# Efficient Methods of Feature Selection Based on Combinatorial Optimization Motivated by the Analysis of Large Biological Datasets

**Mateus Rocha de Paula**

M.Sci. (Computer Science)

B.Sci. (Computer Science)

This dissertation is submitted as a partial requirement for the Degree of **Doctor of Philosophy**

THE UNIVERSITY OF

**NEWCASTLE**

AUSTRALIA

Faculty of Engineering and Built Environment

School of Electrical Engineering and Computer Science

Newcastle, NSW, Australia

August, 2012

# Statement of Originality

The thesis contains no material which has been accepted for the award of any other degree or diploma in any university or other tertiary institution and, to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference has been made in the text. I give consent to this copy of my thesis, when deposited in the University Library[1], being made available for loan and photocopying subject to the provisions of the Copyright Act 1968.

Mateus Rocha de Paula

---

[1]Unless an Embargo has been approved for a determined period.

# Acknowledgements

Matsuo Basho, one of the most celebrated Japanese poets, once wrote *"Every day is a journey, and the journey itself is home"*. When I think about my PhD today, that couldn't sound more true. The degree of Doctor of Philosophy not only represents the completion of this work, but also celebrates a series of meaningful events that allowed me to grow both as a person and professional; and marks the beginning of a new and exciting stage of my life. Every little step I took since I ventured this path has its own story, of immense importance, and none of which would be possible without the help great personalities that I was already fortunate to know, and others that I had the great pleasure to meet on my way. To all those, my many big thanks.

To my mum France, dad José Vicente and sister Rafaella, my eternal gratitude for all the love and support. Every time the light in the end of the tunnel seemed a little distant, it was you who made it brighter. This victory is very much yours too.

Many thanks to my supervisors Pablo and Regina, who so patiently mentored me through all the process, teaching me the skills and giving me the tools that I needed to succeed on this challenge and many yet to come. You not only helped me to be a better researcher, but also taught me how to be truly committed to science, and for that I'll always be thankful. To my friend and colleague Martín for supporting me from the very beginning of my academic career and for opening the doors to this opportunity, to Australia and of his own home to me, making it all possible. It's always been a pleasure to work with such a good friend, always with an open mind, positive attitude and a smile on his face! A big thanks to my friend Alexandre, for all the help both inside and outside the lab. This story would have had a much different ending without your advice, friendship, and many beers and fun times shared. To

*To Vó Cota, Tio Tonico and Shilly.*

*Wherever you are.*

# Contents

# List of Figures

# List of Tables

xiii

# List of Algorithms

# Abstract

Intuitively, the Feature Selection problem is to choose a subset of a given a set of features that best represents the whole in a particular aspect, preserving the original semantics of the variables on the given samples and classes. In practice, the objective of finding such a subset is often to reveal a particular characteristic present in the given samples.

In 2004, a new feature selection approach was proposed. It was based on a combinatorial optimization problem called $(\alpha, \beta)$-k-Feature Set Problem. The main advantage of using this approach over ranking methods is that the features are evaluated as groups, instead of only considering their individual performance.

The main drawback of this approach is the complexity of the combinatorial problems involved. Since some of them are NP-Complete, it is unlikely that there would exist an efficient method to solve them to optimality efficiently. To the best of the author's knowledge at the moment of this research, the available tools to deal with the $(\alpha, \beta)$-k-Feature Set Problem approach can not solve problems of the magnitude required by many practical applications.

Given the big advantage brought by the multivariate characteristic of this method, its successful wide applicability and knowing that its only real known drawback is scalability, further research to overcome such a difficulty is appropriate. Even though the optimal solution of the problem is always desirable, it often is not strictly necessary in the case of many biological applications. Therefore, this work aims to propose fast heuristics to address the $(\alpha, \beta)$-k-Feature Set Problem approach, and propose procedures to obtain dual bounds that do not rely on external optimization packages.