# JCU ePrints

This file is part of the following reference:

**Hancock, Timothy Peter (2006)** *Multivariate consensus trees: tree-based clustering and profiling for mixed data types.* **PhD thesis, James Cook University.**

Access to this file is available from:

http://eprints.jcu.edu.au/17497

JCU
JAMES COOK UNIVERSITY

# Multivariate Consensus Trees:

## Tree-based clustering and profiling for mixed data types

Thesis submitted by

Timothy Peter Hancock BSc(Hons)

in 2006

# STATEMENT ON SOURCES

## *Declaration*

I declare that this thesis is my own work and has not been submitted in any form for another degree or diploma at any university or other institution of tertiary education. Information derived from the published or unpublished work of others has been acknowledged in the text and a list of references is given.

……………………………………….. ………………………

(Signature) (Date)

# STATEMENT OF ACCESS

I, the undersigned, author of this work, understand that James Cook University will make this thesis available for use within the University Library and, via the Australian Digital Theses network, for use elsewhere.  I understand that, as an unpublished work, a thesis has significant protection under the Copyright Act and;

I do not wish to place any further restriction on access to this work.

_____                                    _____
Signature                                                                                    Date

# ELECTRONIC COPY

I, the undersigned, the author of this work, declare that the electronic copy of this thesis provided to the James Cook University Library is an accurate copy of the print thesis submitted, within the limits of the technology available.

_____          _____

Signature                                                              Date

# Statement of Contributions

# Acknowledgements

Studying the same thing for three and a half years and then writing a thesis to prove it attains a PhD. With this in mind I would first like to thank the people who have distracted me and made me realise that there is far more to life than tree-based models for mixed type clustering and profiling. In particular I would like to thank all mu friends and colleagues that have had to bare the load of everything that went even just a little bit wrong over the last three and a half years. My family in particular must also be mentioned as over the course of my PhD they had to put up with even more. Thanks Dad for reading mine opus. I would like to thank my supervisors, Danny Coomans and Bruce Litow, for the time and effort they gave such that my PhD was. Also I would like to thank the people at FABI for food, shelter and support for my trip over to Belgium, Europe and the world.

In short this thesis is the culmination of a lot of people's insight and effort that has been condensed into the next 200 odd pages, thanks all …

# Abstract

Multivariate profiling aims to find groups in a response dataset that are described by relationships with another. Profiling is not predicting each variable within the response set, but finding stable relationships between the two datasets that define common groups. Profiling styles of analysis arise commonly within the context of survey, experimental design and diagnosis type of studies. These studies produce complex multivariate datasets that contain mixed variables often with missing values that require analysis with a flexible, stable statistical technique.

The profiling model under consideration within this thesis is a Classification and Regression Tree (CART). A standard CART model finds groups within a univariate response by building a decision tree from a set of predictor variables. The flexible structure of a CART model allow it to be used for either discriminate or regression analysis whilst also catering for mixed types within the predictor set.

**The goal of this thesis to develop methods that extend CART for a multivariate response dataset involving mixed data types**. Multivariate regression for CART (MRT) has recently been shown to be a powerful profiling and clustering tool. However the same successes in extending CART for multivariate classification and multivariate mixed type analysis is yet to be realised. To begin with thesis explores simple extensions to CART for multivariate mixed type analysis. These are binary substitution of categorical variables within the response set and partitioning of a distance matrix using Db-MRT. These techniques use already existing extensions to

CART methods and are used as comparison methods to gauge the performance of the ensemble and consensus approaches that are the focus of this thesis.

Ensemble models using CART, such as random forests and treeboost, not only improve the overall accuracy of the model predictions but also introduce an ensemble proximity matrix as a measure of similarity between observations of the response set. In this thesis, through MRT, extensions to both random forests and treeboost are developed such that they predict a multivariate response. Furthermore, by binary substitution of the categorical variables within the response set these multivariate ensemble techniques are further extended to mixed type profiling. A result of this extension is that the ensemble proximity matrix now describes the groups found within the multivariate response. In this way multivariate tree-base ensembles can be interpreted as a cluster ensemble method, where the ensemble proximity matrices can be seen as cluster ensemble consensus matrices. In this thesis these proximity matrices are found to be powerful visualisation tools providing improved resolution of group structure found by a multivariate ensemble method. More so, as in cluster ensembles using these matrices as an input in to a clustering method improves the accuracy of the groups found.

**The main work of this thesis is the development of the Multivariate Consensus Tree (MCT) framework for mixed type profiling**. Motivating the MCT approach is the need to further understand which variables relate to the groups observed within the proximity matrix. To do this MCTs describe three methods to intelligently combine the ensemble proximity matrices of individual responses into one overall consensus matrix. This consensus matrix is a summary of the overall group structure

within each individual proximity matrix. As MCTs work solely with proximity matrices they are independent of the data types within the variables of the response set. Furthermore as each response variable is explicitly predicted it is possible to assess the quality of each proximity matrix in terms of predictive accuracy of the corresponding ensemble.

The MCT consensus matrix is a visualisation tool for the groups present within both the response and predictor datasets. As a consensus matrix is a similarity matrix this thesis proposes five new splitting criteria for tree-based models that search for decision rules within variables of the predictor set that partition the consensus matrix into the observed groups. This tree provides a logical decision path that predicts each group. As the groups within the response are now defined by their relationships within the predictor set, the MCT profiling is complete. This thesis proposes two algorithms for building an MCT; global MCTs and local MCTs. Global MCTs construct an overall consensus matrix spanning all observations, and recursively partition on this matrix to build the tree. Local MCTs build a new consensus matrix at each terminal node to evaluate each new split.

As MCTs have the proximity matrices to summarise the group structure within each response variable methods to identify important subgroups within these variables are also proposed. This search for subgroups within the response can be done on two levels. Firstly to identify subgroups of response variables for overall analysis; and secondly to identify subsets of response variables within any specific group found by the MCT. By finding subsets of response variables that relate to specific group structure the understanding of structure within the dataset is greatly improved.

This thesis shows tree-based methods for profiling, in particular MCTs, to be a powerful tool for mixed type analysis. Firstly, the visualisation of the tree, combined with the proximity matrices, provide a unique view of the groups found and allow for their easy interpretation within the context of the analysis. Secondly, MCTs are shown to accurately estimate the number of groups and provide measures on their stability and accuracy. Furthermore, MCTs are found to be resistant to noise variables within the analysis. Finally they provide methods to find subgroups within the response variables and to identify unimportant variables from the analysis. Throughout this thesis these tree-based methods are compared with standard clustering techniques to provide an accurate benchmark for their performance.

# Table of Contents

# List of Figures

# List of Tables