

# Clustering of Proteomics Imaging Mass Spectrometry Data

Annie Conway

*Thesis submitted for the degree of*

*Master of Philosophy*

*in*

*Statistics*

*at*

*The University of Adelaide*

*Faculty of Mathematical and Computer Sciences*

School of Mathematical Sciences



January 12, 2016



# Contents

<b>Signed Statement</b>	<b>vi</b>
<b>Acknowledgements</b>	<b>vii</b>
<b>Abstract</b>	<b>viii</b>
<b>Notation Index</b>	<b>ix</b>
<b>1 Introduction</b>	<b>2</b>
1.1 Chapter overview . . . . .	3
<b>2 Proteomics Imaging Mass Spectrometry</b>	<b>5</b>
2.1 Proteomics and ovarian cancer . . . . .	6
2.2 Mass spectrometry . . . . .	7
2.3 A brief overview of the acquisition of data for this project . . . . .	8
2.4 Conclusion . . . . .	11
<b>3 Functional Data Analysis</b>	<b>14</b>
3.1 Introduction . . . . .	14
3.1.1 The extension of multivariate to functional . . . . .	15
3.2 Smoothing of functional data . . . . .	18
3.2.1 The basis function approach . . . . .	19
3.2.2 Kernel smoothing . . . . .	21
3.3 Preprocessing and interpolation of IMS data . . . . .	22

3.4	Interpolation . . . . .	24
3.5	Conclusion . . . . .	25
<b>4</b>	<b>Normalisation</b>	<b>26</b>
4.1	Motivation for normalisation . . . . .	26
4.2	Some simple normalisation techniques . . . . .	27
4.2.1	Normalisation factors . . . . .	28
4.3	Comparison of normalisation techniques . . . . .	29
4.3.1	Distribution of the calibrants before and after normalisation . . . . .	30
4.4	Peak intensity correction . . . . .	34
4.4.1	Distribution of the calibrants after PIC . . . . .	37
4.5	Conclusion . . . . .	39
<b>5</b>	<b>Clustering</b>	<b>42</b>
5.1	Clustering methods . . . . .	43
5.1.1	Hierarchical clustering . . . . .	43
5.1.2	$k$ -means clustering . . . . .	44
5.2	Comments on the parameters of $k$ -means . . . . .	46
5.3	Clustering raw peak-list data . . . . .	50
5.4	Clustering after a binary transformation . . . . .	54
5.5	Clustering after normalisation . . . . .	56
5.6	Cluster after log transformation . . . . .	56
5.7	Visualisation of closeness to the cancer cluster . . . . .	61
5.8	Conclusion . . . . .	62
<b>6</b>	<b>Principal Component Analysis and Clustering</b>	<b>64</b>
6.1	Principal component analysis . . . . .	65
6.2	The $k$ -means objective function . . . . .	67
6.3	Clustering when $k = 2$ . . . . .	67
6.4	Theorem . . . . .	72

6.4.1	Discussion . . . . .	75
6.5	Application for $k = 2$ and $k > 2$ . . . . .	76
6.6	Further comments . . . . .	80
6.7	A practical application of PCA clustering . . . . .	81
6.8	Conclusion . . . . .	81
<b>7</b>	<b>Comparison of Clustering Methods</b>	<b>82</b>
7.1	Methods of comparison . . . . .	83
7.1.1	Jaccard distance . . . . .	83
7.1.2	Variation of information . . . . .	84
7.1.3	The Rand index . . . . .	86
7.1.4	Prediction strength . . . . .	87
7.2	Application of measures of comparison . . . . .	89
7.2.1	Jaccard distance and variation of information: comparisons with the binary transformation . . . . .	89
7.2.2	Comparing PCA clustering with $k$ -means clustering . . . . .	90
7.2.3	Rand index: compare clustering on full dataset with clustering on principal components . . . . .	93
7.2.4	Prediction strength and prediction error loss: a four-way com- parison . . . . .	95
7.3	Conclusion . . . . .	97
<b>8</b>	<b>Conclusion</b>	<b>98</b>
	<b>Bibliography</b>	<b>99</b>

# Signed Statement

I certify that this work contains no material which has been accepted for the award of any other degree or diploma in my name, 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. In addition, I certify that no part of this work will, in the future, be used in a submission in my name, for any other degree or diploma in any university or other tertiary institution without the prior approval of the University of Adelaide and where applicable, any partner institution responsible for the joint-award of this degree.

I give consent to this copy of my thesis, when deposited in the University Library, being made available for loan and photocopying, subject to the provisions of the Copyright Act 1968. I also give permission for the digital version of my thesis to be made available on the web, via the University's digital research repository, the Library Search and also through web search engines, unless permission has been granted by the University to restrict access for a period of time.

SIGNED: ..... DATE: .....

# Acknowledgements

I give my foremost thanks to my supervisor Inge Koch for her helpful advice and teaching, as well as her unwavering support and encouragement.

I would also like to thank Lyron Winderbaum and the Adelaide Proteomics Centre for providing data, assistance and information. I acknowledge The University of Adelaide for my scholarship and The School of Mathematical Sciences for giving additional support as well as providing me with a place to work.

Finally, I would like to thank my friends and family for their support and encouragement. Many people were a source of insightful conversation which was valuable for both my research and understanding of mathematics.

# Abstract

This thesis presents a toolbox for the exploratory analysis of multivariate data, in particular proteomics imaging mass spectrometry data. Typically such data consist of 15000 - 20000 spectra with a spatial component, and for each spectrum ion intensities are recorded at specific masses. Clustering is a focus of this thesis, with discussion of  $k$ -means clustering and clustering with principal component analysis (PCA). Theoretical results relating PCA and clustering are given based on Ding and He (2004), and detailed and corrected proofs of the authors' results are presented. The benefits of transformations prior to clustering of the data are explored. Transformations include normalisation, peak intensity correction (PIC), binary and log transformations. A number of techniques for comparing different clustering results are also discussed and these include set based comparisons with the Jaccard distance, an information based criterion (variation of information), point-pair comparisons (Rand index) and a modified version of the prediction strength of Tibshirani and Walther (2005).

These exploratory analyses are applied to imaging mass spectrometry data taken from patients with ovarian cancer. The data are taken from slices of cancerous tissue. The analyses in this thesis are primarily focused on data from one patient, with some techniques demonstrated on other patients for comparison.



# Notation Index

$\mathbb{R}$	the real numbers.
$p$	number of variables.
$n$	number of observations.
$x$	(observation of) a functional random variable.
$\mathbf{x}$	(observation of) a multivariate random variable, $p \times 1$ .
$\mathbb{X}$	data matrix, $p \times n$ .
$\boldsymbol{\mu}$	the mean of $\mathbf{x}$ , $p \times 1$ .
$\Sigma$	the covariance matrix of $\mathbf{x}$ , $p \times p$ .
$\bar{\mathbf{x}}$	the sample mean of $\mathbb{X}$ , $p \times 1$ .
$S$	the sample covariance matrix of $\mathbb{X}$ , $p \times p$ .
$d(\mathbf{x}_1, \mathbf{x}_2)$	distance between two vectors $\mathbf{x}_1$ and $\mathbf{x}_2$ .
$\delta(C_1, C_2)$	distance between two sets $C_1$ and $C_2$ .
$k$	number of clusters.
$\mathcal{C} = \{C_1, \dots, C_k\}$	a $k$ -cluster arrangement.
$\mathcal{P}(\mathbb{X})$	the power set of $\mathbb{X}$ , i.e. the set of all subsets of $\mathbb{X}$ .
$m/z$	mass-on-charge.