# Modelling, Data Mining and Visualisation of Genetic Variation Data

A Thesis Submitted for the Degree of
Doctor of Philosophy

By

*Ahmad A. Aloqaily*

in

FACULTY OF ENGINEERING AND INFORMATION TECHNOLOGY
UNIVERSITY OF TECHNOLOGY, SYDNEY
AUSTRALIA
2012

# CERTIFICATE

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

<div align="center">

_____

Signature of Author

</div>

# Acknowledgements

my beloved wife Malak, who has been with me for nearly 5 years already and sacrificed a lot for helping me to pursue my academic pathway. Thank you for your patience and care, and for accompanying me to go through highs and lows accentuated by my research studies.

*To My Family,*

*My Wife and Our Angel, Tala.*

# Table of Contents

# List of Tables

# List of Figures

# Abstract

Data mining and knowledge discovery have been applied to datasets in various industries including biomedical informatics. The major challenges in data mining in the area stem from the fact that biomedical data comes in many forms with a highly dimensional nature. This research thesis focuses on one specific biomedical dataset, termed as genetic variation data in the form of genome-wide single nucleotide polymorphisms (SNPs) datasets.

Advances in single nucleotide polymorphism genotyping technologies have revolutionised our ability to explore the genetic architecture and models underlying complex diseases by conducting studies based on the whole genome. These studies are called genome-wide association studies. The basic strategy used in these studies is to examine the relationship between the disease of interest and genetic markers across the whole genome.

Many association studies have led to the discovery of single genetic variants associated with common diseases. However, complex diseases are not caused by single genes acting alone but are the result of complex non-linear interactions among genetic factors, with each gene having a small effect on disease risk. For this reason there is a critical need to implement new approaches that can take into account non-linear gene-gene interactions in searching for markers that jointly cause complex diseases.

Several computational methods have been developed to deal with the genetic complexity of complex diseases. However, testing each SNP for main effects and different

orders of gene-gene interaction is computationally infeasible for such high-dimensional data. Also, these methods do not scale well. Therefore, there is growing interest in applying non-parametric predictive models including data mining and machine learning approaches to understand genetic variation data.

This thesis constructs models which incorporate genetic variation data in a manner that will alleviate the error induced by the high dimensionality of such data. Data mining approaches, specifically non-parametric ones, are developed for the modelling, exploration and visualization of patient-to-patient relationships based on genome-wide SNP data. This thesis focuses on three main issues in genetic variation studies: (1) feature selection and distance calculations, (2) framework for the task of disease diagnosis and prognosis, and (3) models for the comparison and visualisation of patient-to-patient relationships based on genome-wide SNP profiles.

This thesis proposes efficient feature selection approaches to find an optimal subset of markers with the highest predictive power for the disease of interest, while managing the large search space required. The proposed approaches select genetic markers for marginal effects as well as gene-gene interaction effects. Markers with marginal effects are selected with an iterative random forest (RF) based procedure, called RF-RFE. The importance measure generated by random forest was chosen for estimating the importance of each SNP (weighting) and facilitates the selection of an appropriate set of SNPs. To deal with the large search space involved in detecting gene-gene interactions, putative markers are prioritized in the search using a new measure, called Interaction Effect (IE), that quantifies the potential for a SNP to be involved in gene-gene interaction. This measure can also be used as a splitting criterion in random forest construction to define a cut-off value of a ranked list of SNPs. The prioritized SNP set is used to construct new combined features, which carry the information to account for gene-gene interactions.

This thesis proposes three new methods for calculating distances between genotype

profiles based on a kernel-based weighting function including: RFK, using the RF variable importance measure; MAFK, based on the minor allele frequency measure and EK, using the entropy measure. The distances can be subsequently incorporated for the purpose of disease classifications, cluster analyses and visualizations.

The feasibility of using genetic variation data for disease diagnosis and prognosis is explored with a new computational framework. The framework demonstrates the use of different phases of data processing and modelling to build reliable disease diagnostic and prognostic models using genetic variation data. The proposed feature selection approaches are incorporated in the framework to select an optimal subset of SNPs with the highest predictive power.

The proposed framework is empirically evaluated using two case studies of acute lymphoblastic leukaemia. The results demonstrate that the framework can produce highly accurate diagnosis and prognosis models. This thesis shows that a significant improvement of models' performance requires including interaction markers. The results are consistent with known biology while the accuracy of the produced models is also high.

Finally, several data reduction methods are used to visualize genetic variation data. For unsupervised-based visualization, they are compared based on the trust-worthiness metric. For the supervised-based visualization, the performance is compared based on class discrimination. This thesis finds that the Neighbour Retrieval Visualizer method shows the best results for unsupervised-based visualization. Furthermore, in the supervised-based approach, the results highlight the importance of using feature selection to remove insignificant features. The visualization has the potential to assist clinicians and biomedical researchers in understanding relationships between patients and has the potential to lead to delivery of advanced personalized medicine.

The methodologies and approaches presented in this thesis emphasise the critical role that genetic variation data plays in understanding complex disease. The availability of a flexible framework for the task of disease diagnosis and prognosis, as proposed in this thesis, will play an important role in understanding the genetic basis to common complex diseases. A comprehensive validation of the methods and approaches embedded in the framework is a matter of applying this framework to other complex diseases.

# Publications

**Al-Oqaily, A.**, Tafavogh, S., Catchpoole, D., and Kennedy, P. 'A new computational framework for the task of disease diagnosis and prognosis', In preparation for submission to BMC Bioinformatics journal.

**Al-Oqaily, A.,** Kennedy, P., Catchpoole, D., and S. Simoff, S., (2008). 'Comparison of visualization methods of genome-wide SNP profiles in childhood acute lymphoblastic leukemia', *Data Mining and Analytics 2008: proceedings of the Seventh Australasian Data Mining Conference (AusDM'08)*, **87**, pp. 111-121, Australian Computer Society.

Kennedy, P., Simoff, S., Catchpoole, D., Skillicorn, D., Ubaudi, F. and **Al-Oqaily, A.** (2008). 'Integrative visual data mining of biomedical data: Investigating cases in chronic fatigue syndrome and acute lymphoblastic leukaemia'. In S. J. Simoff, M. H. Boehlen, and A. Mazeika, editors, *Visual Data Mining: Theory, Techniques and Tools for Visual Analytics*. Springer-Verlag New York Inc, pp. 367-388.

The work in chapters 3 and 4 was turned into a project plan for a successful grant application: 'Implementing personalized medicine using global genomic similarity', Cancer Institute NSW Research Innovation Grant 2011, 10/RFG/2-23, Daniel Catchpoole, Paul Kennedy, $50,000.