Lemmas and Their Proofs

Detailed asymptotic analysis is necessary to obtain the asymptotic normality theorems for the maximum likelihood estimators for Gaussian response linear mixed models and maximum quasi-likelihood estimators for generalized linear mixed models in Chapters 3, 4 and 7 respectively. In the process of doing so, we deal with working with population limits of predictor-dependent sample mean quantities and establishing matrix norm asymptotic negligibility between matrix square roots of inverse Fisher information matrices and their simpler asymptotic block diagonal forms. Currently, there are no results available to deal with both these tasks in a simple manner.

Therefore, in this chapter, we introduce three novel lemmas that will act as essential tools required to solve these two tasks.

The appendix contains the proofs for the lemmas introduced in this chapter.

### 2.1 Lemma 1

Certain population quantities appear in the asymptotic normality theorems in Chapters 3, 4 and 7 respectively. These population quantities correspond to the convergence in probability limit of two particular predictor-dependent sample mean quantities each. In this section, we isolate the problem of deriving the population leading term of the first predictor-dependent sample mean quantity in the form of Lemma 1.

**Lemma 1.** Let  $\mathbf{X} \equiv (\mathbf{X}_A^T, \mathbf{X}_B^T)^T$  and  $\mathbf{X}_{ij} \equiv (\mathbf{X}_{Aij}^T, \mathbf{X}_{Bij}^T)^T$ ,  $1 \leq i \leq m$ ,  $1 \leq j \leq n_i$  be independent and identically distributed  $(d_A + d_B) \times 1$  random vectors, with  $d_A \geq 1$  being the number of entries of  $\mathbf{X}_A$  and the  $\mathbf{X}_{Aij}$ s and  $d_B \geq 1$  being the number of entries of  $\mathbf{X}_B$  and the  $\mathbf{X}_{Bij}$ s, with  $d = d_A + d_B$ . Also, let  $\mathbf{U}$  and  $\mathbf{U}_1, \dots, \mathbf{U}_m$  be independent and identically distributed  $N(\mathbf{0}, \mathbf{I})$  ( $d_A \times 1$ ) random vectors, distributed independently of  $\mathbf{X}$  and the  $\mathbf{X}_{ij}$ s, where  $\Sigma$  is symmetric and positive definite. Let  $f$  be a Borel measurable, positive real-valued function on  $\mathbb{R}^{d_A+d_B}$  and assume that

$$E \{|X_k X_{k'}| f(\mathbf{X}, \mathbf{U})\} < \infty \text{ for all } 1 \leq k, k' \leq d$$

where  $X_k$  is the  $k^{\text{th}}$  row of  $\mathbf{X}$ . Then

$$\frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^{n_i} \mathbf{X}_{ij} \mathbf{X}_{ij}^T E \{f(\mathbf{X}_{ij}, \mathbf{U}) | \mathbf{X}_{ij}\} = E \{\mathbf{X} \mathbf{X}^T f(\mathbf{X}, \mathbf{U})\} + o_P(1) \mathbf{1}_d \mathbf{1}_d^T$$

where

$$n \equiv \frac{1}{m} \sum_{i=1}^m n_i.$$

## 2.2 Lemma 2

In this section, we isolate the problem of deriving the population leading term of the second predictor-dependent sample mean quantity in the form of Lemma 2. Using both Lemmas 1 and 2 lead to a full expression for the population leading term in the main Fisher information block, represented by  $\Sigma_{\beta_B}$ ,  $\Lambda_{\beta_B}$  and  $\Lambda_{\beta_1}^*$  in Chapters 3, 4 and 7 respectively.

**Lemma 2.** Let  $\mathbf{X} \equiv (\mathbf{X}_A^T, \mathbf{X}_B^T)^T$  and  $\mathbf{X}_{ij} \equiv (\mathbf{X}_{Aij}^T, \mathbf{X}_{Bij}^T)^T$ ,  $1 \leq i \leq m$ ,  $1 \leq j \leq n_i$  be independent and identically distributed  $(d_A + d_B) \times 1$  random vectors, with  $d_A \geq 1$  being the number of entries of  $\mathbf{X}_A$  and the  $\mathbf{X}_{Aij}$ s and  $d_B \geq 1$  being the number of entries of  $\mathbf{X}_B$  and the  $\mathbf{X}_{Bij}$ s, with  $d = d_A + d_B$ . Also, let  $\mathbf{U}$  and  $\mathbf{U}_1, \dots, \mathbf{U}_m$  be independent and identically distributed random vectors, distributed independently of  $\mathbf{X}$  and the  $\mathbf{X}_{ij}$ s. Let  $f$  be a Borel measurable, positive real-valued function on  $\mathbb{R}^{d_A+d_B}$  and assume that

$$E \left[ \frac{E [\max\{1, \|\mathbf{X}\|\}^8 \max\{1, f(\mathbf{X}, \mathbf{U})\}^4 | \mathbf{U}]}{\min\{1, \lambda_{\min}(E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\})\}^2} \right] < \infty.$$

If  $m$  and  $n_i$  satisfy assumptions that the number of groups  $m$  diverges to  $\infty$  and that the within-group sample sizes  $n_i$  diverge to  $\infty$  in such a way that  $n_i/n \rightarrow C_i$  for constants

$0 < C_i < \infty, 1 \leq i \leq m$ , then

$$\begin{aligned} & E \left[ \frac{1}{mn} \sum_{i=1}^m \left\{ \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T f(\mathbf{X}_{ij}, U_i) \right\} \left\{ \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T f(\mathbf{X}_{ij}, U_i) \right\}^{-1} \right. \\ & \quad \left. \times \left\{ \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T f(\mathbf{X}_{ij}, U_i) \right\}^T \middle| \mathbf{X}_{11}, \dots, \mathbf{X}_{mnm} \right] \\ & \xrightarrow{P} E \left[ E(\mathbf{X}_B \mathbf{X}_A^T f(\mathbf{X}, U) | U) \{ E(\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, U) | U) \}^{-1} E(\mathbf{X}_B \mathbf{X}_A^T f(\mathbf{X}, U) | U)^T \right]. \end{aligned}$$

### 2.3 Lemma 3

The asymptotic normality theorems in Chapters 3, 4 and 7 involve replacement of the matrix square root of the inverse Fisher information matrix by the matrix square root of the asymptotic expression for the inverse Fisher information matrix. This is due to the remainder terms (calculated as a difference between the inverse Fisher information matrix and its asymptotic counterpart) having an asymptotically negligible effect on the relevant matrix square roots. Lemma 3 provides a formalization of this state of affairs, which is used in the final steps of the derivation in the asymptotic normality theorems in Chapters 3, 4 and 7.

**Lemma 3.** *Define the sequences of matrices*

$$\mathbf{M}_n \equiv \begin{bmatrix} \mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2} & R_n \mathbf{1}_p \mathbf{1}_q^T \\ R_n \mathbf{1}_q \mathbf{1}_p^T & \frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2} \end{bmatrix} \quad \text{and} \quad \mathbf{M}_{n,\infty} \equiv \begin{bmatrix} \mathbf{K} & \mathbf{0} \\ \mathbf{0} & \frac{1}{n} \mathbf{L} \end{bmatrix}$$

where  $\mathbf{K}$  and  $\mathbf{L}$  are  $p \times p$  and  $q \times q$  symmetric positive definite matrices and  $Q_n, R_n$  and  $T_n$  are sequences of random variables satisfying  $Q_n = o_P(1)$ ,  $R_n = O_P(n^{-1})$  and  $T_n = o_P(n^{-1})$ . Also note that  $\nu^{\otimes 2} \equiv \nu \nu^T$ . Then, as  $n \rightarrow \infty$ ,

$$\left\| \mathbf{M}_{n,\infty}^{-1/2} \mathbf{M}_n^{1/2} - \mathbf{I} \right\|_F \xrightarrow{P} 0.$$

## 2.4 Appendix

### 2.4.1 Proof of Lemma 1

Let

$$\bar{\mathbf{G}} \equiv \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{ij} \mathbf{X}_{ij}^T E\{f(\mathbf{X}_{ij}, \mathbf{U}) | \mathbf{X}_{ij}\}.$$

Firstly, note that a more explicit form of  $\bar{\mathbf{G}}$  is:

$$\bar{\mathbf{G}} = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{ij} \mathbf{X}_{ij}^T (2\pi)^{-d_A/2} |\boldsymbol{\Sigma}|^{-1/2} \int_{\mathbb{R}^{d_A}} f(\mathbf{X}_{ij}, \mathbf{u}) \exp\left(-\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) d\mathbf{u}.$$

Next, let

$$N \equiv mn = \sum_{i=1}^m n_i,$$

and  $\mathbf{G}_l$ ,  $1 \leq l \leq N$ , be the following ‘‘single subscript’’ re-labelling of the terms inside the double summation of  $\bar{\mathbf{G}}$ . Then we have the following terms where

$$\begin{aligned} \mathbf{G}_1 &= \mathbf{X}_{11} \mathbf{X}_{11}^T (2\pi)^{-d_A/2} |\boldsymbol{\Sigma}|^{-1/2} \int_{\mathbb{R}^{d_A}} f(\mathbf{X}_{ij}, \mathbf{u}) \exp\left(-\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) d\mathbf{u}, \\ \mathbf{G}_2 &= \mathbf{X}_{12} \mathbf{X}_{12}^T (2\pi)^{-d_A/2} |\boldsymbol{\Sigma}|^{-1/2} \int_{\mathbb{R}^{d_A}} f(\mathbf{X}_{ij}, \mathbf{u}) \exp\left(-\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) d\mathbf{u}, \\ &\vdots \\ \mathbf{G}_N &= \mathbf{X}_{mn_m} \mathbf{X}_{mn_m}^T (2\pi)^{-d_A/2} |\boldsymbol{\Sigma}|^{-1/2} \int_{\mathbb{R}^{d_A}} f(\mathbf{X}_{ij}, \mathbf{u}) \exp\left(-\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) d\mathbf{u}. \end{aligned}$$

Then,

$$\bar{\mathbf{G}} = \frac{1}{N} \sum_{l=1}^N \mathbf{G}_l$$

is a sample mean of  $N$  independent and identically distributed  $d \times d$  random matrices with the common distribution

$$\mathbf{G} \equiv \mathbf{X} \mathbf{X}^T E\{f(\mathbf{X}, \mathbf{U}) | \mathbf{X}\}.$$

Using the law of total expectation, the mean of  $\mathbf{G}$  is

$$\begin{aligned} E(\mathbf{G}) &= E[\mathbf{X} \mathbf{X}^T E\{f(\mathbf{X}, \mathbf{U}) | \mathbf{X}\}] \\ &= E[E\{\mathbf{X} \mathbf{X}^T f(\mathbf{X}, \mathbf{U}) | \mathbf{X}\}] \\ &= E\{\mathbf{X} \mathbf{X}^T f(\mathbf{X}, \mathbf{U})\}. \end{aligned}$$

Lastly, we need to impose the following first order moment conditions on the entries of  $\mathbf{G}$  where

$$E\{\mathbf{X}_k \mathbf{X}_{k'} | f(\mathbf{X}, \mathbf{U})\} < \infty \text{ for all } 1 \leq k, k' \leq d, \quad (2.1)$$

where  $\mathbf{X}_k$  is the  $k$ th entry of  $\mathbf{X}$ . Therefore, under moment conditions involving the entries of  $\mathbf{G}$  as stated in (2.1), we have,

$$\bar{\mathbf{G}} = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{ij} \mathbf{X}_{ij}^T E\{f(\mathbf{X}_{ij}, \mathbf{U}) | \mathbf{X}_{ij}\} = E\{\mathbf{X} \mathbf{X}^T f(\mathbf{X}, \mathbf{U})\} + o_P(1) \mathbf{1}_d \mathbf{1}_d^T.$$

## 2.4.2 Proof of Lemma 2

### 2.4.2.1 A Fundamental Inequality for the Spectral Norm of a Vectorised Matrix

Let  $\mathbf{A}$  be a  $d_1 \times d_2$  matrix. This subsection looks into the relationship between

$$\|\mathbf{A}\|_s \quad \text{and} \quad \|\text{vec}(\mathbf{A})\|_s.$$

We start with the following inequality:

$$\|\mathbf{A}\|_s \leq \|\mathbf{A}\|_F \leq \sqrt{\text{rank}(\mathbf{A})} \|\mathbf{A}\|_s. \quad (2.2)$$

However, since  $\text{rank}(\mathbf{A}) \leq \max(d_1, d_2)$ , we then obtain the following

$$\|\mathbf{A}\|_s \leq \|\mathbf{A}\|_F \leq \max(\sqrt{d_1}, \sqrt{d_2}) \|\mathbf{A}\|_s. \quad (2.3)$$

Next note that

$$\|\mathbf{A}\|_F = \sqrt{\text{tr}(\mathbf{A}^T \mathbf{A})} = \sqrt{\text{vec}(\mathbf{A})^T \text{vec}(\mathbf{A})} = \|\text{vec}(\mathbf{A})\|_F. \quad (2.4)$$

Replacement of  $\mathbf{A}$  by  $\text{vec}(\mathbf{A})$  in the first inequality of (2.2) gives

$$\|\text{vec}(\mathbf{A})\|_s \leq \|\text{vec}(\mathbf{A})\|_F.$$

The equality in (2.4) then gives

$$\|\text{vec}(\mathbf{A})\|_s \leq \|\mathbf{A}\|_F.$$

Application of the second inequality of (2.3) leads to

$$\|\text{vec}(\mathbf{A})\|_s \leq \max(\sqrt{d_1}, \sqrt{d_2}) \|\mathbf{A}\|_s \quad (2.5)$$

for all  $d_1 \times d_2$  matrices  $\mathbf{A}$ .

### 2.4.2.2 Notational Definitions

Let us define the following sample and population type quantities as follows:

$$\begin{aligned} \hat{\mathcal{N}}_i(\mathbf{U}) &\equiv \frac{1}{n_i} \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T f(\mathbf{X}_{ij}, \mathbf{U}), & \hat{\mathcal{D}}_i(\mathbf{U}) &\equiv \frac{1}{n_i} \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T f(\mathbf{X}_{ij}, \mathbf{U}), \\ \mathcal{N}(\mathbf{U}) &\equiv E\{\mathbf{X}_B \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} & \text{and} & \quad \mathcal{D}(\mathbf{U}) \equiv E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}. \end{aligned}$$

Next, for  $t \in [0, 1]$ , let

$$\mathcal{N}_{it}^\dagger(\mathbf{U}) \equiv (1-t)\mathcal{N}(\mathbf{U}) + t\widehat{\mathcal{N}}_i(\mathbf{U}) \quad \text{and} \quad \mathcal{D}_{it}^\dagger(\mathbf{U}) \equiv (1-t)\mathcal{D}(\mathbf{U}) + t\widehat{\mathcal{D}}_i(\mathbf{U}).$$

Throughout this proof, we also let

$$\mathcal{X}_i \equiv \{\mathbf{X}_{i1}, \dots, \mathbf{X}_{in_i}\}.$$

In addition, if  $\mathbf{S}$  is a  $d_B \times d_A$  matrix and  $\mathbf{T}$  is a  $d_A \times d_A$  symmetric matrix, define

$$\mathcal{R} \left( \begin{bmatrix} \mathbf{S} \\ \mathbf{T} \end{bmatrix} \right) \equiv \text{vec}(\mathbf{S}\mathbf{T}^{-1}\mathbf{S}^T)^T. \quad (2.6)$$

In the next subsection, we wish to find an explicit expression for

$$\nabla_{\text{vec}(\mathbf{S}), \text{vec}(\mathbf{T})} \mathcal{R} \left( \begin{bmatrix} \mathbf{S} \\ \mathbf{T} \end{bmatrix} \right). \quad (2.7)$$

### 2.4.2.3 Derivation of (2.7)

#### Differentiation with Respect to $\mathbf{S}$

Throughout this subsection it is assumed that the differential operator  $d$  is with respect to  $\mathbf{S}$ . Now, noting that  $\mathbf{T}^T = \mathbf{T}$ ,

$$\begin{aligned} d \text{vec}(\mathbf{S}\mathbf{T}^{-1}\mathbf{S}^T) &= \text{vec}(d(\mathbf{S}\mathbf{T}^{-1}\mathbf{S}^T)) \\ &= \text{vec}(d(\mathbf{S}\mathbf{T}^{-1})\mathbf{S}^T + (\mathbf{S}\mathbf{T}^{-1})d\mathbf{S}^T) \\ &= \text{vec}((d\mathbf{S})\mathbf{T}^{-1}\mathbf{S}^T) + \text{vec}(\mathbf{S}\mathbf{T}^{-1}(d\mathbf{S})^T) \\ &= \text{vec}((d\mathbf{S})\mathbf{T}^{-1}\mathbf{S}^T) + \mathbf{K}_{d_B} \text{vec}((d\mathbf{S})\mathbf{T}^{-1}\mathbf{S}^T) \\ &= (\mathbf{I}_{d_B^2} + \mathbf{K}_{d_B}) \text{vec}((d\mathbf{S})\mathbf{T}^{-1}\mathbf{S}^T) \\ &= (\mathbf{I}_{d_B^2} + \mathbf{K}_{d_B}) \text{vec}(\mathbf{I}_{d_B}(d\mathbf{S})\mathbf{T}^{-1}\mathbf{S}^T) \\ &= (\mathbf{I}_{d_B^2} + \mathbf{K}_{d_B}) \{(\mathbf{T}^{-1}\mathbf{S}^T)^T \otimes \mathbf{I}_{d_B}\} \text{vec}(d\mathbf{S}) \\ &= (\mathbf{I}_{d_B^2} + \mathbf{K}_{d_B}) \{(\mathbf{S}\mathbf{T}^{-1}) \otimes \mathbf{I}_{d_B}\} d\text{vec}(\mathbf{S}). \end{aligned}$$

Hence,

$$\mathbf{D}_{\text{vec}(\mathbf{S})} \text{vec}(\mathbf{S}\mathbf{T}^{-1}\mathbf{S}^T) = (\mathbf{I}_{d_B^2} + \mathbf{K}_{d_B}) \{(\mathbf{S}\mathbf{T}^{-1}) \otimes \mathbf{I}_{d_B}\}.$$

Since  $\mathbf{K}_{d_B}^T = \mathbf{K}_{d_B}$  as in (1.8), we have,

$$\nabla_{\text{vec}(\mathbf{S})} \text{vec}(\mathbf{S}\mathbf{T}^{-1}\mathbf{S}^T)^T = \{(\mathbf{T}^{-1}\mathbf{S}^T) \otimes \mathbf{I}_{d_B}\} (\mathbf{I}_{d_B^2} + \mathbf{K}_{d_B}).$$

Differentiation with Respect to  $\mathbf{T}$ 

Throughout this subsubsection it is assumed that the differential operator  $d$  is with respect to  $\mathbf{T}$ . Now, noting that  $\mathbf{T}^T = \mathbf{T}$ ,

$$\begin{aligned} d \operatorname{vec}(\mathbf{S}\mathbf{T}^{-1}\mathbf{S}^T) &= \operatorname{vec}(\mathbf{S}(d\mathbf{T}^{-1})\mathbf{S}^T) \\ &= -\operatorname{vec}(\mathbf{S}\mathbf{T}^{-1}(d\mathbf{T})\mathbf{T}^{-1}\mathbf{S}^T) \\ &= -\{(\mathbf{T}^{-1}\mathbf{S}^T)^T \otimes (\mathbf{S}\mathbf{T}^{-1})\} \operatorname{vec}(d\mathbf{T}) \\ &= -\{(\mathbf{S}\mathbf{T}^{-T}) \otimes (\mathbf{S}\mathbf{T}^{-1})\} d\operatorname{vec}(\mathbf{T}). \end{aligned}$$

Hence,

$$\mathbf{D}_{\operatorname{vec}(\mathbf{T})} \operatorname{vec}(\mathbf{S}\mathbf{T}^{-1}\mathbf{S}^T) = -(\mathbf{S}\mathbf{T}^{-T}) \otimes (\mathbf{S}\mathbf{T}^{-1}).$$

Therefore, we have,

$$\nabla_{\operatorname{vec}(\mathbf{T})} \operatorname{vec}(\mathbf{S}\mathbf{T}^{-1}\mathbf{S}^T)^T = -(\mathbf{T}^{-1}\mathbf{S}^T) \otimes (\mathbf{T}^{-1}\mathbf{S}^T).$$

Combination of the Two Gradient Vectors

On combining the results of the previous two subsubsubsections, we get an explicit expression for (2.7) as follows

$$\nabla_{\operatorname{vec}(\mathbf{S}), \operatorname{vec}(\mathbf{T})} \mathcal{R} \left( \begin{bmatrix} \mathbf{S} \\ \mathbf{T} \end{bmatrix} \right) = \begin{bmatrix} \{(\mathbf{T}^{-1}\mathbf{S}^T) \otimes \mathbf{I}_{d_B}\}(\mathbf{I}_{d_B^2} + \mathbf{K}_{d_B}) \\ -(\mathbf{T}^{-1}\mathbf{S}^T) \otimes (\mathbf{T}^{-1}\mathbf{S}^T) \end{bmatrix}. \quad (2.8)$$

**2.4.2.4 Expression for (2.6) with Lagrange Form of Remainder**

Using (2.8), a Taylor series expansion of  $\mathcal{R}$  with the Lagrange form of the remainder is

$$\begin{aligned} \mathcal{R} \left( \begin{bmatrix} \mathbf{S} \\ \mathbf{T} \end{bmatrix} \right) &= \mathcal{R} \left( \begin{bmatrix} \mathbf{S}_0 + \mathbf{S} - \mathbf{S}_0 \\ \mathbf{T}_0 + \mathbf{T} - \mathbf{T}_0 \end{bmatrix} \right) \\ &= \mathcal{R} \left( \begin{bmatrix} \mathbf{S}_0 \\ \mathbf{T}_0 \end{bmatrix} \right) + \begin{bmatrix} \operatorname{vec}(\mathbf{S} - \mathbf{S}_0) \\ \operatorname{vec}(\mathbf{T} - \mathbf{T}_0) \end{bmatrix}^T \nabla_{\operatorname{vec}(\mathbf{S}), \operatorname{vec}(\mathbf{T})} \mathcal{R} \left( \begin{bmatrix} \mathbf{S}_t^\dagger \\ \mathbf{T}_t^\dagger \end{bmatrix} \right) \\ &= \mathcal{R} \left( \begin{bmatrix} \mathbf{S}_0 \\ \mathbf{T}_0 \end{bmatrix} \right) + \begin{bmatrix} \operatorname{vec}(\mathbf{S} - \mathbf{S}_0) \\ \operatorname{vec}(\mathbf{T} - \mathbf{T}_0) \end{bmatrix}^T \begin{bmatrix} \{(\mathbf{T}_t^\dagger)^{-1}(\mathbf{S}_t^\dagger)^T\} \otimes \mathbf{I}_{d_B}(\mathbf{I}_{d_B^2} + \mathbf{K}_{d_B}) \\ -\{(\mathbf{T}_t^\dagger)^{-1}(\mathbf{S}_t^\dagger)^T\} \otimes \{(\mathbf{T}_t^\dagger)^{-1}(\mathbf{S}_t^\dagger)^T\} \end{bmatrix} \end{aligned}$$

where

$$\begin{bmatrix} \mathbf{S}_t^\dagger \\ \mathbf{T}_t^\dagger \end{bmatrix} \equiv (1-t) \begin{bmatrix} \mathbf{S}_0 \\ \mathbf{T}_0 \end{bmatrix} + t \begin{bmatrix} \mathbf{S} \\ \mathbf{T} \end{bmatrix} = \begin{bmatrix} (1-t)\mathbf{S}_0 + t\mathbf{S} \\ (1-t)\mathbf{T}_0 + t\mathbf{T} \end{bmatrix}$$



and  $t \in [0, 1]$ . It follows that

$$\left\{ \mathcal{R} \left( \begin{bmatrix} \mathbf{S} \\ \mathbf{T} \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathbf{S}_0 \\ \mathbf{T}_0 \end{bmatrix} \right) \right\}^T = (\mathbf{I}_{d_B^2} + \mathbf{K}_{d_B}) [\{\mathbf{S}_t^\dagger(\mathbf{T}_t^\dagger)^{-1}\} \otimes \mathbf{I}_{d_B}] \text{vec}(\mathbf{S} - \mathbf{S}_0) \\ - [\{\mathbf{S}_t^\dagger(\mathbf{T}_t^\dagger)^{-1}\} \otimes \{\mathbf{S}_t^\dagger(\mathbf{T}_t^\dagger)^{-1}\}] \text{vec}(\mathbf{T} - \mathbf{T}_0). \quad (2.9)$$

#### 2.4.2.5 Spectral Norm Bounding of (2.9)

It follows from (2.9) that

$$\left\| \mathcal{R} \left( \begin{bmatrix} \mathbf{S} \\ \mathbf{T} \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathbf{S}_0 \\ \mathbf{T}_0 \end{bmatrix} \right) \right\|_S \leq \|\mathbf{I}_{d_B^2} + \mathbf{K}_{d_B}\|_S \|\{\mathbf{S}_t^\dagger(\mathbf{T}_t^\dagger)^{-1}\} \otimes \mathbf{I}_{d_B}\|_S \|\mathbf{S} - \mathbf{S}_0\|_F \\ + \|\{\mathbf{S}_t^\dagger(\mathbf{T}_t^\dagger)^{-1}\} \otimes \{\mathbf{S}_t^\dagger(\mathbf{T}_t^\dagger)^{-1}\}\|_S \|\mathbf{T} - \mathbf{T}_0\|_F \\ \leq (\|\mathbf{I}_{d_B^2}\|_S + \|\mathbf{K}_{d_B}\|_S) \|\mathbf{S}_t^\dagger(\mathbf{T}_t^\dagger)^{-1}\|_S \|\mathbf{I}_{d_B}\|_S \|\mathbf{S} - \mathbf{S}_0\|_F \\ + \{\|\mathbf{S}_t^\dagger(\mathbf{T}_t^\dagger)^{-1}\|_S\}^2 \|\mathbf{T} - \mathbf{T}_0\|_F \\ \leq 2\|\mathbf{S}_t^\dagger\|_S \|\mathbf{T}_t^\dagger\|_S^{-1} \|\mathbf{S} - \mathbf{S}_0\|_F + \{\|\mathbf{S}_t^\dagger\|_S \|\mathbf{T}_t^\dagger\|_S^{-1}\}^2 \|\mathbf{T} - \mathbf{T}_0\|_F.$$

Now, in terms of the notation given in the previous subsections, our goal is to show that

$$E \left\{ \frac{1}{m} \sum_{i=1}^m \frac{n_i}{n} \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) \middle| \mathcal{X}_i \right\} \xrightarrow{P} E \left\{ \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}) \\ \mathcal{D}(\mathbf{U}) \end{bmatrix} \right) \right\}. \quad (2.10)$$

in order to prove Lemma 2.

#### 2.4.2.6 Strategy for Proving (2.10)

Result (2.10) is a consequence of

$$E \left\| E \left\{ \frac{1}{m} \sum_{i=1}^m \frac{n_i}{n} \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) \middle| \mathcal{X}_i \right\} - E \left\{ \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}) \\ \mathcal{D}(\mathbf{U}) \end{bmatrix} \right) \right\} \right\|_S \rightarrow 0 \quad (2.11)$$

as  $m, n \rightarrow \infty$ . If we bring the second inner expectation inside the  $i = 1, \dots, m$  summation and replace the  $\mathbf{U}$  of this term by  $\mathbf{U}_i$ , then we can replace (2.11) by

$$E \left\| \frac{1}{m} \sum_{i=1}^m \frac{n_i}{n} E \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \middle| \mathcal{X}_i \right\} \right\|_S. \quad (2.12)$$

By noting that the left-hand side of (2.12) is bounded above by

$$\frac{1}{mn} \sum_{i=1}^m n_i E \left\| E \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \middle| \mathcal{X}_i \right\} \right\|_S,$$

it is sufficient to prove that

$$\frac{1}{mn} \sum_{i=1}^m n_i E \left\| E \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \middle| \mathcal{X}_i \right\} \right\|_S \rightarrow 0 \quad (2.13)$$

as  $m, n \rightarrow \infty$ . Next, for each  $1 \leq i \leq m$ , define the event

$$\mathcal{A}_i \equiv \left\{ \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_S \leq 1, \lambda_{\min}(\widehat{\mathcal{D}}_i(\mathbf{U}_i)) \geq \frac{1}{2} \lambda_{\min}(\mathcal{D}(\mathbf{U}_i)) \right\}. \quad (2.14)$$

Now, note that

$$\begin{aligned} & \frac{1}{mn} \sum_{i=1}^m n_i E \left\| E \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \middle| \mathcal{X}_i \right\} \right\|_S \\ & \leq \frac{1}{mn} \sum_{i=1}^m n_i E \left\| E \left[ \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\} I(\mathcal{A}_i) \middle| \mathcal{X}_i \right] \right\|_S \\ & \quad + \frac{1}{mn} \sum_{i=1}^m n_i E \left\| E \left[ \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\} I(\mathcal{A}_i^C) \middle| \mathcal{X}_i \right] \right\|_S. \end{aligned} \quad (2.15)$$

Then, to prove Lemma 2, it is sufficient to prove that

$$\frac{1}{mn} \sum_{i=1}^m n_i E \left\| E \left[ \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\} I(\mathcal{A}_i) \middle| \mathcal{X}_i \right] \right\|_S \rightarrow 0 \quad (2.16)$$

as  $m, n, \rightarrow \infty$  and

$$\frac{1}{mn} \sum_{i=1}^m n_i E \left\| E \left[ \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\} I(\mathcal{A}_i^C) \middle| \mathcal{X}_i \right] \right\|_S \rightarrow 0 \quad (2.17)$$

as  $m, n \rightarrow \infty$ .

#### 2.4.2.7 Proof of Result (2.16)

Throughout this subsection we are considering

$$(\mathbf{U}_i, \mathcal{X}_i) \text{ such that } \mathcal{A}_i \text{ occurs, } 1 \leq i \leq m. \quad (2.18)$$

Since

$$\mathcal{N}_{it}^\dagger(\mathbf{U}_i) \equiv (1-t)\mathcal{N}(\mathbf{U}_i) + t\widehat{\mathcal{N}}_i(\mathbf{U}_i) = \mathcal{N}(\mathbf{U}_i) + t\{\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\}$$

we have

$$\|\mathcal{N}_{it}^\dagger(\mathbf{U}_i)\|_S \leq \|\mathcal{N}(\mathbf{U}_i)\|_S + t\|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_S,$$

and it follows that, for  $t \in [0, 1]$  and under (2.18), we have,

$$\|\mathcal{N}_{it}^\dagger(\mathbf{U}_i)\|_S \leq \|\mathcal{N}(\mathbf{U}_i)\|_S + 1. \quad (2.19)$$

Next, note that

$$\begin{aligned}
\|\mathcal{D}_{it}^\dagger(\mathbf{U})^{-1}\|_S &= 1/\lambda_{\min}(\mathcal{D}_{it}^\dagger(\mathbf{U})) \\
&= 1/\lambda_{\min}((1-t)\mathcal{D}(\mathbf{U}_i) + t\widehat{\mathcal{D}}_i(\mathbf{U}_i)) \\
&\leq 1/\left\{(1-t)\lambda_{\min}(\mathcal{D}(\mathbf{U}_i)) + t\lambda_{\min}(\widehat{\mathcal{D}}_i(\mathbf{U}_i))\right\}.
\end{aligned} \tag{2.20}$$

Under (2.18), we have,

$$\begin{aligned}
\|\mathcal{D}_{it}^\dagger(\mathbf{U}_i)^{-1}\|_S &\leq 1/\left\{(1-t)\lambda_{\min}(\mathcal{D}(\mathbf{U}_i)) + t\lambda_{\min}(\widehat{\mathcal{D}}_i(\mathbf{U}_i))\right\} \\
&\leq 1/\left\{(1-t)\lambda_{\min}(\mathcal{D}(\mathbf{U}_i)) + \frac{t}{2}\lambda_{\min}(\mathcal{D}(\mathbf{U}_i))\right\} \\
&= \frac{1}{(1-t/2)\lambda_{\min}(\mathcal{D}(\mathbf{U}_i))} \\
&= \left(\frac{2}{2-t}\right)\|\mathcal{D}(\mathbf{U}_i)^{-1}\|_S.
\end{aligned}$$

Since

$$\sup_{t \in [0,1]} \left(\frac{2}{2-t}\right) = 2$$

and  $t \in [0, 1]$ , under (2.18), we have,

$$\|\mathcal{D}_{it}^\dagger(\mathbf{U}_i)^{-1}\|_S \leq 2\|\mathcal{D}(\mathbf{U}_i)^{-1}\|_S. \tag{2.21}$$

Substituting (2.19) and (2.21) into the above discrepancy in (2.16) involving the  $\mathcal{R}$ ,  $\widehat{\mathcal{D}}_i$ ,  $\widehat{\mathcal{N}}_i$ ,  $\mathcal{D}$  and  $\mathcal{N}$  functions we have,

$$\begin{aligned}
&\left\| \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\} I(\mathcal{A}_i) \right\|_S \\
&\leq 4(\|\mathcal{N}(\mathbf{U}_i)\|_S + 1) \|\mathcal{D}(\mathbf{U}_i)^{-1}\|_S \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F \\
&\quad + \left\{ 2(\|\mathcal{N}(\mathbf{U}_i)\|_S + 1) \|\mathcal{D}(\mathbf{U}_i)^{-1}\|_S \right\}^2 \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F \\
&= 4 \left\{ \mathcal{W}(\mathbf{U}_i) \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F + \mathcal{W}(\mathbf{U}_i)^2 \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F \right\}
\end{aligned}$$

where

$$\mathcal{W}(\mathbf{U}_i) \equiv (\|\mathcal{N}(\mathbf{U}_i)\|_S + 1) \|\mathcal{D}(\mathbf{U}_i)^{-1}\|_S = \frac{\|\mathcal{N}(\mathbf{U}_i)\|_S + 1}{\lambda_{\min}(\mathcal{D}(\mathbf{U}_i))}.$$

The left-hand side of (2.16) can now be re-written as

$$\begin{aligned}
& \frac{1}{mn} \sum_{i=1}^m n_i E \left\| E \left[ \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\} I_{(\mathcal{A}_i)} \middle| \mathcal{X}_i \right] \right\|_s \\
& \leq \frac{1}{mn} \sum_{i=1}^m n_i E \left[ E \left\| \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\} I_{(\mathcal{A}_i)} \middle| \mathcal{X}_i \right\|_s \right] \\
& \leq \frac{4}{mn} \sum_{i=1}^m n_i E \left( E \left[ \left\{ \mathcal{W}(\mathbf{U}_i) \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F + \mathcal{W}(\mathbf{U}_i)^2 \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F \right\} \middle| \mathcal{X}_i \right] \right) \\
& \leq \frac{4}{mn} \sum_{i=1}^m n_i E \left( E \left[ \mathcal{W}(\mathbf{U}_i) \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F \middle| \mathcal{X}_i \right] \right) \\
& \quad + \frac{4}{mn} \sum_{i=1}^m n_i E \left( E \left[ \mathcal{W}(\mathbf{U}_i)^2 \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F \middle| \mathcal{X}_i \right] \right) \\
& = \frac{4}{mn} \sum_{i=1}^m n_i E \left\{ \mathcal{W}(\mathbf{U}_i) \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F \right\} \\
& \quad + \frac{4}{mn} \sum_{i=1}^m n_i E \left\{ \mathcal{W}(\mathbf{U}_i)^2 \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F \right\}.
\end{aligned} \tag{2.22}$$

For the first term in the final expression of (2.22), note that,

$$\begin{aligned}
E \left\{ \mathcal{W}(\mathbf{U}_i) \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F \right\} &= E \left[ E \left\{ \mathcal{W}(\mathbf{U}_i) \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F \middle| \mathbf{U}_i \right\} \right] \\
&= E \left[ \mathcal{W}(\mathbf{U}_i) E \left\{ \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F \middle| \mathbf{U}_i \right\} \right].
\end{aligned} \tag{2.23}$$

From a conditional version of the Cauchy-Schwarz inequality,

$$\begin{aligned}
E \left\{ \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F \middle| \mathbf{U}_i \right\} &= \left( \left[ E \left\{ \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F \middle| \mathbf{U}_i \right\} \right]^2 \right)^{1/2} \\
&\leq \left[ E \left\{ \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F^2 \middle| \mathbf{U}_i \right\} \right]^{1/2}.
\end{aligned} \tag{2.24}$$

Using (2.24), observe that

$$\|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F^2 = \sum_{k=1}^{d_B} \sum_{k'=1}^{d_A} [\{\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\}_{kk'}]^2.$$

Then note that

$$\left[ \{\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\}_{kk'} \right]^2 = \left( \frac{1}{n_i} \sum_{j=1}^{n_i} [X_{Bijk} X_{Aijk'} f(\mathbf{X}_{ij}, \mathbf{U}_i) - E\{X_{Bk} X_{Ak'} f(\mathbf{X}, \mathbf{U}_i) | \mathbf{U}_i\}] \right)^2$$

from which it follows that

$$\begin{aligned}
& E\left(\left[\{\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\}_{kk'}\right]^2 \middle| \mathbf{U}_i\right) \\
&= E\left[\left(\frac{1}{n_i} \sum_{j=1}^{n_i} [X_{Bijk} X_{Aijk'} f(\mathbf{X}_{ij}, \mathbf{U}_i) - E\{X_{Bk} X_{Ak'} f(\mathbf{X}, \mathbf{U}_i) | \mathbf{U}_i\}]\right)^2 \middle| \mathbf{U}_i\right] \\
&= \text{Var}\left[\left(\frac{1}{n_i} \sum_{j=1}^{n_i} [X_{Bijk} X_{Aijk'} f(\mathbf{X}_{ij}, \mathbf{U}_i) - E\{X_{Bk} X_{Ak'} f(\mathbf{X}, \mathbf{U}_i) | \mathbf{U}_i\}]\right) \middle| \mathbf{U}_i\right] \\
&= \text{Var}\left[\left(\frac{1}{n_i} \sum_{j=1}^{n_i} X_{Bijk} X_{Aijk'} f(\mathbf{X}_{ij}, \mathbf{U}_i)\right) \middle| \mathbf{U}_i\right] \\
&= \frac{1}{n_i^2} \sum_{j=1}^{n_i} \text{Var}\{X_{Bijk} X_{Aijk'} f(\mathbf{X}_{ij}, \mathbf{U}_i) | \mathbf{U}_i\} \\
&= \frac{1}{n_i} \text{Var}\{X_{Bk} X_{Ak'} f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \\
&= \frac{1}{n_i} \left(E\{X_{Bk}^2 X_{Ak'}^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} - [E\{X_{Bk} X_{Ak'} f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}]^2\right) \\
&\leq \frac{1}{n_i} E\{X_{Bk}^2 X_{Ak'}^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}.
\end{aligned}$$

This implies that

$$\left[E\left\{\|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F^2 \middle| \mathbf{U}_i\right\}\right]^{1/2} \leq \left[\frac{1}{n_i} E\left\{\sum_{k=1}^{d_B} \sum_{k'=1}^{d_A} X_{Bk}^2 X_{Ak'}^2 f(\mathbf{X}, \mathbf{U}) \middle| \mathbf{U}\right\}\right]^{1/2}. \quad (2.25)$$

Substituting (2.25) into (2.24) and substituting (2.24) into (2.23) leads to

$$\begin{aligned}
& \frac{4}{mn} \sum_{i=1}^m n_i E\left\{\mathcal{W}(\mathbf{U}_i) \|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F\right\} \\
&\leq \frac{4}{mn} \left(\sum_{i=1}^m \sqrt{n_i}\right) E\left[\mathcal{W}(\mathbf{U}_i) \left\{E\left(\sum_{k=1}^{d_A} \sum_{k'=1}^{d_B} X_{Ak}^2 X_{Bk'}^2 f(\mathbf{X}, \mathbf{U}) \middle| \mathbf{U}\right)\right\}^{1/2}\right] \\
&\leq \frac{4}{\sqrt{\min_{1 \leq i \leq m} (n_i)}} E\left[\mathcal{W}(\mathbf{U}) \left\{E\left(\|\mathbf{X}_B \mathbf{X}_A^T\|_F^2 f(\mathbf{X}, \mathbf{U}) \middle| \mathbf{U}\right)\right\}^{1/2}\right] \\
&= \frac{4}{\sqrt{\min_{1 \leq i \leq m} (n_i)}} E\left[\mathcal{W}(\mathbf{U}) \left\{E\left(\|\mathbf{X}_A\|^2 \|\mathbf{X}_B\|^2 f(\mathbf{X}, \mathbf{U}) \middle| \mathbf{U}\right)\right\}^{1/2}\right].
\end{aligned} \quad (2.26)$$

The final expression in (2.26) converges to zero provided that we assume the moment condition

$$E\left[\mathcal{W}(\mathbf{U}) \left\{E\left(\|\mathbf{X}_A\|^2 \|\mathbf{X}_B\|^2 f(\mathbf{X}, \mathbf{U}) \middle| \mathbf{U}\right)\right\}^{1/2}\right] < \infty. \quad (2.27)$$

If the spectral norm in  $\mathcal{W}(\mathbf{U}_i)$  is replaced by the Frobenius norm, then (2.27) leads to

$$E \left[ \frac{\left\{ \|E\{\mathbf{X}_B \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}\|_F + 1 \right\} \left\{ E \left( \|\mathbf{X}_A\|^2 \|\mathbf{X}_B\|^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U} \right) \right\}^{1/2}}{\lambda_{\min} \left( E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \right)} \right] < \infty. \quad (2.28)$$

Now note that,

$$\begin{aligned} \|E\{\mathbf{X}_B \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}\|_F &\leq E\{\|\mathbf{X}_B \mathbf{X}_A^T\|_F f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \\ &= E\{\|\mathbf{X}_B\| \|\mathbf{X}_A\| f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}. \end{aligned}$$

This means that we can replace (2.28) by

$$E \left[ \frac{\left\{ E\{\|\mathbf{X}_A\| \|\mathbf{X}_B\| f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} + 1 \right\} \left\{ E \left( \|\mathbf{X}_A\|^2 \|\mathbf{X}_B\|^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U} \right) \right\}^{1/2}}{\lambda_{\min} \left( E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \right)} \right] < \infty. \quad (2.29)$$

Next we treat part of the second term in the final expression of (2.22). Note that,

$$\begin{aligned} E \left\{ \mathcal{W}(\mathbf{U}_i)^2 \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F \right\} &= E \left[ E \left\{ \mathcal{W}(\mathbf{U}_i)^2 \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F | \mathbf{U}_i \right\} \right] \\ &= E \left[ \mathcal{W}(\mathbf{U}_i)^2 E \left\{ \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F | \mathbf{U}_i \right\} \right]. \end{aligned} \quad (2.30)$$

By using a conditional version of the Cauchy-Schwarz inequality, we obtain

$$\begin{aligned} E \left\{ \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F | \mathbf{U}_i \right\} &= \left( \left[ E \left\{ \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F | \mathbf{U}_i \right\} \right]^2 \right)^{1/2} \\ &\leq \left[ E \left\{ \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F^2 | \mathbf{U}_i \right\} \right]^{1/2}. \end{aligned} \quad (2.31)$$

Using (2.31), note that

$$\|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F^2 = \sum_{k=1}^{d_A} \sum_{k'=1}^{d_A} [\{\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\}_{kk'}]^2.$$

Then note that

$$\left[ \{\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\}_{kk'} \right]^2 = \left( \frac{1}{n_i} \sum_{j=1}^{n_i} [X_{A_{ijk}} X_{A_{ijk'}} f(\mathbf{X}_{ij}, \mathbf{U}_i) - E\{X_{Ak} X_{Ak'} f(\mathbf{X}, \mathbf{U}_i) | \mathbf{U}_i\}] \right)^2.$$

Then we have,

$$\begin{aligned}
& E \left\{ \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F \middle| \mathbf{U}_i \right\} \\
& \leq \left[ E \left\{ \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F^2 \middle| \mathbf{U}_i \right\} \right]^{1/2} \\
& = \left\{ E \left( \sum_{k=1}^{d_A} \sum_{k'=1}^{d_A} [\{\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\}_{kk'}]^2 \middle| \mathbf{U}_i \right) \right\}^{1/2} \\
& = \left[ \sum_{k=1}^{d_A} \sum_{k'=1}^{d_A} E \left\{ \left( \frac{1}{n_i} \sum_{j=1}^{n_i} [X_{Aij} X_{Aijk'} f(\mathbf{X}_{ij}, \mathbf{U}_i) - E\{X_{Ak} X_{Ak'} f(\mathbf{X}, \mathbf{U}) | \mathbf{U}_i\}] \right)^2 \middle| \mathbf{U}_i \right\} \right]^{1/2} \\
& = \left[ \frac{1}{n_i} \sum_{k=1}^{d_A} \sum_{k'=1}^{d_A} \text{Var}\{X_{Ak} X_{Ak'} f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \right]^{1/2} \\
& \leq \left[ \frac{1}{n_i} \sum_{k=1}^{d_A} \sum_{k'=1}^{d_A} E\{X_{Ak}^2 X_{Ak'}^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \right]^{1/2}.
\end{aligned} \tag{2.32}$$

By substituting (2.32) into (2.30), we have that

$$\begin{aligned}
& \frac{4}{mn} \sum_{i=1}^m n_i E \left\{ \mathcal{W}(\mathbf{U}_i)^2 \|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F \right\} \\
& \leq \frac{4}{mn} \left( \sum_{i=1}^m \sqrt{n_i} \right) E \left[ \mathcal{W}(\mathbf{U})^2 \left\{ E \left( \sum_{k=1}^{d_A} \sum_{k'=1}^{d_A} X_{Ak}^2 X_{Ak'}^2 f(\mathbf{X}, \mathbf{U}) \middle| \mathbf{U} \right) \right\}^{1/2} \right] \\
& \leq \frac{4}{\sqrt{\min_{1 \leq i \leq m} (n_i)}} E \left[ \mathcal{W}(\mathbf{U})^2 \left\{ E \left( \|\mathbf{X}_A \mathbf{X}_A^T\|_F^2 f(\mathbf{X}, \mathbf{U}) \middle| \mathbf{U} \right) \right\}^{1/2} \right] \\
& = \frac{4}{\sqrt{\min_{1 \leq i \leq m} (n_i)}} E \left[ \mathcal{W}(\mathbf{U})^2 \left\{ E \left( \|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) \middle| \mathbf{U} \right) \right\}^{1/2} \right]
\end{aligned}$$

which converges to zero provided that we assume the moment condition

$$E \left[ \mathcal{W}(\mathbf{U})^2 \left\{ E \left( \|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) \middle| \mathbf{U} \right) \right\}^{1/2} \right] < \infty. \tag{2.33}$$

If the spectral norm in  $\mathcal{W}(\mathbf{U}_i)$  is replaced by the Frobenius norm, then (2.33) leads to

$$E \left[ \frac{\left\{ \|E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}\|_F + 1 \right\}^2 \left\{ E \left( \|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) \middle| \mathbf{U} \right) \right\}^{1/2}}{\lambda_{\min} \left( E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \right)^2} \right] < \infty. \tag{2.34}$$

Now note that

$$\begin{aligned}
\|E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}\|_F & \leq E\{\|\mathbf{X}_A \mathbf{X}_A^T\|_F f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \\
& = E\{\|\mathbf{X}_A\|^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}.
\end{aligned}$$

This means that we can replace (2.34) by

$$E \left[ \frac{[E\{\|\mathbf{X}_A\|^2 f(\mathbf{X}, \mathbf{U})|\mathbf{U}\} + 1]^2 \left\{ E\left(\|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U})|\mathbf{U}\right) \right\}^{1/2}}{\lambda_{\min}\left(E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U})|\mathbf{U}\}\right)^2} \right] < \infty. \quad (2.35)$$

#### 2.4.2.8 Proof of Result (2.17)

In this subsection we aim to prove (2.17), which is

$$\frac{1}{mn} \sum_{i=1}^m n_i E \left\| E \left[ \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\} I(\mathcal{A}_i^C) \middle| \mathcal{X}_i \right] \right\|_S \rightarrow 0$$

as  $m, n \rightarrow \infty$ . Note that throughout this subsection, we are considering

$$(\mathbf{U}_i, \mathcal{X}_i) \text{ such that } \mathcal{A}_i^C \text{ occurs, } 1 \leq i \leq m. \quad (2.36)$$

We first start with

$$\begin{aligned} & E \left\| E \left[ \left\{ \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\} I(\mathcal{A}_i^C) \middle| \mathcal{X}_i \right] \right\|_S \\ & \leq E \left( E \left[ \left\| \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) \right\|_S I(\mathcal{A}_i^C) \middle| \mathcal{X}_i \right] \right) + E \left( E \left[ \left\| \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\|_S I(\mathcal{A}_i^C) \middle| \mathcal{X}_i \right] \right) \\ & = E \left\{ \left\| \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) \right\|_S I(\mathcal{A}_i^C) \right\} + E \left\{ \left\| \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\|_S I(\mathcal{A}_i^C) \right\} \\ & \leq \left[ E \left\{ \left\| \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) \right\|_S^2 \right\} \right]^{1/2} P(\mathcal{A}_i^C)^{1/2} + \left[ E \left\{ \left\| \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\|_S^2 \right\} \right]^{1/2} P(\mathcal{A}_i^C)^{1/2}. \end{aligned}$$

To deal with the following expression

$$E \left\{ \left\| \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) \right\|_S^2 \right\},$$

first recall that

$$\begin{aligned} \left\| \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) \right\|_S & \leq \sqrt{d_B} \left\| \widehat{\mathcal{N}}_i(\mathbf{U}_i) \widehat{\mathcal{D}}_i(\mathbf{U}_i)^{-1} \widehat{\mathcal{N}}_i(\mathbf{U}_i)^T \right\|_S \\ & = \sqrt{d_B} \left\| \left\{ \frac{1}{n_i} \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T f(\mathbf{X}_{ij}, \mathbf{U}_i) \right\} \left\{ \frac{1}{n_i} \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T f(\mathbf{X}_{ij}, \mathbf{U}_i) \right\}^{-1} \right. \\ & \quad \left. \times \left\{ \frac{1}{n_i} \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Bij}^T f(\mathbf{X}_{ij}, \mathbf{U}_i) \right\} \right\|_S. \end{aligned}$$



Now we appeal to the second-last displayed equation on page 1093 of Chipman (1964), which is referred to as the *generalized Schwarz inequality*, to justify:

$$\left\| \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) \right\|_s^2 \leq \sqrt{d_B} \left\| \frac{1}{n_i} \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Bij}^T f(\mathbf{X}_{ij}, \mathbf{U}_i) \right\|_s^2.$$

Hence,

$$\begin{aligned} E \left\{ \left\| \mathcal{R} \left( \begin{bmatrix} \widehat{\mathcal{N}}_i(\mathbf{U}_i) \\ \widehat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) \right\|_s^2 \right\} &\leq \frac{d_B}{n_i} \sum_{j=1}^n E \{ \|\mathbf{X}_{Bij} \mathbf{X}_{Bij}^T f(\mathbf{X}_{ij}, \mathbf{U}_i)\|_F^2 \} \\ &= d_B E \{ \|\mathbf{X}_B\|^4 f(\mathbf{X}, \mathbf{U})^2 \}. \end{aligned}$$

Next, to deal with

$$E \left\{ \left\| \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\|_s^2 \right\}$$

note that

$$\begin{aligned} \left\| \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\|_s &\leq \sqrt{d_B} \|\mathcal{N}(\mathbf{U}_i) \mathcal{D}(\mathbf{U}_i)^{-1} \mathcal{N}(\mathbf{U}_i)^T\|_s \\ &\leq \sqrt{d_B} \|\mathcal{N}(\mathbf{U}_i)\|_s \|\mathcal{D}(\mathbf{U}_i)^{-1}\|_s \|\mathcal{N}(\mathbf{U}_i)\|_s \\ &= \frac{\sqrt{d_B} \|\mathcal{N}(\mathbf{U}_i)\|_s^2}{\lambda_{\min}(\mathcal{D}(\mathbf{U}_i))}. \end{aligned}$$

Hence,

$$E \left\{ \left\| \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\|_s^2 \right\} \leq d_B E \left\{ \frac{\|\mathcal{N}(\mathbf{U})\|_s^4}{\lambda_{\min}(\mathcal{D}(\mathbf{U}))^2} \right\}.$$

Finally, note that

$$\begin{aligned} \|\mathcal{N}(\mathbf{U})\|_s^4 &\leq \|E\{\mathbf{X}_B \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}\|_F^4 \\ &\leq E\{\|\mathbf{X}_B \mathbf{X}_A^T\|_F^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \\ &= E\{\|\mathbf{X}_A\|^4 \|\mathbf{X}_B\|^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \end{aligned}$$

which implies that

$$E \left\{ \frac{\|\mathcal{N}(\mathbf{U})\|_s^4}{\lambda_{\min}(\mathcal{D}(\mathbf{U}))^2} \right\} \leq E \left\{ \frac{E\{\|\mathbf{X}_A\|^4 \|\mathbf{X}_B\|^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}}{\lambda_{\min}(\mathcal{D}(\mathbf{U}))^2} \right\}.$$

From(2.36), note that the following holds

$$\begin{aligned} P(\mathcal{A}_i^C) &\leq P(\|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_s > 1) \\ &\quad + P\left(|\lambda_{\min}(\widehat{\mathcal{D}}_i(\mathbf{U}_i)) - \lambda_{\min}(\mathcal{D}(\mathbf{U}_i))| > \frac{1}{2} \lambda_{\min}(\mathcal{D}(\mathbf{U}_i))\right). \end{aligned} \tag{2.37}$$

Using Markov's inequality and the Cauchy-Schwarz inequality, we have

$$\begin{aligned}
P(\|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_S > 1) &\leq P(\|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_S^2 > 1) \\
&\leq E\{\|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_S^2\} \\
&\leq E\{\|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F^2\} \\
&= E[E\{\|\widehat{\mathcal{N}}_i(\mathbf{U}_i) - \mathcal{N}(\mathbf{U}_i)\|_F^2 | \mathbf{U}_i\}] \\
&\leq \frac{1}{n_i} \sum_{k=1}^{d_B} \sum_{k'=1}^{d_A} E[E\{X_{Bk}^2 X_{Ak'}^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}] \\
&= \frac{1}{n_i} E[E\{\|\mathbf{X}_A\|^2 \|\mathbf{X}_B\|^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}].
\end{aligned} \tag{2.38}$$

Using Markov's inequality again as well as Theorem 8.1.4 (Wielandt-Hoffman) of Golub and Van Loan (2013), we obtain the following expression for the second term in (2.37),

$$\begin{aligned}
&P\left(|\lambda_{\min}(\widehat{\mathcal{D}}_i(\mathbf{U}_i)) - \lambda_{\min}(\mathcal{D}(\mathbf{U}_i))| > \frac{1}{2} \lambda_{\min}(\mathcal{D}(\mathbf{U}_i))\right) \\
&= E\left\{P\left(|\lambda_{\min}(\widehat{\mathcal{D}}_i(\mathbf{U}_i)) - \lambda_{\min}(\mathcal{D}(\mathbf{U}_i))| > \frac{1}{2} \lambda_{\min}(\mathcal{D}(\mathbf{U}_i)) \mid \mathbf{U}_i\right)\right\} \\
&= E\left\{P\left(\{\lambda_{\min}(\widehat{\mathcal{D}}_i(\mathbf{U}_i)) - \lambda_{\min}(\mathcal{D}(\mathbf{U}_i))\}^2 > \frac{1}{4} \lambda_{\min}(\mathcal{D}(\mathbf{U}_i))^2 \mid \mathbf{U}_i\right)\right\} \\
&\leq 4E\left(\frac{E\left[\{\lambda_{\min}(\widehat{\mathcal{D}}_i(\mathbf{U}_i)) - \lambda_{\min}(\mathcal{D}(\mathbf{U}_i))\}^2 \mid \mathbf{U}_i\right]}{\lambda_{\min}(\mathcal{D}(\mathbf{U}_i))^2}\right) \\
&\leq 4E\left(\frac{E\left\{\|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F^2 \mid \mathbf{U}_i\right\}}{\lambda_{\min}(\mathcal{D}(\mathbf{U}_i))^2}\right).
\end{aligned}$$

From earlier calculations, we have that

$$\begin{aligned}
E\left\{\|\widehat{\mathcal{D}}_i(\mathbf{U}_i) - \mathcal{D}(\mathbf{U}_i)\|_F^2 \mid \mathbf{U}_i\right\} &\leq \frac{1}{n_i} \sum_{k=1}^{d_A} \sum_{k'=1}^{d_A} E\{X_{Ak}^2 X_{Ak'}^2 f(\mathbf{X}, \mathbf{U}) \mid \mathbf{U}\} \\
&= \frac{1}{n_i} E\left(\|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) \mid \mathbf{U}\right).
\end{aligned}$$

Therefore,

$$P\left(|\lambda_{\min}(\widehat{\mathcal{D}}_i(\mathbf{U}_i)) - \lambda_{\min}(\mathcal{D}(\mathbf{U}_i))| > \frac{1}{2} \lambda_{\min}(\mathcal{D}(\mathbf{U}_i))\right) \leq \frac{4}{n_i} E\left[\frac{E\{\|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) \mid \mathbf{U}\}}{\lambda_{\min}(\mathcal{D}(\mathbf{U}))^2}\right]. \tag{2.39}$$

Combining (2.38) and (2.39), we obtain

$$P(\mathcal{A}_i^C) \leq \frac{1}{n_i} \left( E\left[E\{\|\mathbf{X}_A\| \|\mathbf{X}_B\| f(\mathbf{X}, \mathbf{U}) \mid \mathbf{U}\}\right] + 4E\left[\frac{E\{\|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) \mid \mathbf{U}\}}{\lambda_{\min}(\mathcal{D}(\mathbf{U}))^2}\right] \right).$$

Using each of the several bounds established in this subsection, we have

$$\begin{aligned}
& \frac{1}{mn} \sum_{i=1}^m n_i E \left\| E \left[ \left\{ \mathcal{R} \left( \begin{bmatrix} \hat{\mathcal{N}}_i(\mathbf{U}_i) \\ \hat{\mathcal{D}}_i(\mathbf{U}_i) \end{bmatrix} \right) - \mathcal{R} \left( \begin{bmatrix} \mathcal{N}(\mathbf{U}_i) \\ \mathcal{D}(\mathbf{U}_i) \end{bmatrix} \right) \right\} I(\mathcal{A}_i^c) \middle| \mathcal{X}_i \right] \right\|_S \\
& \leq \frac{\sqrt{d_B}}{mn} \sum_{i=1}^m \left[ n_i \left( [E\{\|\mathbf{X}_B\|^4 f(\mathbf{X}, \mathbf{U})^2\}]^{1/2} + \left[ E \left\{ \frac{E\{\|\mathbf{X}_A\|^4 \|\mathbf{X}_B\|^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}}{\lambda_{\min}(\mathcal{D}(\mathbf{U}))^2} \right\} \right]^{1/2} \right) \right. \\
& \quad \times \left. \left\{ \frac{1}{n_i} \left( E[E\{\|\mathbf{X}_A\|^2 \|\mathbf{X}_B\|^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}] + 4E \left[ \frac{E\{\|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}}{\lambda_{\min}(\mathcal{D}(\mathbf{U}))^2} \right] \right) \right\}^{1/2} \right] \\
& \leq \frac{\sqrt{d_B}}{\sqrt{\min_{1 \leq i \leq m} (n_i)}} \left[ \left( [E\{\|\mathbf{X}_B\|^4 f(\mathbf{X}, \mathbf{U})^2\}]^{1/2} \right) \right. \\
& \quad \times \left. \left\{ \left( E[E\{\|\mathbf{X}_A\| \|\mathbf{X}_B\| f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}] + 4E \left[ \frac{E\{\|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}}{\lambda_{\min}(\mathcal{D}(\mathbf{U}))^2} \right] \right) \right\}^{1/2} \right] \\
& \quad + \frac{\sqrt{d_B}}{\sqrt{\min_{1 \leq i \leq m} (n_i)}} \left[ \left( \left[ E \left\{ \frac{E\{\|\mathbf{X}_A\|^4 \|\mathbf{X}_B\|^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}}{\lambda_{\min}(\mathcal{D}(\mathbf{U}))^2} \right\} \right]^{1/2} \right) \right. \\
& \quad \times \left. \left\{ \left( E[E\{\|\mathbf{X}_A\| \|\mathbf{X}_B\| f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}] + 4E \left[ \frac{E\{\|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}}{\lambda_{\min}(\mathcal{D}(\mathbf{U}))^2} \right] \right) \right\}^{1/2} \right],
\end{aligned}$$

which converges to zero under Theorems 11, 12 and 14 sample size and moment assumptions.

#### 2.4.2.9 Summary of Moment Assumptions

The moment assumptions used to prove (2.10), are as follows:

$$\begin{aligned}
(\text{MA1}) \quad & E\{\|\mathbf{X}_B\|^4 f(\mathbf{X}, \mathbf{U})^2\} < \infty, \\
(\text{MA2}) \quad & E \left[ \frac{\left\{ E \left( \|\mathbf{X}_A\|^2 \|\mathbf{X}_B\|^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U} \right) \right\}^{1/2} \{ E\{\|\mathbf{X}_A\| \|\mathbf{X}_B\| f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} + 1 \}}{\lambda_{\min} \left( E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \right)} \right] < \infty, \\
(\text{MA3}) \quad & E \left[ \frac{\left\{ E \left( \|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U} \right) \right\}^{1/2} [E\{\|\mathbf{X}_A\|^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} + 1]^2}{\lambda_{\min} \left( E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \right)^2} \right] < \infty, \\
(\text{MA4}) \quad & E[E\{\|\mathbf{X}_A\|^2 \|\mathbf{X}_B\|^2 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}] < \infty, \\
(\text{MA5}) \quad & E \left\{ \frac{E\{\|\mathbf{X}_A\|^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}}{\lambda_{\min} \left( E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \right)^2} \right\} < \infty, \\
(\text{MA6}) \quad & E \left\{ \frac{E\{\|\mathbf{X}_A\|^4 \|\mathbf{X}_B\|^4 f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\}}{\lambda_{\min} \left( E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\} \right)^2} \right\} < \infty.
\end{aligned}$$

### 2.4.2.10 Succinct Expression for Moment Assumptions

We now attempt to find an alternative way to express the moment assumptions in Section 2.4.2.9 in a unified fashion.

#### Inequality to Re-write the Numerator of (MA2)

To re-express the moment assumption (MA2), we intend to apply the following inequality to the numerator of (MA2):

$$(x + 1)y < 1 + x^2 + y^2 \quad \text{for all } x, y \in \mathbb{R}. \quad (2.40)$$

One way to prove (2.40) is to write it as

$$x^2 + y^2 - y - xy + 1 > 0 \quad \text{for all } x, y \in \mathbb{R}. \quad (2.41)$$

If the left-hand side is then written as

$$Ax^2 + By^2 + Cx + Dy + Exy + F$$

where

$$A = 1, \quad B = 1, \quad C = 0, \quad D = -1, \quad E = -1, \quad F = 1.$$

The discriminant-type quantity for the quadratic expression is

$$4AB - E^2 = 4 - (-1)^2 = 3 > 0$$

and since  $A > 0$ , the left-hand side of (2.41) has a minimum at

$$(x, y) = (DE - 2BC, CE - 2AD)/(4AB - E^2) = (1/3, 2/3).$$

Substitution of this point into the left-hand side of (2.41) leads to

$$(1/9) + (4/9) - (2/3) - (2/9) + 1 = (1/9) + (4/9) - (6/9) - (2/9) + (9/9) = 2/3 > 0.$$

Thus, we have established (2.40).

#### Inequality to Re-write the Numerator of (MA3)

Next, to re-express the moment assumption (MA3), we intend to apply the following inequality to the numerator of (MA3):

$$(x + 1)^2y < 1 + x^2 + x^4 + 2y^2 \quad \text{for all } x, y \in \mathbb{R}. \quad (2.42)$$

We first start by expanding the term on the left-hand side of (2.42) which leads to

$$\begin{aligned} (x + 1)^2y &= (x^2 + 2x + y)y \\ &= (x^2 + 1)y + 2xy. \end{aligned} \quad (2.43)$$

Also note that, by using the inequality  $(x - y)^2 \geq 0$ , we obtain the following

$$2xy \leq x^2 + y^2. \quad (2.44)$$

Hence, the result in (2.42) can be obtained using the inequalities in (2.40), (2.43) and (2.44).

### Re-expressing the Complete Set of Moment Conditions

Now we write

$$M_{p_1 p_2} = M_{p_1 p_2}(\mathbf{U}) \equiv E(\|\mathbf{X}_A\|^{p_1} \|\mathbf{X}_B\|^{p_2} f(\mathbf{X}, \mathbf{U}) | \mathbf{U}).$$

Then the application of (2.40) to (MA2) and application of (2.42) to (MA3) leads to

$$M_{22}^{1/2}(M_{11} + 1) \leq 1 + M_{11}^2 + M_{22}$$

and

$$M_{40}^{1/2}(M_{20} + 1)^2 \leq 1 + M_{20}^2 + M_{20}^4 + 2M_{40}.$$

Next, noting that

$$1 \geq \min(1, x), \quad x \geq \min(1, x) \quad \text{and} \quad x^2 \geq \{\min(1, x)\}^2 \quad \text{for all } x > 0,$$

the following inequalities can be applied to the denominators of (MA2)–(MA6):

$$\frac{1}{1} \leq \frac{1}{\{\min(1, x)\}^2}, \quad \frac{1}{x} \leq \frac{1}{\{\min(1, x)\}^2} \quad \text{and} \quad \frac{1}{x^2} \leq \frac{1}{\{\min(1, x)\}^2} \quad \text{for all } x > 0.$$

Combining all of these facts leads to the following alternative set of moment assumptions:

$$E \left[ \frac{\{E(\|\mathbf{X}_A\|^{p_1} \|\mathbf{X}_B\|^{p_2} f(\mathbf{X}, \mathbf{U}) | \mathbf{U})\}^{p_3}}{[\min\{1, \lambda_{\min}(E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\})\}]^2} \right] < \infty$$

for each of  $(p_1, p_2, p_3) \in \{(0, 0, 0), (4, 0, 1), (2, 2, 1), (4, 4, 1), (1, 1, 2), (2, 0, 2), (2, 0, 4)\}$ .

#### 2.4.2.11 A Sufficient Condition for the Moment Assumptions

Let  $\mathbf{X} = (\mathbf{X}_A, \mathbf{X}_B)$ . Then we have  $\|\mathbf{X}\| > \max(\|\mathbf{X}_A\|, \|\mathbf{X}_B\|)$ . Hence,

$$\begin{aligned} \{E(\|\mathbf{X}_A\|^{p_1} \|\mathbf{X}_B\|^{p_2} f(\mathbf{X}, \mathbf{U}) | \mathbf{U})\}^{p_3} &\leq \{E(\|\mathbf{X}\|^{p_1+p_2} f(\mathbf{X}, \mathbf{U}) | \mathbf{U})\}^{p_3} \\ &\leq E(\|\mathbf{X}\|^{(p_1+p_2)p_3} f(\mathbf{X}, \mathbf{U})^{p_3} | \mathbf{U}) \\ &\leq E\left(\max\{1, \|\mathbf{X}\|\}^{(p_1+p_2)p_3} \max\{1, f(\mathbf{X}, \mathbf{U})\}^{p_3} | \mathbf{U}\right) \\ &\leq E(\max\{1, \|\mathbf{X}\|\}^8 \max\{1, f(\mathbf{X}, \mathbf{U})\}^4 | \mathbf{U}) \end{aligned}$$

since  $p_3 \leq 4$  and  $(p_1 + p_2)p_3 \leq 8$  over the set of values that  $(p_1, p_2, p_3)$  takes in Section 2.4.2.10.

It follows that the condition

$$E \left[ \frac{E(\max\{1, \|\mathbf{X}\|\}^8 \max\{1, f(\mathbf{X}, \mathbf{U})\}^4 | \mathbf{U})}{[\min\{1, \lambda_{\min}(E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\})\}]^2} \right] < \infty \quad (2.45)$$

implies each of (MA2)-(MA6). Also, recalling that,

$$\frac{1}{\{\min(1, x)\}^2} \geq 1$$

we have

$$\frac{E(\max\{1, \|\mathbf{X}\|\}^8 \max\{1, f(\mathbf{X}, \mathbf{U})\}^4 | \mathbf{U})}{[\min\{1, \lambda_{\min}(E\{\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}\})\}]^2} \geq E\{\|\mathbf{X}_B\|^8 f(\mathbf{X}, \mathbf{U})^4 | \mathbf{U}\}.$$

Hence, moment condition (2.45) implies that

$$E [E\{\|\mathbf{X}_B\|^8 f(\mathbf{X}, \mathbf{U})^4 | \mathbf{U}\}] < \infty,$$

which is equivalent to

$$E [\{\|\mathbf{X}_B\|^4 f(\mathbf{X}, \mathbf{U})^2\}^2] < \infty. \quad (2.46)$$

However, since for any random variable  $Z$ , it is the case that

$$E(Z^2) < \infty \quad \text{implies} \quad E(|Z|) < \infty$$

condition (2.46) implies (MA1). Thus, we can claim that (2.45) implies each of (MA1)-(MA6).

### 2.4.3 Proof of Lemma 3

#### 2.4.3.1 Matrix Extension of Results Concerning Integrals of Half-Cauchy Forms

Using (1.12), premultiplication on both sides of the equation by  $(\frac{2}{\pi} \mathbf{A}^{1/2})^{-1}$  leads to

$$\frac{\pi}{2} \mathbf{I} = \int_0^\infty \mathbf{A}^{1/2} (\mathbf{A} + x^2 \mathbf{I})^{-1} dx. \quad (2.47)$$

Next, note that by using (2.47), we have,

$$\begin{aligned} \int_0^\infty (\mathbf{I} + x^2 \mathbf{A})^{-1} dx &= \int_0^\infty \{\mathbf{A} (\mathbf{A}^{-1} + x^2 \mathbf{I})\}^{-1} dx \\ &= \int_0^\infty (\mathbf{A}^{-1} + x^2 \mathbf{I})^{-1} \mathbf{A}^{-1} dx \\ &= \mathbf{A}^{1/2} \int_0^\infty \mathbf{A}^{-1/2} (\mathbf{A}^{-1} + x^2 \mathbf{I})^{-1} dx \mathbf{A}^{-1} \\ &= \mathbf{A}^{1/2} \left( \frac{\pi}{2} \mathbf{I} \right) \mathbf{A}^{-1} \\ &= \frac{\pi}{2} \mathbf{A}^{-1/2}. \end{aligned}$$

Therefore,

$$\frac{\pi}{2}\mathbf{I} = \int_0^\infty (\mathbf{I} + x^2\mathbf{A})^{-1} \mathbf{A}^{1/2} dx. \quad (2.48)$$

Combining both (2.47) and (2.48), for  $\mathbf{A} \in \mathbb{C}^{n \times n}$  with no eigenvalues on  $\mathbb{R}^-$ , we have,

$$\frac{\pi}{2}\mathbf{I} = \int_0^\infty \mathbf{A}^{1/2} (\mathbf{A} + x^2\mathbf{I})^{-1} dx = \int_0^\infty (\mathbf{I} + x^2\mathbf{A})^{-1} \mathbf{A}^{1/2} dx. \quad (2.49)$$

From these results, we have,

$$\mathbf{I}_p = \frac{4}{\pi^2} \int_0^\infty \int_0^\infty (\mathbf{I}_p + t^2\mathbf{K})^{-1} \mathbf{K} (\mathbf{K} + u^2\mathbf{I}_p)^{-1} dt du \quad (2.50)$$

and

$$\mathbf{I}_q = \frac{4}{\pi^2} \int_0^\infty \int_0^\infty \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n}\mathbf{L} \right) \right\}^{-1} \left( \frac{1}{n}\mathbf{L} \right) \left\{ \left( \frac{1}{n}\mathbf{L} \right) + u^2\mathbf{I}_q \right\}^{-1} dt du. \quad (2.51)$$

### 2.4.3.2 Derivation of Integrand Expressions

An Integral Expression for  $\mathbf{M}_{n,\infty}^{-1/2} = (\mathbf{M}_{n,\infty}^{-1})^{1/2}$

Note that

$$\mathbf{M}_{n,\infty} \equiv \begin{bmatrix} \mathbf{K} & \mathbf{O} \\ \mathbf{O} & \frac{1}{n}\mathbf{L} \end{bmatrix}$$

and

$$(\mathbf{M}_{n,\infty})^{-1} \equiv \begin{bmatrix} \mathbf{K}^{-1} & \mathbf{O} \\ \mathbf{O} & n\mathbf{L}^{-1} \end{bmatrix}.$$

Application of (2.49) to  $\mathbf{M}_{n,\infty}^{-1}$  leads to

$$\begin{aligned} \mathbf{M}_{n,\infty}^{-1/2} &= (\mathbf{M}_{n,\infty}^{-1})^{1/2} \\ &= \frac{2}{\pi} \int_0^\infty \mathbf{M}_{n,\infty}^{-1} (\mathbf{M}_{n,\infty}^{-1} + t^2\mathbf{I})^{-1} dt \\ &= \frac{2}{\pi} \int_0^\infty \begin{bmatrix} \mathbf{K}^{-1} & \mathbf{O} \\ \mathbf{O} & n\mathbf{L}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{K}^{-1} + t^2\mathbf{I}_p & \mathbf{O} \\ \mathbf{O} & n\mathbf{L}^{-1} + t^2\mathbf{I}_q \end{bmatrix}^{-1} dt \\ &= \frac{2}{\pi} \int_0^\infty \begin{bmatrix} \mathbf{K}^{-1} & \mathbf{O} \\ \mathbf{O} & n\mathbf{L}^{-1} \end{bmatrix} \begin{bmatrix} (\mathbf{K}^{-1} + t^2\mathbf{I}_p)^{-1} & \mathbf{O} \\ \mathbf{O} & (n\mathbf{L}^{-1} + t^2\mathbf{I}_q)^{-1} \end{bmatrix} dt \\ &= \frac{2}{\pi} \int_0^\infty \begin{bmatrix} \mathbf{K}^{-1} & \mathbf{O} \\ \mathbf{O} & n\mathbf{L}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{K} (\mathbf{I}_p + t^2\mathbf{K})^{-1} & \mathbf{O} \\ \mathbf{O} & \left( \frac{1}{n}\mathbf{L} \right) \{ \mathbf{I}_q + t^2 \left( \frac{1}{n}\mathbf{L} \right) \}^{-1} \end{bmatrix} dt \\ &= \frac{2}{\pi} \int_0^\infty \begin{bmatrix} (\mathbf{I}_p + t^2\mathbf{K})^{-1} & \mathbf{O} \\ \mathbf{O} & \{ \mathbf{I}_q + t^2 \left( \frac{1}{n}\mathbf{L} \right) \}^{-1} \end{bmatrix} dt. \end{aligned}$$

An Integral Expression for  $\mathbf{M}_n^{1/2}$

Note that

$$\mathbf{M}_n \equiv \begin{bmatrix} \mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2} & R_n \mathbf{1}_p \mathbf{1}_q^T \\ R_n \mathbf{1}_q \mathbf{1}_p^T & \frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2} \end{bmatrix}$$

where  $Q_n = o_p(1)$ ,  $R_n = O_p(n^{-1})$  and  $T_n = o_p(n^{-1})$ . For all  $n$  sufficiently large so that negative eigenvalues are avoided, application of (2.49) to  $\mathbf{M}_n$  leads to

$$\begin{aligned} \mathbf{M}_n^{1/2} &= \frac{2}{\pi} \int_0^\infty \mathbf{M}_n (\mathbf{M}_n + u^2 \mathbf{I})^{-1} du \\ &= \frac{2}{\pi} \int_0^\infty \begin{bmatrix} \mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2} & R_n \mathbf{1}_p \mathbf{1}_q^T \\ R_n \mathbf{1}_q \mathbf{1}_p^T & \frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2} \end{bmatrix} \begin{bmatrix} \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} & R_n \mathbf{1}_p \mathbf{1}_q^T \\ R_n \mathbf{1}_q \mathbf{1}_p^T & \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \end{bmatrix}^{-1} du. \end{aligned}$$

An Integral Expression for  $\mathbf{M}_{n,\infty}^{-1/2} \mathbf{M}_n^{1/2}$

Firstly, note that

$$\begin{aligned} &\begin{bmatrix} (\mathbf{I}_p + t^2 \mathbf{K})^{-1} & \mathbf{O} \\ \mathbf{O} & \{\mathbf{I}_q + t^2 (\frac{1}{n} \mathbf{L})\}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2} & R_n \mathbf{1}_p \mathbf{1}_q^T \\ R_n \mathbf{1}_q \mathbf{1}_p^T & \frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2} \end{bmatrix} \\ &= \begin{bmatrix} (\mathbf{I}_p + t^2 \mathbf{K})^{-1} (\mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2}) & R_n (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{1}_p \mathbf{1}_q^T \\ R_n \{\mathbf{I}_q + t^2 (\frac{1}{n} \mathbf{L})\}^{-1} \mathbf{1}_q \mathbf{1}_p^T & \{\mathbf{I}_q + t^2 (\frac{1}{n} \mathbf{L})\}^{-1} \frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2} \end{bmatrix}. \end{aligned}$$

Then

$$\begin{aligned} &\mathbf{M}_{n,\infty}^{-1/2} \mathbf{M}_n^{1/2} \\ &= \frac{4}{\pi^2} \int_0^\infty \int_0^\infty \begin{bmatrix} (\mathbf{I}_p + t^2 \mathbf{K})^{-1} (\mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2}) & R_n (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{1}_p \mathbf{1}_q^T \\ R_n \{\mathbf{I}_q + t^2 (\frac{1}{n} \mathbf{L})\}^{-1} \mathbf{1}_q \mathbf{1}_p^T & \{\mathbf{I}_q + t^2 (\frac{1}{n} \mathbf{L})\}^{-1} \frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2} \end{bmatrix} \\ &\quad \times \begin{bmatrix} \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} & R_n \mathbf{1}_p \mathbf{1}_q^T \\ R_n \mathbf{1}_q \mathbf{1}_p^T & \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \end{bmatrix}^{-1} dt du. \end{aligned}$$

An Integral Expression for  $\mathbf{M}_{n,\infty}^{-1/2} \mathbf{M}_n^{1/2} - \mathbf{I}$

For any  $u > 0$  and values of  $n$ ,  $\mathbf{K}$  and  $\mathbf{L}$ , define the  $(p+q) \times (p+q)$  matrix,

$$\mathbf{H}_1(u; n, \mathbf{K}, \mathbf{L}) \equiv \begin{bmatrix} \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} & R_n \mathbf{1}_p \mathbf{1}_q^T \\ R_n \mathbf{1}_q \mathbf{1}_p^T & \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \end{bmatrix}^{-1}. \quad (2.52)$$



Then, for any  $(t, u) \in \mathbb{R}^2$ , define the  $(p+q) \times (p+q)$  matrix as follows,

$$\begin{aligned} & \mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L}) \\ & \equiv \begin{bmatrix} (\mathbf{I}_p + t^2 \mathbf{K})^{-1} (\mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2}) & R_n (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{1}_p \mathbf{1}_q^T \\ R_n \{\mathbf{I}_q + t^2 (\frac{1}{n} \mathbf{L})\}^{-1} \mathbf{1}_q \mathbf{1}_p^T & \{\mathbf{I}_q + t^2 (\frac{1}{n} \mathbf{L})\}^{-1} (\frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2}) \end{bmatrix} \mathbf{H}_1(u; n, \mathbf{K}, \mathbf{L}) \\ & - \begin{bmatrix} (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{K} (\mathbf{K} + u^2 \mathbf{I}_p)^{-1} & \mathbf{O} \\ \mathbf{O} & \{\mathbf{I}_q + t^2 (\frac{1}{n} \mathbf{L})\}^{-1} (\frac{1}{n} \mathbf{L}) \{\frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q\}^{-1} \end{bmatrix}. \end{aligned} \quad (2.53)$$

Then from (2.50) and (2.51) and according to the definitions above,

$$\mathbf{M}_{n, \infty}^{-1/2} \mathbf{M}_n^{1/2} - \mathbf{I}_{p+q} = \frac{4}{\pi^2} \int_0^\infty \int_0^\infty \mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L}) dt du. \quad (2.54)$$

Throughout the rest of this subsection, we aim to find a more explicit expression for  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$ .

#### Inversion of $\mathbf{H}_1(u; n, \mathbf{K}, \mathbf{L})$ Using Block Matrix Inversion

The upper left  $p \times p$  block of  $\mathbf{H}_1(u; n, \mathbf{K}, \mathbf{L})^{-1}$  is

$$\left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1}.$$

The upper right  $p \times q$  block of  $\mathbf{H}_1(u; n, \mathbf{K}, \mathbf{L})^{-1}$  is

$$\begin{aligned} & - R_n \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\ & \times \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1}. \end{aligned}$$

The lower left  $q \times p$  block of  $\mathbf{H}_1(u; n, \mathbf{K}, \mathbf{L})^{-1}$  is

$$\begin{aligned} & - R_n \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \\ & \times \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1}. \end{aligned}$$

The lower right  $q \times q$  block of  $\mathbf{H}_1(u; n, \mathbf{K}, \mathbf{L})^{-1}$  is

$$\left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T \left( \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} \right)^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1}.$$

The Upper Left  $p \times p$  Block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$

The upper left  $p \times p$  block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$  is

$$\begin{aligned}
& (\mathbf{I}_p + t^2 \mathbf{K})^{-1} (\mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2}) \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\
& \quad - R_n^2 (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \\
& \quad \times \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\
& \quad - (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{K} (\mathbf{K} + u^2 \mathbf{I}_p)^{-1} \\
& = \left\{ (\mathbf{I}_p + t^2 \mathbf{K})^{-1} (\mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2}) - R_n^2 (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\} \\
& \quad \times \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\
& \quad - (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{K} (\mathbf{K} + u^2 \mathbf{I}_p)^{-1} \\
& = (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \left[ \left\{ \mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\} \right. \\
& \quad \times \left. \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \right. \\
& \quad \left. - \mathbf{K} (\mathbf{K} + u^2 \mathbf{I}_p)^{-1} \right] \\
& = (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{K} \left[ \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \right. \\
& \quad \left. - (\mathbf{K} + u^2 \mathbf{I}_p)^{-1} \right] + (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \left\{ Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\} \\
& \quad \times \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1}.
\end{aligned}$$

Using (1.1), the first term in the expression can be re-written as

$$\begin{aligned}
& (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{K} \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\
& \quad \times \left\{ R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T - Q_n \mathbf{1}_p^{\otimes 2} \right\} (\mathbf{K} + u^2 \mathbf{I}_p)^{-1}.
\end{aligned}$$

Putting together the results given so far in this subsection, we have that the upper left  $p \times p$  block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$  is:

$$\begin{aligned} & (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{K} \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\ & \times \left\{ R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T - Q_n \mathbf{1}_p^{\otimes 2} \right\} (\mathbf{K} + u^2 \mathbf{I}_p)^{-1} \\ & - (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \left\{ R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T - Q_n \mathbf{1}_p^{\otimes 2} \right\} \\ & \times \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1}. \end{aligned}$$

The Lower Left  $q \times p$  Block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$

Noting that the lower left  $q \times p$  block of  $\mathbf{I}$  is a zero matrix, the lower left  $q \times p$  block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$  is

$$\begin{aligned} & R_n \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \mathbf{1}_q \mathbf{1}_p^T \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\ & - R_n \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \left( \frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2} \right) \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \\ & \times \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\ & = R_n \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \left\{ \mathbf{I}_q - \left( \frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2} \right) \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \right\} \mathbf{1}_q \mathbf{1}_p^T \\ & \times \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\ & = R_n u^2 \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \\ & \times \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1}. \end{aligned}$$

The Upper Right  $p \times q$  Block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$

Noting that the upper right  $p \times q$  block of  $\mathbf{I}$  is a zero matrix, the upper right  $p \times q$

block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$  is

$$\begin{aligned}
& -R_n(\mathbf{I}_p + t^2 \mathbf{K})^{-1}(\mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2}) \\
& \quad \times \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\
& \quad \times \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \\
& \quad + R_n(\mathbf{I}_p + t^2 \mathbf{K})^{-1} \mathbf{1}_p \mathbf{1}_q^T \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1} \\
& = R_n(\mathbf{I}_p + t^2 \mathbf{K})^{-1} \left[ \mathbf{1}_p \mathbf{1}_q^T \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1} \right. \\
& \quad - (\mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2}) \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \mathbf{1}_p \mathbf{1}_q^T \\
& \quad \left. \times \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \right].
\end{aligned}$$

The Lower Right  $q \times q$  Block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$

The lower right  $q \times q$  block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$  is

$$\begin{aligned}
& \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \left( \frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2} \right) \\
& \quad \times \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1} - R_n^2 \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \\
& \quad \times \mathbf{1}_q \mathbf{1}_p^T \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\
& \quad \times \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} - \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \left( \frac{1}{n} \mathbf{L} \right) \left\{ \left( \frac{1}{n} \mathbf{L} \right) + u^2 \mathbf{I}_q \right\}^{-1} \\
& = \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \\
& \quad \times \left[ \left( \frac{1}{n} \mathbf{L} + T_n \mathbf{1}_q^{\otimes 2} \right) \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1} \right. \\
& \quad - \left( \frac{1}{n} \mathbf{L} \right) \left\{ \left( \frac{1}{n} \mathbf{L} \right) + u^2 \mathbf{I}_q \right\}^{-1} \\
& \quad - R_n^2 \mathbf{1}_q \mathbf{1}_p^T \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\
& \quad \left. \times \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \right].
\end{aligned}$$

Therefore, the lower right  $q \times q$  block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$  is

$$\begin{aligned}
&= \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \\
&\quad \times \left( \left( \frac{1}{n} \mathbf{L} \right) \left[ \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1} \right. \right. \\
&\quad \left. \left. - \left\{ \left( \frac{1}{n} \mathbf{L} \right) + u^2 \mathbf{I}_q \right\}^{-1} \right] + T_n \mathbf{1}_q^{\otimes 2} \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q \right. \right. \\
&\quad \left. \left. + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1} \right. \\
&\quad \left. - R_n^2 \mathbf{1}_q \mathbf{1}_p^T \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \right. \\
&\quad \left. \times \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \right).
\end{aligned}$$

Application of (1.1) to the first term in the preceding expression leads to

$$\begin{aligned}
&\left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \left( \frac{1}{n} \mathbf{L} \right) \\
&\quad \times \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1} \\
&\quad \times \left\{ R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T - T_n \mathbf{1}_q^{\otimes 2} \right\} \left\{ \left( \frac{1}{n} \mathbf{L} \right) + u^2 \mathbf{I}_q \right\}^{-1}.
\end{aligned}$$

The next term of the lower right  $q \times q$  block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$  is

$$\begin{aligned}
&T_n \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \mathbf{1}_q^{\otimes 2} \\
&\quad \times \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1}.
\end{aligned}$$

Hence, the final term of the lower right  $q \times q$  block of  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$  is

$$\begin{aligned}
&- R_n^2 \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \mathbf{1}_q \mathbf{1}_p^T \\
&\quad \times \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\
&\quad \times \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1}.
\end{aligned}$$

Putting all these expressions together, we have,

$$\begin{aligned}
& \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \left[ \left( \frac{1}{n} \mathbf{L} \right) \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1} \right. \\
& \quad \times \left\{ R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T - T_n \mathbf{1}_q^{\otimes 2} \right\} \left\{ \left( \frac{1}{n} \mathbf{L} \right) + u^2 \mathbf{I}_q \right\}^{-1} \\
& \quad + T_n \mathbf{1}_q^{\otimes 2} \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T \right\}^{-1} \\
& \quad - R_n^2 \mathbf{1}_q \mathbf{1}_p^T \left\{ \mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2} - R_n^2 \mathbf{1}_q \mathbf{1}_p^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \mathbf{1}_q \mathbf{1}_p^T \right\}^{-1} \\
& \quad \left. \times \mathbf{1}_p \mathbf{1}_q^T \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \right].
\end{aligned}$$

### 2.4.3.3 Succinct Expressions for the Components in (2.53)

Define

$$\begin{aligned}
\Gamma_{1n}(u; T_n, \mathbf{K}, \mathbf{L}) &\equiv \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \\
\Gamma_{2n}(u; Q_n, R_n, \mathbf{K}, \mathbf{L}) &\equiv R_n^2 \mathbf{1}_p \mathbf{1}_q^T (\Gamma_{1n}(u; T_n, \mathbf{K}, \mathbf{L})) \mathbf{1}_q \mathbf{1}_p^T - Q_n \mathbf{1}_p^{\otimes 2} \\
\Gamma_{3n}(u; Q_n, R_n, T_n, \mathbf{K}, \mathbf{L}) &\equiv R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T - T_n \mathbf{1}_q^{\otimes 2} \\
\Gamma_{4n}(u; Q_n, R_n, \mathbf{K}, \mathbf{L}) &\equiv \left\{ \mathbf{K} + u^2 \mathbf{I}_p - (\Gamma_{2n}(u; Q_n, R_n, \mathbf{K}, \mathbf{L})) \right\}^{-1} \quad \text{and} \\
\Gamma_{5n}(u; Q_n, R_n, T_n, \mathbf{K}, \mathbf{L}) &\equiv \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q - (\Gamma_{3n}(u; Q_n, R_n, T_n, \mathbf{K}, \mathbf{L})) \right\}^{-1}.
\end{aligned}$$

From now on, we write our expressions for  $\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})$  in terms of  $\Gamma_{1n}(u), \dots, \Gamma_{5n}(u)$  and suppress all other arguments. Then,

$$\begin{aligned}
\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})_{11} &= (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \left\{ \mathbf{K} \Gamma_{4n}(u) \Gamma_{2n}(u) (\mathbf{K} + u^2 \mathbf{I}_p)^{-1} - \Gamma_{2n}(u) \Gamma_{4n}(u) \right\} \\
\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})_{21} &= R_n u^2 \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \Gamma_{1n}(u) \mathbf{1}_q \mathbf{1}_p^T \Gamma_{4n}(u) \\
\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})_{12} &= R_n (\mathbf{I}_p + t^2 \mathbf{K})^{-1} \left\{ \mathbf{1}_p \mathbf{1}_q^T \Gamma_{5n}(u) - (\mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2}) \Gamma_{4n}(u) \mathbf{1}_p \mathbf{1}_q^T \Gamma_{1n}(u) \right\} \quad \text{and} \\
\mathbf{H}_2(t, u; n, \mathbf{K}, \mathbf{L})_{22} &= \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} \left[ \left( \frac{1}{n} \mathbf{L} \right) \Gamma_{5n}(u) \Gamma_{3n}(u) \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q \right\}^{-1} \right. \\
& \quad \left. + T_n \mathbf{1}_q^{\otimes 2} \Gamma_{5n}(u) - R_n^2 \mathbf{1}_q \mathbf{1}_p^T \Gamma_{4n}(u) \mathbf{1}_p \mathbf{1}_q^T \Gamma_{1n}(u) \right].
\end{aligned}$$

### 2.4.3.4 Simplification of Integrals

It follows from the results in Section 2.4.3.1 that

$$\int_0^\infty (\mathbf{I}_p + t^2 \mathbf{K})^{-1} dt = \frac{\pi}{2} \mathbf{K}^{-1/2}.$$

Also, using (2.49),

$$\int_0^\infty \left\{ \mathbf{I}_q + t^2 \left( \frac{1}{n} \mathbf{L} \right) \right\}^{-1} dt = \frac{\pi}{2} \left( \frac{1}{n} \mathbf{L} \right)^{-1/2} = \frac{\pi n^{1/2}}{2} \mathbf{L}^{-1/2}.$$

### 2.4.3.5 Explicit Expressions for (2.54)

Define

$$\frac{\pi}{2} \left( M_{n,\infty}^{-1/2} M_n^{1/2} - \mathbf{I}_{p+q} \right) = \begin{bmatrix} \mathbf{K}^{-1/2} \int_0^\infty \mathbf{F}_{11n}(u; \mathbf{K}, \mathbf{L}) du & \mathbf{K}^{-1/2} \int_0^\infty \mathbf{F}_{12n}(u; \mathbf{K}, \mathbf{L}) du \\ \mathbf{L}^{-1/2} \int_0^\infty \mathbf{F}_{21n}(u; \mathbf{K}, \mathbf{L}) du & \mathbf{L}^{-1/2} \int_0^\infty \mathbf{F}_{22n}(u; \mathbf{K}, \mathbf{L}) du \end{bmatrix} \quad (2.55)$$

where

$$\begin{aligned} \mathbf{F}_{11n}(u; \mathbf{K}, \mathbf{L}) &\equiv \mathbf{K} \Gamma_{4n}(u) \Gamma_{2n}(u) (\mathbf{K} + u^2 \mathbf{I}_p)^{-1} - \Gamma_{2n}(u) \Gamma_{4n}(u) \\ \mathbf{F}_{21n}(u; \mathbf{K}, \mathbf{L}) &\equiv n^{1/2} R_n u^2 \Gamma_{1n}(u) \mathbf{1}_q \mathbf{1}_p^T \Gamma_{4n}(u) \\ \mathbf{F}_{12n}(u; \mathbf{K}, \mathbf{L}) &\equiv R_n \{ \mathbf{1}_p \mathbf{1}_q^T \Gamma_{5n}(u) - (\mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2}) \Gamma_{4n}(u) \mathbf{1}_p \mathbf{1}_q^T \Gamma_{1n}(u) \} \quad \text{and} \\ \mathbf{F}_{22n}(u; \mathbf{K}, \mathbf{L}) &\equiv n^{1/2} \left[ \left( \frac{1}{n} \mathbf{L} \right) \Gamma_{5n}(u) \Gamma_{3n}(u) \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q \right\}^{-1} \right. \\ &\quad \left. + T_n \mathbf{1}_q^{\otimes 2} \Gamma_{5n}(u) - R_n^2 \mathbf{1}_q \mathbf{1}_p^T \Gamma_{4n}(u) \mathbf{1}_p \mathbf{1}_q^T \Gamma_{1n}(u) \right]. \end{aligned}$$

### 2.4.3.6 Convergence in Probability Limits of the Functions in (2.55)

In Appendix 2.4.4, we establish that

$$\text{plim}_{n \rightarrow \infty} \int_0^\infty \mathbf{F}_{kk'n}(u; \mathbf{K}, \mathbf{L}) du = \mathbf{O}, \quad k, k' = 1, 2.$$

for  $u > 0$ . Hence, the lemma in Section 2.3 is proven.

## 2.4.4 Multivariate Integral Limits for the Matrix Square Root Result

### 2.4.4.1 Overview of this Appendix

In this appendix we deal with the problem of establishing that

$$\text{plim}_{n \rightarrow \infty} \int_0^\infty \mathbf{F}_{kk'n}(u; \mathbf{K}, \mathbf{L}) du = \mathbf{O}, \quad k, k' = 1, 2.$$

The approach involves spectral norm bounds on the integrands  $\mathbf{F}_{kk'n}(u; \mathbf{K}, \mathbf{L})$  uniformly over  $u > 0$  and exact integral results over the positive half-line for functions of  $x$  with factors of the form  $1/(a_j + x^2)$ ,  $a_j > 0$ .

### 2.4.4.2 Computing Spectral Norms

#### Bounding of $\|\mathbf{\Gamma}_{1n}(u)\|_s$

Recall that

$$\mathbf{\Gamma}_{1n}(u) \equiv \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1}.$$

Then, noting that  $\mathbf{\Gamma}_{1n}(u)$  is symmetric and positive definite,

$$\|\mathbf{\Gamma}_{1n}(u)\|_s = \lambda_{\max} \left\{ \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right)^{-1} \right\} = 1 / \lambda_{\min} \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right).$$

Application of (1.9) leads to

$$\begin{aligned} \lambda_{\min} \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right) &\geq \lambda_{\min} \left( \frac{1}{n} \mathbf{L} \right) + \lambda_{\min} (u^2 \mathbf{I}_q) + \lambda_{\min} (T_n \mathbf{1}_q^{\otimes 2}) \\ &= \frac{1}{n} \lambda_{\min} (\mathbf{L}) + u^2 \lambda_{\min} (\mathbf{I}_q) + T_n \lambda_{\min} (\mathbf{1}_q^{\otimes 2}) \\ &\geq \frac{1}{n} \lambda_{\min} (\mathbf{L}) + u^2 + \min (T_n, 0). \end{aligned}$$

Since  $T_n = o_P(n^{-1})$ , for every  $0 < \varepsilon \leq 1$  we can choose  $n_1 \in \mathbb{N}$  such that, for all  $n > n_1$ ,  $|T_n| < \frac{1}{2n} \lambda_{\min} (\mathbf{L})$  with probability exceeding  $1 - \varepsilon$ . For all such large  $n$  we then have

$$T_n > -\frac{1}{2n} \lambda_{\min} (\mathbf{L}).$$

Also,

$$0 > -\frac{1}{2n} \lambda_{\min} (\mathbf{L})$$

and therefore we have,

$$\min (T_n, 0) > -\frac{1}{2n} \lambda_{\min} (\mathbf{L}).$$



Therefore, for all such large  $n$ , we then have

$$\lambda_{\min} \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q + T_n \mathbf{1}_q^{\otimes 2} \right) > \frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2 \quad \text{for all } u > 0.$$

This leads to

$$\|\mathbf{\Gamma}_{1n}(u)\|_s < \frac{1}{\frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2} \quad \text{for all } n > n_1 \text{ and } u > 0$$

with probability exceeding  $1 - \varepsilon$ .

#### Bounding of $\|\mathbf{\Gamma}_{2n}(u)\|_s$

Recall that

$$\begin{aligned} \mathbf{\Gamma}_{2n}(u) &\equiv R_n^2 \mathbf{1}_p \mathbf{1}_q^T \mathbf{\Gamma}_{1n}(u) \mathbf{1}_q \mathbf{1}_p^T - Q_n \mathbf{1}_p^{\otimes 2} \\ &= \{R_n^2 \mathbf{1}_q^T \mathbf{\Gamma}_{1n}(u) \mathbf{1}_q - Q_n\} \mathbf{1}_p^{\otimes 2}. \end{aligned}$$

Also note that for  $n$  large enough and all  $u > 0$

$$\|\mathbf{\Gamma}_{1n}(u)\|_s < \frac{1}{\frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2} < \frac{1}{\frac{1}{2n} \lambda_{\min}(\mathbf{L})} = \frac{2n}{\lambda_{\min}(\mathbf{L})}.$$

Then

$$\begin{aligned} \|\mathbf{\Gamma}_{2n}(u)\|_s &= \|\{R_n^2 \mathbf{1}_q^T \mathbf{\Gamma}_{1n}(u) \mathbf{1}_q - Q_n\} \mathbf{1}_p^{\otimes 2}\|_s \\ &\leq \|R_n^2 \mathbf{1}_q^T \mathbf{\Gamma}_{1n}(u) \mathbf{1}_q - Q_n\|_s \|\mathbf{1}_p^{\otimes 2}\|_s \\ &\leq p \{R_n^2 \|\mathbf{1}_q^T \mathbf{\Gamma}_{1n}(u) \mathbf{1}_q\|_s + \|Q_n\|_s\} \\ &\leq p \{q R_n^2 \|\mathbf{\Gamma}_{1n}(u)\|_s + |Q_n|\} \\ &< p \left\{ \frac{2qn R_n^2}{\lambda_{\min}(\mathbf{L})} + |Q_n| \right\}. \end{aligned}$$

In summary,

$$\|\mathbf{\Gamma}_{2n}(u)\|_s < p \left\{ \frac{2qn R_n^2}{\lambda_{\min}(\mathbf{L})} + |Q_n| \right\} \quad \text{for all } u > 0$$

with probability exceeding  $1 - \varepsilon$ .

#### Bounding of $\|\mathbf{\Gamma}_{3n}(u)\|_s$

Recall that

$$\mathbf{\Gamma}_{3n}(u) \equiv R_n^2 \mathbf{1}_q \mathbf{1}_p^T (\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1} \mathbf{1}_p \mathbf{1}_q^T - T_n \mathbf{1}_q^{\otimes 2}.$$

Then

$$\begin{aligned} \|\mathbf{\Gamma}_{3n}(u)\|_s &\leq R_n^2 \|\mathbf{1}_q\|_s \|\mathbf{1}_p^T\|_s \|(\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})^{-1}\|_s \|\mathbf{1}_p\|_s \|\mathbf{1}_q^T\|_s + |T_n| \|\mathbf{1}_q^{\otimes 2}\|_s \\ &= \frac{pq R_n^2}{\lambda_{\min}(\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2})} + q |T_n|. \end{aligned}$$

Now,

$$\begin{aligned}\lambda_{\min}(\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2}) &\geq \lambda_{\min}(\mathbf{K}) + \lambda_{\min}(u^2 \mathbf{I}_p) + \lambda_{\min}(Q_n \mathbf{1}_p^{\otimes 2}) \\ &\geq \lambda_{\min}(\mathbf{K}) + u^2 + \min(Q_n, 0).\end{aligned}$$

Since  $Q_n = o_P(1)$ , for every  $0 < \varepsilon \leq 1$  we can choose  $n_3 \in \mathbb{N}$  such that, for all  $n > n_3$ ,  $|Q_n| < \frac{1}{2} \lambda_{\min}(\mathbf{K})$  with probability exceeding  $1 - \varepsilon$ . For all such large  $n$  we then have

$$Q_n > -\frac{1}{2} \lambda_{\min}(\mathbf{K}).$$

Also,

$$0 > -\frac{1}{2} \lambda_{\min}(\mathbf{K})$$

and therefore we have,

$$\min(Q_n, 0) > -\frac{1}{2} \lambda_{\min}(\mathbf{K}).$$

Therefore, for all such large  $n$ , we then have

$$\lambda_{\min}(\mathbf{K} + u^2 \mathbf{I}_p + Q_n \mathbf{1}_p^{\otimes 2}) > \frac{1}{2} \lambda_{\min}(\mathbf{K}) + u^2.$$

Therefore,

$$\|\mathbf{\Gamma}_{3n}(u)\|_s < \frac{pqR_n^2}{\frac{1}{2} \lambda_{\min}(\mathbf{K}) + u^2} + q|T_n| \quad \text{for all } n > n_3 \text{ and } u > 0.$$

#### Bounding of $\|\mathbf{\Gamma}_{4n}(u)\|_s$

Recall that

$$\mathbf{\Gamma}_{4n}(u) \equiv \{\mathbf{K} + u^2 \mathbf{I}_p - \mathbf{\Gamma}_{2n}(u)\}^{-1}.$$

Hence,

$$\begin{aligned}\|\mathbf{\Gamma}_{4n}(u)\|_s &= \lambda_{\max} \left\{ (\mathbf{K} + u^2 \mathbf{I}_p - \mathbf{\Gamma}_{2n}(u))^{-1} \right\} \\ &= 1 / \lambda_{\min}(\mathbf{K} + u^2 \mathbf{I}_p - \mathbf{\Gamma}_{2n}(u)).\end{aligned}$$

Then, from (1.9) we have

$$\begin{aligned}\lambda_{\min}(\mathbf{K} + u^2 \mathbf{I}_p - \mathbf{\Gamma}_{2n}(u)) &\geq \lambda_{\min}(\mathbf{K} + u^2 \mathbf{I}_p) + \lambda_{\min}(-\mathbf{\Gamma}_{2n}(u)) \\ &\geq \lambda_{\min}(\mathbf{K}) + \lambda_{\min}(u^2 \mathbf{I}_p) - \lambda_{\max}(\mathbf{\Gamma}_{2n}(u)) \\ &\geq \lambda_{\min}(\mathbf{K}) + u^2 - \lambda_{\max}(\mathbf{\Gamma}_{2n}(u)) \\ &= \lambda_{\min}(\mathbf{K}) + u^2 - \|\mathbf{\Gamma}_{2n}(u)\|_s \\ &\geq \lambda_{\min}(\mathbf{K}) + u^2 - p \left\{ \frac{2qnR_n^2}{\lambda_{\min}(\mathbf{L})} + |Q_n| \right\}.\end{aligned}$$

Since  $R_n = O_P(n^{-1})$  and  $Q_n = o_P(1)$ , for every  $0 < \varepsilon \leq 1$  we can choose  $n_4 \in \mathbb{N}$  such that, for all  $n > n_4$ ,

$$p \left\{ \frac{2qnR_n^2}{\lambda_{\min}(\mathbf{L})} + |Q_n| \right\} < \frac{1}{2} \lambda_{\min}(\mathbf{K}).$$

This is equivalent to

$$-p \left\{ \frac{2qnR_n^2}{\lambda_{\min}(\mathbf{L})} + |Q_n| \right\} > -\frac{1}{2}\lambda_{\min}(\mathbf{K}).$$

This allows us to claim that for all  $n > n_4$ ,

$$\lambda_{\min}(\mathbf{K} + u^2\mathbf{I}_p - \mathbf{\Gamma}_{2n}(u)) > \frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2 \quad \text{for all } u > 0.$$

Hence, we have,

$$\|\mathbf{\Gamma}_{4n}(u)\|_s < \frac{1}{\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2} \quad \text{for all } n > n_4 \text{ and } u > 0$$

with probability exceeding  $1 - \varepsilon$ .

#### Bounding of $\|\mathbf{\Gamma}_{5n}(u)\|_s$

Recall that

$$\mathbf{\Gamma}_{5n}(u) \equiv \left\{ \frac{1}{n}\mathbf{L} + u^2\mathbf{I}_q - \mathbf{\Gamma}_{3n}(u) \right\}^{-1}.$$

Then,

$$\|\mathbf{\Gamma}_{5n}(u)\|_s = 1/\lambda_{\min}\left(\frac{1}{n}\mathbf{L} + u^2\mathbf{I}_q - \mathbf{\Gamma}_{3n}(u)\right).$$

Next,

$$\begin{aligned} \lambda_{\min}\left(\frac{1}{n}\mathbf{L} + u^2\mathbf{I}_q - \mathbf{\Gamma}_{3n}(u)\right) &\geq \lambda_{\min}\left(\frac{1}{n}\mathbf{L}\right) + u^2\lambda_{\min}(\mathbf{I}_q) - \lambda_{\min}(\mathbf{\Gamma}_{3n}(u)) \\ &\geq \frac{1}{n}\lambda_{\min}(\mathbf{L}) + u^2 - \lambda_{\max}(\mathbf{\Gamma}_{3n}(u)) \\ &= \frac{1}{n}\lambda_{\min}(\mathbf{L}) + u^2 - \|\mathbf{\Gamma}_{3n}(u)\|_s \\ &> \frac{1}{n}\lambda_{\min}(\mathbf{L}) + u^2 - \frac{pqR_n^2}{\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2} - q|T_n|. \end{aligned}$$

Since  $R_n = O_P(n^{-1})$  and  $T_n = o_P(n^{-1})$ , for every for every  $0 < \varepsilon \leq 1$  we can choose  $n_5 \in \mathbb{N}$  such that for all sufficiently large  $n$ ,

$$\frac{pqR_n^2}{\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2} + q|T_n| < \frac{1}{2n}\lambda_{\min}(\mathbf{L}).$$

This is equivalent to

$$-\frac{pqR_n^2}{\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2} - q|T_n| > -\frac{1}{2n}\lambda_{\min}(\mathbf{L}).$$

This allows us to claim that for all  $n > n_5$ ,

$$\lambda_{\min}\left(\frac{1}{n}\mathbf{L} + u^2\mathbf{I}_q - \mathbf{\Gamma}_{3n}(u)\right) > \frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2.$$

Therefore, we have,

$$\|\mathbf{\Gamma}_{5n}(u)\|_s < \frac{1}{\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2} \quad \text{for all } n > n_5 \text{ and } u > 0$$

with probability exceeding  $1 - \varepsilon$ .

Bounding of  $\|\mathbf{F}_{11n}(u; \mathbf{K}, \mathbf{L})\|_s$

Recall that

$$\mathbf{F}_{11n}(u; \mathbf{K}, \mathbf{L}) \equiv \mathbf{K}\mathbf{\Gamma}_{4n}(u)\mathbf{\Gamma}_{2n}(u) (\mathbf{K} + u^2\mathbf{I}_p)^{-1} - \mathbf{\Gamma}_{2n}(u)\mathbf{\Gamma}_{4n}(u).$$

Hence

$$\|\mathbf{F}_{11n}(u; \mathbf{K}, \mathbf{L})\|_s \leq \|\mathbf{K}\|_s \|\mathbf{\Gamma}_{4n}(u)\|_s \|\mathbf{\Gamma}_{2n}(u)\|_s \|(\mathbf{K} + u^2\mathbf{I}_p)^{-1}\|_s + \|\mathbf{\Gamma}_{2n}(u)\|_s \|\mathbf{\Gamma}_{4n}(u)\|_s.$$

Now,

$$\|\mathbf{K}\|_s = \lambda_{\max}(\mathbf{K})$$

and

$$\|(\mathbf{K} + u^2\mathbf{I}_p)^{-1}\|_s = \frac{1}{\lambda_{\min}(\mathbf{K} + u^2\mathbf{I}_p)} \leq \frac{1}{\lambda_{\min}(\mathbf{K}) + u^2}.$$

Hence, for all sufficiently large  $n$  and  $u > 0$ ,

$$\|\mathbf{F}_{11n}(u; \mathbf{K}, \mathbf{L})\|_s \leq \left\{ \frac{\lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K}) + u^2} + 1 \right\} \|\mathbf{\Gamma}_{2n}(u)\|_s \|\mathbf{\Gamma}_{4n}(u)\|_s.$$

Therefore, for all sufficiently large  $n$  and  $u > 0$ ,

$$\|\mathbf{F}_{11n}(u; \mathbf{K}, \mathbf{L})\|_s < p \left\{ \frac{2qnR_n^2}{\lambda_{\min}(\mathbf{L})} + |Q_n| \right\} \left\{ \frac{\lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K}) + u^2} + 1 \right\} \left\{ \frac{1}{\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2} \right\} \quad (2.56)$$

with probability exceeding  $1 - \varepsilon$ .

Bounding of  $\|\mathbf{F}_{21n}(u; \mathbf{K}, \mathbf{L})\|_s$

Recall that

$$\mathbf{F}_{21n}(u; \mathbf{K}, \mathbf{L}) \equiv n^{1/2}R_n u^2 \mathbf{\Gamma}_{1n}(u) \mathbf{1}_q \mathbf{1}_p^T \mathbf{\Gamma}_{4n}(u).$$

Then,

$$\begin{aligned} \|\mathbf{F}_{21n}(u; \mathbf{K}, \mathbf{L})\|_s &\leq n^{1/2}|R_n|u^2 \|\mathbf{\Gamma}_{1n}(u)\|_s \|\mathbf{1}_q\|_s \|\mathbf{1}_p^T\|_s \|\mathbf{\Gamma}_{4n}(u)\|_s \\ &= \sqrt{pq}n^{1/2}|R_n|u^2 \|\mathbf{\Gamma}_{1n}(u)\|_s \|\mathbf{\Gamma}_{4n}(u)\|_s. \end{aligned}$$

Hence, for all sufficiently large  $n$  we have,

$$\|\mathbf{F}_{21n}(u; \mathbf{K}, \mathbf{L})\|_s < \frac{\sqrt{pq}n^{1/2}|R_n|u^2}{\left\{ \frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2 \right\} \left\{ \frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2 \right\}} \quad \text{for all } u > 0 \quad (2.57)$$

with probability exceeding  $1 - \varepsilon$ .

Bounding of  $\|\mathbf{F}_{12n}(u; \mathbf{K}, \mathbf{L})\|_S$

Recall that

$$\mathbf{F}_{12n}(u; \mathbf{K}, \mathbf{L}) \equiv R_n \{ \mathbf{1}_p \mathbf{1}_q^T \mathbf{\Gamma}_{5n}(u) - (\mathbf{K} + Q_n \mathbf{1}_p^{\otimes 2}) \mathbf{\Gamma}_{4n}(u) \mathbf{1}_p \mathbf{1}_q^T \mathbf{\Gamma}_{1n}(u) \}.$$

Then

$$\begin{aligned} \|\mathbf{F}_{12n}(u; \mathbf{K}, \mathbf{L})\|_S &\leq |R_n| \left\{ \|\mathbf{1}_p\|_S \|\mathbf{1}_q^T\|_S \|\mathbf{\Gamma}_{5n}(u)\|_S \right. \\ &\quad \left. + (\|\mathbf{K}\|_S + |Q_n| \|\mathbf{1}_p^{\otimes 2}\|_S) \|\mathbf{\Gamma}_{4n}(u)\|_S \|\mathbf{1}_p\|_S \|\mathbf{1}_q^T\|_S \|\mathbf{\Gamma}_{1n}(u)\|_S \right\} \\ &= |R_n| \{ \sqrt{pq} \|\mathbf{\Gamma}_{5n}(u)\|_S + \sqrt{pq} (\lambda_{\max}(\mathbf{K}) + p|Q_n|) \|\mathbf{\Gamma}_{4n}(u)\|_S \|\mathbf{\Gamma}_{1n}(u)\|_S \} \\ &= \sqrt{pq} |R_n| \{ \|\mathbf{\Gamma}_{5n}(u)\|_S + (\lambda_{\max}(\mathbf{K}) + p|Q_n|) \|\mathbf{\Gamma}_{4n}(u)\|_S \|\mathbf{\Gamma}_{1n}(u)\|_S \}. \end{aligned}$$

For all sufficiently large  $n$  and  $u > 0$ , we have,

$$\begin{aligned} \|\mathbf{F}_{12n}(u; \mathbf{K}, \mathbf{L})\|_S &< \sqrt{pq} |R_n| \left\{ \frac{1}{\frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2} \right. \\ &\quad \left. + (\lambda_{\max}(\mathbf{K}) + p|Q_n|) \left( \frac{1}{\frac{1}{2} \lambda_{\min}(\mathbf{K}) + u^2} \right) \left( \frac{1}{\frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2} \right) \right\} \end{aligned} \quad (2.58)$$

with probability exceeding  $1 - \varepsilon$ .

Bounding of  $\|\mathbf{F}_{22n}(u; \mathbf{K}, \mathbf{L})\|_S$

Recall that

$$\begin{aligned} \mathbf{F}_{22n}(u; \mathbf{K}, \mathbf{L}) &\equiv n^{1/2} \left[ \left( \frac{1}{n} \mathbf{L} \right) \mathbf{\Gamma}_{5n}(u) \mathbf{\Gamma}_{3n}(u) \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q \right\}^{-1} \right. \\ &\quad \left. + T_n \mathbf{1}_q^{\otimes 2} \mathbf{\Gamma}_{5n}(u) - R_n^2 \mathbf{1}_q \mathbf{1}_p^T \mathbf{\Gamma}_{4n}(u) \mathbf{1}_p \mathbf{1}_q^T \mathbf{\Gamma}_{1n}(u) \right]. \end{aligned}$$

We then have

$$\begin{aligned} \|\mathbf{F}_{22n}(u; \mathbf{K}, \mathbf{L})\|_S &\leq n^{1/2} \left[ \left\| \frac{1}{n} \mathbf{L} \right\|_S \|\mathbf{\Gamma}_{5n}(u)\|_S \|\mathbf{\Gamma}_{3n}(u)\|_S \left\| \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q \right\}^{-1} \right\|_S \right. \\ &\quad \left. + |T_n| \|\mathbf{1}_q^{\otimes 2}\|_S \|\mathbf{\Gamma}_{5n}(u)\|_S + R_n^2 \|\mathbf{1}_q\|_S \|\mathbf{1}_p^T\|_S \|\mathbf{\Gamma}_{4n}(u)\|_S \|\mathbf{1}_p\|_S \|\mathbf{1}_q^T\|_S \|\mathbf{\Gamma}_{1n}(u)\|_S \right] \\ &= n^{1/2} \left[ \frac{1}{n} \lambda_{\max}(\mathbf{L}) \|\mathbf{\Gamma}_{5n}(u)\|_S \|\mathbf{\Gamma}_{3n}(u)\|_S \left\| \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q \right\}^{-1} \right\|_S \right. \\ &\quad \left. + q |T_n| \|\mathbf{\Gamma}_{5n}(u)\|_S + pq R_n^2 \|\mathbf{\Gamma}_{4n}(u)\|_S \|\mathbf{\Gamma}_{1n}(u)\|_S \right]. \end{aligned} \quad (2.59)$$

Next, note that

$$\left\| \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q \right\}^{-1} \right\|_s \leq 1 / \lambda_{\min} \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q \right).$$

Since

$$\lambda_{\min} \left( \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q \right) \geq \frac{1}{n} \lambda_{\min}(\mathbf{L}) + u^2$$

we have the bound

$$\left\| \left\{ \frac{1}{n} \mathbf{L} + u^2 \mathbf{I}_q \right\}^{-1} \right\|_s \leq \frac{1}{\frac{1}{n} \lambda_{\min}(\mathbf{L}) + u^2} \quad \text{for all } u > 0.$$

It follows that for all sufficiently large  $n$ , the first term on the right-hand side of (2.59) is bounded above by

$$\begin{aligned} & \frac{n^{-1/2} \lambda_{\max}(\mathbf{L})}{\frac{1}{n} \lambda_{\min}(\mathbf{L}) + u^2} \|\mathbf{\Gamma}_{5n}(u)\|_s \|\mathbf{\Gamma}_{3n}(u)\|_s \\ & < \left\{ \frac{n^{-1/2} \lambda_{\max}(\mathbf{L})}{\frac{1}{n} \lambda_{\min}(\mathbf{L}) + u^2} \right\} \left\{ \frac{1}{\frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2} \right\} \left\{ \frac{pqR_n^2}{\frac{1}{2} \lambda_{\min}(\mathbf{K}) + u^2} + q|T_n| \right\}. \end{aligned}$$

It follows that for all sufficiently large  $n$ , the second term on the right-hand side of (2.59) is bounded above by

$$\frac{qn^{1/2}|T_n|}{\frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2}.$$

Finally, it follows that for all sufficiently large  $n$ , the third term on the right-hand side of (2.59) is bounded above by

$$\frac{pqn^{1/2}R_n^2}{\left(\frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2\right) \left(\frac{1}{2} \lambda_{\min}(\mathbf{K}) + u^2\right)}.$$

Combining all of these bounds, we have

$$\begin{aligned} \|\mathbf{F}_{22n}(u; \mathbf{K}, \mathbf{L})\|_s & < \left\{ \frac{n^{-1/2} \lambda_{\max}(\mathbf{L})}{\frac{1}{n} \lambda_{\min}(\mathbf{L}) + u^2} \right\} \left\{ \frac{1}{\frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2} \right\} \left\{ \frac{pqR_n^2}{\frac{1}{2} \lambda_{\min}(\mathbf{K}) + u^2} + q|T_n| \right\} \\ & + \frac{qn^{1/2}|T_n|}{\frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2} + \frac{pqn^{1/2}R_n^2}{\left(\frac{1}{2n} \lambda_{\min}(\mathbf{L}) + u^2\right) \left(\frac{1}{2} \lambda_{\min}(\mathbf{K}) + u^2\right)} \end{aligned} \quad (2.60)$$

for all sufficiently large  $n$  and  $u > 0$ , with probability exceeding  $1 - \varepsilon$ .

### 2.4.4.3 Verifying Convergence in Probability Limits of the Functions in (2.55)

The Convergence in Probability Limit of  $\int_0^\infty \mathbf{F}_{11n}(u; \boldsymbol{\kappa}, \mathbf{K}) du$

Noting that

$$\left\| \int_0^\infty \mathbf{F}_{11n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \right\|_s \leq \int_0^\infty \|\mathbf{F}_{11n}(u; \boldsymbol{\kappa}, \mathbf{K})\|_s du$$

and from (2.56) we have, for all sufficiently large  $n$

$$\begin{aligned} & \int_0^\infty \|\mathbf{F}_{11n}(u; \boldsymbol{\kappa}, \mathbf{K})\|_s du \\ & < p \left\{ \frac{2qnR_n^2}{\lambda_{\min}(\mathbf{L})} + |Q_n| \right\} \int_0^\infty \left\{ \frac{\lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K}) + u^2} + 1 \right\} \left\{ \frac{1}{\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2} \right\} du \\ & < p \left\{ \frac{2qnR_n^2}{\lambda_{\min}(\mathbf{L})} + |Q_n| \right\} \left\{ \int_0^\infty \frac{\lambda_{\max}(\mathbf{K})}{(\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2)^2} du + \int_0^\infty \frac{1}{\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2} du \right\} \\ & = p \left\{ \frac{2qnR_n^2}{\lambda_{\min}(\mathbf{L})} + |Q_n| \right\} \left( \frac{\pi \lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K}) \sqrt{2\lambda_{\min}(\mathbf{K})}} + \frac{\pi}{\sqrt{2\lambda_{\min}(\mathbf{K})}} \right) \end{aligned}$$

with probability exceeding  $1 - \varepsilon$ . Since  $R_n = O_P(n^{-1})$ ,  $Q_n = o_P(1)$  and  $\varepsilon$  is arbitrary, we must have

$$\left\| \int_0^\infty \mathbf{F}_{11n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \right\|_s \xrightarrow{P} 0 \quad \text{as } n \rightarrow \infty$$

and therefore,

$$\int_0^\infty \mathbf{F}_{11n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \xrightarrow{P} \mathbf{O} \quad \text{as } n \rightarrow \infty.$$

The Convergence in Probability Limit of  $\int_0^\infty \mathbf{F}_{21n}(u; \boldsymbol{\kappa}, \mathbf{K}) du$

Noting that

$$\left\| \int_0^\infty \mathbf{F}_{21n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \right\|_s \leq \int_0^\infty \|\mathbf{F}_{21n}(u; \boldsymbol{\kappa}, \mathbf{K})\|_s du$$

and from (2.57) we have, for all sufficiently large  $n$

$$\begin{aligned} & \int_0^\infty \|\mathbf{F}_{21n}(u; \boldsymbol{\kappa}, \mathbf{K})\|_s du \\ & < \sqrt{pq}n^{1/2}|R_n| \int_0^\infty \frac{u^2}{\left\{ \frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2 \right\} \left\{ \frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2 \right\}} du \\ & = \frac{\sqrt{pq}\pi n^{1/2}|R_n|}{\sqrt{2} \left( \sqrt{\lambda_{\min}(\mathbf{K})} + \sqrt{\lambda_{\min}(\mathbf{L})/n} \right)} \end{aligned}$$

with probability exceeding  $1 - \varepsilon$ . Since  $R_n = O_P(n^{-1})$  and  $\varepsilon$  is arbitrary, we must have

$$\left\| \int_0^\infty \mathbf{F}_{21n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \right\|_s \xrightarrow{P} 0 \quad \text{as } n \rightarrow \infty$$

and therefore,

$$\int_0^\infty \mathbf{F}_{21n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \xrightarrow{P} \mathbf{O} \quad \text{as } n \rightarrow \infty.$$

The Convergence in Probability Limit of  $\int_0^\infty \mathbf{F}_{12n}(u; \boldsymbol{\kappa}, \mathbf{K}) du$

Noting that

$$\left\| \int_0^\infty \mathbf{F}_{12n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \right\|_s \leq \int_0^\infty \|\mathbf{F}_{12n}(u; \boldsymbol{\kappa}, \mathbf{K})\|_s du$$

and from (2.58) we have, for all sufficiently large  $n$

$$\begin{aligned} & \int_0^\infty \|\mathbf{F}_{12n}(u; \boldsymbol{\kappa}, \mathbf{K})\|_s du \\ & < \sqrt{pq}|R_n| \left\{ \int_0^\infty \frac{1}{\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2} du \right. \\ & \quad \left. + (\lambda_{\max}(\mathbf{K}) + p|Q_n|) \int_0^\infty \left( \frac{1}{\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2} \right) \left( \frac{1}{\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2} \right) du \right\} \\ & = \sqrt{pq}|R_n| \left[ \frac{\pi}{2\sqrt{\frac{1}{2n}\lambda_{\min}(\mathbf{L})}} \right. \\ & \quad \left. + (\lambda_{\max}(\mathbf{K}) + p|Q_n|) \left\{ \frac{\pi}{2\sqrt{\left(\frac{1}{2}\lambda_{\min}(\mathbf{K})\right)\left(\frac{1}{2n}\lambda_{\min}(\mathbf{L})\right)\left(\sqrt{\frac{1}{2}\lambda_{\min}(\mathbf{K})} + \sqrt{\frac{1}{2n}\lambda_{\min}(\mathbf{L})}\right)}} \right\} \right] \\ & = \sqrt{pq}|R_n| \left[ \frac{n^{1/2}\pi}{\sqrt{2\lambda_{\min}(\mathbf{L})}} \right. \\ & \quad \left. + (\lambda_{\max}(\mathbf{K}) + p|Q_n|) \left\{ \frac{\sqrt{2}n^{1/2}\pi}{\sqrt{\lambda_{\min}(\mathbf{K})\lambda_{\min}(\mathbf{L})}\left(\sqrt{\lambda_{\min}(\mathbf{K})} + \sqrt{\frac{1}{n}\lambda_{\min}(\mathbf{L})}\right)} \right\} \right] \end{aligned}$$

with probability exceeding  $1 - \varepsilon$ . Since  $R_n = O_P(n^{-1})$ ,  $Q_n = o_P(1)$  and  $\varepsilon$  is arbitrary, we must have

$$\left\| \int_0^\infty \mathbf{F}_{12n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \right\|_s \xrightarrow{P} 0 \quad \text{as } n \rightarrow \infty$$

and therefore,

$$\int_0^\infty \mathbf{F}_{12n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \xrightarrow{P} \mathbf{O} \quad \text{as } n \rightarrow \infty.$$



The Convergence in Probability Limit of  $\int_0^\infty \mathbf{F}_{22n}(u; \boldsymbol{\kappa}, \mathbf{K}) du$

Noting that

$$\left\| \int_0^\infty \mathbf{F}_{22n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \right\|_s \leq \int_0^\infty \|\mathbf{F}_{22n}(u; \boldsymbol{\kappa}, \mathbf{K})\|_s du$$

and from (2.60) we have, for all sufficiently large  $n$

$$\begin{aligned} & \int_0^\infty \|\mathbf{F}_{22n}(u; \boldsymbol{\kappa}, \mathbf{K})\|_s du \\ & < \lambda_{\max}(\mathbf{L})pqn^{-1/2}R_n^2 \int_0^\infty \frac{1}{\left(\frac{1}{n}\lambda_{\min}(\mathbf{L}) + u^2\right) \left(\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2\right) \left(\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2\right)} du \\ & \quad + n^{-1/2}\lambda_{\max}(\mathbf{L})q|T_n| \int_0^\infty \frac{1}{\left(\frac{1}{n}\lambda_{\min}(\mathbf{L}) + u^2\right) \left(\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2\right)} du \\ & \quad + qn^{1/2}|T_n| \int_0^\infty \frac{1}{\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2} du \\ & \quad + pqn^{1/2}R_n^2 \int_0^\infty \frac{1}{\left(\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2\right) \left(\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2\right)} du \\ & < \lambda_{\max}(\mathbf{L})pqn^{-1/2}R_n^2 \int_0^\infty \frac{1}{\left(\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2\right)^2 \left(\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2\right)} du \\ & \quad + n^{-1/2}\lambda_{\max}(\mathbf{L})q|T_n| \int_0^\infty \frac{1}{\left(\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2\right)^2} du \\ & \quad + qn^{1/2}|T_n| \int_0^\infty \frac{1}{\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2} du \\ & \quad + pqn^{1/2}R_n^2 \int_0^\infty \frac{1}{\left(\frac{1}{2n}\lambda_{\min}(\mathbf{L}) + u^2\right) \left(\frac{1}{2}\lambda_{\min}(\mathbf{K}) + u^2\right)} du \\ & < \frac{2\sqrt{2}pq\pi n\lambda_{\max}(\mathbf{L})R_n^2}{\lambda_{\min}(\mathbf{L})\sqrt{\lambda_{\min}(\mathbf{L})\lambda_{\min}(\mathbf{K})} \left(\sqrt{\frac{1}{2}\lambda_{\min}(\mathbf{L})} + \sqrt{\lambda_{\min}(\mathbf{K})}\right)} \\ & \quad + \frac{q\pi n\lambda_{\max}(\mathbf{K})|T_n|}{\sqrt{2}\lambda_{\min}(\mathbf{L})\sqrt{\lambda_{\min}(\mathbf{L})}} + \frac{q\pi n|T_n|}{\sqrt{2}\lambda_{\min}(\mathbf{L})} \\ & \quad + \frac{pq\pi n\sqrt{2}R_n^2}{\sqrt{\lambda_{\min}(\mathbf{L})\lambda_{\min}(\mathbf{K})} \left(\sqrt{\frac{1}{n}\lambda_{\min}(\mathbf{L})} + \sqrt{\lambda_{\min}(\mathbf{K})}\right)} \end{aligned}$$

with probability exceeding  $1 - \varepsilon$ . Since  $R_n = O_P(n^{-1})$ ,  $T_n = o_P(n^{-1})$  and  $\varepsilon$  is arbitrary, it follows that

$$\left\| \int_0^\infty \mathbf{F}_{22n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \right\|_s \xrightarrow{P} 0 \quad \text{as } n \rightarrow \infty$$

and therefore,

$$\int_0^\infty \mathbf{F}_{22n}(u; \boldsymbol{\kappa}, \mathbf{K}) du \xrightarrow{P} \mathbf{O} \quad \text{as } n \rightarrow \infty.$$

#### 2.4.4.4 Conclusion for Multivariate Integral Limits for the Matrix Square Root Result

Hence, we have shown that

$$\text{plim}_{n \rightarrow \infty} \int_0^\infty \mathbf{F}_{kk'n}(u; \mathbf{K}, \mathbf{L}) du = \mathbf{O}, \quad k, k' = 1, 2$$

for  $u > 0$ .

## Chapter 3

# Usable Asymptotic Normality Results and Inference for Gaussian Response Linear Mixed Models

Even though estimation by maximum likelihood for linear mixed models and generalized linear mixed models is well established, asymptotic normality results that can be used to construct confidence intervals and Wald tests via Studentization are currently unavailable in the existing linear mixed model and generalized linear mixed model literature.

Asymptotic normality results for maximum likelihood estimators for Gaussian response linear mixed models have been presented in literature such as McCulloch et al. (2008), Miller (1973), Miller (1977), Jiang (1996) and Jiang and Nguyen (2021). As discussed in Section 1.9.3, Nie (2007) investigated an extension to a generalized linear mixed model setting, but did not give explicit forms void of limits or expectations with respect to the response. Hence, the existing *generalized* linear mixed model literature lacks asymptotic covariance results that are amenable to practical purposes such as confidence interval construction.

Other recent related literature by Lyu and Welsh (2022) derive explicit asymptotic normality results for both maximum likelihood estimators and restricted maximum likelihood estimators for model parameters in a nested regression model (random intercept model) for clustered data. The authors, as done in this thesis, considered the scenario where both the number of independent clusters and number of observations within each cluster go to infinity. When restricted to Gaussian responses, the work by

Lyu and Welsh (2022) is closely related to the work in this chapter. However, while Lyu and Welsh (2022) consider regression models with a random intercept, this chapter considers Gaussian response linear mixed models that allow for both a random intercept and slope to be included. On the other hand, Westfall (1986) considers a linear mixed model set-up with vectors of fixed effects, nonerror random effects and error random effects and develop asymptotic distribution theory for the corresponding analysis of variance estimates. In contrast, this chapter presents the explicit first order (leading term) asymptotic approximations used for the asymptotic normality results in this chapter, serving as a prequel to the results presented for the generalized linear mixed model case in Chapter 4.

In this chapter, we aim to derive asymptotic normality results that are directly usable for asymptotically valid confidence intervals and Wald tests for analysis concerning linear mixed models with Gaussian responses. The main theorem in this chapter concerns the joint asymptotic normality of all maximum likelihood estimators for a Gaussian response mixed model and elegantly shows faster rates of convergence of fixed effects that are not accompanied by a random effect as compared to fixed effects that have partnering random effects.

The results presented in this chapter will then be extended for generalized linear mixed models with multivariate fixed and random effects in Chapter 4.

### 3.1 Model Description

In this section, we study Gaussian response linear mixed models of the following form, for observations of the random triples  $(\mathbf{X}_{Aij}, \mathbf{X}_{Bij}, Y_{ij})$ ,  $1 \leq i \leq m, 1 \leq j \leq n_i$ ,

$$Y_{ij} | \mathbf{X}_{Aij}, \mathbf{X}_{Bij}, \mathbf{U}_i \text{ are independent } N((\boldsymbol{\beta}_A^0 + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B^0)^T \mathbf{X}_{Bij}, \sigma_\varepsilon^2). \quad (3.1)$$

The  $\mathbf{U}_i$  are  $d_A \times 1$  unobserved random vectors for each  $1 \leq i \leq m$ . The  $\mathbf{X}_{Aij}$  are  $d_A \times 1$  random vectors corresponding to predictors that are partnered by both a fixed effect and a random effect. The  $\mathbf{X}_{Bij}$  are  $d_B \times 1$  random vectors are predictors that have a fixed effect only. Let  $\mathbf{X}_{ij} \equiv (\mathbf{X}_{Aij}^T, \mathbf{X}_{Bij}^T)^T$  denote the combined predictor vectors such that  $d_A + d_B = d$ . We also assumed that the  $\mathbf{X}_{ij}$  and  $\mathbf{U}_i$ , for  $1 \leq i \leq m$  and  $1 \leq j \leq n_i$  are independent, with the  $\mathbf{X}_{ij}$  each having the same distribution as the  $(d_A + d_B) \times 1$  random vector  $\mathbf{X} = (\mathbf{X}_A^T, \mathbf{X}_B^T)^T$  and the  $\mathbf{U}_i$  are independent  $N(0, (\boldsymbol{\Sigma})^0)$ , each having the same distribution as the random vector  $\mathbf{U}$ .

Then, for any  $\beta_A(d_A \times 1)$ ,  $\beta_B(d_B \times 1)$ ,  $\sigma_\varepsilon^2$  and  $\Sigma(d_A \times d_A)$ , the maximum likelihood estimator of  $(\beta_A^0, \beta_B^0, \Sigma^0, (\sigma^2)^0)$  is,

$$(\hat{\beta}_A, \hat{\beta}_B, \hat{\Sigma}, \hat{\sigma}^2) = \underset{\beta_A, \beta_B, \Sigma, \sigma_\varepsilon^2}{\operatorname{argmax}} \ell(\beta_A, \beta_B, \Sigma, \sigma_\varepsilon^2)$$

where the conditional log-likelihood is

$$\begin{aligned} & \ell(\beta_A, \beta_B, \Sigma, \sigma_\varepsilon^2) \\ &= \sum_{i=1}^m \sum_{j=1}^n \left[ \left\{ Y_{ij} (\beta_A^T \mathbf{X}_{Aij} + \beta_B^T \mathbf{X}_{Bij}) - \frac{1}{2} Y_{ij}^2 \right\} / \sigma_\varepsilon^2 - \frac{1}{2} \log(2\pi\sigma_\varepsilon^2) \right] - \frac{m}{2} \log |2\pi\Sigma| \\ &+ \sum_{i=1}^m \log \int_{\mathbb{R}^{d_A}} \exp \left( \sum_{j=1}^n \left[ Y_{ij} \mathbf{u}^T \mathbf{X}_{Aij} - \frac{1}{2} \{ (\beta_A + \mathbf{u})^T \mathbf{X}_{Aij} + \beta_B^T \mathbf{X}_{Bij} \} \right] / \sigma_\varepsilon^2 - \frac{1}{2} \mathbf{u}^T \Sigma^{-1} \mathbf{u} \right) d\mathbf{u}. \end{aligned}$$

## 3.2 Notation Required for Fisher Information Calculations

Let

$$\mathbf{X}_i = \begin{bmatrix} \mathbf{X}_{i1}^T \\ \vdots \\ \mathbf{X}_{in_i}^T \end{bmatrix} \quad \text{and} \quad \mathbf{Y}_i = \begin{bmatrix} Y_{i1} \\ \vdots \\ Y_{in_i} \end{bmatrix}.$$

Also define

$$n \equiv \frac{1}{m} \sum_{i=1}^m n_i = \text{average of the within-group sample sizes,}$$

and

$$\Sigma_{\beta_B} = \text{lower right } d_B \times d_B \text{ block of } E \left( \begin{bmatrix} \mathbf{X}_A \mathbf{X}_A^T & \mathbf{X}_A \mathbf{X}_B^T \\ \mathbf{X}_B \mathbf{X}_A^T & \mathbf{X}_B \mathbf{X}_B^T \end{bmatrix} \right)^{-1}.$$

## 3.3 Asymptotic Normality Theorem

The main theoretical contribution of this chapter is an asymptotic normality theorem for the maximum likelihood estimators for a Gaussian response mixed model as described in Section 3.1.

The theorem relies on the following assumptions:

- (A1) The number of groups  $m$  diverges to  $\infty$ .
- (A2) The within-group sample sizes  $n_i$  diverge to  $\infty$  in such a way that  $n_i/n \rightarrow C_i$  for constants  $0 < C_i < \infty$ ,  $1 \leq i \leq m$ .
- (A3) The distribution of  $\mathbf{X}$  is such that

$$E(\|\mathbf{X}\|^8) < \infty$$

and none of the entries in  $\mathbf{X}_A$  are zero degenerate random variables.

**Theorem 11.** *Assume that conditions (A1) - (A3) hold. Then we have the following*

$$\sqrt{m} \begin{bmatrix} \hat{\beta}_A - \beta_A^0 \\ \sqrt{n} (\hat{\beta}_B - \beta_B^0) \\ \text{vech}(\hat{\Sigma} - \Sigma^0) \\ \sqrt{n} (\hat{\sigma}^2 - (\sigma^2)^0) \end{bmatrix} \xrightarrow{\mathcal{D}} N \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \Sigma^0 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\sigma^2)^0 \Sigma_{\beta_B} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & 2\mathbf{D}_{d_A}^+ (\Sigma^0 \otimes \Sigma^0) \mathbf{D}_{d_A}^{+T} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & 2\{(\sigma^2)^0\}^2 \end{bmatrix} \right),$$

where  $\mathbf{D}_{d_A}^+$  is the Moore-Penrose inverse of  $\mathbf{D}_{d_A}$ .

The proof for Theorem 11 is in the appendix.

## 3.4 Appendix

### 3.4.1 Linear Mixed Models with Multivariate Fixed and Random Effects

Let

$$\mathbf{X}_i = \begin{bmatrix} \mathbf{X}_{Ai} & \mathbf{X}_{Bi} \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_A \\ \beta_B \end{bmatrix} \quad \text{and} \quad \theta = \begin{bmatrix} \text{vec}(\Sigma) \\ \sigma_\varepsilon^2 \end{bmatrix}.$$

Then the model presented in Section 3.1 can be rewritten in the following form

$$\mathbf{y}_i \stackrel{\text{ind.}}{\sim} N(\mathbf{X}_i \beta, \mathbf{V}_{\theta i}), \quad (3.2)$$

with

$$\mathbf{V}_{\theta i} = \mathbf{X}_{Ai} \Sigma \mathbf{X}_{Ai}^T + \sigma_\varepsilon^2 I_{n_i},$$

where  $\mathbf{V}_{\theta i}$  is the  $n_i \times n_i$  covariance matrix of  $\mathbf{y}_i$  parametrized by  $\theta$ . The log-likelihood of  $(\beta, \theta)$  can then be expressed as

$$\ell(\beta, \theta) = -\frac{1}{2} \sum_{i=1}^m \left\{ \log |\mathbf{V}_{\theta i}| + (\mathbf{y}_i - \mathbf{X}_i \beta)^T \mathbf{V}_{\theta i}^{-1} (\mathbf{y}_i - \mathbf{X}_i \beta) + n_i \log(2\pi) \right\}.$$

The model in (3.2) is a special case of a general Gaussian variance regression model presented in Section 4.3 of Wand (2002). Using steps similar to those in Wand (2002), an expression for the Fisher information matrix for a Gaussian response linear mixed model with multivariate fixed and random effects can be obtained. Note that the model in (3.2) belongs to a common class of submodels with

$$\mathbf{V}_{\boldsymbol{\theta}_i} = \sum_{h=1}^c \theta_h \mathbf{K}_h, \quad \boldsymbol{\theta} = [\theta_1, \dots, \theta_c], \quad (3.3)$$

for a set of  $n_i \times n_i$  matrices  $\mathbf{K}_1, \dots, \mathbf{K}_c$ , which leads to considerable simplifications in obtaining the expression for the Fisher information matrix. An alternative expression for equation (3.3) is as follows, where

$$\text{vec}(\mathbf{V}_{\boldsymbol{\theta}_i}) = \boldsymbol{\kappa}_i \boldsymbol{\theta}, \quad \boldsymbol{\kappa} = [\text{vec}(\mathbf{K}_1) | \dots | \text{vec}(\mathbf{K}_c)].$$

It remains to solve for  $\boldsymbol{\kappa}_i$ . By making use of the fourth property in Subsection 1.4.5, it can be done as follows

$$\begin{aligned} \text{vec}(\mathbf{V}_{\boldsymbol{\theta}_i}) &= \text{vec}(\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_n) \\ &= \text{vec}(\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T) + \text{vec}(\sigma_\varepsilon^2 \mathbf{I}_n) \\ &= (\mathbf{X}_{A_i} \otimes \mathbf{X}_{A_i}) \text{vec}(\boldsymbol{\Sigma}) + \text{vec}(\mathbf{I}_{n_i}) \sigma_\varepsilon^2 \\ &= (\mathbf{X}_{A_i} \otimes \mathbf{X}_{A_i}) \mathbf{D}_{d_A} \text{vech}(\boldsymbol{\Sigma}) + \text{vec}(\mathbf{I}_{n_i}) \sigma_\varepsilon^2 \\ &= \left[ (\mathbf{X}_{A_i} \otimes \mathbf{X}_{A_i}) \mathbf{D}_{d_A} \quad \text{vec}(\mathbf{I}_{n_i}) \right] \begin{bmatrix} \text{vech}(\boldsymbol{\Sigma}) \\ \sigma_\varepsilon^2 \end{bmatrix} \\ &= \boldsymbol{\kappa}_i \boldsymbol{\theta}, \end{aligned}$$

with

$$\boldsymbol{\kappa}_i = \left[ (\mathbf{X}_{A_i} \otimes \mathbf{X}_{A_i}) \mathbf{D}_{d_A} \quad \text{vec}(\mathbf{I}_{n_i}) \right].$$

The expression for the full Fisher information matrix is subsequently

$$\begin{aligned} I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2) &= \begin{bmatrix} I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B) & 0 \\ 0 & I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2) \end{bmatrix} \\ &= \sum_{i=1}^m \begin{bmatrix} \mathbf{X}_i^T \mathbf{V}_{\boldsymbol{\theta}_i}^{-1} \mathbf{X}_i & 0 \\ 0 & \frac{1}{2} \boldsymbol{\kappa}_i^T (\mathbf{V}_{\boldsymbol{\theta}_i}^{-1} \otimes \mathbf{V}_{\boldsymbol{\theta}_i}^{-1}) \boldsymbol{\kappa}_i \end{bmatrix} \\ &= \begin{bmatrix} \sum_{i=1}^m \mathbf{X}_i^T \mathbf{V}_{\boldsymbol{\theta}_i}^{-1} \mathbf{X}_i & 0 \\ 0 & \frac{1}{2} \sum_{i=1}^m \boldsymbol{\kappa}_i^T (\mathbf{V}_{\boldsymbol{\theta}_i}^{-1} \otimes \mathbf{V}_{\boldsymbol{\theta}_i}^{-1}) \boldsymbol{\kappa}_i \end{bmatrix}. \end{aligned} \quad (3.4)$$

To find an explicit expression for  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)$ , it suffices to find explicit expressions for the block matrices  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$  and  $I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)$  separately.

### 3.4.2 Expression for Top Left Block of Fisher Information Matrix

Expanding the current expression for  $I(\beta_A, \beta_B)$  leads to

$$\begin{aligned}
I(\beta_A, \beta_B) &= \sum_{i=1}^m \mathbf{X}_i^T \mathbf{V}_{\theta_i}^{-1} \mathbf{X}_i \\
&= \sum_{i=1}^m \begin{bmatrix} \mathbf{X}_{A_i} & \mathbf{X}_{B_i} \end{bmatrix}^T (\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \begin{bmatrix} \mathbf{X}_{A_i} & \mathbf{X}_{B_i} \end{bmatrix} \\
&= \sum_{i=1}^m \begin{bmatrix} \mathbf{X}_{A_i}^T \\ \mathbf{X}_{B_i}^T \end{bmatrix} (\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \begin{bmatrix} \mathbf{X}_{A_i} & \mathbf{X}_{B_i} \end{bmatrix} \\
&= \sum_{i=1}^m \begin{bmatrix} \mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{A_i} & \mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{B_i} \\ \mathbf{X}_{B_i}^T (\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{A_i} & \mathbf{X}_{B_i}^T (\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{B_i} \end{bmatrix}.
\end{aligned} \tag{3.5}$$

It remains to evaluate and simplify the expressions in the matrix, which can be done by making use of the matrix identities from Harville (1977), highlighted in Subsubsection 1.4.7.1.

#### 3.4.2.1 Top Left Block of (3.5)

Now, we solve a part of the expression of the top left block of  $I(\beta_A, \beta_B)$ , specifically

$$\mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{A_i}.$$

Using (1.10b) and (1.3) by setting  $\mathbf{A} = \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \mathbf{X}_{A_i}$  and  $\mathbf{B} = \boldsymbol{\Sigma}$ , we have,

$$\begin{aligned}
\mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{A_i} &= \{ \mathbf{I}_{d_A} + \mathbf{X}_{A_i}^T (\sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{A_i} \boldsymbol{\Sigma} \}^{-1} \mathbf{X}_{A_i}^T (\sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{A_i} \\
&= \left( \mathbf{I}_{d_A} + \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \mathbf{X}_{A_i} \boldsymbol{\Sigma} \right)^{-1} \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \mathbf{X}_{A_i} \right) \\
&= \left\{ \mathbf{I}_{d_A} - \boldsymbol{\Sigma}^{-1} \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \mathbf{X}_{A_i} \right)^{-1} + \dots \right\} \boldsymbol{\Sigma}^{-1} \\
&= \left\{ \mathbf{I}_{d_A} - \sigma_\varepsilon^2 \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} + \dots \right\} \boldsymbol{\Sigma}^{-1} \\
&= \boldsymbol{\Sigma}^{-1} - \sigma_\varepsilon^2 \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \boldsymbol{\Sigma}^{-1} + \dots \\
&= \boldsymbol{\Sigma}^{-1} + O_P(n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T.
\end{aligned}$$

Therefore, the top left block of  $I(\beta_A, \beta_B)$  can be computed as

$$\sum_{i=1}^m (\boldsymbol{\Sigma}^{-1} + O_P(n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T) = m \boldsymbol{\Sigma}^{-1} + O_P(mn^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T.$$

#### 3.4.2.2 Top Right Block of (3.5)

Now, we solve a part of the expression of the top right block of  $I(\beta_A, \beta_B)$ , specifically

$$\mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{B_i}.$$



Using (1.10b) and (1.3) by setting  $\mathbf{A} = \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{Ai}^T \mathbf{X}_{Ai}$ ,  $\mathbf{B} = \boldsymbol{\Sigma}$  and  $\mathbf{C} = \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{Ai}^T \mathbf{X}_{Bi}$  we have,

$$\begin{aligned} \mathbf{X}_{Ai}^T (\mathbf{X}_{Ai} \boldsymbol{\Sigma} \mathbf{X}_{Ai}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{Bi} &= \{\mathbf{I}_{d_A} + \mathbf{X}_{Ai}^T (\sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{Ai} \boldsymbol{\Sigma}\}^{-1} \mathbf{X}_{Ai}^T (\sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{Bi} \\ &= \left( \mathbf{I}_{d_A} + \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{Ai}^T \mathbf{X}_{Ai} \boldsymbol{\Sigma} \right)^{-1} \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{Ai}^T \mathbf{X}_{Bi} \right) \\ &= \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Bi}) \\ &\quad - \sigma_\varepsilon^2 \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Bi}) + \dots \\ &= \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Bi}) + O_P(n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T. \end{aligned}$$

Hence, the top right block of  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$  can be computed as

$$\begin{aligned} &\sum_{i=1}^m \left\{ \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Bi}) + O_P(n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T \right\} \\ &= \sum_{i=1}^m \left\{ \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Bi}) \right\} + O_P(mn^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T. \\ &= O_P(m) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T. \end{aligned}$$

### 3.4.2.3 Bottom Left Block of (3.5)

Subsequently, the expression for the bottom left block of  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$  can be easily computed as

$$\begin{aligned} &\left[ \sum_{i=1}^m \left\{ \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} (\mathbf{X}_{Ai}^T \mathbf{X}_{Bi}) \right\} + O_P(mn^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T \right]^T \\ &= \sum_{i=1}^m \left\{ (\mathbf{X}_{Bi}^T \mathbf{X}_{Ai}) (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \boldsymbol{\Sigma}^{-1} \right\} + O_P(mn^{-1}) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T \\ &= O_P(m) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T. \end{aligned}$$

### 3.4.2.4 Bottom Right Block of (3.5)

Now, we solve a part of the expression of the bottom right block of  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$ , specifically

$$\mathbf{X}_{Bi}^T (\mathbf{X}_{Ai} \boldsymbol{\Sigma} \mathbf{X}_{Ai}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{Bi}.$$

Using (1.10a) and (1.5) by setting  $\mathbf{A} = \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \mathbf{X}_{A_i}$  and  $\mathbf{B} = \boldsymbol{\Sigma}$ , we have,

$$\begin{aligned}
& \mathbf{X}_{B_i}^T \{ \mathbf{X}_{A_i} \boldsymbol{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i} \}^{-1} \mathbf{X}_{B_i} \\
&= \mathbf{X}_{B_i}^T \left[ (\sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} - (\sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{A_i} \boldsymbol{\Sigma} \{ \mathbf{I}_{d_A} + \mathbf{X}_{A_i}^T (\sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{A_i} \boldsymbol{\Sigma} \}^{-1} \mathbf{X}_{A_i}^T (\sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \right] \mathbf{X}_{B_i} \\
&= \mathbf{X}_{B_i}^T \left\{ \frac{1}{\sigma_\varepsilon^2} \mathbf{I}_{n_i} - \left( \frac{1}{\sigma_\varepsilon^2} \right)^2 \mathbf{X}_{A_i} \boldsymbol{\Sigma} \left( \mathbf{I}_{d_A} + \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \mathbf{X}_{A_i} \boldsymbol{\Sigma} \right)^{-1} \mathbf{X}_{A_i}^T \right\} \mathbf{X}_{B_i} \\
&= \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{B_i}^T \mathbf{X}_{B_i} - \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{B_i}^T \mathbf{X}_{A_i} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{X}_{A_i}^T \mathbf{X}_{B_i} \\
&\quad + \mathbf{X}_{B_i}^T \mathbf{X}_{A_i} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \boldsymbol{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{X}_{A_i}^T \mathbf{X}_{B_i} + \dots \\
&= \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{B_i}^T \mathbf{X}_{B_i} - \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{B_i}^T \mathbf{X}_{A_i} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{X}_{A_i}^T \mathbf{X}_{B_i} + O_P(1) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T.
\end{aligned}$$

Hence, the bottom right block of  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$  can be computed as

$$\begin{aligned}
& \sum_{i=1}^m \left\{ \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{B_i}^T \mathbf{X}_{B_i} - \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{B_i}^T \mathbf{X}_{A_i} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{X}_{A_i}^T \mathbf{X}_{B_i} + O_P(1) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T \right\} \\
&= \frac{1}{\sigma_\varepsilon^2} \sum_{i=1}^m \mathbf{X}_{B_i}^T \mathbf{X}_{B_i} - \frac{1}{\sigma_\varepsilon^2} \sum_{i=1}^m \mathbf{X}_{B_i}^T \mathbf{X}_{A_i} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{X}_{A_i}^T \mathbf{X}_{B_i} + O_P(m) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T \\
&= \frac{1}{\sigma_\varepsilon^2} \sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{B_{ij}} \mathbf{X}_{B_{ij}}^T - \frac{1}{\sigma_\varepsilon^2} \sum_{i=1}^m \left( \sum_{j=1}^n \mathbf{X}_{B_{ij}} \mathbf{X}_{A_{ij}}^T \right) \left( \sum_{j=1}^n \mathbf{X}_{A_{ij}} \mathbf{X}_{A_{ij}}^T \right)^{-1} \left( \sum_{j=1}^n \mathbf{X}_{A_{ij}} \mathbf{X}_{B_{ij}}^T \right) \\
&\quad + O_P(m) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T.
\end{aligned}$$

The bottom right block of  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$  can then be re-expressed as follows

$$\begin{aligned}
& \frac{mn}{\sigma_\varepsilon^2} \left( \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{B_{ij}} \mathbf{X}_{B_{ij}}^T \right) \\
&\quad - \frac{mn}{\sigma_\varepsilon^2} \left\{ \frac{1}{mn} \sum_{i=1}^m \left( \sum_{j=1}^n \mathbf{X}_{B_{ij}} \mathbf{X}_{A_{ij}}^T \right) \left( \sum_{j=1}^n \mathbf{X}_{A_{ij}} \mathbf{X}_{A_{ij}}^T \right)^{-1} \left( \sum_{j=1}^n \mathbf{X}_{A_{ij}} \mathbf{X}_{B_{ij}}^T \right) \right\} \\
&\quad + O_P(m) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T.
\end{aligned}$$

Both of the leading terms in the previous expression are of order  $O_P(mn) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T$  but it remains to determine the overall order of magnitude of the bottom right block of  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$ . In order to do this, we shall treat both of the leading terms separately.

#### Treatment of First Leading Term

Note that for a Gaussian mixed model with multivariate fixed and random effects as described in Section 3.1,

$$b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{A_{ij}} + \boldsymbol{\beta}_B^T \mathbf{X}_{B_{ij}}) = 1, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n_i.$$

Then from using Lemma 1, we have

$$\frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{B_{ij}} \mathbf{X}_{B_{ij}}^T = E(\mathbf{X}_B \mathbf{X}_B^T) + o_P(1) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T.$$

Finally we obtain

$$\frac{mn}{\sigma_\varepsilon^2} \left( \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Bij}^T \right) = \frac{mn}{\sigma_\varepsilon^2} E(\mathbf{X}_B \mathbf{X}_B^T) + o_P(mn) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T. \quad (3.6)$$

#### Treatment of Second Leading Term

Note that for a Gaussian mixed model with multivariate fixed and random effects as described in Section 3.1,

$$f(\mathbf{X}_{ij}, \mathbf{U}_i) = 1, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n_i,$$

where  $f(\mathbf{X}_{ij}, \mathbf{U}_i)$  is as defined in Lemma 2. Assuming that the conditions and assumptions in Lemma 2 are met, in the Gaussian mixed model case, we can simplify the left side of Lemma 2 as follows

$$\begin{aligned} & E \left[ \frac{1}{mn} \sum_{i=1}^m \left\{ \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T f(\mathbf{X}_{ij}, \mathbf{U}_i) \right\} \left\{ \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T f(\mathbf{X}_{ij}, \mathbf{U}_i) \right\}^{-1} \times \right. \\ & \quad \left. \left\{ \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T f(\mathbf{X}_{ij}, \mathbf{U}_i) \right\}^T \middle| \mathbf{X}_{11}, \dots, \mathbf{X}_{mn_m} \right] \\ &= E \left\{ \frac{1}{mn} \sum_{i=1}^m \left( \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T \right) \left( \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T \right)^{-1} \left( \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T \right)^T \middle| \mathbf{X}_{11}, \dots, \mathbf{X}_{mn_m} \right\} \\ &= \frac{1}{mn} \sum_{i=1}^m \left( \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T \right) \left( \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T \right)^{-1} \left( \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Bij}^T \right). \end{aligned}$$

We can also simplify the right hand side of Lemma 2 as follows

$$\begin{aligned} & E \left[ E(\mathbf{X}_B \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}) \{ E(\mathbf{X}_A \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U}) \}^{-1} E(\mathbf{X}_B \mathbf{X}_A^T f(\mathbf{X}, \mathbf{U}) | \mathbf{U})^T \right] \\ &= E \left[ E(\mathbf{X}_B \mathbf{X}_A^T) \{ E(\mathbf{X}_A \mathbf{X}_A^T) \}^{-1} E(\mathbf{X}_B \mathbf{X}_A^T)^T \right] \\ &= E(\mathbf{X}_B \mathbf{X}_A^T) \{ E(\mathbf{X}_A \mathbf{X}_A^T) \}^{-1} E(\mathbf{X}_A \mathbf{X}_B^T). \end{aligned}$$

Now we have a simplified version of Lemma 2 for the Gaussian response linear mixed model which states

$$\begin{aligned} & \frac{1}{mn} \sum_{i=1}^m \left( \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T \right) \left( \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T \right)^{-1} \left( \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Bij}^T \right) \\ & \xrightarrow{P} E(\mathbf{X}_B \mathbf{X}_A^T) \{ E(\mathbf{X}_A \mathbf{X}_A^T) \}^{-1} E(\mathbf{X}_A \mathbf{X}_B^T). \end{aligned} \quad (3.7)$$

Equation (3.7) can also be written as

$$\begin{aligned} & \frac{1}{mn} \sum_{i=1}^m \left( \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T \right) \left( \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T \right)^{-1} \left( \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Bij}^T \right) \\ &= E(\mathbf{X}_B \mathbf{X}_A^T) \{ E(\mathbf{X}_A \mathbf{X}_A^T) \}^{-1} E(\mathbf{X}_A \mathbf{X}_B^T) + o_P(1) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T. \end{aligned} \quad (3.8)$$

Therefore we have,

$$\begin{aligned} & \frac{mn}{\sigma_\varepsilon^2} \left\{ \frac{1}{mn} \sum_{i=1}^m \left( \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T \right) \left( \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T \right)^{-1} \left( \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Bij}^T \right) \right\} \\ &= \frac{mn}{\sigma_\varepsilon^2} E(\mathbf{X}_B \mathbf{X}_A^T) \{E(\mathbf{X}_A \mathbf{X}_A^T)\}^{-1} E(\mathbf{X}_A \mathbf{X}_B^T) + o_P(mn) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T. \end{aligned}$$

Hence, from equations (3.6) and (3.7) we have the following expression for the bottom right block of  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$

$$\begin{aligned} & \sum_{i=1}^m \mathbf{X}_{Bi}^T (\mathbf{X}_{Ai} \boldsymbol{\Sigma} \mathbf{X}_{Ai}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{Bi} \\ &= \frac{mn}{\sigma_\varepsilon^2} \left[ E(\mathbf{X}_B \mathbf{X}_B^T) - E(\mathbf{X}_B \mathbf{X}_A^T) \{E(\mathbf{X}_A \mathbf{X}_A^T)\}^{-1} E(\mathbf{X}_A \mathbf{X}_B^T) \right] + o_P(mn) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T. \end{aligned}$$

Consider the following

$$E \left( \begin{bmatrix} \mathbf{X}_A \mathbf{X}_A^T & \mathbf{X}_A \mathbf{X}_B^T \\ \mathbf{X}_B \mathbf{X}_A^T & \mathbf{X}_B \mathbf{X}_B^T \end{bmatrix} \right)^{-1} = \begin{bmatrix} E(\mathbf{X}_A \mathbf{X}_A^T) & E(\mathbf{X}_A \mathbf{X}_B^T) \\ E(\mathbf{X}_B \mathbf{X}_A^T) & E(\mathbf{X}_B \mathbf{X}_B^T) \end{bmatrix}^{-1}.$$

Then let

$$\begin{aligned} \boldsymbol{\Sigma}_{\boldsymbol{\beta}_B} &= \text{lower right } d_B \times d_B \text{ block of } E \left( \begin{bmatrix} \mathbf{X}_A \mathbf{X}_A^T & \mathbf{X}_A \mathbf{X}_B^T \\ \mathbf{X}_B \mathbf{X}_A^T & \mathbf{X}_B \mathbf{X}_B^T \end{bmatrix} \right)^{-1} \\ &= \left[ E(\mathbf{X}_B \mathbf{X}_B^T) - E(\mathbf{X}_B \mathbf{X}_A^T) \{E(\mathbf{X}_A \mathbf{X}_A^T)\}^{-1} E(\mathbf{X}_A \mathbf{X}_B^T) \right]^{-1}. \end{aligned}$$

The expression for the bottom right block of  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$  can then be re-defined as

$$\frac{mn \boldsymbol{\Sigma}_{\boldsymbol{\beta}_B}^{-1}}{\sigma_\varepsilon^2} + o_P(mn) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T.$$

Putting the expressions obtained for the sub-blocks of  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$  together, we have

$$I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B) = \begin{bmatrix} m \boldsymbol{\Sigma}^{-1} + O_P(mn^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T & O_P(m) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T \\ O_P(m) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T & \frac{mn \boldsymbol{\Sigma}_{\boldsymbol{\beta}_B}^{-1}}{\sigma_\varepsilon^2} + o_P(mn) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T \end{bmatrix}. \quad (3.9)$$

### 3.4.3 Expression for Bottom Right Block of Fisher Information Matrix

Using similar steps to Section 3.4.2, an explicit expression for  $I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2)$  can be obtained. Expanding the current expression for  $I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2)$  leads to

$$\begin{aligned}
& I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2) \\
&= \frac{1}{2} \sum_{i=1}^m \mathcal{K}_i^T (\mathbf{V}_{\theta_i}^{-1} \otimes \mathbf{V}_{\theta_i}^{-1}) \mathcal{K}_i \\
&= \frac{1}{2} \sum_{i=1}^m \left[ (\mathbf{X}_{A_i} \otimes \mathbf{X}_{A_i}) \mathbf{D}_{d_A} \quad \text{vec}(\mathbf{I}_{n_i}) \right]^T \left\{ (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \right\} \times \\
&\quad \left[ (\mathbf{X}_{A_i} \otimes \mathbf{X}_{A_i}) \mathbf{D}_{d_A} \quad \text{vec}(\mathbf{I}_{n_i}) \right] \\
&= \left[ \begin{aligned} & \frac{1}{2} \sum_{i=1}^m \mathbf{D}_{d_A}^T (\mathbf{X}_{A_i}^T \otimes \mathbf{X}_{A_i}^T) \left\{ (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \right\} (\mathbf{X}_{A_i} \otimes \mathbf{X}_{A_i}) \mathbf{D}_{d_A} \\ & \frac{1}{2} \sum_{i=1}^m \{ \text{vec}(\mathbf{I}_{n_i}) \}^T \left\{ (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \right\} (\mathbf{X}_{A_i} \otimes \mathbf{X}_{A_i}) \mathbf{D}_{d_A} \\ & \frac{1}{2} \sum_{i=1}^m \mathbf{D}_{d_A}^T (\mathbf{X}_{A_i}^T \otimes \mathbf{X}_{A_i}^T) \left\{ (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \right\} \text{vec}(\mathbf{I}_{n_i}) \\ & \frac{1}{2} \sum_{i=1}^m \{ \text{vec}(\mathbf{I}_{n_i}) \}^T \left\{ (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \right\} \text{vec}(\mathbf{I}_{n_i}) \end{aligned} \right]. \tag{3.10}
\end{aligned}$$

Once again, it remains to evaluate and simplify the expressions in the matrix, which can be done by making use of the matrix identities from Harville (1977) highlighted in Section 1.4.7.1.

#### 3.4.3.1 Top Left Block of (3.10)

Now, we solve a part of the expression of the top left block of  $I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2)$ , specifically

$$\mathbf{D}_{d_A}^T (\mathbf{X}_{A_i}^T \otimes \mathbf{X}_{A_i}^T) \left\{ (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \right\} (\mathbf{X}_{A_i} \otimes \mathbf{X}_{A_i}) \mathbf{D}_{d_A}.$$

Using the properties of Kronecker products listed under Subsection 1.4.5 we have,

$$\begin{aligned}
& \mathbf{D}_{d_A}^T (\mathbf{X}_{A_i}^T \otimes \mathbf{X}_{A_i}^T) \left\{ (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \right\} (\mathbf{X}_{A_i} \otimes \mathbf{X}_{A_i}) \mathbf{D}_{d_A} \\
&= \mathbf{D}_{d_A}^T \left[ \left\{ \mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{A_i} \right\} \otimes \left\{ \mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{X}_{A_i} \right\} \right] \mathbf{D}_{d_A} \\
&= \mathbf{D}_{d_A}^T \left[ \left\{ \mathbf{\Sigma}^{-1} + O_P(n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T \right\} \otimes \left\{ \mathbf{\Sigma}^{-1} + O_P(n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T \right\} \right] \mathbf{D}_{d_A} \\
&= \mathbf{D}_{d_A}^T (\mathbf{\Sigma}^{-1} \otimes \mathbf{\Sigma}^{-1}) \mathbf{D}_{d_A} + O_P(n^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T.
\end{aligned}$$

Therefore, the top left block of  $I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2)$  can be computed as

$$\begin{aligned}
& \frac{1}{2} \sum_{i=1}^m \left\{ \mathbf{D}_{d_A}^T (\mathbf{\Sigma}^{-1} \otimes \mathbf{\Sigma}^{-1}) \mathbf{D}_{d_A} + O_P(n^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T \right\} \\
&= \frac{m}{2} \left\{ \mathbf{D}_{d_A}^T (\mathbf{\Sigma}^{-1} \otimes \mathbf{\Sigma}^{-1}) \mathbf{D}_{d_A} \right\} + O_P(mn^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T.
\end{aligned}$$

### 3.4.3.2 Top Right Block of (3.10)

Now, we solve a part of the expression of the top right block of  $I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2)$ , specifically

$$D_{d_A}^T (\mathbf{X}_{A_i}^T \otimes \mathbf{X}_{A_i}^T) \{(\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1}\} \text{vec}(\mathbf{I}_{n_i}).$$

It follows that

$$\begin{aligned} & D_{d_A}^T (\mathbf{X}_{A_i}^T \otimes \mathbf{X}_{A_i}^T) \{(\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1}\} \text{vec}(\mathbf{I}_{n_i}) \\ &= D_{d_A}^T \{ \mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \otimes \mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \} \text{vec}(\mathbf{I}_{n_i}) \\ &= D_{d_A}^T \text{vec} \left[ \mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \mathbf{I}_{n_i} \{ \mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \}^T \right] \\ &= D_{d_A}^T \text{vec} \left[ \mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \{ \mathbf{X}_{A_i}^T (\mathbf{X}_{A_i} \mathbf{\Sigma} \mathbf{X}_{A_i}^T + \sigma_\varepsilon^2 \mathbf{I}_{n_i})^{-1} \}^T \right] \\ &= D_{d_A}^T \text{vec} \left[ \left\{ \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} + \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} + \dots \right\}^{-1} \right. \\ &\quad \times \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \right) \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i} \right) \left\{ \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} + \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \right. \\ &\quad \left. \left. \times (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} + \dots \right\}^{-1} \right]. \end{aligned}$$

By expanding and simplifying the last line of the previous expression, we have,

$$\begin{aligned} & D_{d_A}^T \text{vec} \left[ \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \right) \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i} \right) \sigma_\varepsilon^2 (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \right. \\ &\quad - \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \right) \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i} \right) \sigma_\varepsilon^2 (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \\ &\quad - \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \right) \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i} \right) \sigma_\varepsilon^2 (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \\ &\quad - \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i}^T \right) \left( \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{A_i} \right) \sigma_\varepsilon^2 (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \\ &\quad \left. \times (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} + \dots \right] \\ &= D_{d_A}^T \text{vec} \left\{ \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} - \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \right. \\ &\quad \left. - \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \right. \\ &\quad \left. + (\sigma_\varepsilon^2)^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} + \dots \right\} \\ &= D_{d_A}^T \left[ \text{vec} \{ \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \} - 2 \text{vec} \{ \sigma_\varepsilon^2 \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \} + \dots \right] \\ &= D_{d_A}^T \text{vec} \{ \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \} + O_P(n^{-2}) \mathbf{1}_{d_A(d_A+1)/2} \\ &= O_P(n^{-1}) \mathbf{1}_{d_A(d_A+1)/2}. \end{aligned}$$

Therefore, the top right block of  $I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2)$  can be computed as

$$\frac{1}{2} \sum_{i=1}^m D_{d_A}^T \text{vec} \{ \mathbf{\Sigma}^{-1} (\mathbf{X}_{A_i}^T \mathbf{X}_{A_i})^{-1} \mathbf{\Sigma}^{-1} \} + O_P(mn^{-2}) \mathbf{1}_{d_A(d_A+1)/2} = O_P(mn^{-1}) \mathbf{1}_{d_A(d_A+1)/2}.$$

### 3.4.3.3 Bottom Left Block of (3.10)

Now, we solve a part of the expression of the bottom left block of  $I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2)$ , specifically

$$\{\text{vec}(\mathbf{I}_{n_i})\}^T \{(\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1}\} (\mathbf{X}_{Ai} \otimes \mathbf{X}_{Ai}) \mathbf{D}_{d_A}.$$

Using the properties for Kronecker products highlighted in Subsection 1.4.5, we have,

$$\begin{aligned} & \{\text{vec}(\mathbf{I}_{n_i})\}^T \{(\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1}\} (\mathbf{X}_{Ai} \otimes \mathbf{X}_{Ai}) \mathbf{D}_{d_A} \\ &= [\mathbf{D}_{d_A}^T (\mathbf{X}_{Ai}^T \otimes \mathbf{X}_{Ai}^T) \{(\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1}\} \text{vec}(\mathbf{I}_{n_i})]^T \\ &= [\mathbf{D}_{d_A}^T \text{vec}\{\mathbf{\Sigma}^{-1}(\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \mathbf{\Sigma}^{-1}\} + O_P(n^{-2})\mathbf{1}_{d_A(d_A+1)/2}]^T \\ &= [\mathbf{D}_{d_A}^T \text{vec}\{\mathbf{\Sigma}^{-1}(\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \mathbf{\Sigma}^{-1}\}]^T + O_P(n^{-2})\mathbf{1}_{d_A(d_A+1)/2}^T \\ &= O_P(n^{-1})\mathbf{1}_{d_A(d_A+1)/2}^T. \end{aligned}$$

Therefore, the bottom left block of  $I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2)$  can be computed as

$$\frac{1}{2} \sum_{i=1}^m [\mathbf{D}_{d_A}^T \text{vec}\{\mathbf{\Sigma}^{-1}(\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \mathbf{\Sigma}^{-1}\}]^T + O_P(mn^{-2})\mathbf{1}_{d_A(d_A+1)/2}^T = O_P(mn^{-1})\mathbf{1}_{d_A(d_A+1)/2}^T.$$

### 3.4.3.4 Bottom Right Block of (3.10)

Now, we solve a part of the expression of the bottom right block of  $I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2)$ , specifically

$$\{\text{vec}(\mathbf{I}_{n_i})\}^T \{(\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1}\} \{\text{vec}(\mathbf{I}_{n_i})\}.$$

Using (1.10a) and (1.5), we have

$$\begin{aligned} & \{\text{vec}(\mathbf{I}_{n_i})\}^T \{(\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1} \otimes (\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1}\} \{\text{vec}(\mathbf{I}_{n_i})\} \\ &= \text{tr} \{(\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1} (\mathbf{X}_{Ai}\mathbf{\Sigma}\mathbf{X}_{Ai}^T + \sigma_\varepsilon^2\mathbf{I}_{n_i})^{-1}\} \\ &= \text{tr} \left[ \left\{ \frac{1}{\sigma_\varepsilon^2} \mathbf{I}_{n_i} - \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{Ai} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \mathbf{X}_{Ai}^T + \dots \right\} \left\{ \frac{1}{\sigma_\varepsilon^2} \mathbf{I}_{n_i} - \frac{1}{\sigma_\varepsilon^2} \mathbf{X}_{Ai} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \mathbf{X}_{Ai}^T + \dots \right\} \right] \\ &= \text{tr} \left\{ \left( \frac{1}{\sigma_\varepsilon^2} \right)^2 \mathbf{I}_{n_i} - 2 \left( \frac{1}{\sigma_\varepsilon^2} \right)^2 \mathbf{X}_{Ai} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \mathbf{X}_{Ai}^T \right. \\ & \quad \left. + \left( \frac{1}{\sigma_\varepsilon^2} \right)^2 \mathbf{X}_{Ai} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \mathbf{X}_{Ai}^T \mathbf{X}_{Ai} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \mathbf{X}_{Ai}^T \right\} + \dots \\ &= \left( \frac{1}{\sigma_\varepsilon^2} \right)^2 \text{tr}(\mathbf{I}_{n_i}) - \left( \frac{1}{\sigma_\varepsilon^2} \right)^2 \text{tr} \{ \mathbf{X}_{Ai} (\mathbf{X}_{Ai}^T \mathbf{X}_{Ai})^{-1} \mathbf{X}_{Ai}^T \} + \dots \\ &= \frac{n}{(\sigma_\varepsilon^2)^2} + O_P(1). \end{aligned}$$

Therefore, the bottom right block of  $I(\text{vech}(\mathbf{\Sigma}), \sigma_\varepsilon^2)$  can be computed as

$$\frac{1}{2} \sum_{i=1}^m \left\{ \frac{n}{(\sigma_\varepsilon^2)^2} + O_P(1) \right\} = \frac{mn}{2(\sigma_\varepsilon^2)^2} + O_P(m).$$

Putting the expressions obtained for the sub-blocks of  $I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)$  together, we have

$$I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2) = \begin{bmatrix} \frac{m}{2} \left\{ \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} \right\} + O_P(mn^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T & O_P(mn^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \\ O_P(mn^{-1}) \mathbf{1}_{d_A(d_A+1)/2}^T & \frac{mn}{2(\sigma_\varepsilon^2)^2} + O_P(m) \end{bmatrix}. \quad (3.11)$$

### 3.4.4 The Inverse of the Fisher Information Matrix

Since  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)$  is a block diagonal matrix, we can invert  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$  and  $I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)$  separately and put the expressions together to find  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)^{-1}$  as shown below.

$$I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)^{-1} = \begin{bmatrix} I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)^{-1} & 0 \\ 0 & I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)^{-1} \end{bmatrix}$$

#### 3.4.4.1 Expression for Top Left Block of Inverse Fisher Information Matrix

Firstly, let us partition  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)$  as follows

$$I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \quad \text{where } \mathbf{A}_{21} = \mathbf{A}_{12}^T.$$

The expressions for  $\mathbf{A}_{11}$ ,  $\mathbf{A}_{12}$ ,  $\mathbf{A}_{21}$  and  $\mathbf{A}_{22}$  are currently as follows

$$\begin{aligned} \mathbf{A}_{11} &= m\boldsymbol{\Sigma}^{-1} + O_P(mn^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T, \\ \mathbf{A}_{12} &= O_P(m) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T, \\ \mathbf{A}_{21} &= O_P(m) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T, \\ \mathbf{A}_{22} &= \frac{mn\boldsymbol{\Sigma}_B^{-1}}{\sigma_\varepsilon^2} + O_P(mn) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T. \end{aligned}$$

Let  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)^{-1}$  assume the following form

$$\begin{aligned} I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)^{-1} &= \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}^{-1} \\ &= \begin{bmatrix} \mathbf{A}^{11} & \mathbf{A}^{12} \\ \mathbf{A}^{21} & \mathbf{A}^{22} \end{bmatrix} \quad \text{where } \mathbf{A}^{21} = (\mathbf{A}^{12})^T. \end{aligned}$$

Firstly note that

$$\begin{aligned} \mathbf{A}_{11}^{-1} &= \{m\boldsymbol{\Sigma}^{-1} + O_P(mn^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T\}^{-1} \\ &= [m\boldsymbol{\Sigma}^{-1} \{ \mathbf{I}_{d_A} + O_P(n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T \}]^{-1} \\ &= m\boldsymbol{\Sigma}^{-1} \{ \mathbf{I}_{d_A} - O_P(n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T + \dots \} \\ &= m^{-1} \boldsymbol{\Sigma} + O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T. \end{aligned} \quad (3.12)$$



Also note that

$$\begin{aligned}
\mathbf{A}_{22}^{-1} &= \left\{ \frac{mn\boldsymbol{\Sigma}_{\beta_B}^{-1}}{\sigma_\varepsilon^2} + o_P(mn)\mathbf{1}_{d_B}\mathbf{1}_{d_B}^T \right\}^{-1} \\
&= \left[ \frac{mn\boldsymbol{\Sigma}_{\beta_B}^{-1}}{\sigma_\varepsilon^2} \{ \mathbf{I}_{d_B} + o_P(1)\boldsymbol{\Sigma}_{\beta_B}\mathbf{1}_{d_B}\mathbf{1}_{d_B}^T \} \right]^{-1} \\
&= \frac{\sigma_\varepsilon^2\boldsymbol{\Sigma}_{\beta_B}}{mn} \{ \mathbf{I}_{d_B} + o_P(1)\boldsymbol{\Sigma}_{\beta_B}\mathbf{1}_{d_B}\mathbf{1}_{d_B}^T \}^{-1} \\
&= \frac{\sigma_\varepsilon^2\boldsymbol{\Sigma}_{\beta_B}}{mn} \{ \mathbf{I}_{d_B} + o_P(1)\boldsymbol{\Sigma}_{\beta_B}\mathbf{1}_{d_B}\mathbf{1}_{d_B}^T + \dots \} \\
&= \frac{\sigma_\varepsilon^2\boldsymbol{\Sigma}_{\beta_B}}{mn} + o_P(m^{-1}n^{-1})\mathbf{1}_{d_B}\mathbf{1}_{d_B}^T.
\end{aligned} \tag{3.13}$$

Using the result for carrying out block matrix inversion presented in Section 1,  $\mathbf{A}^{11}$ ,  $\mathbf{A}^{12}$ ,  $\mathbf{A}^{21}$  and  $\mathbf{A}^{22}$  can be calculated as follows. Firstly we have,

$$\mathbf{A}^{11} = \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1}\mathbf{A}_{12}(\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1}.$$

Note that

$$(\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12})^{-1} = \mathbf{A}_{22}^{-1} + \mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} + \dots$$

It follows that

$$\mathbf{A}_{11}^{-1}\mathbf{A}_{12}(\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} = o_P(m^{-1}n^{-1})\mathbf{1}_{d_A}\mathbf{1}_{d_A}^T.$$

Therefore by making use of (3.12), we get,

$$\mathbf{A}^{11} = m^{-1}\boldsymbol{\Sigma} + o_P(m^{-1}n^{-1})\mathbf{1}_{d_A}\mathbf{1}_{d_A}^T.$$

Next we have,

$$\mathbf{A}^{12} = -(\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1}.$$

Note that

$$(\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})^{-1} = \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} + \dots$$

Therefore, by making use of (3.12) and (3.13), we have

$$\mathbf{A}^{12} = o_P(m^{-1}n^{-1})\mathbf{1}_{d_A}\mathbf{1}_{d_B}^T.$$

Note that

$$\mathbf{A}^{21} = (\mathbf{A}^{12})^T.$$

Therefore, we have

$$\mathbf{A}^{21} = o_P(m^{-1}n^{-1})\mathbf{1}_{d_B}\mathbf{1}_{d_A}^T.$$

Finally, making use of (3.13), we get

$$\mathbf{A}^{22} = \mathbf{A}_{22}^{-1} + \mathbf{A}_{22}^{-1}\mathbf{A}_{21}(\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1}.$$

Also note that

$$\mathbf{A}_{22}^{-1} \mathbf{A}_{21} (\mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21})^{-1} \mathbf{A}_{12} \mathbf{A}_{22}^{-1} = O_P(m^{-1} n^{-2}) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T.$$

This leads to

$$\begin{aligned} \mathbf{A}^{22} &= \frac{\sigma_\varepsilon^2 \boldsymbol{\Sigma} \boldsymbol{\beta}_B}{mn} + o_P(m^{-1} n^{-1}) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T + O_P(m^{-1} n^{-2}) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T \\ &= \frac{\sigma_\varepsilon^2 \boldsymbol{\Sigma} \boldsymbol{\beta}_B}{mn} + o_P(m^{-1} n^{-1}) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T. \end{aligned}$$

Using the expressions for  $\mathbf{A}^{11}$ ,  $\mathbf{A}^{12}$ ,  $\mathbf{A}^{21}$  and  $\mathbf{A}^{22}$ , we have the following expression for  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)^{-1}$

$$I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)^{-1} = \begin{bmatrix} m^{-1} \boldsymbol{\Sigma} + O_P(m^{-1} n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T & O_P(m^{-1} n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T \\ O_P(m^{-1} n^{-1}) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T & \frac{\sigma_\varepsilon^2 \boldsymbol{\Sigma} \boldsymbol{\beta}_B}{mn} + o_P(m^{-1} n^{-1}) \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T \end{bmatrix}.$$

#### 3.4.4.2 Expression for Bottom Right Block of Inverse Fisher Information Matrix

Similarly, let us partition  $I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)$  as follows

$$I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2) = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} \text{ where } \mathbf{B}_{21} = \mathbf{B}_{12}^T.$$

The expressions for  $\mathbf{B}_{11}$ ,  $\mathbf{B}_{12}$ ,  $\mathbf{B}_{21}$  and  $\mathbf{B}_{22}$  are currently as follows

$$\begin{aligned} \mathbf{B}_{11} &= \frac{m}{2} \{ \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} \} + O_P(mn^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T, \\ \mathbf{B}_{12} &= O_P(mn^{-1}) \mathbf{1}_{d_A(d_A+1)/2}, \\ \mathbf{B}_{21} &= O_P(mn^{-1}) \mathbf{1}_{d_A(d_A+1)/2}^T, \\ \mathbf{B}_{22} &= \frac{mn}{2(\sigma_\varepsilon^2)^2} + O_P(m). \end{aligned}$$

Let  $I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)^{-1}$  assume the following form

$$\begin{aligned} I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)^{-1} &= \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix}^{-1} \\ &= \begin{bmatrix} \mathbf{B}^{11} & \mathbf{B}^{12} \\ \mathbf{B}^{21} & \mathbf{B}^{22} \end{bmatrix} \text{ where } \mathbf{B}^{21} = (\mathbf{B}^{12})^T. \end{aligned}$$

Firstly note that

$$\begin{aligned}
\mathbf{B}_{11}^{-1} &= \left[ \frac{m}{2} \left\{ \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} \right\} + O_P(mn^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T \right]^{-1} \\
&= \left[ \frac{m}{2} \left\{ \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} \right\} \left\{ \mathbf{I}_{d_A(d_A+1)/2} + O_P(n^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T \right\} \right]^{-1} \\
&= \frac{2}{m} \left\{ \mathbf{I}_{d_A(d_A+1)/2} + O_P(n^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T \right\}^{-1} \left\{ \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} \right\}^{-1} \\
&= \frac{2}{m} \left\{ \mathbf{I}_{d_A(d_A+1)/2} + O_P(n^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T + \dots \right\} \mathbf{D}_{d_A}^+ (\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}) \mathbf{D}_{d_A}^{+T} \\
&= \frac{2\mathbf{D}_{d_A}^+ (\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}) \mathbf{D}_{d_A}^{+T}}{m} + O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T.
\end{aligned} \tag{3.14}$$

Also note that

$$\begin{aligned}
\mathbf{B}_{22}^{-1} &= \left\{ \frac{mn}{2(\sigma_\varepsilon^2)^2} + O_P(m) \right\}^{-1} \\
&= \left[ \frac{mn}{2(\sigma_\varepsilon^2)^2} \{1 + O_P(n^{-1})\} \right]^{-1} \\
&= \frac{2(\sigma_\varepsilon^2)^2}{mn} \{1 + O_P(n^{-1}) + \dots\} \\
&= \frac{2(\sigma_\varepsilon^2)^2}{mn} + O_P(m^{-1}n^{-2}).
\end{aligned} \tag{3.15}$$

Using the result for carrying out block matrix inversion presented in Section 1,  $\mathbf{B}^{11}$ ,  $\mathbf{B}^{12}$ ,  $\mathbf{B}^{21}$  and  $\mathbf{B}^{22}$  can be calculated as follows. Firstly, we have,

$$\mathbf{B}^{11} = \mathbf{B}_{11}^{-1} + \mathbf{B}_{11}^{-1} \mathbf{B}_{12} (\mathbf{B}_{22} - \mathbf{B}_{21} \mathbf{B}_{11}^{-1} \mathbf{B}_{12})^{-1} \mathbf{B}_{21} \mathbf{B}_{11}^{-1}.$$

Note that

$$(\mathbf{B}_{22} - \mathbf{B}_{21} \mathbf{B}_{11}^{-1} \mathbf{B}_{12})^{-1} = \mathbf{B}_{22}^{-1} + \mathbf{B}_{22}^{-1} \mathbf{B}_{21} \mathbf{A}_{11}^{-1} \mathbf{B}_{12} \mathbf{B}_{22}^{-1} + \dots$$

It follows that

$$\mathbf{B}_{11}^{-1} \mathbf{B}_{12} (\mathbf{B}_{22} - \mathbf{B}_{21} \mathbf{B}_{11}^{-1} \mathbf{B}_{12})^{-1} \mathbf{B}_{21} \mathbf{B}_{11}^{-1} = O_P(m^{-1}n^{-3}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T.$$

Therefore by making use of (3.14), we get,

$$\mathbf{B}^{11} = \frac{2\mathbf{D}_{d_A}^+ (\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}) \mathbf{D}_{d_A}^{+T}}{m} + O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A(d_A+1)/2} \mathbf{1}_{d_A(d_A+1)/2}^T.$$

Next we have,

$$\mathbf{B}^{12} = -(\mathbf{B}_{11} - \mathbf{B}_{12} \mathbf{B}_{22}^{-1} \mathbf{B}_{21})^{-1} \mathbf{B}_{12} \mathbf{B}_{22}^{-1}.$$

Note that

$$(\mathbf{B}_{11} - \mathbf{B}_{12} \mathbf{B}_{22}^{-1} \mathbf{B}_{21})^{-1} = \mathbf{B}_{11}^{-1} + \mathbf{B}_{11}^{-1} \mathbf{B}_{12} \mathbf{B}_{22}^{-1} \mathbf{B}_{21} \mathbf{B}_{11}^{-1} + \dots$$

Therefore, by making use of (3.14) and (3.15), we have

$$\mathbf{B}^{12} = O_P(m^{-1}n^{-2}) \mathbf{1}_{d_A(d_A+1)/2}.$$

Subsequently

$$\mathbf{B}^{21} = (\mathbf{B}^{12})^T.$$

Therefore we have

$$\mathbf{B}^{21} = O_P(m^{-1}n^{-2})\mathbf{1}_{d_A(d_A+1)/2}^T.$$

Finally, making use of (3.15), we get

$$\mathbf{B}^{22} = \mathbf{B}_{22}^{-1} + \mathbf{B}_{22}^{-1}\mathbf{B}_{21}(\mathbf{B}_{11} - \mathbf{B}_{12}\mathbf{B}_{22}^{-1}\mathbf{B}_{21})^{-1}\mathbf{B}_{12}\mathbf{B}_{22}^{-1}.$$

Note that

$$\mathbf{B}_{22}^{-1}\mathbf{B}_{21}(\mathbf{B}_{11} - \mathbf{B}_{12}\mathbf{B}_{22}^{-1}\mathbf{B}_{21})^{-1}\mathbf{B}_{12}\mathbf{B}_{22}^{-1} = O_P(m^{-1}n^{-4}).$$

This leads to

$$\begin{aligned} \mathbf{B}^{22} &= \frac{2(\sigma_\varepsilon^2)^2}{mn} + O_P(m^{-1}n^{-2}) + O_P(m^{-1}n^{-4}) \\ &= \frac{2(\sigma_\varepsilon^2)^2}{mn} + O_P(m^{-1}n^{-2}). \end{aligned}$$

Using the expressions for  $\mathbf{B}^{11}$ ,  $\mathbf{B}^{12}$ ,  $\mathbf{B}^{21}$  and  $\mathbf{B}^{22}$ , we have the following expression for  $I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)^{-1}$

$$\begin{aligned} &I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)^{-1} \\ &= \begin{bmatrix} \frac{2\mathbf{D}_{d_A}^+(\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma})\mathbf{D}_{d_A}^{+T}}{m} + O_P(m^{-1}n^{-1})\mathbf{1}_{d_A(d_A+1)/2}\mathbf{1}_{d_A(d_A+1)/2}^T & O_P(m^{-1}n^{-2})\mathbf{1}_{d_A(d_A+1)/2} \\ O_P(m^{-1}n^{-2})\mathbf{1}_{d_A(d_A+1)/2}^T & \frac{2(\sigma_\varepsilon^2)^2}{mn} + O_P(m^{-1}n^{-2}) \end{bmatrix}. \end{aligned}$$

The resultant expression for  $I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)$  is

$$\begin{aligned} &I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)^{-1} \\ &= \begin{bmatrix} I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B)^{-1} & \mathbf{0} \\ \mathbf{0} & I(\text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)^{-1} \end{bmatrix} \\ &= I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)_\infty^{-1} \\ &\quad + \frac{1}{mn} \begin{bmatrix} O_P(1)\mathbf{1}_{d_A}\mathbf{1}_{d_A}^T & O_P(1)\mathbf{1}_{d_A}\mathbf{1}_{d_B}^T & \mathbf{0} & \mathbf{0} \\ O_P(1)\mathbf{1}_{d_B}\mathbf{1}_{d_A}^T & O_P(1)\mathbf{1}_{d_B}\mathbf{1}_{d_B}^T & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & O_P(1)\mathbf{1}_{d_A(d_A+1)/2}\mathbf{1}_{d_A(d_A+1)/2}^T & O_P(n^{-1})\mathbf{1}_{d_A(d_A+1)/2} \\ \mathbf{0} & \mathbf{0} & O_P(n^{-1})\mathbf{1}_{d_A(d_A+1)/2}^T & O_P(n^{-1}) \end{bmatrix}, \end{aligned}$$

where

$$I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}), \sigma_\varepsilon^2)_\infty^{-1} = \begin{bmatrix} \frac{\boldsymbol{\Sigma}}{m} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \frac{\sigma_\varepsilon^2 \boldsymbol{\Sigma} \boldsymbol{\beta}_B}{mn} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \frac{2\mathbf{D}_{d_A}^+(\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma})\mathbf{D}_{d_A}^{+T}}{m} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \frac{2(\sigma_\varepsilon^2)^2}{mn} \end{bmatrix}.$$

### 3.4.5 Derivation of the Final Asymptotic Normality Result for Gaussian Response Linear Mixed Models

For a matrix  $\mathbf{M}$  let

$$\|\mathbf{M}\|_F = \sqrt{\text{tr}(\mathbf{M}^T \mathbf{M})}$$

denote the *Frobenius norm* of  $\mathbf{M}$ .

Consider working with the order  $(\beta_A, \text{vech}(\Sigma), \beta_B, \sigma_\varepsilon^2)$  rather than the order in  $(\beta_A, \beta_B, \text{vech}(\Sigma), \sigma_\varepsilon^2)$ , when working with the Fisher information matrix for the derivation of the final asymptotic normality result. From standard results concerning asymptotic normality of maximum likelihood estimators we have

$$\{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2}(\hat{\theta} - \theta^0) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{I})$$

where  $\hat{\theta} = [\hat{\beta}_A^T \text{vech}(\hat{\Sigma})^T \hat{\beta}_B^T \hat{\sigma}^2]^T$  and  $\theta^0 = [(\beta_A^0)^T \{\text{vech}(\Sigma^0)\}^T \beta_B^0]^T ((\sigma^2)^0)^T]^T$ . Therefore, for all  $(d_A + d_A(d_A + 1)/2 + d_B + 1) \times 1$  vectors  $\mathbf{a} \neq \mathbf{0}$  we have

$$\mathbf{a}^T \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2}(\hat{\theta} - \theta^0) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{a}^T \mathbf{a}).$$

Note that

$$\begin{aligned} & \mathbf{a}^T \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2}(\hat{\theta} - \theta^0) \\ &= \mathbf{a}^T \left[ \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)_\infty^{-1}\}^{-1/2} + \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2} \right. \\ & \quad \left. - \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)_\infty^{-1}\}^{-1/2} \right] (\hat{\theta} - \theta^0) \\ &= \mathbf{a}^T \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)_\infty^{-1}\}^{-1/2}(\hat{\theta} - \theta^0) \\ & \quad + \mathbf{a}^T \left[ \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2} \right. \\ & \quad \left. - \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)_\infty^{-1}\}^{-1/2} \right] (\hat{\theta} - \theta^0). \end{aligned}$$

As a consequence

$$\mathbf{a}^T \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)_\infty^{-1}\}^{-1/2}(\hat{\theta} - \theta^0) + r_{mn}(\mathbf{a}) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{a}^T \mathbf{a}) \quad (3.16)$$

with

$$\begin{aligned} & r_{mn}(\mathbf{a}) \\ &= \mathbf{a}^T \left[ \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2} - \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)_\infty^{-1}\}^{-1/2} \right] (\hat{\theta} - \theta^0) \\ &= \mathbf{a}^T \left[ \mathbf{I} - \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)_\infty^{-1}\}^{-1/2} \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{1/2} \right] \\ & \quad \times \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2}(\hat{\theta} - \theta^0) \\ &= \left( - \left[ \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)_\infty^{-1}\}^{-1/2} \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{1/2} - \mathbf{I} \right]^T \mathbf{a} \right)^T \mathbf{Z} \end{aligned}$$

where  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I}_{d_A+d_A(d_A+1)/2+d_B+1})$ . Next, note that using the matrix norm properties  $\|-\mathbf{A}\| = \|\mathbf{A}\|$  and  $\|\mathbf{AB}\| \leq \|\mathbf{A}\|\|\mathbf{B}\|$  for any matrices  $\mathbf{A}$  and  $\mathbf{B}$  and the fact that  $\|\mathbf{M}^T\|_F = \|\mathbf{M}\|$  for any matrix  $\mathbf{M}$ , we have

$$\begin{aligned} & \left\| - \left[ \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2} \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{1/2} - \mathbf{I} \right]^T \mathbf{a} \right\|_F \\ & \leq \left\| \left[ \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2} \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{1/2} - \mathbf{I} \right]^T \right\|_F \|\mathbf{a}\|_F \\ & = \left\| \left[ \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2} \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{1/2} - \mathbf{I} \right] \right\|_F \|\mathbf{a}\|_F. \end{aligned} \quad (3.17)$$

Our next aim is to establish that

$$\left\| \left[ \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2} \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{1/2} - \mathbf{I} \right] \right\|_F \xrightarrow{P} 0 \quad (3.18)$$

Recall that

$$\begin{aligned} & I(\beta_A, \text{vech}(\Sigma), \beta_B, \sigma_\varepsilon^2)^{-1} \\ & = I(\beta_A, \text{vech}(\Sigma), \beta_B, \sigma_\varepsilon^2)_\infty^{-1} \\ & + \frac{1}{m} \begin{bmatrix} O_P(n^{-1})\mathbf{1}_{d_A}\mathbf{1}_{d_A}^T & \mathbf{0} & O_P(n^{-1})\mathbf{1}_{d_A}\mathbf{1}_{d_B}^T & \mathbf{0} \\ \mathbf{0} & O_P(n^{-1})\mathbf{1}_{d_A(d_A+1)/2}\mathbf{1}_{d_A(d_A+1)/2}^T & \mathbf{0} & O_P(n^{-2})\mathbf{1}_{d_A(d_A+1)/2} \\ O_P(n^{-1})\mathbf{1}_{d_B}\mathbf{1}_{d_A}^T & \mathbf{0} & O_P(n^{-1})\mathbf{1}_{d_B}\mathbf{1}_{d_B}^T & \mathbf{0} \\ \mathbf{0} & O_P(n^{-2})\mathbf{1}_{d_A(d_A+1)/2}^T & \mathbf{0} & O_P(n^{-2}) \end{bmatrix} \end{aligned}$$

where

$$I(\beta_A, \text{vech}(\Sigma), \beta_B, \sigma_\varepsilon^2)_\infty^{-1} = \frac{1}{m} \begin{bmatrix} \Sigma & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 2D_{d_A}^+ (\Sigma \otimes \Sigma) D_{d_A}^{+T} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \frac{\sigma_\varepsilon^2 \Sigma \beta_B}{n} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \frac{2(\sigma_\varepsilon^2)^2}{n} \end{bmatrix},$$

so that

$$\left\{ I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1} \right\}^{-1/2} \left\{ I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1} \right\}^{1/2} = \mathbf{M}_{n,\infty}^{-1/2} \mathbf{M}_n^{1/2}$$

where

$$\mathbf{M}_{n,\infty} = m \left\{ I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1} \right\} \quad \text{and} \quad \mathbf{M}_n = \left\{ I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1} \right\}.$$

Therefore, we can apply Lemma 3 with the following

$$p = d_A + d_A(d_A + 1)/2, \quad \mathbf{K} = \begin{bmatrix} \Sigma & \mathbf{0} \\ \mathbf{0} & 2D_{d_A}^+ (\Sigma \otimes \Sigma) D_{d_A}^{+T} \end{bmatrix}, \quad q = d_B + 1$$

and

$$\mathbf{L} = \begin{bmatrix} \frac{\sigma_\varepsilon^2 \Sigma \beta_B}{n} & \mathbf{0} \\ \mathbf{0} & \frac{2(\sigma_\varepsilon^2)^2}{n} \end{bmatrix}$$

in order to show that (3.18) holds. It then follows from (3.17) and (3.18) that

$$\left[ \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{-1/2} \{I(\beta_A^0, \text{vech}(\Sigma^0), \beta_B^0, (\sigma^2)^0)^{-1}\}^{1/2} - \mathbf{I} \right] \mathbf{a} \xrightarrow{P} \mathbf{0}.$$

Application of Slutsky's Theorem then gives  $r_{mn}(\mathbf{a}) \xrightarrow{P} 0$ . From (3.16) and another application of Slutsky's Theorem we have

$$\mathbf{a}^T \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0, (\sigma^2)^0)_{\infty}^{-1}\}^{-1/2} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{a}^T \mathbf{a}).$$

It then follows from the Cramér-Wold Device and the Continuous Mapping Theorem that

$$\sqrt{m} \begin{bmatrix} \hat{\boldsymbol{\beta}}_A - \boldsymbol{\beta}_A^0 \\ \sqrt{n} (\hat{\boldsymbol{\beta}}_B - \boldsymbol{\beta}_B^0) \\ \text{vech}(\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^0) \\ \sqrt{n} (\hat{\sigma}^2 - (\sigma^2)^0) \end{bmatrix} \xrightarrow{\mathcal{D}} N \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}^0 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\sigma^2)^0 \boldsymbol{\Sigma} \boldsymbol{\beta}_B & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & 2\mathbf{D}_{d_A}^+ (\boldsymbol{\Sigma}^0 \otimes \boldsymbol{\Sigma}^0) \mathbf{D}_{d_A}^{+T} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & 2((\sigma^2)^0)^2 \end{bmatrix} \right)$$

as shown in Theorem 11.

## Chapter 4

# Usable Asymptotic Normality Results and Inference for Generalized Linear Mixed Models

In this chapter, we aim to derive asymptotic normality results beyond those that have been derived in Chapter 3 for Gaussian response linear mixed models. The main theorem in this chapter concerns the joint asymptotic normality of all of the maximum quasi-likelihood estimators for a generalized linear mixed model. Once again, it elegantly shows faster rates of convergence for fixed effects that are not accompanied by a random effect compared to fixed effects that have a partnering random effect. The results derived in this chapter can also be used for the construction of asymptotically valid confidence intervals and Wald tests for generalized linear mixed model analysis, which will be discussed in the next chapter.

In addition, we extend this theorem under certain circumstances to dispersion parameters, introduced to account for overdispersion, as well. For the class of two-parameter exponential families, maximum likelihood estimation is possible for all model parameters including the dispersion parameter. Thus, we extend Theorem 12 and derive the asymptotic normality results for the maximum likelihood estimator for a dispersion parameter in the Gaussian and Gamma response cases.

The appendix contains the proofs for the theorem introduced in this chapter.



## 4.1 Model Description

Consider the following density, or probability mass, function for the class of one-parameter exponential families

$$p(y; \eta) = \exp \{y\eta - b(\eta) + c(y)\} h(y) \quad (4.1)$$

where  $\eta$  is the natural parameter. For example, the Bernoulli probability mass function has  $b(x) = \log(1 + e^x)$ ,  $c(x) = 0$  and  $h(x) = I(x \in \{0, 1\})$ . Whereas for the Poisson mass function,  $b(x) = e^x$ ,  $c(x) = -\log(x!)$  and  $h(x) = I(x \in \{0\} \cup \mathbb{N})$ . Here,  $I(\mathcal{P}) = 1$  if the condition  $\mathcal{P}$  is true and  $I(\mathcal{P}) = 0$  if  $\mathcal{P}$  is false. If the random variable  $Y$  has density, or probability mass, function as in (4.1), then  $E(Y) = b'(\eta)$  and  $\text{Var}(Y) = b''(\eta)$ . To account for overdispersion in data and to allow one to model the variance flexibly, a common modelling extension is implemented such that  $\text{Var}(Y) = \phi b''(\eta)$ , where  $\phi > 0$  represents the dispersion parameter. This involves the replacement of  $\log\{p(y; \eta)\}$  by the following quasi-likelihood function

$$\{y\eta - b(\eta) + c(y)\}/\phi + d(y, \phi) \quad (4.2)$$

where  $d(y, \phi)$  is a function of  $y$  and  $\phi$  only. Note that for ordinary binomial and Poisson response models,  $\phi$  is fixed at 1. For Gaussian and gamma response models, (4.2) corresponds to the expression of  $\log\{p(y; \eta)\}$  for a two-parameter exponential family density function and ordinary likelihood applies. In this section, we study generalized linear mixed models of the following form, for observations of the random triples  $(\mathbf{X}_{Aij}, \mathbf{X}_{Bij}, Y_{ij})$ ,  $1 \leq i \leq m, 1 \leq j \leq n_i$ ,

$$Y_{ij} | \mathbf{X}_{Aij}, \mathbf{X}_{Bij}, \mathbf{U}_i \text{ are independent having quasi-likelihood function (4.2) with natural parameter } (\boldsymbol{\beta}_A^0 + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B^0)^T \mathbf{X}_{Bij} \text{ such that the } \mathbf{U}_i \text{ are independent } N(\mathbf{0}, \boldsymbol{\Sigma}^0) \text{ random vectors.} \quad (4.3)$$

The  $\mathbf{U}_i$  are  $d_A \times 1$  unobserved random vectors for each  $1 \leq i \leq m$ . The  $\mathbf{X}_{Aij}$  are  $d_A \times 1$  random vectors corresponding to predictors that are partnered by both a fixed effect and a random effect. The  $\mathbf{X}_{Bij}$  are  $d_B \times 1$  random vectors which are predictors that have a fixed effect only. Let  $\mathbf{X}_{ij} \equiv (\mathbf{X}_{Aij}^T, \mathbf{X}_{Bij}^T)^T$  denote the combined predictor vectors such that  $d_A + d_B = d$ . We also assumed that the  $\mathbf{X}_{ij}$  and  $\mathbf{U}_i$ , for  $1 \leq i \leq m$  and  $1 \leq j \leq n_i$ , are independent. The  $\mathbf{X}_{ij}$  are each assumed as having the same distribution as the  $(d_A + d_B) \times 1$  random vector  $\mathbf{X} = (\mathbf{X}_A^T, \mathbf{X}_B^T)^T$  and the  $\mathbf{U}_i$  are assumed to be independent  $N(\mathbf{0}, \boldsymbol{\Sigma}^0)$ , each having the same distribution as the random vector  $\mathbf{U}$ .

Then, for any  $\beta_A(d_A \times 1)$ ,  $\beta_B(d_B \times 1)$  and  $\Sigma(d_A \times d_A)$  that is symmetric, positive definite and conditional on the  $\mathbf{X}_{ij}$  data, the maximum likelihood estimator of  $(\beta_A^0, \beta_B^0, \Sigma^0)$  is,

$$(\widehat{\beta}_A, \widehat{\beta}_B, \widehat{\Sigma}) = \underset{\beta_A, \beta_B, \Sigma}{\operatorname{argmax}} \ell(\beta_A, \beta_B, \Sigma)$$

where  $\ell(\beta_A, \beta_B, \Sigma)$  is the conditional log-likelihood and has the expression

$$\begin{aligned} \ell(\beta_A, \beta_B, \Sigma) &= \sum_{i=1}^m \sum_{j=1}^n [\{Y_{ij} (\beta_A^T \mathbf{X}_{Aij} + \beta_B^T \mathbf{X}_{Bij}) + c(Y_{ij})\} / \phi + d(Y_{ij}, \phi)] - \frac{m}{2} \log |2\pi \Sigma| \\ &\quad + \sum_{i=1}^m \log \int_{\mathbb{R}^{d_A}} \exp \left[ \sum_{j=1}^n \{Y_{ij} \mathbf{u}^T \mathbf{X}_{Aij} - b((\beta_A + \mathbf{u})^T \mathbf{X}_{Aij} + \beta_B^T \mathbf{X}_{Bij})\} / \phi \right. \\ &\quad \left. - \frac{1}{2} \mathbf{u}^T \Sigma^{-1} \mathbf{u} \right] d\mathbf{u}. \end{aligned}$$

Note that one-parameter exponential family densities have a variety of desirable properties, so that the regularity conditions for Theorem 7 are met. For example, in the class of one-parameter exponential families, the support of  $f(\mathbf{y}|\boldsymbol{\theta})$  does not depend on  $\boldsymbol{\theta}$ . Additional assumptions, for instance, having  $\Sigma^0$  being positive definite, ensures that the true values for all model parameters are interior to the parameter space. Thus, the model description above ensures that the regularity conditions required of the density function is met for the convergence in distribution result for MLEs. Similar explanations apply for the model descriptions in Chapters 6 and 7 as well.

## 4.2 Notation

Define

$$n \equiv \frac{1}{m} \sum_{i=1}^m n_i = \text{average of the within-group sample sizes,}$$

$$\boldsymbol{\Omega}_{\beta_B}(\mathbf{U}) \equiv E \left\{ b''((\beta_A^0 + \mathbf{U})^T \mathbf{X}_A + (\beta_B^0)^T \mathbf{X}_B) \begin{bmatrix} \mathbf{X}_A \mathbf{X}_A^T & \mathbf{X}_A \mathbf{X}_B^T \\ \mathbf{X}_B \mathbf{X}_A^T & \mathbf{X}_B \mathbf{X}_B^T \end{bmatrix} \middle| \mathbf{U} \right\}$$

and

$$\boldsymbol{\Lambda}_{\beta_B} \equiv \left( E \left[ \left\{ \text{lower right } d_B \times d_B \text{ block of } \boldsymbol{\Omega}_{\beta_B}(\mathbf{U})^{-1} \right\}^{-1} \right] \right)^{-1}.$$

### 4.3 Asymptotic Normality Theorem

The main theoretical contribution of this chapter is an asymptotic normality theorem for the maximum quasi-likelihood estimators for a generalized linear mixed model as described in Section 4.1.

The theorem relies on the following assumptions:

- (A4) The number of groups  $m$  diverges to  $\infty$ .
- (A5) The within-group sample sizes  $n_i$  diverge to  $\infty$  in such a way that  $n_i/n \rightarrow C_i$  for constants  $0 < C_i < \infty$ ,  $1 \leq i \leq m$ . Also,  $n/m \rightarrow 0$  as  $m$  and  $n$  diverge.
- (A6) The distribution of  $\mathbf{X}$  is such that

$$E \left[ \frac{E [\max\{1, \|\mathbf{X}\|\}^8 \max\{1, b''((\beta_A + \mathbf{U})^T \mathbf{X}_A + \beta_B^T \mathbf{X}_B)\}^4 | \mathbf{U}]}{\min\{1, \lambda_{\min}(E\{\mathbf{X}_A \mathbf{X}_A^T b''((\beta_A + \mathbf{U})^T \mathbf{X}_A + \beta_B^T \mathbf{X}_B) | \mathbf{U})\})^2} \right] < \infty$$

for all  $\beta_A \in \mathbb{R}^{d_A}$ ,  $\beta_B \in \mathbb{R}^{d_B}$  and  $\Sigma$  that is a  $d_A \times d_A$  symmetric and positive definite matrix.

**Theorem 12.** *Assume that conditions (A4) - (A6) hold. Then we have the following*

$$\sqrt{m} \begin{bmatrix} \hat{\beta}_A - \beta_A^0 \\ \sqrt{n} (\hat{\beta}_B - \beta_B^0) \\ \text{vech}(\hat{\Sigma}) - \text{vech}(\Sigma^0) \end{bmatrix} \xrightarrow{\mathcal{D}} N \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma^0 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \phi \Lambda_{\beta_B} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & 2D_{d_A}^+ (\Sigma^0 \otimes \Sigma^0) D_{d_A}^{+T} \end{bmatrix} \right).$$

Proof of Theorem 12 is in the appendices. Some remarks concerning Theorem 12 are:

1. Firstly, note that the asymptotic variances of the maximum quasi-likelihood estimator of the fixed effects that are accompanied by random effects and the maximum quasi-likelihood estimator of the variance and covariance parameters of the random effects, both have a convergence rate of  $m^{-1}$ . On the other hand, the asymptotic variance of the estimator of the fixed effects unaccompanied by random effects has a much faster convergence rate of  $(mn)^{-1}$ .
2. The off-diagonal entries of the asymptotic variance-covariance matrix of the maximum quasi-likelihood estimators of the model parameters are zero matrices. Hence, this implies asymptotic orthogonality between the  $\beta_A$ ,  $\beta_B$  and  $\Sigma$  model parameters.

3. In existing literature, results such as Theorem 3 of Nie (2007) contains limits and expectations over the response distribution in their Fisher information approximations. In Theorem 12, we avoid such limits and expectations over the response distribution. The  $\Lambda_{\beta_B}$  matrix only involves expectations over the simpler random effects distribution. Therefore the results in this thesis, including Theorem 12, provide expressions that are easy to implement for practical tasks such as confidence interval construction, making them “usable” in practice compared to other theoretical asymptotic normality results (rather than results that make use of observed Fisher information) available for maximum likelihood estimators for generalized linear mixed models.

## 4.4 Dispersion Parameter Extension

In Theorem 12, we treat the dispersion parameter  $\phi$  as being fixed. When considering the Gaussian and Gamma response cases, all of the model parameters in (4.2), including  $\phi$ , can be estimated using ordinary maximum likelihood. Theorem 12 can then be extended for maximum likelihood estimation of  $\phi^0$  and involves the addition of

$$\sqrt{mn} \left( \hat{\phi} - \phi^0 \right) \xrightarrow{D} N(0, v(\phi^0))$$

where  $v(x) \equiv 2x^2$  for Gaussian responses and  $v(x) \equiv x^4 / \{\text{trigamma}(1/x) - x\}$ . Proof of the asymptotic variance expression for the maximum likelihood estimator of the dispersion parameter in the Gaussian response case is in Chapter 3 and proof of the asymptotic variance expression for the maximum likelihood estimator of the dispersion parameter in the Gamma response case is in the appendix for this chapter.

Note that from our results, for the Gaussian and gamma response cases, we can conclude that exact orthogonality exists between  $\phi$  and  $(\beta_A, \beta_B)$  and asymptotic orthogonality exists between  $\phi$  and  $\Sigma$ . Thus, the covariance matrices of Theorem 12 still hold for  $\hat{\beta}_A, \hat{\beta}_B$  and  $\text{vech}(\hat{\Sigma})$ . On the other hand, when implementing the quasi-likelihood extension of the Binomial and Poisson response cases,  $\phi$  cannot be estimated via maximum quasi-likelihood and is typically estimated via a method of moments approach. Note that the values of the maximum quasi-likelihood estimates of  $\beta_A, \beta_B$  and  $\Sigma$  asymptotically do not depend on  $\phi$ . This can be deduced from the likelihood equations formed using the first-order asymptotic approximations of the scores that only have  $\phi$  as a constant. Hence, Theorem 12 is unaffected by the estimation of  $\phi$  for the response cases under the umbrella of one-parameter exponential families too.

## 4.5 Appendix

### 4.5.1 Multivariate Extension of (2.6) of Tierney et al. (1989)

To carry out the derivations in the next appendix, one has to deal with solving ratios of intractable integrals. In this appendix, we show how to deal with such ratios of intractable integrals containing  $d$ -variate arguments by working with its equivalent multi-term Laplace's method expansion instead. This is accomplished by considering the multivariate extension of (2.6) of Tierney et al. (1989), which follows from results in Appendix A of Miyata (2004).

#### 4.5.1.1 Overview

For smooth real-valued functions  $g$ ,  $c$  and  $h$ , Equation (2.6) of Tierney et al. (1989) states that

$$\frac{\int_{-\infty}^{\infty} b_N(x) \exp\{-nh(x)\} dx}{\int_{-\infty}^{\infty} b_D(x) \exp\{-nh(x)\} dx} = g(x^*) + \frac{b'_D(x^*)g'(x^*)}{nb_D(x^*)h''(x^*)} + \frac{g''(x^*)}{2nh''(x^*)} - \frac{g'(x^*)h'''(x^*)}{2nh''(x^*)^2} + O(n^{-2}). \quad (4.4)$$

where

$$g \equiv b_N/b_D \quad (4.5)$$

and

$$x^* \equiv \text{value of } x \text{ that minimises } h \text{ over } \mathbb{R}.$$

Now consider the first two equations in Appendix A of Miyata (2004). Suppose that in the right-hand side of the first equation we set  $\Theta = \mathbb{R}^d$ , replace the  $\rho$  symbol by  $b_D$ , replace the  $h_n$  symbol by  $h$  and replace the integral dummy variable  $\theta$  by  $\mathbf{x}$ . Then we get

$$\frac{\int_{\mathbb{R}^d} g(\mathbf{x}) b_D(\mathbf{x}) \exp\{-nh(\mathbf{x})\} d\mathbf{x}}{\int_{\mathbb{R}^d} b_D(\mathbf{x}) \exp\{-nh(\mathbf{x})\} d\mathbf{x}}.$$

If the function  $b_N$  is defined according to (7.3.6), then this quantity becomes

$$\frac{\int_{\mathbb{R}^d} b_N(\mathbf{x}) \exp\{-nh(\mathbf{x})\} d\mathbf{x}}{\int_{\mathbb{R}^d} b_D(\mathbf{x}) \exp\{-nh(\mathbf{x})\} d\mathbf{x}},$$

which is equivalent to both the right-hand sides of the first and second equations in Appendix A of Miyata (2004). So the asymptotic expansion is given by the right-hand side of the second equation in Appendix A of Miyata (2004) with  $\hat{\theta}$  replaced by  $\mathbf{x}^*$  and  $\rho$  replaced  $b_D$ .

### 4.5.1.2 Multivariate Derivative Notation

Let  $f$  be a smooth real-valued  $d$ -variate function of the  $d$ -variate argument  $\mathbf{x} \equiv (x_1, \dots, x_d)$ . The gradient vector of  $f$  is

$$\nabla f(x_1, \dots, x_d) = d \times 1 \text{ vector with } i\text{th entry } \frac{\partial f(x_1, \dots, x_d)}{\partial x_i}.$$

The Hessian matrix of  $f$  is

$$\nabla^2 f(x_1, \dots, x_d) = d \times d \text{ matrix with } (i, j)\text{th entry } \frac{\partial^2 f(x_1, \dots, x_d)}{\partial x_i \partial x_j}.$$

The third derivatives three-dimensional array of  $f$  is

$$\nabla^3 f(x_1, \dots, x_d) = d \times d \times d \text{ array with } (i, j, k)\text{th entry } \frac{\partial^3 f(x_1, \dots, x_d)}{\partial x_i \partial x_j \partial x_k}.$$

### 4.5.1.3 Check of the Miyata (2004) Appendix A Result for the Univariate Case

We first make sure that the right-hand side of the second equation in Appendix A of Miyata (2004) matches (2.6) of Tierney et al. (1989) when  $d = 1$ . Recall that we have to replace  $\rho$  by  $b_D$ . We first deal with the term that is the one just before the  $O(n^{-2})$  term. The author points out that this term is, in fact (with some trivial re-arrangement),

$$\frac{1}{2n} \text{tr}[\{\nabla^2 h(\mathbf{x}^*)\}^{-1} \nabla^2 g(\mathbf{x}^*)]$$

In the univariate case this becomes

$$\frac{g''(x^*)}{2n h''(x^*)}$$

which is one of the terms in (7.3.6).

Consider the second term involving the  $\sum_{ij}$  symbol. Next note that Miyata (2004) uses  $h^{ij}$  to denote the components of  $\{\nabla^2 h(\mathbf{x}^*)\}^{-1}$ . In the univariate case this is simply  $1/h''(x^*)$ . In this univariate case the summation over  $ij$  collapses to a scalar and the first component is

$$\frac{1}{n} g'(x^*) \{1/h''(x^*)\} b'_D(x^*)/b_D(x^*) = \frac{b'_D(x^*) g'(x^*)}{n b_D(x^*) h''(x^*)},$$

which is also one of the terms in (7.3.6).

It remains to show that the second component of the main  $\sum_{ij}$  expression reduces to

$$-\frac{g'(x^*) h'''(x^*)}{2n h''(x^*)^2}.$$

For a general set of model parameters  $\boldsymbol{\theta}$ , let  $h_{j_1 \dots j_d}(\hat{\boldsymbol{\theta}})$  denote the  $d$ th partial derivative  $\partial^d h_n(\boldsymbol{\theta}) / \partial \theta_{j_1} \dots \partial \theta_{j_d}$  with respect to  $\boldsymbol{\theta}$  evaluated  $\hat{\boldsymbol{\theta}}$ . Then from this definition of  $h_{j_1 \dots j_d}$  provided in Miyata (2004), it is apparent that, in the  $d = 1$  case, we can replace  $h_{rsj}$  by  $h_{111}$  and then set it to  $h'''(x^*)$ . Also, in this  $d = 1$  case the summation of  $rs$  collapses to a scalar. Combining we get

$$\frac{1}{n} g'(x^*) \{1/h''(x^*)\} (-\frac{1}{2}) \{1/h''(x^*)\} h'''(x^*) = -\frac{g'(x^*) h'''(x^*)}{2n h''(x^*)^2}$$

as required.

In summary, the result in Appendix A of Miyata (2004) does reduce to (2.6) of Tierney et al. (1989) in the  $d = 1$  case.

#### 4.5.1.4 The Multivariate Case

The second equation of Appendix A of Miyata (2004) gives an expansion for

$$\frac{\int_{\mathbb{R}^d} g(\mathbf{x}) b_D(\mathbf{x}) \exp\{-nh(\mathbf{x})\} d\mathbf{x}}{\int_{\mathbb{R}^d} b_D(\mathbf{x}) \exp\{-nh(\mathbf{x})\} d\mathbf{x}}.$$

It is relatively easy to show that

$$\begin{aligned} \frac{\int_{\mathbb{R}^d} g(\mathbf{x}) b_D(\mathbf{x}) \exp\{-nh(\mathbf{x})\} d\mathbf{x}}{\int_{\mathbb{R}^d} b_D(\mathbf{x}) \exp\{-nh(\mathbf{x})\} d\mathbf{x}} &= g(\mathbf{x}^*) + \frac{\nabla g(\mathbf{x}^*)^T \{\nabla^2 h(\mathbf{x}^*)\}^{-1} \nabla b_D(\mathbf{x}^*)}{n b_D(\mathbf{x}^*)} \\ &\quad + \frac{\text{tr}[\{\nabla^2 h(\mathbf{x}^*)\}^{-1} \nabla^2 g(\mathbf{x}^*)]}{2n} + \Upsilon(\mathbf{x}^*) + O(n^{-2}) \end{aligned}$$

for smooth real-valued  $d$ -variate functions  $g$ ,  $c$  and  $h$ , where  $\Upsilon(\mathbf{x}^*)$  denotes the term involving third order derivatives and therefore is more challenging when it comes to getting succinct matrix algebraic expressions.

In terms of the subscript derivative notation used in Miyata (2004):

$$\Upsilon(\mathbf{x}^*) = -\frac{1}{2n} \sum_{i=1}^d \sum_{j=1}^d g_i h^{ij} \sum_{k=1}^d \sum_{\ell=1}^d h^{k\ell} h_{k\ell j}$$

where

$$g_i = i\text{th entry of } \nabla g(\mathbf{x}^*),$$

$$h^{ij} = (i, j) \text{ entry of } \{\nabla^2 h(\mathbf{x}^*)\}^{-1}$$

$$\text{and } h_{ijk} = (i, j, k) \text{ entry of } \nabla^3 h(\mathbf{x}^*).$$

Note that if  $\boldsymbol{\omega}(\mathbf{x}^*)$  is the  $d \times 1$  vector with  $k$ th entry equal to

$$\sum_{i=1}^d \sum_{j=1}^d h^{ij} h_{ijk} = \sum_{i=1}^d \sum_{j=1}^d [\{\nabla^2 h(\mathbf{x}^*)\}^{-1}]_{ij} \{\nabla^3 h(\mathbf{x}^*)\}_{ijk}$$

then we have

$$\Upsilon(\mathbf{x}^*) = -\frac{\nabla g(\mathbf{x}^*)^T \{\nabla^2 h(\mathbf{x}^*)\}^{-1} \boldsymbol{\omega}(\mathbf{x}^*)}{2n}.$$

We could also define

$$\nabla^3 h(\mathbf{x})_{[k]} \equiv d \times d \text{ matrix with } (i, j) \text{ entry equal to the } (i, j, k) \text{ entry of } \nabla^3 h(\mathbf{x})$$

and then the  $k$ th entry of  $\boldsymbol{\omega}(\mathbf{x}^*)$  is

$$\text{tr}[\{\nabla^2 h(\mathbf{x}^*)\}^{-1} \nabla^3 h(\mathbf{x}^*)_{[k]}].$$

#### 4.5.1.5 Final Expression for the Multivariate Extension of (2.6) of Tierney et al. (1989)

Putting everything together from the previous subsection, the multivariate extension of (2.6) of Tierney et al. (1989) is:

$$\begin{aligned} \frac{\int_{\mathbb{R}^d} g(\mathbf{x}) b_D(\mathbf{x}) \exp\{-nh(\mathbf{x})\} d\mathbf{x}}{\int_{\mathbb{R}^d} b_D(\mathbf{x}) \exp\{-nh(\mathbf{x})\} d\mathbf{x}} &= g(\mathbf{x}^*) + \frac{\nabla g(\mathbf{x}^*)^T \{\nabla^2 h(\mathbf{x}^*)\}^{-1} \nabla b_D(\mathbf{x}^*)}{nb_D(\mathbf{x}^*)} \\ &+ \frac{\text{tr}[\{\nabla^2 h(\mathbf{x}^*)\}^{-1} \nabla^2 g(\mathbf{x}^*)]}{2n} \\ &- \frac{\nabla g(\mathbf{x}^*)^T \{\nabla^2 h(\mathbf{x}^*)\}^{-1} \boldsymbol{\omega}(\mathbf{x}^*)}{2n} + O(n^{-2}) \end{aligned} \quad (4.6)$$

where

$$\boldsymbol{\omega}(\mathbf{x}) \equiv \begin{bmatrix} \text{tr}[\{\nabla^2 h(\mathbf{x})\}^{-1} \nabla^3 h(\mathbf{x})_{[1]}] \\ \vdots \\ \text{tr}[\{\nabla^2 h(\mathbf{x})\}^{-1} \nabla^3 h(\mathbf{x})_{[d]}] \end{bmatrix}$$

and

$$\nabla^3 h(\mathbf{x})_{[k]} \equiv d \times d \text{ matrix with } (i, j) \text{ entry equal to the } (i, j, k) \text{ entry of } \nabla^3 h(\mathbf{x}).$$

## 4.5.2 Proof of Theorem 12

This appendix contains the details for the derivations leading up to Theorem 12.

### 4.5.2.1 Constructing the Fisher Information Matrix

In order to compute the asymptotic covariance matrix for the maximum quasi-likelihood estimators, we would first need to compute the Fisher information matrix for the model



parameters as per the model description in (4.3). To do so, let

$$\mathbf{S}_i \equiv \begin{bmatrix} \mathbf{S}_{Ai} \\ \mathbf{S}_{Bi} \\ \mathbf{S}_{Ci} \end{bmatrix} = \begin{bmatrix} \nabla_{\beta_A} \log p_{\mathbf{Y}_i|\mathbf{X}_i}(\mathbf{Y}_i|\mathbf{X}_i) \\ \nabla_{\beta_B} \log p_{\mathbf{Y}_i|\mathbf{X}_i}(\mathbf{Y}_i|\mathbf{X}_i) \\ \nabla_{\text{vech}(\boldsymbol{\Sigma})} \log p_{\mathbf{Y}_i|\mathbf{X}_i}(\mathbf{Y}_i|\mathbf{X}_i) \end{bmatrix} \quad (4.7)$$

denote the  $i$ th contribution to the scores for each of the model parameters. Then the Fisher information matrix can be computed as

$$I(\beta_A, \beta_B, \text{vech}(\boldsymbol{\Sigma})) = \sum_{i=1}^m E(\mathbf{S}_i \mathbf{S}_i^T | \mathbf{X}_i).$$

The next few sections then focus on obtaining the expressions for the scores and the quadratic conditional expectations that are required to construct the final Fisher information matrix.

#### 4.5.2.2 Expression for Conditional Density Function

The expression for  $p_{\mathbf{Y}_i|\mathbf{X}_i}(\mathbf{Y}_i|\mathbf{X}_i)$  as per the model description in (4.3) is as follows

$$\begin{aligned} & p_{\mathbf{Y}_i|\mathbf{X}_i}(\mathbf{Y}_i|\mathbf{X}_i) \\ &= \int_{\mathbb{R}^{d_A}} \prod_{j=1}^{n_i} \{p(Y_{ij}|\mathbf{X}_{Aij}, \mathbf{X}_{Bij}, \mathbf{U}_i)\} p(\mathbf{U}_i) d\mathbf{U}_i \\ &= \int_{\mathbb{R}^{d_A}} \exp \left\{ \sum_{j=1}^n \left( \left[ Y_{ij} \left\{ (\beta_A + \mathbf{u})^T \mathbf{X}_{Aij} + \beta_B^T \mathbf{X}_{Bij} \right\} - b \left( (\beta_A + \mathbf{u})^T \mathbf{X}_{Aij} + \beta_B^T \mathbf{X}_{Bij} \right) \right. \right. \right. \\ &\quad \left. \left. \left. + c(Y_{ij}) \right] / \phi + d(Y_{ij}, \phi) \right) \right\} \times (2\pi)^{-d_A/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) d\mathbf{u} \\ &= \int_{\mathbb{R}^{d_A}} |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left\{ \sum_{j=1}^n \left( \left[ Y_{ij} \left\{ (\beta_A + \mathbf{u})^T \mathbf{X}_{Aij} + \beta_B^T \mathbf{X}_{Bij} \right\} \right. \right. \right. \\ &\quad \left. \left. \left. - b \left( (\beta_A + \mathbf{u})^T \mathbf{X}_{Aij} + \beta_B^T \mathbf{X}_{Bij} \right) + c(Y_{ij}) \right] / \phi + d(Y_{ij}, \phi) \right) - \frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} \right\} d\mathbf{u} \\ &= |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left( \sum_{j=1}^n \left[ \{ Y_{ij} (\beta_A^T \mathbf{X}_{Aij} + \beta_B^T \mathbf{X}_{Bij}) + c(Y_{ij}) \} / \phi + d(Y_{ij}, \phi) \right] \right) \\ &\quad \times \int_{\mathbb{R}^{d_A}} \exp \left[ \sum_{j=1}^n \left\{ Y_{ij} \mathbf{u}^T \mathbf{X}_{Aij} - b \left( (\beta_A + \mathbf{u})^T \mathbf{X}_{Aij} + \beta_B^T \mathbf{X}_{Bij} \right) \right\} / \phi - \frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} \right] d\mathbf{u}. \end{aligned}$$

### 4.5.2.3 Introduction of Useful Notation and its Properties

Here we introduce expressions that will be useful in summarising the derivations throughout this appendix. Firstly,

$$\mathbf{1}_{d_A^{\boxplus}} \equiv \mathbf{1}_{d_A(d_A+1)/2}$$

and

$$\mathbf{1}_{d_A}^{\otimes 2} \equiv \mathbf{1}_{d_A} \mathbf{1}_{d_A}^T, \quad \mathbf{1}_{d_B}^{\otimes 2} \equiv \mathbf{1}_{d_B} \mathbf{1}_{d_B}^T \quad \text{and} \quad \mathbf{1}_{d_A^{\boxplus}}^{\otimes 2} \equiv \mathbf{1}_{d_A^{\boxplus}} \mathbf{1}_{d_A^{\boxplus}}^T.$$

Then note that

$$\mathcal{G}_{Ai} \equiv \sum_{j=1}^n \{Y_{ij} - b'((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Aij},$$

$$\mathcal{G}_{Bi} \equiv \sum_{j=1}^n \{Y_{ij} - b'((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Bij},$$

$$\mathcal{H}_{AAi} \equiv \sum_{j=1}^n b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T,$$

$$\mathcal{H}_{ABi} = \sum_{j=1}^n b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \mathbf{X}_{Aij} \mathbf{X}_{Bij}^T,$$

$$\mathcal{H}_{BBi} \equiv \sum_{j=1}^n b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \mathbf{X}_{Bij} \mathbf{X}_{Bij}^T.$$

$\mathcal{H}'_{AAAi}$  is the  $d_A \times d_A \times d_A$  array with  $(r, s, t)$  entry equal to

$$\sum_{j=1}^n b'''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) (\mathbf{X}_{Aij})_r (\mathbf{X}_{Aij})_s (\mathbf{X}_{Aij})_t,$$

and  $\mathcal{H}'_{AAAi[t]}$  is the  $d_A \times d_A$  matrix with  $(r, s)$  entry equal to the  $(r, s, t)$  entry of  $\mathcal{H}'_{AAAi}$ . Note that the expressions listed above have the following probabilistic orders where

$$\mathcal{G}_{Ai} = O_P(n^{1/2}) \mathbf{1}_{d_A}, \quad \mathcal{G}_{Bi} = O_P(n^{1/2}) \mathbf{1}_{d_B},$$

$$\mathcal{H}_{AAi} = O_P(n) \mathbf{1}_{d_A}^{\otimes 2}, \quad \mathcal{H}_{ABi} = O_P(n) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T \quad \text{and} \quad \mathcal{H}_{BBi} = O_P(n) \mathbf{1}_{d_B}^{\otimes 2}.$$

The Conditional Expectation of  $\mathcal{G}_{Ai}$  and  $\mathcal{G}_{Bi}$  Given  $(\mathbf{X}_i, \mathbf{U}_i)$

$$\begin{aligned} E(\mathcal{G}_{Ai}|\mathbf{X}_i, \mathbf{U}_i) &= E\left(\left[\sum_{j=1}^n \left\{Y_{ij} - b' \left((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij}\right)\right\} \mathbf{X}_{Aij}\right] \middle| \mathbf{X}_i, \mathbf{U}_i\right) \\ &= \sum_{j=1}^n \left\{E(Y_{ij}|\mathbf{X}_i, \mathbf{U}_i) - b' \left((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij}\right)\right\} \mathbf{X}_{Aij} \\ &= \mathbf{0}. \end{aligned}$$

Similarly we also have that

$$E(\mathcal{G}_{Bi}|\mathbf{X}_i, \mathbf{U}_i) = \mathbf{0}.$$

The Conditional Expectations of  $\mathcal{G}_{Ai} \mathcal{G}_{Ai}^T$ ,  $\mathcal{G}_{Ai} \mathcal{G}_{Bi}^T$  and  $\mathcal{G}_{Bi} \mathcal{G}_{Bi}^T$  Given  $(\mathbf{X}_i, \mathbf{U}_i)$

$$\begin{aligned} &E(\mathcal{G}_{Ai} \mathcal{G}_{Ai}^T | \mathbf{X}_i, \mathbf{U}_i) \\ &= E\left(\left[\sum_{j=1}^n \left\{Y_{ij} - b' \left((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij}\right)\right\} \mathbf{X}_{Aij}\right] \right. \\ &\quad \left. \times \left[\sum_{j=1}^n \left\{Y_{ij'} - b' \left((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij'} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij'}\right)\right\} \mathbf{X}_{Aij'}\right]^T \middle| \mathbf{X}_i, \mathbf{U}_i\right) \\ &= \sum_{j \neq j'} E\left(\left[\left\{Y_{ij} - b' \left((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij}\right)\right\} \mathbf{X}_{Aij}\right] \right. \\ &\quad \left. \times \left[\left\{Y_{ij'} - b' \left((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij'} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij'}\right)\right\} \mathbf{X}_{Aij'}\right]^T \middle| \mathbf{X}_i, \mathbf{U}_i\right) \\ &\quad + \sum_{j=1}^n E\left(\left[\left\{Y_{ij} - b' \left((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij}\right)\right\} \mathbf{X}_{Aij}\right] \right. \\ &\quad \left. \times \left[\left\{Y_{ij} - b' \left((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij}\right)\right\} \mathbf{X}_{Aij}\right]^T \middle| \mathbf{X}_i, \mathbf{U}_i\right) \\ &= \sum_{j \neq j'} E\left\{\left[\left\{Y_{ij} - b' \left((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij}\right)\right\} \mathbf{X}_{Aij} \middle| \mathbf{X}_i, \mathbf{U}_i\right] \right. \\ &\quad \left. \times E\left(\left[\left\{Y_{ij'} - b' \left((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij'} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij'}\right)\right\} \mathbf{X}_{Aij'}\right]^T \middle| \mathbf{X}_i, \mathbf{U}_i\right)\right\} \\ &\quad + \sum_{j=1}^n \text{Var}(Y_{ij} \mathbf{X}_{Aij} | \mathbf{X}_i, \mathbf{U}_i) \end{aligned}$$

The equation above simplifies to

$$\begin{aligned}
& \sum_{j \neq j'} \sum \left[ \left\{ E(Y_{ij} | \mathbf{X}_i, \mathbf{U}_i) - b' \left( (\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij} \right) \right\} \mathbf{X}_{Aij} \right] \\
& \quad \times \left[ \left\{ E(Y_{ij'} | \mathbf{X}_i, \mathbf{U}_i) - b' \left( (\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij'} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij'} \right) \right\} \mathbf{X}_{Aij'} \right]^T \\
& \quad + \sum_{j=1}^n \mathbf{X}_{Aij} \text{Var}(Y_{ij} | \mathbf{X}_i, \mathbf{U}_i) \mathbf{X}_{Aij}^T \\
& = \phi \sum_{j=1}^n b'' \left( (\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij} \right) \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T \\
& = \phi \mathcal{H}_{AAi}.
\end{aligned}$$

Similarly we also have that

$$E(\mathcal{G}_{Ai} \mathcal{G}_{Bi}^T | \mathbf{X}_i, \mathbf{U}_i) = \phi \mathcal{H}_{ABi} \text{ and } E(\mathcal{G}_{Bi} \mathcal{G}_{Bi}^T | \mathbf{X}_i, \mathbf{U}_i) = \phi \mathcal{H}_{BBi}.$$

#### 4.5.2.4 Computing an Asymptotic Approximation for the First Entry in (4.7)

To overcome the intractability of the ratio of integrals present when deriving the scores with respect to each of the model parameters, we will work with an asymptotic approximation of the ratio of integrals by using a multi-term Laplace's method expansion as described in the appendix in Subsection 4.5.1. Hence, the  $i$ th contribution to the score of  $\boldsymbol{\beta}_A$  is

$$\begin{aligned}
\mathbf{S}_{Ai} &= \nabla_{\boldsymbol{\beta}_A} \log p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i) \\
&= \frac{\int_{\mathbb{R}^{d_A}} \mathbf{b}_N^{\text{1st}}(\mathbf{u}) \exp\{-nh_N(\mathbf{u})\} d\mathbf{u}}{\int_{\mathbb{R}^{d_A}} b_D^{\text{1st}}(\mathbf{u}) \exp\{-nh_N(\mathbf{u})\} d\mathbf{u}}
\end{aligned}$$

where

$$\begin{aligned}
\mathbf{b}_N^{\text{1st}}(\mathbf{u}) &\equiv \exp\left(-\frac{1}{2}\mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) \frac{1}{\phi} \sum_{j=1}^n \left\{ Y_{ij} - b' \left( (\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij} \right) \right\} \mathbf{X}_{Aij}, \\
b_D^{\text{1st}}(\mathbf{u}) &\equiv \exp\left(-\frac{1}{2}\mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) \quad \text{and} \\
h_N(\mathbf{u}) &\equiv -\frac{1}{n\phi} \sum_{j=1}^n \left\{ Y_{ij} \mathbf{u}^T \mathbf{X}_{Aij} - b \left( (\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij} \right) \right\}.
\end{aligned}$$

Now define

$$\begin{aligned}
\mathbf{U}_i^* &\equiv \text{value of } \mathbf{u} \text{ that minimises } h_N(\mathbf{u}) \\
&= \text{value of } \mathbf{u} \text{ such that } \nabla_{\mathbf{u}} h_N(\mathbf{u}) = \mathbf{0} \\
&= \text{value of } \mathbf{u} \text{ such that } \sum_{j=1}^n \left\{ Y_{ij} - b' \left( (\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij} \right) \right\} \mathbf{X}_{Aij} = \mathbf{0}.
\end{aligned}$$

However,

$$\mathbf{b}_N^{1st}(\mathbf{U}_i^*) = \mathbf{0}.$$

This violates the condition in Hsu (1948), in the sense that  $\mathbf{b}_N^{1st}(\mathbf{U}_i^*) \neq \mathbf{0}$  is required in order for the Laplace approximation to hold. To counter this issue, firstly note that the numerator of  $\mathbf{S}_{Ai}$  is

$$\int_{\mathbb{R}^{d_A}} \{\nabla_{\mathbf{u}} s(\mathbf{u})\} t(\mathbf{u}) d\mathbf{u} \quad (4.8)$$

where

$$s(\mathbf{u}) \equiv \exp \left( \frac{n}{\phi} \left[ \frac{1}{n} \sum_{j=1}^n \left\{ Y_{ij} \mathbf{u}^T \mathbf{X}_{Aij} - b((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \right\} \right] \right)$$

and

$$t(\mathbf{u}) \equiv \exp \left( -\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} \right).$$

Using an  $\mathbb{R}^{d_A}$  extension of integration by parts, we can replace (4.8) by

$$- \int_{\mathbb{R}^{d_A}} s(\mathbf{u}) \{\nabla_{\mathbf{u}} t(\mathbf{u})\} d\mathbf{u}.$$

We then obtain that

$$\begin{aligned} dt(\mathbf{u}) &= \exp \left( -\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} \right) d \left\{ -\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} \right\} \\ &= - \exp \left( -\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} \right) \frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} d\mathbf{u}, \end{aligned}$$

which leads to

$$\nabla_{\mathbf{u}} t(\mathbf{u}) = - \exp \left( -\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} \right) \boldsymbol{\Sigma}^{-1} \mathbf{u}.$$

Now by rewriting the numerator of  $\mathbf{S}_{Ai}$ , we have,

$$\mathbf{S}_{Ai} = \frac{\int_{\mathbb{R}^{d_A}} \mathbf{b}_N(\mathbf{u}) \exp\{-nh_N(\mathbf{u})\} d\mathbf{u}}{\int_{\mathbb{R}^{d_A}} b_D(\mathbf{u}) \exp\{-nh_N(\mathbf{u})\} d\mathbf{u}}$$

where

$$\begin{aligned} \mathbf{b}_N(\mathbf{u}) &\equiv \exp \left( -\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} \right) \boldsymbol{\Sigma}^{-1} \mathbf{u} \\ b_D(\mathbf{u}) &\equiv \exp \left( -\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} \right) \quad \text{and} \\ h_N(\mathbf{u}) &\equiv -\frac{1}{n\phi} \sum_{j=1}^n \left\{ Y_{ij} \mathbf{u}^T \mathbf{X}_{Aij} - b((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \right\}. \end{aligned}$$

Expansion of  $\mathbf{U}_i^*$ 

Here we find an asymptotic expression for  $\mathbf{U}_i^*$ . We have that

$$\begin{aligned}
0 &= \nabla_{\mathbf{u}} h_N(\mathbf{u}) \\
&= \sum_{j=1}^n \{Y_{ij} - b'((\boldsymbol{\beta}_A + \mathbf{U}_i^*)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Aij} \\
&= \sum_{j=1}^n \{Y_{ij} - b'((\boldsymbol{\beta}_A^0 + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B^0)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Aij} \\
&\quad - \sum_{j=1}^n b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T (\mathbf{U}_i^* - \mathbf{U}_i) + r_{it} \\
&= \mathcal{G}_{Ai} - \mathcal{H}_{AAi}(\mathbf{U}_i^* - \mathbf{U}_i) + r_{it}
\end{aligned}$$

where  $r_{it}$  is the Lagrange form of the remainder and is a quadratic form in  $\mathbf{U}_i^* - \mathbf{U}_i$  and a smooth function of  $\mathbf{U}_{it}^\dagger \equiv (1-t)\mathbf{U}_i + t\mathbf{U}_i^*$  for some  $t \in [0, 1]$ . Inversion of this asymptotic series leads to

$$\mathbf{U}_i^* = \mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} + O_P(n^{-1}) \mathbf{1}_{d_A}.$$

The  $\mathbf{e}_k$  Notation

The following notation is useful for the upcoming calculations. For each  $1 \leq k_A \leq d_A$ ,  $1 \leq k_B \leq d_B$ ,  $1 \leq k_C \leq d_A^\boxplus$ , let

$\mathbf{e}_{k_A}$  = the  $d_A \times 1$  vector with 1 in the  $k_A$ th entry and all other entries equal to zero,

$\mathbf{e}_{k_B}$  = the  $d_B \times 1$  vector with 1 in the  $k_B$ th entry and all other entries equal to zero

and

$\mathbf{e}_{k_C}$  = the  $d_A^\boxplus \times 1$  vector with 1 in the  $k_C$ th entry and all other entries equal to zero.

The First Term of  $\mathbf{S}_{Ai}$ 

For each  $1 \leq k_A \leq d_A$ , the  $k_A$ th entry of the first term of  $\mathbf{S}_{Ai}$  depends on the function

$$g(\mathbf{u}) = \frac{\mathbf{e}_{k_A}^T \mathbf{b}_N(\mathbf{u})}{b_D(\mathbf{u})} = \mathbf{e}_{k_A}^T \boldsymbol{\Sigma}^{-1} \mathbf{u},$$

and is

$$g(\mathbf{U}_i^*) = \mathbf{e}_{k_A}^T \boldsymbol{\Sigma}^{-1} \mathbf{U}_i^* = \mathbf{e}_{k_A}^T \boldsymbol{\Sigma}^{-1} (\mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(n^{-1}).$$

Therefore, the first term of  $\mathbf{S}_{Ai}$  is

$$\boldsymbol{\Sigma}^{-1} (\mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(n^{-1}) \mathbf{1}_{d_A}.$$

The Second Term of  $\mathbf{S}_{Ai}$

For each  $1 \leq k_A \leq d_A$ , the  $k_A$ th entry of the second term of  $\mathbf{S}_{Ai}$  depends on

$$\nabla g(\mathbf{u}) = \boldsymbol{\Sigma}^{-1} \mathbf{e}_{k_A}$$

and

$$db_D(\mathbf{u}) = d \exp\left(-\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) = -\exp\left(-\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) \mathbf{u}^T \boldsymbol{\Sigma}^{-1} d\mathbf{u},$$

$$\nabla b_D(\mathbf{u}) = -b_D(\mathbf{u}) \boldsymbol{\Sigma}^{-1} \mathbf{u}.$$

It also depends on  $\nabla^2 h_N(\mathbf{u})$  which can be evaluated as follows

$$\begin{aligned} dh_N(\mathbf{u}) &= -\frac{1}{n\phi} \sum_{j=1}^n d \{Y_{ij} \mathbf{u}^T \mathbf{X}_{Aij} - b'((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \\ &= -\frac{1}{n\phi} \sum_{j=1}^n d \{Y_{ij} \mathbf{u}^T \mathbf{X}_{Aij} - b'((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\}^T \\ &= -\frac{1}{n\phi} \sum_{j=1}^n \{Y_{ij} d(\mathbf{u})^T \mathbf{X}_{Aij} - b'((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) d(\mathbf{u})^T \mathbf{X}_{Aij}\}^T \\ &= -\frac{1}{n\phi} \sum_{j=1}^n d \{Y_{ij} - b'((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Aij}^T d\mathbf{u}, \\ d^2 h_N(\mathbf{u}) &= \frac{1}{n\phi} \sum_{j=1}^n (d\mathbf{u})^T \mathbf{X}_{Aij} \{b''((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Aij}^T d\mathbf{u}, \\ \nabla^2 h_N(\mathbf{u}) &= \frac{1}{n\phi} \sum_{j=1}^n \{b''((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T. \end{aligned}$$

By using a stochastic Taylor series approximation, one can show that

$$\{\nabla^2 h_N(\mathbf{U}_i^*)\}^{-1} = n\phi \mathcal{H}_{AAi}^{-1} + O_P(n^{-1/2}) \mathbf{1}_{d_A}^{\otimes 2}.$$

Hence, for each  $1 \leq k_A \leq d_A$ , the  $k_A$ th entry of the second term of  $\mathbf{S}_{Ai}$  is

$$\begin{aligned} & \frac{\nabla g(\mathbf{U}_i^*)^T \{\nabla^2 h_N(\mathbf{U}_i^*)\}^{-1} \nabla b_D(\mathbf{U}_i^*)}{nb_D(\mathbf{U}_i^*)} \\ &= \frac{-\mathbf{e}_{k_A}^T \boldsymbol{\Sigma}^{-1} \left\{ n\phi \mathcal{H}_{AAi}^{-1} + O_P(n^{-1/2}) \mathbf{1}_{d_A}^{\otimes 2} \right\} b_D(\mathbf{U}_i^*) \boldsymbol{\Sigma}^{-1} \mathbf{U}_i^*}{nb_D(\mathbf{U}_i^*)} \\ &= -\frac{1}{n} \mathbf{e}_{k_A}^T \boldsymbol{\Sigma}^{-1} \left\{ n\phi \mathcal{H}_{AAi}^{-1} + O_P(n^{-1/2}) \mathbf{1}_{d_A}^{\otimes 2} \right\} \boldsymbol{\Sigma}^{-1} \{\mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} + O_P(n^{-1} \mathbf{1}_{d_A})\} \\ &= -\phi \mathbf{e}_{k_A}^T \boldsymbol{\Sigma}^{-1} \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{U}_i + O_P(n^{-3/2}). \end{aligned}$$

Then the leading term behaviour of the second term of  $\mathbf{S}_{Bi}$  is as follows

$$-\phi \boldsymbol{\Sigma}^{-1} \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{U}_i + O_P(n^{-3/2}) \mathbf{1}_{d_A} = O_P(n^{-1}) \mathbf{1}_{d_A}.$$

The Third Term of  $\mathbf{S}_{Ai}$ 

For each  $1 \leq k_A \leq d_A$ , the  $k_A$ th entry of the third term of  $\mathbf{S}_{Ai}$  depends on

$$\nabla^2 g(\mathbf{u}) = \nabla^2 (\mathbf{e}_{k_A}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}) = \mathbf{O},$$

the  $d_A \times d_A$  matrix of zeroes. Since this matrix appears in a trace expression, the third term of  $\mathbf{S}_{Ai}$  is  $\mathbf{0}_{d_A}$ .

The Fourth Term of  $\mathbf{S}_{Ai}$ 

The contribution from the fourth term of  $\mathbf{S}_{Ai}$  is  $O_P(1)\mathbf{1}_{d_A}$  but does not have a concise matrix algebraic expression. It is of lower order compared to the leading term of  $\mathbf{S}_{Ai}$ .

Overall Leading Term Expression for  $\mathbf{S}_{Ai}$ 

Putting the terms of  $\mathbf{S}_{Ai}$  together, we can assert that

$$\mathbf{S}_{Ai} = \boldsymbol{\Sigma}^{-1} (\mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(n^{-1})\mathbf{1}_{d_A}.$$

**4.5.2.5 Computing an Asymptotic Approximation for the Second Entry in (4.7)**

The  $i$ th contribution to the score of  $\boldsymbol{\beta}_B$  is

$$\begin{aligned} \mathbf{S}_{Bi} &= \nabla_{\boldsymbol{\beta}_B} \log p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i) \\ &= \frac{\int_{\mathbb{R}^{d_A}} b_N(\mathbf{u}) \exp\{-nh_N(\mathbf{u})\} d\mathbf{u}}{\int_{\mathbb{R}^{d_A}} b_D(\mathbf{u}) \exp\{-nh_N(\mathbf{u})\} d\mathbf{u}} \end{aligned}$$

where

$$\begin{aligned} b_N(\mathbf{u}) &\equiv \exp\left(-\frac{1}{2}\mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) \frac{1}{\phi} \sum_{j=1}^n [\{Y_{ij} - b'((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Bij}] \\ b_D(\mathbf{u}) &\equiv \exp\left(-\frac{1}{2}\mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) \quad \text{and} \\ h_N(\mathbf{u}) &\equiv -\frac{1}{\phi n} \sum_{j=1}^n \{Y_{ij} \mathbf{u}^T \mathbf{X}_{Aij} - b((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\}. \end{aligned}$$

The First Term of  $\mathbf{S}_{Bi}$ 

For each  $1 \leq k_B \leq d_B$ , the  $k_B$ th entry of the first term of  $\mathbf{S}_{Bi}$  depends on the function

$$g(\mathbf{u}) = \frac{\mathbf{e}_{k_B}^T b_N(\mathbf{u})}{b_D(\mathbf{u})} = \frac{1}{\phi} \sum_{j=1}^n \mathbf{e}_{k_B}^T [\{Y_{ij} - b'((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Bij}]$$



and is

$$\begin{aligned}
g(\mathbf{U}_i^*) &= \frac{1}{\phi} \sum_{j=1}^n \mathbf{e}_{k_B}^T [\{Y_{ij} - b'((\boldsymbol{\beta}_A + \mathbf{U}_i^*)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Bij}] \\
&= \frac{1}{\phi} \sum_{j=1}^n \mathbf{e}_{k_B}^T \left[ \{Y_{ij} - b'((\boldsymbol{\beta}_A + \mathbf{U}_i + \mathbf{U}_i^* - \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij})\} \mathbf{X}_{Bij} \right] \\
&= \frac{1}{\phi} \sum_{j=1}^n \mathbf{e}_{k_B}^T \left[ \{Y_{ij} - b'((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij})\} \mathbf{X}_{Bij} \right. \\
&\quad \left. - b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + \boldsymbol{\beta}_B^T \mathbf{X}_{Bij}) \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T (\mathbf{U}_i^* - \mathbf{U}_i) + O_P(n^{-1}) \mathbf{1}_{d_B} \right] \\
&= \frac{1}{\phi} \mathbf{e}_{k_B}^T (\mathcal{G}_{Bi} - \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(1).
\end{aligned}$$

Therefore, the first term of  $\mathbf{S}_{Bi}$  is

$$\frac{1}{\phi} (\mathcal{G}_{Bi} - \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(1) \mathbf{1}_{d_B}.$$

#### The Second Term of $\mathbf{S}_{Bi}$

For each  $1 \leq k_B \leq d_B$ , the  $k_B$ th entry of the second term of  $\mathbf{S}_{Bi}$  depends on

$$\begin{aligned}
dg(\mathbf{u}) &= \frac{1}{\phi} \sum_{j=1}^n \mathbf{e}_{k_B}^T \mathbf{X}_{Bij} [d \{Y_{ij} - b'((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\}] \\
&= -\frac{1}{\phi} \sum_{j=1}^n \mathbf{e}_{k_B}^T \mathbf{X}_{Bij} b''((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \mathbf{X}_{Aij}^T (d\mathbf{u}), \\
\nabla g(\mathbf{u}) &= -\frac{1}{\phi} \sum_{j=1}^n b''((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \mathbf{X}_{Aij} \mathbf{X}_{Bij}^T \mathbf{e}_{k_B}
\end{aligned}$$

and

$$\nabla b_D(\mathbf{u}) = -b_D(\mathbf{u}) \boldsymbol{\Sigma}^{-1} \mathbf{u}.$$

It also depends on

$$\nabla^2 h_N(\mathbf{u}) = \frac{1}{n\phi} \sum_{j=1}^n \{b''((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})\} \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T.$$

By using a stochastic Taylor series approximation, note that

$$\{\nabla^2 h_N(\mathbf{U}_i^*)\}^{-1} = n\phi \mathcal{H}_{AAi}^{-1} + O_P(n^{-1/2}) \mathbf{1}_{d_A}^{\otimes 2}.$$

Hence, for each  $1 \leq k_B \leq d_B$ , the  $k_B$ th entry of the second term of  $\mathbf{S}_{B_i}$  is

$$\begin{aligned}
& \frac{\nabla g(\mathbf{U}_i^*)^T \{\nabla^2 h_N(\mathbf{U}_i^*)\}^{-1} \nabla b_D(\mathbf{U}_i^*)}{nb_D(\mathbf{U}_i^*)} \\
&= \frac{\frac{1}{\phi} \mathbf{e}_{k_B}^T \sum_{j=1}^n b''((\beta_A + \mathbf{U}_i^*)^T \mathbf{X}_{A_{ij}} + (\beta_B)^T \mathbf{X}_{B_{ij}}) \mathbf{X}_{B_{ij}} \mathbf{X}_{A_{ij}}^T \left\{ n\phi \mathcal{H}_{AA_i}^{-1} + O_P(n^{-1/2}) \mathbf{1}_{d_A}^{\otimes 2} \right\}}{nb_D(\mathbf{U}_i^*)}} \\
&\quad \times \frac{b_D(\mathbf{U}_i^*) \boldsymbol{\Sigma}^{-1} \mathbf{U}_i^*}{nb_D(\mathbf{U}_i^*)} \\
&= \frac{1}{n\phi} \mathbf{e}_{k_B}^T \mathcal{H}_{AB_i}^T \left\{ n\phi \mathcal{H}_{AA_i}^{-1} + O_P(n^{-1/2}) \mathbf{1}_{d_A}^{\otimes 2} \right\} \boldsymbol{\Sigma}^{-1} \left\{ \mathbf{U}_i + \mathcal{H}_{AA_i}^{-1} \mathcal{G}_{A_i} + O_P(n^{-1} \mathbf{1}_{d_A}) \right\} \\
&= \mathbf{e}_{k_B}^T \mathcal{H}_{AB_i}^T \mathcal{H}_{AA_i}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{U}_i + O_P(n^{-1/2}).
\end{aligned}$$

Then the leading term behaviour of the second term of  $\mathbf{S}_{B_i}$  is as follows

$$\mathcal{H}_{AB_i}^T \mathcal{H}_{AA_i}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{U}_i + O_P(n^{-1/2}) \mathbf{1}_{d_B} = O_P(1) \mathbf{1}_{d_B}.$$

#### The Third Term of $\mathbf{S}_{B_i}$

The contribution from the third term of  $\mathbf{S}_{B_i}$  is  $O_P(1) \mathbf{1}_{d_B}$  but does not have a concise matrix algebraic expression. It is of lower order compared to the leading term of  $\mathbf{S}_{B_i}$ .

#### The Fourth Term of $\mathbf{S}_{B_i}$

Similarly, the contribution from the fourth term of  $\mathbf{S}_{B_i}$  is also  $O_P(1) \mathbf{1}_{d_B}$  and does not have a concise matrix algebraic expression. It is of lower order compared to the leading term of  $\mathbf{S}_{B_i}$ .

#### Overall Leading Term Expression for $\mathbf{S}_{B_i}$

Combining all four asymptotic approximations of the terms of  $\mathbf{S}_{B_i}$  together, we have that

$$\mathbf{S}_{B_i} = \frac{1}{\phi} (\mathcal{G}_{B_i} - \mathcal{H}_{AB_i}^T \mathcal{H}_{AA_i}^{-1} \mathcal{G}_{A_i}) + O_P(1) \mathbf{1}_{d_B}.$$

### 4.5.2.6 Computing an Asymptotic Approximation for the Third Entry in (4.7)

First note that,

$$\begin{aligned}
d \log |\boldsymbol{\Sigma}| &= \frac{d|\boldsymbol{\Sigma}|}{|\boldsymbol{\Sigma}|} \\
&= \frac{|\boldsymbol{\Sigma}| \operatorname{tr}(\boldsymbol{\Sigma}^{-1} d\boldsymbol{\Sigma})}{|\boldsymbol{\Sigma}|} \\
&= \operatorname{vec}(\boldsymbol{\Sigma}^{-1})^T \operatorname{vec}(d\boldsymbol{\Sigma}) \\
&= \operatorname{vec}(\boldsymbol{\Sigma}^{-1})^T d\operatorname{vec}(\boldsymbol{\Sigma}) \\
&= \operatorname{vec}(\boldsymbol{\Sigma}^{-1})^T d\mathbf{D}_{d_A} \operatorname{vech}(\boldsymbol{\Sigma}) \\
&= \operatorname{vec}(\boldsymbol{\Sigma}^{-1})^T \mathbf{D}_{d_A} d\operatorname{vech}(\boldsymbol{\Sigma}).
\end{aligned}$$

Next note that,

$$\begin{aligned}
d\mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} &= -\mathbf{u}^T \boldsymbol{\Sigma}^{-1} (d\boldsymbol{\Sigma}) \boldsymbol{\Sigma}^{-1} \mathbf{u} \\
&= -\operatorname{tr}\{\boldsymbol{\Sigma}^{-1} \mathbf{u} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} (d\boldsymbol{\Sigma})\} \\
&= -\operatorname{vec}(\boldsymbol{\Sigma}^{-1} \mathbf{u} \mathbf{u}^T \boldsymbol{\Sigma}^{-1})^T \operatorname{vec}(d\boldsymbol{\Sigma}) \\
&= -\{(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \operatorname{vec}(\mathbf{u} \mathbf{u}^T)\}^T \mathbf{D}_{d_A} d\operatorname{vech}(\boldsymbol{\Sigma}).
\end{aligned}$$

With further calculations, the  $i$ th contribution to the score of  $\operatorname{vech}(\boldsymbol{\Sigma})$  can be computed as

$$\begin{aligned}
\mathbf{S}_{C_i} &= \nabla_{\operatorname{vech}(\boldsymbol{\Sigma})} \log p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i) \\
&= -\frac{1}{2} \mathbf{D}_{d_A}^T \operatorname{vec}(\boldsymbol{\Sigma}^{-1}) + \frac{1}{2} \frac{\int_{\mathbb{R}^{d_A}} b_N(\mathbf{u}) \exp\{-nh_N(\mathbf{u})\} d\mathbf{u}}{\int_{\mathbb{R}^{d_A}} b_D(\mathbf{u}) \exp\{-nh_N(\mathbf{u})\} d\mathbf{u}}
\end{aligned}$$

where

$$\begin{aligned}
b_N(\mathbf{u}) &\equiv \exp\left(-\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \operatorname{vec}(\mathbf{u} \mathbf{u}^T), \\
b_D(\mathbf{u}) &\equiv \exp\left(-\frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}\right) \quad \text{and} \\
h_N(\mathbf{u}) &\equiv -\frac{1}{n\phi} \sum_{j=1}^n \left\{ Y_{ij} \mathbf{u}^T \mathbf{X}_{Aij} - b((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \right\}.
\end{aligned}$$

#### The First Term of the Integral Ratio Component of $\mathbf{S}_{C_i}$

For each  $1 \leq k_C \leq d_A^{\#}$ , the  $k_C$ th entry of the first term of the integral ratio component of  $\mathbf{S}_{C_i}$  depends on the function

$$g(\mathbf{u}) = \frac{e_{k_C}^T \mathbf{b}_N(\mathbf{u})}{b_D(\mathbf{u})} = e_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \operatorname{vec}(\mathbf{u} \mathbf{u}^T)$$

and is

$$\begin{aligned}
g(\mathbf{U}_i^*) &= \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec} \{ \mathbf{U}_i^* (\mathbf{U}_i^*)^T \} \\
&= \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec} \left[ \{ \mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} + O_P(n^{-1} \mathbf{1}_{d_A}) \} \right. \\
&\quad \left. \times \{ \mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} + O_P(n^{-1} \mathbf{1}_{d_A}) \}^T \right] \\
&= \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec} (\mathbf{U}_i \mathbf{U}_i^T + \mathbf{U}_i \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathbf{U}_i^T) + O_P(n^{-1}).
\end{aligned}$$

Therefore, the first term of the integral ratio component of  $\mathbf{S}_{C_i}$  is

$$\begin{aligned}
&\mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec} (\mathbf{U}_i \mathbf{U}_i^T + \mathbf{U}_i \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathbf{U}_i^T) + O_P(n^{-1}) \mathbf{1}_{d_A^{\boxplus}} \\
&= \mathbf{D}_{d_A}^T \text{vec} \{ \boldsymbol{\Sigma}^{-1} (\mathbf{U}_i \mathbf{U}_i^T + \mathbf{U}_i \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathbf{U}_i^T) \boldsymbol{\Sigma}^{-1} \} + O_P(n^{-1}) \mathbf{1}_{d_A^{\boxplus}}.
\end{aligned}$$

The Second Term of the Integral Ratio Component of  $\mathbf{S}_{C_i}$

For each  $1 \leq k_C \leq d_A^{\boxplus}$ , the  $k_C$ th entry of the second term of the integral ratio component of  $\mathbf{S}_{C_i}$  depends on the following function. By making use of the first property in Theorem 12 of Magnus and Neudecker (1999), we have

$$\begin{aligned}
dg(\mathbf{u}) &= \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) d\text{vec} (\mathbf{u} \mathbf{u}^T) \\
&= \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec} \{ (d\mathbf{u}) \mathbf{u}^T + \mathbf{u} (d\mathbf{u})^T \} \\
&= \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec} \{ \mathbf{I}_{d_A} (d\mathbf{u}) \mathbf{u}^T \} + \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec} \{ \mathbf{u} (d\mathbf{u})^T \mathbf{I}_{d_A} \} \\
&= \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) (\mathbf{u} \otimes \mathbf{I}_{d_A}) \text{vec} (d\mathbf{u}) \\
&\quad + \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) (\mathbf{I}_{d_A} \otimes \mathbf{u}) \text{vec} \{ (d\mathbf{u})^T \} \\
&= \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \{ (\mathbf{u} \otimes \mathbf{I}_{d_A}) + (\mathbf{I}_{d_A} \otimes \mathbf{u}) \} d\mathbf{u} \\
&= \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T \{ (\boldsymbol{\Sigma}^{-1} \mathbf{u} \otimes \boldsymbol{\Sigma}^{-1}) + (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1} \mathbf{u}) \} d\mathbf{u} \\
&= \mathbf{e}_{k_C}^T \left\{ \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \mathbf{u} \otimes \boldsymbol{\Sigma}^{-1}) + (\mathbf{K}_{d_A} \mathbf{D}_{d_A})^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1} \mathbf{u}) \right\} d\mathbf{u} \\
&= \mathbf{e}_{k_C}^T \left\{ \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \mathbf{u} \otimes \boldsymbol{\Sigma}^{-1}) + \mathbf{D}_{d_A}^T \mathbf{K}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1} \mathbf{u}) \right\} d\mathbf{u} \\
&= \mathbf{e}_{k_C}^T \left\{ \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \mathbf{u} \otimes \boldsymbol{\Sigma}^{-1}) + \mathbf{D}_{d_A}^T \mathbf{K}_{d_A} (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1} \mathbf{u}) \right\} d\mathbf{u} \\
&= \mathbf{e}_{k_C}^T \left\{ \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \mathbf{u} \otimes \boldsymbol{\Sigma}^{-1}) + \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \mathbf{u} \otimes \boldsymbol{\Sigma}^{-1}) \right\} d\mathbf{u} \\
&= 2\mathbf{e}_{k_C}^T \left\{ \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \mathbf{u} \otimes \boldsymbol{\Sigma}^{-1}) \right\} d\mathbf{u}, \\
\nabla g(\mathbf{u}) &= 2(\boldsymbol{\Sigma}^{-1} \mathbf{u} \otimes \boldsymbol{\Sigma}^{-1})^T \mathbf{D}_{d_A} \mathbf{e}_{k_C}
\end{aligned}$$

and

$$\nabla b_D(\mathbf{u}) = -b_D(\mathbf{u}) \boldsymbol{\Sigma}^{-1} \mathbf{u}.$$

It also depends on

$$\nabla^2 h_N(\mathbf{u}) = \frac{1}{n\phi} \sum_{j=1}^n \{ b''((\boldsymbol{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \} \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T.$$

By using a stochastic Taylor series approximation, note that

$$\{\nabla^2 h_N(\mathbf{U}_i^*)\}^{-1} = n\phi\mathcal{H}_{AAi}^{-1} + O_P(n^{-1/2})\mathbf{1}_{d_A}^{\otimes 2}.$$

Hence, for each  $1 \leq k_C \leq d_A^{\boxplus}$ , the  $k_C$ th entry of the second term of the integral component of  $\mathbf{S}_{Ci}$  is

$$\begin{aligned} & \frac{\nabla g(\mathbf{U}_i^*)^T \{\nabla^2 h_N(\mathbf{U}_i^*)\}^{-1} \nabla b_D(\mathbf{U}_i^*)}{nb_D(\mathbf{U}_i^*)} \\ &= \frac{-2\mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \mathbf{U}_i^* \otimes \boldsymbol{\Sigma}^{-1}) \left\{ n\phi\mathcal{H}_{AAi}^{-1} + O_P(n^{-1/2})\mathbf{1}_{d_A}^{\otimes 2} \right\} b_D(\mathbf{U}_i^*) \boldsymbol{\Sigma}^{-1} \mathbf{U}_i^*}{nb_D(\mathbf{U}_i^*)} \\ &= -2\phi \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T \left[ \boldsymbol{\Sigma}^{-1} \left\{ \mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} + O_P(n^{-1} \mathbf{1}_{d_A}) \right\} \otimes \boldsymbol{\Sigma}^{-1} \right] \left\{ \mathcal{H}_{AAi}^{-1} + O_P(n^{-3/2})\mathbf{1}_{d_A}^{\otimes 2} \right\} \boldsymbol{\Sigma}^{-1} \\ &\quad \times \left\{ \mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} + O_P(n^{-1} \mathbf{1}_{d_A}) \right\} \\ &= -2\phi \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \mathbf{U}_i \otimes \boldsymbol{\Sigma}^{-1}) \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{U}_i + O_P(n^{-3/2}) \\ &= -2\phi \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \mathbf{U}_i \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{U}_i) + O_P(n^{-3/2}) \\ &= -2\phi \mathbf{e}_{k_C}^T \mathbf{D}_{d_A}^T \text{vec}(\boldsymbol{\Sigma}^{-1} \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{U}_i \mathbf{U}_i^T \boldsymbol{\Sigma}^{-1}) + O_P(n^{-3/2}). \end{aligned}$$

Then the leading term behaviour of the second term of the integral component of  $\mathbf{S}_{Ci}$  is as follows

$$-2\phi \mathbf{D}_{d_A}^T \text{vec}(\boldsymbol{\Sigma}^{-1} \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{U}_i \mathbf{U}_i^T \boldsymbol{\Sigma}^{-1}) + O_P(n^{-3/2})\mathbf{1}_{d_A^{\boxplus}} = O_P(n^{-1})\mathbf{1}_{d_A^{\boxplus}}.$$

#### The Third Term of the Integral Ratio Component of $\mathbf{S}_{Ci}$

The leading term behaviour of the third term of the integral component of  $\mathbf{S}_{Ci}$  is

$$\phi \mathbf{D}_{d_A}^T \text{vec}(\boldsymbol{\Sigma}^{-1} \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1}) + O_P(n^{-3/2})\mathbf{1}_{d_A^{\boxplus}} = O_P(n^{-1})\mathbf{1}_{d_A^{\boxplus}}.$$

#### The Fourth Term of the Integral Ratio Component of $\mathbf{S}_{Ci}$

The contribution from the fourth term of  $\mathbf{S}_{Bi}$  is also  $O_P(n^{-1})\mathbf{1}_{d_A^{\boxplus}}$  but does not have a concise matrix algebraic expression. It is of lower order compared to the leading term of the integral ratio component of  $\mathbf{S}_{Ci}$ .

#### Overall Leading Term Expression for $\mathbf{S}_{Ci}$

Combining all four asymptotic approximations of the terms of the integral ratio component of  $\mathbf{S}_{Ci}$  together, we have that

$$\begin{aligned} \mathbf{S}_{Ci} &= \frac{1}{2} \mathbf{D}_{d_A}^T \left[ \text{vec} \left\{ \boldsymbol{\Sigma}^{-1} (\mathbf{U}_i \mathbf{U}_i^T + \mathbf{U}_i \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathbf{U}_i^T) \boldsymbol{\Sigma}^{-1} \right\} - \text{vec}(\boldsymbol{\Sigma}^{-1}) \right] \\ &\quad + O_P(n^{-1})\mathbf{1}_{d_A^{\boxplus}}. \end{aligned}$$

### 4.5.2.7 The Quadratic Conditional Expectations of the Scores

In this subsection we find the conditional expectations required to compute the Fisher information matrix of  $(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}))$ .

#### The Expectation of $\mathbf{S}_{Ai}\mathbf{S}_{Ai}^T$ Given $\mathbf{X}_i$

From the previous sections, we have the following approximation

$$\mathbf{S}_{Ai} = \boldsymbol{\Sigma}^{-1} (\mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(n^{-1}) \mathbf{1}_{d_A}.$$

Therefore,

$$\begin{aligned} & E(\mathbf{S}_{Ai}\mathbf{S}_{Ai}^T | \mathbf{X}_i) \\ &= E \left[ \left\{ \boldsymbol{\Sigma}^{-1} (\mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(n^{-1}) \mathbf{1}_{d_A} \right\} \left\{ \boldsymbol{\Sigma}^{-1} (\mathbf{U}_i + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(n^{-1}) \mathbf{1}_{d_A} \right\}^T \middle| \mathbf{X}_i \right] \\ &= \boldsymbol{\Sigma}^{-1} E(\mathbf{U}_i \mathbf{U}_i^T) \boldsymbol{\Sigma}^{-1} + E(\boldsymbol{\Sigma}^{-1} \mathbf{U}_i \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} | \mathbf{X}_i) + E(\boldsymbol{\Sigma}^{-1} \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathbf{U}_i^T \boldsymbol{\Sigma}^{-1} | \mathbf{X}_i) \\ &\quad + E(\boldsymbol{\Sigma}^{-1} \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} | \mathbf{X}_i) + O_P(n^{-1}) \mathbf{1}_{d_A}^{\otimes 2} \\ &= \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \boldsymbol{\Sigma}^{-1} + E \left\{ \boldsymbol{\Sigma}^{-1} \mathbf{U}_i E(\mathcal{G}_{Ai}^T | \mathbf{X}_i, \mathbf{U}_i) \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} | \mathbf{X}_i \right\} \\ &\quad + E \left\{ \boldsymbol{\Sigma}^{-1} \mathcal{H}_{AAi}^{-1} E(\mathcal{G}_{Ai} | \mathbf{X}_i, \mathbf{U}_i) \mathbf{U}_i^T \boldsymbol{\Sigma}^{-1} | \mathbf{X}_i \right\} \\ &\quad + E \left\{ \boldsymbol{\Sigma}^{-1} \mathcal{H}_{AAi}^{-1} E(\mathcal{G}_{Ai} \mathcal{G}_{Ai}^T | \mathbf{X}_i, \mathbf{U}_i) \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} | \mathbf{X}_i \right\} + O_P(n^{-1}) \mathbf{1}_{d_A}^{\otimes 2} \\ &= \boldsymbol{\Sigma}^{-1} + \phi E \left\{ \boldsymbol{\Sigma}^{-1} \mathcal{H}_{AAi}^{-1} \boldsymbol{\Sigma}^{-1} | \mathbf{X}_i \right\} + O_P(n^{-1}) \mathbf{1}_{d_A}^{\otimes 2} \\ &= \boldsymbol{\Sigma}^{-1} + O_P(n^{-1}) \mathbf{1}_{d_A}^{\otimes 2}. \end{aligned}$$

#### The Expectation of $\mathbf{S}_{Bi}\mathbf{S}_{Bi}^T$ Given $\mathbf{X}_i$

From the previous sections, we have the following approximation

$$\mathbf{S}_{Bi} = \frac{1}{\phi} (\mathcal{G}_{Bi} - \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(1) \mathbf{1}_{d_B}.$$

Therefore,

$$\begin{aligned} & E(\mathbf{S}_{Bi}\mathbf{S}_{Bi}^T | \mathbf{X}_i) \\ &= E \left[ \left\{ \frac{1}{\phi} (\mathcal{G}_{Bi} - \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(1) \mathbf{1}_{d_B} \right\} \left\{ \frac{1}{\phi} (\mathcal{G}_{Bi} - \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai}) + O_P(1) \mathbf{1}_{d_B} \right\}^T \middle| \mathbf{X}_i \right] \\ &= \frac{1}{\phi^2} \left\{ E(\mathcal{G}_{Bi} \mathcal{G}_{Bi}^T | \mathbf{X}_i) - E(\mathcal{G}_{Bi} \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1} \mathcal{H}_{ABi} | \mathbf{X}_i) - E(\mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathcal{G}_{Bi}^T | \mathbf{X}_i) \right. \\ &\quad \left. + E(\mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1} \mathcal{H}_{ABi} | \mathbf{X}_i) \right\} + O_P(1) \mathbf{1}_{d_B}^{\otimes 2} \\ &= \frac{1}{\phi^2} \left[ E \left\{ E(\mathcal{G}_{Bi} \mathcal{G}_{Bi}^T | \mathbf{X}_i, \mathbf{U}_i) \middle| \mathbf{X}_i \right\} - E \left\{ E(\mathcal{G}_{Bi} \mathcal{G}_{Ai}^T | \mathbf{X}_i, \mathbf{U}_i) \mathcal{H}_{AAi}^{-1} \mathcal{H}_{ABi} \middle| \mathbf{X}_i \right\} \right. \\ &\quad \left. - E \left\{ \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} E(\mathcal{G}_{Ai} \mathcal{G}_{Bi}^T | \mathbf{X}_i, \mathbf{U}_i) \middle| \mathbf{X}_i \right\} \right. \\ &\quad \left. + E \left\{ \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} E(\mathcal{G}_{Ai} \mathcal{G}_{Ai}^T | \mathbf{X}_i, \mathbf{U}_i) \mathcal{H}_{AAi}^{-1} \mathcal{H}_{ABi} \middle| \mathbf{X}_i \right\} \right] + O_P(1) \mathbf{1}_{d_B}^{\otimes 2}. \end{aligned}$$

Simplifying the previous expression, we have,

$$\begin{aligned} & \frac{1}{\phi^2} \left\{ \phi E(\mathcal{H}_{\text{BB}i} | \mathbf{X}_i) - \phi E(\mathcal{H}_{\text{AB}i}^T \mathcal{H}_{\text{AA}i}^{-1} \mathcal{H}_{\text{AB}i} | \mathbf{X}_i) - \phi E(\mathcal{H}_{\text{AB}i}^T \mathcal{H}_{\text{AA}i}^{-1} \mathcal{H}_{\text{AB}i} | \mathbf{X}_i) \right. \\ & \quad \left. + \phi E(\mathcal{H}_{\text{AB}i}^T \mathcal{H}_{\text{AA}i}^{-1} \mathcal{H}_{\text{AB}i} | \mathbf{X}_i) \right\} + O_P(1) \mathbf{1}_{d_{\text{B}}}^{\otimes 2} \\ & = \frac{1}{\phi} E(\mathcal{H}_{\text{BB}i} - \mathcal{H}_{\text{AB}i}^T \mathcal{H}_{\text{AA}i}^{-1} \mathcal{H}_{\text{AB}i} | \mathbf{X}_i) + O_P(1) \mathbf{1}_{d_{\text{B}}}^{\otimes 2}. \end{aligned}$$

The Expectation of  $\mathbf{S}_{\text{Ci}} \mathbf{S}_{\text{Ci}}^T$  Given  $\mathbf{X}_i$

From the previous sections, we have the following approximation

$$\begin{aligned} \mathbf{S}_{\text{Ci}} & = \frac{1}{2} \mathbf{D}_{d_{\text{A}}}^T [\text{vec} \{ \boldsymbol{\Sigma}^{-1} (\mathbf{U}_i \mathbf{U}_i^T + \mathbf{U}_i \mathcal{G}_{\text{Ai}}^T \mathcal{H}_{\text{AA}i}^{-1} + \mathcal{H}_{\text{AA}i}^{-1} \mathcal{G}_{\text{Ai}} \mathbf{U}_i^T) \boldsymbol{\Sigma}^{-1} \} - \text{vec}(\boldsymbol{\Sigma}^{-1})] \\ & \quad + O_P(n^{-1}) \mathbf{1}_{d_{\text{A}}}^{\text{H}} \\ & = \frac{1}{2} \mathbf{D}_{d_{\text{A}}}^T \{ \text{vec}(\boldsymbol{\Sigma}^{-1} \mathbf{U}_i \mathbf{U}_i^T \boldsymbol{\Sigma}^{-1}) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \} \\ & \quad + \frac{1}{2} \mathbf{D}_{d_{\text{A}}}^T [\text{vec} \{ \boldsymbol{\Sigma}^{-1} (\mathbf{U}_i \mathcal{G}_{\text{Ai}}^T \mathcal{H}_{\text{AA}i}^{-1} + \mathcal{H}_{\text{AA}i}^{-1} \mathcal{G}_{\text{Ai}} \mathbf{U}_i^T) \boldsymbol{\Sigma}^{-1} \}] + O_P(n^{-1}) \mathbf{1}_{d_{\text{A}}}^{\text{H}} \\ & = \frac{1}{2} \mathbf{D}_{d_{\text{A}}}^T \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \} \\ & \quad + \frac{1}{2} \mathbf{D}_{d_{\text{A}}}^T \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathcal{G}_{\text{Ai}}^T \mathcal{H}_{\text{AA}i}^{-1} + \mathcal{H}_{\text{AA}i}^{-1} \mathcal{G}_{\text{Ai}} \mathbf{U}_i^T) \} + O_P(n^{-1}) \mathbf{1}_{d_{\text{A}}}^{\text{H}}. \end{aligned}$$

We will deal with each term arising in  $E(\mathbf{S}_{\text{Ci}} \mathbf{S}_{\text{Ci}}^T | \mathbf{X}_i)$  separately. The first term in  $E(\mathbf{S}_{\text{Ci}} \mathbf{S}_{\text{Ci}}^T | \mathbf{X}_i)$  is

$$\frac{1}{4} \mathbf{D}_{d_{\text{A}}}^T E \{ [(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1})] \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \}^T ] \mathbf{D}_{d_{\text{A}}}.$$

Next note that

$$\begin{aligned} & E \{ [(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1})] \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \}^T ] \\ & = E \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T)^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \} \\ & \quad - E \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) \text{vec}(\boldsymbol{\Sigma}^{-1})^T \} - E \{ \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T)^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \} \\ & \quad + \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}^{-1})^T. \end{aligned} \tag{4.9}$$

Then note that

$$E \{ \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) \} = \text{vec}(E(\mathbf{U}_i \mathbf{U}_i^T)) = \text{vec}(\boldsymbol{\Sigma}).$$

With this we can simplify the following

$$\begin{aligned} & - E \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) \text{vec}(\boldsymbol{\Sigma}^{-1})^T \} \\ & = - (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}) \text{vec}(\boldsymbol{\Sigma}^{-1})^T \\ & = - \text{vec}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}^{-1})^T \\ & = - \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}^{-1})^T \end{aligned}$$

and

$$\begin{aligned}
& - E \{ \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T)^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \} \\
& = - \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma})^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \\
& = - \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \boldsymbol{\Sigma}^{-1})^T \\
& = - \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}^{-1})^T.
\end{aligned}$$

These calculations lead to the right-hand side of (4.9) to simplify as

$$(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) E \{ \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T)^T \} (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}^{-1})^T. \quad (4.10)$$

Next, we appeal to Theorem 4.3 (iv) Magnus and Neudecker (1979) to get

$$\begin{aligned}
E \{ \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T)^T \} & = E \{ (\mathbf{U}_i \otimes \mathbf{U}_i) (\mathbf{U}_i \otimes \mathbf{U}_i)^T \} \\
& = \text{Cov}(\mathbf{U}_i \otimes \mathbf{U}_i) + E(\mathbf{U}_i \otimes \mathbf{U}_i) \{ E(\mathbf{U}_i \otimes \mathbf{U}_i) \}^T \\
& = (\mathbf{I}_{d_A^2} + \mathbf{K}_{d_A}) (\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}) + E \{ \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) \} E \{ \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) \}^T \\
& = (\mathbf{I}_{d_A^2} + \mathbf{K}_{d_A}) (\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}) + \text{vec}(\boldsymbol{\Sigma}) \text{vec}(\boldsymbol{\Sigma})^T.
\end{aligned} \quad (4.11)$$

Substitution of (4.11) into (4.10) leads to

$$\begin{aligned}
& (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) (\mathbf{I}_{d_A^2} + \mathbf{K}_{d_A}) (\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}) (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \\
& + (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}) \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}) \}^T - \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}^{-1})^T \\
& = (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) (\mathbf{I}_{d_A^2} + \mathbf{K}_{d_A}) (\mathbf{I}_d \otimes \mathbf{I}_d) + \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}^{-1})^T - \text{vec}(\boldsymbol{\Sigma}^{-1}) \text{vec}(\boldsymbol{\Sigma}^{-1})^T \\
& = (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) (\mathbf{I}_{d_A^2} + \mathbf{K}_{d_A}) (\mathbf{I}_d \otimes \mathbf{I}_d) \\
& = (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) + (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{K}_{d_A} (\mathbf{I}_d \otimes \mathbf{I}_d) \\
& = (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) + (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) (\mathbf{I}_d \otimes \mathbf{I}_d) \mathbf{K}_{d_A} \\
& = (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) + (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{K}_{d_A} \\
& = (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) + \mathbf{K}_{d_A} (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \\
& = (\mathbf{I}_{d_A^2} + \mathbf{K}_{d_A}) (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}).
\end{aligned}$$

Therefore, we have the first term of  $E(\mathbf{S}_{C_i} \mathbf{S}_{C_i}^T | \mathbf{X}_i)$  being equal to

$$\begin{aligned}
& \frac{1}{4} \mathbf{D}_{d_A}^T (\mathbf{I}_{d_A^2} + \mathbf{K}_{d_A}) (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} \\
& = \frac{1}{4} \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} + \frac{1}{4} (\mathbf{K}_{d_A} \mathbf{D}_{d_A})^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} \\
& = \frac{1}{2} \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A}.
\end{aligned}$$

We can then show that the second term of  $E(\mathbf{S}_{C_i} \mathbf{S}_{C_i}^T | \mathbf{X}_i)$  is

$$\begin{aligned}
& \frac{1}{4} \mathbf{D}_{d_A}^T E \{ \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \} \\
& \quad \times \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathcal{G}_{A_i}^T \mathcal{H}_{AA_i}^{-1} + \mathcal{H}_{AA_i}^{-1} \mathcal{G}_{A_i} \mathbf{U}_i^T) \}^T | \mathbf{X}_i \} \mathbf{D}_{d_A} \\
& = \frac{1}{4} \mathbf{D}_{d_A}^T E \left( \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \} \right. \\
& \quad \left. \times [ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec} \{ \mathbf{U}_i E(\mathcal{G}_{A_i}^T | \mathbf{X}_i, \mathbf{U}_i) \mathcal{H}_{AA_i}^{-1} + \mathcal{H}_{AA_i}^{-1} E(\mathcal{G}_{A_i} | \mathbf{X}_i, \mathbf{U}_i) \mathbf{U}_i^T \} \}^T | \mathbf{X}_i \right) \mathbf{D}_{d_A} \\
& = \mathbf{O}.
\end{aligned}$$



Similarly, we can then show that the third term of  $E(\mathbf{S}_{Ci}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$  is

$$\begin{aligned} & \frac{1}{4}\mathbf{D}_{d_A}^T E\left[\{(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\text{vec}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\} \right. \\ & \quad \times \left. \{(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\text{vec}(\mathbf{U}_i\mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1})\}^T|\mathbf{X}_i\right]\mathbf{D}_{d_A} \\ &= \frac{1}{4}\mathbf{D}_{d_A}^T E\left(\left[(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\text{vec}\left\{\mathbf{U}_i E\left(\mathcal{G}_{Ai}^T|\mathbf{X}_i, \mathbf{U}_i\right)\mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1}E\left(\mathcal{G}_{Ai}|\mathbf{X}_i, \mathbf{U}_i\right)\mathbf{U}_i^T\right\}\right] \right. \\ & \quad \times \left. \{(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\text{vec}(\mathbf{U}_i\mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1})\}^T|\mathbf{X}_i\right)\mathbf{D}_{d_A} \\ &= \mathbf{O}. \end{aligned}$$

Lastly, the fourth term of  $E(\mathbf{S}_{Ci}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$  can be shown to simplify as follows

$$\begin{aligned} & \frac{1}{4}\mathbf{D}_{d_A}^T E\left[\{(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\text{vec}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\} \right. \\ & \quad \times \left. \{(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\text{vec}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\}^T|\mathbf{X}_i\right]\mathbf{D}_{d_A} \\ &= \frac{1}{4}\mathbf{D}_{d_A}^T E\left[(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\{\text{vec}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1}) + \text{vec}(\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\} \right. \\ & \quad \times \left. \{\text{vec}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1}) + \text{vec}(\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\}^T(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})|\mathbf{X}_i\right]\mathbf{D}_{d_A} \\ &= \frac{1}{4}\mathbf{D}_{d_A}^T E\left[(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\{(\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i)\text{vec}(\mathcal{G}_{Ai}^T) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1})\text{vec}(\mathcal{G}_{Ai})\} \right. \\ & \quad \times \left. \{(\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i)\text{vec}(\mathcal{G}_{Ai}^T) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1})\text{vec}(\mathcal{G}_{Ai})\}^T(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})|\mathbf{X}_i\right]\mathbf{D}_{d_A} \\ &= \frac{1}{4}\mathbf{D}_{d_A}^T E\left[(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\{(\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1})\}\mathcal{G}_{Ai}\mathcal{G}_{Ai}^T \right. \\ & \quad \times \left. \{(\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1})\}^T(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})|\mathbf{X}_i\right]\mathbf{D}_{d_A} \\ &= \frac{1}{4}\mathbf{D}_{d_A}^T E\left[(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\{(\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1})\}E(\mathcal{G}_{Ai}\mathcal{G}_{Ai}^T|\mathbf{X}_i, \mathbf{U}_i) \right. \\ & \quad \times \left. \{(\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1})\}^T(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})|\mathbf{X}_i\right]\mathbf{D}_{d_A} \\ &= \frac{\phi}{4}\mathbf{D}_{d_A}^T E\left[(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\{(\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1})\}\mathcal{H}_{AAi} \right. \\ & \quad \times \left. \{(\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1})\}^T(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})|\mathbf{X}_i\right]\mathbf{D}_{d_A} \\ &= O_P(n^{-1})\mathbf{1}_{d_A^{\boxplus}}\mathbf{1}_{d_A^{\boxplus}}^T. \end{aligned}$$

Putting together all the terms of  $E(\mathbf{S}_{Ci}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$ , we have that

$$E(\mathbf{S}_{Ci}\mathbf{S}_{Ci}^T|\mathbf{X}_i) = \frac{1}{2}\mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1})\mathbf{D}_{d_A} + O_P(n^{-1})\mathbf{1}_{d_A^{\boxplus}}\mathbf{1}_{d_A^{\boxplus}}^T.$$

#### The Expectation of $\mathbf{S}_{Ai}\mathbf{S}_{Bi}^T$ Given $\mathbf{X}_i$

From the previous sections, we have the following approximations

$$\begin{aligned} \mathbf{S}_{Ai} &= \boldsymbol{\Sigma}^{-1} (\mathbf{U}_i + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}) + O_P(n^{-1})\mathbf{1}_{d_A}, \\ \mathbf{S}_{Bi} &= \frac{1}{\phi} (\mathcal{G}_{Bi} - \mathcal{H}_{ABi}^T\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}) + O_P(1)\mathbf{1}_{d_B}. \end{aligned}$$

Therefore, by using the law of total expectation,

$$\begin{aligned}
& E(\mathbf{S}_{Ai}\mathbf{S}_{Bi}^T) \\
&= E \left[ \left\{ \boldsymbol{\Sigma}^{-1}(\mathbf{U}_i + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}) + O_P(n^{-1})\mathbf{1}_{d_A} \right\} \left\{ \frac{1}{\phi}(\mathcal{G}_{Bi} - \mathcal{H}_{ABi}^T\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}) + O_P(1)\mathbf{1}_{d_B} \right\}^T \middle| \mathbf{X}_i \right] \\
&= \frac{1}{\phi} \left\{ E(\boldsymbol{\Sigma}^{-1}\mathbf{U}_i\mathcal{G}_{Bi}^T | \mathbf{X}_i) + E(\boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathcal{G}_{Bi}^T | \mathbf{X}_i) - E(\boldsymbol{\Sigma}^{-1}\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1}\mathcal{H}_{ABi} | \mathbf{X}_i) \right. \\
&\quad \left. - E(\boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1}\mathcal{H}_{ABi} | \mathbf{X}_i) \right\} + O_P(1)\mathbf{1}_{d_A}\mathbf{1}_{d_B}^T \\
&= \frac{1}{\phi} \left\{ E(\boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}\mathcal{H}_{ABi} | \mathbf{X}_i) - E(\boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}\mathcal{H}_{ABi} | \mathbf{X}_i) \right\} + O_P(1)\mathbf{1}_{d_A}\mathbf{1}_{d_B}^T \\
&= O_P(1)\mathbf{1}_{d_A}\mathbf{1}_{d_B}^T
\end{aligned}$$

The Expectation of  $\mathbf{S}_{Ai}\mathbf{S}_{Ci}^T$  Given  $\mathbf{X}_i$

From the previous sections, we have the following approximations where

$$\mathbf{S}_{Ai} = \boldsymbol{\Sigma}^{-1}(\mathbf{U}_i + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}) + O_P(n^{-1})\mathbf{1}_{d_A}$$

and

$$\begin{aligned}
\mathbf{S}_{Ci} &= \frac{1}{2}\mathbf{D}_{d_A}^T [\text{vec} \{ \boldsymbol{\Sigma}^{-1}(\mathbf{U}_i\mathbf{U}_i^T + \mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T) \boldsymbol{\Sigma}^{-1} \} - \text{vec}(\boldsymbol{\Sigma}^{-1})] \\
&\quad + O_P(n^{-1})\mathbf{1}_{d_A^{\#}} \\
&= \frac{1}{2}\mathbf{D}_{d_A}^T \{ \text{vec}(\boldsymbol{\Sigma}^{-1}\mathbf{U}_i\mathbf{U}_i^T\boldsymbol{\Sigma}^{-1}) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \} \\
&\quad + \frac{1}{2}\mathbf{D}_{d_A}^T [\text{vec} \{ \boldsymbol{\Sigma}^{-1}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T) \boldsymbol{\Sigma}^{-1} \}] + O_P(n^{-1})\mathbf{1}_{d_A^{\#}} \\
&= \frac{1}{2}\mathbf{D}_{d_A}^T \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i\mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \} \\
&\quad + \frac{1}{2}\mathbf{D}_{d_A}^T \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T) \} + O_P(n^{-1})\mathbf{1}_{d_A^{\#}}.
\end{aligned}$$

We will deal with each term arising in  $E(\mathbf{S}_{Ai}\mathbf{S}_{Ci}^T | \mathbf{X}_i)$  separately. The first term in  $E(\mathbf{S}_{Ai}\mathbf{S}_{Ci}^T | \mathbf{X}_i)$  can be simplified as follows

$$\begin{aligned}
& \frac{1}{2}E \left[ \boldsymbol{\Sigma}^{-1}\mathbf{U}_i \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i\mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \}^T \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\
&= \frac{1}{2}E \left\{ \boldsymbol{\Sigma}^{-1}\mathbf{U}_i \text{vec}(\mathbf{U}_i\mathbf{U}_i^T)^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} - \boldsymbol{\Sigma}^{-1}\mathbf{U}_i \text{vec}(\boldsymbol{\Sigma}^{-1})^T \mathbf{D}_{d_A} | \mathbf{X}_i \right\} \\
&= \mathbf{O}.
\end{aligned}$$

The second term in  $E(\mathbf{S}_{Ai}\mathbf{S}_{Ci}^T | \mathbf{X}_i)$  can be shown to be

$$\begin{aligned}
& \frac{1}{2}E \left[ \boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai} \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i\mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \}^T \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\
&= \frac{1}{2}E \left[ \boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1} E(\mathcal{G}_{Ai} | \mathbf{X}_i, \mathbf{U}_i) \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i\mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \}^T \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\
&= \mathbf{O}.
\end{aligned}$$

Similarly, the third term in  $E(\mathbf{S}_{Ai}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$  can be shown to be

$$\begin{aligned} & \frac{1}{2}E\left[\boldsymbol{\Sigma}^{-1}\mathbf{U}_i\left\{\text{vec}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1}+\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\right\}^T(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\mathbf{D}_{d_A}|\mathbf{X}_i\right] \\ &= \frac{1}{2}E\left(\boldsymbol{\Sigma}^{-1}\mathbf{U}_i\left[\text{vec}\left\{\mathbf{U}_iE(\mathcal{G}_{Ai}^T|\mathbf{X}_i,\mathbf{U}_i)\mathcal{H}_{AAi}^{-1}+\mathcal{H}_{AAi}^{-1}E(\mathcal{G}_{Ai}|\mathbf{X}_i,\mathbf{U}_i)\mathbf{U}_i^T\right\}\right]^T\right. \\ & \quad \left.\times(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\mathbf{D}_{d_A}|\mathbf{X}_i\right) \\ &= \mathbf{O}. \end{aligned}$$

Lastly, the fourth term in  $E(\mathbf{S}_{Ai}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$  simplifies as follows

$$\begin{aligned} & \frac{1}{2}E\left[\boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\left\{\text{vec}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1}+\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\right\}^T(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\mathbf{D}_{d_A}|\mathbf{X}_i\right] \\ &= \frac{1}{2}E\left[\boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\left\{\text{vec}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1})+\text{vec}(\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\right\}^T(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\mathbf{D}_{d_A}|\mathbf{X}_i\right] \\ &= \frac{1}{2}E\left[\boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\left\{(\mathcal{H}_{AAi}^{-1}\otimes\mathbf{U}_i)\text{vec}(\mathcal{G}_{Ai}^T)+(\mathbf{U}_i\otimes\mathcal{H}_{AAi}^{-1})\text{vec}(\mathcal{G}_{Ai})\right\}^T(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\mathbf{D}_{d_A}|\mathbf{X}_i\right] \\ &= \frac{1}{2}E\left[\boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathcal{G}_{Ai}^T\left\{(\mathcal{H}_{AAi}^{-1}\otimes\mathbf{U}_i)+(\mathbf{U}_i\otimes\mathcal{H}_{AAi}^{-1})\right\}^T(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\mathbf{D}_{d_A}|\mathbf{X}_i\right] \\ &= \frac{1}{2}E\left[\boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}E(\mathcal{G}_{Ai}\mathcal{G}_{Ai}^T|\mathbf{X}_i,\mathbf{U}_i)\left\{(\mathcal{H}_{AAi}^{-1}\otimes\mathbf{U}_i)+(\mathbf{U}_i\otimes\mathcal{H}_{AAi}^{-1})\right\}^T(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\mathbf{D}_{d_A}|\mathbf{X}_i\right] \\ &= \frac{\phi}{2}E\left[\boldsymbol{\Sigma}^{-1}\mathcal{H}_{AAi}^{-1}\mathcal{H}_{AAi}\left\{(\mathcal{H}_{AAi}^{-1}\otimes\mathbf{U}_i)+(\mathbf{U}_i\otimes\mathcal{H}_{AAi}^{-1})\right\}^T(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\mathbf{D}_{d_A}|\mathbf{X}_i\right] \\ &= \frac{\phi}{2}E\left[\boldsymbol{\Sigma}^{-1}\left\{(\mathcal{H}_{AAi}^{-1}\otimes\mathbf{U}_i)+(\mathbf{U}_i\otimes\mathcal{H}_{AAi}^{-1})\right\}^T(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\mathbf{D}_{d_A}|\mathbf{X}_i\right] \\ &= O_P(n^{-1})\mathbf{1}_{d_A}\mathbf{1}_{d_A^\#}^T. \end{aligned}$$

Putting together all the terms of  $E(\mathbf{S}_{Ai}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$ , we have that

$$E(\mathbf{S}_{Ai}\mathbf{S}_{Ci}^T|\mathbf{X}_i) = O_P(n^{-1})\mathbf{1}_{d_A}\mathbf{1}_{d_A^\#}^T.$$

### The Expectation of $\mathbf{S}_{Bi}\mathbf{S}_{Ci}^T$ Given $\mathbf{X}_i$

From the previous sections, we have the following approximations

$$\mathbf{S}_{Bi} = \frac{1}{\phi}(\mathcal{G}_{Bi} - \mathcal{H}_{ABi}^T\mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}) + O_P(1)\mathbf{1}_{d_B}$$

and

$$\begin{aligned} \mathbf{S}_{Ci} &= \frac{1}{2}\mathbf{D}_{d_A}^T\left[\text{vec}\left\{\boldsymbol{\Sigma}^{-1}(\mathbf{U}_i\mathbf{U}_i^T + \mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\boldsymbol{\Sigma}^{-1}\right\} - \text{vec}(\boldsymbol{\Sigma}^{-1})\right] \\ & \quad + O_P(n^{-1})\mathbf{1}_{d_A^\#} \\ &= \frac{1}{2}\mathbf{D}_{d_A}^T\left\{\text{vec}(\boldsymbol{\Sigma}^{-1}\mathbf{U}_i\mathbf{U}_i^T\boldsymbol{\Sigma}^{-1}) - \text{vec}(\boldsymbol{\Sigma}^{-1})\right\} \\ & \quad + \frac{1}{2}\mathbf{D}_{d_A}^T\left[\text{vec}\left\{\boldsymbol{\Sigma}^{-1}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\boldsymbol{\Sigma}^{-1}\right\}\right] + O_P(n^{-1})\mathbf{1}_{d_A^\#} \\ &= \frac{1}{2}\mathbf{D}_{d_A}^T\left\{(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\text{vec}(\mathbf{U}_i\mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1})\right\} \\ & \quad + \frac{1}{2}\mathbf{D}_{d_A}^T\left\{(\boldsymbol{\Sigma}^{-1}\otimes\boldsymbol{\Sigma}^{-1})\text{vec}(\mathbf{U}_i\mathcal{G}_{Ai}^T\mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1}\mathcal{G}_{Ai}\mathbf{U}_i^T)\right\} + O_P(n^{-1})\mathbf{1}_{d_A^\#}. \end{aligned}$$

We will deal with each term arising in  $E(\mathbf{S}_{Bi}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$  separately. The first term in  $E(\mathbf{S}_{Bi}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$  can be simplified as follows

$$\begin{aligned} & \frac{1}{2\phi} E \left[ \mathcal{G}_{Bi} \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \}^T \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= \frac{1}{2\phi} E \left[ E(\mathcal{G}_{Bi} | \mathbf{X}_i, \mathbf{U}_i) \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \}^T \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= \mathbf{O}. \end{aligned}$$

. The second term in  $E(\mathbf{S}_{Bi}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$  can be shown to be

$$\begin{aligned} & -\frac{1}{2\phi} E \left[ \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \}^T \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= -\frac{1}{2\phi} E \left[ \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} E(\mathcal{G}_{Ai} | \mathbf{X}_i, \mathbf{U}_i) \{ (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{U}_i \mathbf{U}_i^T) - \text{vec}(\boldsymbol{\Sigma}^{-1}) \}^T \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= \mathbf{O}. \end{aligned}$$

The third term in  $E(\mathbf{S}_{Bi}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$  can be shown to be

$$\begin{aligned} & \frac{1}{2\phi} E \left[ \mathcal{G}_{Bi} \{ \text{vec}(\mathbf{U}_i \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathbf{U}_i^T) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= \frac{1}{2\phi} E \left[ \mathcal{G}_{Bi} \{ \text{vec}(\mathbf{U}_i \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1}) + \text{vec}(\mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathbf{U}_i^T) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= \frac{1}{2\phi} E \left[ \mathcal{G}_{Bi} \{ (\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) \text{vec}(\mathcal{G}_{Ai}^T) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1}) \text{vec}(\mathcal{G}_{Ai}) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= \frac{1}{2\phi} E \left[ \mathcal{G}_{Bi} \mathcal{G}_{Ai}^T \{ (\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1}) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= \frac{1}{2\phi} E \left[ E(\mathcal{G}_{Bi} \mathcal{G}_{Ai}^T | \mathbf{X}_i, \mathbf{U}_i) \{ (\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1}) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= \frac{1}{2\phi^2} E \left[ \mathcal{H}_{ABi}^T \{ (\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1}) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= O_P(1) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T. \end{aligned}$$

Lastly, the fourth term in  $E(\mathbf{S}_{Bi}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$  simplifies as follows

$$\begin{aligned} & -\frac{1}{2\phi} E \left[ \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \{ \text{vec}(\mathbf{U}_i \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1} + \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathbf{U}_i^T) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= -\frac{1}{2\phi} E \left[ \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \{ \text{vec}(\mathbf{U}_i \mathcal{G}_{Ai}^T \mathcal{H}_{AAi}^{-1}) + \text{vec}(\mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathbf{U}_i^T) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= -\frac{1}{2\phi} E \left[ \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \{ (\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) \text{vec}(\mathcal{G}_{Ai}^T) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1}) \text{vec}(\mathcal{G}_{Ai}) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= -\frac{1}{2\phi} E \left[ \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{G}_{Ai} \mathcal{G}_{Ai}^T \{ (\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1}) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= -\frac{1}{2\phi} E \left[ \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} E(\mathcal{G}_{Ai} \mathcal{G}_{Ai}^T | \mathbf{X}_i, \mathbf{U}_i) \{ (\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1}) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= -\frac{1}{2\phi^2} E \left[ \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{H}_{AAi} \{ (\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1}) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= \frac{1}{2\phi^2} E \left[ \mathcal{H}_{ABi}^T \{ (\mathcal{H}_{AAi}^{-1} \otimes \mathbf{U}_i) + (\mathbf{U}_i \otimes \mathcal{H}_{AAi}^{-1}) \}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A} | \mathbf{X}_i \right] \\ &= O_P(1) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T. \end{aligned}$$

Putting together all the terms of  $E(\mathbf{S}_{Bi}\mathbf{S}_{Ci}^T|\mathbf{X}_i)$ , we have that

$$E(\mathbf{S}_{Ai}\mathbf{S}_{Ci}^T|\mathbf{X}_i) = O_P(1)\mathbf{1}_{d_B}\mathbf{1}_{d_A}^T.$$

#### 4.5.2.8 Treating the Leading Term of the (2,2)-Entry of the Fisher Information Matrix

Note that the leading term of  $\sum_{i=1}^m E(\mathbf{S}_{Bi}\mathbf{S}_{Bi}^T|\mathbf{X}_i)$  is

$$\begin{aligned} & \frac{1}{\phi} \sum_{i=1}^m E(\mathcal{H}_{BBi} - \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{H}_{ABi} | \mathbf{X}_i) \\ &= \frac{mn}{\phi} E \left\{ \frac{1}{mn} \sum_{i=1}^m (\mathcal{H}_{BBi} - \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{H}_{ABi} | \mathbf{X}_i) \right\} \\ &= \frac{mn}{\phi} E \left\{ \frac{1}{mn} \sum_{i=1}^m (\mathcal{H}_{BBi} | \mathbf{X}_i) - \frac{1}{mn} \sum_{i=1}^m (\mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{H}_{ABi} | \mathbf{X}_i) \right\}. \end{aligned} \quad (4.12)$$

Using Lemma 1 from Chapter 2 with  $f(\mathbf{X}_{ij}, \mathbf{U}_i) = b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})$ , we have that the first term in (4.12) can be re-expressed as follows

$$\begin{aligned} & \frac{mn}{\phi} E \left\{ \frac{1}{mn} \sum_{i=1}^m (\mathcal{H}_{BBi} | \mathbf{X}_i) \right\} \\ &= \frac{mn}{\phi} E \left[ \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Bij}^T E \left\{ b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) | \mathbf{X}_i \right\} \right] \\ &= \frac{mn}{\phi} E \left\{ \mathbf{X}_B \mathbf{X}_B^T b''((\boldsymbol{\beta}_A + \mathbf{U})^T \mathbf{X}_A + \boldsymbol{\beta}_B^T \mathbf{X}_B) \right\} + o_P(mn) \mathbf{1}_{d_B}^{\otimes 2}. \end{aligned} \quad (4.13)$$

Now, using Lemma 2 from Chapter 2 with  $f(\mathbf{X}_{ij}, \mathbf{U}_i) = b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij})$ , we have that the second term in (4.12) can be re-expressed as follows

$$\begin{aligned} & \frac{mn}{\phi} E \left\{ \frac{1}{mn} \sum_{i=1}^m (\mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{H}_{ABi} | \mathbf{X}_i) \right\} \\ &= \frac{mn}{\phi} E \left[ \frac{1}{mn} \sum_{i=1}^m \left\{ \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \right\} \right. \\ & \quad \times \left. \left\{ \sum_{j=1}^n \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \right\}^{-1} \right. \\ & \quad \times \left. \left. \left\{ \sum_{j=1}^n \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T b''((\boldsymbol{\beta}_A + \mathbf{U}_i)^T \mathbf{X}_{Aij} + (\boldsymbol{\beta}_B)^T \mathbf{X}_{Bij}) \right\}^T \right| \mathbf{X}_i \right] \end{aligned}$$

The previous expression simplifies to

$$\begin{aligned} & \frac{mn}{\phi} E \left( E \{ \mathbf{X}_B \mathbf{X}_A^T b'' ((\boldsymbol{\beta}_A + \mathbf{U})^T \mathbf{X}_A + \boldsymbol{\beta}_B^T \mathbf{X}_B) \} [E \{ \mathbf{X}_A \mathbf{X}_A^T b'' ((\boldsymbol{\beta}_A + \mathbf{U})^T \mathbf{X}_A + \boldsymbol{\beta}_B^T \mathbf{X}_B) \}]^{-1} \right. \\ & \quad \left. \times E \{ \mathbf{X}_B \mathbf{X}_A^T b'' ((\boldsymbol{\beta}_A + \mathbf{U})^T \mathbf{X}_A + \boldsymbol{\beta}_B^T \mathbf{X}_B) \}^T \right) + o_P(mn) \mathbf{1}_{d_B}^{\otimes 2}. \end{aligned} \quad (4.14)$$

Combining (4.13) and (4.5.2.8), we have

$$\begin{aligned} & \frac{1}{\phi} \sum_{i=1}^m E (\mathcal{H}_{BBi} - \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{H}_{ABi} | \mathbf{X}_i) \\ &= \frac{mn}{\phi} E \left( E \{ \mathbf{X}_B \mathbf{X}_B^T b'' ((\boldsymbol{\beta}_A + \mathbf{U})^T \mathbf{X}_A + \boldsymbol{\beta}_B^T \mathbf{X}_B) \} - E \{ \mathbf{X}_B \mathbf{X}_A^T b'' ((\boldsymbol{\beta}_A + \mathbf{U})^T \mathbf{X}_A + \boldsymbol{\beta}_B^T \mathbf{X}_B) \} \right. \\ & \quad \left. \times [E \{ \mathbf{X}_B \mathbf{X}_A^T b'' ((\boldsymbol{\beta}_A + \mathbf{U})^T \mathbf{X}_A + \boldsymbol{\beta}_B^T \mathbf{X}_B) \}]^{-1} E \{ \mathbf{X}_B \mathbf{X}_A^T b'' ((\boldsymbol{\beta}_A + \mathbf{U})^T \mathbf{X}_A + \boldsymbol{\beta}_B^T \mathbf{X}_B) \}^T \right) \\ & \quad + o_P(mn) \mathbf{1}_{d_B}^{\otimes 2}. \end{aligned}$$

Now since

$$\boldsymbol{\Omega}_{\beta_B}(\mathbf{U}) \equiv E \left\{ b'' ((\boldsymbol{\beta}_A^0 + \mathbf{U})^T \mathbf{X}_A + (\boldsymbol{\beta}_B^0)^T \mathbf{X}_B) \begin{bmatrix} \mathbf{X}_A \mathbf{X}_A^T & \mathbf{X}_A \mathbf{X}_B^T \\ \mathbf{X}_B \mathbf{X}_A^T & \mathbf{X}_B \mathbf{X}_B^T \end{bmatrix} \middle| \mathbf{U} \right\}$$

and

$$\boldsymbol{\Lambda}_{\beta_B} \equiv \left( E \left[ \left\{ \text{lower right } d_B \times d_B \text{ block of } \boldsymbol{\Omega}_{\beta_B}(\mathbf{U})^{-1} \right\}^{-1} \right] \right)^{-1},$$

then we have,

$$\frac{1}{\phi} \sum_{i=1}^m E (\mathcal{H}_{BBi} - \mathcal{H}_{ABi}^T \mathcal{H}_{AAi}^{-1} \mathcal{H}_{ABi} | \mathbf{X}_i) = \frac{mn \boldsymbol{\Lambda}_{\beta_B}^{-1}}{\phi} + o_P(mn) \mathbf{1}_{d_B}^{\otimes 2}.$$

#### 4.5.2.9 The Fisher Information Matrix

Putting together the expressions for the quadratic conditional expectations of the scores from the earlier sections, we have

$$\begin{aligned} & I(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma})) \\ &= \begin{bmatrix} m \boldsymbol{\Sigma}^{-1} + O_P(mn^{-1}) \mathbf{1}_{d_A}^{\otimes 2} & O_P(m) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T & O_P(mn^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T \\ O_P(m) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T & \frac{mn \boldsymbol{\Lambda}_{\beta_B}^{-1}}{\phi} + o_P(mn) \mathbf{1}_{d_B}^{\otimes 2} & O_P(m) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T \\ O_P(mn^{-1}) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T & O_P(m) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T & \frac{m \mathbf{D}_{d_A}^T (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{D}_{d_A}}{2} + O_P(mn^{-1}) \mathbf{1}_{d_B}^{\otimes 2} \end{bmatrix}. \end{aligned}$$

#### 4.5.2.10 The Inverse of the Fisher Information Matrix

To invert the Fisher information matrix, we choose to work with the  $(\boldsymbol{\beta}_A, \text{vech}(\boldsymbol{\Sigma}), \boldsymbol{\beta}_B)$  ordering instead of  $(\boldsymbol{\beta}_A, \boldsymbol{\beta}_B, \text{vech}(\boldsymbol{\Sigma}))$ . A trivial rearrangement of the matrix entries

leads to

$$I(\beta_A, \beta_B, \text{vech}(\Sigma)) = \begin{bmatrix} m\Sigma^{-1} + O_P(mn^{-1})\mathbf{1}_{d_A}^{\otimes 2} & O_P(mn^{-1})\mathbf{1}_{d_A}\mathbf{1}_{d_A^\#}^T & O_P(m)\mathbf{1}_{d_A}\mathbf{1}_{d_B}^T \\ O_P(mn^{-1})\mathbf{1}_{d_A^\#}\mathbf{1}_{d_A}^T & \frac{mD_{d_A}^T(\Sigma^{-1} \otimes \Sigma^{-1})D_{d_A}}{2} + O_P(mn^{-1})\mathbf{1}_{d_A^\#}^{\otimes 2} & O_P(m)\mathbf{1}_{d_A^\#}\mathbf{1}_{d_B}^T \\ O_P(m)\mathbf{1}_{d_B}\mathbf{1}_{d_A}^T & O_P(m)\mathbf{1}_{d_B}\mathbf{1}_{d_A^\#}^T & \frac{mn\Lambda\beta_B^{-1}}{\phi} + O_P(mn)\mathbf{1}_{d_B}^{\otimes 2} \end{bmatrix}.$$

Now, let us partition  $I(\beta_A, \text{vech}(\Sigma), \beta_B)$  as follows

$$I(\beta_A, \text{vech}(\Sigma), \beta_B) = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \text{ where } \mathbf{A}_{21} = \mathbf{A}_{12}^T.$$

The expressions for  $\mathbf{A}_{11}$ ,  $\mathbf{A}_{12}$ ,  $\mathbf{A}_{21}$  and  $\mathbf{A}_{22}$  are currently as follows

$$\begin{aligned} \mathbf{A}_{11} &= \begin{bmatrix} m\Sigma^{-1} + O_P(mn^{-1})\mathbf{1}_{d_A}^{\otimes 2} & O_P(mn^{-1})\mathbf{1}_{d_A}\mathbf{1}_{d_A^\#}^T \\ O_P(mn^{-1})\mathbf{1}_{d_A^\#}\mathbf{1}_{d_A}^T & \frac{mD_{d_A}^T(\Sigma^{-1} \otimes \Sigma^{-1})D_{d_A}}{2} + O_P(mn^{-1})\mathbf{1}_{d_A^\#}^{\otimes 2} \end{bmatrix}, \\ \mathbf{A}_{12} &= O_P(m) \begin{bmatrix} \mathbf{1}_{d_A}\mathbf{1}_{d_B}^T \\ \mathbf{1}_{d_A^\#}\mathbf{1}_{d_B}^T \end{bmatrix}, \\ \mathbf{A}_{21} &= O_P(m) \begin{bmatrix} \mathbf{1}_{d_B}\mathbf{1}_{d_A}^T & \mathbf{1}_{d_B}\mathbf{1}_{d_A^\#}^T \end{bmatrix}, \\ \mathbf{A}_{22} &= \frac{mn\Lambda\beta_B^{-1}}{\phi} + O_P(mn)\mathbf{1}_{d_B}^{\otimes 2}. \end{aligned}$$

Let  $I(\beta_A, \text{vech}(\Sigma), \beta_B)^{-1}$  assume the following form

$$I(\beta_A, \text{vech}(\Sigma), \beta_B)^{-1} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{A}^{11} & \mathbf{A}^{12} \\ \mathbf{A}^{21} & \mathbf{A}^{22} \end{bmatrix} \text{ where } \mathbf{A}^{21} = (\mathbf{A}^{12})^T.$$

Firstly note that

$$\mathbf{A}_{11}^{-1} = \begin{bmatrix} \frac{\Sigma}{m} + O_P(m^{-1}n^{-1})\mathbf{1}_{d_A}^{\otimes 2} & O_P(m^{-1}n^{-1})\mathbf{1}_{d_A}\mathbf{1}_{d_A^\#}^T \\ O_P(m^{-1}n^{-1})\mathbf{1}_{d_A^\#}\mathbf{1}_{d_A}^T & \frac{2D_{d_A}^+(\Sigma \otimes \Sigma)D_{d_A}^+{}^T}{m} + O_P(m^{-1}n^{-1})\mathbf{1}_{d_A^\#}^{\otimes 2} \end{bmatrix}.$$

Also note that

$$\mathbf{A}_{22}^{-1} = \frac{\phi\Lambda\beta_B}{mn} + O_P(m^{-1}n^{-1})\mathbf{1}_{d_B}^{\otimes 2}.$$

Using the result for carrying out block matrix inversion under Result 1, the quantities  $\mathbf{A}^{11}$ ,  $\mathbf{A}^{12}$ ,  $\mathbf{A}^{21}$  and  $\mathbf{A}^{22}$  can be calculated as follows. Firstly we have

$$\mathbf{A}^{11} = \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1}\mathbf{A}_{12}(\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1}.$$

Note that

$$(\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12})^{-1} = \mathbf{A}_{22}^{-1} + \mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} + \dots$$

It follows that

$$\mathbf{A}_{11}^{-1} \mathbf{A}_{12} (\mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12})^{-1} \mathbf{A}_{21} \mathbf{A}_{11}^{-1} = O_P(m^{-1}n^{-1}) \begin{bmatrix} \mathbf{1}_{d_A}^{\otimes 2} & \mathbf{1}_{d_A} \mathbf{1}_{d_A^{\boxplus}}^T \\ \mathbf{1}_{d_A^{\boxplus}} \mathbf{1}_{d_A}^T & \mathbf{1}_{d_A^{\boxplus}}^{\otimes 2} \end{bmatrix}.$$

Therefore, it follows that

$$\mathbf{A}^{11} = \begin{bmatrix} \frac{\Sigma}{m} + O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A}^{\otimes 2} & O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A^{\boxplus}}^T \\ O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A^{\boxplus}} \mathbf{1}_{d_A}^T & \frac{2D_{d_A}^+(\Sigma \otimes \Sigma) D_{d_A}^{+T}}{m} + O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A^{\boxplus}}^{\otimes 2} \end{bmatrix}.$$

Next we have,

$$\mathbf{A}^{12} = -(\mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21})^{-1} \mathbf{A}_{12} \mathbf{A}_{22}^{-1}.$$

Note that

$$(\mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21})^{-1} = \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1} \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21} \mathbf{A}_{11}^{-1} + \dots$$

Therefore we have

$$\mathbf{A}^{12} = O_P(m^{-1}n^{-1}) \begin{bmatrix} \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T \\ \mathbf{1}_{d_A^{\boxplus}} \mathbf{1}_{d_B}^T \end{bmatrix}.$$

Subsequently,

$$\mathbf{A}^{21} = (\mathbf{A}^{12})^T.$$

Therefore we have

$$\mathbf{A}^{21} = O_P(m^{-1}n^{-1}) \begin{bmatrix} \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T & \mathbf{1}_{d_B} \mathbf{1}_{d_A^{\boxplus}}^T \end{bmatrix}.$$

Finally, we obtain

$$\mathbf{A}^{22} = \mathbf{A}_{22}^{-1} + \mathbf{A}_{22}^{-1} \mathbf{A}_{21} (\mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21})^{-1} \mathbf{A}_{12} \mathbf{A}_{22}^{-1}.$$

Note that

$$\mathbf{A}_{22}^{-1} \mathbf{A}_{21} (\mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21})^{-1} \mathbf{A}_{12} \mathbf{A}_{22}^{-1} = O_P(m^{-1}n^{-2}) \mathbf{1}_{d_B}^{\otimes 2}.$$

This leads to

$$\mathbf{A}^{22} = \frac{\phi \Lambda \beta_B}{mn} + o_P(m^{-1}n^{-1}) \mathbf{1}_{d_B}^{\otimes 2}.$$

Using the expressions for  $\mathbf{A}^{11}$ ,  $\mathbf{A}^{12}$ ,  $\mathbf{A}^{21}$  and  $\mathbf{A}^{22}$ , we have the following expression for  $I(\beta_A, \text{vech}(\Sigma), \beta_B)^{-1}$  where

$$I(\beta_A, \text{vech}(\Sigma), \beta_B)^{-1} = \begin{bmatrix} \frac{\Sigma}{m} + O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A}^{\otimes 2} & O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_A^{\boxplus}}^T & O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T \\ O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A^{\boxplus}} \mathbf{1}_{d_A}^T & \frac{2D_{d_A}^+(\Sigma \otimes \Sigma) D_{d_A}^{+T}}{m} + O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A^{\boxplus}}^{\otimes 2} & O_P(m^{-1}n^{-1}) \mathbf{1}_{d_A^{\boxplus}} \mathbf{1}_{d_B}^T \\ O_P(m^{-1}n^{-1}) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T & O_P(m^{-1}n^{-1}) \mathbf{1}_{d_B} \mathbf{1}_{d_A^{\boxplus}}^T & \frac{\phi \Lambda \beta_B}{mn} + o_P(m^{-1}n^{-1}) \mathbf{1}_{d_B}^{\otimes 2} \end{bmatrix}.$$



The expression for the inverse of the Fisher information matrix can be also written as follows

$$\begin{aligned} & I(\boldsymbol{\beta}_A, \text{vech}(\boldsymbol{\Sigma}), \boldsymbol{\beta}_B)^{-1} \\ &= I(\boldsymbol{\beta}_A, \text{vech}(\boldsymbol{\Sigma}), \boldsymbol{\beta}_B)_\infty^{-1} + \frac{1}{mn} \begin{bmatrix} O_P(1)\mathbf{1}_{d_A}^{\otimes 2} & O_P(1)\mathbf{1}_{d_A}\mathbf{1}_{d_A^{\text{ff}}}^T & O_P(1)\mathbf{1}_{d_A}\mathbf{1}_{d_B}^T \\ O_P(1)\mathbf{1}_{d_A^{\text{ff}}}\mathbf{1}_{d_A}^T & O_P(1)\mathbf{1}_{d_A^{\text{ff}}}^{\otimes 2} & O_P(1)\mathbf{1}_{d_A^{\text{ff}}}\mathbf{1}_{d_B}^T \\ O_P(1)\mathbf{1}_{d_B}\mathbf{1}_{d_A}^T & O_P(1)\mathbf{1}_{d_B}\mathbf{1}_{d_A^{\text{ff}}}^T & O_P(1)\mathbf{1}_{d_B}^{\otimes 2} \end{bmatrix}, \end{aligned}$$

where

$$I(\boldsymbol{\beta}_A, \text{vech}(\boldsymbol{\Sigma}), \boldsymbol{\beta}_B)_\infty^{-1} = \begin{bmatrix} \frac{\boldsymbol{\Sigma}}{m} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \frac{2D_{d_A}^+(\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma})D_{d_A}^{+T}}{m} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \frac{\phi \boldsymbol{\Lambda} \boldsymbol{\beta}_B}{mn} \end{bmatrix}.$$

#### 4.5.2.11 Derivation of the Final Asymptotic Normality Result for Generalized Response Linear Mixed Models

For a matrix  $\mathbf{M}$  let

$$\|\mathbf{M}\|_F = \sqrt{\text{tr}(\mathbf{M}^T \mathbf{M})}$$

denote the *Frobenius norm* of  $\mathbf{M}$ .

For regular likelihood situations, from standard results concerning asymptotic normality of maximum likelihood estimators we have

$$\{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)^{-1}\}^{-1/2}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{I})$$

where  $\widehat{\boldsymbol{\theta}} = [\widehat{\boldsymbol{\beta}}_A^T \text{vech}(\widehat{\boldsymbol{\Sigma}})^T \widehat{\boldsymbol{\beta}}_B^T]^T$  and  $\boldsymbol{\theta}^0 = [(\boldsymbol{\beta}_A^0)^T \{\text{vech}(\boldsymbol{\Sigma}^0)\}^T (\boldsymbol{\beta}_B^0)^T]^T$ . On the other hand, general quasi-likelihood situations require asymptotic normality theory as treated in, for example, Section 5.3 of the book *Asymptotic Statistics* (see van der Vaart (1998)). Therefore, for all  $(d_A + d_A^{\text{ff}} + d_B) \times 1$  vectors  $\mathbf{a} \neq \mathbf{0}$  we have

$$\mathbf{a}^T \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)^{-1}\}^{-1/2}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{a}^T \mathbf{a}).$$

Note that

$$\begin{aligned} & \mathbf{a}^T \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)^{-1}\}^{-1/2}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0) \\ &= \mathbf{a}^T \left[ \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty^{-1}\}^{-1/2} + \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)^{-1}\}^{-1/2} \right. \\ & \quad \left. - \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty^{-1}\}^{-1/2} \right] (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0) \\ &= \mathbf{a}^T \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty^{-1}\}^{-1/2}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0) \\ & \quad + \mathbf{a}^T \left[ \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)^{-1}\}^{-1/2} - \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty^{-1}\}^{-1/2} \right] (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0). \end{aligned}$$

As a consequence

$$\mathbf{a}^T \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{-1/2} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0) + r_{mn}(\mathbf{a}) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{a}^T \mathbf{a}) \quad (4.15)$$

with

$$\begin{aligned} r_{mn}(\mathbf{a}) &= \mathbf{a}^T \left[ \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{-1/2} - \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{-1/2} \right] (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0) \\ &= \mathbf{a}^T \left[ \mathbf{I} - \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{-1/2} \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{1/2} \right] \\ &\quad \times \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{-1/2} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0) \\ &= \left( - \left[ \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{-1/2} \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{1/2} - \mathbf{I} \right]^T \mathbf{a} \right)^T \mathbf{Z} \end{aligned}$$

where  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I}_{d_A + d_A^\# + d_B})$ . Next note that using the matrix norm properties  $\|-\mathbf{A}\| = \|\mathbf{A}\|$  and  $\|\mathbf{A}\mathbf{B}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$  for any matrices  $\mathbf{A}$  and  $\mathbf{B}$  and the fact that  $\|\mathbf{M}^T\|_F = \|\mathbf{M}\|_F$  for any matrix  $\mathbf{M}$ , we have

$$\begin{aligned} &\left\| - \left[ \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{-1/2} \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{1/2} - \mathbf{I} \right]^T \mathbf{a} \right\|_F \\ &\leq \left\| \left[ \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{-1/2} \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{1/2} - \mathbf{I} \right]^T \right\|_F \|\mathbf{a}\|_F \quad (4.16) \\ &= \left\| \left[ \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{-1/2} \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{1/2} - \mathbf{I} \right] \right\|_F \|\mathbf{a}\|_F. \end{aligned}$$

Our next aim is to establish that

$$\left\| \left[ \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{-1/2} \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty\}^{1/2} - \mathbf{I} \right] \right\|_F \xrightarrow{P} 0 \quad (4.17)$$

Recall that

$$\begin{aligned} &I(\boldsymbol{\beta}_A, \text{vech}(\boldsymbol{\Sigma}), \boldsymbol{\beta}_B)^{-1} \\ &= I(\boldsymbol{\beta}_A, \text{vech}(\boldsymbol{\Sigma}), \boldsymbol{\beta}_B)_\infty^{-1} + \frac{1}{mn} \begin{bmatrix} O_P(1) \mathbf{1}_{d_A}^{\otimes 2} & O_P(1) \mathbf{1}_{d_A} \mathbf{1}_{d_A^\#}^T & O_P(1) \mathbf{1}_{d_A} \mathbf{1}_{d_B}^T \\ O_P(1) \mathbf{1}_{d_A^\#} \mathbf{1}_{d_A}^T & O_P(1) \mathbf{1}_{d_A^\#}^{\otimes 2} & O_P(1) \mathbf{1}_{d_A^\#} \mathbf{1}_{d_B}^T \\ O_P(1) \mathbf{1}_{d_B} \mathbf{1}_{d_A}^T & O_P(1) \mathbf{1}_{d_B} \mathbf{1}_{d_A^\#}^T & o_P(1) \mathbf{1}_{d_B}^{\otimes 2} \end{bmatrix}, \end{aligned}$$

where

$$I(\boldsymbol{\beta}_A, \text{vech}(\boldsymbol{\Sigma}), \boldsymbol{\beta}_B)_\infty^{-1} = \begin{bmatrix} \frac{\boldsymbol{\Sigma}}{m} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \frac{2D_{d_A}^+(\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}) D_{d_A}^+{}^T}{m} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \frac{\phi \boldsymbol{\Lambda}_{\boldsymbol{\beta}_B}}{mn} \end{bmatrix}.$$

so that

$$\left\{ I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty \right\}^{-1/2} \left\{ I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty \right\}^{1/2} = \mathbf{M}_{n,\infty}^{-1/2} \mathbf{M}_n^{1/2}$$

with

$$\mathbf{M}_{n,\infty} = m \left\{ I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty \right\} \quad \text{and} \quad \mathbf{M}_n = m \left\{ I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_\infty \right\}^{-1}.$$

Therefore, Lemma 3 from Chapter 2 applies with

$$p = d_A + d_A^{\boxplus}, \quad \mathbf{K} = \begin{bmatrix} \boldsymbol{\Sigma} & \mathbf{0} \\ \mathbf{0} & 2\mathbf{D}_{d_A}^+ (\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}) \mathbf{D}_{d_A}^{+T} \end{bmatrix}, \quad q = d_B \quad \text{and} \quad \mathbf{L} = \frac{\phi \boldsymbol{\Lambda} \boldsymbol{\beta}_B}{n}.$$

Therefore (4.17) holds. It then follows from (4.16) and (4.17) that

$$\left[ \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_{\infty}^{-1}\}^{-1/2} \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)^{-1}\}^{1/2} - \mathbf{I} \right] \mathbf{a} \xrightarrow{P} 0.$$

Application of Slutsky's Theorem then gives  $r_{mn}(\mathbf{a}) \xrightarrow{P} 0$ . From (4.15) and another application of Slutsky's Theorem we have

$$\mathbf{a}^T \{I(\boldsymbol{\beta}_A^0, \text{vech}(\boldsymbol{\Sigma}^0), \boldsymbol{\beta}_B^0)_{\infty}^{-1}\}^{-1/2} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{a}^T \mathbf{a}).$$

It then follows from the Cramér-Wold Device and the Continuous Mapping Theorem that

$$\sqrt{m} \begin{bmatrix} \hat{\boldsymbol{\beta}}_A - \boldsymbol{\beta}_A^0 \\ \sqrt{n} (\hat{\boldsymbol{\beta}}_B - \boldsymbol{\beta}_B^0) \\ \text{vech}(\hat{\boldsymbol{\Sigma}}) - \text{vech}(\boldsymbol{\Sigma}^0) \end{bmatrix} \xrightarrow{\mathcal{D}} N \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}^0 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \phi \boldsymbol{\Lambda} \boldsymbol{\beta}_B & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & 2\mathbf{D}_{d_A}^+ (\boldsymbol{\Sigma}^0 \otimes \boldsymbol{\Sigma}^0) \mathbf{D}_{d_A}^{+T} \end{bmatrix} \right).$$

### 4.5.3 The Reciprocal Dispersion Parameter Fisher Information Block for Gamma Responses

In this appendix, we derive the block of the Fisher information matrix for the parameter  $\psi \equiv 1/\phi$  where  $\phi$  is the dispersion parameter. We start with the general response situation and, later, focus on the Gamma case. With notational simplicity in mind we treat the  $d_A = d_B = 1$  case with  $\boldsymbol{\Sigma} = \sigma^2$  and  $n_i = n$ ,  $1 \leq i \leq m$ . These restrictions will not affect the  $\psi$  contribution to the Fisher information matrix.

#### 4.5.3.1 The Conditional Density Function

The conditional density function for the  $i$ th group is

$$\begin{aligned} p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i) &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp \left[ -\frac{u^2}{2\sigma^2} + \sum_{j=1}^n \left\{ \psi \left( Y_{ij}(\beta_0 + \beta_1 X_{ij} + u) \right. \right. \right. \\ &\quad \left. \left. \left. - b(\beta_0 + \beta_1 X_{ij} + u) + c(Y_{ij}) \right) + d(Y_{ij}, \psi) \right\} \right] du \\ &= (2\pi\sigma^2)^{-1/2} \exp \left( \sum_{j=1}^n \left[ \psi \{ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + c(Y_{ij}) \} + d(Y_{ij}, \psi) \right] \right) \\ &\quad \times \int_{-\infty}^{\infty} \exp \left[ -\frac{u^2}{2\sigma^2} + \psi \sum_{j=1}^n \{ Y_{ij}u - b(\beta_0 + \beta_1 X_{ij} + u) \} \right] du. \end{aligned}$$

In the Gamma response case,

$$c(y) = \log(y) \quad \text{and} \quad d(y, \psi) = \psi \log(\psi) - \log\{\Gamma(\psi)\} - \log(y).$$

Therefore

$$\frac{\partial d(y, \psi)}{\partial \psi} = \mathbf{q}(\psi) \quad \text{and} \quad \frac{\partial^2 d(y, \psi)}{\partial \psi^2} = (1/\psi) - \text{trigamma}(\psi),$$

where

$$\mathbf{q}(x) \equiv 1 + \log(x) - \text{digamma}(x).$$

Note that the model is such that

$$Y_{ij} | \mathbf{X}_i, U_i \stackrel{\text{ind.}}{\sim} \text{Gamma}(\psi, \psi/b'(\beta_0 + \beta_1 X_{ij} + U_i)). \quad (4.18)$$

Since  $b'(x) = -1/x$  the statement (4.18) is equivalent to

$$Y_{ij} | \mathbf{X}_i, U_i \stackrel{\text{ind.}}{\sim} \text{Gamma}(\psi, -\psi(\beta_0 + \beta_1 X_{ij} + U_i)). \quad (4.19)$$

Using the result that

$$X \sim \text{Gamma}(\kappa, \lambda) \quad \text{implies} \quad E\{\log(X)\} = \text{digamma}(\kappa) - \log(\lambda)$$

and using  $b(x) = -\log(-x)$ , we have,

$$\begin{aligned} E\{c(Y_{ij}) | U_i, \mathbf{X}_i\} &= E\{\log(Y_{ij}) | U_i, \mathbf{X}_i\} \\ &= \text{digamma}(\psi) - \log[\{-(\beta_0 + \beta_1 X_{ij} + U_i)\}\psi] \\ &= -\log\{-(\beta_0 + \beta_1 X_{ij} + U_i)\} + \text{digamma}(\psi) - \log(\psi) \\ &= b(\beta_0 + \beta_1 X_{ij} + U_i) + \text{digamma}(\psi) - \log(\psi). \end{aligned}$$

Therefore, if we define

$$\mathcal{A}_i \equiv \sum_{j=1}^n b(\beta_0 + \beta_1 X_{ij} + U_i)$$

then

$$\sum_{j=1}^n E\{\log(Y_{ij}) | U_i, \mathbf{X}_i\} = \mathcal{A}_i - n\{\log(\psi) - \text{digamma}(\psi)\}.$$

Also, if we define

$$\mathbf{q}(x) \equiv 1 + \log(x) - \text{digamma}(x)$$

then

$$E\{\log(Y_{ij}) | U_i, \mathbf{X}_i\} = b(\beta_0 + \beta_1 X_{ij} + U_i) + 1 - \mathbf{q}(\psi)$$

and

$$\sum_{j=1}^n E\{\log(Y_{ij}) | U_i, \mathbf{X}_i\} = \mathcal{A}_i + n\{1 - \mathbf{q}(\psi)\}.$$

In the Gamma case

$$d(y, \psi) = \psi \log(\psi) - \log\{\Gamma(\psi)\} - \log(y) \quad \text{implies} \quad \frac{\partial d(y, \psi)}{\partial \psi} = \mathbf{q}(\psi).$$

Therefore,

$$\sum_{j=1}^n \frac{\partial}{\partial \psi} d(Y_{ij}, \psi) = n\mathbf{q}(\psi).$$

Next consider the problem of obtaining an expression for  $E(Y_{ij}^2 | \mathbf{X}_i, U_i)$ .

$$\begin{aligned} E(Y_{ij}^2 | \mathbf{X}_i, U_i) &= \text{Var}(Y_{ij} | \mathbf{X}_i, U_i) + \{E(Y_{ij} | \mathbf{X}_i, U_i)\}^2 \\ &= \frac{1}{\psi(\beta_0 + \beta_1 X_{ij} + U_i)^2} + \left( \frac{1}{-(\beta_0 + \beta_1 X_{ij} + U_i)} \right)^2 \\ &= \left( \frac{1}{\psi} + 1 \right) \frac{1}{(\beta_0 + \beta_1 X_{ij} + U_i)^2} \\ &= \left( \frac{1}{\psi} + 1 \right) \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2. \end{aligned}$$

Later on we also need expressions for

$$E\{Y_{ij} \log(Y_{ij}) | \mathbf{X}_i, U_i\} \quad \text{and} \quad E\{[\log(Y_{ij})]^2 | \mathbf{X}_i, U_i\}.$$

As a prelude to obtaining these expressions we consider

$$X \sim \text{Gamma}(\kappa, \lambda).$$

Then note that

$$E\{X \log(X)\} = \{\kappa \text{digamma}(\kappa) + 1 - \kappa \log(\lambda)\} / \lambda$$

and

$$E\{[\log(X)]^2\} = \text{trigamma}(\kappa) + \{\text{digamma}(\kappa) - \log(\lambda)\}^2.$$

The first of these results leads to (with  $\kappa = \psi$  and  $\lambda = -\psi(\beta_0 + \beta_1 X_{ij} + U_i)$ ),

$$\begin{aligned} &E\{Y_{ij} \log(Y_{ij}) | \mathbf{X}_i, U_i\} \\ &= \frac{1}{\psi\{-(\beta_0 + \beta_1 X_{ij} + U_i)\}} \left( \psi \text{digamma}(\psi) + 1 - \psi \log[\psi\{-(\beta_0 + \beta_1 X_{ij} + U_i)\}] \right) \\ &= \frac{1}{\psi\{-(\beta_0 + \beta_1 X_{ij} + U_i)\}} \left[ \psi \text{digamma}(\psi) + 1 - \psi \log(\psi) - \psi \log\{-(\beta_0 + \beta_1 X_{ij} + U_i)\} \right] \\ &= \frac{1}{\psi\{-(\beta_0 + \beta_1 X_{ij} + U_i)\}} \left\{ \psi \text{digamma}(\psi) + 1 - \psi \log(\psi) + \psi b(\beta_0 + \beta_1 X_{ij} + U_i) \right\} \\ &= b'(\beta_0 + \beta_1 X_{ij} + U_i) \left\{ \frac{1}{\psi} + \text{digamma}(\psi) - \log(\psi) + b(\beta_0 + \beta_1 X_{ij} + U_i) \right\} \\ &= b'(\beta_0 + \beta_1 X_{ij} + U_i) \left\{ \frac{1}{\psi} + 1 - \mathbf{q}(\psi) + b(\beta_0 + \beta_1 X_{ij} + U_i) \right\}. \end{aligned}$$

The second of these results leads to (with  $\kappa = \psi$  and  $\lambda = -\psi(\beta_0 + \beta_1 X_{ij} + U_i)$ )

$$\begin{aligned} E[\{\log(Y_{ij})\}^2 | \mathbf{X}_i, U_i] &= \text{trigamma}(\psi) + \left( \text{digamma}(\psi) - \log[\psi\{-(\beta_0 + \beta_1 X_{ij} + U_i)\}] \right)^2 \\ &= \text{trigamma}(\psi) + [\text{digamma}(\psi) - \log(\psi) - \log\{-(\beta_0 + \beta_1 X_{ij} + U_i)\}]^2 \\ &= \text{trigamma}(\psi) + \{1 - \mathfrak{q}(\psi) + b(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \\ &= \text{trigamma}(\psi) + \{b(\beta_0 + \beta_1 X_{ij} + U_i) + 1 - \mathfrak{q}(\psi)\}^2. \end{aligned}$$

#### 4.5.3.2 The Score of the Reciprocal Dispersion Parameter

The score of  $\psi$  is

$$\begin{aligned} S_{3i} &\equiv \frac{\partial \log p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i)}{\partial \psi} \\ &= \sum_{j=1}^n \left[ \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij})\} + \mathfrak{q}(\psi) \right] \\ &\quad + \frac{\int_{-\infty}^{\infty} \exp \left[ -\frac{u^2}{2\sigma^2} + \psi \sum_{j=1}^n \{Y_{ij}u - b(\beta_0 + \beta_1 X_{ij} + u)\} \right] \sum_{j=1}^n \{Y_{ij}u - b(\beta_0 + \beta_1 X_{ij} + u)\} du}{\int_{-\infty}^{\infty} \exp \left[ -\frac{u^2}{2\sigma^2} + \psi \sum_{j=1}^n \{Y_{ij}u - b(\beta_0 + \beta_1 X_{ij} + u)\} \right] du} \\ &= \sum_{j=1}^n \left[ \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij})\} + \mathfrak{q}(\psi) \right] + \frac{\int_{-\infty}^{\infty} b_N(u) \exp\{-nh_N(u)\} du}{\int_{-\infty}^{\infty} b_D(u) \exp\{-nh_N(u)\} du} \end{aligned}$$

where

$$\begin{aligned} h_N(u) &\equiv -\frac{\psi}{n} \sum_{j=1}^n \{Y_{ij}u - b(\beta_0 + \beta_1 X_{ij} + u)\}, \\ b_N(u) &\equiv \sum_{j=1}^n \{Y_{ij}u - b(\beta_0 + \beta_1 X_{ij} + u)\} \exp\left(-\frac{u^2}{2\sigma^2}\right) \quad \text{and} \quad b_D(u) \equiv \exp\left(-\frac{u^2}{2\sigma^2}\right). \end{aligned}$$

#### The First Term of the $S_{3i}$ Integral

The first term on the right-hand side of (2.6) of Tierney et al. (1989) is

$$g(U_i^*)$$

where

$$g(u) \equiv b_N(u)/b_D(u) = \sum_{j=1}^n \{Y_{ij}u - b(\beta_0 + \beta_1 X_{ij} + u)\}.$$

Using steps similar to those given for approximating the score of  $\beta_B$  in Subsubsection 4.5.2.5, we get

$$g(U_i^*) = U_i Y_{i\bullet} - \mathcal{A}_i + O_P(1)$$

where

$$Y_{i\bullet} \equiv \sum_{j=1}^n Y_{ij} \quad \text{and} \quad \mathcal{A}_i \equiv \sum_{j=1}^n b(\beta_0 + \beta_1 X_{ij} + U_i).$$

Also define

$$\mathcal{A}'_i \equiv \sum_{j=1}^n b'(\beta_0 + \beta_1 X_{ij} + U_i).$$

#### 4.5.3.3 Computing the Fisher Information Block for the Reciprocal Dispersion Parameter

Define

$$\mathcal{B}_i \equiv \sum_{j=1}^n \left\{ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathfrak{q}(\psi) \right\}.$$

Then

$$S_{3i} = \mathcal{B}_i + U_i Y_{i\bullet} - \mathcal{A}_i + O_P(1).$$

Therefore,

$$E(S_{3i}^2 | \mathbf{X}_i) \approx T_A + T_B + T_C + T_D + T_E + T_F$$

where

$$T_A = E(\mathcal{B}_i^2 | \mathbf{X}_i),$$

$$T_B = E(U_i^2 Y_{i\bullet}^2 | \mathbf{X}_i)$$

$$T_C = E(\mathcal{A}_i^2 | \mathbf{X}_i),$$

$$T_D = 2E(\mathcal{B}_i U_i Y_{i\bullet} | \mathbf{X}_i),$$

$$T_E = -2E(\mathcal{A}_i \mathcal{B}_i | \mathbf{X}_i),$$

$$\text{and} \quad T_F = -2E(U_i \mathcal{A}_i Y_{i\bullet} | \mathbf{X}_i).$$

#### Treatment of $T_A$

In the Gamma case

$$b'(x) = -1/x \quad \text{and} \quad c(x) = \log(x).$$

Also,

$$\frac{\partial d(Y_{ij}, \psi)}{\partial \psi} = \mathfrak{q}(\psi).$$

Hence

$$\mathcal{B}_i = \sum_{j=1}^n \left\{ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathfrak{q}(\psi) \right\}$$

and

$$\begin{aligned} \mathcal{B}_i^2 &= \sum_{j=1}^n \sum_{j'=1}^n \left[ \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij})\} + \mathbf{q}(\psi) \right] \left[ \{Y_{ij'}(\beta_0 + \beta_1 X_{ij'}) + \log(Y_{ij'})\} + \mathbf{q}(\psi) \right] \\ &= \sum_{j=1}^n \left[ \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij})\} + \mathbf{q}(\psi) \right]^2 \\ &\quad + \sum_{j \neq j'}^n \left[ \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij})\} + \mathbf{q}(\psi) \right] \left[ \{Y_{ij'}(\beta_0 + \beta_1 X_{ij'}) + \log(Y_{ij'})\} + \mathbf{q}(\psi) \right]. \end{aligned}$$

Therefore

$$T_A = E(\mathcal{B}_i^2 | \mathbf{X}_i) = \mathbf{r}_1(\mathbf{X}_i) + \mathbf{r}_2(\mathbf{X}_i)$$

where

$$\mathbf{r}_1(\mathbf{X}_i) \equiv \sum_{j \neq j'}^n E \left( \left[ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right] \left[ Y_{ij'}(\beta_0 + \beta_1 X_{ij'}) + \log(Y_{ij'}) + \mathbf{q}(\psi) \right] \middle| \mathbf{X}_i \right)$$

and

$$\mathbf{r}_2(\mathbf{X}_i) \equiv \sum_{j=1}^n E \left[ \left\{ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right\}^2 \middle| \mathbf{X}_i \right].$$

#### Treatment of $\mathbf{r}_1(\mathbf{X}_i)$

The  $(j, j')$ th term in the  $\mathbf{r}_1(\mathbf{X}_i)$  summation is

$$\begin{aligned} &E \left\{ E \left( \left[ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right] \left[ Y_{ij'}(\beta_0 + \beta_1 X_{ij'}) + \log(Y_{ij'}) + \mathbf{q}(\psi) \right] \middle| \mathbf{X}_i, U_i \right) \middle| \mathbf{X}_i \right\} \\ &= E \left\{ E \left( \left[ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right] \middle| \mathbf{X}_i, U_i \right) \right. \\ &\quad \left. \times E \left( \left[ Y_{ij'}(\beta_0 + \beta_1 X_{ij'}) + \log(Y_{ij'}) + \mathbf{q}(\psi) \right] \middle| \mathbf{X}_i, U_i \right) \middle| \mathbf{X}_i \right\} \end{aligned}$$

Next note that

$$\begin{aligned} &E \left( \left[ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right] \middle| \mathbf{X}_i, U_i \right) \\ &= b'(\beta_0 + \beta_1 X_{ij} + U_i)(\beta_0 + \beta_1 X_{ij}) + b(\beta_0 + \beta_1 X_{ij} + U_i) + 1 - \mathbf{q}(\psi) + \mathbf{q}(\psi) \\ &= b'(\beta_0 + \beta_1 X_{ij} + U_i)(\beta_0 + \beta_1 X_{ij} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij} + U_i) + b(\beta_0 + \beta_1 X_{ij} + U_i) + 1 \\ &= b(\beta_0 + \beta_1 X_{ij} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij} + U_i). \end{aligned}$$

Hence,

$$\begin{aligned} \mathbf{r}_1(\mathbf{X}_i) &= \sum_{j \neq j'}^n E \left[ \{b(\beta_0 + \beta_1 X_{ij} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\} \right. \\ &\quad \left. \times \{b(\beta_0 + \beta_1 X_{ij'} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij'} + U_i)\} \middle| \mathbf{X}_i \right]. \end{aligned}$$



Treatment of  $\tau_2(\mathbf{X}_i)$ 

The  $j$ th term in the  $\tau_2(\mathbf{X}_i)$  summation is

$$E \left( E \left[ \left\{ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right\}^2 \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right).$$

Next we expand out the inner conditional expectation:

$$\begin{aligned} & E \left[ \left\{ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right\}^2 \middle| \mathbf{X}_i, U_i \right] \\ &= (\beta_0 + \beta_1 X_{ij})^2 E(Y_{ij}^2 | \mathbf{X}_i, U_i) + E[\{\log(Y_{ij})\}^2 | \mathbf{X}_i, U_i] + \mathbf{q}(\psi)^2 \\ &\quad + 2(\beta_0 + \beta_1 X_{ij}) E\{Y_{ij} \log(Y_{ij}) | \mathbf{X}_i, U_i\} \\ &\quad + 2(\beta_0 + \beta_1 X_{ij}) \mathbf{q}(\psi) E(Y_{ij} | \mathbf{X}_i, U_i) + 2\mathbf{q}(\psi) E\{\log(Y_{ij}) | \mathbf{X}_i, U_i\} \\ &= (\beta_0 + \beta_1 X_{ij})^2 \left( \frac{1}{\psi} + 1 \right) \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \\ &\quad + \text{trigamma}(\psi) + \{b(\beta_0 + \beta_1 X_{ij} + U_i) + 1 - \mathbf{q}(\psi)\}^2 \\ &\quad + \mathbf{q}(\psi)^2 + 2(\beta_0 + \beta_1 X_{ij}) b'(\beta_0 + \beta_1 X_{ij} + U_i) \left\{ \frac{1}{\psi} + 1 - \mathbf{q}(\psi) + b(\beta_0 + \beta_1 X_{ij} + U_i) \right\} \\ &\quad + 2(\beta_0 + \beta_1 X_{ij}) \mathbf{q}(\psi) b'(\beta_0 + \beta_1 X_{ij} + U_i) + 2\mathbf{q}(\psi) \{b(\beta_0 + \beta_1 X_{ij} + U_i) + 1 - \mathbf{q}(\psi)\}. \end{aligned}$$

The first of these terms is

$$\begin{aligned} & \left( \frac{1}{\psi} + 1 \right) (\beta_0 + \beta_1 X_{ij} + U_i - U_i)^2 \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \\ &= \left( \frac{1}{\psi} + 1 \right) \{(\beta_0 + \beta_1 X_{ij} + U_i)^2 - 2U_i(\beta_0 + \beta_1 X_{ij} + U_i) + U_i^2\} \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2. \\ &= \left( \frac{1}{\psi} + 1 \right) \left[ 1 + 2U_i b'(\beta_0 + \beta_1 X_{ij} + U_i) + U_i^2 \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \right]. \end{aligned}$$

The sixth of these terms is

$$2(\beta_0 + \beta_1 X_{ij} + U_i - U_i) \mathbf{q}(\psi) b'(\beta_0 + \beta_1 X_{ij} + U_i) = -2\mathbf{q}(\psi) \{1 + U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\}.$$

With similar steps, the fifth of these terms is

$$-2 \left\{ \frac{1}{\psi} + 1 - \mathbf{q}(\psi) + b(\beta_0 + \beta_1 X_{ij} + U_i) \right\} \{1 + U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\}.$$

The sum of the fifth and sixth terms is then

$$-2 \left\{ \frac{1}{\psi} + 1 + b(\beta_0 + \beta_1 X_{ij} + U_i) \right\} \{1 + U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\}.$$

Assembling the above results we have

$$\begin{aligned}
T_A &= n \operatorname{trigamma}(\psi) + n\{\mathfrak{q}(\psi)\}^2 \\
&+ \sum_{j \neq j'}^n E \left[ \{b(\beta_0 + \beta_1 X_{ij} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\} \{b(\beta_0 + \beta_1 X_{ij'} + U_i) \right. \\
&\quad \left. - U_i b'(\beta_0 + \beta_1 X_{ij'} + U_i)\} | \mathbf{X}_i \right] \\
&+ \left( \frac{1}{\psi} + 1 \right) \sum_{j=1}^n E \left( [1 + 2U_i b'(\beta_0 + \beta_1 X_{ij} + U_i) + U_i^2 \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2] | \mathbf{X}_i \right) \\
&+ \sum_{j=1}^n E \left( \{b(\beta_0 + \beta_1 X_{ij} + U_i) + 1 - \mathfrak{q}(\psi)\}^2 | \mathbf{X}_i \right) \\
&- 2 \sum_{j=1}^n E \left[ \left\{ \left( \frac{1}{\psi} + 1 \right) + b(\beta_0 + \beta_1 X_{ij} + U_i) \right\} \{1 + U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\} | \mathbf{X}_i \right] \\
&+ 2\mathfrak{q}(\psi) \sum_{j=1}^n E \left[ \{b(\beta_0 + \beta_1 X_{ij} + U_i) + 1 - \mathfrak{q}(\psi)\} | \mathbf{X}_i \right].
\end{aligned}$$

A next useful step (for cancellation purposes) is to expand out the second, third, fourth, fifth and sixth terms of this expression for  $T_A$ .

Expansion of the second term of  $T_A$

For the second term of  $T_A$  we have

$$\begin{aligned}
&\sum_{j \neq j'}^n E \left[ \{b(\beta_0 + \beta_1 X_{ij} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\} \right. \\
&\quad \left. \times \{b(\beta_0 + \beta_1 X_{ij'} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij'} + U_i)\} | \mathbf{X}_i \right] \\
&= \sum_{j \neq j'}^n E \left\{ b(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} \\
&\quad - \sum_{j \neq j'}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} \\
&\quad - \sum_{j \neq j'}^n E \left\{ U_i b'(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} \\
&\quad + \sum_{j \neq j'}^n E \left\{ U_i^2 b'(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\}.
\end{aligned}$$

Expansion of the third term of  $T_A$ 

For the third term of  $T_A$  we have

$$\begin{aligned} & \left(\frac{1}{\psi} + 1\right) \sum_{j=1}^n E \left( [1 + 2U_i b'(\beta_0 + \beta_1 X_{ij} + U_i) + U_i^2 \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2] \middle| \mathbf{X}_i \right) \\ &= n \left(\frac{1}{\psi} + 1\right) + 2 \left(\frac{1}{\psi} + 1\right) \sum_{j=1}^n E \left\{ U_i b'(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right\} \\ & \quad + \left(\frac{1}{\psi} + 1\right) \sum_{j=1}^n E \left[ U_i^2 \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \middle| \mathbf{X}_i \right]. \end{aligned}$$

Expansion of the fourth term of  $T_A$ 

For the fourth term of  $T_A$  we have

$$\begin{aligned} & \sum_{j=1}^n E \left( \{b(\beta_0 + \beta_1 X_{ij} + U_i) + 1 - \mathbf{q}(\psi)\}^2 \middle| \mathbf{X}_i \right) \\ &= \sum_{j=1}^n E \left( \{b(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \middle| \mathbf{X}_i \right) + 2\{1 - \mathbf{q}(\psi)\} \sum_{j=1}^n E \left( b(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right) \\ & \quad + n\{1 - \mathbf{q}(\psi)\}^2 \\ &= \sum_{j=1}^n E \left( \{b(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \middle| \mathbf{X}_i \right) + 2 \sum_{j=1}^n E \left( b(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right) \\ & \quad - 2\mathbf{q}(\psi) \sum_{j=1}^n E \left\{ b(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right\} + n\{\mathbf{q}(\psi)\}^2 - 2n\mathbf{q}(\psi) + n. \end{aligned}$$

Expansion of the fifth term of  $T_A$ 

For the fifth term of  $T_A$  we have

$$\begin{aligned} & -2 \sum_{j=1}^n E \left[ \left\{ \left(\frac{1}{\psi} + 1\right) + b(\beta_0 + \beta_1 X_{ij} + U_i) \right\} \{1 + U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\} \middle| \mathbf{X}_i \right] \\ &= -2n \left(\frac{1}{\psi} + 1\right) - 2 \left(\frac{1}{\psi} + 1\right) \sum_{j=1}^n E \left\{ U_i b'(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right\} \\ & \quad - 2 \sum_{j=1}^n E \left\{ b(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right\} - 2 \sum_{j=1}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right\}. \end{aligned}$$

Expansion of the sixth term of  $T_A$ 

For the sixth term of  $T_A$  we have

$$\begin{aligned}
& 2q(\psi) \sum_{j=1}^n E \left[ \{b(\beta_0 + \beta_1 X_{ij} + U_i) + 1 - q(\psi)\} \middle| \mathbf{X}_i \right] \\
&= 2q(\psi) \sum_{j=1}^n E [b(\beta_0 + \beta_1 X_{ij} + U_i) | \mathbf{X}_i] + 2nq(\psi)\{1 - q(\psi)\} \\
&= 2q(\psi) \sum_{j=1}^n E [b(\beta_0 + \beta_1 X_{ij} + U_i) | \mathbf{X}_i] + 2nq(\psi) - 2n\{q(\psi)\}^2.
\end{aligned}$$

The fully expanded version of  $T_A$  is as follows:

$$\begin{aligned}
T_A &= n \left\{ \text{trigamma}(\psi) - \frac{1}{\psi} \right\} \\
&+ \sum_{j \neq j'}^n E \left\{ b(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right\} \\
&- \sum_{j \neq j'}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right\} \\
&- \sum_{j \neq j'}^n E \left\{ U_i b'(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right\} \\
&+ \sum_{j \neq j'}^n E \left\{ U_i^2 b'(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right\} \\
&+ \left( \frac{1}{\psi} + 1 \right) \sum_{j=1}^n E \left[ U_i^2 \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \middle| \mathbf{X}_i \right] \\
&+ \sum_{j=1}^n E \left( \{b(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \middle| \mathbf{X}_i \right) \\
&- 2 \sum_{j=1}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right\}.
\end{aligned}$$

Treatment of  $T_B$ 

We have

$$E \left( U_i^2 Y_{i\bullet}^2 \middle| \mathbf{X}_i \right) = E \left\{ E \left( U_i^2 Y_{i\bullet}^2 \middle| \mathbf{X}_i, U_i \right) \middle| \mathbf{X}_i \right\} = E \left\{ U_i^2 E \left( Y_{i\bullet}^2 \middle| \mathbf{X}_i, U_i \right) \middle| \mathbf{X}_i \right\}.$$

Recall that

$$Y_{i\bullet} = \sum_{j=1}^n Y_{ij} \quad \text{implying that} \quad Y_{i\bullet}^2 = \sum_{j \neq j'}^n Y_{ij} Y_{ij'} + \sum_{j=1}^n Y_{ij}^2.$$

Therefore

$$\begin{aligned} & E(Y_{i\bullet}^2 | \mathbf{X}_i, U_i) \\ &= \sum_{j \neq j'}^n E(Y_{ij} Y_{ij'} | \mathbf{X}_i, U_i) + \sum_{j=1}^n E(Y_{ij}^2 | \mathbf{X}_i, U_i) \\ &= \sum_{j \neq j'}^n b'(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) + \left(\frac{1}{\psi} + 1\right) \sum_{j=1}^n \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2. \end{aligned}$$

This implies that

$$\begin{aligned} T_B &= \sum_{j \neq j'}^n E\{U_i^2 b'(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i\} \\ &\quad + \left(\frac{1}{\psi} + 1\right) \sum_{j=1}^n E[U_i^2 \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 | \mathbf{X}_i]. \end{aligned}$$

#### Treatment of $T_C$

We have

$$T_C = E(\mathcal{A}_i^2 | \mathbf{X}_i)$$

which implies that

$$T_C = \sum_{j=1}^n \sum_{j'=1}^n E\{b(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i\}.$$

Breaking this up into “diagonal” and “off-diagonal” components, we get

$$\begin{aligned} T_C &= \sum_{j \neq j'}^n E\{b(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i\} \\ &\quad + \sum_{j=1}^n E[\{b(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 | \mathbf{X}_i]. \end{aligned}$$

#### Treatment of $T_D$

First recall that

$$2E(\mathcal{B}_i U_i Y_{i\bullet} | \mathbf{X}_i)$$

where

$$\mathcal{B}_i \equiv \sum_{j=1}^n \left[ \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij})\} + \mathbf{q}(\psi) \right].$$

Therefore

$$\begin{aligned} 2\mathcal{B}_i U_i Y_{i\bullet} &= 2U_i \sum_{j=1}^n \sum_{j'=1}^n \left[ \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij})\} + \mathbf{q}(\psi) \right] Y_{ij'} \\ &= 2U_i \sum_{j \neq j'}^n \left[ \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij})\} + \mathbf{q}(\psi) \right] Y_{ij'} \\ &\quad + 2U_i \sum_{j=1}^n \left[ \{Y_{ij}^2(\beta_0 + \beta_1 X_{ij}) + Y_{ij} \log(Y_{ij})\} + Y_{ij} \mathbf{q}(\psi) \right]. \end{aligned}$$

It follows that

$$2E(\mathcal{B}_i U_i Y_{i\bullet} | \mathbf{X}_i) = \mathbf{r}_3(\mathbf{X}_i) + \mathbf{r}_4(\mathbf{X}_i)$$

where

$$\mathbf{r}_3(\mathbf{X}_i) \equiv 2 \sum_{j \neq j'}^n E \left( U_i \left[ \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij})\} + \mathbf{q}(\psi) \right] Y_{ij'} \middle| \mathbf{X}_i \right)$$

and

$$\mathbf{r}_4(\mathbf{X}_i) \equiv 2 \sum_{j=1}^n E \left( U_i \left[ \{Y_{ij}^2(\beta_0 + \beta_1 X_{ij}) + Y_{ij} \log(Y_{ij})\} + Y_{ij} \mathbf{q}(\psi) \right] \middle| \mathbf{X}_i \right).$$

#### Treatment of $\mathbf{r}_3(\mathbf{X}_i)$

The  $(j, j')$ th term in the  $\mathbf{r}_3(\mathbf{X}_i)$  summation is

$$\begin{aligned} &2E \left\{ E \left( U_i \left[ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right] Y_{ij'} \middle| \mathbf{X}_i, U_i \right) \middle| \mathbf{X}_i \right\} \\ &= 2E \left\{ U_i E \left( \left[ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right] Y_{ij'} \middle| \mathbf{X}_i, U_i \right) \middle| \mathbf{X}_i \right\} \\ &= 2E \left\{ U_i E \left( \left[ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right] \middle| \mathbf{X}_i, U_i \right) E \left( Y_{ij'} \middle| \mathbf{X}_i, U_i \right) \middle| \mathbf{X}_i \right\}. \end{aligned}$$

As stated earlier in this document,

$$E \left( \left[ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right] \middle| \mathbf{X}_i, U_i \right) = b(\beta_0 + \beta_1 X_{ij} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij} + U_i).$$

Also,

$$E\left(Y_{ij'} \mid \mathbf{X}_i, U_i\right) = b'(\beta_0 + \beta_1 X_{ij'} + U_i)$$

from which it follows that the  $(j, j')$ th term in the  $\mathbf{r}_3(\mathbf{X}_i)$  summation is

$$2E\left[U_i \{b(\beta_0 + \beta_1 X_{ij} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\} b'(\beta_0 + \beta_1 X_{ij'} + U_i) \mid \mathbf{X}_i\right].$$

Hence,

$$\mathbf{r}_3(\mathbf{X}_i) = 2 \sum_{j \neq j'}^n E\left[U_i \{b(\beta_0 + \beta_1 X_{ij} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\} b'(\beta_0 + \beta_1 X_{ij'} + U_i) \mid \mathbf{X}_i\right]$$

#### Treatment of $\mathbf{r}_4(\mathbf{X}_i)$

The  $j$ th term in the  $\mathbf{r}_4(\mathbf{X}_i)$  summation is

$$\begin{aligned} & 2E\left(E\left[U_i \left\{Y_{ij}^2(\beta_0 + \beta_1 X_{ij}) + Y_{ij} \log(Y_{ij}) + Y_{ij} \mathbf{q}(\psi)\right\} \mid \mathbf{X}_i, U_i\right] \mid \mathbf{X}_i\right) \\ &= 2E\left(U_i E\left[\left\{Y_{ij}^2(\beta_0 + \beta_1 X_{ij}) + Y_{ij} \log(Y_{ij}) + Y_{ij} \mathbf{q}(\psi)\right\} \mid \mathbf{X}_i, U_i\right] \mid \mathbf{X}_i\right). \end{aligned}$$

Next note that

$$\begin{aligned} E\left[Y_{ij}^2(\beta_0 + \beta_1 X_{ij}) \mid \mathbf{X}_i, U_i\right] &= \left(\frac{1}{\psi} + 1\right) (\beta_0 + \beta_1 X_{ij}) \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \\ &= \left(\frac{1}{\psi} + 1\right) (\beta_0 + \beta_1 X_{ij} + U_i) \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \\ &\quad - \left(\frac{1}{\psi} + 1\right) U_i \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \\ &= -\left(\frac{1}{\psi} + 1\right) [b'(\beta_0 + \beta_1 X_{ij} + U_i) + U_i \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2]. \end{aligned}$$

Also, remember that

$$E\{Y_{ij} \log(Y_{ij}) \mid \mathbf{X}_i, U_i\} = b'(\beta_0 + \beta_1 X_{ij} + U_i) \left\{ \frac{1}{\psi} + 1 - \mathbf{q}(\psi) + b(\beta_0 + \beta_1 X_{ij} + U_i) \right\}$$

and

$$\mathbf{q}(\psi) E\{Y_{ij} \mid \mathbf{X}_i, U_i\} = \mathbf{q}(\psi) b'(\beta_0 + \beta_1 X_{ij} + U_i),$$

which means that

$$E\{Y_{ij} \log(Y_{ij}) + \mathbf{q}(\psi) Y_{ij} \mid \mathbf{X}_i, U_i\} = b'(\beta_0 + \beta_1 X_{ij} + U_i) \left\{ \frac{1}{\psi} + 1 + b(\beta_0 + \beta_1 X_{ij} + U_i) \right\}.$$

This means that

$$\begin{aligned}
& E \left[ \left\{ Y_{ij}^2(\beta_0 + \beta_1 X_{ij}) + Y_{ij} \log(Y_{ij}) + Y_{ij} \mathbf{q}(\psi) \right\} \middle| \mathbf{X}_i, U_i \right] \\
&= - \left( \frac{1}{\psi} + 1 \right) b'(\beta_0 + \beta_1 X_{ij} + U_i) - \left( \frac{1}{\psi} + 1 \right) U_i \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \\
&\quad + \left( \frac{1}{\psi} + 1 \right) b'(\beta_0 + \beta_1 X_{ij} + U_i) + b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij} + U_i) \\
&= b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij} + U_i) - \left( \frac{1}{\psi} + 1 \right) U_i \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2.
\end{aligned}$$

Hence, the  $j$ th term in the  $\mathbf{r}_4(\mathbf{X}_i)$  summation is

$$2E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right\} - 2 \left( \frac{1}{\psi} + 1 \right) E \left[ U_i^2 \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \middle| \mathbf{X}_i \right].$$

Putting all of this together, we arrive at

$$\begin{aligned}
T_D &= 2 \sum_{j \neq j'}^n E \left[ U_i \{b(\beta_0 + \beta_1 X_{ij} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\} b'(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right] \\
&\quad + 2 \sum_{j=1}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right\} \\
&\quad - 2 \left( \frac{1}{\psi} + 1 \right) \sum_{j=1}^n E \left[ U_i^2 \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \middle| \mathbf{X}_i \right].
\end{aligned}$$

However, note that the first term is

$$\begin{aligned}
& 2 \sum_{j \neq j'}^n E \left[ U_i \{b(\beta_0 + \beta_1 X_{ij} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij} + U_i)\} b'(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right] \\
&= 2 \sum_{j \neq j'}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right\} \\
&\quad - 2 \sum_{j \neq j'}^n E \left\{ U_i^2 b'(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right\}.
\end{aligned}$$



Therefore

$$\begin{aligned}
T_D &= 2 \sum_{j \neq j'}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right\} \\
&\quad - 2 \sum_{j \neq j'}^n E \left\{ U_i^2 b'(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right\} \\
&\quad + 2 \sum_{j=1}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij} + U_i) \middle| \mathbf{X}_i \right\} \\
&\quad - 2 \left( \frac{1}{\psi} + 1 \right) \sum_{j=1}^n E \left\{ U_i^2 \{b'(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 \middle| \mathbf{X}_i \right\}.
\end{aligned}$$

Treatment of  $T_E$

Recall that

$$T_E \equiv -2E(\mathcal{A}_i \mathcal{B}_i | \mathbf{X}_i)$$

where

$$\mathcal{A}_i \equiv \sum_{j=1}^n b(\beta_0 + \beta_1 X_{ij} + U_i) \quad \text{and} \quad \mathcal{B}_i \equiv \sum_{j=1}^n \left\{ Y_{ij}(\beta_0 + \beta_1 X_{ij}) + \log(Y_{ij}) + \mathbf{q}(\psi) \right\}.$$

Therefore

$$\begin{aligned}
T_E &= -2 \sum_{j=1}^n \sum_{j'=1}^n E \left[ b(\beta_0 + \beta_1 X_{ij} + U_i) \left\{ Y_{ij'}(\beta_0 + \beta_1 X_{ij'}) + \log(Y_{ij'}) + \mathbf{q}(\psi) \right\} \middle| \mathbf{X}_i \right] \\
&= -2 \sum_{j=1}^n \sum_{j'=1}^n E \left( E \left[ b(\beta_0 + \beta_1 X_{ij} + U_i) \left\{ Y_{ij'}(\beta_0 + \beta_1 X_{ij'}) + \log(Y_{ij'}) + \mathbf{q}(\psi) \right\} \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) \\
&= -2 \sum_{j=1}^n \sum_{j'=1}^n E \left( b(\beta_0 + \beta_1 X_{ij} + U_i) E \left[ Y_{ij'}(\beta_0 + \beta_1 X_{ij'}) + \log(Y_{ij'}) + \mathbf{q}(\psi) \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) \\
&= -2 \sum_{j=1}^n \sum_{j'=1}^n E \left( b(\beta_0 + \beta_1 X_{ij} + U_i) \{ b(\beta_0 + \beta_1 X_{ij'} + U_i) - U_i b'(\beta_0 + \beta_1 X_{ij'} + U_i) \} \middle| \mathbf{X}_i \right) \\
&= -2 \sum_{j=1}^n \sum_{j'=1}^n E \left\{ b(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right\} \\
&\quad + 2 \sum_{j=1}^n \sum_{j'=1}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) \middle| \mathbf{X}_i \right\}
\end{aligned}$$

Simplifying the previous expression, we have,

$$\begin{aligned}
& -2 \sum_{j \neq j'}^n E[b(\beta_0 + \beta_1 X_{ij} + U_i)b(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i] \\
& + 2 \sum_{j \neq j'}^n E[U_i b(\beta_0 + \beta_1 X_{ij} + U_i)b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i] \\
& - 2 \sum_{j=1}^n E[\{b(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 | \mathbf{X}_i] \\
& + 2 \sum_{j=1}^n E[U_i b(\beta_0 + \beta_1 X_{ij} + U_i)b'(\beta_0 + \beta_1 X_{ij} + U_i) | \mathbf{X}_i].
\end{aligned}$$

#### Treatment of $T_F$

Using steps similar to those given above for  $T_B$  we have

$$T_F = -2 E(U_i \mathcal{A}_i \mathcal{A}'_i | \mathbf{X}_i).$$

Therefore,

$$T_F = -2 \sum_{j=1}^n \sum_{j'=1}^n E\{U_i b(\beta_0 + \beta_1 X_{ij} + U_i)b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i\}.$$

Breaking this up into “diagonal” and “off-diagonal” components, we get

$$\begin{aligned}
T_F &= -2 \sum_{j \neq j'}^n E\{U_i b(\beta_0 + \beta_1 X_{ij} + U_i)b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i\} \\
& - 2 \sum_{j=1}^n E\{U_i b(\beta_0 + \beta_1 X_{ij} + U_i)b'(\beta_0 + \beta_1 X_{ij} + U_i) | \mathbf{X}_i\}.
\end{aligned}$$

#### The Sum of $T_E$ and $T_F$

Inspection of the fully expanded versions of  $T_E$  and  $T_F$  reveals that

$$\begin{aligned}
T_E + T_F &= -2 \sum_{j \neq j'}^n E[b(\beta_0 + \beta_1 X_{ij} + U_i)b(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i] \\
& - 2 \sum_{j=1}^n E[\{b(\beta_0 + \beta_1 X_{ij} + U_i)\}^2 | \mathbf{X}_i].
\end{aligned}$$

Combining  $T_A, T_B, T_C, T_D, T_E$  and  $T_F$ 

We now combine the  $T_A, T_B, T_C, T_D, T_E$  and  $T_F$  terms to get the full approximation of  $E(S_{3i}^2 | \mathbf{X}_i)$ .

$$\begin{aligned}
& E(S_{3i}^2 | \mathbf{X}_i) \\
& \approx T_A + T_B + T_C + T_D + T_E + T_F \\
& = n \left\{ \text{trigamma}(\psi) - \frac{1}{\psi} \right\} + \sum_{j \neq j'}^n E \left\{ b(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} \\
& \quad - \sum_{j \neq j'}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} \\
& \quad - \sum_{j \neq j'}^n E \left\{ U_i b'(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} \\
& \quad + \sum_{j \neq j'}^n E \left\{ U_i^2 b'(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} \\
& \quad + \left( \frac{1}{\psi} + 1 \right) \sum_{j=1}^n E \left[ U_i^2 \{ b'(\beta_0 + \beta_1 X_{ij} + U_i) \}^2 | \mathbf{X}_i \right] + \sum_{j=1}^n E \left[ \{ b(\beta_0 + \beta_1 X_{ij} + U_i) \}^2 | \mathbf{X}_i \right] \\
& \quad - 2 \sum_{j=1}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij} + U_i) | \mathbf{X}_i \right\} \\
& \quad + \sum_{j \neq j'}^n E \left\{ U_i^2 b'(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} \\
& \quad + \left( \frac{1}{\psi} + 1 \right) \sum_{j=1}^n E \left[ U_i^2 \{ b'(\beta_0 + \beta_1 X_{ij} + U_i) \}^2 | \mathbf{X}_i \right] \\
& \quad + \sum_{j \neq j'}^n E \left\{ b(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} + \sum_{j=1}^n E \left[ \{ b(\beta_0 + \beta_1 X_{ij} + U_i) \}^2 | \mathbf{X}_i \right] \\
& \quad + 2 \sum_{j \neq j'}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} \\
& \quad - 2 \sum_{j \neq j'}^n E \left\{ U_i^2 b'(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right\} \\
& \quad + 2 \sum_{j=1}^n E \left\{ U_i b(\beta_0 + \beta_1 X_{ij} + U_i) b'(\beta_0 + \beta_1 X_{ij} + U_i) | \mathbf{X}_i \right\} \\
& \quad - 2 \left( \frac{1}{\psi} + 1 \right) \sum_{j=1}^n E \left\{ U_i^2 \{ b'(\beta_0 + \beta_1 X_{ij} + U_i) \}^2 | \mathbf{X}_i \right\} \\
& \quad - 2 \sum_{j \neq j'}^n E \left[ b(\beta_0 + \beta_1 X_{ij} + U_i) b(\beta_0 + \beta_1 X_{ij'} + U_i) | \mathbf{X}_i \right] - 2 \sum_{j=1}^n E \left[ \{ b(\beta_0 + \beta_1 X_{ij} + U_i) \}^2 | \mathbf{X}_i \right].
\end{aligned}$$

Many of the terms cancel with each other, and we are left with

$$E(S_{3i}^2 | \mathbf{X}_i) \approx n \left\{ \text{trigamma}(\psi) - \frac{1}{\psi} \right\}.$$

#### 4.5.3.4 Asymptotic Normality and Variance Results for the Maximum Likelihood Estimator of the Reciprocal Dispersion Parameter

We can show that (using results given in Wand (2007))  $\psi$  is orthogonal to the other model parameters. Therefore, results in the previous sections of this document lead to

$$\sqrt{mn}(\hat{\psi} - \psi) \xrightarrow{\mathcal{D}} N\left(0, \frac{1}{\text{trigamma}(\psi) - \frac{1}{\psi}}\right)$$

and we have

$$\text{Asy.Var}(\hat{\psi}) = \frac{1}{mn \left\{ \text{trigamma}(\psi) - \frac{1}{\psi} \right\}}.$$

#### 4.5.3.5 Asymptotic Normality and Variance Results for the Maximum Likelihood Estimator of the Dispersion Parameter

If

$$\phi = 1/\psi = g(\psi) \quad \text{and} \quad g(x) = x^{-1}$$

then using the delta method leads to

$$\sqrt{mn}\{g(\hat{\psi}) - g(\psi)\} \xrightarrow{\mathcal{D}} N\left(0, \frac{g'(\psi)^2}{\text{trigamma}(\psi) - \frac{1}{\psi}}\right).$$

Noting that

$$g'(x) = -x^{-2} \quad \text{and} \quad g'(x)^2 = x^{-4}$$

, we have

$$\frac{g'(\psi)^2}{\text{trigamma}(\psi) - \frac{1}{\psi}} = \frac{\psi^{-4}}{\text{trigamma}(\psi) - \frac{1}{\psi}} = \frac{\phi^4}{\text{trigamma}(1/\phi) - \phi}.$$

Hence

$$\sqrt{mn}(\hat{\phi} - \phi) \xrightarrow{\mathcal{D}} N\left(0, \frac{\phi^4}{\text{trigamma}(1/\phi) - \phi}\right).$$

Lastly, we have that,

$$\text{Asy.Var}(\hat{\phi}) = \frac{\phi^4}{mn \left\{ \text{trigamma}(1/\phi) - \phi \right\}}.$$

## Chapter 5

# Consequences and Applications of Asymptotic Normality Results

In this chapter, we discuss the consequences and applications of the novel asymptotic normality results presented under Theorem 12.

Firstly, we present how Theorem 12 can be used to carry out asymptotically valid statistical inference in generalized linear mixed model analysis. This is done using confidence intervals constructed via the studentization process. Following that, to assess the efficacy of the Theorem 12-based confidence intervals, we ran two simulation studies and investigated the performance of our confidence intervals when the samples are finite.

Next, we move onto the implications of Theorem 12 on optimal design theory. Optimal designs contain values of the covariates in the design matrix such that these designs give the smallest standard errors of the estimators of the model parameters. These in turn lead to narrower confidence intervals and result in higher powers for hypothesis tests as compared to non-optimal designs. However, when dealing with generalized linear mixed models, choosing optimal designs can be complicated with most optimality criteria being based on the Fisher information matrix, which is computationally expensive to compute. Therefore, we then demonstrate how the derivations leading to Theorem 12 that involve large sample approximations of the Fisher information can be used in approximate optimal design settings.

This chapter is broken up into two main parts. Section 5.1 details how Theorem 12-based confidence intervals can be constructed using the studentization process. Next, Section 5.2 applies the derivations leading up to Theorem 12 in the design setting and briefly demonstrates how approximate locally D-optimal designs can be constructed

using the D-optimality criterion as presented in Theorem 13.

## 5.1 Asymptotically Valid Inference

In the first part of this section, we present details regarding the construction of asymptotically valid confidence intervals using Theorem 12. It is followed up by two simulation studies that are used to assess the efficacy of these asymptotically valid confidence intervals.

### 5.1.1 Construction of Asymptotically Valid Confidence Intervals

The asymptotic normality results for maximum quasi-likelihood estimators presented in Theorem 12 can be used to construct asymptotically valid  $100(1 - \alpha)\%$  confidence intervals. The confidence intervals for  $\beta_A^0, \beta_B^0$  and  $\Sigma^0$  are developed using the studentization process, which involve replacing the true quantities in the asymptotic variances in Theorem 12 by their consistent estimators.

#### 100(1 - $\alpha$ )% Confidence Interval for the Entries of $\beta_A^0$

For  $\hat{\beta}_A$ , its asymptotic covariance matrix only involves  $\Sigma^0$ . Hence, studentization simply involves replacing  $\Sigma^0$  by  $\hat{\Sigma}$  which leads to the following asymptotic normality result

$$\sqrt{m}\hat{\Sigma}^{-1/2} \left( \hat{\beta}_A - \beta_A^0 \right) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{I}).$$

Let  $\hat{\sigma}_k^2$  denote the  $k$ th diagonal entry of  $\hat{\Sigma}$ ,  $1 \leq k \leq d_A$ . It follows that the asymptotically valid  $100(1 - \alpha)\%$  confidence interval for the  $k$ th entry of  $\beta_A^0$  is

$$(\hat{\beta}_A)_k \pm \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right) \sqrt{\frac{\hat{\sigma}_k^2}{m}}, \quad (5.1)$$

where  $\Phi$  represents the standard normal cumulative distribution function.

#### 100(1 - $\alpha$ )% Confidence Intervals for the Entries of $\beta_B^0$

Constructing asymptotically valid confidence intervals for the entries of  $\beta_B^0$  is less straightforward compared to constructing confidence intervals for  $\beta_A^0$  or  $\Sigma^0$ . Studentiza-

tion involves replacing  $\mathbf{\Lambda}_{\beta_B}$  by  $\widehat{\mathbf{\Lambda}}_{\beta_B}$  which leads to the following asymptotic normality result

$$\sqrt{mn} \left( \phi \widehat{\mathbf{\Lambda}}_{\beta_B} \right)^{-1/2} \left( \widehat{\beta}_B - \beta_B^0 \right) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{I}),$$

where  $\widehat{\mathbf{\Lambda}}_{\beta_B}$  is defined as follows

$$\widehat{\mathbf{\Lambda}}_{\beta_B} \equiv \left[ |2\pi \widehat{\Sigma}|^{-1/2} \int_{\mathbb{R}^{d_A}} \left\{ \text{lower right } d_B \times d_B \text{ block of } \widehat{\Omega}_{\beta_B}(\mathbf{u})^{-1} \right\}^{-1} \exp \left( -\frac{1}{2} \mathbf{u}^T \widehat{\Sigma}^{-1} \mathbf{u} \right) d\mathbf{u} \right]^{-1}$$

and

$$\widehat{\Omega}_{\beta_B}(\mathbf{u}) \equiv \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \left\{ b'' \left( (\widehat{\beta}_A + \mathbf{u})^T \mathbf{X}_{Aij} + (\widehat{\beta}_B)^T \mathbf{X}_{Bij} \right) \begin{bmatrix} \mathbf{X}_{Aij} \mathbf{X}_{Aij}^T & \mathbf{X}_{Aij} \mathbf{X}_{Bij}^T \\ \mathbf{X}_{Bij} \mathbf{X}_{Aij}^T & \mathbf{X}_{Bij} \mathbf{X}_{Bij}^T \end{bmatrix} \right\}.$$

Then, the asymptotically valid  $100(1 - \alpha)\%$  confidence interval for the  $k$ th entry of  $\beta_B^0$  is

$$\left( \widehat{\beta}_B \right)_k \pm \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right) \sqrt{\frac{\left( \phi \widehat{\mathbf{\Lambda}}_{\beta_B} \right)_{kk}}{m}}. \quad (5.2)$$

#### 100(1 - $\alpha$ )% Confidence Intervals for the Entries of $\Sigma^0$

For  $\widehat{\Sigma}$ , in a similar manner to the case involving  $\widehat{\beta}_A$ , its asymptotic covariance matrix only consists of  $\Sigma^0$ . We apply the studentization process and replace  $\Sigma^0$  by  $\widehat{\Sigma}$  which leads to the following asymptotic normality result

$$\sqrt{m} \left\{ 2\mathbf{D}_{d_A}^+ \left( \widehat{\Sigma} \otimes \widehat{\Sigma} \right) \mathbf{D}_{d_A}^{+T} \right\}^{-1/2} \text{vech} \left( \widehat{\Sigma} - \Sigma^0 \right) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{I}).$$

It then follows that the asymptotically valid  $100(1 - \alpha)\%$  confidence interval for the  $(k, k)$ th entry of  $\Sigma^0$  is

$$\widehat{\sigma}_k^2 \pm \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right) \sqrt{\frac{2 \left( \widehat{\sigma}_k^2 \right)^2}{m}}. \quad (5.3)$$

### 5.1.2 Simulation Study

Two simulation studies were run to assess the efficacy of the confidence intervals constructed using the asymptotic normality results presented in Theorem 12. In this study, confidence intervals for the following  $d_A = d_B = 1$  Poisson mixed model were

generated

$$Y_{ij}|X_{ij}, U_i \stackrel{\text{ind.}}{\sim} \text{Poisson}(\exp(\beta_0^0 + \beta_B^0 X_{ij} + U_i)),$$

$$U_i \stackrel{\text{ind.}}{\sim} N(0, (\sigma^2)^0), \quad 1 \leq i \leq m, \quad 1 \leq j \leq n,$$

with  $\phi = 1$ . To simplify the notation involved,  $\beta_A, \beta_B$  and  $\Sigma$  have been replaced by the scalar parameter symbols  $\beta_0, \beta_B$  and  $\sigma^2$ . The values for the true parameter vector  $(\beta_0^0, \beta_B^0, (\sigma^2)^0)$  were chosen from the following possible set of pre-determined values

$$\{(-0.3, 0.2, 0.5), (2.2, -0.1, 0.16), (1.2, 0.4, 0.1), (0.02, 1.3, 1), (-0.3, 0.2, 0.1)\},$$

such that the data was well-behaved and led to fewer singularity issues. The distribution for  $X_{ij}$  was also chosen to be either  $N(0, 1)$  or  $\text{Uniform}(-1, 1)$ . The number of groups in the simulated data,  $m$ , varied over the set  $\{100, 200, \dots, 1000\}$  and the number of observations present within each group,  $n$ , was fixed at  $m/10$ . A total of 1000 replications were simulated for every possible combination of the true parameter vector, chosen  $X_{ij}$  distribution and value of  $(m, n)$  pair.

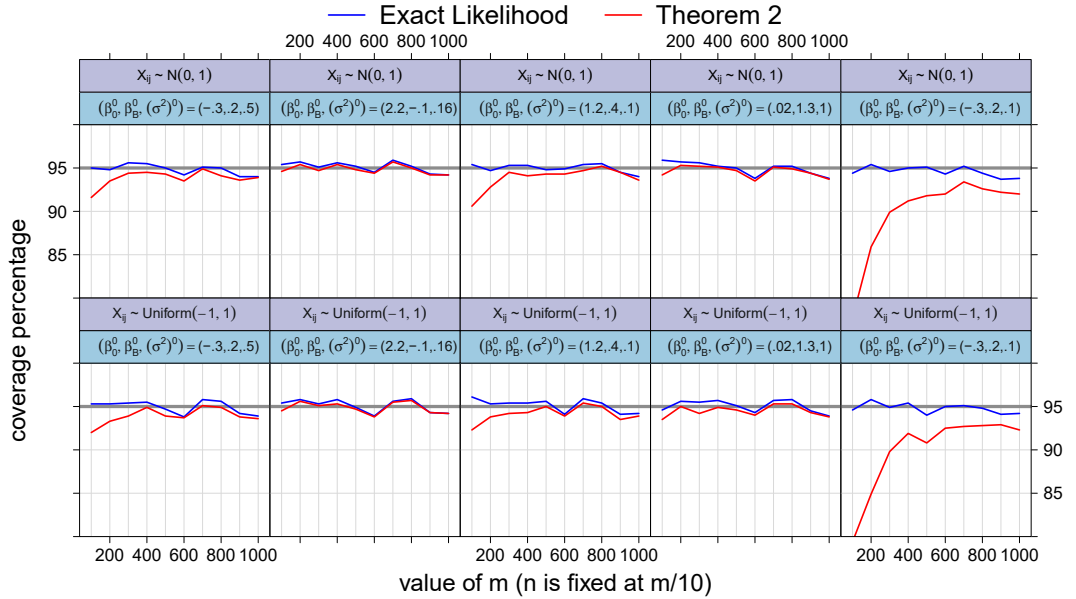


Figure 5.1: Actual coverage percentage of nominally 95% confidence intervals for  $\beta_0^0$  in a  $d_A = d_B = 1$  Poisson mixed model. The confidence intervals are obtained using the exact observed Fisher information computations provided by the function `glmer()` in the R package `lme4` (blue lines) and Theorem 12 with studentization according to (5.1) (red lines). The nominal percentage is shown as a thick grey horizontal line. The percentages are based on 1000 replications. The values of  $m$  are 100, 200,  $\dots$ , 1000. The value of  $n$  is fixed at  $m/10$ .



Using the `glmer()` function in the R package `lme4` (Bates et al., 2015), maximum likelihood estimates of  $\beta_0^0, \beta_B^0$  and  $(\sigma^2)^0$  were obtained for every sample of simulated data generated. 95% confidence intervals based on (5.1) and (5.2) were computed using the maximum likelihood estimates obtained. 95% confidence intervals based on exact observed Fisher information were also obtained using `glmer()`. The coverage percentages corresponding to the percentage of times the true value landed in the confidence intervals were also calculated for both approaches.

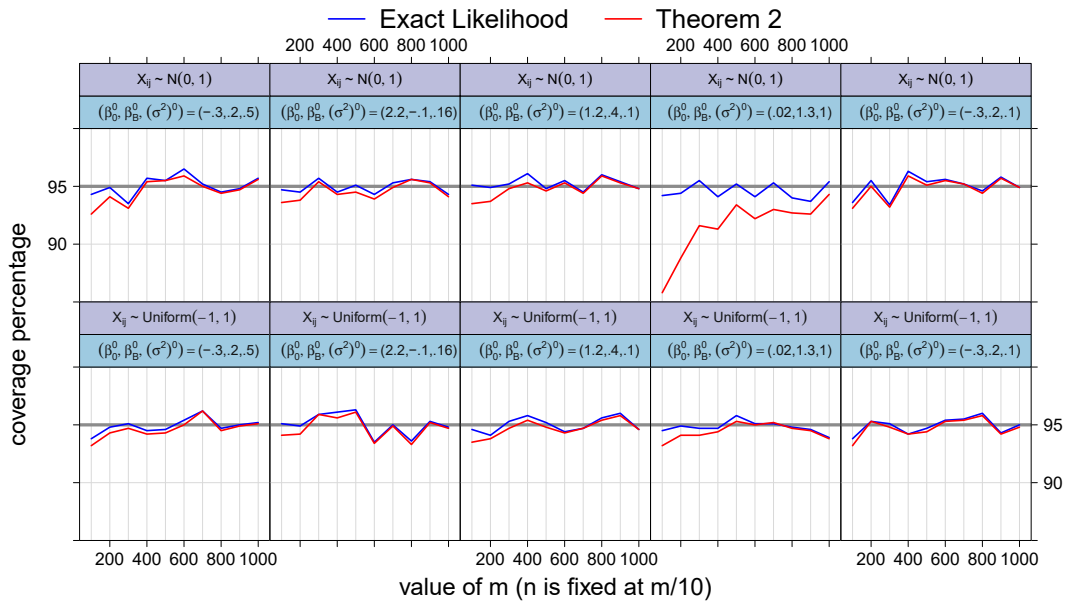


Figure 5.2: Actual coverage percentage of nominally 95% confidence intervals for  $\beta_B^0$  in a  $d_A = d_B = 1$  Poisson mixed model. The confidence intervals are obtained using the exact observed Fisher information computations provided by the function `glmer()` in the R package `lme4` (blue lines) and Theorem 12 with studentization according to (5.2) (red lines). The nominal percentage is shown as a thick grey horizontal line. The percentages are based on 1000 replications. The values of  $m$  are 100, 200,  $\dots$ , 1000. The value of  $n$  is fixed at  $m/10$ .

Figure 5.1 displays the coverage percentages for the 95% confidence intervals computed using the two approaches for various simulation settings. Both the Theorem 12-based approach and the approach based on exact observed Fisher information give almost identical coverage percentages for all response distribution and sample size combinations across all values of  $m$ , for the first four true parameter vector settings. For the last setting, while the Theorem 12-based asymptotic approach does not perform well for smaller values of  $m$ , the asymptotic properties it is based on gives similar coverage

percentages to that of the exact observed Fisher information approach once the value of  $m$  exceeds 500 and continues to get larger. This suggests that the asymptotic variance of  $\hat{\beta}_0$  being  $\sigma^2/m$  is a very good approximation to the variance of  $\hat{\beta}_0$  that arises from using exact observed Fisher information, especially with larger values of  $m$ .

In Figure 5.2, once again, both the Theorem 12-based approach and the approach based on exact observed Fisher information, for most of the true parameter vector settings, give almost identical coverage percentages for all response distribution and sample size combinations across all values of  $m$ . For the fourth setting, when the  $X_{ij}$ s are generated from a standard normal distribution, the Theorem 12-based approach does not perform well for smaller values of  $m$  as compared to the exact observed Fisher information approach. However, both approaches have similar performances for larger values of  $m$ . This suggests that the asymptotic variance of  $\hat{\beta}_B$  being  $\Lambda_{\beta_B}/mn$  is a very good approximation to the variance of  $\hat{\beta}_B$  that arises from using exact observed Fisher information.

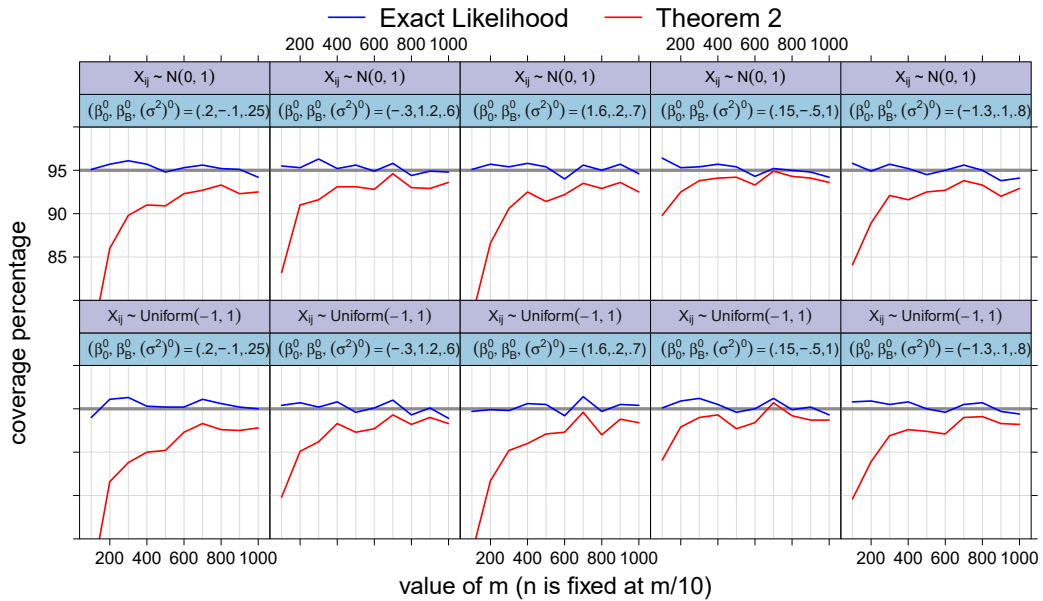


Figure 5.3: Actual coverage percentage of nominally 95% confidence intervals for  $\beta_0^0$  in a  $d_A = d_B = 1$  logistic mixed model. The confidence intervals are obtained using the exact observed Fisher information computations provided by the function `glmer()` in the R package `lme4` (blue lines) and Theorem 12 with studentization according to (5.1) (red lines). The nominal percentage is shown as a thick grey horizontal line. The percentages are based on 1000 replications. The values of  $m$  are 100, 200,  $\dots$ , 1000. The value of  $n$  is fixed at  $m/10$ .

In the next study, confidence intervals for the following  $d_A = d_B = 1$  logistic mixed model were generated

$$Y_{ij}|X_{ij}, U_i \stackrel{\text{ind.}}{\sim} \text{Bernoulli}(\text{expit}(\beta_0^0 + \beta_B^0 X_{ij} + U_i)), \\ U_i \stackrel{\text{ind.}}{\sim} N(0, (\sigma^2)^0), \quad 1 \leq i \leq m, \quad 1 \leq j \leq n,$$

with  $\phi = 1$ .

The values for the true parameter vector  $(\beta_0^0, \beta_B^0, (\sigma^2)^0)$  were chosen from the following possible set of pre-determined values

$$\{(0.2, -0.1, 0.25), (-0.3, 1.2, 0.6), (1.6, 0.2, 0.7), (0.15, -0.5, 1), (-1.3, 0.1, 0.8)\}.$$

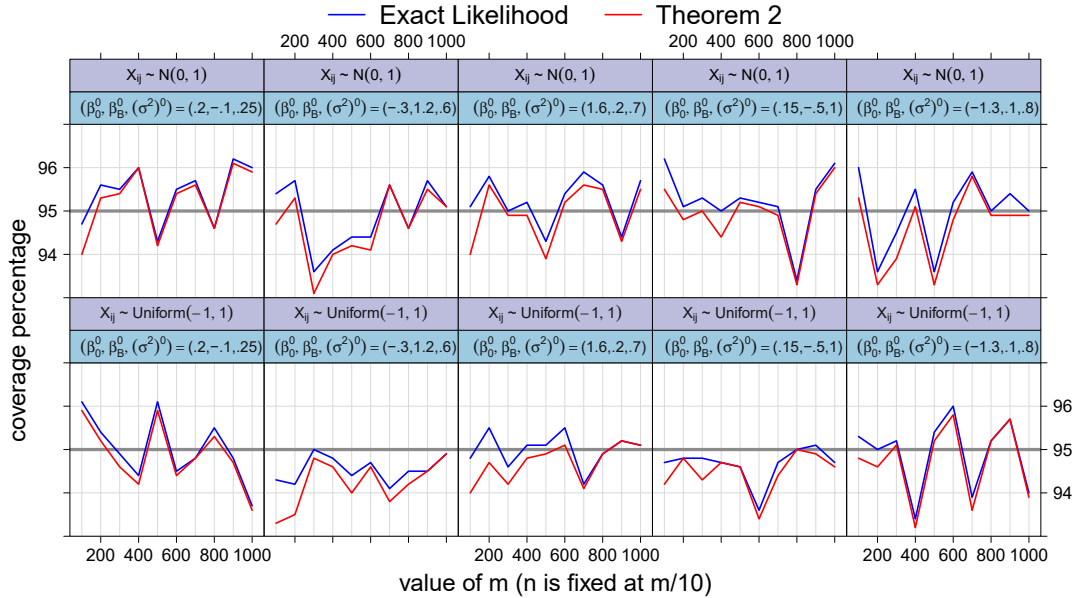


Figure 5.4: Actual coverage percentage of nominally 95% confidence intervals for  $\beta_B^0$  in a  $d_A = d_B = 1$  logistic mixed model. The confidence intervals are obtained using the exact observed Fisher information computations provided by the function `glmer()` in the R package `lme4` (blue lines) and Theorem 12 with studentization according to (5.2) (red lines). The nominal percentage is shown as a thick grey horizontal line. The percentages are based on 1000 replications. The values of  $m$  are 100, 200,  $\dots$ , 1000. The value of  $n$  is fixed at  $m/10$ .

The remaining simulation settings match that of the Poisson simulation study conducted. Once again, the maximum likelihood estimates of  $\beta_0^0, \beta_B^0$  and  $(\sigma^2)^0$  were

obtained and 95% confidence intervals based on both Theorem 12 and exact observed Fisher information were computed.

In Figure 5.3, using exact observed Fisher information leads to better coverage percentages across all values of  $m$ . While the confidence intervals constructed using Theorem 12 do not perform as well for lower values of  $m$ , their performances match those of exact observed Fisher information when the values of  $m$  grows beyond  $m = 500$ .

In Figure 5.4, both approaches give almost identical coverage percentages for all response distribution and sample size combinations across all values of  $m$ , for all five true parameter vector settings. Note that the asymptotic variance of  $\hat{\beta}_B$  has a convergence rate of  $(mn)^{-1}$  while the asymptotic variance of  $\hat{\beta}_0$  has a slower convergence rate of  $m^{-1}$ . This attributes to why the coverage percentages computed from Theorem 12-based confidence intervals achieve percentages closer to 95% for smaller values of  $m$  and  $n$  in Figure 5.4 as compared to the results shown in Figure 5.3 in the simple logistic mixed model simulation study.

Across all four figures, the simulation results indicate that sometimes, certain true values of the model parameters and the chosen distribution of  $\mathbf{X}$  may result in Theorem 12 performing worse in comparison to the other simulation cases presented. This is particularly evident in Figures 5.3 and 5.4. To explain this, we need to look beyond the first-order asymptotic covariances presented in this thesis and analyse the second-order asymptotic covariances instead (Maestrini et al., 2023). Let  $\mathcal{X} \equiv \{\mathbf{X}_{ij} : 1 \leq i \leq m, 1 \leq j \leq n_i\}$ . As an example, for the  $d_A = 1, d_B = 1$  Poisson quasi-likelihood special case of (4.3), a solution to the two-term asymptotic covariance problem can be expressed relatively simply with parameters

$$\beta_A = \beta_0, \quad \beta_B = \beta_1 \quad \text{and} \quad \Sigma = \sigma^2 \quad \text{and predictor variable} \quad \mathbf{X} = \begin{bmatrix} 1 \\ X \end{bmatrix}$$

for a scalar random variable  $X$ . Define

$$a_1(\beta_0, \beta_1, \sigma^2) \equiv e^{\beta_0 + \sigma^2/2} [E(X^2 e^{\beta_1 X}) E(e^{\beta_1 X}) - \{E(X e^{\beta_1 X})\}^2]$$

and

$$a_2(\beta_1, \sigma^2) \equiv \frac{e^{\sigma^2} E(X^2 e^{\beta_1 X}) E(e^{\beta_1 X}) + (1 - e^{\sigma^2}) E\{(X e^{\beta_1 X})\}^2}{E(e^{\beta_1 X})}.$$

Then the two-term asymptotic covariance matrix of  $(\widehat{\beta}_0, \widehat{\beta}_1)$  is

$$\text{Cov} \left( \begin{bmatrix} \widehat{\beta}_0 \\ \widehat{\beta}_1 \end{bmatrix} \middle| \mathcal{X} \right) = \frac{1}{m} \begin{bmatrix} (\sigma^2)^0 & 0 \\ 0 & 0 \end{bmatrix} + \frac{\phi\{1 + o_p(1)\}}{a_1(\beta_0^0, \beta_1^0, (\sigma^2)^0) mn} \begin{bmatrix} a_2(\beta_1^0, (\sigma^2)^0) & -E(Xe^{\beta_1^0 X}) \\ -E(Xe^{\beta_1^0 X}) & E(e^{\beta_1^0 X}) \end{bmatrix}. \quad (5.4)$$

As is apparent from (5.4), the differences between one-term and two-term asymptotic variances depend on  $m$ ,  $n$ ,  $\sigma_0^2$  and particular moments of the  $X$  distribution in a complicated way. Theoretically speaking, it is possible to make these differences arbitrarily small or large by appropriate choices of  $m$ ,  $n$ ,  $\sigma_0^2$  and the distribution of  $X$ . The 4th column of Figure 5.2 contrasts with the other columns in terms of the empirical coverage because the one-term and two-term approximations differ more, with exact likelihood being closer to the two-term approximation.

Therefore, for certain chosen values of the true model parameters and distribution of  $\mathbf{X}$ , the first order asymptotic variance terms are insufficient to obtain good coverage percentages as the second term in the asymptotic variance is also significant, especially for small values of  $m$  and  $n$ . In these cases, the confidence intervals should be constructed using two-term asymptotic variances instead rather than the asymptotic variances presented in Theorem 12, which will lead to better coverage percentages.

With regards to computing these confidence intervals, the Theorem 12 and studentization based approach provides the analyst with a quicker and simpler option, especially for large  $m$ . When computing asymptotically valid confidence intervals using the exact approach, numerical integration is required to compute the ratios of integrals involved when computing the exact observed Fisher information matrix. Note that for  $d_A > 1$ , multivariate numerical integration is needed for the exact approach. On the other hand, for constructing asymptotically valid confidence intervals for  $\beta_0^0$  using Theorem 12, when  $m$  is in the several hundreds or thousands, the closed form confidence interval arising from Theorem 12 and studentization is an attractive alternative to the numerical integration-based exact approach. For constructing asymptotically valid confidence intervals for  $\beta_B^0$  using Theorem 12, the logistic case requires simple numerical integration as compared to the exact approach while the Poisson case does not require any numerical integration.

## 5.2 Approximate Optimal Design

In this section, we demonstrate how the derivations leading to Theorem 12 involving large sample approximations of the Fisher information can be used in approximate optimal design settings.

### 5.2.1 Background and Model Description

In previous sections, we assumed that the data was observed in accordance to the model described in (4.3). Now, let us consider the use of a generalized linear mixed model as in (4.3) with  $d_A = 1$ ,  $\mathbf{X}_A = \mathbf{1}$ ,  $\beta_A = \beta_0$ ,  $\Sigma = \sigma^2$  and with the same number of observations in each group. Also consider the case where the data is yet to be observed. This simplifies to a random intercept generalized linear mixed model as follows

$$\begin{aligned} Y_{ij}|U_i \text{ are independent having quasi-likelihood function (4.2) with} \\ \text{natural parameter } \beta_0^0 + (\beta_B^0)^T \mathbf{x}_{Bij} + U_i \text{ such that the } U_i \text{ are independent} \quad (5.5) \\ N(0, (\sigma^2)^0) \text{ random variables,} \end{aligned}$$

where  $1 \leq i \leq m$  and  $1 \leq j \leq N$ . The unique values of the non-random  $\mathbf{x}_{Bij}$  predictor vectors in the  $i$ th group can be viewed as a finite set of points in  $\mathbb{R}^{d_B}$  and can be denoted as  $\mathbf{x}_1, \dots, \mathbf{x}_s$ . Here,  $s$  denotes the number of unique predictor vectors.

In optimal design, with the help of an optimality criterion, one is firstly required to select the possible values of  $\mathbf{x}$  at which the observations of  $Y_{ij}$  will be made. One also has to determine the fraction of occurrences of independent observations made at each value of  $\mathbf{x}$  (Russell, 2018). Each  $\mathbf{x}$  used in the design is a support point. Hence, there are  $s$  unique support points included in the design.

Now let  $\mathcal{X} \subseteq \mathbb{R}^{d_B}$  denote the set to which the support points are restricted to. For example, if  $d_B = 2$  with the first predictor being binary and the second predictor being a proportion then  $\mathcal{X} = \{(x_1, x_2) : x_1 \in \{0, 1\}, 0 \leq x_2 \leq 1\}$ . Also, denote the number of independent observations made at  $\mathbf{x}_k$ ,  $1 \leq k \leq s$ , as  $n_k$  and define

$$\delta_k \equiv \frac{n_k}{N}, \quad 1 \leq k \leq s,$$

where  $N \equiv n_1 + \dots + n_s$ . Hence, the  $\delta_k$  are known as design weights and represent the fraction of data in the  $i$ th group associated to each support point. Note that  $\delta_1 + \dots + \delta_s = 1$ . Our working assumption throughout the rest of this section is that the

asymptotically valid D-optimal designs are such that each of the  $m$  groups have exactly the same support points and design weights.

In this thesis, we restrict our attention to approximate optimal designs, which are common when the design weights are represented as decimals or fractions if there are recurring decimals. For example, a design could contain 4 support points with the design weights for  $(\delta_1, \delta_2, \delta_3, \delta_4)$  being  $\{0.127, 0.2378, 0.452, 0.1832\}$ . This is an example of an ideal design that is rarely exactly attainable. In such cases, exact designs with all  $n_k \in \mathbb{Z}^+$  are not always possible, especially for small values of  $N$ , and approximate designs are achievable instead.

Out of the optimality-criteria available, we also restrict our attention to D-optimality. This involves choosing the design that maximises the determinant of the Fisher information matrix. The derivations leading to Theorem 12 involve large sample expressions for the Fisher information for the class of generalized linear mixed models. By using the D-optimality criterion, these analogous large sample approximations of the Fisher information allow for approximate locally optimal design determination. Since the approximations of the Fisher information matrix are asymptotic approximations, we only considered designed experiments for which large sample sizes are feasible.

In addition, note that in non-Gaussian generalized linear mixed models, the Fisher information matrix contains entries dependent on the model parameters. Hence, we work with designs that maximise the determinant of the Fisher information matrix with fixed values for the model parameters, known as locally D-optimal designs.

### 5.2.2 Approximate Locally D-Optimal Design Determination

Define

$$n \equiv \frac{1}{s} \sum_{k=1}^s n_k = \text{average of the support point replication sizes within each group.}$$

The theorem relies on the following assumption:

- (A7) The design sample sizes  $n_k$  diverge to  $\infty$  in such a way that  $n_k/(sn) \rightarrow \delta_k$  for constants  $0 \leq \delta_1 \leq 1, 1 \leq k \leq s$ .

**Theorem 13.** *Consider the random intercept generalized linear mixed model described in (5.5) with design weights  $\delta_k$  and corresponding support points  $x_k \in \mathcal{X} \subseteq \mathbb{R}^{d_B}, 1 \leq k \leq s$ . Assume that condition (A7) holds. Then, based on the exact leading term behaviour of*

the determinant of the Fisher information matrix, approximate locally D-optimal designs at the parameter vector  $(\beta_0, \boldsymbol{\beta}_B, \sigma^2)$  are those for which

$$\left| \int_{-\infty}^{\infty} \left\{ \text{lower right } d_B \times d_B \text{ block of } \left( \sum_{k=1}^s \delta_k b''(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u) \begin{bmatrix} 1 & \mathbf{x}_k^T \\ \mathbf{x}_k & \mathbf{x}_k \mathbf{x}_k^T \end{bmatrix} \right)^{-1} \right\}^{-1} \times \exp\{-u^2/(2\sigma^2)\} du \right|$$

is maximal over  $\{\delta_k : \delta_k \geq 0, \sum_{k=1}^s \delta_k = 1, 1 \leq k \leq s\}$  and  $\{\mathbf{x}_k \in \mathcal{X} : 1 \leq k \leq s\}$ .

The proof of Theorem 13 is in the appendix. Some remarks regarding Theorem 13 are as follows:

1. Existing literature that consider the selection of optimal designs using the D-optimality criterion for classes of generalized linear mixed models similar to those considered in this thesis include Waite and Woods (2015) and Zhang et al. (2017). Waite and Woods (2015) used generalized estimating equations to obtain approximations of the mixed model Fisher Information Matrix. On the other hand, Zhang et al. (2017) explored three methods to approximate the Fisher information matrix, namely importance sampling, Laplace approximation and joint sampling. In contrast, Theorem 13 facilitates approximate locally D-optimal design determination in a more direct manner and is based on the precise leading term behaviour of the Fisher information matrix.
2. In this chapter, the theory and discussion has been restricted to D-optimality. Other optimality criteria such as A-optimality, which requires minimization of the trace of the inverse of the information matrix, also benefit from our precise asymptotic approximations of the Fisher information matrix for generalized linear mixed models.
3. Note that in the Gaussian case,  $b''(x) = 1$ , and the determinant in Theorem 13 is proportional to

$$\left| \sum_{k=1}^s \delta_k \begin{bmatrix} 1 & \mathbf{x}_k^T \\ \mathbf{x}_k & \mathbf{x}_k \mathbf{x}_k^T \end{bmatrix} \right|. \quad (5.6)$$

Since the expression above does not depend on any of the model parameters, designs that maximize (5.6) are globally D-optimal.

4. In the Poisson linear mixed model case,  $b''(x) = \exp(x)$  leads to the simplification



of Theorem 13. The D-optimality criterion then reduces to

$$\left| \sum_{k=1}^s \delta_k \exp(\boldsymbol{\beta}_B^T \mathbf{x}_k) \begin{bmatrix} 1 & \mathbf{x}_k^T \\ \mathbf{x}_k & \mathbf{x}_k \mathbf{x}_k^T \end{bmatrix} \right| / \sum_{k=1}^s \delta_k \exp(\boldsymbol{\beta}_B^T \mathbf{x}_k).$$

Note that generating approximate locally D-optimal designs for the Poisson linear mixed model case is only dependent on  $\boldsymbol{\beta}_B$ .

5. When considering logistic mixed models,  $b''(x) = 1/[2\{1 + \cosh(x)\}]$  and Theorem 13 does not simplify further. Hence, in the logistic mixed model case, approximate locally D-optimal designs depend on  $\beta_0, \boldsymbol{\beta}_B$  and  $\sigma^2$ . Although Theorem 13 does not admit an explicit form in this case, each of the entries of the approximate Fisher information matrix can be computed using univariate numerical integration.

### 5.2.3 Illustration of Theorem 13

In this section, we illustrate the use of Theorem 13 in determining the approximate optimal design when  $d_B = 2$  and both predictors are binary, taking on values of either 0 or 1. In this scenario, there are at most  $s = 4$  support points with the only possible support points being  $\mathbf{x}_k \in \{(0, 0), (0, 1), (1, 0), (1, 1)\}$ . Since all the possible support points are known, one would only need to maximise the expression in Theorem 13 over the design weights. Figure 5.5 shows the approximate locally D-optimal designs, for the situation where  $\beta_0 = -0.3$ ,  $\boldsymbol{\beta}_B = (1.7, 2.1)$  and the values of  $\sigma$  take on a value from  $\{0.6, 0.76, 0.97, 1.24, 1.57, 2.00\}$ , by displaying the optimal design weights for the possible corresponding support points. To obtain Figure 5.5, we used code similar to that provided in Section 4.5 of Russell (2018), which is based on the `optim()` function in R (R Core Team (2022)) and Nelder-Mead searches with 100 random initial values.

We noted that the choices for the initial values of the design weights did not impact the results. From Figure 5.5, we see that for the two lowest values of  $\sigma$ , the optimal designs have only three support points, with the point  $(1, 1)$  being excluded from the design. However, as the value of  $\sigma$  increases, the corresponding design weight for the support point  $(1, 1)$  becomes positive and larger.

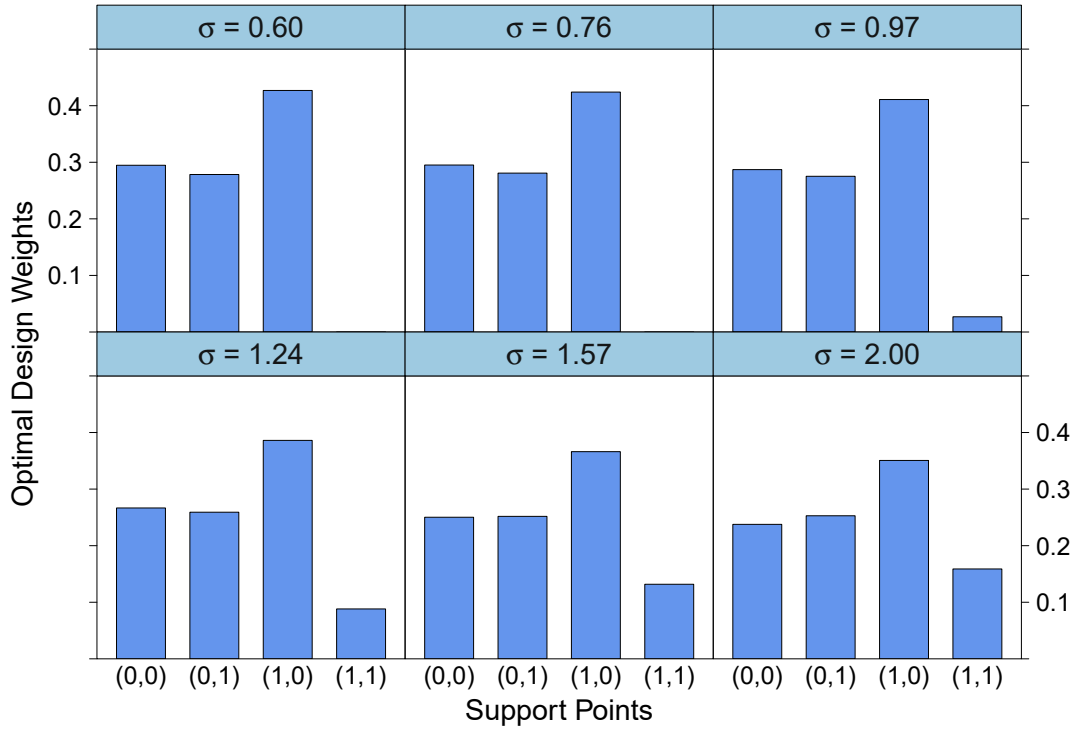


Figure 5.5: Approximate locally D-optimal designs for logistic mixed models with two binary predictors when  $\beta_0 = -0.3$ ,  $\beta_B = (1.7, 2.1)$  and the values of  $\sigma$  take on a value from  $\{0.6, 0.76, 0.97, 1.24, 1.57, 2.00\}$ .

## 5.3 Appendix

This appendix contains the details for the derivations leading up to Theorem 13.

### 5.3.1 Model Description

Our working assumption throughout this section is that the asymptotically D-optimal designs have exactly the same support points and design weights for each of the  $m$  groups. Let  $\mathbf{x}_k$ ,  $1 \leq k \leq s$ , be the support points that are common to each group. For each  $1 \leq i \leq m$ , let

$$n_k = \text{number of } \mathbf{x}_k \text{ values in the design, } 1 \leq k \leq s.$$

Following that, the full data for the  $i$ th group can be expressed as:

$$\begin{aligned} & (\mathbf{x}_1, Y_{i1}^{[1]}), (\mathbf{x}_1, Y_{i2}^{[1]}), \dots, (\mathbf{x}_1, Y_{in_1}^{[1]}), \\ & (\mathbf{x}_2, Y_{i1}^{[2]}), (\mathbf{x}_2, Y_{i2}^{[2]}), \dots, (\mathbf{x}_2, Y_{in_2}^{[2]}), \\ & \quad \vdots \\ & (\mathbf{x}_s, Y_{i1}^{[s]}), (\mathbf{x}_s, Y_{is}^{[s]}), \dots, (\mathbf{x}_s, Y_{in_s}^{[s]}). \end{aligned}$$

Then, the conditional quasi-probability mass function or quasi-density function of  $Y_{ij}^{[k]}$  given  $U_i$  is

$$p_{Y_{ij}^{[k]}|U_i}(y|U_i = u_i) = \exp \left[ \frac{1}{\phi} \{y(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u_i) - b(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u_i)\} + c(y) \right],$$

where  $U_i \sim N(0, \sigma^2)$ ,  $1 \leq i \leq m$ ,  $1 \leq j \leq n_k$  and  $1 \leq k \leq s$ . Also, for each  $1 \leq i \leq m$ , conditional on  $U_i$  the

$$Y_{ij}^{[k]}, \quad 1 \leq j \leq n_k, \quad 1 \leq k \leq s$$

are independent. Therefore, the quasi-likelihood of  $(\beta_0, \boldsymbol{\beta}_B, \sigma^2)$  is

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}_B, \sigma^2) = \prod_{i=1}^m \int_{-\infty}^{\infty} \left\{ \prod_{k=1}^s \prod_{j=1}^{n_k} p_{Y_{ij}^{[k]}|U_i}(Y_{ij}^{[k]}|U_i = u) \right\} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\{-(u^2)/(2\sigma^2)\} du,$$

and the conditional quasi-log-likelihood is

$$\begin{aligned} \ell(\beta_0, \boldsymbol{\beta}_B, \sigma^2) &= \sum_{i=1}^m \log \int_{-\infty}^{\infty} \left\{ \prod_{k=1}^s \prod_{j=1}^{n_k} p_{Y_{ij}^{[k]}|U_i}(Y_{ij}^{[k]}|U_i = u) \right\} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\{-(u^2)/(2\sigma^2)\} du \\ &= \sum_{i=1}^m \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp \left[ \sum_{k=1}^s \sum_{j=1}^{n_k} \left\{ \frac{1}{\phi} Y_{ij}^{[k]} (\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u) \right. \right. \right. \\ &\quad \left. \left. \left. - \frac{1}{\phi} b(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u) + c(Y_{ij}^{[k]}) \right\} - (u^2)/(2\sigma^2) \right] du \right). \end{aligned}$$

Note that,

$$\begin{aligned} & \int_{-\infty}^{\infty} \exp \left[ \sum_{k=1}^s \sum_{j=1}^{n_k} \frac{1}{\phi} \left\{ Y_{ij}^{[k]} (\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u) - b(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u) \right\} - \frac{u^2}{2\sigma^2} \right] du \\ &= \exp \left[ \sum_{k=1}^s \sum_{j=1}^{n_k} \frac{1}{\phi} \left\{ Y_{ij}^{[k]} (\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k) \right\} \right] \\ &\quad \times \int_{-\infty}^{\infty} \exp \left[ \sum_{k=1}^s \sum_{j=1}^{n_k} \frac{1}{\phi} \left\{ Y_{ij}^{[k]} u - b(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u) \right\} - \frac{u^2}{2\sigma^2} \right] du. \end{aligned}$$

Hence,

$$\begin{aligned} & \log \int_{-\infty}^{\infty} \exp \left[ \sum_{k=1}^s \sum_{j=1}^{n_k} \frac{1}{\phi} \left\{ Y_{ij}^{[k]} (\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u) - b(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u) \right\} - \frac{u^2}{2\sigma^2} \right] du \\ &= \sum_{k=1}^s \sum_{j=1}^{n_k} \frac{1}{\phi} Y_{ij}^{[k]} (\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k) \\ & \quad + \log \int_{-\infty}^{\infty} \exp \left[ \sum_{k=1}^s \sum_{j=1}^{n_k} \frac{1}{\phi} \left\{ Y_{ij}^{[k]} u - b(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u) \right\} - \frac{u^2}{2\sigma^2} \right] du. \end{aligned}$$

Therefore, by rewriting the conditional quasi-log-likelihood, we have,

$$\begin{aligned} \ell(\beta_0, \boldsymbol{\beta}_B, \sigma^2) &= -\frac{m}{2} \log(2\pi\sigma^2) + \sum_{i=1}^m \sum_{k=1}^s \sum_{j=1}^{n_k} \frac{1}{\phi} Y_{ij}^{[k]} (\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k) \\ & \quad + \sum_{i=1}^m \log \int_{-\infty}^{\infty} \exp \left[ \sum_{k=1}^s \sum_{j=1}^{n_k} \frac{1}{\phi} \left\{ Y_{ij}^{[k]} u - b(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + u) \right\} - \frac{u^2}{2\sigma^2} \right] du \\ & \quad + C, \end{aligned}$$

where  $C$  denotes a constant term independent of  $(\beta_0, \boldsymbol{\beta}_B, \sigma^2)$ .

### 5.3.2 Asymptotic Assumption for Support Point Sample Sizes

Define

$$n \equiv \frac{1}{s} \sum_{k=1}^s n_k = \text{average of the support point replication sizes within each group.}$$

For the upcoming asymptotic analysis we assume:

- (A4) The design sample sizes  $n_k$  diverge to  $\infty$  in such a way that  $n_k/(sn) \rightarrow \delta_k$  for constants  $0 < \delta_k < 1$ ,  $1 \leq k \leq s$ .

The  $\delta_k$  correspond to the so-called *design weights*.

### 5.3.3 Useful Notation

Let  $\mathbf{v}$  be a generic  $d \times 1$  vector. Then for  $r = 0, 1, 2$  we define

$$\mathbf{v}^{\otimes r} \equiv \begin{cases} 1 & \text{for } r = 0 \\ \mathbf{v} & \text{for } r = 1 \\ \mathbf{v}\mathbf{v}^T & \text{for } r = 2. \end{cases}$$

Note that, according to this notation, for all  $(r, r') \in \{(0, 0), (1, 0), (1, 1)\}$ ,

$$\mathbf{v}^{\otimes r}(\mathbf{v}^{\otimes r'})^T = \mathbf{v}^{\otimes(r+r')}.$$

For  $r = 0, 1$  and  $1 \leq i \leq m$ , also define

$$\tilde{\mathcal{G}}_{ri} \equiv \sum_{k=1}^s \sum_{j=1}^{n_k} \mathbf{x}_k^{\otimes r} \{Y_{ij}^{[k]} - b'(\beta_0 + \beta_B^T \mathbf{x}_k + U_i)\} \quad \text{and} \quad \tilde{\mathcal{H}}_{ri} \equiv \sum_{k=1}^s n_k \mathbf{x}_k^{\otimes r} b''(\beta_0 + \beta_B^T \mathbf{x}_k + U_i).$$

Also, let  $\tilde{\mathcal{H}}'_{ir}$  be the same as  $\tilde{\mathcal{H}}_{ir}$  but with  $b''$  replaced by  $b'''$ .

### 5.3.4 Key Moment Results

In this subsection, we present some key moment results required for the asymptotic derivations in the subsections to follow.

#### The Expectation of $\tilde{\mathcal{G}}_{ri}$ Given $U_i$

Note that for  $r, r' = 0, 1$ :

$$\begin{aligned} & E[\mathbf{x}_k^{\otimes r}(\mathbf{x}_k^{\otimes r'})^T \{Y_{ij}^{[k]} - b'(\beta_0 + \beta_B^T \mathbf{x}_k + U_i)\} | U_i] \\ &= \mathbf{x}_k^{\otimes r}(\mathbf{x}_k^{\otimes r'})^T E\{Y_{ij}^{[k]} - b'(\beta_0 + \beta_B^T \mathbf{x}_k + U_i) | U_i\} \\ &= \mathbf{x}_k^{\otimes r}(\mathbf{x}_k^{\otimes r'})^T E\{b'(\beta_0 + \beta_B^T \mathbf{x}_k + U_i) - b'(\beta_0 + \beta_B^T \mathbf{x}_k + U_i)\} \\ &= \mathbf{O}. \end{aligned}$$

Hence

$$E[\mathbf{x}_k^{\otimes r}(\mathbf{x}_k^{\otimes r'})^T \{Y_{ij}^{[k]} - b'(\beta_0 + \beta_B^T \mathbf{x}_k + U_i)\} | U_i] = \mathbf{O}. \quad (5.7)$$

It follows immediately from (5.7) that

$$E(\tilde{\mathcal{G}}_{0i} | U_i) = \mathbf{0}, \quad E(\tilde{\mathcal{G}}_{1i} | U_i) = \mathbf{0} \quad \text{and} \quad E(\tilde{\mathcal{G}}_{2i} | U_i) = \mathbf{O}.$$

#### The Expectation of $\tilde{\mathcal{G}}_{ri} \tilde{\mathcal{G}}_{r'i}^T$ Given $U_i$

Note that for all  $(r, r') \in \{(0, 0), (1, 0), (1, 1)\}$ ,

$$\begin{aligned} E(\tilde{\mathcal{G}}_{ri} \tilde{\mathcal{G}}_{r'i}^T | U_i) &= E \left( \left[ \sum_{k=1}^s \sum_{j=1}^{n_k} \mathbf{x}_k^{\otimes r} \{Y_{ij}^{[k]} - b'(\beta_0 + \beta_B^T \mathbf{x}_k + U_i)\} \right] \right. \\ &\quad \times \left. \left[ \sum_{k'=1}^s \sum_{j'=1}^{n_{k'}} \mathbf{x}_{k'}^{\otimes r'} \{Y_{ij'}^{[k']} - b'(\beta_0 + \beta_B^T \mathbf{x}_{k'} + U_i)\} \right]^T \middle| U_i \right) \\ &= \sum_{k=1}^s \sum_{k'=1}^s \sum_{j=1}^{n_k} \sum_{j'=1}^{n_{k'}} \mathbf{x}_k^{\otimes r}(\mathbf{x}_{k'}^{\otimes r'})^T E \left[ \{Y_{ij}^{[k]} - b'(\beta_0 + \beta_B^T \mathbf{x}_k + U_i)\} \right. \\ &\quad \left. \times \{Y_{ij'}^{[k']} - b'(\beta_0 + \beta_B^T \mathbf{x}_{k'} + U_i)\} | U_i \right]. \end{aligned} \quad (5.8)$$

First consider those terms for which

$$k = k' \quad \text{and} \quad j = j'. \quad (5.9)$$

Such terms have the form

$$\mathbf{x}_k^{\otimes r} (\mathbf{x}_k^{\otimes r'})^T \text{Var} \left[ \{Y_{ij}^{[k]} - b'(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + U_i)\} | U_i \right] = \phi b''(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + U_i).$$

Therefore, the contribution to (5.8) from the terms satisfying (5.9) is

$$\phi \sum_{k=1}^s n_k \mathbf{x}_k^{\otimes(r+r')} b''(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + U_i) = \phi \tilde{\mathcal{H}}_{(r+r')i}.$$

Next consider those terms for which

$$k = k' \quad \text{and} \quad j \neq j'. \quad (5.10)$$

Since

$$\begin{aligned} & E[\{Y_{ij}^{[k]} - b'(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + U_i)\} \{Y_{ij'}^{[k]} - b'(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + U_i)\} | U_i] \\ &= E[\{Y_{ij}^{[k]} - b'(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + U_i)\} | U_i] E[\{Y_{ij'}^{[k]} - b'(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + U_i)\} | U_i] \\ &= 0, \end{aligned}$$

the contribution to (5.8) from the terms satisfying (5.10) is 0. Next consider those terms for which

$$k \neq k'. \quad (5.11)$$

Then  $Y_{ij}^{[k]}$  and  $Y_{ij}^{[k']}$  must be distinct random variables, which implies that

$$E[\{Y_{ij}^{[k]} - b'(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_k + U_i)\} \{Y_{ij}^{[k']} - b'(\beta_0 + \boldsymbol{\beta}_B^T \mathbf{x}_{k'} + U_i)\} | U_i] = 0 \quad \text{for all } k \neq k'.$$

Therefore, for  $(r, r') \in \{(0, 0), (1, 0), (1, 1)\}$  we have

$$E(\tilde{\mathcal{G}}_{ri} \tilde{\mathcal{G}}_{r'i}^T | U_i) = \phi \tilde{\mathcal{H}}_{(r+r')i}.$$

Lastly, there is the issue when  $(r, r') = (0, 1)$ . Note the special case

$$E(\tilde{\mathcal{G}}_{1i} \tilde{\mathcal{G}}_{0i}^T | U_i) = \phi \tilde{\mathcal{H}}_{1i}.$$

Taking the transpose on each side of this equation we obtain

$$E(\tilde{\mathcal{G}}_{0i} \tilde{\mathcal{G}}_{1i}^T | U_i) = \phi \tilde{\mathcal{H}}_{1i}^T.$$

Then, the full set of results for  $E(\tilde{\mathcal{G}}_{ri} \tilde{\mathcal{G}}_{r'i}^T | U_i)$  is:

$$\begin{aligned} E(\tilde{\mathcal{G}}_{0i} \tilde{\mathcal{G}}_{0i}^T | U_i) &= \phi \tilde{\mathcal{H}}_{0i}, & E(\tilde{\mathcal{G}}_{0i} \tilde{\mathcal{G}}_{1i}^T | U_i) &= \phi \tilde{\mathcal{H}}_{1i}^T, \\ E(\tilde{\mathcal{G}}_{1i} \tilde{\mathcal{G}}_{0i}^T | U_i) &= \phi \tilde{\mathcal{H}}_{1i}, & E(\tilde{\mathcal{G}}_{1i} \tilde{\mathcal{G}}_{1i}^T | U_i) &= \phi \tilde{\mathcal{H}}_{2i}. \end{aligned}$$

But since  $\tilde{\mathcal{G}}_{0i}$  is a scalar, these results can be simplified to

$$E(\tilde{\mathcal{G}}_{0i}^2 | U_i) = \phi \tilde{\mathcal{H}}_{0i}, \quad E(\tilde{\mathcal{G}}_{0i} \tilde{\mathcal{G}}_{1i} | U_i) = \phi \tilde{\mathcal{H}}_{1i} \quad \text{and} \quad E(\tilde{\mathcal{G}}_{1i} \tilde{\mathcal{G}}_{1i}^T | U_i) = \phi \tilde{\mathcal{H}}_{2i}. \quad (5.12)$$

### 5.3.5 The Fisher Information Matrix

To construct an asymptotic approximation of the Fisher information matrix, we require asymptotic expansions of the scores and their quadratic expectations. Hence, we follow steps similar to those detailed under the proof for Theorem 12 for these computations as well. Putting together the resultant expressions from these computations, we have the following expression for the Fisher information

$$I(\beta_0, \boldsymbol{\beta}_B, \sigma^2) = \begin{bmatrix} \frac{m}{\sigma^2} + O(mn^{-1}) & O(m)\mathbf{1}_{d_B}^T & O(mn^{-1}) \\ O(m)\mathbf{1}_{d_B} & \frac{1}{\phi} \sum_{i=1}^m E \left( \tilde{\mathcal{H}}_{2i} - \frac{\tilde{\mathcal{H}}_{1i}^{\otimes 2}}{\tilde{\mathcal{H}}_{0i}} \right) + O(m)\mathbf{1}_{d_B}^{\otimes 2} & O(m)\mathbf{1}_{d_B} \\ O(mn^{-1}) & O(m)\mathbf{1}_{d_B}^T & \frac{m}{2\sigma^4} + O(mn^{-1}) \end{bmatrix}.$$

### 5.3.6 The Asymptotic D-Optimality Criterion

Next, we change the ordering of the parameters from  $(\beta_0, \boldsymbol{\beta}_B, \sigma^2)$  to  $(\beta_0, \sigma^2, \boldsymbol{\beta}_B)$ . Then partition  $I(\beta_0, \sigma^2, \boldsymbol{\beta}_B)$  according to

$$I(\beta_0, \sigma^2, \boldsymbol{\beta}_B) = \begin{bmatrix} \tilde{\mathbf{A}}_{11} & \tilde{\mathbf{A}}_{12}^T \\ \tilde{\mathbf{A}}_{12} & \tilde{\mathbf{A}}_{22} \end{bmatrix}$$

where

$$\tilde{\mathbf{A}}_{11} \equiv \begin{bmatrix} \frac{m}{\sigma^2} + O(mn^{-1}) & O(mn^{-1}) \\ O(mn^{-1}) & \frac{m}{2\sigma^4} + O(mn^{-1}) \end{bmatrix}, \quad \tilde{\mathbf{A}}_{12} \equiv O(m)[\mathbf{1}_{d_B} \quad \mathbf{1}_{d_B}],$$

and

$$\tilde{\mathbf{A}}_{22} \equiv \frac{m}{\phi} E \left( \tilde{\mathcal{H}}_{21} - \frac{\tilde{\mathcal{H}}_{11}^{\otimes 2}}{\tilde{\mathcal{H}}_{01}} \right) + O(m)\mathbf{1}_{d_B}^{\otimes 2}.$$

Now, we apply a standard result concerning the determinant of a  $2 \times 2$  block-partitioned matrix (e.g. Harville, 2008; Theorem 13.3.8) to obtain

$$|I(\beta_0, \boldsymbol{\beta}_B, \sigma^2)| = \left| \tilde{\mathbf{A}}_{11} \right| \left| \tilde{\mathbf{A}}_{22} - \tilde{\mathbf{A}}_{12}^T \tilde{\mathbf{A}}_{11}^{-1} \tilde{\mathbf{A}}_{12} \right|.$$

It is easily verified that  $\left| \tilde{\mathbf{A}}_{11} \right| = m^2/(2\sigma^6) + O(mn^{-1})$  and  $\tilde{\mathbf{A}}_{12}^T \tilde{\mathbf{A}}_{11}^{-1} \tilde{\mathbf{A}}_{12} = O(m)\mathbf{1}_{d_B}^{\otimes 2}$ . It follows that

$$\frac{2\phi^{d_B} \sigma^6}{m^{d_B+2}} |I(\beta_0, \boldsymbol{\beta}_B, \sigma^2)| = \left| E \left( \tilde{\mathcal{H}}_{21} - \frac{\tilde{\mathcal{H}}_{11}^{\otimes 2}}{\tilde{\mathcal{H}}_{01}} \right) + O(1)\mathbf{1}_{d_B}^{\otimes 2} \right| = \left| \boldsymbol{\Psi}_n + O(1)\mathbf{1}_{d_B}^{\otimes 2} \right| \quad (5.13)$$

where  $\Psi_n \equiv E \left( \tilde{\mathcal{H}}_{21} - \tilde{\mathcal{H}}_{11}^{\otimes 2} / \tilde{\mathcal{H}}_{01} \right)$ . Since  $\tilde{\mathcal{H}}_{01} = O_P(n)$ ,  $\tilde{\mathcal{H}}_{11} = O_P(n) \mathbf{1}_{d_B}$  and  $\tilde{\mathcal{H}}_{21} = O_P(n) \mathbf{1}_{d_B}^{\otimes 2}$ , we have  $\Psi_n = O(n) \mathbf{1}_{d_B}^{\otimes 2}$ . Let  $\lambda_1(\mathbf{M}), \dots, \lambda_{d_B}(\mathbf{M})$  denote the eigenvalues of a generic  $d_B \times d_B$  matrix  $\mathbf{M}$ . Then

$$\left| \Psi_n + O(1) \mathbf{1}_{d_B}^{\otimes 2} \right| = \prod_{j=1}^{d_B} \lambda_j \left( \Psi_n + O(1) \mathbf{1}_{d_B}^{\otimes 2} \right).$$

As a consequence of Theorem 8.1.4 (Wielandt-Hoffman) of Golub and Van Loan (2013),

$$\lambda_j \left( \Psi_n + O(1) \mathbf{1}_{d_B}^{\otimes 2} \right) = \lambda_j(\Psi_n) + O(1)$$

for each  $1 \leq j \leq d_B$ . Hence

$$\left| \Psi_n + O(1) \mathbf{1}_{d_B}^{\otimes 2} \right| = |\Psi_n| + O(1) \sum_{j=1}^{d_B} |\Psi_n| / \lambda_j(\Psi_n). \quad (5.14)$$

To obtain the order of magnitude of the  $\lambda_j(\Psi_n)$  we appeal to Theorem 8.1.3 (Gershgorin) of Golub and Van Loan (2013). Since all entries of  $\Psi_n$  are  $O(n)$ , the same is true for the lower and upper limits of each of the Gershgorin discs of  $\Psi_n$ . Since each eigenvalue of  $\Psi_n$  is inside at least one Gershgorin disc, we have  $\lambda_j(\Psi_n) = O(n)$ ,  $1 \leq j \leq d_B$ . It follows from this fact and (5.14) that

$$\left| \Psi_n + O(1) \mathbf{1}_{d_B}^{\otimes 2} \right| = |\Psi_n| \{1 + o(1)\}.$$

In view of (5.13), the determinant of  $|I(\beta_0, \beta_B, \sigma^2)|$  is proportional to a quantity with leading term  $|\Psi_n|$  as  $n \rightarrow \infty$ . Recalling that  $n_k = ns\delta_k$  and dividing through by  $ns$  we can assert that approximate locally D-optimal designs, based on the exact leading term behaviour of the determinant of the Fisher information matrix, are those which maximize

$$\left| E \left[ \sum_{k=1}^s \delta_k \mathbf{x}_k^{\otimes r} b''(\beta_0 + \beta_B^T \mathbf{x}_k + U) - \frac{\left( \sum_{k=1}^s \delta_k \mathbf{x}_k b''(\beta_0 + \beta_B^T \mathbf{x}_k + U) \right)^{\otimes 2}}{\sum_{k=1}^s \delta_k b''(\beta_0 + \beta_B^T \mathbf{x}_k + U)} \right] \right|, \quad U \sim N(0, \sigma^2),$$

over the design weights  $\delta_k$  and support points  $\mathbf{x}_k$ ,  $1 \leq k \leq s$ .

The following quantity can then be used to obtain asymptotic D-optimal designs:

$$\left| \sum_{k=1}^s \delta_k E \{ b''(\beta_0 + \beta_B^T \mathbf{x}_k + U) \} \mathbf{x}_k^{\otimes 2} - E \left\{ \frac{\left( \sum_{k=1}^s \delta_k b''(\beta_0 + \beta_B^T \mathbf{x}_k + U) \mathbf{x}_k \right)^{\otimes 2}}{\sum_{k=1}^s \delta_k b''(\beta_0 + \beta_B^T \mathbf{x}_k + U)} \right\} \right| \quad (5.15)$$



where  $U \sim N(0, \sigma^2)$ . An equivalent integral form is:

$$\left| \begin{array}{c} \sum_{k=1}^s \delta_k \mathbf{x}_k^{\otimes 2} \int_{-\infty}^{\infty} b''(\beta_0 + \beta_B^T \mathbf{x}_k + u) \exp\{-u^2/(2\sigma^2)\} du \\ - \int_{-\infty}^{\infty} \frac{\left( \sum_{k=1}^s \delta_k b''(\beta_0 + \beta_B^T \mathbf{x}_k + u) \mathbf{x}_k \right)^{\otimes 2} \exp\{-u^2/(2\sigma^2)\} du}{\sum_{k=1}^s \delta_k b''(\beta_0 + \beta_B^T \mathbf{x}_k + u)} \end{array} \right|. \quad (5.16)$$

### 5.3.7 Alternative Final Asymptotic D-optimality Criterion

Consider the matrix-valued function  $\Omega(u)$  given by

$$\Omega(u) \equiv \sum_{k=1}^s \delta_k b''(\beta_0 + \beta_B^T + u) \begin{bmatrix} 1 & \mathbf{x}_k \\ \mathbf{x}_k^T & \mathbf{x}_k \mathbf{x}_k^T \end{bmatrix}.$$

Then

$$\Omega(u) = \begin{bmatrix} \sum_{k=1}^s \delta_k b''(\beta_0 + \beta_B^T + u) & \sum_{k=1}^s \delta_k \mathbf{x}_k b''(\beta_0 + \beta_B^T + u) \\ \sum_{k=1}^s \delta_k \mathbf{x}_k^T b''(\beta_0 + \beta_B^T + u) & \sum_{k=1}^s \delta_k \mathbf{x}_k \mathbf{x}_k^T b''(\beta_0 + \beta_B^T + u) \end{bmatrix}$$

and the lower right  $d_B \times d_B$  block of  $\Omega(u)^{-1}$  is

$$\left\{ \sum_{k=1}^s \delta_k \mathbf{x}_k \mathbf{x}_k^T b''(\beta_0 + \beta_B^T + u) - \frac{\left( \sum_{k=1}^s \delta_k \mathbf{x}_k b''(\beta_0 + \beta_B^T + u) \right)^{\otimes 2}}{\sum_{k=1}^s \delta_k b''(\beta_0 + \beta_B^T + u)} \right\}^{-1}.$$

The inverse of this function is the function of  $u$  that is multiplied by  $\exp\{-u^2/(2\sigma^2)\}$  in the determinant of (5.16). Therefore, an alternative expression for (5.16) is

$$\left| \int_{-\infty}^{\infty} \left\{ \text{lower right } d \times d \text{ block of } \left( \sum_{k=1}^s \delta_k b''(\beta_0 + \beta_B^T \mathbf{x}_k + u) \begin{bmatrix} 1 & \mathbf{x}_k \\ \mathbf{x}_k^T & \mathbf{x}_k \mathbf{x}_k^T \end{bmatrix} \right)^{-1} \right\}^{-1} \times \exp\{-u^2/(2\sigma^2)\} du \right|$$

and is the form of the asymptotic D-optimality criterion presented in Theorem 13.

### 5.3.8 Special Distribution Cases

#### Poisson Special Case

In the Poisson special case,  $b'' = \exp$  and therefore, we have,

$$b''(\beta_0 + \beta_B^T \mathbf{x}_k + U) = \exp(U + \beta_0) \exp(\beta_B^T \mathbf{x}_k).$$

We can show that  $\exp(U + \beta_0)$  comes out as a multiplicative factor in (5.15). This is the only random factor in the expectation expression. Therefore, an equivalent D-optimality criterion is

$$\left| \frac{\sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \mathbf{x}_k^{\otimes 2} - \frac{\left( \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \mathbf{x}_k \right)^{\otimes 2}}{\sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k)}}{\sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k)} \right|. \quad (5.17)$$

Define

$$\mathbf{\Lambda} \equiv \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \begin{bmatrix} 1 & \mathbf{x}_k^T \\ \mathbf{x}_k & \mathbf{x}_k \mathbf{x}_k^T \end{bmatrix}.$$

Then

$$\mathbf{\Lambda} = \begin{bmatrix} \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) & \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \mathbf{x}_k^T \\ \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \mathbf{x}_k & \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \mathbf{x}_k^{\otimes 2} \end{bmatrix}.$$

From the result concerning the determinant of a  $2 \times 2$  block partitioned matrix, we obtain

$$|\mathbf{\Lambda}| = \left| \begin{bmatrix} \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \end{bmatrix} \right| \left| \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \mathbf{x}_k^{\otimes 2} - \frac{\left( \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \mathbf{x}_k \right)^{\otimes 2}}{\sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k)} \right|$$

It follows that (5.17) is equivalent to

$$|\mathbf{\Lambda}| / \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k).$$

In other words, an equivalent asymptotic D-optimality criterion for Poisson mixed models is

$$\left| \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \begin{bmatrix} 1 & \mathbf{x}_k^T \\ \mathbf{x}_k & \mathbf{x}_k^{\otimes 2} \end{bmatrix} \right| / \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k).$$

Note that, in the Poisson case, this means that the asymptotic locally D-optimal designs do not depend on  $\beta_0$  or  $\sigma^2$ . They only depend on  $\beta_B$ . Also, for Poisson regression models, the D-optimality criterion is

$$\left| \sum_{k=1}^s \delta_k \exp(\beta_B^T \mathbf{x}_k) \begin{bmatrix} 1 & \mathbf{x}_k^T \\ \mathbf{x}_k & \mathbf{x}_k^{\otimes 2} \end{bmatrix} \right|.$$

Therefore, asymptotic D-optimality for Poisson mixed models involves maximisation of a quantity that is similar, but not identical, to that for Poisson regression models.

#### Logistic Special Case

In the logistic mixed model special case, after taking out multiplicative factors, the quantity (5.16) becomes:

$$\left| \frac{\sum_{k=1}^s \delta_k \mathbf{x}_k^{\otimes 2} \int_{-\infty}^{\infty} \frac{\exp\{-u^2/(2\sigma^2)\} du}{1 + \cosh(\beta_0 + \beta_B^T \mathbf{x}_k + u)}}{\int_{-\infty}^{\infty} \frac{\left( \sum_{k=1}^s \delta_k \mathbf{x}_k / \{1 + \cosh(\beta_0 + \beta_B^T \mathbf{x}_k + u)\} \right)^{\otimes 2} \exp\{-u^2/(2\sigma^2)\} du}{\sum_{k=1}^s \delta_k / \{1 + \cosh(\beta_0 + \beta_B^T \mathbf{x}_k + u)\}}} \right|.$$

## Chapter 6

# Thouless-Anderson-Palmer Enhancement of Generalized Linear Mixed Models

Frequentist inference for GLMMs is hindered by integral intractability problems. In machine learning contexts, Thouless-Anderson-Palmer approaches, highlighted in Subsection 1.11.1 in Chapter 1, can not only help overcome issues involving intractable integrals but theoretically also provide better approximations. However, statistical applications such as longitudinal data analysis and multilevel models analysis have not been investigated at all. Therefore, in this chapter, the goal is to apply the Thouless-Anderson-Palmer frequentist variational approach to generalized linear mixed models with canonical links.

Firstly, the TAP enhancement approach is explained and we obtain a result detailing the explicit form of the TAP approximate negative log-likelihood expression for GLMMs with canonical links, which can then be locally minimized to obtain TAP estimates of the true model parameters.

We then carry out simulation studies to investigate the use of the TAP enhancement approach in practical settings and compare it to the popular Gaussian variational approximation approach for simplified Poisson generalized linear mixed model set-ups.

This chapter is broken up into several parts. Section 6.1 details the model set-up used throughout sections 6.2 to 6.4. Section 6.2 then provides an explicit expression for the GVA log-likelihood. Following that, Section 6.3 provide details regarding the TAP enhancement approach which builds on the GVA approach. An explicit result

for the TAP negative approximate log-likelihood for GLMMs with canonical links is then provided in Section 6.4. Lastly, Section 6.5 delves into two simulation studies constructed for Poisson GLMMs and compares the quality of estimates obtained from the TAP approach against the GVA approach.

## 6.1 Model Description

Consider the use of a simple canonical link generalized linear mixed model as follows where

$$\begin{aligned} Y_{ij}|X_{ij}, U_i \text{ are independent having density function (4.1) with} \\ \text{natural parameter } \beta_0^0 + \beta_1^0 X_{ij} + U_i \text{ such that the } U_i \text{ are independent} \\ N\left(0, (\sigma^2)^0\right) \text{ random variables.} \end{aligned} \quad (6.1)$$

Here, the values of  $(X_{ij}, Y_{ij})$  are observed for  $1 \leq i \leq m$  and  $1 \leq j \leq n$ . We have assumed that the  $X_{ij}$  and  $U_i$  are independent random variables. In addition, the  $X_{ij}$  are each assumed as having the same distribution as the random vector  $X$ . The  $U_i$  are the unobserved random effects variables and are assumed to be having the same distribution as the random vector  $U$ .

Let  $\boldsymbol{\beta} \equiv (\beta_0, \beta_1)$  be the vector of fixed parameters. Then the model parameters for this set-up are  $(\boldsymbol{\beta}, \sigma^2)$ .

## 6.2 The Gaussian Variational Approximate Log-Likelihood

Using the model description in (6.1), one can then obtain  $\ell(\boldsymbol{\beta}, \sigma^2)$ , the conditional log-likelihood of  $(\boldsymbol{\beta}, \sigma^2)$ , where

$$\begin{aligned} \ell(\boldsymbol{\beta}, \sigma^2) = \sum_{i=1}^m \sum_{j=1}^n \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + c(Y_{ij})\} - \frac{m}{2} \log(2\pi\sigma^2) \\ + \sum_{i=1}^m \log \int_{-\infty}^{\infty} \exp \left\{ \sum_{j=1}^n (Y_{ij}u - b(\beta_0 + \beta_1 X_{ij} + u)) - \frac{u^2}{2\sigma^2} \right\} du. \end{aligned} \quad (6.2)$$

Maximum likelihood estimation is hindered due to the  $m$  intractable integrals arising in the expression in (6.2). However, each of the  $m$  integrals can be re-written to overcome

this obstacle by re-expressing the  $i$ th integral as follows (Hall et al., 2011)

$$\begin{aligned} & \int_{-\infty}^{\infty} \exp \left\{ \sum_{j=1}^n (Y_{ij}u - b(\beta_0 + \beta_1 X_{ij} + u)) - \frac{u^2}{2\sigma^2} \right\} \frac{e^{-(1/2)(u-\mu_i)^2\lambda_i/\sqrt{2\pi\lambda_i}}}{e^{-(1/2)(u-\mu_i)^2\lambda_i/\sqrt{2\pi\lambda_i}}} du \\ &= \sqrt{2\pi\lambda_i} E_{\tilde{U}_i} \left[ \exp \left\{ \sum_{j=1}^n (Y_{ij}\tilde{U}_i - b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)) - \frac{\tilde{U}_i^2}{2\sigma^2} + \frac{(\tilde{U}_i - \mu_i)^2}{2\lambda_i} \right\} \right] \end{aligned}$$

where  $E_{\tilde{U}_i}$  denotes the expectation with respect to the random variable  $\tilde{U}_i \sim N(\mu_i, \lambda_i)$  with  $\lambda_i > 0$ , for  $1 \leq i \leq m$ . Here,  $(\boldsymbol{\mu}, \boldsymbol{\lambda})$  are known as the variational parameters where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$  and  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_m)$ . Using Jensen's inequality, one can then obtain the following lower bound

$$\begin{aligned} & \log E_{\tilde{U}_i} \left[ \exp \left\{ \sum_{j=1}^n (Y_{ij}\tilde{U}_i - b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)) - \frac{\tilde{U}_i^2}{2\sigma^2} + \frac{(\tilde{U}_i - \mu_i)^2}{2\lambda_i} \right\} \right] \\ & \geq E_{\tilde{U}_i} \left\{ \sum_{j=1}^n (Y_{ij}\tilde{U}_i - b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)) - \frac{\tilde{U}_i^2}{2\sigma^2} + \frac{(\tilde{U}_i - \mu_i)^2}{2\lambda_i} \right\}, \end{aligned}$$

which is now tractable. Then, the Gaussian variational approximation to  $\ell(\boldsymbol{\beta}, \sigma^2)$  is derived as,

$$\begin{aligned} & \underline{\ell}_{\text{GVA}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) \\ &= \sum_{i=1}^m \sum_{j=1}^n \{Y_{ij}(\beta_0 + \beta_1 X_{ij}) + c(Y_{ij})\} - \frac{m}{2} \log(2\pi\sigma^2) + \sum_{i=1}^m \log(\sqrt{2\pi\lambda_i}) \\ & \quad + \sum_{i=1}^m E_{\tilde{U}_i} \left\{ \sum_{j=1}^n (Y_{ij}\tilde{U}_i - b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)) - \frac{\tilde{U}_i^2}{2\sigma^2} + \frac{(\tilde{U}_i - \mu_i)^2}{2\lambda_i} \right\} \quad (6.3) \\ &= \sum_{i=1}^m E_{\tilde{U}_i} \left[ \sum_{j=1}^n \{Y_{ij}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) - b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) + c(Y_{ij})\} \right. \\ & \quad \left. - \frac{\tilde{U}_i^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right] + \frac{1}{2} \sum_{i=1}^m \{1 + \log(2\pi\lambda_i)\}. \end{aligned}$$

Note that

$$\ell(\boldsymbol{\beta}, \sigma^2) \geq \underline{\ell}_{\text{GVA}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})$$

for all vectors  $\boldsymbol{\mu}$  and  $\boldsymbol{\lambda}$ .

### 6.3 Overview of Thouless-Anderson-Palmer Enhancement

In this section, we provide details on how the TAP enhancement approach builds upon the GVA approach. Firstly, for each  $1 \leq i \leq m$ , define the following data vectors:

$$\mathbf{Y}_i \equiv (Y_{i1}, \dots, Y_{in}) \text{ and } \mathbf{X}_i \equiv (X_{i1}, \dots, X_{in}).$$

The GVA negative log-likelihood can then be expressed as follows (Johnstone, 2022)

$$-\underline{\ell}_{\text{GVA}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) = -\frac{1}{2} \sum_{i=1}^m \{1 + \log(2\pi\lambda_i)\} + \sum_{i=1}^m E_{\tilde{U}_i} \left\{ \Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i \right\},$$

where

$$\Psi_i(\tilde{U}_i) = \sum_{j=1}^n \left\{ -Y_{ij}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) + b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) - c(Y_{ij}) \right\} + \frac{\tilde{U}_i^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2).$$

The TAP enhancement approach theoretically obtains better approximations than the GVA approach by enhancing the expression in  $-\underline{\ell}_{\text{GVA}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})$  through the addition of the Onsager's correction term to  $-\underline{\ell}_{\text{GVA}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})$ , which was first introduced in Thouless et al. (1977). The Onsager's correction term is defined as follows

$$-1/2 \sum_{i=1}^m \xi_i(\mu_i, \lambda_i; \mathbf{X}_i, \mathbf{Y}_i, \beta_0, \beta_1, \sigma^2),$$

with

$$\begin{aligned} & \xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2) \\ & \equiv \text{Var}\{\Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} - \lambda_i \left[ E\{\Psi_i'(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} \right]^2 - \frac{\lambda_i^2}{2} \left[ E\{\Psi_i''(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} \right]^2, \end{aligned} \quad (6.4)$$

where (6.4) was obtained based on personal communication with Professor Song Mei (University of California, Berkeley, U.S.A) and Professor Iain Johnstone (Stanford University, U.S.A). They derived a working expression for the main quantity in the Onsager's correction term for density functions from exponential families.

Following that, the TAP approximate negative log-likelihood can be obtained as,

$$\begin{aligned}
 -\ell_{\text{TAP}}(\beta, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) &= -\ell_{\text{GVA}}(\beta, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) - \frac{1}{2} \sum_{i=1}^m \xi_i(\mu_i, \lambda_i; \mathbf{X}_i, \mathbf{Y}_i, \beta_0, \beta_1, \sigma^2) \\
 &= -\frac{1}{2} \sum_{i=1}^m \{1 + \log(2\pi\lambda_i)\} + \sum_{i=1}^m E \left\{ \Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i \right\} \\
 &\quad - \frac{1}{2} \sum_{i=1}^m \left( \text{Var}\{\Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} - \lambda_i \left[ E\{\Psi_i'(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} \right]^2 \right. \\
 &\quad \left. - \frac{\lambda_i^2}{2} \left[ E\{\Psi_i''(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} \right]^2 \right). \tag{6.5}
 \end{aligned}$$

Now define

$$Y_{i\bullet} \equiv \sum_{j=1}^n Y_{ij} \quad \text{and} \quad \mathcal{A}_i(\tilde{U}_i) \equiv \sum_{j=1}^n b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i).$$

Then, we can re-express  $\Psi_i(\tilde{U}_i)$  as

$$\Psi_i(\tilde{U}_i) = - \sum_{j=1}^n Y_{ij}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) + \mathcal{A}_i(\tilde{U}_i) - \sum_{j=1}^n c(Y_{ij}) + \frac{\tilde{U}_i^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2).$$

It follows that

$$\Psi_i'(\tilde{U}_i) = -Y_{i\bullet} + \mathcal{A}_i'(\tilde{U}_i) + \frac{\tilde{U}_i}{\sigma^2}$$

and

$$\Psi_i''(\tilde{U}_i) = \mathcal{A}_i''(\tilde{U}_i) + \frac{1}{\sigma^2}.$$

Also note that

$$\mathcal{A}_i'(\tilde{U}_i) = \sum_{j=1}^n b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \quad \text{and} \quad \mathcal{A}_i''(\tilde{U}_i) = \sum_{j=1}^n b''(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i).$$

## 6.4 The Thouless-Anderson-Palmer Approximate Negative Log-Likelihood

In this section, we present a result detailing the Thouless-Anderson-Palmer approximate negative log-likelihood expression for canonical link GLMMs after having solved for the equation in (6.5). First, for  $p, q \in \{0, 1, 2\}$ ,  $r > 0$  and  $s, t \in \mathbb{R}$ , define

$$\mathcal{Q}(p, q, r, s, t) \equiv (2\pi)^{-1/2} \int_{-\infty}^{\infty} (s + rx)^p b^{(q)}(t + rx) \exp\left(-\frac{x^2}{2}\right) dx. \tag{6.6}$$



Also, for  $r > 0$  and  $s, t \in \mathbb{R}$ , define

$$\mathcal{R}(r, s, t) \equiv (2\pi)^{-1/2} \int_{-\infty}^{\infty} b(s + rx)b(t + rx) \exp\left(-\frac{x^2}{2}\right) dx. \quad (6.7)$$

We then have the following result.

**Result 2.** *Consider the model set-up as in (6.1). Then the Thouless-Anderson-Palmer approximate negative likelihood is*

$$-\underline{\ell}_{TAP}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) = -\underline{\ell}_{GVA}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) - \frac{1}{2} \sum_{i=1}^m \xi_i(\mu_i, \lambda_i; \mathbf{X}_i, \mathbf{Y}_i, \beta_0, \beta_1, \sigma^2),$$

where

$$\begin{aligned} \underline{\ell}_{GVA}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) &= \sum_{i=1}^m E_{\tilde{U}_i} \left[ \sum_{j=1}^n \left\{ Y_{ij}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) - b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) + c(Y_{ij}) \right\} \right. \\ &\quad \left. - \frac{\tilde{U}_i^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right] + \frac{1}{2} \sum_{i=1}^m \{1 + \log(2\pi\lambda_i)\} \end{aligned}$$

and

$$\begin{aligned} &\xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2) \\ &= \sum_{j=1}^n \left\{ \left( 2Y_{i\bullet} \mu_i - \frac{\lambda_i + \mu_i^2}{\sigma^2} \right) \mathcal{Q}(0, 0, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i) \right. \\ &\quad + 2\lambda_i \left( Y_{i\bullet} - \frac{\mu_i}{\sigma^2} \right) \mathcal{Q}(0, 1, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i) \\ &\quad - 2Y_{i\bullet} \mathcal{Q}(1, 0, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i) \\ &\quad - \frac{\lambda_i^2}{\sigma^2} \mathcal{Q}(0, 2, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i) + \frac{1}{\sigma^2} \mathcal{Q}(2, 0, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i) \\ &\quad - \mathcal{Q}(0, 0, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i)^2 - \lambda_i \mathcal{Q}(0, 1, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i)^2 \\ &\quad - \frac{\lambda_i^2}{2} \mathcal{Q}(0, 2, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i)^2 \\ &\quad \left. + \sum_{j'=1}^n \mathcal{R}(\sqrt{\lambda_i}, \beta_0 + \beta_1 X_{ij} + \mu_i, \beta_0 + \beta_1 X_{ij'} + \mu_i) \right\}. \end{aligned}$$

The proof for Result 2 is provided in the appendix under Subsection 6.6.1.

## 6.5 Thouless-Anderson-Palmer Enhancement for Poisson Generalized Linear Mixed Models

In the Poisson case, when  $b(x) = \exp(x)$ ,  $Q(p, q, r, s, t)$  and  $R(r, s, t)$  admit exact expressions when evaluating Result 2. However, for general  $b$  functions, numerical integration is required for evaluating (6.6) and (6.7). Hence in this section, we work with a Poisson generalized linear mixed model and capitalize on having exact expressions to work with. Here we consider the model set-up of Hall et al. (2011) where

$$Y_{ij}|X_{ij}, U_i \text{ independent Poisson with mean } \exp(\beta_0^0 + \beta_1^0 X_{ij} + U_i), \quad (6.8)$$

such that the  $U_i$  are independent  $N(0, (\sigma^2)^0)$ .

In Subsection 6.5.1, we detail the expression for the Gaussian variational approximate log-likelihood for the model set-up in (6.8). Similarly, in Subsection 6.5.2, the expression for the Thouless-Anderson-Palmer approximate negative log-likelihood is presented. Subsection 6.5.3 moves on to study the optimisation issues present when using the TAP approach by exploring a simpler version of the set-up in (6.8) with  $m = 1$ . Lastly, in Subsection 6.5.4, a full simulation study is carried out to assess and compare the accuracy of the estimates of the model parameters across the GVA and TAP approaches for the model set-up in (6.8).

### 6.5.1 The Gaussian Variational Approximate Log-Likelihood for Simulation Set-Up

Substituting  $b(x) = \exp(x)$  and  $c(x) = -\log(x!)$  into (7.7), the Gaussian variational approximate log-likelihood for the model set-up in (6.8) is

$$\begin{aligned} \underline{\ell}_{\text{GVA}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) &= \sum_{i=1}^m E_{\tilde{U}_i} \left[ \sum_{j=1}^n \left\{ Y_{ij}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) - e^{\beta_0 + \beta_1 X_{ij} + \tilde{U}_i} - \log(Y_{ij}!) \right\} \right. \\ &\quad \left. - \frac{\tilde{U}_i^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right] + \frac{1}{2} \sum_{i=1}^m \{1 + \log(2\pi\lambda_i)\} \\ &= \sum_{i=1}^m \sum_{j=1}^n \left\{ Y_{ij}(\beta_0 + \beta_1 X_{ij} + \mu_i) - e^{\beta_0 + \beta_1 X_{ij} + \mu_i + \frac{1}{2}\lambda_i} \right\} \\ &\quad - \sum_{i=1}^m \frac{\mu_i^2 + \lambda_i}{2\sigma^2} - \frac{m}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^m \log(\lambda_i) + C, \end{aligned}$$

with vectors  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$  and  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_m)$  as the variational parameters. Also, note that  $C$  denotes the constant term independent of any model or variational parameters.

### 6.5.2 The Thouless-Anderson-Palmer Negative Approximate Log-Likelihood for Simulation Set-Up

We start out with the following result:

**Result 3.** *Consider the model set-up as in (6.8). Then the Onsager's correction term in the Thouless-Anderson-Palmer negative approximate likelihood is*

$$\begin{aligned} & -\frac{1}{2} \sum_{i=1}^m \xi_i(\mu_i, \lambda_i; \mathbf{X}_i, \mathbf{Y}_i, \beta_0, \beta_1, \sigma^2) \\ &= -\frac{1}{2} \sum_{i=1}^m \left[ \left\{ \exp(\lambda_i) - 1 - \lambda_i - \frac{1}{2} \lambda_i^2 \right\} \exp(2\mu_i + \lambda_i) \left\{ \sum_{j=1}^n \exp(\beta_0 + \beta_1 X_{ij}) \right\}^2 \right]. \end{aligned}$$

The details for the derivations leading to Result 3 is in the appendix under Subsection 6.6.3. By using Results 2 and 3 together, we have that the TAP negative approximate log-likelihood for the model set-up in (6.8) is

$$\begin{aligned} & -\underline{\ell}_{\text{TAP}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) \\ &= \sum_{i=1}^m \sum_{j=1}^n \left\{ -Y_{ij}(\beta_0 + \beta_1 X_{ij} + \mu_i) + e^{\beta_0 + \beta_1 X_{ij} + \mu_i + \frac{1}{2} \lambda_i} \right\} + \sum_{i=1}^m \frac{\mu_i^2 + \lambda_i}{2\sigma^2} + \frac{m}{2} \log(\sigma^2) \\ & -\frac{1}{2} \sum_{i=1}^m \log(\lambda_i) - \frac{1}{2} \sum_{i=1}^m \left[ \left\{ \exp(\lambda_i) - 1 - \lambda_i - \frac{1}{2} \lambda_i^2 \right\} \exp(2\mu_i + \lambda_i) \left\{ \sum_{j=1}^n \exp(\beta_0 + \beta_1 X_{ij}) \right\}^2 \right] \\ & + C, \end{aligned}$$

where  $C$  denotes a constant term independent of any model or variational parameters.

### 6.5.3 Optimisation Issues

In this section, we look into the problem of obtaining a local minimum for the Thouless-Anderson-Palmer approximate negative log-likelihood. We have chosen to focus on a simplified version of the model set-up in (6.8) to highlight the optimisation problems encountered.

### 6.5.3.1 A Simplified Version of the Optimisation Problem

First, consider the case where

$$\beta_1 = 0, \quad \sigma^2 = \sigma_{\text{fixed}}^2, \quad m = 1.$$

This situation corresponds to the following simplified Poisson mixed model:

$$Y_{1j}|U_1 \sim \text{Poisson}\left(e^{\beta_0+U_1}\right), \quad 1 \leq j \leq n, \quad U_1 \sim N(0, \sigma_{\text{fixed}}^2). \quad (6.9)$$

It follows that in this case, the Thouless-Anderson-Palmer approximate negative log-likelihood can be simplified to

$$\begin{aligned} -\underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1) &= -(\beta_0 + \mu_1) \sum_{j=1}^n Y_{ij} + ne^{\beta_0+\mu_1+\frac{1}{2}\lambda_1} + \frac{\mu_1^2 + \lambda_1}{2\sigma_{\text{fixed}}^2} - \frac{1}{2} \log(\lambda_1) \\ &\quad - \frac{1}{2} \left\{ \exp(\lambda_1) - 1 - \lambda_1 - \frac{1}{2}\lambda_1^2 \right\} n^2 \exp(2\mu_1 + \lambda_1 + 2\beta_0) + C, \end{aligned}$$

where  $C$  denotes a constant term independent of any model or variational parameters.

To facilitate the optimization process, we compute the partial derivatives required for minimising the Thouless-Anderson-Palmer approximate negative log-likelihood. The first order partial derivatives are as follows:

$$\begin{aligned} \frac{-\partial \underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)}{\partial \beta_0} &= -\sum_{j=1}^n Y_{ij} + n \exp\left(\beta_0 + \mu_1 + \frac{1}{2}\lambda_1\right) \\ &\quad - n^2 \left\{ \exp(\lambda_1) - 1 - \lambda_1 - \frac{1}{2}\lambda_1^2 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0), \\ \frac{-\partial \underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)}{\partial \mu_1} &= -\sum_{j=1}^n Y_{ij} + n \exp\left(\beta_0 + \mu_1 + \frac{1}{2}\lambda_1\right) + \frac{\mu_1}{\sigma_{\text{fixed}}^2} \\ &\quad - n^2 \left\{ \exp(\lambda_1) - 1 - \lambda_1 - \frac{1}{2}\lambda_1^2 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0) \quad \text{and} \\ \frac{-\partial \underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)}{\partial \lambda_1} &= \frac{n}{2} \exp\left(\beta_0 + \mu_1 + \frac{1}{2}\lambda_1\right) + \frac{1}{2\sigma_{\text{fixed}}^2} - \frac{1}{2\lambda_1} \\ &\quad - \frac{n^2}{2} \left\{ \exp(\lambda_1) - 1 - \lambda_1 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0) \\ &\quad - \frac{n^2}{2} \left\{ \exp(\lambda_1) - 1 - \lambda_1 - \frac{1}{2}\lambda_1^2 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0) \\ &= \frac{n}{2} \exp\left(\beta_0 + \mu_1 + \frac{1}{2}\lambda_1\right) + \frac{1}{2\sigma_{\text{fixed}}^2} - \frac{1}{2\lambda_1} \\ &\quad - \frac{n^2}{2} \left\{ 2 \exp(\lambda_1) - 2 - 2\lambda_1 - \frac{1}{2}\lambda_1^2 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0). \end{aligned}$$

Next, the second order partial derivatives in the diagonal of the Hessian matrix are as follows:

$$\begin{aligned} \frac{-\partial \underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)}{\partial \beta_0^2} &= n \exp\left(\beta_0 + \mu_1 + \frac{1}{2}\lambda_1\right) \\ &\quad - 2n^2 \left\{ \exp(\lambda_1) - 1 - \lambda_1 - \frac{1}{2}\lambda_1^2 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0), \\ \frac{-\partial \underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)}{\partial \mu_1^2} &= n \exp\left(\beta_0 + \mu_1 + \frac{1}{2}\lambda_1\right) + \frac{1}{\sigma_{\text{fixed}}^2} \\ &\quad - 2n^2 \left\{ \exp(\lambda_1) - 1 - \lambda_1 - \frac{1}{2}\lambda_1^2 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0) \quad \text{and} \\ \frac{-\partial \underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)}{\partial \lambda_1^2} &= \frac{n}{4} \exp\left(\beta_0 + \mu_1 + \frac{1}{2}\lambda_1\right) + \frac{1}{2\lambda_1^2} \\ &\quad - \frac{n^2}{2} \left\{ 4 \exp(\lambda_1) - 4 - 3\lambda_1 - \frac{1}{2}\lambda_1^2 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0). \end{aligned}$$

Lastly, the second order partial derivatives in the off-diagonals of the Hessian matrix are as follows:

$$\begin{aligned} \frac{-\partial \underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)}{\partial \beta_0 \partial \mu_1} &= n \exp\left(\beta_0 + \mu_1 + \frac{1}{2}\lambda_1\right) \\ &\quad - 2n^2 \left\{ \exp(\lambda_1) - 1 - \lambda_1 - \frac{1}{2}\lambda_1^2 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0), \\ \frac{-\partial \underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)}{\partial \beta_0 \partial \lambda_1} &= \frac{n}{2} \exp\left(\beta_0 + \mu_1 + \frac{1}{2}\lambda_1\right) \\ &\quad - n^2 \left\{ 2 \exp(\lambda_1) - 2 - 2\lambda_1 - \frac{1}{2}\lambda_1^2 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0) \quad \text{and} \\ \frac{-\partial \underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)}{\partial \mu_1 \partial \lambda_1} &= \frac{n}{2} \exp\left(\beta_0 + \mu_1 + \frac{1}{2}\lambda_1\right) \\ &\quad - n^2 \left\{ 2 \exp(\lambda_1) - 2 - 2\lambda_1 - \frac{1}{2}\lambda_1^2 \right\} \exp(2\mu_1 + \lambda_1 + 2\beta_0). \end{aligned}$$

### 6.5.3.2 Simplified Simulation Study

A simplified simulation study with the model set-up in (6.9) was run with the following settings where

$$\beta_0^0 = -0.2, \quad \sigma_{\text{fixed}} = 0.3, \quad \text{and} \quad n = 20$$

with 10000 replications. The search for a local minimiser of  $-\underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)$  in a three-dimensional space can be challenging. Hence, we use the following strategy involving the `optim()` function in R:

1. First, initialise values for  $(\beta_0, \mu_1, \lambda_1)$ .
2. Next, based on the initial values specified in (1), carry out a large number of Nelder-Mead iterations, via `optim()`, to search for a local minimiser of  $-\underline{\ell}_{\text{TAP}}(\beta_0, \mu_1, \lambda_1)$ .
3. Lastly, use the results from step 2 to obtain starting values for the Broyden-Fletcher-Goldberg-Shanno quasi-Newton approach, via `optim()`, to improve on the result from implementing step 2.

Whilst this strategy may seem reasonable, it turns out that it is prone to erratic behaviour if the initial value in step 1 is a poor choice. For the simplified model set-up in (6.9), the expressions for the partial derivatives can be analysed to develop stationary point equations and determine suitable starting values.

### 6.5.3.3 Results and Conclusion

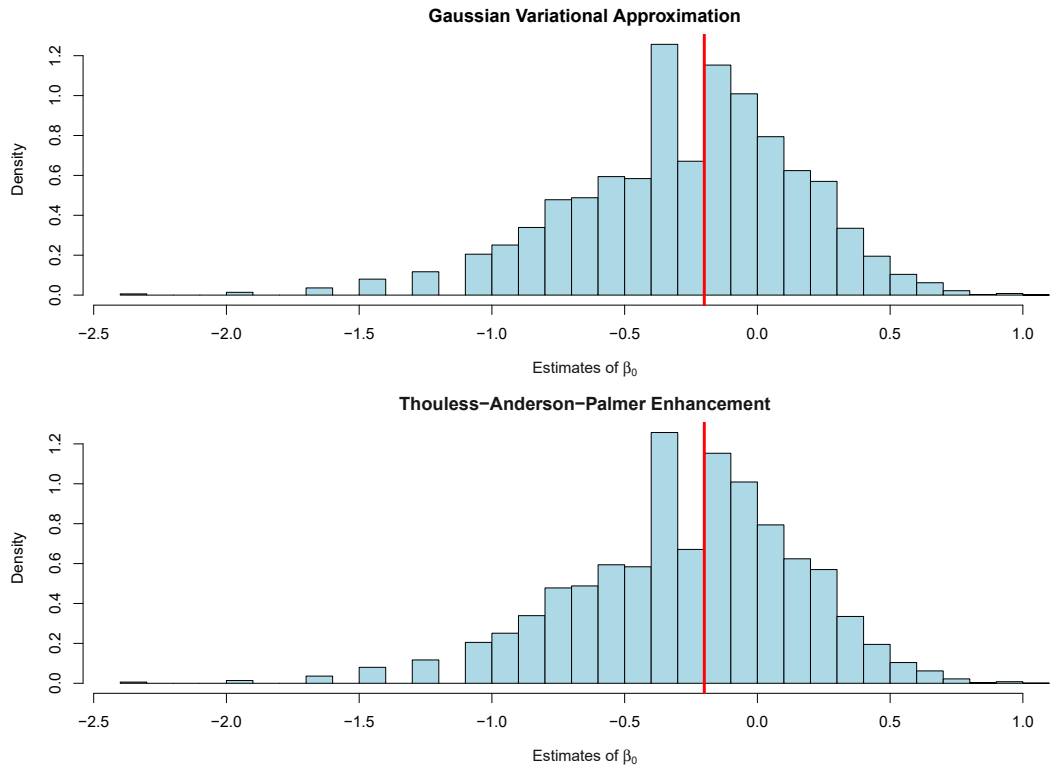


Figure 6.1: Histograms of the Gaussian variational approximation estimates of  $\beta_0$  and the Thouless-Anderson-Palmer enhancement estimates of  $\beta_0$ .

After having found finite local minima of the Thouless-Anderson-Palmer approximate log-likelihood surface for this simulation study, the estimates of  $\beta_0$  were gathered from implementing both the GVA and TAP approaches. Figure 6.1 displays the histograms of the estimates of  $\beta_0$  obtained from the GVA approach and TAP enhancement approach. Note that the vertical red line is situated at the true value of  $\beta_0^0 = -0.2$ . We see that for both approximation approaches, the estimates are distributed about the true value.

We then computed the absolute error values for the estimates. For a generic estimate  $\widehat{\beta}_0$ , the absolute error is computed as

$$|\widehat{\beta}_0 - \beta_0^0|.$$

Applying this definition to the vectors of the estimates of  $\beta_0$  obtained from both the GVA and TAP enhancement approaches, we obtain vectors of the GVA absolute errors and the TAP enhancement absolute errors of length 10000 each. Next, we obtained the pairwise differences with ordering as follows:

$$\text{Pairwise difference} = (\text{TAP enhancement absolute error}) - (\text{GVA absolute error}).$$

Figure 6.2 shows a histogram of the pairwise differences.

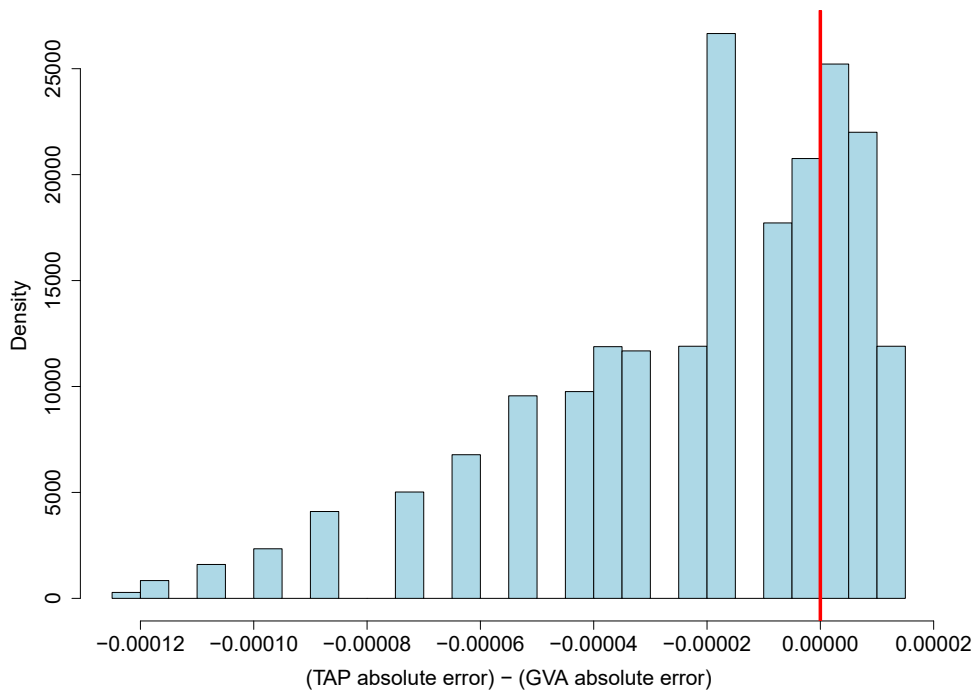


Figure 6.2: Histogram of the pairwise differences between the TAP enhancement absolute errors and GVA absolute errors.

Figure 6.2 suggests a statistically significant improvement due to using the TAP enhancement. However, the improvement seems to be quite slight from a practical standpoint. On the positive side, the results from this simulation study shows that with some care it is possible to get estimates of Poisson mixed model parameters that are improved by the TAP enhancement. On the negative side, the numerical problem is quite challenging even for the  $m = 1$  and  $\sigma^2 = \sigma_{\text{fixed}}^2$  case. In the more practically relevant case where  $m$  is in the hundreds and  $\sigma^2$  is estimated instead of being fixed, obtaining local minima for the Thouless-Anderson-Palmer approximate negative log-likelihood could be quite challenging. We will investigate this further in the next subsection.

#### 6.5.4 Simulation Study

A full simulation study was run to investigate and compare the accuracy of approximations from the GVA approach and the TAP enhancement approach. In this simulation study, a Poisson linear mixed model was used following the set-up in (6.8). The values for the true parameter vector  $(\beta_0^0, \beta_1^0, (\sigma^2)^0)$  were chosen from the following possible set of pre-determined values

$$\{(-0.3, 0.2, 0.5), (2.2, -0.1, 0.16), (1.2, 0.4, 0.1), (0.02, 1.3, 1), (-0.3, 0.2, 0.1)\}.$$

and the distribution of the  $X_{ij}$  was taken to be either  $N(0, 1)$  or Uniform  $(-1, 1)$ . The number of groups in the simulated data,  $m$ , varied over the set  $\{100, 200, \dots, 1000\}$  and the number of observations present within each group,  $n$ , was fixed at  $m/10$ . 100 replications were simulated for every possible combination of the true parameter vector, chosen  $X_{ij}$  distribution and value of  $(m, n)$  pair. For each sample, estimates for the true parameter vector  $(\beta_0^0, \beta_1^0, (\sigma^2)^0)$ , were obtained using similar steps to those outlined in Section 6.5.3.2 for the TAP enhancement approach. Estimates for  $(\beta_0^0, \beta_1^0, (\sigma^2)^0)$  via the GVA approach were also obtained. For generic estimates  $\hat{\beta}_0, \hat{\beta}_1$  and  $\hat{\sigma}^2$ , their absolute error values were computed as

$$|\hat{\beta}_0 - \beta_0^0|, \quad |\hat{\beta}_1 - \beta_1^0| \quad \text{and} \quad |\hat{\sigma}^2 - (\sigma^2)^0|$$

respectively. Applying this definition to the vectors of the estimates of  $\beta_0, \beta_1$  and  $\sigma^2$  obtained from both the GVA and TAP enhancement approaches, we obtain vectors of the GVA absolute errors and the TAP enhancement absolute errors of length 100 each. Like in the previous simulation study, we obtained the pairwise differences with



ordering as follows:

$$\text{Pairwise difference} = (\text{TAP enhancement absolute error}) - (\text{GVA absolute error}).$$

Grouped boxplots for each model parameter were then produced to graphically demonstrate the mean and spread for each vector of pairwise differences.

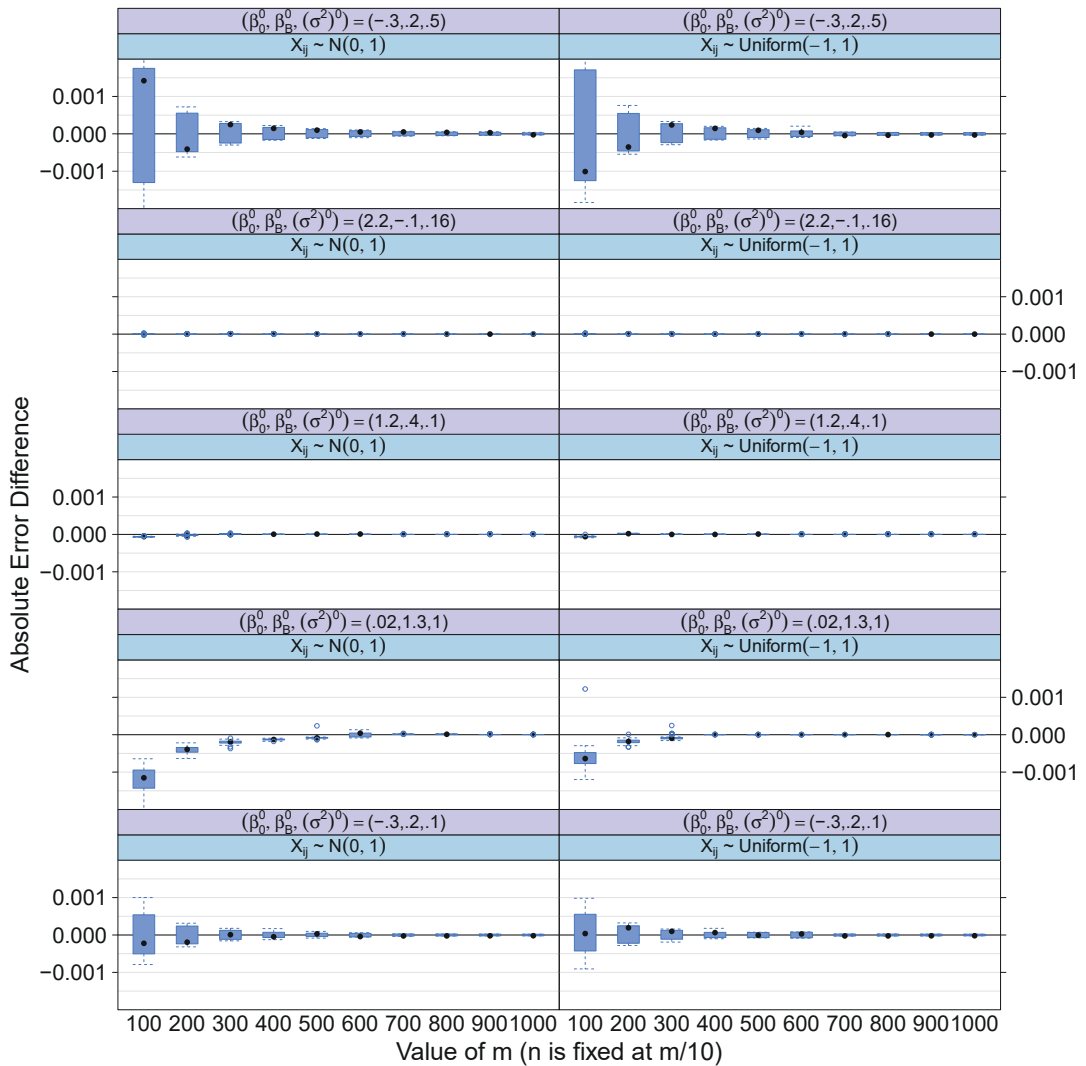


Figure 6.3: Grouped boxplots representing the pairwise difference in absolute errors (TAP enhancement absolute error - GVA absolute error) in estimating  $\beta_0$ . The values of  $m$  are 100, 200, . . . , 1000 while the value of  $n$  is fixed at  $m/10$ .

Figure 6.3 shows that for values of  $m$  in the lower hundreds, if there is a difference in accuracy in estimating  $\beta_0$  between the GVA approach and TAP enhancement approach,

the TAP enhancement approach seems to give slightly better estimates in most of those cases. There seems to be a more significant improvement when the true value  $\beta_0^0$  is close to 0 coupled with a large  $(\sigma^2)^0$  value. However, as the value of  $m$  increases, the difference in absolute errors between the TAP enhancement approach and GVA approach diminishes.

Inspection of Figure 6.4 reveals almost no observable difference in the accuracy of estimates of  $\beta_1$  between the GVA approach and TAP enhancement approach, for all values of  $m$ .

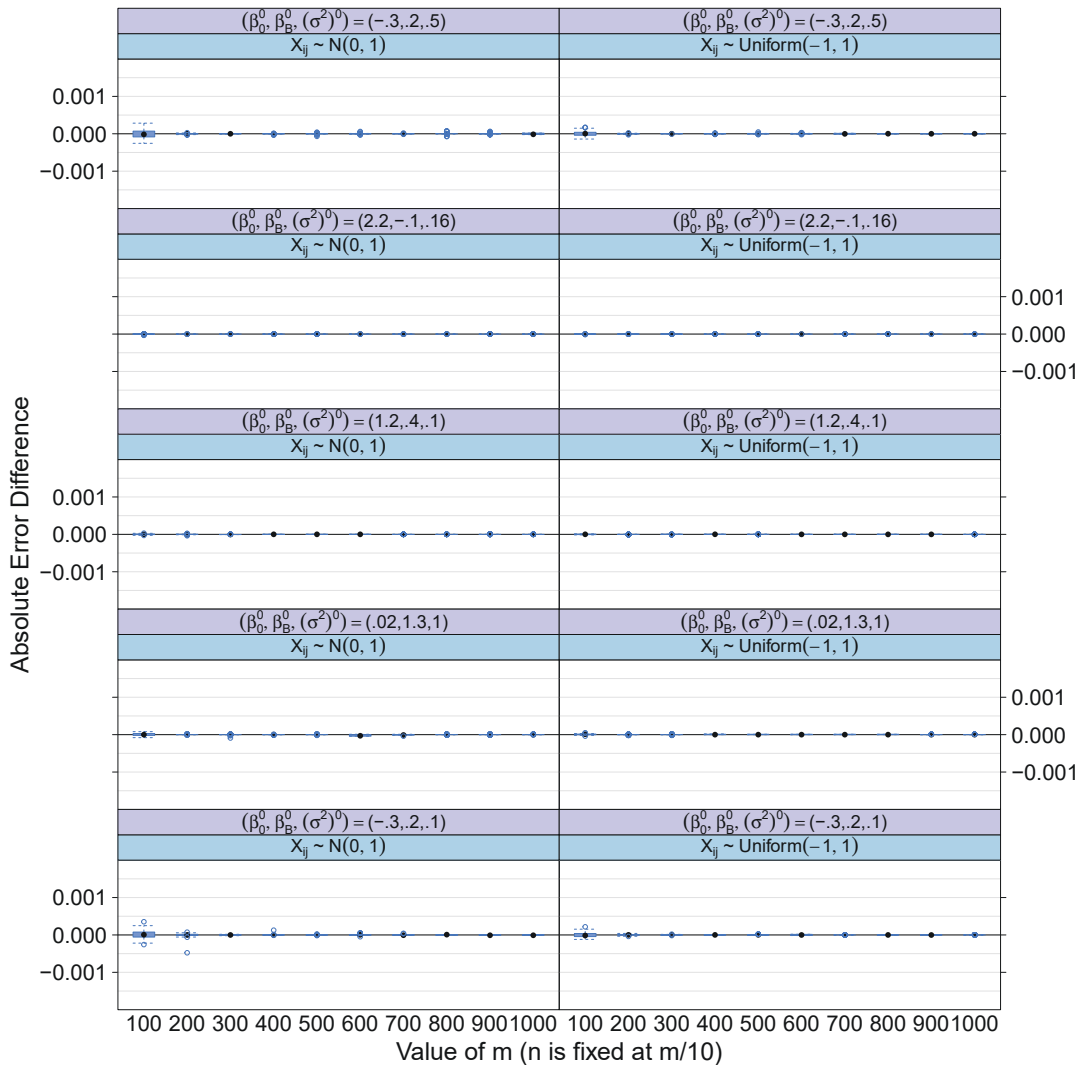


Figure 6.4: Grouped boxplots representing the pairwise difference in absolute errors (TAP enhancement absolute error - GVA absolute error) in estimating  $\beta_1$ . The values of  $m$  are 100, 200, ..., 1000 while the value of  $n$  is fixed at  $m/10$ .

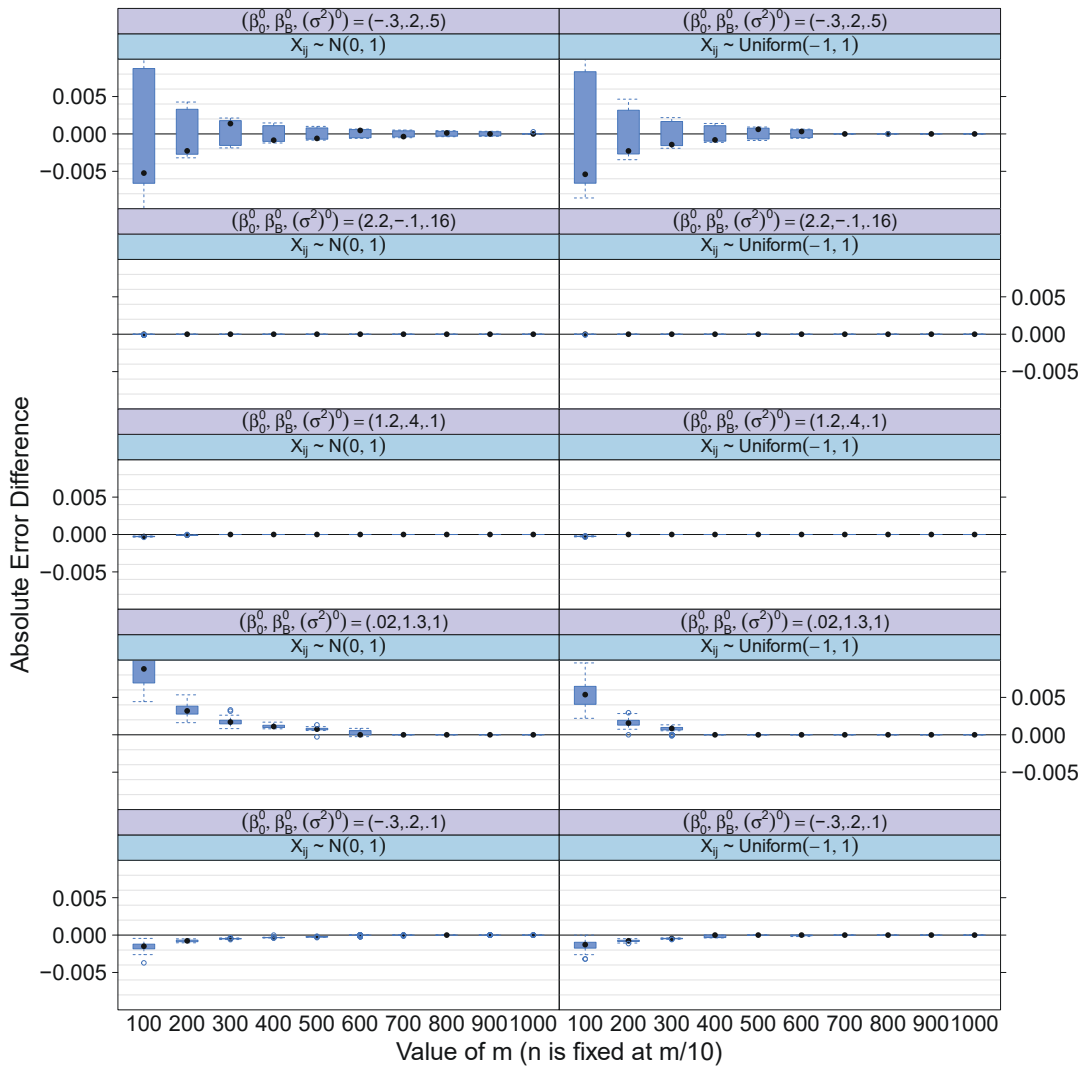


Figure 6.5: Grouped boxplots representing the pairwise difference in absolute errors (TAP enhancement absolute error - GVA absolute error) in estimating  $\sigma^2$ . The values of  $m$  are 100, 200, ..., 1000 while the value of  $n$  is fixed at  $m/10$ .

In Figure 6.5, for values of  $m$  in the lower hundreds, there are 6 subplots where a visible difference in the accuracy in estimating  $\sigma^2$  between the GVA approach and TAP enhancement approach is present. The TAP enhancement approach seems to give slightly better estimates in 4 out of 6 such cases. In contrast to Figure 6.3, the TAP estimates are worse when the true value  $\beta_0^0$  is close to 0 coupled with a large  $(\sigma^2)^0$  value. Once again, as the value of  $m$  increases, there is no significant difference in the absolute errors between the TAP enhancement approach and GVA approach.

To conclude, the Thouless-Anderson-Palmer enhancement approach suggests a slight yet statistically significant improvement as compared to using the Gaussian variational approximation approach for small datasets ( $m \leq 200$ ). However, although the TAP estimates for  $\beta_0^0$  fare better than the GVA estimates when the true value  $\beta_0^0$  is close to 0 coupled with a large  $(\sigma^2)^0$  value, that is not the case for the TAP estimates of  $(\sigma^2)^0$ . It is also worth exploring how the TAP approach performs against the GVA approach for generalized linear mixed models with response distributions other than the Poisson family.

## 6.6 Appendix

### 6.6.1 Proof of Result 2

This appendix contains the details for the derivations leading up to Result 2.

#### 6.6.1.1 Main Quantity in Onsager's Correction Term

Note that the main quantity in the Onsager's correction term is

$$\begin{aligned} & \xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2) \\ &= \text{Var}\{\Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} - \lambda_i \left[ E\{\Psi_i'(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} \right]^2 - \frac{\lambda_i^2}{2} \left[ E\{\Psi_i''(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} \right]^2. \end{aligned} \quad (6.10)$$

This quantity consists of three terms which we will work with in the next three subsections.

#### 6.6.1.2 An Explicit Expression for the First Term in (6.10)

In this subsection, we find an explicit expression for  $\text{Var}\{\Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\}$ . Firstly, note that

$$\begin{aligned} & \text{Var}\{\Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} \\ &= \text{Var}\left( \frac{\tilde{U}_i^2}{2\sigma^2} - Y_{i\bullet} \tilde{U}_i + \mathcal{A}_i(\tilde{U}_i) \middle| \mathbf{Y}_i, \mathbf{X}_i \right) \\ &= \text{Var}\left( \frac{\tilde{U}_i^2}{2\sigma^2} - Y_{i\bullet} \tilde{U}_i \middle| \mathbf{Y}_i, \mathbf{X}_i \right) + 2\text{Cov}\left( \frac{\tilde{U}_i^2}{2\sigma^2} - Y_{i\bullet} \tilde{U}_i, \mathcal{A}_i(\tilde{U}_i) \middle| \mathbf{Y}_i, \mathbf{X}_i \right) + \text{Var}\left( \mathcal{A}_i(\tilde{U}_i) \middle| \mathbf{Y}_i, \mathbf{X}_i \right). \end{aligned}$$

Now we treat each term in the expression for  $\text{Var}\{\Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\}$  individually.

Treatment of the  $\text{Var}\left(\frac{\tilde{U}_i^2}{2\sigma^2} - Y_{i\bullet}\tilde{U}_i \mid \mathbf{Y}_i, \mathbf{X}_i\right)$  Term

We first obtain an expression for

$$\text{Var}(\mathbf{a}\tilde{U}_i^2 + \mathbf{b}\tilde{U}_i + \mathbf{c}) = \text{Var}(\mathbf{a}\tilde{U}_i^2 + \mathbf{b}\tilde{U}_i)$$

for general  $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}$ . Then, note that

$$\begin{aligned} \text{Var}_{\tilde{U}_i}(\tilde{U}_i^2) &= E(\tilde{U}_i^4) - \{E(\tilde{U}_i^2)\}^2 \\ &= (\mu_i^4 + 6\mu_i^2\lambda_i + 3\lambda_i^2) - (\mu_i^2 + \lambda_i)^2 \\ &= 2\lambda_i^2 + 4\mu_i^2\lambda_i \\ &= 2\lambda_i(2\mu_i^2 + \lambda_i) \end{aligned}$$

and

$$\begin{aligned} \text{Cov}_{\tilde{U}_i}(\tilde{U}_i, \tilde{U}_i^2) &= E(\tilde{U}_i^3) - E(\tilde{U}_i)E(\tilde{U}_i^2) \\ &= (\mu_i + 3\mu_i\lambda_i) - (\mu_i)(\mu_i^2 + \lambda_i) \\ &= 2\mu_i\lambda_i. \end{aligned}$$

Then we have

$$\begin{aligned} \text{Var}(\mathbf{a}\tilde{U}_i^2 + \mathbf{b}\tilde{U}_i) &= \mathbf{a}\text{Var}(\tilde{U}_i^2) + 2\mathbf{a}\mathbf{b}\text{Cov}(\tilde{U}_i^2, \tilde{U}_i) + \mathbf{b}^2\text{Var}(\tilde{U}_i) \\ &= 2\mathbf{a}^2\lambda_i(2\mu_i^2 + \lambda_i) + 4\mathbf{a}\mathbf{b}\lambda_i\mu_i + \mathbf{b}^2\lambda_i \\ &= 4\mathbf{a}^2\lambda_i\mu_i^2 + 2\mathbf{a}^2\lambda_i^2 + 4\mathbf{a}\mathbf{b}\lambda_i\mu_i + \mathbf{b}^2\lambda_i. \end{aligned}$$

Setting

$$\mathbf{a} = \frac{1}{2\sigma^2} \quad \text{and} \quad \mathbf{b} = -Y_{i\bullet},$$

we get

$$\text{Var}\left(\frac{\tilde{U}_i^2}{2\sigma^2} - Y_{i\bullet}\tilde{U}_i \mid \mathbf{Y}_i, \mathbf{X}_i\right) = \frac{2\lambda_i\mu_i^2 + \lambda_i^2}{2\sigma^4} - \frac{2Y_{i\bullet}\lambda_i\mu_i}{\sigma^2} + \lambda_i(Y_{i\bullet})^2. \quad (6.11)$$

Treatment of the  $2\text{Cov}\left(\frac{\tilde{U}_i^2}{2\sigma^2} - Y_{i\bullet}\tilde{U}_i, \mathcal{A}_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i\right)$  Term

Next, we obtain an expression for

$$\text{Cov}(\mathbf{a}\tilde{U}_i^2 + \mathbf{b}\tilde{U}_i, \mathcal{A}_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i) = \mathbf{a}\text{Cov}(\tilde{U}_i^2, \mathcal{A}_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i) + \mathbf{b}\text{Cov}(\tilde{U}_i, \mathcal{A}_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i)$$

for general  $\mathbf{a}, \mathbf{b} \in \mathbb{R}$ . Next note that, for  $k \in \{1, 2\}$ ,

$$\begin{aligned} &\text{Cov}\left(\tilde{U}_i^k, \mathcal{A}_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i\right) \\ &= \sum_{j=1}^n \text{Cov}\left(\tilde{U}_i^k, b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right) \\ &= \sum_{j=1}^n E\left(\tilde{U}_i^k b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right) - E(\tilde{U}_i^k) \sum_{j=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right). \end{aligned}$$

Therefore,

$$\begin{aligned} & \text{Cov}\left(\tilde{U}_i, \mathcal{A}_i(\tilde{U}_i) \mid \mathbf{X}_i\right) \\ &= \sum_{j=1}^n E\left(\tilde{U}_i b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right) - \mu_i \sum_{j=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right), \end{aligned}$$

and

$$\begin{aligned} & \text{Cov}\left(\tilde{U}_i^2, \mathcal{A}_i(\tilde{U}_i) \mid \mathbf{X}_i\right) \\ &= \sum_{j=1}^n E\left(\tilde{U}_i^2 b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right) - (\lambda_i + \mu_i^2) \sum_{j=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right). \end{aligned}$$

It follows from these results that

$$\begin{aligned} & \text{Cov}(\mathbf{a}\tilde{U}_i^2 + \mathbf{b}\tilde{U}_i, \mathcal{A}_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i) \\ &= 2\mathbf{a} \sum_{j=1}^n E\left(\tilde{U}_i^2 b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right) + 2\mathbf{b} \sum_{j=1}^n E\left(\tilde{U}_i b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right) \\ &\quad - 2\left\{\mathbf{a}(\lambda_i + \mu_i^2) + \mathbf{b}\mu_i\right\} \sum_{j=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right). \end{aligned}$$

Setting

$$\mathbf{a} = \frac{1}{2\sigma^2} \quad \text{and} \quad \mathbf{b} = -Y_{i\bullet},$$

we obtain

$$\begin{aligned} & 2\text{Cov}\left(\frac{\tilde{U}_i^2}{2\sigma^2} - Y_{i\bullet}\tilde{U}_i, \mathcal{A}_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i\right) \\ &= \frac{1}{\sigma^2} \sum_{j=1}^n E\left(\tilde{U}_i^2 b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right) - 2Y_{i\bullet} \sum_{j=1}^n E\left(\tilde{U}_i b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right) \\ &\quad + \left(2Y_{i\bullet}\mu_i - \frac{\lambda_i + \mu_i^2}{\sigma^2}\right) \sum_{j=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i\right). \end{aligned}$$

(6.12)

Treatment of the  $\text{Var}(\mathcal{A}_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i)$  Term

It follows that

$$\begin{aligned}
& \text{Var}(\mathcal{A}_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i) \\
&= \text{Var}\left(\sum_{j=1}^n b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right) \\
&= \sum_{j=1}^n \text{Var}(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i) + \sum_{j \neq j'} \text{Cov}(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i), b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) | \mathbf{X}_i) \\
&= \sum_{j=1}^n E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)^2 | \mathbf{X}_i) - \sum_{j=1}^n \left\{E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i)\right\}^2 \\
&\quad + \sum_{j \neq j'} E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) | \mathbf{X}_i) \\
&\quad - \sum_{j \neq j'} \left\{E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i) E(b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) | \mathbf{X}_i)\right\}
\end{aligned}$$

The previous expression then simplifies to

$$\begin{aligned}
& \sum_{j=1}^n E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)^2 | \mathbf{X}_i) + \sum_{j \neq j'} E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) | \mathbf{X}_i) \\
& \quad - \left\{ \sum_{j=1}^n E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i) \right\}^2. \tag{6.13}
\end{aligned}$$

The Resultant  $\text{Var}\{\Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\}$  Expression

Putting together the expressions from (6.11), (6.12) and (6.13), we have

$$\begin{aligned}
& \text{Var}\{\Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} \\
&= \frac{2\lambda_i \mu_i^2 + \lambda_i^2}{2\sigma^4} - \frac{2Y_{i\bullet} \lambda_i \mu_i}{\sigma^2} + \lambda_i (Y_{i\bullet})^2 + \frac{1}{\sigma^2} \sum_{j=1}^n E(\tilde{U}_i^2 b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i) \\
&\quad - 2Y_{i\bullet} \sum_{j=1}^n E(\tilde{U}_i b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i) \\
&\quad + \left(2Y_{i\bullet} \mu_i - \frac{\lambda_i + \mu_i^2}{\sigma^2}\right) \sum_{j=1}^n E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i) \\
&\quad + \sum_{j=1}^n E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)^2 | \mathbf{X}_i) + \sum_{j \neq j'} E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) | \mathbf{X}_i) \\
&\quad - \left\{ \sum_{j=1}^n E(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i) \right\}^2.
\end{aligned}$$

However, note that,

$$\begin{aligned} & \sum_{j=1}^n E \left( b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)^2 \mid \mathbf{X}_i \right) + \sum_{j \neq j'} E \left( b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) \mid \mathbf{X}_i \right) \\ &= \sum_{j=1}^n \sum_{j'=1}^n E \left( b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) \mid \mathbf{X}_i \right). \end{aligned}$$

Therefore,

$$\begin{aligned} & \text{Var}\{\Psi_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i\} \\ &= \frac{2\lambda_i \mu_i^2 + \lambda_i^2}{2\sigma^4} - \frac{2Y_{i\bullet} \lambda_i \mu_i}{\sigma^2} + \lambda_i (Y_{i\bullet})^2 + \frac{1}{\sigma^2} \sum_{j=1}^n E \left( \tilde{U}_i^2 b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i \right) \\ & \quad - 2Y_{i\bullet} \sum_{j=1}^n E \left( \tilde{U}_i b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i \right) + \left( 2Y_{i\bullet} \mu_i - \frac{\lambda_i + \mu_i^2}{\sigma^2} \right) \sum_{j=1}^n E \left( b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i \right) \\ & \quad + \sum_{j=1}^n \sum_{j'=1}^n E \left( b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) \mid \mathbf{X}_i \right) - \left\{ \sum_{j=1}^n E \left( b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i \right) \right\}^2. \end{aligned}$$

### 6.6.1.3 An Explicit Expression for the Second Term in (6.10)

In this subsection, we find an explicit expression for  $-\lambda_i \left[ E\{\Psi'_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i\} \right]^2$ .

First note that

$$\Psi'_i(\tilde{U}_i) = -Y_{i\bullet} + \mathcal{A}'_i(\tilde{U}_i) + \frac{\tilde{U}_i}{\sigma^2}.$$

Hence,

$$E\{\Psi'_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i\} = -Y_{i\bullet} + \sum_{j=1}^n E \left( b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i \right) + \frac{\mu_i}{\sigma^2}.$$

Then we have,

$$\begin{aligned} -\lambda_i \left[ E\{\Psi'_i(\tilde{U}_i) \mid \mathbf{Y}_i, \mathbf{X}_i\} \right]^2 &= \lambda_i \left\{ -Y_{i\bullet} + \sum_{j=1}^n E \left( b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i \right) + \frac{\mu_i}{\sigma^2} \right\}^2 \\ &= -\lambda_i (Y_{i\bullet})^2 - \lambda_i \left\{ \sum_{j=1}^n E \left( b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i \right) \right\}^2 \\ & \quad - \frac{\lambda_i \mu_i^2}{\sigma^4} + 2\lambda_i Y_{i\bullet} \left\{ \sum_{j=1}^n E \left( b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i \right) \right\} \\ & \quad + \frac{2Y_{i\bullet} \lambda_i \mu_i}{\sigma^2} - \frac{2\mu_i \lambda_i}{\sigma^2} \left\{ \sum_{j=1}^n E \left( b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \mid \mathbf{X}_i \right) \right\}. \end{aligned}$$



### 6.6.1.4 An Explicit Expression for the Third Term in (6.10)

In this subsection, we find an explicit expression for  $-\frac{\lambda_i^2}{2} \left[ E\{\Psi_i''(\tilde{U}_i)|\mathbf{Y}_i, \mathbf{X}_i\} \right]^2$ . Note that

$$\Psi_i''(\tilde{U}_i) = \mathcal{A}_i''(\tilde{U}_i) + \frac{1}{\sigma^2}.$$

Therefore,

$$E\{\Psi_i''(\tilde{U}_i)|\mathbf{Y}_i, \mathbf{X}_i\} = \sum_{j=1}^n E\left(b''(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) + \frac{1}{\sigma^2}.$$

Following that, we have,

$$\begin{aligned} -\frac{\lambda_i^2}{2} \left[ E\{\Psi_i''(\tilde{U}_i)|\mathbf{Y}_i, \mathbf{X}_i\} \right]^2 &= -\frac{\lambda_i^2}{2} \left( \sum_{j=1}^n E\left(b''(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) + \frac{1}{\sigma^2} \right)^2 \\ &= -\frac{\lambda_i^2}{2} \left( \sum_{j=1}^n E\left(b''(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) \right)^2 \\ &\quad - \frac{\lambda_i^2}{\sigma^2} \sum_{j=1}^n E\left(b''(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) - \frac{\lambda_i^2}{2\sigma^4}. \end{aligned}$$

### 6.6.1.5 The Resultant Expression for the Main Quantity in the Onsager's Correction Term

Based on the results from the previous three subsections, we have the following

$$\begin{aligned} \xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2) &= \frac{2\lambda_i\mu_i^2 + \lambda_i^2}{2\sigma^4} - \frac{2Y_{i\bullet}\lambda_i\mu_i}{\sigma^2} + \lambda_i(Y_{i\bullet})^2 + \frac{1}{\sigma^2} \sum_{j=1}^n E\left(\tilde{U}_i^2 b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) \\ &\quad - 2Y_{i\bullet} \sum_{j=1}^n E\left(\tilde{U}_i b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) + \left(2Y_{i\bullet}\mu_i - \frac{\lambda_i + \mu_i^2}{\sigma^2}\right) \sum_{j=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) \\ &\quad + \sum_{j=1}^n \sum_{j'=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i)|\mathbf{X}_i\right) - \left\{ \sum_{j=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) \right\}^2 - \lambda_i(Y_{i\bullet})^2 \\ &\quad - \lambda_i \left\{ \sum_{j=1}^n E\left(b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) \right\}^2 - \frac{\lambda_i\mu_i^2}{\sigma^4} + 2\lambda_i Y_{i\bullet} \left\{ \sum_{j=1}^n E\left(b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) \right\} + \frac{2Y_{i\bullet}\lambda_i\mu_i}{\sigma^2} \\ &\quad - \frac{2\mu_i\lambda_i}{\sigma^2} \left\{ \sum_{j=1}^n E\left(b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) \right\} - \frac{\lambda_i^2}{2\sigma^4} - \frac{\lambda_i^2}{2} \left( \sum_{j=1}^n E\left(b''(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right) \right)^2 \\ &\quad - \frac{\lambda_i^2}{\sigma^2} \sum_{j=1}^n E\left(b''(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)|\mathbf{X}_i\right). \end{aligned}$$

By simplifying this expression further, we have,

$$\begin{aligned}
& \xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2) \\
&= \frac{1}{\sigma^2} \sum_{j=1}^n E\left(\tilde{U}_i^2 b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right) - 2Y_{i\bullet} \sum_{j=1}^n E\left(\tilde{U}_i b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right) \\
&+ \left(2Y_{i\bullet} \mu_i - \frac{\lambda_i + \mu_i^2}{\sigma^2}\right) \sum_{j=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right) \\
&+ \sum_{j=1}^n \sum_{j'=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) | \mathbf{X}_i\right) \\
&- \left\{ \sum_{j=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right) \right\}^2 - \lambda_i \left\{ \sum_{j=1}^n E\left(b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right) \right\}^2 \\
&+ 2\lambda_i \left(Y_{i\bullet} - \frac{\mu_i}{\sigma^2}\right) \left\{ \sum_{j=1}^n E\left(b'(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right) \right\} \\
&- \frac{\lambda_i^2}{2} \left( \sum_{j=1}^n E\left(b''(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right) \right)^2 - \frac{\lambda_i^2}{\sigma^2} \sum_{j=1}^n E\left(b''(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right).
\end{aligned}$$

For any  $n \times 1$  random vector  $\mathbf{X}_i = (X_{i1}, \dots, X_{in})$  and  $p, q \in \mathbb{Z}^+$ , define

$$\mathcal{B}(p, q, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) \equiv \sum_{j=1}^n E\left(\tilde{U}_i^p b^{(q)}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right),$$

for  $1 \leq i \leq m$ . Also define

$$\mathcal{C}(\beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) \equiv \sum_{j=1}^n \sum_{j'=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) | \mathbf{X}_i\right),$$

for  $1 \leq i \leq m$ . We then have

$$\begin{aligned}
& \xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2) \\
&= \left(2Y_{i\bullet} \mu_i - \frac{\lambda_i + \mu_i^2}{\sigma^2}\right) \mathcal{B}(0, 0, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) \\
&+ 2\lambda_i \left(Y_{i\bullet} - \frac{\mu_i}{\sigma^2}\right) \mathcal{B}(0, 1, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) - 2Y_{i\bullet} \mathcal{B}(1, 0, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) \\
&- \frac{\lambda_i^2}{\sigma^2} \mathcal{B}(0, 2, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) + \frac{1}{\sigma^2} \mathcal{B}(2, 0, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) \\
&- \mathcal{B}(0, 0, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i)^2 - \lambda_i \mathcal{B}(0, 1, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i)^2 \\
&- \frac{\lambda_i^2}{2} \mathcal{B}(0, 2, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i)^2 + \mathcal{C}(\beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i).
\end{aligned} \tag{6.14}$$

### 6.6.2 Expressing the Main Quantity in the Onsager's Correction Term Using Integral Families

The aim of this final section in this appendix is to obtain an expression for the quantity  $\xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2)$  that does not involve any expectation operators. Rather, we use specially tailored families of integrals for the desired expressions.

We first define the integral family definitions that we will work with. For  $p, q \in \{0, 1, 2\}, r > 0$  and  $s, t \in \mathbb{R}$ , define

$$\mathcal{Q}(p, q, r, s, t) \equiv (2\pi)^{-1/2} \int_{-\infty}^{\infty} (s + rx)^p b^{(q)}(t + rx) \exp\left(-\frac{x^2}{2}\right) dx. \quad (6.15)$$

Also, for  $r > 0$  and  $s, t \in \mathbb{R}$ , define

$$\mathcal{R}(r, s, t) \equiv (2\pi)^{-1/2} \int_{-\infty}^{\infty} b(s + rx)b(t + rx) \exp\left(-\frac{x^2}{2}\right) dx. \quad (6.16)$$

In the Poisson special case, since  $b(x) = \exp(x)$ ,  $\mathcal{Q}(p, q, r, s, t)$  and  $\mathcal{R}(r, s, t)$  admit exact expressions. However, for general  $b$  functions, we are stuck with expressions as in (6.15) and (6.16). Thus, numerical integration is required for evaluation purposes. Then note that in (6.14),

$$\begin{aligned} & \mathcal{B}(p, q, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) \\ &= \sum_{j=1}^n E\left(\tilde{U}_i^p b^{(q)}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i\right) \\ &= \sum_{j=1}^n \int_{-\infty}^{\infty} u^p b^{(q)}(\beta_0 + \beta_1 X_{ij} + u) (2\pi\lambda_i)^{-1/2} \exp\left\{-\frac{(u - \mu_i)^2}{2\lambda_i}\right\} du \\ &= \sum_{j=1}^n (2\pi)^{-1/2} \int_{-\infty}^{\infty} (\mu_i + \sqrt{\lambda_i}z)^p b^{(q)}(\beta_0 + \beta_1 X_{ij} + \mu_i + \sqrt{\lambda_i}z) \exp\left(-\frac{z^2}{2}\right) dz \\ &= \sum_{j=1}^n \mathcal{Q}(p, q, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i). \end{aligned}$$

Also note that

$$\begin{aligned} & \mathcal{C}(\beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) \\ &= \sum_{j=1}^n E\left(b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i)b(\beta_0 + \beta_1 X_{ij'} + \tilde{U}_i) | \mathbf{X}_i\right) \\ &= \sum_{j=1}^n \sum_{j'=1}^n \int_{-\infty}^{\infty} b(\beta_0 + \beta_1 X_{ij} + u)b(\beta_0 + \beta_1 X_{ij'} + u) (2\pi\lambda_i)^{-1/2} \exp\left\{-\frac{(u - \mu_i)^2}{2\lambda_i}\right\} du \\ &= \sum_{j=1}^n \sum_{j'=1}^n (2\pi)^{-1/2} \int_{-\infty}^{\infty} b(\beta_0 + \beta_1 X_{ij} + \mu_i + \sqrt{\lambda_i}z)b(\beta_0 + \beta_1 X_{ij'} + \mu_i + \sqrt{\lambda_i}z) \\ & \quad \times \exp\left(-\frac{z^2}{2}\right) dz \\ &= \sum_{j=1}^n \sum_{j'=1}^n \mathcal{R}(\sqrt{\lambda_i}, \beta_0 + \beta_1 X_{ij} + \mu_i, \beta_0 + \beta_1 X_{ij'} + \mu_i). \end{aligned}$$

Hence, the final expression for  $\xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2)$  used in Result 2 is as follows

$$\begin{aligned}
& \xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2) \\
&= \sum_{j=1}^n \left\{ \left( 2Y_{i\bullet} \mu_i - \frac{\lambda_i + \mu_i^2}{\sigma^2} \right) \mathcal{Q}(0, 0, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i) \right. \\
&\quad + 2\lambda_i \left( Y_{i\bullet} - \frac{\mu_i}{\sigma^2} \right) \mathcal{Q}(0, 1, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i) \\
&\quad - 2Y_{i\bullet} \mathcal{Q}(1, 0, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i) \\
&\quad - \frac{\lambda_i^2}{\sigma^2} \mathcal{Q}(0, 2, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i) + \frac{1}{\sigma^2} \mathcal{Q}(2, 0, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i) \\
&\quad - \mathcal{Q}(0, 0, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i)^2 - \lambda_i \mathcal{Q}(0, 1, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i)^2 \\
&\quad - \frac{\lambda_i^2}{2} \mathcal{Q}(0, 2, \sqrt{\lambda_i}, \mu_i, \beta_0 + \beta_1 X_{ij} + \mu_i)^2 \\
&\quad \left. + \sum_{j'=1}^n \mathcal{R}(\sqrt{\lambda_i}, \beta_0 + \beta_1 X_{ij} + \mu_i, \beta_0 + \beta_1 X_{ij'} + \mu_i) \right\}.
\end{aligned}$$

### 6.6.3 Proof of Result 3

This appendix contains the details leading to Result 3. In the Poisson case,  $b(x) = \exp(x)$ . Hence for  $p, q \in \{0, 1, 2\}$ , each term in (6.14) can be simplified greatly, leading to a reduced expression for  $\xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2)$ .

#### 6.6.3.1 Simplifications in Poisson Case

When  $p = 0, q \in \{0, 1, 2\}$ ,

$$\begin{aligned}
\mathcal{B}(0, q, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) &= \sum_{j=1}^n E \left( b^{(q)}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i \right) \\
&= \sum_{j=1}^n E \left\{ \exp(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i \right\} \\
&= E \left\{ \exp(\tilde{U}_i) \right\} \sum_{j=1}^n \exp(\beta_0 + \beta_1 X_{ij}) \\
&= \exp(\mu_i + \frac{1}{2} \lambda_i) \sum_{j=1}^n \exp(\beta_0 + \beta_1 X_{ij}).
\end{aligned} \tag{6.17}$$

When  $p = 1, q = 0$ ,

$$\begin{aligned}
\mathcal{B}(1, 0, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) &= \sum_{j=1}^n E \left( \tilde{U}_i b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i \right) \\
&= \sum_{j=1}^n E \left\{ \tilde{U}_i \exp(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i \right\} \\
&= E \left\{ \tilde{U}_i \exp(\tilde{U}_i) \right\} \sum_{j=1}^n \exp(\beta_0 + \beta_1 X_{ij}) \\
&= (\mu_i + \lambda_i) \exp(\mu_i + \frac{1}{2} \lambda_i) \sum_{j=1}^n \exp(\beta_0 + \beta_1 X_{ij}).
\end{aligned} \tag{6.18}$$

When  $p = 2, q = 0$ ,

$$\begin{aligned}
\mathcal{B}(2, 0, \beta_0, \beta_1, \mu_i, \lambda_i, \mathbf{X}_i) &= \sum_{j=1}^n E \left( \tilde{U}_i^2 b(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i \right) \\
&= \sum_{j=1}^n E \left\{ \tilde{U}_i^2 \exp(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) | \mathbf{X}_i \right\} \\
&= E \left\{ \tilde{U}_i^2 \exp(\tilde{U}_i) \right\} \sum_{j=1}^n \exp(\beta_0 + \beta_1 X_{ij}) \\
&= \{ \lambda_i + (\mu_i + \lambda_i)^2 \} \exp(\mu_i + \frac{1}{2} \lambda_i) \sum_{j=1}^n \exp(\beta_0 + \beta_1 X_{ij}).
\end{aligned} \tag{6.19}$$

Substituting (6.17), (6.18) and (6.19) into (6.14), several terms cancel with each other and we have

$$\begin{aligned}
&\xi_i(\mu_i, \lambda_i; \mathbf{Y}_i, \mathbf{X}_i, \beta_0, \beta_1, \sigma^2) \\
&= \left\{ \exp(\lambda_i) - 1 - \lambda_i - \frac{1}{2} \lambda_i^2 \right\} \exp(2\mu_i + \lambda_i) \left\{ \sum_{j=1}^n \exp(\beta_0 + \beta_1 X_{ij}) \right\}^2,
\end{aligned} \tag{6.20}$$

leading to Result 3.

## Chapter 7

# Extensions to Noncanonical Link Generalized Linear Mixed Models

In some situations, for a better data fit, using a noncanonical link might be beneficial over using a canonical link. To cater for these situations, we present an extension of the asymptotic normality results derived in Chapter 4 for generalized linear mixed models with noncanonical links. Thus, this chapter presents the final asymptotic normality theorem in this thesis, that concerns the joint asymptotic normality of all of the maximum quasi-likelihood estimators for a generalized linear mixed model with noncanonical links. As in Chapter 4, it elegantly shows faster rates of convergence for fixed effects that are not accompanied by a random effect compared to fixed effects that have a partnering random effect.

Finally, to wrap up the thesis, we build on the theory presented in Chapter 6 and present some details concerning the Thouless-Palmer-Anderson enhancement approach for improving statistical inference for generalized linear mixed models when noncanonical links are involved.

This chapter starts off by providing asymptotic normality results concerning generalized linear mixed models with noncanonical links in Section 7.1. Section 7.1.1 presents the model being used while Section 7.1.2 presents the notation required for the asymptotic normality theorem presented in Section 7.1.3. The chapter concludes with Section 7.2 which presents an introduction into the usage of the Thouless-Anderson-Palmer enhancement approach when noncanonical links are involved. Section 7.2.1 provides the model description for this section. Section 7.2.2 then provides an explicit expression for the GVA log-likelihood. Finally, Section 7.2.3 provide some introductory details

regarding the TAP enhancement approach for noncanonical links.

## 7.1 Asymptotic Normality Results Involving Noncanonical Links

In this section, we present asymptotic normality results for maximum quasi-likelihood estimators for a generalized linear mixed model with noncanonical links.

### 7.1.1 Model Description

To accommodate the use of noncanonical links, consider the following density, or probability mass, function for the class of one-parameter exponential families as in Fan et al. (1995) where

$$p(y; \eta) = \exp [y(g \circ b')^{-1}(\eta) - \{b \circ (g \circ b')^{-1}\}(\eta) + c(y)] h(y) \quad (7.1)$$

where  $g$  is the link function and  $\eta$  is the natural parameter. Here,  $I(\mathcal{P}) = 1$  if the condition  $\mathcal{P}$  is true and  $I(\mathcal{P}) = 0$  if  $\mathcal{P}$  is false. If the random variable  $Y$  has density, or probability mass, function as in (7.1), then  $E(Y) = g^{-1}(\eta)$  and  $\text{Var}(Y) = \{b'' \circ (b')^{-1} \circ g^{-1}\}(\eta)$ . To account for overdispersion in the data and to allow one to model the variance flexibly, a common modelling extension is implemented such that  $\text{Var}(Y) = \phi \{b'' \circ (b')^{-1} \circ g^{-1}\}(\eta)$ , where  $\phi > 0$  represents the dispersion parameter. This involves replacement of  $\log\{p(y; \eta)\}$  by the following quasi-likelihood function

$$[y(g \circ b')^{-1}(\eta) - \{b \circ (g \circ b')^{-1}\}(\eta) + c(y)] / \phi + d(y, \phi) \quad (7.2)$$

where  $d(y, \phi)$  is a function of  $y$  and  $\phi$  only. Note that for ordinary Binomial and Poisson response models,  $\phi$  is fixed at 1. For Gaussian and Gamma response models, (7.2) corresponds to the expression of  $\log\{p(y; \eta)\}$  for a two-parameter exponential family density function and ordinary likelihood applies. In this section, we study generalized linear mixed models of the following form, for observations of the random pairs  $(\mathbf{X}_{ij}, Y_{ij})$ ,  $1 \leq i \leq m, 1 \leq j \leq n_i$ ,

$$Y_{ij} | \mathbf{X}_{ij}, \mathbf{U}_i \text{ are independent having quasi-likelihood function (7.2) with natural parameter } \beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i \text{ such that the } U_i \text{ are independent } N(0, \sigma^2) \text{ random vectors.} \quad (7.3)$$

The  $U_i$  are unobserved random effects variables. We assumed that the  $\mathbf{X}_{ij}$  and  $U_i$ , for  $1 \leq i \leq m$  and  $1 \leq j \leq n_i$ , are totally independent, with the  $\mathbf{X}_{ij}$  each having the same distribution as the  $d \times 1$  random vector  $\mathbf{X}$  and the  $U_i$  each having the same distribution as the random variable  $U$ .

Then, for any  $\beta_0 \in \mathbb{R}$ ,  $\beta_1 \in \mathbb{R}^d$  and  $\sigma^2 > 0$ , the conditional log-quasi-likelihood is the joint density function of the  $Y_{ij}$ , given the  $\mathbf{X}_{ij}$ , as a function of the parameters  $(\beta_0^0, \beta_1^0, (\sigma^2)^0)$  is,

$$(\hat{\beta}_0, \hat{\beta}_1, \hat{\sigma}^2) = \operatorname{argmax}_{\beta_0, \beta_1, \sigma^2} \ell(\beta_0, \beta_1, \sigma^2)$$

where  $\ell(\beta_0, \beta_1, \sigma^2)$  is the conditional log-likelihood and has the expression

$$\begin{aligned} \ell(\beta_0, \beta_1, \sigma^2) = & -\frac{m}{2} \log(2\pi\sigma^2) + \sum_{i=1}^m \sum_{j=1}^{n_i} (c(Y_{ij})/\phi + d(Y_{ij}, \phi)) \\ & + \sum_{i=1}^m \log \int_{-\infty}^{\infty} \exp \left( \sum_{j=1}^{n_i} \left[ Y_{ij} (g \circ b')^{-1} (\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) \right. \right. \\ & \left. \left. - \{b \circ (g \circ b')^{-1}\} (\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) \right] / \phi - \frac{u^2}{2\sigma^2} \right) du. \end{aligned} \quad (7.4)$$

### 7.1.2 Notation

Define

$$n \equiv \frac{1}{m} \sum_{i=1}^m n_i = \text{average of the within-group sample sizes.}$$

Also define

$$\Omega_{\beta_1}^*(U) \equiv E \left\{ w(\beta_0 + \beta_1 + U) \begin{bmatrix} 1 & \mathbf{X}^T \\ \mathbf{X} & \mathbf{X}\mathbf{X}^T \end{bmatrix} \middle| U \right\},$$

where

$$w \equiv \frac{\{(g^{-1})'\}^2}{\{b'' \circ (b')^{-1} \circ g^{-1}\}}$$

and let

$$\Lambda_{\beta_1}^* \equiv \left( E \left[ \left\{ \text{lower right } d \times d \text{ block of } \Omega_{\beta_1}^*(U)^{-1} \right\}^{-1} \right] \right)^{-1}.$$



### 7.1.3 Asymptotic Normality Theorem

The main theoretical contribution of this chapter is an asymptotic normality theorem for the maximum quasi-likelihood estimators for a generalized linear mixed model with noncanonical links as described in Section 7.1.1.

The theorem relies on the following assumptions:

- (A8) The number of groups  $m$  diverges to  $\infty$ .
- (A9) The within-group sample sizes  $n_i$  diverge to  $\infty$  in such a way that  $n_i/n \rightarrow C_i$  for constants  $0 < C_i < \infty$ ,  $1 \leq i \leq m$ . Also,  $n/m \rightarrow 0$  as  $m$  and  $n$  diverge.
- (A10) The distribution of  $\mathbf{X}$  is such that

$$E \left[ \frac{E [\max\{1, \|\mathbf{X}\|\}^8 \max\{1, w(\beta_0 + \beta_1^T \mathbf{X} + U)\}^4 | U]}{\min\{1, \lambda_{\min}(E\{w(\beta_0 + \beta_1^T \mathbf{X} + U) | U\})\}^2} \right] < \infty$$

for all  $\beta_0 \in \mathbb{R}$ ,  $\beta_1 \in \mathbb{R}^d$  and  $\sigma^2 > 0$ .

**Theorem 14.** *Assume that conditions (A8) - (A10) hold. Then we have the following*

$$\sqrt{m} \begin{bmatrix} \hat{\beta}_0 - \beta_0^0 \\ \sqrt{n} (\hat{\beta}_1 - \beta_1^0) \\ \hat{\sigma}^2 - (\sigma^2)^0 \end{bmatrix} \xrightarrow{D} N \left( \begin{bmatrix} 0 \\ \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} (\sigma^2)^0 & \mathbf{0} & 0 \\ \mathbf{0} & \phi \mathbf{\Lambda}_{\beta_1}^* & \mathbf{0} \\ 0 & \mathbf{0} & 2 \{(\sigma^2)^0\}^2 \end{bmatrix} \right).$$

The proof of Theorem 14 is in the appendix. A remark concerning Theorem 14 is as follows:

1. When working with binary responses, there are two common alternative noncanonical link functions used to the canonical logit link, namely the probit link and the complementary log-log link. Both noncanonical links map the mean response restricted to the  $(0, 1)$  interval to the  $(-\infty, \infty)$  interval. Also in both cases,  $w(x)$  in Section 7.1.2 can be expressed as an explicit function to work with. Note that  $b(x) = \log\{1 + \exp(x)\}$  for binary responses. Then the usage of a probit link where  $g(x) = \Phi^{-1}(x)$  leads to  $w(x)$  taking on the following expression

$$w = \frac{\varphi^2}{\Phi(1 - \Phi)},$$

where  $\varphi$  represents the standard normal density function.

Meanwhile, the usage of a complementary log-log link where

$$g(x) = \ln\{-\ln(1-x)\}$$

results in

$$w(x) = \frac{\exp(2x)}{\exp\{\exp(x)\} - 1}.$$

## 7.2 Thouless-Anderson-Palmer Approach Involving Non-canonical Links

In this section, we present a primer for those who wish to carry out the Thouless-Anderson-Palmer approach for generalized linear mixed models with noncanonical links.

### 7.2.1 Model Description

Now, consider the use of a simple noncanonical link generalized linear mixed model as follows

$$\begin{aligned} Y_{ij}|X_{ij}, U_i \text{ are independent having density function (7.1) with} \\ \text{natural parameter } \beta_0^0 + \beta_1^0 X_{ij} + U_i \text{ such that the } U_i \text{ are independent} \\ N\left(0, (\sigma^2)^0\right) \text{ random variables.} \end{aligned} \quad (7.5)$$

Here, the values of  $(X_{ij}, Y_{ij})$  are observed for  $1 \leq i \leq m$  and  $1 \leq j \leq n$ . We have assumed that the  $X_{ij}$  and  $U_i$  are independent random variables. In addition, the  $X_{ij}$  are each assumed as having the same distribution as the random vector  $X$ . The  $U_i$  are the unobserved random effects variables and are assumed to be having the same distribution as the random vector  $U$ .

Let  $\boldsymbol{\beta} = (\beta_0, \beta_1)$  be the vector of fixed parameters. Then the model parameters for this set-up are  $(\boldsymbol{\beta}, \sigma^2)$ . Following that  $\ell(\boldsymbol{\beta}, \sigma^2)$ , the conditional log-likelihood of  $(\boldsymbol{\beta}, \sigma^2)$ ,

is

$$\begin{aligned}
 \ell(\boldsymbol{\beta}, \sigma^2) &= -\frac{m}{2} \log(2\pi\sigma^2) + \sum_{i=1}^m \sum_{j=1}^n c(Y_{ij}) \\
 &+ \sum_{i=1}^m \log \int_{-\infty}^{\infty} \exp \left( \sum_{j=1}^n \left[ Y_{ij} (g \circ b')^{-1}(\beta_0 + \beta_1 X_{ij} + u) \right. \right. \\
 &\left. \left. - \{b \circ (g \circ b')^{-1}\}(\beta_0 + \beta_1 X_{ij} + u) \right] - \frac{u^2}{2\sigma^2} \right) du.
 \end{aligned} \tag{7.6}$$

### 7.2.2 The Gaussian Variational Approximate Log-Likelihood

Following steps similar to those in Section 6.2, the Gaussian variational approximation to  $\ell(\boldsymbol{\beta}, \sigma^2)$  is derived as,

$$\begin{aligned}
 &\underline{\ell}_{\text{GVA}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) \\
 &= \sum_{i=1}^m E_{\tilde{U}_i} \left[ \sum_{j=1}^n \left\{ Y_{ij} (g \circ b')^{-1}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) - \{b \circ (g \circ b')^{-1}\}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \right. \right. \\
 &\quad \left. \left. + c(Y_{ij}) \right\} - \frac{\tilde{U}_i^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right] + \frac{1}{2} \sum_{i=1}^m \{1 + \log(2\pi\lambda_i)\},
 \end{aligned} \tag{7.7}$$

where  $E_{\tilde{U}_i}$  denotes the expectation with respect to the random variable  $\tilde{U}_i \sim N(\mu_i, \lambda_i)$  with  $\lambda_i > 0$ , for  $1 \leq i \leq m$ . The variational parameters are  $(\boldsymbol{\mu}, \boldsymbol{\lambda})$  where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$  and  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_m)$ . Note that

$$\ell(\boldsymbol{\beta}, \sigma^2) \geq \underline{\ell}_{\text{GVA}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda})$$

for all vectors  $\boldsymbol{\mu}$  and  $\boldsymbol{\lambda}$ .

### 7.2.3 Overview of Thouless-Anderson-Palmer Enhancement

Next, we provide details on how the TAP enhancement approach builds upon the GVA approach when using noncanonical links. Firstly, for each  $1 \leq i \leq m$ , define the following data vectors:

$$\mathbf{Y}_i \equiv (Y_{i1}, \dots, Y_{in}) \text{ and } \mathbf{X}_i \equiv (X_{i1}, \dots, X_{in}).$$

The Gaussian variational approximate negative log-likelihood can then be expressed as follows

$$-\underline{\ell}_{\text{GVA}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) = -\frac{1}{2} \sum_{i=1}^m \{1 + \log(2\pi\lambda_i)\} + \sum_{i=1}^m E_{\tilde{U}_i} \left\{ \Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i \right\},$$

where

$$\begin{aligned} \Psi_i(\tilde{U}_i) = \sum_{j=1}^n & \left[ -Y_{ij}(g \circ b')^{-1}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) + \{b \circ (g \circ b')^{-1}\}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) - c(Y_{ij}) \right] \\ & + \frac{\tilde{U}_i^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2). \end{aligned}$$

Using theory detailed in Section 6.3, the TAP approximate negative log-likelihood can be obtained as,

$$\begin{aligned} -\underline{\ell}_{\text{TAP}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) &= -\underline{\ell}_{\text{GVA}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\mu}, \boldsymbol{\lambda}) - \frac{1}{2} \sum_{i=1}^m \xi_i(\mu_i, \lambda_i; \mathbf{X}_i, \mathbf{Y}_i, \beta_0, \beta_1, \sigma^2) \\ &= -\frac{1}{2} \sum_{i=1}^m \{1 + \log(2\pi\lambda_i)\} + \sum_{i=1}^m E \left\{ \Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i \right\} \\ &\quad - \frac{1}{2} \sum_{i=1}^m \left( \text{Var}\{\Psi_i(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} - \lambda_i \left[ E\{\Psi_i'(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} \right]^2 \right. \\ &\quad \left. - \frac{\lambda_i^2}{2} \left[ E\{\Psi_i''(\tilde{U}_i) | \mathbf{Y}_i, \mathbf{X}_i\} \right]^2 \right). \end{aligned} \tag{7.8}$$

Now define

$$\begin{aligned} \mathcal{E}_i(\tilde{U}_i) &\equiv \sum_{j=1}^n Y_{ij}(g \circ b')^{-1}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i) \quad \text{and} \\ \mathcal{F}_i(\tilde{U}_i) &\equiv \sum_{j=1}^n \{b \circ (g \circ b')^{-1}\}(\beta_0 + \beta_1 X_{ij} + \tilde{U}_i). \end{aligned}$$

Then, we can re-express  $\Psi_i(\tilde{U}_i)$  as

$$\Psi_i(\tilde{U}_i) = -\mathcal{E}_i(\tilde{U}_i) + \mathcal{F}_i(\tilde{U}_i) - \sum_{j=1}^n c(Y_{ij}) + \frac{\tilde{U}_i^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2).$$

It follows that

$$\Psi_i'(\tilde{U}_i) = -\mathcal{E}_i'(\tilde{U}_i) + \mathcal{F}_i'(\tilde{U}_i) + \frac{\tilde{U}_i}{\sigma^2}$$

and

$$\Psi_i''(\tilde{U}_i) = -\mathcal{E}_i''(\tilde{U}_i) + \mathcal{F}_i''(\tilde{U}_i) + \frac{1}{\sigma^2}.$$

By solving for and the substituting explicit expressions for  $\Psi_i(\tilde{U}_i)$ ,  $\Psi_i'(\tilde{U}_i)$  and  $\Psi_i''(\tilde{U}_i)$

into (7.8), the TAP approximate negative log-likelihood can be obtained. This expression can then be optimized to find a local minima and obtain TAP estimates of the model parameters.

## 7.3 Appendix

This appendix contained the details for the derivations leading up to Theorem 14.

### 7.3.1 Constructing the Fisher Information Matrix

In order to compute the asymptotic covariance matrix for the maximum quasi-likelihood estimators, we would first need to compute the Fisher information matrix for the model parameters as per the model description in 7.1.1. To do so, let

$$\mathbf{S}_i \equiv \begin{bmatrix} \mathbf{S}_{0i} \\ \mathbf{S}_{1i} \\ \mathbf{S}_{2i} \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial \beta_0} \log p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i) \\ \nabla_{\beta_1} \log p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i) \\ \frac{\partial}{\partial \sigma^2} \log p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i) \end{bmatrix} \quad (7.9)$$

denote the  $i$ th contribution to the scores for each of the model parameters. Then the Fisher information matrix can be computed as

$$I(\beta_0, \beta_1, \sigma^2) = \sum_{i=1}^m E(\mathbf{S}_i \mathbf{S}_i^T | \mathbf{X}_i).$$

The next few sections then focus on obtaining the expressions for the scores and the quadratic conditional expectations that are required to construct the final Fisher information matrix.

### 7.3.2 Expression for Conditional Density Function

The expression for  $p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i)$  as per the model description in (7.3) is

$$\begin{aligned} & p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i) \\ &= \int_{-\infty}^{\infty} \prod_{j=1}^{n_i} \{p(Y_{ij} | \mathbf{X}_{ij}, U_i)\} p(U_i) dU_i \\ &= \int_{-\infty}^{\infty} \exp \left\{ \sum_{j=1}^n \left( [Y_{ij}(g \circ b')^{-1}(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) - \{b \circ (g \circ b')^{-1}\}(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) \right. \right. \\ & \quad \left. \left. + c(Y_{ij}) \right] / \phi + d(Y_{ij}, \phi) \right\} \times (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{u^2}{2\sigma^2}\right) du \end{aligned}$$

The previous expression can be further simplified as follows

$$\begin{aligned}
& \int_{-\infty}^{\infty} (2\pi\sigma^2)^{-1/2} \exp \left\{ \sum_{j=1}^n \left( \left[ Y_{ij}(g \circ b')^{-1} (\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) \right. \right. \right. \\
& \quad \left. \left. \left. - \{b \circ (g \circ b')^{-1}\} (\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) + c(Y_{ij}, \phi) \right] / \phi + d(Y_{ij}, \phi) \right) - \frac{u^2}{2\sigma^2} \right\} du \\
&= (2\pi\sigma^2)^{-1/2} \exp \left\{ \sum_{j=1}^n (c(Y_{ij}) / \phi + d(Y_{ij}, \phi)) \right\} \\
& \quad \times \int_{-\infty}^{\infty} \exp \left[ \sum_{j=1}^n \left\{ Y_{ij}(g \circ b')^{-1} (\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) - \{b \circ (g \circ b')^{-1}\} (\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) \right\} / \phi \right. \\
& \quad \left. - \frac{u^2}{2\sigma^2} \right] du.
\end{aligned}$$

### 7.3.3 Deriving Expressions for the Expectation and Variance of the Response Variable

Next, using the Bartlett identities, we can obtain expressions for  $E(Y)$  and  $\text{Var}(Y)$ . Firstly let,

$$a_1(\eta) \equiv (g \circ b')^{-1}(\eta) \quad \text{and} \quad a_2(\eta) \equiv \{b \circ (g \circ b')^{-1}\}(\eta).$$

Then note that the noncanonical extension of the one-parameter exponential family of density, or probability mass, functions takes on the following form

$$\begin{aligned}
p(y; \eta, \phi) &\propto \exp \left[ \frac{y(g \circ b')^{-1}(\eta) - \{b \circ (g \circ b')^{-1}\}(\eta) + c(y)}{\phi} \right] \\
&= \exp \left\{ \frac{a_1(\eta)y - a_2(\eta) + c(y)}{\phi} \right\}
\end{aligned}$$

where  $\eta$  is the natural parameter. Then, we have the following

$$\log p(y; \eta, \phi) = \frac{a_1(\eta)y - a_2(\eta) + c(y)}{\phi} + C$$

where  $C$  is a constant independent of  $\eta$ . Let  $\ell(\eta) \equiv \log p(y; \eta, \phi)$ . Then,

$$\frac{\partial \ell(\eta)}{\partial \eta} = \frac{a'_1(\eta)y - a'_2(\eta)}{\phi}. \quad (7.10)$$

The first Bartlett identity states that

$$E \left( \frac{\partial \ell(\eta)}{\partial \eta} \right) = 0.$$

Hence, by substituting (7.10) into the first Bartlett identity, we have,

$$E \left\{ \frac{a'_1(\eta)y - a'_2(\eta)}{\phi} \right\} = 0.$$

This leads to

$$\frac{a'_1(\eta)E(y) - a'_2(\eta)}{\phi} = 0.$$

Hence, now we have an initial expression for  $E(Y)$  which is as follows

$$E(Y) = \frac{a'_2(\eta)}{a'_1(\eta)}. \quad (7.11)$$

Now note that,

$$\frac{\partial^2 \ell}{\partial \eta^2} = \frac{a''_1(\eta)(y) - a''_2(\eta)}{\phi}. \quad (7.12)$$

The second Bartlett identity states that

$$E\left(\frac{\partial^2 \ell}{\partial \eta^2}\right) + E\left\{\left(\frac{\partial \ell}{\partial \eta}\right)^2\right\} = 0.$$

Using the first Bartlett identity, we can re-write the second Bartlett identity as follows

$$E\left(\frac{\partial^2 \ell}{\partial \eta^2}\right) + \text{Var}\left(\frac{\partial \ell}{\partial \eta}\right) = 0.$$

Substituting (7.10) and (7.12) into the second Bartlett identity, we have,

$$E\left\{\frac{a''_1(\eta)(y) - a''_2(\eta)}{\phi}\right\} + \text{Var}\left\{\frac{a'_1(\eta)y - a'_2(\eta)}{\phi}\right\} = 0. \quad (7.13)$$

This leads to

$$\frac{a''_1(\eta)E(y) - a''_2(\eta)}{\phi} + \frac{(a'_1(\eta))^2}{\phi^2} \text{Var}(y) = 0.$$

Therefore, the initial expression for  $\text{Var}(Y)$  is as follows

$$\begin{aligned} \text{Var}(Y) &= \left[ \frac{a''_2(\eta) - a''_1(\eta) \left\{ \frac{a'_2(\eta)}{a'_1(\eta)} \right\}}{\phi} \right] \left[ \frac{\phi^2}{\{a'_1(\eta)\}^2} \right] \\ &= \frac{\phi \{a'_1(\eta)a''_2(\eta) - a''_1(\eta)a'_2(\eta)\}}{\{a'_1(\eta)\}^3} \\ &= \frac{\phi \{a'_2(\eta)/a'_1(\eta)\}'}{a'_1(\eta)}. \end{aligned} \quad (7.14)$$

Now let us simplify the expressions for  $E(Y)$  and  $\text{Var}(Y)$ . Note that

$$\begin{aligned} a_2(\eta) &= \{b \circ (g \circ b')^{-1}\}(\eta) = (b \circ a_1)(\eta) \\ a'_2(\eta) &= (b \circ a_1)'(\eta) = (b' \circ a_1)(\eta) \times a'_1(\eta). \end{aligned}$$

Hence, the expression for  $E(Y)$  simplifies to

$$\begin{aligned} E(Y) &= \frac{a'_2(\eta)}{a'_1(\eta)} \\ &= (b' \circ a_1)(\eta) \\ &= \{b' \circ (g \circ b')^{-1}\}(\eta) \\ &= \{b' \circ (b')^{-1} \circ g^{-1}\}(\eta) \\ &= g^{-1}(\eta). \end{aligned}$$

The expression for  $\text{Var}(Y)$  simplifies to

$$\begin{aligned}\text{Var}(Y) &= \frac{\phi\{a'_2(\eta)/a'_1(\eta)\}'}{a'_1(\eta)} \\ &= \frac{\phi\{(b' \circ a_1)'(\eta)\}}{a'_1(\eta)} \\ &= \frac{\phi\{(b'' \circ a_1)(\eta)\}\{a'_1(\eta)\}}{a'_1(\eta)} \\ &= \phi\{b'' \circ (g \circ b')^{-1}\}(\eta) \\ &= \phi\{b'' \circ (b')^{-1} \circ g^{-1}\}(\eta).\end{aligned}$$

### 7.3.4 Introduction of Useful Notation and Its Properties

Let  $\mathbf{v}$  be a generic  $d \times 1$  vector. Then for  $r = 0, 1, 2$  we define

$$\mathbf{v}^{\otimes r} \equiv \begin{cases} 1 & \text{for } r = 0 \\ \mathbf{v} & \text{for } r = 1 \\ \mathbf{v}\mathbf{v}^T & \text{for } r = 2. \end{cases}$$

For  $r \in \{0, 1, 2\}$ , let

$$\begin{aligned}\mathcal{G}_{ri}^* &\equiv \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes r} \{Y_{ij}a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} \quad \text{and} \\ \mathcal{H}_{ri}^* &\equiv \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes r} \{a''_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - Y_{ij}a''_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\}.\end{aligned}$$

Note that the expressions listed above have the following probabilistic orders where

$$\mathcal{G}_{ri}^* = O_P(n^{1/2})\mathbf{1}_d^{\otimes r} \quad \text{and} \quad \mathcal{H}_{ri}^* = O_P(n)\mathbf{1}_d^{\otimes r}.$$

### 7.3.5 Key Conditional Moment Results

We now compute the conditional expectations of  $\mathcal{G}_{0i}^*$  and  $\mathcal{G}_{1i}^*$  given  $(\mathbf{X}_i, U_i)$ . Note that,

$$\begin{aligned}&E(\mathcal{G}_{0i}^* | \mathbf{X}_i, U_i) \\ &= E \left[ \sum_{j=1}^n \{Y_{ij}a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} \middle| \mathbf{X}_i, U_i \right] \\ &= \sum_{j=1}^n \{E(Y_{ij} | \mathbf{X}_i, U_i)a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} \\ &= \sum_{j=1}^n \left\{ \frac{a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)}{a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) \right\} \\ &= 0.\end{aligned}$$



Similarly, we have that

$$E(\mathcal{G}_{1i}^* | \mathbf{X}_i, U_i) = \mathbf{0}.$$

Next, we compute the conditional expectations of  $\mathcal{G}_{0i}^{*2}$ ,  $\mathcal{G}_{0i}^* \mathcal{G}_{1i}^*$  and  $\mathcal{G}_{1i}^{*2}$  given  $(\mathbf{X}_i, U_i)$ . Note that,

$$\begin{aligned} & E(\mathcal{G}_{0i}^{*2} | \mathbf{X}_i, U_i) \\ &= E \left( \left[ \sum_{j=1}^n \{Y_{ij} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} \right] \right. \\ &\quad \left. \times \left[ \sum_{j'=1}^n \{Y_{ij'} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij'} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij'} + U_i)\} \right] \middle| \mathbf{X}_i, U_i \right) \\ &= \sum_{j \neq j'} E \left[ \{Y_{ij} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} \right. \\ &\quad \left. \times \{Y_{ij'} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij'} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij'} + U_i)\} \middle| \mathbf{X}_i, U_i \right] \\ &\quad + \sum_{j=1}^n E \left[ \{Y_{ij} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} \right. \\ &\quad \left. \times \{Y_{ij} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} \middle| \mathbf{X}_i, U_i \right] \\ &= \sum_{j \neq j'} E \left( E \left[ \{Y_{ij} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} \middle| \mathbf{X}_i, U_i \right] \right. \\ &\quad \left. \times E \left[ \{Y_{ij'} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij'} + U_i) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij'} + U_i)\} \middle| \mathbf{X}_i, U_i \right] \right) \\ &\quad + \sum_{j=1}^n \text{Var} \left\{ Y_{ij} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) \middle| \mathbf{X}_i, U_i \right\} \\ &= \sum_{j=1}^n a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)^2 \text{Var} \left( Y_{ij} \middle| \mathbf{X}_i, U_i \right). \end{aligned}$$

Now let  $\eta_{ij} \equiv \beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i$ . By substituting the expression for  $\text{Var}(Y)$  in line 2 of (7.14), we have,

$$\begin{aligned} E(\mathcal{G}_{0i}^{*2} | \mathbf{X}_i, U_i) &= \phi \sum_{j=1}^n a'_1(\eta_{ij})^2 \left[ \frac{\{a'_1(\eta_{ij}) a''_2(\eta_{ij}) - a''_1(\eta_{ij}) a'_2(\eta_{ij})\}}{\{a'_1(\eta_{ij})\}^3} \right] \\ &= \phi \sum_{j=1}^n \left\{ a''_2(\eta_{ij}) - \frac{a''_1(\eta_{ij}) a'_2(\eta_{ij})}{a'_1(\eta_{ij})} \right\} \\ &= \phi \sum_{j=1}^n \{a''_2(\eta_{ij}) - a''_1(\eta_{ij}) E(Y_{ij} | \mathbf{X}_i, U_i)\} \\ &= \phi E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i). \end{aligned}$$

Similarly, we have that

$$E(\mathcal{G}_{0i}^* \mathcal{G}_{1i}^* | \mathbf{X}_i, U_i) = \phi E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i) \quad \text{and} \quad E(\mathcal{G}_{1i}^{*2} | \mathbf{X}_i, U_i) = \phi E(\mathcal{H}_{2i}^* | \mathbf{X}_i, U_i).$$

Now let us simplify the expression for  $E(\mathcal{H}_{ri}^* | \mathbf{X}_i, U_i)$  for  $r \in \{0, 1, 2\}$ . Note that,

$$\begin{aligned}
E(\mathcal{H}_{ri}^* | \mathbf{X}_i, U_i) &= \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes r} \{a_2''(\eta_{ij}) - a_1''(\eta_{ij})E(Y_{ij} | \mathbf{X}_i, U_i)\} \\
&= \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes r} \left\{ a_2''(\eta_{ij}) - \frac{a_1''(\eta_{ij})a_2'(\eta_{ij})}{a_1'(\eta_{ij})} \right\} \\
&= \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes r} a_1'(\eta_{ij})^2 \left[ \frac{\{a_1'(\eta_{ij})a_2''(\eta_{ij}) - a_1''(\eta_{ij})a_2'(\eta_{ij})\}}{\{a_1'(\eta_{ij})\}^3} \right] \\
&= \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes r} a_1'(\eta_{ij})^2 \{b'' \circ (g \circ b')^{-1}\}(\eta_{ij}).
\end{aligned}$$

The term  $a_1'(\eta_{ij})^2$  can be simplified as follows

$$\begin{aligned}
a_1'(\eta_{ij})^2 &= \left[ \{(g \circ b')^{-1}\}' \right]^2 (\eta_{ij}) \\
&= \left\{ \frac{1}{(g \circ b')' \circ (g \circ b')^{-1}} \right\}^2 (\eta_{ij}) \\
&= \left[ \frac{1}{\{(g \circ b')b''\} \circ \{(b')^{-1} \circ g^{-1}\}} \right]^2 (\eta_{ij}) \\
&= \left[ \frac{1}{(g' \circ g^{-1}) \{b'' \circ (b')^{-1} \circ g^{-1}\}} \right]^2 (\eta_{ij}) \\
&= \left[ \frac{(g^{-1})'}{\{b'' \circ (b')^{-1} \circ g^{-1}\}} \right]^2 (\eta_{ij}).
\end{aligned}$$

Hence,

$$\begin{aligned}
E(\mathcal{H}_{ri}^* | \mathbf{X}_i, U_i) &= \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes r} \left[ \frac{\{(g^{-1})'\}^2}{\{b'' \circ (b')^{-1} \circ g^{-1}\}} \right] (\eta_{ij}) \\
&= \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes r} w(\eta_{ij}),
\end{aligned} \tag{7.15}$$

where

$$w \equiv \frac{\{(g^{-1})'\}^2}{b'' \circ (b')^{-1} \circ g^{-1}}.$$

### 7.3.6 Computing an Asymptotic Approximation for the First Entry in (7.9)

Once again, to overcome the intractability of the ratio of integrals present when deriving the scores with respect to each of the model parameters, we will work with an asymptotic approximation of the ratio of integrals by using a multi-term Laplace's method expansion.

For smooth real-valued functions  $b$ ,  $g$  and  $h$ , Equation (2.6) of Tierney et al. (1989) states that

$$\frac{\int_{-\infty}^{\infty} b_N(x) \exp\{-nh(x)\} dx}{\int_{-\infty}^{\infty} b_D(x) \exp\{-nh(x)\} dx} = g(x^*) + \frac{b'_D(x^*)g'(x^*)}{nb_D(x^*)h''(x^*)} + \frac{g''(x^*)}{2nh''(x^*)} - \frac{g'(x^*)h'''(x^*)}{2nh''(x^*)^2} + O(n^{-2}).$$

where

$$g \equiv b_N/b_D$$

and

$$x^* \equiv \text{value of } x \text{ that minimises } h \text{ over } \mathbb{R}.$$

Hence, the  $i$ th contribution to the score of  $\beta_0$  can be expressed as

$$\begin{aligned} S_{0i} &\equiv \frac{\partial \log p_{\mathbf{Y}_i|\mathbf{X}_i}(\mathbf{Y}_i|\mathbf{X}_i)}{\partial \beta_0} \\ &= \frac{\int_{-\infty}^{\infty} b_N^{1st}(u) \exp\{-nh_N(u)\} du}{\int_{-\infty}^{\infty} b_D^{1st}(u) \exp\{-nh_N(u)\} du} \end{aligned}$$

where

$$\begin{aligned} b_N^{1st}(u) &= \exp\left(-\frac{u^2}{2\sigma^2}\right) \frac{1}{\phi} \sum_{j=1}^n \{Y_{ij}a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u)\}, \\ b_D^{1st}(u) &= \exp\left(-\frac{u^2}{2\sigma^2}\right) \quad \text{and} \\ h_N(u) &= -\frac{1}{n\phi} \sum_{j=1}^n \{Y_{ij}a_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) - a_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u)\}. \end{aligned}$$

Now define

$$\begin{aligned} U_i^* &\equiv \text{value of } u \text{ that minimises } h_N(u) \\ &= \text{value of } u \text{ such that } \frac{d}{du} h_N(u) = 0 \\ &= \text{value of } u \text{ such that } \sum_{j=1}^n \{Y_{ij}a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u)\} = 0. \end{aligned}$$

However,

$$b_N^{1st}(U_i^*) = 0.$$

This violates the necessary condition in Hsu (1948), in the sense that  $b_N^{1st}(U_i^*) \neq 0$  is required in order for the Laplace approximation to hold. To counter this issue, firstly note that the numerator of  $S_{0i}$  is

$$\int_{-\infty}^{\infty} s'(u)t(u)du$$

where

$$s(u) \equiv \exp\left(\frac{n}{\phi} \left[ \frac{1}{n} \sum_{j=1}^n \{Y_{ij}a_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) - a_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u)\} \right]\right)$$

and

$$t(u) \equiv \exp\left(-\frac{u^2}{2\sigma^2}\right).$$

Application of integration by parts leads to the integral being equal to

$$-\int_{-\infty}^{\infty} s(u)t'(u)du.$$

Note that

$$t'(u) = (-u/\sigma^2) \exp\left(-\frac{u^2}{2\sigma^2}\right).$$

Now by rewriting the numerator of  $S_{0i}$ , we have,

$$\begin{aligned} S_{0i} &\equiv \frac{\partial \log p_{\mathbf{Y}_i|\mathbf{X}_i}(\mathbf{Y}_i|\mathbf{X}_i)}{\partial \beta_0} \\ &= \frac{\int_{-\infty}^{\infty} b_N(u) \exp\{-nh_N(u)\} du}{\int_{-\infty}^{\infty} b_D(u) \exp\{-nh_N(u)\} du} \end{aligned}$$

where

$$\begin{aligned} b_N(u) &= (u/\sigma^2) \exp\left(-\frac{u^2}{2\sigma^2}\right), \\ b_D(u) &= \exp\left(-\frac{u^2}{2\sigma^2}\right) \quad \text{and} \\ h_N(u) &= -\frac{1}{n\phi} \sum_{j=1}^n \{Y_{ij}a_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) - a_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u)\}. \end{aligned}$$

### Expansion of $U_i^*$

Here we find an asymptotic expression for  $U_i^*$ . We have that

$$\begin{aligned} 0 &= \frac{d}{du} h_N(\mathbf{u}) \\ &= \sum_{j=1}^n \{Y_{ij}a_1'(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i^*) - a_2'(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i^*)\} \\ &= \sum_{j=1}^n \{Y_{ij}a_1'(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a_2'(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} \\ &\quad - (U_i^* - U_i) \sum_{j=1}^n \{a_2''(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - Y_{ij}a_1''(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} + r_{it} \\ &= \mathcal{G}_{0i}^* - (U_i^* - U_i)\mathcal{H}_{0i}^* + r_{it} \end{aligned}$$

where  $r_{it}$  is the Lagrange form of the remainder and is a quadratic form in  $U_i^* - U_i$  and a smooth function of  $U_{it}^\dagger \equiv (1-t)U_i + tU_i^*$  for some  $t \in [0, 1]$ . Inversion of this asymptotic series leads to

$$U_i^* = U_i + \frac{\mathcal{G}_{0i}^*}{\mathcal{H}_{0i}^*} + O_P(n^{-1}).$$

### 7.3.6.1 The First Term of the First Score

The first term of  $S_{0i}$  is

$$g(U_i^*) = \frac{b_N(U_i^*)}{b_D(U_i^*)} = \frac{U_i^*}{\sigma^2} = \frac{1}{\sigma^2} \left( U_i + \frac{\mathcal{G}_{0i}^*}{\mathcal{H}_{0i}^*} \right) + O_P(n^{-1}).$$

### 7.3.6.2 The Other Terms of the First Score

The second, third and fourth terms of  $S_{0i}$  are  $O_P(n^{-1})$ , 0 and  $O_P(n^{-1})$  as in the canonical case in Chapter 4.

### 7.3.6.3 Overall Leading Term Expression for the First Score

The first term of  $S_{0i}$  is  $O_P(1)$ , the second term of  $S_{0i}$  is  $O_P(n^{-1})$ , the third term of  $S_{0i}$  is 0 and the fourth term of  $S_{0i}$  is  $O_P(n^{-1})$ . Putting these together we obtain the following asymptotic expansion for  $S_{0i}$  such that

$$S_{0i} = \frac{1}{\sigma^2} \left( U_i + \frac{\mathcal{G}_{0i}^*}{\mathcal{H}_{0i}^*} \right) + O_P(n^{-1}).$$

### 7.3.7 Computing an Asymptotic Approximation for the Second Entry in (7.9)

The  $i$ th contribution to the score of  $\beta_1$  is

$$\begin{aligned} S_{1i} &\equiv \nabla_{\beta_1} \log p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i) \\ &= \frac{\int_{-\infty}^{\infty} b_N(u) \exp\{-nh_N(u)\} du}{\int_{-\infty}^{\infty} b_D(u) \exp\{-nh_N(u)\} du} \end{aligned}$$

where

$$b_N(u) = \exp\left(-\frac{u^2}{2\sigma^2}\right) \frac{1}{\phi} \sum_{j=1}^n \{Y_{ij} a'_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) - a'_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u)\} \mathbf{X}_{ij},$$

$$b_D(u) = \exp\left(-\frac{u^2}{2\sigma^2}\right) \text{ and}$$

$$h_N(u) = -\frac{1}{n\phi} \sum_{j=1}^n \{Y_{ij} a_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) - a_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u)\}.$$

### 7.3.7.1 The First Term of the Second Score

The first term of  $S_{1i}$  is

$$\begin{aligned} g(U_i^*) &= \frac{b_N(U_i^*)}{b_D(U_i^*)} \\ &= \frac{1}{\phi} \sum_{j=1}^n \{Y_{ij} a_1'(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i^*) - a_2'(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i^*)\} \mathbf{X}_{ij}. \end{aligned}$$

Next note that,

$$\begin{aligned} &Y_{ij} a_1'(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i^*) - a_2'(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i^*) \\ &= Y_{ij} a_1'(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a_2'(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) \\ &\quad + (U_i^* - U_i) \{Y_{ij} a_1''(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i) - a_2''(\beta_0 + \beta_1^T \mathbf{X}_{ij} + U_i)\} + O_P(n^{-1}). \end{aligned}$$

Therefore,

$$g(U_i^*) = \frac{1}{\phi} \left( \mathcal{G}_{1i}^* - \frac{\mathcal{G}_{0i}^* \mathcal{H}_{1i}^*}{\mathcal{H}_{0i}^*} \right) + O_P(1) \mathbf{1}_d.$$

### 7.3.7.2 The Other Terms of the Second Score

The second, third and fourth terms are  $O_P(1) \mathbf{1}_d$  as in the canonical case in Chapter 4.

### 7.3.7.3 Overall Leading Term Expression for the Second Score

Putting the terms of the score together, we obtain the following asymptotic expansion for  $S_{1i}$  where

$$S_{1i} = \frac{1}{\phi} \left( \mathcal{G}_{1i}^* - \frac{\mathcal{G}_{0i}^* \mathcal{H}_{1i}^*}{\mathcal{H}_{0i}^*} \right) + O_P(1) \mathbf{1}_d.$$

## 7.3.8 Computing an Asymptotic Approximation for the Third Entry in (7.9)

The  $i$ th contribution to the score of  $\sigma^2$  is

$$\begin{aligned} S_{2i} &\equiv \frac{\partial \log p_{\mathbf{Y}_i | \mathbf{X}_i}(\mathbf{Y}_i | \mathbf{X}_i)}{\partial \sigma^2} \\ &= -\frac{1}{2\sigma^2} + \frac{\int_{-\infty}^{\infty} b_N(u) \exp\{-nh_N(u)\} du}{\int_{-\infty}^{\infty} b_D(u) \exp\{-nh_N(u)\} du} \end{aligned}$$

where

$$\begin{aligned} b_N(u) &\equiv \frac{u^2}{2\sigma^4} \exp\left(-\frac{u^2}{2\sigma^2}\right), \\ b_D(u) &\equiv \exp\left(-\frac{u^2}{2\sigma^2}\right) \text{ and} \\ h_N(u) &\equiv -\frac{1}{n\phi} \sum_{j=1}^n \{Y_{ij}a_1(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u) - a_2(\beta_0 + \beta_1^T \mathbf{X}_{ij} + u)\}. \end{aligned}$$

### 7.3.8.1 The First Term of the Third Score

The first term of  $S_{2i}$  is

$$\begin{aligned} g(U_i^*) &= \frac{b_N(U_i^*)}{b_D(U_i^*)} = \frac{(U_i^*)^2}{2\sigma^4} \\ &= \frac{1}{2\sigma^4} \left( U_i + \frac{\mathcal{G}_{0i}^*}{\mathcal{H}_{0i}^*} + O_P(n^{-1}) \right)^2 \\ &= \frac{1}{2\sigma^4} \left( U_i^2 + \frac{2U_i\mathcal{G}_{0i}^*}{\mathcal{H}_{0i}^*} + O_P(n^{-1}) \right). \end{aligned}$$

### 7.3.8.2 The Other Terms of the Third Score

The second, third and fourth terms of  $S_{2i}$  are all  $O_P(n^{-1})$  as in the canonical case in Chapter 4.

### 7.3.8.3 Overall Leading Term Expression for the Third Score

The first term of  $S_{2i}$  is  $O_P(1)$ , the second term of  $S_{2i}$  is  $O_P(n^{-1})$ , the third term of  $S_{2i}$  is  $O_P(n^{-1})$  and the fourth term of  $S_{2i}$  is  $O_P(n^{-1})$ . Putting these together we obtain the following asymptotic expansion for  $S_{2i}$  such that

$$S_{2i} = -\frac{1}{2\sigma^2} + \frac{U_i^2}{2\sigma^4} + \frac{U_i\mathcal{G}_{0i}^*}{\sigma^4\mathcal{H}_{0i}^*} + O_P(n^{-1}).$$

## 7.3.9 The Quadratic Conditional Expectations of the Scores

In this section we find the conditional expectations required for the diagonal entries of the Fisher information matrix of  $(\beta_0, \beta_1, \sigma^2)$ .

### 7.3.9.1 The Conditional Expectation of the Square of the First Score

From the previous sections, we have the following approximation where

$$S_{0i} = \frac{U_i}{\sigma^2} + O_P(n^{-1/2}).$$

Therefore,

$$E(S_{0i}^2 | \mathbf{X}_i) = E \left\{ \left( \frac{U_i}{\sigma^2} + O_P(n^{-1/2}) \right)^2 \middle| \mathbf{X}_i \right\} = \frac{1}{\sigma^2} + O_P(n^{-1}).$$

### 7.3.9.2 The Conditional Expectation of the Square of the Second Score

We have the following approximation

$$\phi S_{1i} = \mathcal{G}_{1i}^* - \frac{\mathcal{G}_{0i}^* \mathcal{H}_{1i}^*}{\mathcal{H}_{0i}^*} + O_P(1) \mathbf{1}_d.$$

Therefore,

$$\begin{aligned} \phi^2 E(S_{1i} S_{1i}^T | \mathbf{X}_i) &= E \left\{ \left( \mathcal{G}_{1i}^* - \frac{\mathcal{G}_{0i}^* \mathcal{H}_{1i}^*}{\mathcal{H}_{0i}^*} + O_P(1) \mathbf{1}_d \right)^{\otimes 2} \middle| \mathbf{X}_i \right\} \\ &= E(\mathcal{G}_{1i}^{*\otimes 2} | \mathbf{X}_i) + E \left( \frac{\mathcal{G}_{0i}^{*2} \mathcal{H}_{1i}^{*\otimes 2}}{\mathcal{H}_{0i}^{*2}} \middle| \mathbf{X}_i \right) - E \left( \frac{\mathcal{G}_{0i}^* \mathcal{G}_{1i}^* \mathcal{H}_{1i}^{*T}}{\mathcal{H}_{0i}^*} \middle| \mathbf{X}_i \right) \\ &\quad - E \left( \frac{\mathcal{G}_{0i}^* \mathcal{H}_{1i}^* \mathcal{G}_{1i}^{*T}}{\mathcal{H}_{0i}^*} \middle| \mathbf{X}_i \right) + O_P(1) \mathbf{1}_d^{\otimes 2}. \end{aligned}$$

We will now solve for the conditional expectations occurring in  $\phi^2 E(S_{1i} S_{1i}^T | \mathbf{X}_i)$ . Firstly,

$$\begin{aligned} E(\mathcal{G}_{1i}^{*\otimes 2} | \mathbf{X}_i) &= E \left\{ E(\mathcal{G}_{1i}^{*\otimes 2} | \mathbf{X}_i, U_i) \middle| \mathbf{X}_i \right\} \\ &= E \left\{ \phi E(\mathcal{H}_{2i}^* | \mathbf{X}_i, U_i) \middle| \mathbf{X}_i \right\} \\ &= \phi E \left\{ E(\mathcal{H}_{2i}^* | \mathbf{X}_i, U_i) \middle| \mathbf{X}_i \right\}. \end{aligned}$$

Next,

$$\begin{aligned} &E \left( \frac{\mathcal{G}_{0i}^{*2} \mathcal{H}_{1i}^{*\otimes 2}}{\mathcal{H}_{0i}^{*2}} \middle| \mathbf{X}_i \right) \\ &= E \left\{ E \left( \frac{\mathcal{G}_{0i}^{*2} \mathcal{H}_{1i}^{*\otimes 2}}{\mathcal{H}_{0i}^{*2}} \middle| \mathbf{X}_i, U_i \right) \middle| \mathbf{X}_i \right\} \\ &= E \left( E \left[ \frac{\mathcal{G}_{0i}^{*2} \{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i) + \mathcal{H}_{1i}^* - E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^{\otimes 2}}{\{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i) + \mathcal{H}_{0i}^* - E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)\}^2} \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right). \end{aligned}$$



The previous expression can be further evaluated as follows

$$\begin{aligned}
& E \left( E \left[ \frac{\mathcal{G}_{0i}^{*2} \{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i) + \mathcal{H}_{1i}^* - E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^{\otimes 2}}{\{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)\}^2 \left\{1 - \frac{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i) - \mathcal{H}_{0i}^*}{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)}\right\}^2} \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) \\
&= E \left( E \left[ \frac{\mathcal{G}_{0i}^{*2} \{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i) + \mathcal{H}_{1i}^* - E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^{\otimes 2}}{\{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)\}^2} \right. \right. \\
&\quad \left. \left. \times \left\{1 - \frac{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i) - \mathcal{H}_{0i}^*}{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)}\right\}^2 \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) \\
&= E \left( E \left[ \frac{\mathcal{G}_{0i}^{*2} \{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^{\otimes 2}}{\{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)\}^2} \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) + \mathcal{R}_i \\
&= \phi E \left[ \frac{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i) \{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^{\otimes 2}}{\{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)\}^2} \middle| \mathbf{X}_i \right] + \mathcal{R}_i \\
&= \phi E \left[ \frac{\{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^{\otimes 2}}{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)} \middle| \mathbf{X}_i \right] + \mathcal{R}_i
\end{aligned}$$

where  $\mathcal{R}_i$  comprises of remainder terms that are of lower order as compared to the leading terms in the final steps of the derivation. Now note that,

$$\begin{aligned}
& E \left( \frac{\mathcal{G}_{0i}^* \mathcal{G}_{1i}^* \mathcal{H}_{1i}^{*T}}{\mathcal{H}_{0i}^*} \middle| \mathbf{X}_i \right) \\
&= E \left\{ E \left( \frac{\mathcal{G}_{0i}^* \mathcal{G}_{1i}^* \mathcal{H}_{1i}^{*T}}{\mathcal{H}_{0i}^*} \middle| \mathbf{X}_i, U_i \right) \middle| \mathbf{X}_i \right\} \\
&= E \left( E \left[ \frac{\mathcal{G}_{0i}^* \mathcal{G}_{1i}^* \{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i) + \mathcal{H}_{1i}^* - E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^T}{\{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i) + \mathcal{H}_{0i}^* - E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)\}} \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) \\
&= E \left( E \left[ \frac{\mathcal{G}_{0i}^* \mathcal{G}_{1i}^* \{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i) + \mathcal{H}_{1i}^* - E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^T}{\{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)\} \left\{1 - \frac{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i) - \mathcal{H}_{0i}^*}{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)}\right\}} \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) \\
&= E \left( E \left[ \frac{\mathcal{G}_{0i}^* \mathcal{G}_{1i}^* \{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i) + \mathcal{H}_{1i}^* - E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^T}{\{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)\}} \right. \right. \\
&\quad \left. \left. \times \left\{1 - \frac{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i) - \mathcal{H}_{0i}^*}{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)}\right\} \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) \\
&= E \left[ E \left\{ \frac{\mathcal{G}_{0i}^* \mathcal{G}_{1i}^* E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)^T}{\{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)\}} \middle| \mathbf{X}_i, U_i \right\} \middle| \mathbf{X}_i \right] + \mathcal{R}_i \\
&= \phi E \left[ \frac{\{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^{\otimes 2}}{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)} \middle| \mathbf{X}_i \right] + \mathcal{R}_i.
\end{aligned}$$

Similarly,

$$E \left( \frac{\mathcal{G}_{0i}^* \mathcal{H}_{1i}^* \mathcal{G}_{1i}^{*T}}{\mathcal{H}_{0i}^*} \middle| \mathbf{X}_i \right) = \phi E \left[ \frac{\{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^{\otimes 2}}{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)} \middle| \mathbf{X}_i \right] + \mathcal{R}_i,$$

where  $\mathcal{R}_i$  comprises of remainder terms that are of lower order as compared to the leading term in the expression. Combining the results so far in this subsection, we obtain

$$E(S_{1i}S_{1i}^T|\mathbf{X}_i) = \frac{1}{\phi} E \left[ E(\mathcal{H}_{2i}^*|\mathbf{X}_i, U_i) - \frac{\{E(\mathcal{H}_{1i}^*|\mathbf{X}_i, U_i)\}^{\otimes 2}}{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i)} \middle| \mathbf{X}_i \right] + \mathcal{R}.$$

### 7.3.9.3 The Conditional Expectation of the Square of the Third Score

We have the following approximation

$$S_{2i} = -\frac{1}{2\sigma^2} + \frac{U_i^2}{2\sigma^4} + \frac{U_i\mathcal{G}_{0i}^*}{\sigma^4\mathcal{H}_{0i}^*} + O_P(n^{-1}).$$

Therefore,

$$\begin{aligned} 4\sigma^8 E(S_{2i}^2|\mathbf{X}_i) &= E \left\{ \left( U_i^2 - \sigma^2 + \frac{2U_i\mathcal{G}_{0i}^*}{\mathcal{H}_{0i}^*} + O_P(n^{-1}) \right)^2 \middle| \mathbf{X}_i \right\} \\ &= E(U_i^4) + \sigma^4 + 4E \left( \frac{U_i^2\mathcal{G}_{0i}^{*2}}{\mathcal{H}_{0i}^{*2}} \middle| \mathbf{X}_i \right) - 2\sigma^2 E(U_i^2) \\ &\quad + 2E \left( \frac{U_i^3\mathcal{G}_{0i}^*}{\mathcal{H}_{0i}^*} \middle| \mathbf{X}_i \right) - 4\sigma^2 \left( \frac{U_i\mathcal{G}_{0i}^*}{\mathcal{H}_{0i}^*} \middle| \mathbf{X}_i \right) + O_P(n^{-1}). \end{aligned}$$

Now note that,

$$\begin{aligned} &E \left( \frac{U_i^2\mathcal{G}_{0i}^{*2}}{\mathcal{H}_{0i}^{*2}} \middle| \mathbf{X}_i \right) \\ &= E \left\{ E \left( \frac{U_i^2\mathcal{G}_{0i}^{*2}}{\mathcal{H}_{0i}^{*2}} \middle| \mathbf{X}_i, U_i \right) \middle| \mathbf{X}_i \right\} \\ &= E \left( E \left[ \frac{U_i^2\mathcal{G}_{0i}^{*2}}{\{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i) + \mathcal{H}_{0i}^* - E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i)\}^2} \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) \\ &= E \left( E \left[ \frac{U_i^2\mathcal{G}_{0i}^{*2}}{\{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i)\}^2 \left\{ 1 - \frac{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i) - \mathcal{H}_{0i}^*}{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i)} \right\}^2} \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) \\ &= E \left( E \left[ \left\{ \frac{U_i^2\mathcal{G}_{0i}^{*2}}{\{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i)\}^2} \right\} \left\{ 1 - \frac{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i) - \mathcal{H}_{0i}^*}{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i)} \right\}^2 \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) \\ &= E \left( E \left[ \frac{U_i^2\mathcal{G}_{0i}^{*2}}{\{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i)\}^2} \middle| \mathbf{X}_i, U_i \right] \middle| \mathbf{X}_i \right) + \mathcal{R}_i \\ &= \phi E \left[ \frac{U_i^2 E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i)}{\{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i)\}^2} \middle| \mathbf{X}_i \right] + \mathcal{R}_i \\ &= \phi E \left[ \frac{U_i^2}{E(\mathcal{H}_{0i}^*|\mathbf{X}_i, U_i)} \middle| \mathbf{X}_i \right] + \mathcal{R}_i \\ &= O_P(n^{-1}). \end{aligned}$$

Also, for  $k = 1, 3$ ,

$$E\left(\frac{U_i^k \mathcal{G}_{0i}^*}{\mathcal{H}_{0i}^*} \middle| \mathbf{X}_i\right) = 0.$$

Hence,

$$\begin{aligned} 4\sigma^8 E(S_{2i}^2 | \mathbf{X}_i) &= E(U_i^4) + \sigma^4 - 2\sigma^2 E(U_i^2) + O_P(n^{-1}) \\ &= 3\sigma^4 + \sigma^4 - 2\sigma^4 + O_P(n^{-1}) \\ &= 2\sigma^4 + O_P(n^{-1}). \end{aligned}$$

Therefore we have,

$$E(S_{2i}^2 | \mathbf{X}_i) = \frac{1}{2\sigma^4} + O_P(n^{-1}).$$

### 7.3.10 The Fisher Information Matrix

The conditional expectations for the off-diagonal entries of the Fisher information matrix of  $(\beta_0, \beta_1, \sigma^2)$  expressed in terms of order notation are the same as in the canonical case. Putting together the expressions for the quadratic conditional expectations of the scores from the earlier three subsections, we have

$$I(\beta_0, \beta_1, \sigma^2) = \begin{bmatrix} \frac{m}{\sigma^2} + O_P(mn^{-1}) & O_P(m)\mathbf{1}_d^T & O_P(mn^{-1}) \\ O_P(m)\mathbf{1}_d & \frac{m}{\phi} E\left[E(\mathcal{H}_{2i}^* | \mathbf{X}_i, U_i) - \frac{\{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^{\otimes 2}}{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)} \middle| \mathbf{X}_i\right] + \mathcal{R} & O_P(m)\mathbf{1}_d \\ O_P(mn^{-1}) & O_P(m)\mathbf{1}_d^T & \frac{m}{2\sigma^4} + O_P(mn^{-1}) \end{bmatrix}.$$

However note that by using (7.15), we can rewrite the expression below as follows,

$$\begin{aligned} E(S_{1i} S_{1i}^T | \mathbf{X}_i) &= \frac{1}{\phi} E\left[E(\mathcal{H}_{2i}^* | \mathbf{X}_i, U_i) - \frac{\{E(\mathcal{H}_{1i}^* | \mathbf{X}_i, U_i)\}^{\otimes 2}}{E(\mathcal{H}_{0i}^* | \mathbf{X}_i, U_i)} \middle| \mathbf{X}_i\right] + \mathcal{R} \\ &= \frac{1}{\phi} E\left[\sum_{j=1}^n \mathbf{X}_{ij}^{\otimes 2} w(\eta_{ij}) - \frac{\left\{\sum_{j=1}^n \mathbf{X}_{ij} w(\eta_{ij})\right\}^{\otimes 2}}{\sum_{j=1}^n w(\eta_{ij})} \middle| \mathbf{X}_i\right] + \mathcal{R}. \end{aligned}$$

It follows that the leading term in (2, 2) block of  $I(\beta_0, \beta_1, \sigma^2)$  is,

$$\frac{mn}{\phi} \left( \frac{1}{mn} \sum_{i=1}^m E\left[\sum_{j=1}^n \mathbf{X}_{ij}^{\otimes 2} w(\eta_{ij}) - \frac{\left\{\sum_{j=1}^n \mathbf{X}_{ij} w(\eta_{ij})\right\}^{\otimes 2}}{\sum_{j=1}^n w(\eta_{ij})} \middle| \mathbf{X}_i\right] \right) = \frac{mn \Sigma_{\beta_1}}{\phi}$$

where

$$\begin{aligned}
\Sigma_{\beta_1} &\equiv \frac{1}{mn} \sum_{i=1}^m E \left[ \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes 2} w(\eta_{ij}) - \frac{\left\{ \sum_{j=1}^n \mathbf{X}_{ij} w(\eta_{ij}) \right\}^{\otimes 2}}{\sum_{j=1}^n w(\eta_{ij})} \middle| \mathbf{X}_i \right] \\
&= \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes 2} E\{w(\eta_{ij}) | \mathbf{X}_i\} \\
&\quad - E \left[ \frac{1}{mn} \sum_{i=1}^m \left\{ \sum_{j=1}^n \mathbf{X}_{ij} w(\eta_{ij}) \right\} \left\{ \sum_{j=1}^n w(\eta_{ij}) \right\}^{-1} \left\{ \sum_{j=1}^n \mathbf{X}_{ij} w(\eta_{ij}) \right\}^T \right].
\end{aligned} \tag{7.16}$$

Using Lemma 1 from Chapter 2 with  $f(\mathbf{X}_{ij}, \mathbf{U}_i) = w(\eta_{ij})$ , we have that the first term in (7.16) can be re-expressed as follows

$$\frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{ij}^{\otimes 2} E\{w(\eta_{ij}) | \mathbf{X}_i\} = E \{ \mathbf{X} \mathbf{X}^T w(\beta_0 + \beta_1 + U) \} + o_P(1) \mathbf{1}_d^{\otimes 2}.$$

Now, using Lemma 2 from Chapter 2 with  $f(\mathbf{X}_{ij}, \mathbf{U}_i) = w(\eta_{ij})$ , we have that the second term in (7.16) can be re-expressed as follows

$$\begin{aligned}
&E \left[ \frac{1}{mn} \sum_{i=1}^m \left\{ \sum_{j=1}^n \mathbf{X}_{ij} w(\eta_{ij}) \right\} \left\{ \sum_{j=1}^n w(\eta_{ij}) \right\}^{-1} \left\{ \sum_{j=1}^n \mathbf{X}_{ij} w(\eta_{ij}) \right\}^T \right] \\
&= E \left( E \{ \mathbf{X} w(\beta_0 + \beta_1 + U) | U \} [E \{ w(\beta_0 + \beta_1 + U) | U \}]^{-1} E \{ \mathbf{X} w(\beta_0 + \beta_1 + U) | U \}^T \right) \\
&\quad + o_P(1) \mathbf{1}_d^{\otimes 2}.
\end{aligned}$$

Now,

$$\Sigma_{\beta_1} = E \left[ E \{ \mathbf{X} \mathbf{X}^T w(\beta_0 + \beta_1 + U) | U \} - \frac{E \{ \mathbf{X} w(\beta_0 + \beta_1 + U) | U \}^{\otimes 2}}{E \{ w(\beta_0 + \beta_1 + U) | U \}} \right] + o_P(1) \mathbf{1}_d^{\otimes 2}.$$

Let

$$\Omega_{\beta_1}^*(U) \equiv E \left\{ w(\beta_0 + \beta_1 + U) \begin{bmatrix} 1 & \mathbf{X}^T \\ \mathbf{X} & \mathbf{X} \mathbf{X}^T \end{bmatrix} \middle| U \right\}.$$

The inverse of the lower right  $d \times d$  block of  $\Omega_{\beta_1}(U)^{-1}$  is

$$E \{ \mathbf{X} \mathbf{X}^T w(\beta_0 + \beta_1 + U) | U \} - \frac{E \{ \mathbf{X} w(\beta_0 + \beta_1 + U) | U \}^{\otimes 2}}{E \{ w(\beta_0 + \beta_1 + U) | U \}}.$$

Therefore, we conclude that that the (2, 2) block of  $I(\beta_0, \beta_1, \sigma^2)$  is as follows

$$\frac{mn(\Lambda_{\beta_1}^*)^{-1}}{\phi} + o_P(mn) \mathbf{1}_d^{\otimes 2}$$

where

$$\Lambda_{\beta_1}^* \equiv \left( E \left[ \left\{ \text{lower right } d \times d \text{ block of } \Omega_{\beta_1}^*(U)^{-1} \right\}^{-1} \right] \right)^{-1}.$$

Hence, the final Fisher information matrix can be expressed as

$$I(\beta_0, \beta_1, \sigma^2) = \begin{bmatrix} \frac{m}{\sigma^2} + O_P(mn^{-1}) & O_P(m)\mathbf{1}_d^T & O_P(mn^{-1}) \\ O_P(m)\mathbf{1}_d & \frac{mn(\Lambda_{\beta_1}^*)^{-1}}{\phi} + o_P(mn)\mathbf{1}_d^{\otimes 2} & O_P(m)\mathbf{1}_d \\ O_P(mn^{-1}) & O_P(m)\mathbf{1}_d^T & \frac{m}{2\sigma^4} + O_P(mn^{-1}) \end{bmatrix}.$$

### 7.3.11 The Inverse of the Fisher Information Matrix

To invert the Fisher information matrix, we choose to work with the  $(\beta_0, \sigma^2, \beta_1)$  ordering instead of  $(\beta_0, \beta_1, \sigma^2)$ . A trivial rearrangement of the matrix entries leads to

$$I(\beta_0, \beta_1, \sigma^2) = \begin{bmatrix} \frac{m}{\sigma^2} + O_P(mn^{-1}) & O_P(mn^{-1}) & O_P(m)\mathbf{1}_d^T \\ O_P(mn^{-1}) & \frac{m}{2\sigma^4} + O_P(mn^{-1}) & O_P(m)\mathbf{1}_d^T \\ O_P(m)\mathbf{1}_d & O_P(m)\mathbf{1}_d & \frac{mn(\Lambda_{\beta_1}^*)^{-1}}{\phi} + o_P(mn)\mathbf{1}_d^{\otimes 2} \end{bmatrix}.$$

Following similar steps to that in Subsubsection 4.5.2.10, we obtain the following expression for  $I(\beta_0, \sigma^2, \beta_1)^{-1}$  where

$$I(\beta_0, \sigma^2, \beta_1)^{-1} = \begin{bmatrix} \frac{\sigma^2}{m} + O_P(m^{-1}n^{-1}) & O_P(m^{-1}n^{-1}) & O_P(m^{-1}n^{-1})\mathbf{1}_d^T \\ O_P(m^{-1}n^{-1}) & \frac{2\sigma^4}{m} + O_P(m^{-1}n^{-1}) & O_P(m^{-1}n^{-1})\mathbf{1}_d^T \\ O_P(m^{-1}n^{-1})\mathbf{1}_d & O_P(m^{-1}n^{-1})\mathbf{1}_d & \frac{\phi\Lambda_{\beta_1}^*}{mn} + o_P(m^{-1}n^{-1})\mathbf{1}_d^{\otimes 2} \end{bmatrix}.$$

The expression for the inverse of the Fisher information matrix can be also written as follows

$$I(\beta_0, \sigma^2, \beta_1)^{-1} = I(\beta_0, \sigma^2, \beta_1)^{-1} + \frac{1}{mn} \begin{bmatrix} O_P(1) & O_P(1) & O_P(1)\mathbf{1}_d^T \\ O_P(1) & O_P(1) & O_P(1)\mathbf{1}_d^T \\ O_P(1)\mathbf{1}_d & O_P(1)\mathbf{1}_d & o_P(1)\mathbf{1}_d^{\otimes 2} \end{bmatrix}.$$

where

$$I(\beta_0, \sigma^2, \beta_1)^{-1} = \begin{bmatrix} \frac{\sigma^2}{m} & 0 & \mathbf{0} \\ 0 & \frac{2\sigma^4}{m} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \frac{\phi\Lambda_{\beta_1}^*}{mn} \end{bmatrix}.$$

### 7.3.12 Final Asymptotic Normality Result

The derivations leading to the final asymptotic normality result follows steps similar to those in Subsubsection 4.5.2.11j. It follows that

$$\sqrt{m} \begin{bmatrix} \widehat{\beta}_0 - \beta_0^0 \\ \sqrt{n}(\widehat{\beta}_1 - \beta_1^0) \\ \widehat{\sigma}^2 - (\sigma^2) \end{bmatrix} \xrightarrow{\mathcal{D}} N \left( \begin{bmatrix} 0 \\ \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} (\sigma^2)^0 & \mathbf{0} & 0 \\ \mathbf{0} & \phi\Lambda_{\beta_1}^* & \mathbf{0} \\ 0 & \mathbf{0} & 2\{(\sigma^2)^0\}^2 \end{bmatrix} \right).$$

## Chapter 8

# Discussion and Conclusion

The aim of this thesis was to address some of the current gaps present in the literature available for GLMMs. Specifically, we aimed to develop asymptotic theory for maximum likelihood and maximum quasi-likelihood estimators for GLMMs. Once these results were derived, we wished to assess the efficacy of the studentized confidence results constructed based upon these asymptotic normality results and also explore the implications of such results on optimal design theory. Lastly, with regards to methodology for GLMMs, we aimed to implement the Thouless-Anderson-Palmer approach and analyse the accuracy of the variational estimates obtained.

In order to carry out detailed asymptotic analysis on maximum likelihood estimators and maximum quasi-likelihood estimators for GLMMs, we first required results for two important tasks. The first of which required deriving population limits of particular predictor-dependent sample mean quantities. The second task involved establishing matrix norm asymptotic negligibility between matrix square roots of inverse Fisher information matrices and their simpler asymptotic block diagonal forms. Currently, there are no results available to deal with either of these tasks in a simple manner. Hence in Chapter 2, we provided the necessary tools and results in the forms of Lemmas 1, 2 and 3.

Next, in Chapter 3, we derived asymptotic normality results for GLMMs involving Gaussian responses. Existing theory for Gaussian linear mixed models by Wand (2002) and Harville (1977) provided base expressions for the required Fisher information matrix. Thereafter, the leading terms in the entries of the Fisher information matrix were retained, as the number of groups and the number of observations within each group diverged. (This approach was repeated throughout Chapters 4 and 7 in this thesis

as well). The resulting theorem as a result of our work concerns the joint asymptotic normality of all maximum likelihood estimators for a Gaussian response mixed model and elegantly shows faster rates of convergence of fixed effects unaccompanied by random effects as compared to fixed effects that have partnering random effects.

The results in Chapter 3 were extended to the class of all GLMMs, including a model extension for overdispersion, in Chapter 4. However, frequentist inference for GLMMs is hindered by the existence of intractable integrals due to the inclusion of random effects in these models. To overcome this obstacle, we used a multi-term Laplace's method expansion for ratios of intractable integrals (Miyata, 2004; Tierney et al., 1989). The resulting asymptotic normality theorem concerns the joint asymptotic normality of all of the maximum quasi-likelihood estimators, for fixed values of the dispersion parameter, for a generalized linear mixed model. Once again, the results derived in this chapter show faster rates of convergence for fixed effects that are not accompanied by random effects as compared to fixed effects accompanied by random effects. We also noted that for the class of two-parameter exponential families, maximum likelihood estimation is possible for all model parameters including the dispersion parameter. Based on this, we also derived asymptotic normality results for the maximum likelihood estimator for the dispersion parameter in the Gaussian and Gamma response cases.

Chapter 5 presents the consequences and applications of the asymptotic normality results derived in Chapter 4. First, we present how studentized confidence intervals can be constructed based on our asymptotic normality results in order to carry out asymptotically valid inference. The efficacy of the confidence intervals were then assessed, which showed that the Theorem 12-based approach in Section 4.3 as an attractive alternative to the exact observed Fisher information approach. The Theorem 12-based approach required simpler or no numerical integration at all compared to the exact observed Fisher information approach and gave similar coverage properties, especially for larger values of the number of groups and number of observations within each group. Next, we looked into the implications of Theorem 12 on optimal design theory. Most optimality criteria are based on the Fisher information matrix, which is computationally expensive to evaluate. Hence, we presented a simple approach to constructing approximate locally D-optimal designs based on large sample approximations of the Fisher information matrix.

In Chapter 6, we tackle the implementation of the Thouless-Anderson-Palmer approach for GLMMs. First, we presented a general result for deriving the Thouless-Anderson-Palmer approximate negative log-likelihood for GLMMs. This expression can then be locally minimized to obtain TAP estimates of the true model parameters. Then,

we carried out several simulation studies using a simple Poisson linear mixed model and analysed the Thouless-Anderson-Palmer variational estimates obtained against the estimates obtained from implementing Gaussian variational approximation. Based on the simulation studies, the Thouless-Anderson-Palmer enhancement approach suggests a slight yet statistically significant improvement as compared to using the Gaussian variational approximation approach, especially for small datasets.

Last but not least, in Chapter 7, we developed theory to consider the usage of noncanonical links for both the development of asymptotic normality results for maximum quasi-likelihood estimators for GLMMs and also the implementation of the Thouless-Anderson-Palmer approach for GLMMs.

In conclusion, this thesis presents important theoretical and methodological work that concerns the asymptotic distributions of maximum quasi-likelihood estimators for GLMMs and the implementation of the Thouless-Anderson-Palmer variational method for GLMMs. We believe that the work in this thesis will make a significant and novel contribution to the area of GLMMs, which have been a mainstay of regression-type statistical analyses in important areas such as longitudinal data analysis, multilevel modelling, panel data analysis and small area estimation.

Potential future work involves deriving second-order asymptotic approximations of the Fisher information matrix. In this thesis, we only retained the leading terms in the entries of the Fisher information matrix, hence leading to a first-order asymptotic approximation. Deriving second-order asymptotic approximations can give more accurate expressions for the asymptotic variance-covariance matrix for the maximum quasi-likelihood estimators for GLMMs, which is especially useful when implementing such results for smaller finite samples. For example, with second-order approximations, better coverages can be achieved when constructing studentized confidence intervals for smaller values for the number of groups and number of observations within each group as compared to those in the simulation study in Subsection 5.1.2. Deriving second-order approximations will also allow us to determine approximate locally D-optimal designs when considering multivariate random effects, which is not met by the theory presented in this thesis.

The techniques used in this thesis could also be used to study statistical models other than GLMMs. An example of a potential class of such models would be generalized additive mixed models. This class of models extends the GLMM framework by allowing for continuous predictors impacting the mean response to be modelled by additive non-parametric functions. This provides additional flexibility for modelling the actual



relationship between the response (specified by an exponential family distribution) and predictors, It also potentially provides better fits to the data as well. The techniques used in this thesis could be used to derive the first-order asymptotic approximations of the ratios of intractable integrals that arise when calculating the scores required for the Fisher information matrix, and subsequently the asymptotic variance-covariance matrix.

Last but not least, we can explore how the TAP approach performs against the GVA approach for GLMMs with response distributions other than the Poisson family.

## References

- Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015). Fitting Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software, Articles*, 67:1–48.
- Breslow, N. E. and Clayton, D. G. (1993). Approximate Inference in Generalized Linear Mixed Models. *Journal of the American Statistical Association*, 88:9–25.
- Chipman, J. S. (1964). On Least Squares with Insufficient Observations. *Journal of the American Statistical Association*, 59(308):1078–1111.
- Fan, J., Heckman, N. E., and Wand, M. P. (1995). Local Polynomial Kernel Regression for Generalized Linear Models and Quasi-Likelihood Functions. *Journal of the American Statistical Association*, 90:141–150.
- Fan, Z., Mei, S., and Montanari, A. (2021). TAP Free Energy, Spin Glasses and Variational Inference. *The Annals of Probability*, 49:1 – 45.
- Golub, G. and Van Loan, C. (2013). *Matrix Computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press.
- Hall, P., Pham, T., Wand, M. P., and Wang, S. S. J. (2011). Asymptotic Normality and Valid Inference for Gaussian Variational Approximation. *Annals of Statistics*, 39:2502–2532.
- Harville, D. A. (1977). Maximum Likelihood Approaches to Variance Component Estimation and to Related Problems. *Journal of the American Statistical Association*, 72:320–338.
- Higham, N. J. (2008). *Functions of Matrices: Theory and Computation*. Society for Industrial and Applied Mathematics.
- Hsu, L. C. (1948). Approximations to a Class of Double Integrals of Functions of Large Numbers. *American Journal of Mathematics*, 70:698–708.
- Jiang, J. (1996). REML Estimation: Asymptotic Behavior and Related Topics. *The Annals of Statistics*, 24:255–286.

- Jiang, J. and Nguyen, T. (2021). *Linear and Generalized Linear Mixed Models and Their Applications, Second Edition*. New York:Springer.
- Johnstone, I. (2022). Expectation Propagation in Mixed Models. Presentation in the conference “Statistics in the Big Data Era”. Simons Institute, California, U.S.A, June 2022.
- Knight, K. (2000). *Mathematical Statistics*. Chapman and Hall/CRC.
- Lyu, Z. and Welsh, A. (2022). Increasing Cluster Size Asymptotics for Nested Error Regression Models. *Journal of Statistical Planning and Inference*, 217:52–68.
- Maestrini, L., Bhaskaran, A., and Wand, M. P. (2023). Second term improvement to generalised linear mixed model asymptotics.
- Magnus, J. and Neudecker, H. (1999). *Matrix Differential Calculus with Applications in Statistics and Econometrics*. John Wiley & Sons.
- Magnus, J. R. and Neudecker, H. (1979). The Commutation Matrix: Some Properties and Applications. *The Annals of Statistics*, 7:381 – 394.
- McCulloch, C., Searle, S., and Neuhaus, J. (2008). *Generalized, Linear, and Mixed Models, Second Edition*. New York: John Wiley & Sons.
- Miller, J. J. (1973). Asymptotic properties of maximum likelihood estimates in the mixed model of the analysis of variance. *Technical Report No. 12, Department of Statistics*.
- Miller, J. J. (1977). Asymptotic properties of maximum likelihood estimates in the mixed model of the analysis of variance. *The Annals of Statistics*, 5:746–762.
- Miyata, Y. (2004). Fully Exponential Laplace Approximations Using Asymptotic Modes. *Journal of the American Statistical Association*, 99:1037–1049.
- Nie, L. (2007). Convergence Rate of MLE in Generalized Linear and Nonlinear Mixed-Effects Models: Theory and Applications. *Journal of Statistical Planning and Inference*, 137:1787–1804.
- Ormerod, J. T. and Wand, M. P. (2010). Explaining Variational Approximations. *The American Statistician*, 64:140–153.
- Pace, L. and Salvan, A. (1997). *Principles of Statistical Inference from a Neo-Fisherian Perspective*. World Scientific.
- Plefka, T. (1982). Convergence Condition of the TAP Equation for the Infinite-Ranged Ising Spin Glass Model. *Journal of Physics A: Mathematical and General*, 15:1971–1978.
- R Core Team (2022). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Russell, K. (2018). *Design of Experiments for Generalized Linear Models*. Chapman and Hall/CRC.

- Sherrington, D. and Kirkpatrick, S. (1975). Solvable model of a spin-glass. *Phys. Rev. Lett.*, 35:1792–1796.
- Thouless, D. J., Anderson, P. W., and Palmer, R. G. (1977). Solution of 'Solvable Model of a Spin Glass'. *The Philosophical Magazine: A Journal of Theoretical Experimental and Applied Physics*, 35:593–601.
- Tierney, L., Kass, R. E., and Kadane, J. B. (1989). Fully Exponential Laplace Approximations to Expectations and Variances of Nonpositive Functions. *Journal of the American Statistical Association*, 84:710–716.
- van der Vaart, A. W. (1998). *Asymptotic Statistics*. Cambridge, U.K.: Cambridge University Press.
- Waite, T. W. and Woods, D. C. (2015). Designs for Generalized Linear Models with Random Block Effects Via Information Matrix Approximations. *Biometrika*, 102:677–693.
- Wand, M. (2007). Fisher Information for Generalised Linear Mixed Models. *Journal of Multivariate Analysis*, 98:1412–1416.
- Wand, M. P. (2002). Vector Differential Calculus in Statistics. *The American Statistician*, 56:55–62.
- Westfall, P. H. (1986). Asymptotic Normality of the Anova Estimates of Components of Variance in the Nonnormal, Unbalanced Hierarchical Mixed Model. *The Annals of Statistics*, 14:1572 – 1582.
- Wolfram Research Inc. (2022). *Mathematica, Version 12.0*. Champaign, Illinois, U.S.A.
- Zhang, W., Mandal, A., and Stufken, J. (2017). Approximations of the Information Matrix for a Panel Mixed Logit Model. *Journal of Statistical Theory and Practice*, 11:269–295.