# Statistical analysis of proteomic mass spectrometry data for the identification of biomarkers and disease diagnosis

Tyman Stanford

**Discipline of Statistics**
**School of Mathematical Sciences**

THE UNIVERSITY
*of* ADELAIDE

# 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: .......................

# Contents

iii

# List of Figures

# List of Tables

# Acknowledgements

I must emphasise my gratitude to my supervisors, Professor Patty Solomon and Dr Chris Bagley. Patty, I have tremendous admiration for your statistical knowledge and thank you for your fantastic insights, guidance and wisdom along the way. Chris, your incredibly sharp eye combined with your patience when answering my questions is thoroughly appreciated.

Mum, Dad, Mel and Liana, I hope we can do away with the "are you done yet?" question. Thank you to you all for your encouragement from near and far, recently and formerly. I am under no illusions that without your combined support I would not be writing this now. I will be a better son, brother and partner now I promise. I extended this to friends who have been very understanding of my absenteeism. Thank you for your support as well.

Thank you to Chris Davies also for his initial work in his honours thesis that allowed me to hit the ground running. Many thanks to The University of Adelaide, specifically to the School of Mathematical Sciences and those I have had the most contact with at the Adelaide Proteomics Centre: Megan Penno, Vicki Clifton and Peter Hoffmann.

Since I have the floor, there are some more general sentiments I would like to make. I am grateful to exist in the time and location I do, and to be able to do what I love. My exclamation of "what a time to be alive!" is rarely sarcastic, albeit a poor attempt at humour. It would be remiss of me not to reference 'standing on the shoulders of giants' (but to complete the metaphor, in my case, rather than standing I might be sitting or even sliding off). I also wish to thank others that I have not met, those who get insufficient acknowledgement for what is a tremendous service to society: people that create and maintain publicly available software. Particularly the authors and contributors of R and TeX/LaTeX, software I have used extensively in this thesis.

# Abstract

Proteomic spectra obtained from matrix-assisted laser desorption ionisation (MALDI) time-of-flight mass spectrometry (TOF-MS) are generated from the proteins and peptides present in serum obtained from blood. By ionising the proteins and resolving them in the mass spectrometer, data on the expression of proteins can be obtained, realised from the amplitude of signal for different mass to charge ratios. Of primary interest is the biological signal, in particular, the expression of proteins related to disease. In common with many 'omic' technologies, the raw spectra suffer from systematic errors due to technological artefacts and batch-effects, in addition to sample and biological variability. To negate these effects, novel application of genetic microarray pre-processing and analysis methods to proteomic TOF-MS data are presented. However, there are important differences between microarray and TOF-MS data which require consideration and non-trivial modifications to be successfully applied. One important difference between MALDI TOF-MS data and other high-throughput data, seldom addressed, is the high proportion of missing values.

The pre-processing of raw proteomic TOF-MS data needs to be undertaken prior to analysis and remains a mathematical and statistical challenge. Performed in distinct steps, pre-processing consists of signal smoothing, baseline correction, spectra normalisation, peak detection and peak alignment. An argument as to why the order of these steps is highly important is presented. Standard and novel data pre-processing methods are investigated and compared to optimise the process. Each step is given due consideration since the cumulative effects of substandard pre-processing can render subsequent statistical analysis highly unreliable.

Ultimately, the aim of proteomic MS is to analyse the protein profiles. Two different but related approaches to the analysis are undertaken. The first approach is to identify biological markers (biomarkers) that exhibit differential expression between disease groups. Identifying potential biomarkers for further research requires appropriate exploratory, visual and statistical modelling which is addressed in detail here. The second approach is to perform statistical discrimination between groups, a classical supervised learning problem. The ability of mathematical models to predict

disease groups using differential biological signal provides insight into the plausibility of diagnostic tests. Methodologically, supervised learning is a multifaceted problem given that feature selection, model parameter optimisation, and the handling of the training and test data all contribute to the inference that can be made from the results. Empirical appraisal of the methods applied to the proteomic data are provided with the outcome of discrimination error as a quantitative benchmark.

A number of proteomic TOF-MS datasets with differing characteristics are used throughout this thesis to assess the validity of the methods presented. The detailed analysis of a murine model MALDI TOF-MS dataset has facilitated the discovery of potential biomarkers for gastric cancer. Correct classification of spectra to their respective disease group (gastric cancer or control mice) as high as 97.4% was achieved using supervised learning. The thorough treatment of all the differently behaved datasets contained in this thesis, starting from the raw data pre-processing steps through to the challenging process of identifying potential biomarkers, provides a comprehensive and best-practice pipeline to analyse real-world proteomic MS data.

# Acronyms and abbreviations

For simplicity, many abbreviations will be used throughout this thesis. The abbreviation/acronym will appear in parentheses at the first occurrence of the phrase but the table below provides a comprehensive list for quick reference.

| Abbreviation | Meaning |
| --- | --- |
| APC | Adelaide Proteomics Centre |
| C | The portable and compiled programming language |
| C8 beads | Alkyl group beads used in proteomic sample fractionation |
| CLSA | Continuous line segment algorithm |
| CLN | Cyclic LOESS normalisation |
| CRC | Colorectal cancer |
| $CV$ | Coefficient of variation |
| (k)Da | (kilo)Daltons; $^1/_{12}$th of a carbon-12 atom's mass ($\sim 1.7 \times 10^{-27}$kg) |
| DNA | Deoxyribonucleic acid |
| DP | Dynamic programming |
| EQN | Empirical quantile normalisation |
| FDR | False discovery rate |
| FS | Fisher score |
| FWHM | Full-width at half-maximum |
| GC | Gastric cancer |
| GC-MS | Gas chromatography-mass spectrometry |
| GEE | Generalised estimating equation |
| $G$FCV | $G$-fold cross-validation, traditionally denoted $k$-fold |
| GLM | Generalized linear model |
| HM | Harmonic mean |
| IMAC-Cu | Immobilised metal affinity chromatography - copper |
| $k$NN | $k$-nearest neighbours |
| LC-MS | Liquid chromatography?mass spectrometry |
| LDA | Linear discriminant analysis |
| LME | Linear mixed effects |
| LOESS | Locally weighted scatterplot smoothing (local regression) |
| LSA | Line segment algorithm |

| Abbreviation | Meaning |
| --- | --- |
| MA | A transformation of paired minus vs. average log intensities |
| MAR | Missing at random |
| MALDI | Matrix-assisted laser desorption/ionisation |
| MCAR | Missing completely at random |
| MS | Mass spectrometry |
| $m/z$ | Mass divided by charge: the $x$-axis of TOF-MS |
| $\mu$m | Micrometre ($10^{-6}$ metres) |
| Nd:YAG | Neodymium-yttrium aluminium garnet (laser) |
| $n_k$ | The number of patients/subjects in $k = 1, \ldots, K$ groups |
| nm | Nanometre ($10^{-9}$ metres) |
| NW | Needleman and Wunsch (algorithm) |
| OOB | Out-of-bag |
| OLS | Ordinary linear least-squares (regression) |
| PC | Prostate cancer |
| PCA | Principal component analysis |
| PF | Pareto Front |
| PFDA | Pairwise fusion discriminant analysis |
| PLS | Penalised least squares (regression) |
| pH | Acidity/akalinity scale; hydrogen ion concentration metric |
| pmol/$\mu$L | Molecular concentration/microlitre; pmol $\approx 6 \times 10^{11}$ molecules |
| QDA | Quadratic discriminant analysis |
| R | The statistical programming environment |
| RDA | Regularised discriminant analysis |
| REML | Restricted maximum likelihood |
| RF | RandomForest |
| RNA | Ribonucleic acid |
| RUV | Remove unwanted variation |
| S2N | Signal to noise (ratio) |
| SAX | Strong anion exchange |
| SE | Structuring element |
| SELDI | Surface-enhanced laser desorption/ionisation |
| S-G | Savitzky-Golay |
| S$n$L$p$ | Small-$n$ Large-$p$ (problem) |
| SVA | Surrogate variable analysis |
| SVD | Singular value decomposition |
| SVM | Support vector machine |
| SW | Smith and Waterman (algorithm) |
| TCN | TIC normalisation |
| TIC | Total ion current |
| TOF | Time-of-flight |
| T$_x$ | Treatment |
| UV | Ultraviolet |
| WCX | Weak cation exchange |