# Analyzing volatile compound measurements using traditional Multivariate techniques and Bayesian networks

A thesis presented in partial fulfillment of the requirements

for the degree of

Master of Arts

in

Statistics

at Massey University, Albany, New Zealand

Shweta Baldawa

2009

# Abstract

The purpose of this project is to compare two statistical approaches, traditional multivariate analysis and Bayesian networks, for representing the relationship between volatile compounds in kiwifruit. Compound measurements were for individual vines which were progeny of an intercross. It was expected that groupings in the data (or compounds) would give some indication of the generic nature of the biochemical pathways. Data for this project was provided by the Flavour Biotech team at Plant and Food Research. This data contained many non-detected observations which were treated as zero and to deal with them, we looked for appropriate value of $c$ for data transformation in $\log(x+c)$. The data is 'large $p$ small $n$' paradigm – and has much in common with data, although it is not as extreme as microarray. Principal component analysis was done to select a subset of compounds that retained most of the multivariate structure for further analysis. The reduced set of data was analyzed by Cluster analysis and Bayesian network techniques. A heat map produced by Cluster analysis and a graphical representation of Bayesian networks were presented to scientists for their comments. According to them, the two graphs complemented each other; both graphs were useful in their own unique way. Along with clusters of compounds, clusters of genotypes were represented by the heat map which showed by how much a particular compound is present in each genotype while the relation among different compounds was seen from the Bayesian networks.

# Acknowledgments

# Table of Contents

# List of Figures

# List of Tables