# A hierarchical model to detect differential gene expression distributions, and their investigation as a reflection of dysregulation in cancer

**by Aedan Roberts**

Thesis submitted in fulfilment of the requirements for the degree of

**Doctor of Philosophy**

under the supervision of Professor Paul Kennedy and Associate Professor Daniel Catchpoole

# Certificate of Original Authorship

I, Aedan Roberts, declare that this thesis is submitted in fulfilment of the requirements for the award of Doctor of Philosophy in the School of Computer Science, Faculty of Engineering and Information Technology at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

SIGNATURE:

Production Note:
Signature removed prior to publication.

[Aedan Roberts]

DATE: 12$^{\text{th}}$ August, 2021

PLACE: Sydney, Australia

# Acknowledgements

I would like to thank my supervisors, Prof Paul Kennedy and A/Prof Daniel Catchpoole, for their support and advice throughout the course of my PhD. Their encouragement, guidance, and the many discussions on the direction of my work have been invaluable.

My endless gratitude goes to my wife Beatrice, for supporting, encouraging and tolerating me, not necessarily in that order. A nearly equal amount of gratitude goes to Waffles and Clio, who really didn't offer any useful advice to be honest, but whose presence was immensely helpful nonetheless.

I would also like to thank the support staff at the School of Computer Science for their help throughout the process, in particular Margot Kopel and Janet Stack. Thanks also go to Thomas Lysaght and George Mundackal, who helped to optimise my R code as part of their Computing Science Studio 1 project.

I am grateful for the financial support provided by the NSW Health PhD Scholarship and the Australian Government Research Training Program.

Finally, I thank the anonymous donors whose tissue donations to GTEx, TCGA and the GEO repository make this and so much other research possible.

# Contents

# List of Figures

# List of Tables

# List of Publications

Listed below are the publications and other outputs associated with the research presented in this thesis.

**Roberts, A. G. K.**, Catchpoole, D. R. & Kennedy, P. J., 2021, 'Identification of differentially distributed gene expression and distinct sets of cancer-related genes identified by changes in expression mean and variability'. (Submitted) Available as a preprint: https://www.biorxiv.org/content/10.1101/2021.02.15.431343v2.

**Roberts, A. G. K.**, 2021, 'DiffDist', https://github.com/aedanr/DiffDist. R package.

**Roberts, A. G. K.**, Catchpoole, D.R. & Kennedy, P.J., 2018, 'Variance-based Feature Selection for Classification of Cancer Subtypes Using Gene Expression Data', *2018 International Joint Conference on Neural Networks (IJCNN)*, Rio de Janeiro.

# List of Abbreviations and Symbols

| Abbreviation | Description |
|---|---|
| ALL | Acute lymphoblastic leukaemia |
| AP | Anti-profiles |
| AUC | Area under the ROC curve |
| BFDR | Bayesian false discovery rate |
| BRCA | Breast invasive adenocarcinoma |
| CGC | Cancer Gene Census |
| COAD | Colon adenocarcinoma |
| CV | Coefficient of variation |
| DE | Differential expression |
| DMD | Differences in means and deviations |
| $D\phi$ | Differential dispersion |
| FDR | False discovery rate |
| FPR | False positive rate |
| GAMLSS | Generalised additive models for location, scale and shape |
| GEO | Gene Expression Omnibus |
| GO | Gene Ontology |
| GTEx | Genotype–Tissue Expression |
| HM | Hierarchical model |
| HMM | Hierarchical mixture model |
| HPD | Highest posterior density |
| KEGG | Kyoto Encyclopedia of Genes and Genomes |
| KIRC | Kidney renal clear cell carcinoma |
| LDA | Linear discriminant analysis |
| LFC | Log fold change |
| LIHC | Liver hepatocellular carcinoma |
| LSVM | Linear SVM |
| LUAD | Lung adenocarcinoma |
| LUSC | Lung squamous cell carcinoma |

| Abbreviation | Description |
|---|---|
| MAD | Median absolute deviation |
| MCMC | Markov chain Monte Carlo |
| mRNA | messenger RNA |
| MSE | Mean squared error |
| NB | Negative binomial |
| NCBI | National Center for Biotechnology Information |
| NOS | Not otherwise specified |
| NSCLC | Non-small cell lung carcinoma |
| PRAD | Prostate adenocarcinoma |
| PSVM | SVM with polynomial kernel |
| RBF | Radial basis function |
| RF | Random forest |
| RIN | RNA integrity number |
| RLE | Relative log expression |
| RNA-seq | RNA sequencing |
| ROC | Receiver operating characteristic |
| RSVM | SVM with RBF kernel |
| SAGE | Serial analysis of gene expression |
| SAM | Significance Analysis of Microarrays |
| SVM | Support vector machine |
| TCGA | The Cancer Genome Atlas |
| THCA | Thyroid carcinoma |
| TMM | Trimmed mean of M-values |
| TPR | True positive rate |
| TRF | Random forest with distance-to-median transformed features |

| Symbol | Description |
| --- | --- |
| $g$ | Number of genes represented in a gene expression dataset |
| $m_\mu$ | Location hyperparameter for log-normal prior on mean |
| $m_\phi$ | Location hyperparameter for log-normal prior on dispersion |
| $n$ | Number of samples in a dataset |
| $n_A$ | Number of samples in group $A$ |
| $n_B$ | Number of samples in group $B$ |
| $\hat{R}$ | Gelman–Rubin diagnostic |
| $s$ | Sample standard deviation |
| $s^2$ | Sample variance |
| $v_\mu$ | Scale hyperparameter for log-normal prior on mean |
| $v_\phi$ | Scale hyperparameter for log-normal prior on dispersion |
| $y$ | Set of observed RNA-seq counts for all genes and samples in a dataset |
| $\bar{y}$ | Sample mean of $y$ |
| $y_{ij}$ | Observed count for gene $j$ in sample $i$ |
| $z_j$ | Mixture component indicator for gene $j$ |
| $\Gamma(\cdot)$ | Gamma function |
| $\gamma$ | Set of all hyperparameters |
| $\gamma_\mu$ | Set of hyperparameters for prior on mean |
| $\gamma_\phi$ | Set of hyperparameters for prior on dispersion |
| $\theta$ | Set of means and dispersions for all genes: $(\mu_j, \phi_j),\ j = 1, \dots, g$ |
| $\theta_j$ | Set of mean and dispersion for gene $j$: $(\mu_j, \phi_j)$ |
| $\lambda$ | Proportion of differentially distributed genes in HMM or Poisson rate parameter |
| $\mu$ | Mean |
| $\phi$ | Negative binomial dispersion |
| $\sigma$ | Standard deviation |
| $\sigma^2$ | Variance |

# Abstract

Data from genome-wide gene expression studies provides a wealth of information
on diseases such as cancer, which can lead to insights into disease mechanisms
and advances in diagnosis and treatment. Analysis of expression data is most
commonly aimed at identifying genes whose mean expression levels are increased
or decreased in disease compared to normal tissue, or between disease subtypes –
differential expression analysis. However, there is strong evidence that changes in
the variability of gene expression, without a difference in mean, can also be relevant.
Genes related to cancer have been shown to have changes in the variability of their
expression between normal and tumour tissue, and these differentially variable
genes have also been found to be informative for diagnostic and prognostic cancer
classification. The research presented in this thesis addresses several aspects of
research on differential gene expression variability, and the broader concept of
differential distribution, defined as any difference in the distribution of expression
values between groups.

   This work makes three contributions to knowledge, relating to cancer clas-
sification, identification of differentially variable or distributed genes, and the
biology of differential variability and distribution in cancer. Contribution 1 extends
previous work by demonstrating that genes identified by differential variability
or distribution can be used to classify closely related cancer subtypes, rather
than purely diagnostic or prognostic classification. Contribution 2 is a Bayesian
hierarchical model for RNA-seq data that provides tests for differential expression,
variability and distribution. The performance of each test is compared with existing
methods on simulated data and on real RNA-seq datasets modified to artificially
introduce changes in expression between groups. The differential expression test
is competitive with state-of-the-art methods, and the differential variability test
improves on existing methods, particularly for small sample sizes. The differential
distribution test is the first such test available for RNA-seq data. Contribution 3

builds on previous work by providing the first clear demonstration that differential variability and differential distribution analyses can identify cancer-related genes, and that differential expression and differential variability analyses identify distinct sets of cancer-related genes, each with different biological functions.

Overall, this research confirms and extends previous findings showing that changes in expression variability and distribution in cancer are both of biological significance and informative for classification. As well as further demonstrating the need to look beyond differential expression to a comprehensive assessment of changes in gene expression distributions, this work provides a method that enables the identification of these differentially distributed genes.