

# **Semiparametric and Nonparametric Density Deconvolution for Data with Measurement Er- ror; Applications to Nutrition Data**

Xiao Chen Yu

Supervisors:Ray Carroll, James Brown and Stephen Woodcock

PhD Thesis Mathematics  
University of Technology Sydney

Faculty of Science

Year of Submission:2021

CERTIFICATE OF ORIGINAL AUTHORSHIP

I, Xiao Chen Yu declare that this thesis, is submitted in fulfillment of the requirements for the award of PhD thesis: 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 the Australian Government Research Training Program.

Signature:

Production Note:

Signature removed prior to publication.

Date: Mar 2021

# Contents

<b>I</b>	<b>Abstract</b>	<b>i</b>
<b>II</b>	<b>Chapter 1: Introduction</b>	<b>iii</b>
<b>III</b>	<b>Chapter 2: Literature review</b>	<b>v</b>
<b>1</b>	<b>Methods for analysing measurement error models</b>	<b>vi</b>
1.1	Bayesian . . . . .	ix
1.2	Regression Calibration . . . . .	xii
1.3	Simulation Extrapolation (SIMEX) . . . . .	xiii
<b>2</b>	<b>Methods for estimating the density analysing only the error model</b>	<b>xvii</b>
2.1	Kernel Density Deconvolution . . . . .	xviii
2.2	Sieve Maximum Likelihood . . . . .	xxii
<b>3</b>	<b>Nutritional data</b>	<b>xxiv</b>
<b>IV</b>	<b>Chapter 3: Semiparametric Density Deconvolution for Continuous Data without Additional Information</b>	<b>xxvi</b>
<b>1</b>	<b>Introduction</b>	<b>xxvi</b>
<b>2</b>	<b>Hermite Polynomials</b>	<b>xxviii</b>
2.1	How to express an unknown density using Hermite polynomial . . . . .	xxx
<b>3</b>	<b>Methodology</b>	<b>xxxv</b>
<b>4</b>	<b>Simulations</b>	<b>xxxvi</b>
4.1	HePD: choosing the number of smoothing parameters . . . . .	xxxvi
4.2	Generating observed variable $W$ for simulation . . . . .	xxxvii
4.3	Assuming known parameters for $U$ . . . . .	xxxviii
4.4	Estimating both the density of $T$ and the parameters of $U$ . . . . .	lii
<b>5</b>	<b>Discussion</b>	<b>lvi</b>
<b>6</b>	<b>Computation details and errors</b>	<b>lvi</b>
<b>V</b>	<b>Chapter 4: Semiparametric and Nonparametric Density Deconvolution for Continuous Data with Replicates</b>	<b>lviii</b>
<b>1</b>	<b>Introduction</b>	<b>lviii</b>
<b>2</b>	<b>methodology</b>	<b>lix</b>

<b>3</b>	<b>Simulations</b>	<b>lxi</b>
3.1	Estimate $T$ when we know both the distribution type and standard deviation of $U$ . . . . .	lxi
3.2	Estimating the density of $T$ when we only know the distribution type for $U$ .	lxxvi
3.3	Estimating the density of $T$ with no information on the density of $U$ . . . . .	lxxi
<b>4</b>	<b>Comparisons</b>	<b>lxxi</b>
4.1	Knowing $\sigma_U$ vs estimating $\sigma_U$ . . . . .	lxxvi
4.2	No replicates vs multiple replicates . . . . .	lxxix
4.3	MAE and MSE results . . . . .	lxxix
<b>5</b>	<b>Real Data</b>	<b>lxxxii</b>
5.1	Quick Introduction to EATS . . . . .	lxxxii
<b>VI</b>	<b>Chapter 5: Semiparametric Density Deconvolution for Data with Excess Zeros</b>	<b>lxxxvii</b>
<b>1</b>	<b>Methodology</b>	<b>lxxxviii</b>
1.1	Using Zero Inflated Data . . . . .	lxxxviii
1.1.1	Usual Intake . . . . .	xc
1.1.2	BoxCox transformation parameter selection . . . . .	xciii
1.2	Bootstrapping . . . . .	xcv
<b>2</b>	<b>Simulations</b>	<b>xcv</b>
2.1	Simulation 1 results . . . . .	xcv
2.2	Simulation 2 results . . . . .	ciii
<b>VII</b>	<b>Chapter 6: EATS Data application</b>	<b>cviii</b>
<b>1</b>	<b>Introduction</b>	<b>cviii</b>
1.1	Energy . . . . .	cx
1.2	Alcohol . . . . .	cx
1.3	Total Fruits . . . . .	cxiii
1.4	Total Vegetables . . . . .	cxvi
<b>2</b>	<b>Bootstrapping</b>	<b>cxviii</b>
<b>VIII</b>	<b>Future plans</b>	<b>cxxi</b>

## Part I

# Abstract

Our inspiration behind this thesis is nutritional data, more specifically nutritional data collected through short term methods such as the 24HR recall. These collection methods obtain results that are quite accurate in what a subject consumed in a day, but is not an accurate representation of what a subject's consumption pattern looks like in long term. This leads to many statistician using measurement error models to adjust for the difference.

As our society as a whole becomes more aware of our health and how our eating pattern may effect it, more and more studies have come to focus on such ideas. And more recently, studies have come to focus on just understanding what the distribution of a populations consumption pattern looks like, in hopes to answer questions such as how does our society as a general consume a nutrition of interest, are we over or under consuming a certain food group or nutrition, has our consumption pattern changed as time passes, and so on.

So far in studies that use measurement error models to help obtain a density curve that represents a populations consumption patterns, most studies require additional information or additional assumptions that are given without specifying a reason such as assuming a certain distribution for the error terms of the model. For our thesis, we wish to develop a method that can obtain an unbiased distribution of a populations long term consumption pattern without additional information and minimal assumptions.

In this thesis, we start with a simple classical error model that will work well for continuous data, this may be good with nutrition data such as protein, fat and fiber. We then move on to allowing replicates in our observed variable, in doing so, we can let go of most assumptions on the error term. We do this because most 24HR recalls collect multiple entries from the same subject, which can work as replicates. We then move on to a more complex error model that is designed for zero-inflated data. We are interested in such a model since data collection methods such as the 24HR recall also collects information on what food we eat in a day, since it is very rare that we will eat every type of food in a 24 hour period, the 24HR recall will contain a large amount of zero. We hope to develop a method that

can help with estimating a populations long term consumption pattern using data collected using this short term method that contains excess zero.