

A Two-Sample Distribution-Free Test with Applications to Correlated Genomic Data



Alison Jane Telford
Department of Statistics
University of Leeds

Submitted in accordance with the requirements for the degree of

Doctor of Philosophy

January 2019

The candidate confirms that the work submitted is his/her/their own and that appropriate credit has been given where reference has been made to the work of others.

This copy has been supplied on the understanding that it is copyright material and that no quotation from the thesis may be published without proper acknowledgement

© 2019 The University of Leeds and Alison Telford

I would like to dedicate this thesis to my supervisors, for being incredibly patient and helpful and for teaching me to never give up; no matter how hard things get.

Acknowledgements

I'd firstly like to thank my family, without whom I'd have never embarked on this incredible journey. They gave me the push and the confidence I needed to apply for a PhD and I've never looked back since. I should also thank Joe for sticking by me when things got hard and giving me the motivation to keep going. You always had faith in me even when my own faith wavered.

I thank my supervisors Dr Arief Gusnanto, Professor Charles Taylor and Dr Henry Wood for all the help they've given me throughout my PhD, and giving me the gentle encouragements to help me gain success. They have turned a masters graduate into a confident human being with amazing career prospects. I seriously cannot thank them enough.

I'd also like to thank Professor John Kent for being my assessor for the end of second and third year reviews. Whilst his comments were critical and always brought me back to reality, I began to ask myself WWJKS (what would John Kent say) when writing this thesis. By thinking this way, I was able to produce a piece of work to be proud of.

I want to add a special mention to Alastair Droop for spending hours with me sorting out my R code into an R package, without your help I would have gotten nowhere!

Abstract

This thesis focuses on the identification of genomic regions that exhibit significant differences of Copy Number Alterations (CNA) between two clinical groups. CNA are a structural variation in the human genome where some regions have more or less copy number than the normal two copies. CNA patterns in some genomic regions across patients have been shown to be associated with disease phenotypes. Our interest is in testing which genomic regions exhibit different distributions between two clinical groups to aid classification of patients on their subtype of cancer and discover new genomic markers for phenotypic identification. To do this we apply a two-sample test on each genomic region to test the null hypothesis that two distributions are equal.

Standard statistical tests are not adequate to deal with the characteristics of the data where the differences between the two groups lie in any one of the following aspects of the distribution: mean, variance, skewness, and multi-modality. When the null hypothesis is that two distributions are equal, the Anderson-Darling (AD) test is generally employed. The AD test was developed from the Cramer-von Mises (CvM) test statistic, which was originally proposed for a goodness-of-fit test. In the case of multi-modality, we find that the AD test often fails to identify true differences. We show, however, that the Cramer test - another modification to the CvM test - does not fail in the case of multi-modality. We have obtained the first four moments of the Cramer test statistic, which are not available previously. We also propose a new method for obtaining a p -value without using resampling techniques by approximating the distribution of the test statistic by a Generalised Pareto Distribution (GPD). By approximating the null distribution in this way, the calculation of the p -value is much faster than current methods, especially for large n . A simulation study indicates that the Cramer test is as powerful as other tests in simple cases and more powerful in more complicated cases.

To test our method, we applied the Cramer test on each genomic region to compare two groups of 76 lung cancer patients - 38 of which have adenocarcinoma type lung cancer and the other 38 have squamous carcinoma type lung cancer. Comparisons with the current method for identifying genomic regions of interest, KC Smart, also indicate that our method works well and is arguably preferable.

When the genome is split into separate regions, we show that adjacent (in genomic location) regions can exhibit very high correlation of CNA. High correlation between genomic locations suggests dependencies between the simultaneously performed tests. Because of these dependencies, multiplicity correction techniques for independent tests cannot be used alone as the number of independent tests performed is unknown. Methods exist to estimate the effective number of independent tests, however we find that these methods are slow and computationally expensive. Because of this, we extend work done on Fisher's method to combine dependent p -values. We compare this method to using a multivariate version of the Cramer test and show that the method produces similar results when performed on the lung cancer data set.

Contents

1	Introduction	1
1.1	Motivation and Background	1
1.2	Objectives	2
1.2.1	Identify a Two-Sample Test for use on Complex Data	2
1.2.2	Inference in Highly-Correlated Data	3
1.3	What are Copy Number Alterations?	3
1.4	Measuring CNA	4
1.4.1	Comparative Genomic Hybridization (CGH)	4
1.4.2	Array CGH (aCGH)	5
1.4.3	Next Generation Sequencing Technologies	5
1.5	Lung Cancer Data Set	7
1.5.1	Data Cleaning	7
1.5.2	Exploratory Data Analysis	8
1.6	Analysing CNA	12
1.6.1	Data Pre-processing	13
1.6.2	Analysis Per Sequence	14
1.6.3	Analysis Across Sequences	18
1.7	Comparative KC Smart	20
1.7.1	Exploratory Data Analysis	20
1.7.2	KC Smart Methodology	25
1.7.3	Critical Assessment of KC Smart	27
1.8	Hypothesis Testing for Identifying Genomic Regions of Interest	29
1.9	Thesis Overview	31
2	Identifying a Two-Sample Test to Locate Genomic Regions of Interest	35
2.1	Introduction	35
2.2	Parametric Tests	36

CONTENTS

2.3	Two Sample Tests Based on Empirical Cumulative Distribution Functions	36
2.3.1	Kolmogorov-Smirnov	37
2.3.2	Cramer-von Mises and Anderson-Darling	39
2.3.3	Cramer Test	42
2.4	Comparison to Current Literature	45
2.5	Further Two Sample Tests	45
2.6	Discussion	46
3	Properties of the Cramer Test	47
3.1	Introduction	47
3.2	Cramer Test Statistic	47
3.3	Moments	47
3.3.1	Expectation	48
3.3.2	Variance	49
3.3.3	Skewness	50
3.3.4	Kurtosis	51
3.4	Transformation of data	53
3.5	Expectation and Variance for Known Distribution	55
3.5.1	Example - Continuous Uniform Distribution	56
3.5.2	Example - Standard Normal Distribution	57
3.5.3	Results for Other Unimodal Distributions Z	58
3.5.4	Results for Multi-Modal Distributions Z	59
3.6	Discussion	59
4	A Faster Approach to Estimate the p-value of the Cramer Test	61
4.1	Introduction	61
4.2	Resampling Approaches	62
4.2.1	Permutation Test	62
4.2.2	Bootstrapping the Limiting Distribution	62
4.3	Empirical Approximations	63
4.3.1	ECD Function of $T_{n,m}$ for Various Distributions of X and Y	63
4.3.2	Generalised Pareto Distribution	66
4.3.3	Measuring the Accuracy in the Right Tail	68
4.3.4	Comparing Empirical and Theoretical Quantiles	69
4.4	Alternative Approaches	71
4.5	Discussion	72

5	Application of Two Sample Test	75
5.1	Introduction	75
5.2	Computational Considerations	76
5.2.1	Test Statistic	76
5.2.2	Moments	77
5.2.3	Reducing Number of Grid Points	79
5.3	Atest - an R Package	81
5.3.1	The Functions	82
5.3.2	An Example	83
5.4	Comparing GPD Method to Bootstrap and Permutation Approach	84
5.4.1	Speed	84
5.4.2	Accuracy	85
5.5	Simulation Study	87
5.5.1	Type-I error control	87
5.5.2	Sensitivity	89
5.6	Genomic Results	92
5.6.1	Results of Cramer Test	93
5.6.2	Results of KC Smart	95
5.7	Application to KC Smart Data	97
5.7.1	Segmenting the Data	98
5.8	Discussion	100
6	Examination of Correlation Structure in the Data	101
6.1	Introduction	101
6.1.1	Notation	101
6.2	Modelling the Correlation Structure Between Variables	102
6.2.1	Autocorrelation	102
6.2.2	Autoregressive Models	103
6.2.3	Multivariate Normal Distribution	106
6.3	Correlation between Patients, ρ^{pa} , ρ^{ps} and ρ^p	107
6.3.1	Lung Cancer Data Set Application	108
6.3.2	Correlation Across all Patients	109
6.4	Correlation between Windows, ρ^{wa} , ρ^{ws} and ρ^w	110
6.4.1	Lung Cancer Data Set Application	111
6.4.2	Correlation Between all Windows in Two Chromosomes	114
6.5	Correlation between p -values, ρ^{pval}	115
6.5.1	Finding a Relationship between ρ^{wa} , ρ^{ws} and ρ^{pval}	115
6.6	Discussion	119

CONTENTS

7	Multiple Testing for Dependent p-values	121
7.1	Multiplicity Burden	122
7.1.1	Estimating m	122
7.1.2	Example	123
7.1.3	Lung Cancer Data Set	126
7.1.4	KC Smart Data Set	127
7.1.5	Estimating Multiplicity Burden	130
7.2	Fisher's Combined Probability Test for Dependent p -values	131
7.2.1	Fisher's Combined Probability Test for Independent p -values	131
7.2.2	Adapting Fisher's Method for Dependent p -values	132
7.2.3	Calculating the Joint Expectation	133
7.2.4	Using Fisher's Combined Probability Test on Dependent p - values	137
7.2.5	Example	138
7.2.6	Lung Cancer Data Set	139
7.2.7	KC Smart Data Set	141
7.2.8	Using the Multivariate Version of the Cramer Test	144
7.3	Segmentation Methods	145
7.4	Discussion	147
8	Discussion	151
A	Alternative Choices of Hypothesis Tests	157
A.0.1	Skew-Normal Distribution	157
A.1	Skew-Adjusted t test	159
A.1.1	Asymmetric, Skew-Adjusted t -test	159
A.1.2	Using Welch's t -test	160
B	Alternative Methods to Obtain a Suitable Null Distribution	167
B.0.1	The (Scaled) Chi-Square Distribution	167
B.0.2	The Gamma Distribution	168
B.0.3	The Log-Normal Distribution	168
B.0.4	Other Two-Parameter Distributions	170
B.1	Transformation of $T_{n,m}$	171
B.1.1	Finding the Optimal ψ	171
B.1.2	Transforming $T_{n,m}$ to a Log-Normal Distribution	171
B.1.3	Transforming $T_{n,m}$ to a Gamma Distribution	172
B.2	Extreme Value Theorem	175

C	Miscellaneous Propositions	179
D	Variance of $T_{n,m}$ Proof	185
E	Third Moment of $T_{n,m}$ Proof	193
F	Fourth Moment of $T_{n,m}$ Proof	203
	References	264

List of Figures

1.1	The estimated CNA for each window along the genome for a patient with adenocarcinoma type lung cancer (top) and squamous carcinoma type lung cancer (bottom). The alternating colouring scheme indicates chromosomes 1–22.	9
1.2	Histograms of the estimated CNA across patients with adenocarcinoma type lung cancer (left) and squamous carcinoma type lung cancer (right). Each window is located at a specific position in a specific chromosome: Window 448 (67.05 – 67.2 Mbp, chromosome 1), Window 2023 (53.85 – 54 Mbp, chromosome 2) and Window 9546 (38.1 – 38.25 Mbp, chromosome 8).	11
1.3	The process of analysing CNA data. Each blue square represents a step in the analysis process. An arrow indicates a potential next step in the process once the previous step has been completed. . . .	13
1.4	The copy number alterations across the genome for a single patient with colorectal cancer before normalisation (top) and after normalisation (bottom). The bottom graph is the output of CNAnorm, with tumCont referring to the percentage of tumour content within the sampled cell. Red points refer to a chromosomal gain and blue points refer to a chromosomal loss. Vertical lines are used to show the separation of the chromosomes.	15
1.5	The estimated CNA for each probe along the genome for a sample from group 1 (top) and group 2 (bottom). The alternating colouring scheme indicates chromosomes 1–22.	22
1.6	Histograms of the estimated CNA across samples from group 1 (left) and group 2 (right) of the sample data set. Probe 1 is taken from chromosome 1, probe 800 is taken from chromosome 4 and probe 3000 is taken from chromosome 22.	23

LIST OF FIGURES

1.7	For one mouse or sample, the log2 ratios of probe intensities plotted against the genomic position of the probes for positive (top) and negative (bottom) gains (grey). The black lines represent the KC score calculated for arbitrary positions along the genome.	26
1.8	The output of the KC Smart analysis with $\sigma = 10^6$ on the sample data.	28
1.9	The output of the KC Smart analysis with $\sigma = 10^7$ on the sample data.	28
1.10	The estimated CNA in window 4170 across patients with adenocarcinoma type lung cancer (left) and squamous carcinoma type lung cancer (right). Both histograms have been plotted on the same x-axis scale to enable easy comparison.	30
2.1	The empirical cumulative distribution functions of two samples drawn from distributions X and Y such that $X \sim N(0,1)$ and $Y \sim N(1,1)$ and $n = m = 100$	37
2.2	The empirical cumulative distribution functions of two samples drawn from distributions X and Y such that $X \sim N(0,2)$ and $Y \sim N(0,1)$ and $n = m = 100$	38
2.3	The empirical cumulative distribution functions of two samples drawn from distributions X and Y such that $X \sim N(0,1)$ and $Y \sim N(1,1)$ and $n = m = 40$	39
2.4	The histograms of the samples X (left) and Y (right).	41
2.5	The empirical cumulative distribution functions of two samples drawn from distributions X and Y such that $X \sim N(1, 0.25)$ and Y which follows a mixture of normals distribution with probability density function $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$. and $n = m = 100$	42
2.6	The graphs of $\eta(t)$ plotted against t' where t' is the linearly transformed version of t so that $t' \in [0, 1]$ (left), and $\eta(t)$ plotted against $H_{n+m}(t)$ (right).	44
4.1	The histogram (left) and the ECD function (right) of the 10000 test statistics where $X \sim N(0,1)$, $Y \sim N(0,1)$, $n = m = 10000$ and $k = 10000$	64
4.2	The histogram (left) and the ECD function (right) of the 10000 test statistics where X and Y are both mixture distributions defined by $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$ for $n = m = 10000$ and $k = 10000$	64

4.3	The histogram (left) and the ECD function (right) of the 10000 test statistics where X and Y are both mixture distributions defined by $\frac{5}{9}N(1, 0.25) + \frac{3}{9}N(3, 0.25) + \frac{1}{9}N(5, 0.25)$ for $n = m = 10000$ and $k = 10000$	65
4.4	The ECDF curves of test statistics when X and Y are sampled from each distributional form.	65
4.5	The percentiles of the 10,000 sampled test statistics when X and Y are distribution as $N(0, 1)$ plotted against the percentiles of the fitted GPD with $\mu = 0.121$, $\sigma = 0.473$, and $\xi = -0.025$	67
4.6	The percentiles of the 10,000 sampled test statistics when X and Y both follow mixture distributions defined by $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$ plotted against the percentiles of the GPD with $\mu = 0.052$, $\sigma = 0.260$ and $\xi = 0.071$	67
4.7	The percentiles of the 10,000 sampled test statistics plotted against the percentiles of the fitted GPD with $\mu = 0.083$, $\sigma = 0.628$ and $\xi = 0.083$	68
5.1	The integrand of the test statistic $\mathcal{J}(t)$ plotted against $t \in [-3, 3]$	77
5.2	The time in seconds for calculating the mean (top left), variance (top right) and third moment (bottom) of the test statistic when $n = m \in [2, 250]$ using R and C++.	80
5.3	The ratios $\frac{E[T_{n,m}N_g]}{E[T_{n,m}z]}$ (top left), $\frac{\text{Var}[T_{n,m}N_g]}{\text{Var}[T_{n,m}z]}$ (top right) and $\frac{E[T_{n,m}^3N_g]}{E[T_{n,m}^3z]}$ (bottom) plotted against $N_g \in [2, 250]$	81
5.4	The average speed over 5 replications if 200,000 simultaneous hypothesis tests are performed using the permutation method, the bootstrap approach and the GPD method with $N_g = 50$ to calculate the p -value. In this scenario two samples are drawn from a mixture of normals distribution with probability density function $\frac{3}{4}N(0, 0.5) + \frac{1}{4}N(2, 0.5)$ and $n = m \in [2, 200]$	84
5.5	Top Left: The ratio of the p -values which are less than 0.10 calculated using the GPD method over the permutation approach. Top Right: The ratio of the p -values which are less than 0.10 calculated using the GPD method over the bootstrap approach. Bottom: The ratio of the p -values which are less than 0.10 calculated using the bootstrap approach over the permutation approach. Here, the sample size is $n = m = 50$, the number of replicates for the permutation and bootstrap approach is 10,000 and for the GPD method $N_g = 50$	86

LIST OF FIGURES

- 5.6 False positive rates for the Cramer test, (two-sample) t -test, Cramer-von Mises test, Anderson Darling test, F -test, and Kolmogorov-Smirnov test at different simulation settings: varying μ (top left panel), σ (top right panel), α (bottom left panel), and both α and σ with $\alpha = \sigma$ (bottom right panel), from skew-normal distribution (see Section A.0.1) $SN(\mu, \sigma, \alpha)$. In the bottom row figures, the values of α are within the interval $[0, 1]$ (bottom horizontal axis), which have been accompanied by the corresponding values of skewness γ (top horizontal axis). 88
- 5.7 False positive rates for the Cramer test, (two-sample) t -test, Cramer-von Mises test, Anderson Darling test, F -test, and Kolmogorov-Smirnov test at different simulation settings: varying d (left panel) and p_1 (right panel) from a multi-modal mixture distribution which follows $p_1N(1, 1) + p_2N(1 + d, 1)$ 89
- 5.8 Sensitivity for the Cramer test, (two-sample) t -test, Cramer-von Mises test, Anderson Darling test, F -test, and Kolmogorov-Smirnov test at different simulation settings: varying μ (top left panel), σ (top right panel), α (bottom left panel), and both α and σ with $\alpha = \sigma$ (bottom right panel), from skew-normal distribution $SN(\mu, \sigma, \alpha)$ in the first sample. In the second sample, the observations are drawn from $SN(0, 1, 0)$. In the bottom row figures, the values of α are within the interval $[0, 1]$ (bottom horizontal axis), which have been accompanied by the corresponding values of skewness γ (top horizontal axis). 91
- 5.9 Sensitivity for the Cramer test, (two-sample) t -test, Cramer-von Mises test, Anderson Darling test, F -test, and Kolmogorov-Smirnov test, where 100 observations are drawn from $N(1, 1)$ in the first sample, and 100 observations from a multi-modal mixture distribution which follows $(1 - p_1)N(1, 1) + p_1N(1 + d, 1)$ in the second sample. The left panel is the setting where d varies in the range $[0, 9]$ and p_1 is fixed at $\frac{1}{8}$. The right panel is the setting where p_1 varies in the range $[0.01, 0.5]$, d is fixed at 2. 92
- 5.10 Sensitivity for the Cramer test, (two-sample) t -test, Cramer-von Mises test, Anderson Darling test, F -test, and Kolmogorov-Smirnov test, where 100 observations are drawn from $N(0, 10)$ in the first sample, and 100 observations from $N(-d, \sigma)^\pi \cdot N(d, \sigma)^{1-\pi}$ distribution where $\pi \sim \text{Bernoulli}(\frac{1}{2})$, $d \sim [7, 10)$ and $\sigma = \sqrt{100 - d^2}$, in the second sample. The sensitivity figures are plotted as a function of d . 93

5.11	The p -values of individual genomic regions across the genome using the Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey line represents the Bonferroni corrected significance threshold for which 669 regions pass this threshold. The alternating colouring scheme indicates chromosomes 1–22 from the left. Sex chromosomes are excluded from the analysis.	94
5.12	The output after applying KC Smart to the lung cancer data set. The same mirror locations were used as the artificial data set from the KC Smart vignette.	96
5.13	The p -values across regions in the genome using the Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey line represents the Bonferroni corrected significance threshold. The alternating colouring scheme indicates chromosomes 1–22. Sex chromosomes are excluded from the analysis.	97
5.14	The first sequence in the artificial KC Smart data. The black line represents the segment means. The alternating colouring scheme indicates chromosomes 1–22, X, Y.	98
5.15	The p -values of each segment using the Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey line represents the Bonferroni corrected significance threshold. The alternating colouring scheme indicates chromosomes 1–22, X, Y.	99
6.1	The autocorrelation function evaluated at $\tau \in [1, 100]$ for a single patient with adenocarcinoma type lung cancer (top) and squamous carcinoma type lung cancer (bottom).	103
6.2	The autocorrelation function of an AR(1) process with $n = 17, 613$, $E[X_t] = 0$ and $\alpha = 0.99$	104
6.3	Realisation of an AR(22) (top) and AR(23) (bottom) model with parameters fitted using the CNA of two patients with adenocarcinoma and squamous carcinoma type lung cancer respectively.	105
6.4	Estimated CNA across windows for two patients with adenocarcinoma type lung cancer (left) and two patients with squamous carcinoma type lung cancer (right).	108
6.5	Estimated CNA across windows for patient 1 with adenocarcinoma type lung cancer and patient 1 with squamous carcinoma type lung cancer.	109

LIST OF FIGURES

6.6	Heat map showing the value of ρ^{pa} , ρ^{ps} and ρ^p calculated for each pair of patients with either adenocarcinoma or squamous carcinoma type lung cancer. The letters on the diagonal represent the subtype of cancer, i.e. “a” represents a patient with adenocarcinoma type lung cancer and “s” represents a patient with squamous carcinoma type lung cancer.	110
6.7	Estimated CNA plotted for patients with adenocarcinoma type lung cancer for window 6 against window 7 (top left), window 1662 against window 1663 (top right), window 6 against window 1663 (bottom left) and window 810 against window 951 (bottom right).	112
6.8	Estimated CNA for patients with squamous carcinoma type lung cancer plotted for window 6 against window 7 (top left), window 1662 against window 1663 (top right), window 6 against window 1663 (bottom left) and window 810 against window 951 (bottom right).	113
6.9	Heat map showing the value of ρ^{wa} calculated for each pair of windows in Chromosome 21 and 22.	114
6.10	Heat map showing the value of ρ^{ws} calculated for each pair of windows in Chromosome 21 and 22.	115
7.1	The p -values for each variable p when comparing two samples from a multivariate normal distribution with mean vectors defined in Equation (7.3) and covariance matrix defined by Equation (7.4). The horizontal grey line represents the 5% Bonferroni corrected significance threshold with $m = 50$	124
7.2	The histogram of the minimum p -values calculated for each 10,000 permutations. The dashed black line represents the fitted Beta(1,10.68) distribution.	125
7.3	The p -values for each variable p when comparing two samples from a multivariate normal distribution with mean vectors defined in Equation (7.3) and covariance matrix defined by Equation (7.4). The horizontal grey line represents the 5% Bonferroni corrected significance threshold with $m = 10$	125
7.4	The histogram (left) and the ECDF (right) of the 10,000 minimum p -values when the method proposed by Dudbridge and Gusnanto (2008) is used on chromosome 1. The dashed grey line represents the fitted Beta(1, m) distribution.	126

- 7.5 The p -values of individual genomic regions across the genome using the Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey lines represent the Bonferroni corrected significance threshold for each chromosome taking into account the percentage of effective independent tests from Table 7.1. The alternating colouring scheme indicates different chromosomes, starting with chromosome 1, 2, \dots , 22 from the left. Sex chromosomes are excluded from the analysis. 128
- 7.6 The histogram (left) and the ECDF (right) of the 10,000 minimum p -values. The dashed grey line represents the fitted Beta(1, m) distribution. 129
- 7.7 The p -values of individual genomic regions across the genome using the Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey line represents the Bonferroni corrected significance threshold with $m = 551$ effective independent tests. The alternating colouring scheme indicates different chromosomes. . . . 129
- 7.8 The relationship between ρ^{pv} and a for $\rho^{pv} \in [-1, 1]$ 135
- 7.9 Using a sliding block of 5 p -values, the Fisher's combined probability test p -values are plotted when the test is performed on each block B_i , $i = 1, \dots, 46$, of highly correlated p -values. The horizontal grey dashed line represents the 5% Bonferroni corrected significance threshold. 138
- 7.10 Using non-overlapping blocks of 5 p -values, the Fisher's combined probability test p -values are plotted using when the test is performed on each block B_i , $i = 1, \dots, 10$, of highly correlated p -values. The horizontal grey dashed line represents the 5% Bonferroni corrected significance threshold. 139
- 7.11 The p -values of each sliding block of 20 Cramer test p -values using the adjusted Fisher's method. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey line represents the Bonferroni corrected significance threshold with $m = 17173$. The alternating colouring scheme indicates different chromosomes. 140

LIST OF FIGURES

- 7.12 The p -values using the univariate Cramer test and the GPD to calculate the p -value for each window of the genome. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. Larger circles are plotted at windows which are identified as significant by both Fisher's combined probability test and using the Bonferroni correction (with the estimated number of effective independent tests), triangles are plotted at windows which are identified as significant by using the Bonferroni correction (with the estimated number of effective independent tests) only and plus signs are plotted at windows which are identified as significant by Fisher's combined probability test only. The alternating colouring scheme indicates different chromosomes. . . . 142
- 7.13 The p -values of applying Fishers combined probability test to each block B_i of p -values for each chromosome. 142
- 7.14 The p -values of each sliding block of 20 windows using the multivariate Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The alternating colouring scheme indicates different chromosomes. . 144
- 7.15 The p -values using the univariate Cramer test and the GPD to calculate the p -value for each window of the genome. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. Larger circles are plotted at windows which are identified as significant by both the multivariate Cramer test and Fisher's combined probability test, triangles are plotted at windows which are identified as significant by the multivariate Cramer test only and plus signs are plotted at windows which are identified as significant by Fisher's combined probability test only. The alternating colouring scheme indicates different chromosomes. 146
- 7.16 The results after applying the CBS segmentation technique to the p -values of the Cramer test applied to each window of the lung cancer data set. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The black lines represent the segments obtained after applying CBS. The horizontal grey line represents the significance threshold after dividing 0.05 by 2389 - the total number of segments. The alternating colouring scheme indicates different chromosomes. 147

A.1	Density scaled histograms of the estimated CNA in various windows across patients with adenocarcinoma type lung cancer (left) and squamous carcinoma type lung cancer (right). Each window is located at a specific position in a specific chromosome: Window 448 (67.05 – 67.2 Mbp, chromosome 1), Window 2023 (53.85 – 54 Mbp, chromosome 2) and Window 9546 (38.1 – 38.25 Mbp, chromosome 8). The solid black line represents the probability density function of the fitted skew-normal distribution.	158
B.1	The percentiles of the 10000 sampled test statistics plotted against the percentiles of the fitted scaled Chi-Square distribution with $c = 0.174$ and $f = 3.353$	168
B.2	The percentiles of the 10000 sampled test statistics plotted against the percentiles of the fitted Gamma distribution with $\alpha = 1.63$ and $\beta = 2.93$	169
B.3	The percentiles of the 10000 sampled test statistics plotted against the percentiles of the fitted log-normal distribution with $\mu = -0.82$ and $\sigma = 0.69$	169
B.4	The percentiles of the 10000 sampled test statistics plotted against the percentiles of the fitted log-normal distribution with $\mu = -1.47$ and $\sigma = 0.78$	170
B.5	The values of RT_{acc} plotted against $\psi \in (0, 1]$	174
B.6	The values of LT_{acc} plotted against $\psi \in [-1, 0)$	174
B.7	The QQ-plots comparing the percentiles of the fitted Gamma distributions against U_1 (left) and U_2 (right).	175

List of Tables

1.1	The number of windows in each chromosome when the window size is 150kbp.	8
1.2	The number of windows which 1) pass the Shapiro-Wilks normality test, 2) have evidence of skewness and 3) have evidence of multimodality for patients with adenocarcinoma and squamous carcinoma type lung cancer respectively.	10
1.3	The quantity of windows which have i number of peaks, $i = 1, \dots, 9$ for patients with adenocarcinoma and squamous carcinoma type lung cancer respectively after applying the clustering algorithm in R.	12
1.4	The hidden states for which QuantiSNP, PennCNV and GenoCN use within their hidden Markov models.	16
1.5	The number of probes in each chromosome for the sample data set.	21
1.6	The number of windows which 1) pass the Shapiro-Wilks normality test, 2) have evidence of skewness and 3) have evidence of multimodality for samples in group 1 and group 2 of the artificial data set respectively.	24
1.7	The quantity of windows which have i number of peaks, $i = 1, \dots, 9$ for samples in group 1 and group 2 of the artificial data set respectively after applying the clustering algorithm in R.	24
1.8	A tabular output of KC Smart showing the locations of the significant regions for a kernel width of $\sigma = 10^6$	27
1.9	A tabular output of KC Smart showing the locations of the significant regions for a kernel width of $\sigma = 10^7$	28
2.1	Number of rejections of H_0 out of 100 simulated datasets in which the Cramer-von Mises test, the Anderson-Darling test and the Cramer test is performed on 100 observations from X and Y such that $X \sim N(1, 0.25)$ and Y which follows a mixture of normals distribution with probability density function $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$	43

LIST OF TABLES

3.1 The values of \tilde{a} , \tilde{b} , and \tilde{c} in Eq. (3.13) for different underlying distributions of Z , in the expectation and variance of the test statistic $T_{n,m}$ under the null hypothesis. The symbol * indicates the value is obtained numerically. 59

3.2 The values of \tilde{a} , \tilde{b} , and \tilde{c} in Eq. (3.13) for different multi-modal underlying distributions of Z , in the expectation and variance of the test statistic $T_{n,m}$ under the null hypothesis. The symbol * indicates the value is obtained numerically. 59

4.1 For X and Y distributed as mixture distributions which follow $p_1N(1, 0.25)+p_2N(1+d, 0.25)+p_3N(1+2d, 0.25)$ with $n = m = 100$, the 95th percentile of the empirical cumulative distribution function for $k = 10,000$, the 95th percentile of the fitted GPD, and the probability p_k of obtaining a value from the fitted GPD larger than the 95th quantile of the empirical cumulative distribution function. ϕ denotes the normal probability density function. 69

4.2 For X and Y distributed as mixture distributions which follow $p_1N(1, 0.25)+p_2N(1+d, 0.25)+p_3N(1+2d, 0.25)$ with $n = m = 100$, the 99.5th percentile of the empirical cumulative distribution function for $k = 10,000$, the 99.5th percentile of the fitted GPD, and the probability p_k of obtaining a value from the fitted GPD larger than the 99.5th quantile of the empirical cumulative distribution function, for various distributions. ϕ denotes the normal probability density function. 70

4.3 For X and Y distributed as mixture distributions which follow $p_1N(1, 0.25)+p_2N(1+d, 0.25)+p_3N(1+2d, 0.25)$ with $n = m = 30$, the 95th percentile of the empirical cumulative distribution function for $k = 10,000$, the 95th percentile of the fitted GPD, and the probability p_k of obtaining a value from the fitted GPD larger than the 95th quantile of the empirical cumulative distribution function, for various distributions. ϕ denotes the normal probability density function. 71

4.4	For X and Y distributed as mixture distributions which follow $p_1N(1, 0.25) + p_2N(1+d, 0.25) + p_3N(1+2d, 0.25)$ with $n = m = 30$, the 99.5th percentile of the empirical cumulative distribution function for $k = 10000$, the 99.5th percentile of the fitted GPD, and the probability p_k of obtaining a value from the fitted GPD larger than the 99.5th quantile of the empirical cumulative distribution function, for various distributions. ϕ denotes the normal probability density function.	71
5.1	The number of genomic regions out of 17,613 with (unadjusted) p -values less than 0.05 under the Cramer test, the t -test, the F -test, the KS test, the AD test and the CvM test in our lung cancer dataset. The (i, j) th entry indicates the number of significant windows in both the i th and j th test.	94
5.2	The number of significant genomic regions after Bonferroni correction out of 17,613 when using the Cramer test and KC Smart.	96
6.1	The value of ρ^{pval} for each $\rho^{wa} \in [-1, -0.7]$ and $\rho^{ws} \in [-1, -0.7]$ when X and Y are simulated from a multivariate normal distribution.	116
6.2	The value of ρ^{pval} for each $\rho^{wa} \in [0.7, 1]$ and $\rho^{ws} \in [0.7, 1]$ when X and Y are simulated from a multivariate normal distribution.	116
6.3	The value of ρ^{pval} for each $\rho^{wa} \in [0.7, 1]$ and $\rho^{ws} \in [0.7, 1]$ when X and Y are simulated from a multivariate normal distribution.	117
6.4	The value of ρ^{pval} for each $\rho^{wa} \in [-1, -0.7]$ and $\rho^{ws} \in [-1, -0.7]$ when X and Y are simulated from correlated mixture distributions.	118
6.5	The value of ρ^{pval} for each $\rho^{wa} \in [0.7, 1]$ and $\rho^{ws} \in [0.7, 1]$ when X and Y are simulated from correlated mixture distributions.	118
6.6	The value of ρ^{pval} for each $\rho^{wa} \in [0.7, 1]$ and $\rho^{ws} \in [0.7, 1]$ when X and Y are simulated from correlated mixture distributions.	118
7.1	The number of effective independent tests m , and the percentage of effective number of independent tests when Dudbridge and Gusnanto (2008) 's method is performed on each chromosome.	127
7.2	The estimated multiplicity burden for each $r_1 \in [0.7, 1]$ and $r_2 \in [0.7, 1]$ for $p = 1000$ variables and $n = 100$ observations.	130
7.3	The number of significant windows (after Bonferroni correction for Fisher's combined probability test) when using Fisher's combined probability test and using the Bonferroni correction (with the estimated number of effective independent tests).	140

LIST OF TABLES

7.4	The Fishers combined probability test p -values when applied to blocks of p -values B_i from each chromosome i	143
7.5	The number of significant windows (after Bonferroni correction for Fisher's combined probability test) when using Fisher's combined probability test and the multivariate version of the Cramer test. . .	145
D.1	The number of occurrences for each case of i, j, k, l for the first summation in equation (D.4). Here i and j are interchangeable, similarly, k and l are interchangeable.	187
D.2	The number of occurrences for each case of i, j, k, l for the second summation in equation (D.4). Here only i and j are interchangeable.	188

Chapter 1

Introduction

1.1 Motivation and Background

Cancer is a complex disease. Certain types of cancer, like lung cancer, have multiple subtypes that require different treatments. Understanding why the subtypes of cancer require different treatment is important when determining the best course of treatment. [Smith and Sheltzer \(2018\)](#) identifies specific genetic alterations that are more common for certain subtypes of cancer. These alterations are Copy Number Alterations (CNA) or Copy Number Variations (CNV) which are the duplications and deletions of chromosomes which occur along the genome, see Sections [1.3](#) and [1.4](#) for a more detailed explanation on CNA, it's biological significance and how its measured.

CNA are extremely common in cancer ([Beroukhim et al., 2010](#)) and are biologically significant when detecting tumour subtypes ([Gusnanto et al., 2015](#)). Many studies, e.g. [Loo et al. \(2011\)](#), [Choi et al. \(2017\)](#) and [Wang et al. \(2016\)](#), have been done to identify genomic markers which display a difference in CNA between patients with different subtypes of cancer. Performing analysis on patients with known subtype of cancer will therefore help identify the subtype of cancer for a patient with unknown subtype. For example, if a genomic marker is identified and the differences between CNA understood, then a new patients subtype of cancer can be discovered by observing the behaviour of CNA at that genomic marker and matching it to the behaviour we expect to see for a subtype of cancer. The behaviour of CNA expected for each subtype of cancer at genomic markers are also identified through analysis. We should therefore be able to use this analysis for the purpose of classifying new patients on their subtype of cancer. The classification of patients is not discussed in detail in this thesis, however it is considered

1. Introduction

future work. In this thesis we instead focus on a different approach for identifying genomic regions for which CNA differs significantly between subtypes of cancer.

Statistical learning and classification techniques for analysing these data are currently a popular topic (Lin et al. (2018), Wang et al. (2017) and Lu et al. (2018)). Usually, CNA data will have a large amount of correlated variables p ($>10,000$) and only a small number of observations n (<100), therefore more care is needed when performing statistical learning and classification.

In order to locate genomic regions of significance, two-sample testing can be used on different genomic locations where each sample of CNA comes from a different subtype of cancer, see Section 1.8 for a more detailed motivation on identifying genomic regions of significance using two-sample testing. Choosing a suitable two-sample test for this purpose however is not a straightforward task and provides one of the main focuses for this thesis. For a chapter by chapter overview of the research carried out in this thesis see Section 1.9.

1.2 Objectives

There are many problems arising from the identification of genomic regions of interest using CNA to determine tumour subtypes. We will address two related problems in this thesis.

1.2.1 Identify a Two-Sample Test for use on Complex Data

When comparing CNA data between two different subtypes of cancer, we find that the data differs not only in the mean, but also in the variance, skewness and even multi-modality. Because of this, we require a two-sample test which is able to not only deal with this form of complex data but is also sensitive at identifying these differences between the two groups of data. A two-sample test that can meet these specifications has not yet been identified.

Further to this, the high-dimensional nature of the data will require the test to be performed simultaneously on a large amount of variables ($> 10,000$). To enable a fast, user-friendly computation, each test will need to be performed efficiently. Thus, a method which does not require resampling to calculate the p -value of the test is also needed.

1.2.2 Inference in Highly-Correlated Data

There are three ways which we can incorporate correlation into our approach to calculate the regions which have significantly different CNA between groups of patients, namely

1. Incorporate the correlation into the approach before the hypothesis tests are performed,
2. Incorporate the correlation into the approach after the hypothesis tests are performed,
3. Consider a multivariate version of the test.

As producing a multivariate version of a test can prove challenging, we focus our attention to a univariate version of the test. Now, correlation can be incorporated into the model using hidden Markov model methods (Section 1.6.2) to segment the data prior to analysis, however it is still likely that these methods cannot completely correct for the correlation in the data. We therefore attempt to research methods into incorporating correlation into our approach after the hypothesis tests are performed.

Usually, after performing many simultaneous hypothesis tests multiplicity corrections are applied. Whilst many multiplicity corrections exist when performing simultaneous independent tests, further research is required when simultaneous dependent tests are performed. For CNA data, the correlation between adjacent genomic locations is very high (> 0.9). Because of this high correlation, we cannot assume our tests are independent, therefore using multiplicity corrections for independent tests alone will not be enough. Whilst various methods exist already to correct for multiplicity in the case where the tests are dependent, further research is needed to extend these methods when the distribution of test statistics is unknown.

1.3 What are Copy Number Alterations?

Wu et al. (2014) describes copy number alterations as “gains and losses of large segments of the genome - ranging in size from a few kilobases to whole chromosomes”, where a kilobase is the unit in which the length of DNA is measured. One kilobase refers to a section of double stranded DNA that makes up one thousand nucleotides.

1. Introduction

In general, each human cell has two copies of every chromosome apart from chromosomes X and Y which vary between males and females. However certain diseases, especially cancer, are caused by alterations to the number of copies along the genome. For example some sections of the genome may have experienced a loss in genetic code meaning there are less than two copies of a chromosome in certain locations. Alternatively some sections of the genome may experience duplications of genetic code, causing the number of copies of the chromosome at certain locations to be larger than two.

CNA has been identified as causes for certain diseases and development abnormalities (Tang and Amon, 2013). For example, duplication of the gene SNCA is associated with Parkinson's disease (Singleton et al., 2003), and duplication of the gene GSKb is associated with bipolar disorder (Lachman et al., 2007). Of course, alterations in regions of the genome which are not gene specific can also be associated with disease phenotypes and many studies have been outlined by Tang and Amon (2013). Many cancers are also a consequence of CNA including breast cancer (Pollack et al., 2002) and prostate cancer (Wolf et al., 2004).

Whilst it is known that CNA play a role in the cause of disease phenotypes, identifying the type and location of the alterations in the genome will help determine how and where in the genome to target treatment. Thereby using information about a cause for a disease to help treat them.

1.4 Measuring CNA

Estimating the CNA at various locations along the genome has recently had many advances. Different technologies are being created to do this more efficiently and less costly. We discuss three main technologies used to measure CNA. We briefly explain how each technology works and provide the advantages and disadvantages for each.

1.4.1 Comparative Genomic Hybridization (CGH)

First developed by Kallioniemi et al. (1992) to detect copy number changes in solid tumours, CGH hybridises differentially labelled test DNA and control DNA to metaphase chromosomes (Theisen, 2008). The test or tumour DNA is labelled using a green fluorochrome and the control is labelled using a red fluorochrome (Weiss et al., 1999).

Hybridization, first introduced by Schildkraut et al. (1961), involves taking DNA from two different sources and combining them to create a single hybrid

DNA. The purpose of hybridization is to detect similarities and differences between the two strands of DNA which will be found through the ratio of fluorescence intensities. A higher intensity of red fluorescence will indicate a loss of DNA, a higher intensity of green fluorescence will indicate a gain of DNA and an equal intensity of green and red fluorescence will indicate no change between the two DNA.

For CGH, the metaphase chromosomes created after hybridization are investigated by looking at the fluorescence intensities along the chromosome. Duplications and deletions occurring along chromosomes can be found by plotting the ratio of fluorescence intensities across each chromosome (Kallioniemi et al., 1992).

CGH has an advantage in that it can quickly scan an entire genome for differences (Theisen, 2008). However a big disadvantage to the CGH method is that genomic abnormalities can only be detected if the size of the region altered is sufficiently large (10Mbp to 20Mbp) (Ostroverkhova et al., 2002). Also CGH requires the use of fresh cells, so will not work on most samples stored in hospitals, suggesting that this method is unsuitable for research purposes.

1.4.2 Array CGH (aCGH)

To overcome the disadvantages of the CGH method, Pinkel et al. (1998) combined the method with the use of microarrays. Using a microarray, tens of thousands of different DNA segments, called probes, are arranged in rows and columns on a glass slide (Govindarajan et al., 2012). The location of the DNA segment on the glass slide is known and can therefore be referred to.

For the method of aCGH, the test or tumour DNA is once again dyed with green fluorescence and the control DNA dyed with red fluorescence (Theisen, 2008). Segments of DNA from both the test and control DNA are mixed together and placed on the microarray and then hybridized. The microarray containing all the hybridized probes can then be analysed by looking at the fluorescence intensities of each probe and thus duplications and deletions can be identified in the same way as for CGH.

This method is preferable to using CGH as the segments of DNA making up the probes are a lot smaller than the metaphase chromosomes (Theisen, 2008) and thus the resolution is much higher than CGH.

1.4.3 Next Generation Sequencing Technologies

It has been quoted by Schuster (2007) that next-generation sequencing technology is transforming biology. These sequencing technologies are increasing the speed of

1. Introduction

sequencing (Xie and Tammi, 2009). Xie and Tammi (2009) also goes on to explain further advantages of sequencing over the aCGH method for detecting CNA.

Firstly test or tumour DNA and control DNA is split into smaller DNA sequences called “reads” (Muzzey et al., 2015). The reads are mapped to the human reference genome (Yoon et al., 2009) which is then split into equal sized non-overlapping windows. The window size can be chosen according to the users specifications. As discussed by Gusnanto et al. (2014), if the window size is too small, there may be windows which have no reads and thus observing a pattern is difficult. Alternatively if the window size is too wide, the genomic features will be smoothed out. Gusnanto et al. (2014) provides a method for estimating the optimal window size based on the read density per window using AIC and cross validation techniques.

Once a window size is chosen, the number of reads are counted (determined by its starting point) in each window and a ratio of read counts for the test or tumour DNA over the control DNA can be calculated (Wood et al., 2010). Mathematically speaking, let u_{ij} be the observed number of reads from a tumour sample in chromosome i , $i = 1, \dots, h$ where h is the total number of chromosomes in the study, and window j , $j = 1, \dots, \omega_i$ where ω_i is the total number of windows in chromosome i . Also let v_{ij} be the observed reads from a normal sample in chromosome i and window j . To estimate the copy number alteration in chromosome i , window j , the ratio of the the number of reads in each window and chromosome for the tumour sample is taken over the normal sample, i.e if r_{ij} represents the estimated CNA in chromosome i , window j , then

$$r_{ij} = \frac{u_{ij}}{v_{ij}}. \quad (1.1)$$

The ratio in Equation (1.1) can then be plotted against the windows j for each chromosome i to identify regions of duplication or deletion.

Gusnanto et al. (2011) states that an advantage of using sequencing compared to array technology is that the signal does not show saturation or background noise which is typical of hybridization techniques. Hurd and Nelson (2009) discusses further advantages of the sequencing technique compared to the microarray technology. Behjati and Tarpey (2013) states that a disadvantage of this technology is cost. However, as the technology is improving, the cost of sequencing is reducing.

1.5 Lung Cancer Data Set

To motivate and test our methods, we use a data set of 76 patients with two subtypes of lung cancer (Belvedere et al., 2012). Belvedere et al. (2012) has already identified genomic regions of significance between the two types of lung cancer. Because of this, any methods and conclusions we produce can therefore be compared to the conclusions found by Belvedere et al. (2012).

In the dataset, half of the patients have adenocarcinoma type lung cancer and the other half has squamous carcinoma type lung cancer. The data is collected using next generation sequencing technology (Section 1.4.3) with a window size of 150kbp (kilobase pairs). The total number of windows along the genome is 20,652 windows. Due to the size of the chromosomes being different e.g. Chromosome 1 is of length 248,956,422 base pairs (bp) compared to chromosome 22 which is of length 50,818,468bp (Genome Reference Consortium, 2017), the number of windows in each chromosome varies. Table 1.1 displays the number of windows in each chromosome for this data set. Note that the length of chromosomes are not always a multiple of the window size. Because of this, the last window in a chromosome is usually smaller than the window size. However, in practice the last few windows are often removed from analysis because the ends of the chromosomes are repetitive and thus too few reads are aligned.

1.5.1 Data Cleaning

As in Gusnanto et al. (2015), we remove the sex chromosomes, the mitochondria chromosome and the centromere regions as missing data can be problematic. The regions of missing data have very repetitive genomic sequences, so DNA cannot be reliably mapped. For better analysis and comparisons, we will remove any windows from the analysis in which the CNA is not recorded for more than one patient. In the windows where the CNA is not recorded for a single patient, we will replace the missing value with the mean CNA across the other patients in the same subtype of cancer for that window. Note that replacing single missing values with the mean CNA across the other patients is not necessary to do if the tests performed does not rely on the number of observations from each group to be equivalent. However, for simplicity and to avoid removing observations given the already small sample size in each group we decided to make the number of observations in each group equal. Note however that this could lead to potential biases and future recommendations if the sample sizes are large enough would be to ignore this step and perform the tests on unequal group sizes. Within the data

1. Introduction

Chromosome	Number of Windows
1	1662
2	1622
3	1321
4	1275
5	1207
6	1141
7	1061
8	976
9	942
10	904
11	901
12	893
13	768
14	716
15	684
16	603
17	542
18	521
19	395
20	421
21	321
22	343
X	1036
Y	396
M	1

Table 1.1: The number of windows in each chromosome when the window size is 150kbp.

set there are 3,039 windows across both types of lung cancer in which the CNA is not recorded for more than one patient. These windows are removed from the analysis for both types of lung cancer. This leaves 17,613 windows for analysis.

1.5.2 Exploratory Data Analysis

Figure 1.1 shows the estimated CNA for each window along the genome for a patient with adenocarcinoma type lung cancer and a patient with squamous carcinoma type lung cancer. Recall that for next generation sequencing, the estimated CNA is calculated by taking the ratio of the number of reads from a tumour cell divided by the number of reads from a normal cell per window. This means that an estimated CNA of 1 suggests that there are no duplications or deletions in that region. Note that for the patient with adenocarcinoma type lung cancer (Figure 1.1 top) there doesn't appear to be any duplications or deletions in many of the

regions of the genome as the estimated CNA is centred on the $y = 1$ line. Now compare this to the patient with squamous carcinoma type lung cancer (Figure 1.1 bottom) whose genome displays signs of many duplications and deletions. This therefore suggests that squamous carcinoma type lung cancer is associated with more sporadic CNA across the genome, whereas adenocarcinoma type lung cancer is associated with CNA in specific regions of the genome. Of course these conclusions are based on only two patients within the data set, and to get a clearer view of what the associated patterns of CNA are for each subtype of cancer, all patients need to be examined.

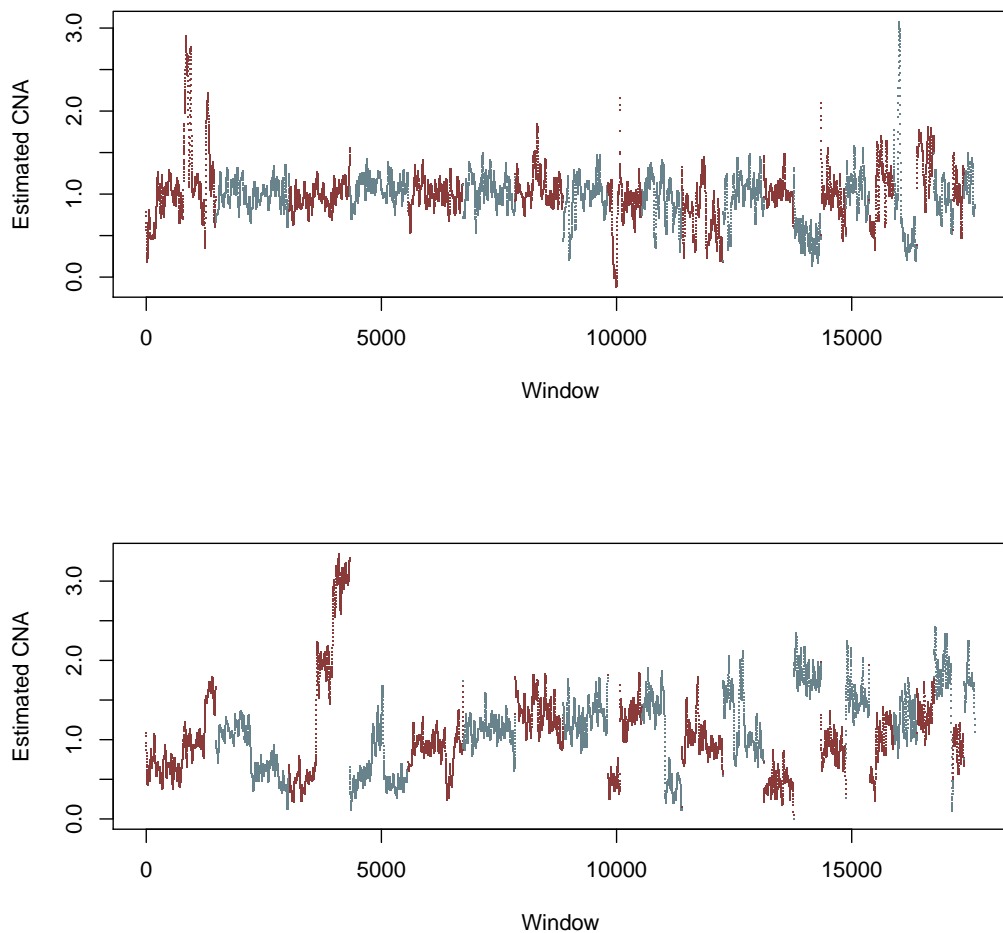


Figure 1.1: The estimated CNA for each window along the genome for a patient with adenocarcinoma type lung cancer (top) and squamous carcinoma type lung cancer (bottom). The alternating colouring scheme indicates chromosomes 1–22.

A selection of three windows is chosen to display the shape of the distributions of estimated CNA within the same patient group in Figure 1.2. It can

1. Introduction

be seen from Figure 1.2 that multi-modality is a common feature of the distributions of estimated CNA within the same patient group. As the estimated CNA $r_{ij} \in \{0.5, 1, 1.5, 2, \dots\}$, there will be peaks centred around these values with very few values in between. There does however exist special cases where $r_{ij} \notin \{0.5, 1, 1.5, 2, \dots\}$, these come from errors in cell division which result in regions of the genome being gained or lost, or they come from whole genome duplications, which increase the ratio across the entire sample. For certain regions, many patients within the same group may have common CNA. In this case it is expected that a larger peak is present centred around the “most common CNA”, with smaller peaks centred around other CNAs. Because outside factors like age or other biological influences could affect CNA it is possible that some patients within the same group will have different CNA for certain regions, this also gives rise to a potential multi-modal distribution with more than two modes - see for example Figure 1.2 (bottom right) which displays three main peaks.

It can also be seen from Figure 1.2 that there is larger variability of estimated CNA within patients with squamous carcinoma type lung cancer compared to the patients with adenocarcinoma type lung cancer. Given that the estimated CNA for the patient with squamous carcinoma from Figure 1.1 was very sporadic, the higher variability within patients with this type of lung cancer would suggest that the sporadic behaviour is consistent across all patients with squamous carcinoma type lung cancer and no clear average pattern exists.

Table 1.2 shows the number of windows which 1) pass the Shapiro-Wilks normality test, 2) have evidence of skewness and 3) have evidence of multi-modality for patients with adenocarcinoma and squamous carcinoma respectively.

Evidence of	Adenocarcinoma	Squamous Carcinoma
Normality	11396	4139
Skewness	6059	13294
Multi-Modality	5517	12605

Table 1.2: The number of windows which 1) pass the Shapiro-Wilks normality test, 2) have evidence of skewness and 3) have evidence of multi-modality for patients with adenocarcinoma and squamous carcinoma type lung cancer respectively.

Table 1.2 therefore shows that the majority of windows for patients with adenocarcinoma type lung cancer could be considered normally distributed and just under a third considered multi-modally distributed. The converse however is true for patients with squamous carcinoma type lung cancer with well over a third of windows considered multi-modally distributed. Out of the windows which had evidence of multi-modality, Table 1.3 shows the quantity which have i number of

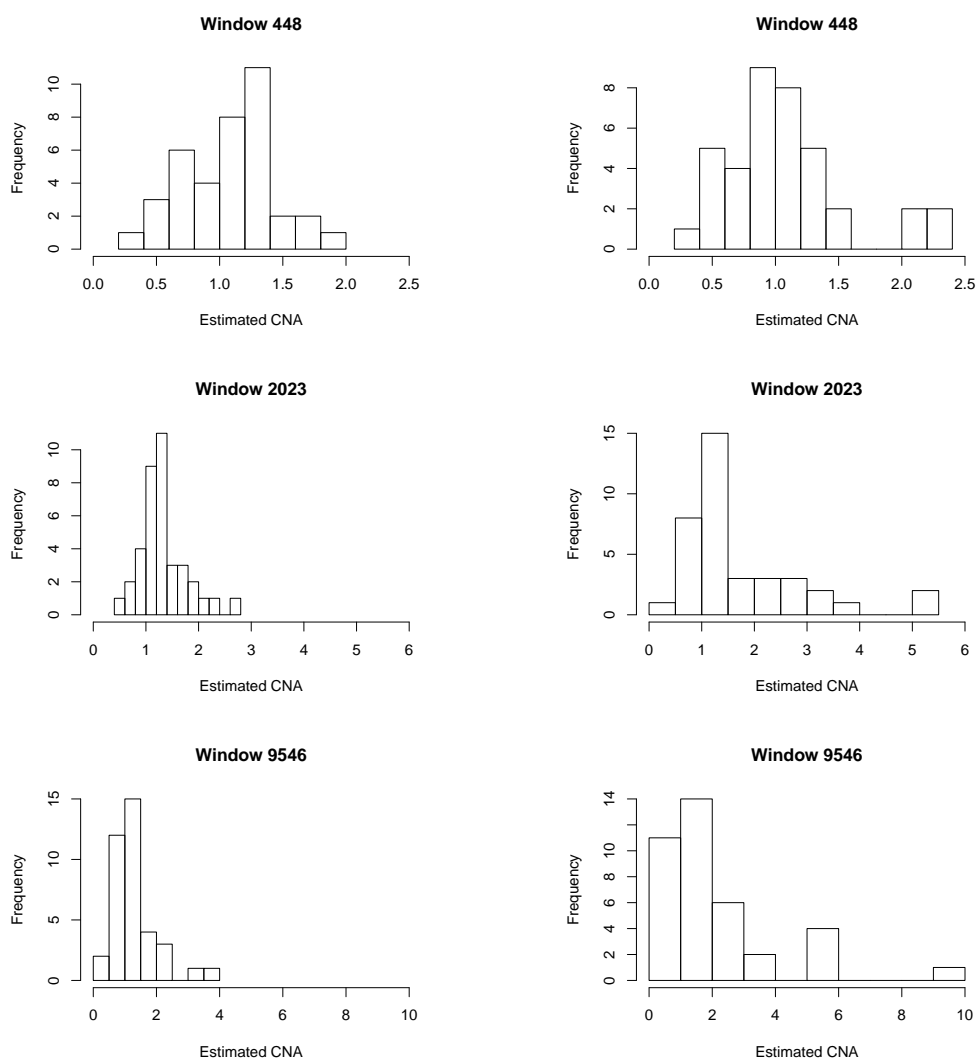


Figure 1.2: Histograms of the estimated CNA across patients with adenocarcinoma type lung cancer (left) and squamous carcinoma type lung cancer (right). Each window is located at a specific position in a specific chromosome: Window 448 (67.05 – 67.2 Mbp, chromosome 1), Window 2023 (53.85 – 54 Mbp, chromosome 2) and Window 9546 (38.1 – 38.25 Mbp, chromosome 8).

1. Introduction

peaks, $i = 1, \dots, 9$ for patients with adenocarcinoma and squamous carcinoma type lung cancer respectively. To obtain the estimated number of peaks, a clustering algorithm was applied in R.

i	Adenocarcinoma	Squamous Carcinoma
2	4927	11732
3	465	719
4	88	112
5	26	25
6	9	9
7	2	5
8	0	2
9	0	1

Table 1.3: The quantity of windows which have i number of peaks, $i = 1, \dots, 9$ for patients with adenocarcinoma and squamous carcinoma type lung cancer respectively after applying the clustering algorithm in R.

1.6 Analysing CNA

The process of analysing CNA data can be described by the flowchart in Figure 1.3. This process starts with the collection of CNA data using the techniques described in Section 1.4. The next step is data pre-processing, we describe this step in Section 1.6.1. After data pre-processing the analysis of CNA can generally be classified into two main types of analysis: analysis per sequence, and analysis across sequences. In the flow chart it can be seen that one has a choice of whether to analyse the data per sequence and then perform analysis across sequences, or skip the analysis per sequence and perform analysis across sequences. Analysis of CNA per sequence aims to estimate the CNA level at each genomic location per sequence and can further be classified into three groups of segmentation methods. The first approach is HMM-based methods, the second is binary segmentation, and the last is smooth segmentation based on a random effects model. We discuss methods which fit into these three groups in Section 1.6.2. Alternatively, analysis across sequences aims to identify differential pattern of CNA for a given genomic location between two groups of sequences or patients, which is the main focus of our study in this thesis. Currently, there exists one commonly used method for the analysis across sequences, namely KC Smart, and we discuss this method in detail in Section 1.7.

Note that, as shown in the flow chart, the analysis across sequences can either be performed on the unsegmented CNA data or on the segmented data, i.e. after

performing analysis per sequence. If the analysis across sequences is performed on the unsegmented data, then this has the advantage of granularity. This means results will be obtained for many small regions of the genome and any features that only exist at a very granular level will be identified. If the analysis across sequences is instead performed on the segmented data this has the advantage that fewer - potentially correlated - simultaneous tests need to be performed.

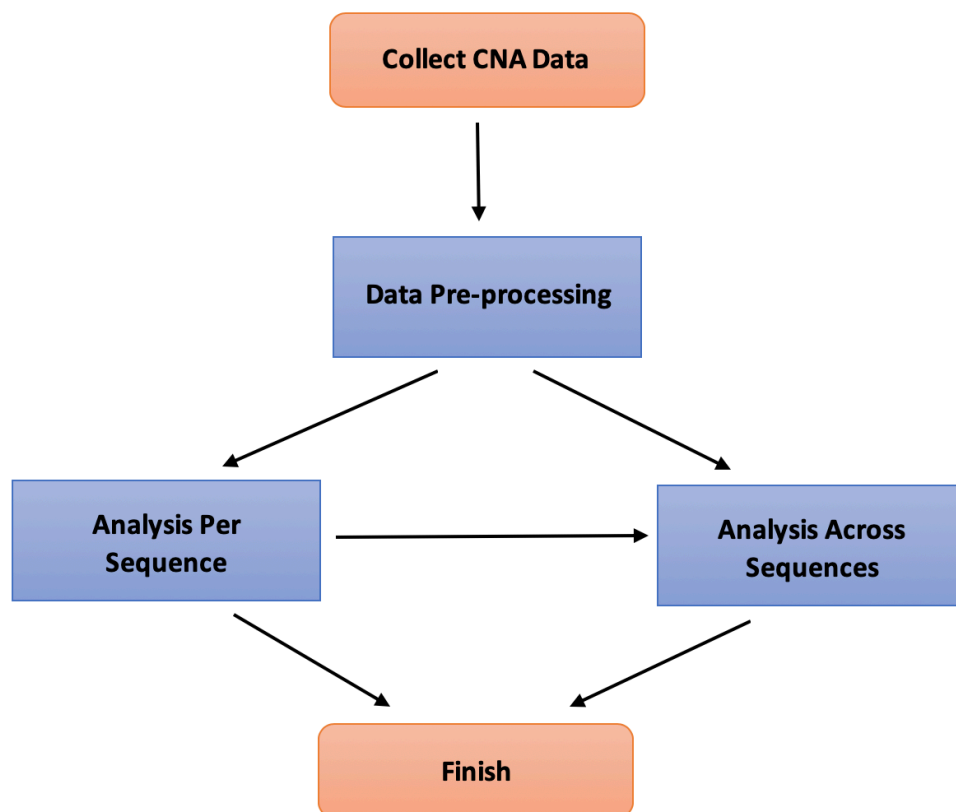


Figure 1.3: The process of analysing CNA data. Each blue square represents a step in the analysis process. An arrow indicates a potential next step in the process once the previous step has been completed.

1.6.1 Data Pre-processing

Data pre-processing concerns the data quality and normalisation. Normalisation is a crucial step in CNA analysis (Gusnanto et al., 2011) to correct for variations in the data affected by factors other than the copy number. Many methodologies for analysing CNA data also include a pre-processing normalisation step, for example CNV-seq (Xie and Tammi, 2009), CNVnator (Abyzov et al., 2011), FREEC (Boeva et al., 2010) and CNAnorm (Gusnanto et al., 2011). In the case of CNAnorm, the

1. Introduction

normalisation procedure is done to correct for various aspects which can affect the accuracy of the CNA approximations. For example, correcting for GC-content which affects staining intensity (Furey and Haussler, 2003), smoothing the data to reduce noise, and correcting for tumour sample contamination. Tumour sample contamination refers to the contamination of normal cells when the sample is taken from the cancerous tumour.

The ploidy of a cell is the number of chromosome copies in that cell. As humans generally have two copies of every chromosome, for a normal cell, the ploidy is two. If there exists copy number alterations within a cell, the ploidy of a genomic location can take any value from the set $\{0, 1, 2, 3, \dots\}$. Hence we would expect the ratio $r_{i,j}$ (1.1) to take any value from the set $\{0, 0.5, 1, 1.5, \dots\}$. Thus for normal genomic regions in the tumour cell, $r_{i,j} \approx 1$. When there exists contamination of normal cells with the tumour cells, $r_{i,j}$ will be shrunk towards ratio 1 (Gusnanto et al., 2011).

Another normalisation step which Gusnanto et al. (2011) considers is to align the estimated CNA so the most common genomic regions are centred at ratio 1. To illustrate the use of CNAnorm, Figure 1.4 gives the copy number alterations across the genome of a single patient with colorectal cancer before and after the normalisation step. The points represent the CNA measured for each window or region and the lines represent the output of the circular binary segmentation procedure DNACopy (see Section 1.6.2).

The proportion of tumour content is estimated by CNAnorm to be 44.07%, meaning that over half of the sample is contaminated by normal cells. This can be seen in Figure 1.4 (top) as the ratios are all shrunk towards a ploidy of 2. In Figure 1.4 (bottom) the ratios are centred so each segment has an integer ploidy. Note that it can be seen that regions of chromosome 20 do not align to an integer ploidy. This is because not every cell in a tumour contains every genomic change. Changes acquired late in a tumour's development may not be present in all cells, and will therefore appear as non-integer using this method.

1.6.2 Analysis Per Sequence

As mentioned at the start of this section, analysis of CNA per sequence can be classified into three groups of segmentation methods. Here we discuss some of the methods which fall into the three segmentation groups.

Circular binary segmentation (CBS) is a method introduced by Olshen et al. (2004) which attempts to segment the data using a modification of binary segmentation. This method was originally created for the analysis of aCGH data

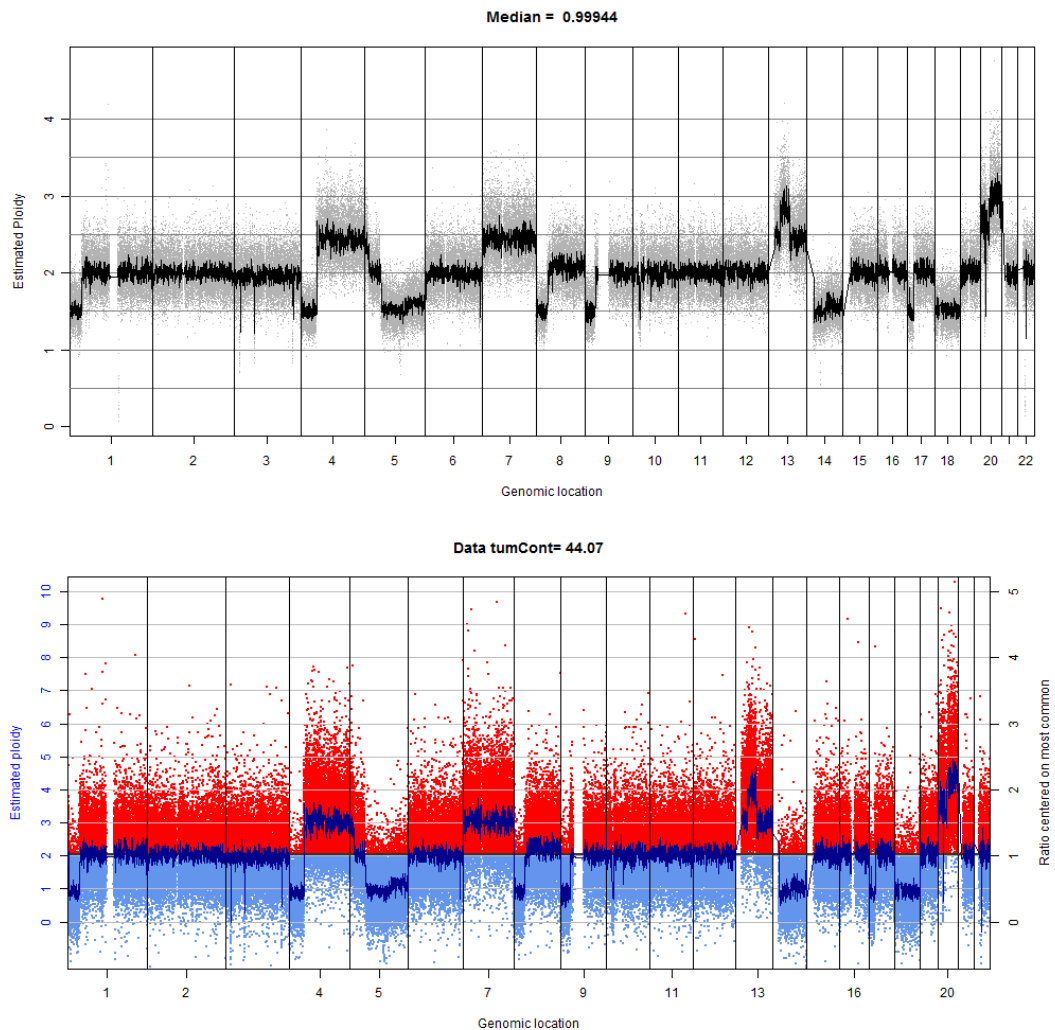


Figure 1.4: The copy number alterations across the genome for a single patient with colorectal cancer before normalisation (top) and after normalisation (bottom). The bottom graph is the output of CNAnorm, with tumCont referring to the percentage of tumour content within the sampled cell. Red points refer to a chromosomal gain and blue points refer to a chromosomal loss. Vertical lines are used to show the separation of the chromosomes.

1. Introduction

(Section 1.4.2) but can be easily adapted for use on next generation sequencing data. Let X_1, \dots, X_n represent the estimated CNA across all regions n and let $S_i = X_1 + \dots + X_i$, $1 \leq i \leq n$ be the partial sums. The method, which has been called DNACopy, calculates a statistic based on the partial sums to determine the locations of breakpoints. A breakpoint is defined as a location ν in which X_1, \dots, X_ν have a common distribution F_0 and $X_{\nu+1}, \dots$ have a common distribution F_1 where $F_0 \neq F_1$ until the next breakpoint. This method is used in the normalisation procedure CNAnorm described in Section 1.6.1 and the results of performing such a segmentation are shown in Figure 1.4.

Lai et al. (2005) shows that under various conditions CBS outperforms many other segmentation methods by comparing the Receiver Operating Characteristic (ROC) curves which quantifies sensitivity and specificity. However, CBS does not come without its limitations. CBS is not sensitive to short segments and often fails to detect them (Ben-Yaacov and Eldar, 2008), this means that any duplications or deletions which may span across only a small number of windows may be overlooked. Another limitation of this method is the computational time taken to segment high density arrays (Ben-Yaacov and Eldar, 2008).

Fridlyand et al. (2004) introduces the use of hidden Markov models to segment the data. They use an unsupervised hidden Markov model (HMM) approach to segment regions into sets with the same underlying copy number. They fit k -state hidden Markov models for $k = 1, \dots, K$, where K is the maximum number of states in the model. In the case of QuantiSNP (Colella et al., 2007), PennCNV (Wang et al., 2007) and GenoCN (Sun et al., 2009), $K = 6$ with Table 1.4 showing the hidden states, associated copy numbers and biological interpretation.

Hidden State	Copy Number	Description
1	0	Full Deletion
2	1	Single Copy Deletion
3	2	Normal (heterozygote)
4	2	Normal (homozygote)
5	3	Single Copy Duplication
6	4	Double Copy Duplication

Table 1.4: The hidden states for which QuantiSNP, PennCNV and GenoCN use within their hidden Markov models.

All three methods identify copy number states based on the log R ratio (LRR) and B allele frequencies (BAF). The log R ratio is defined as $\log_2(R_{\text{tumour}} / R_{\text{normal}})$ where R_{tumour} and R_{normal} are the copy numbers obtained from the tumour sample and normal sample respectively. The B allele frequency measures the normalised allelic contrast.

Each method also has its own way of defining the transition probabilities between each state and emission probabilities for the LRR and BAF. In the case of PennCNV, the method assumes that the mean and standard deviation of the LRR and BAF are known. PennCNV also incorporates the use of family information in CNV calling and validation, an advantage over other methods. QuantiSNP imposes some common priors for the LRR and BAF so that only a few hyperparameters need to be estimated. Finally, [Sun et al. \(2009\)](#) claims that GenoCN can provide output on allele-specific information whereas PennCNV and QuantiSNP cannot.

Further HMM methods for the analysis of CNA data per sequence include hsegHMM ([Choo-Wosoba et al., 2018](#)), GPHMM ([Li et al., 2009](#)) and MixHMM ([Liu et al., 2010](#)). [Choo-Wosoba et al. \(2018\)](#) states that PennCNV and QuantiSNP are based on the assumption of 100% tumour purity whereas [Van Loo et al. \(2010\)](#), [Li et al. \(2009\)](#) and [Liu et al. \(2010\)](#) account for both the tumour purity and ploidy. [Choo-Wosoba et al. \(2018\)](#) also states that all the aforementioned methods use the BAF which is sensitive to mapping bias, therefore instead of using BAF, hsegHMM uses the log Odds Ratio (LogOR).

HMM methods are powerful tools for analysing CNA data, and as many papers in the literature focus on these methods, are a popular method of choice to perform the analysis. However, these methods are not without limitations. Before analysis can be performed using HMM methods, the states are required to be defined. The methods have a very rigid structure for which the data can follow - a data point must fit into one of the predefined states - and are therefore less flexible than other methods.

[Huang et al. \(2007\)](#) also observes that binary segmentation techniques as well as hidden Markov model methods are based on modelling data as a series of discrete segments. It can be argued that in reality most segments are not discrete, providing further limitations of the CBS and HMM methods approaches. Even though the true underlying biological process is discrete, in reality many biological and environmental factors cause the CNA signal to deviate from a stepwise function ([Engler et al., 2006](#)), ([Picard et al., 2005](#)). Because of this, [Huang et al. \(2007\)](#) introduces smoothseg, a method which uses a smooth segmentation algorithm to analyse CNA per sequence. The statistical model used is based on a correlated random effects model. The estimation of the random effects model is carried out using the maximum likelihood method. [Huang et al. \(2007\)](#) also shows that smoothseg performs better than CBS as well as methods which use wavelet smoothing ([Hsu et al., 2005](#)).

1. Introduction

However, [Hsu et al. \(2005\)](#) argues that nonparametric regression techniques ([Jianqing and Gijbels, 1996](#)), ([Percival and Walden, 2006](#)) are suitable for data denoising as they “do not impose any parametric model in finding structures in the data sets”. Because sharp discontinuities of copy number changes can often occur in the tumour DNA and the sizes of the aberrations can vary between being very small (spanning a couple of windows) or very large (whole chromosome arms), [Hsu et al. \(2005\)](#) suggests that applying wavelet analysis is a desirable choice.

The idea of wavelet analysis is to represent the data as a linear combination of wavelets ([Hsu et al., 2005](#)). Consider n windows in a chromosome where for each window the relative copy number is measured. Denote $y(x_i)$ as the observed copy number alteration for the i th genomic location x_i . Then $y(x_i)$ can be expressed as

$$y(x_i) = f(x_i) + \epsilon_i$$

where $f(x_i)$ represents the “true” CNA signal and ϵ_i represents the errors and are independently and identically distributed as $N(0, \sigma^2)$. [Hsu et al. \(2005\)](#) therefore attempts to estimate the signal f by performing discrete wavelet transformations on the data. In particular, they use the maximal overlap discrete wavelet transformation (MODWT) together with the Haar wavelet family for their analysis. They treat the data as being equally spaced and use the SURE ([Donoho and Johnstone, 1995](#)) method for denoising the data. [Wang and Wang \(2007\)](#) later shows that assuming the data is equally spaced leads to “suboptimal” results and therefore considers a wavelet approach to analyse aCGH data which does not have this assumption. [Nguyen et al. \(2007\)](#) also considers dual tree complex wavelet transforms for analysing aCGH data.

As discussed, there exists many methods for the purpose of analysing CNA data per sequence. However the analysis we wish to carry out within this thesis is more concerned with the analysis across sequences, which we discuss further now.

1.6.3 Analysis Across Sequences

As mentioned previously, the purpose for analysing CNA per sequence is to identify common or recurrent regions of CNA ([Rueda and Diaz-Uriarte, 2010](#)). A common or recurrent region is defined by [Rueda and Diaz-Uriarte \(2010\)](#) as “a set of contiguous regions which, as a group, shows a high enough probability (or evidence) of being altered in at least some samples or arrays”. This means that methods attempt to locate “segments” of chromosomes which are consistently duplicated or deleted across patients within a single group.

In the literature, there exists many methods which aims to locate recurrent regions of duplication or deletion by analysing data across multiple sequences within a single group of patients. For example CGHregions (Van De Wiel and Van Wieringen, 2007) uses dimension reduction techniques, STAC (Diskin et al., 2006) which calculates two statistics based on the frequency of occurrence of the regions and the alignment of the regions, MSeq-CNV (Malekpour et al., 2018) which applies a mixture density to model the distribution of read counts and estimates the model parameters using the EM algorithm, and GISTIC (Beroukhim et al., 2007) which calculates a statistic based on both the frequency of occurrence as well as the amplitude of aberration. As well as these methods, work has been done to extend the use of HMMs to identify recurrent regions across multiple sequences. Some HMM methods include H-HMM (Shah et al., 2007) and pREC-A and pREC-S Rueda and Diaz-Uriarte (2009).

In this thesis, we are mainly concerned with the analysis across sequences which aims to identify differential pattern of CNA for a given genomic location between two groups of sequences or patients. There are very few methods in the literature which firstly aims to identify differential pattern of CNA and secondly enables the comparison of CNA between two groups of patients. More recently, a comparative version of KC Smart has been published (de Ronde et al., 2010) and is currently the method which is used in application due to its accessible nature. We provide more details and disadvantages to this method in Section 1.7 as we will use this method as a comparison to the methods we create.

It should be noted that other methods, for example as described in Diskin et al. (2006) use two-way agglomerative hierarchical clustering to determine possible subtypes and therefore classify patients into groups. As our analysis is done when the groups of patients are already known, these methods are not comparable. A further method to distinguish between groups of patients is to plot a heat map of the correlations between patients. Usually, patients within the same clinical group will have higher correlations between their CNAs and thus heat maps can be plotted to visually distinguish between subgroups.

Few papers, such as Wilting et al. (2006), Van De Wiel and Van Wieringen (2007), Huang et al. (2007) and Smeets et al. (2006) perform a hypothesis test on each of their identified regions to compare the CNA from the two groups. Wilting et al. (2006), Van De Wiel and Van Wieringen (2007) and Smeets et al. (2006) apply the Wilcoxon test whereas Huang et al. (2007) use the t -test. In this case a significant p -value will therefore correspond to a region which can be considered to have significantly different CNA between groups of patients - the main aim of our study. They then correct the p -values afterwards accounting for multiple testing.

1. Introduction

In the case of [Huang et al. \(2007\)](#), the t -test is performed on the unsegmented data then the smoothseg algorithm is applied to the t statistics to make comparisons to calculating the t -statistics on the segmented data. It was found that segmenting the t statistics yield a smaller FDR than performing the segmentation on the data. This suggests that performing the test on the unsegmented data is more preferable.

As there currently exists only one method for performing analysis across sequences to identify differential pattern of CNA for a given location between two groups of patients, we aim to create an alternative approach following the hypothesis testing approach done by [Wilting et al. \(2006\)](#), [Van De Wiel and Van Wieringen \(2007\)](#), [Huang et al. \(2007\)](#) and [Smeets et al. \(2006\)](#). We will then compare the method to the comparative KC Smart method - the only method available for this type of analysis.

The next section reviews the comparative KC Smart method, we then motivate the use of hypothesis testing to locate genomic regions which display a significant difference of CNA between groups of patients.

1.7 Comparative KC Smart

At the time of the study, KC Smart was the only package which attempted to look for differences between sequence groups. There are lots of packages which look for things common to one disease (see Section 1.4), but to our knowledge only KCsmart compares groups of patients and attempts to locate differences in the patterns of CNA between sequence groups. The KC Smart Vignette [De Ronde and Klijn \(2013\)](#) describes the implementation of this method on a sample data set included within the R package.

1.7.1 Exploratory Data Analysis

The data is an artificial aCGH data set created by "permuting a BAC data set consisting of 20 samples and introducing an artificial gain" in chromosome 4 ([de Ronde et al., 2019](#)). The data set contains 3268 probe intensities for 20 sequences. Despite the data being artificial, we include the use of it within this thesis and another way of comparing this approach to our methods. Table 1.5 shows the number of probes for each chromosome in this sample data set.

After the removal of chromosomes X and Y, there are 3096 probes for the 20 samples. In this sample data set no missing data is present thus no further data cleaning is required. The data has been split into two separate groups with 10 samples in each group. In the first group of 10 samples there exists a large

Chromosome	Number of Probes
1	210
2	256
3	250
4	208
5	188
6	168
7	190
8	156
9	152
10	146
11	169
12	165
13	99
14	109
15	96
16	104
17	122
18	88
19	75
20	65
21	28
22	52
X	143
Y	29

Table 1.5: The number of probes in each chromosome for the sample data set.

increase in copy number for chromosome 4. The difference of CNA in chromosome 4 between the two groups of samples exist purely to demonstrate the power of the KC Smart method. Figure 1.5 shows the estimated CNA across all probes for a sample from group 1 and group 2 of this sample data set. The increase in CNA in chromosome 4 for a sample in group 1 can clearly be seen in Figure 1.5 (top).

A selection of three probes has is chosen to display the shape of the distributions of estimated CNA within the same sample group in Figure 1.6. The difference between estimated CNA between the two groups of samples in probe 800 is obvious from the histograms. Despite the data being artificial, the shape of the distributions still display an evidence of multi-modality. As multi-modality is a feature we expect to see for these kinds of data sets, it is good that this feature is present in the artificial data set. However, the artificial nature of this data set is made clear in Figure 1.6 (bottom left) as the distribution is similar to that of a normal distribution. For this kind of data it is very rare to see data which can be described as easily using a normal distribution.

1. Introduction

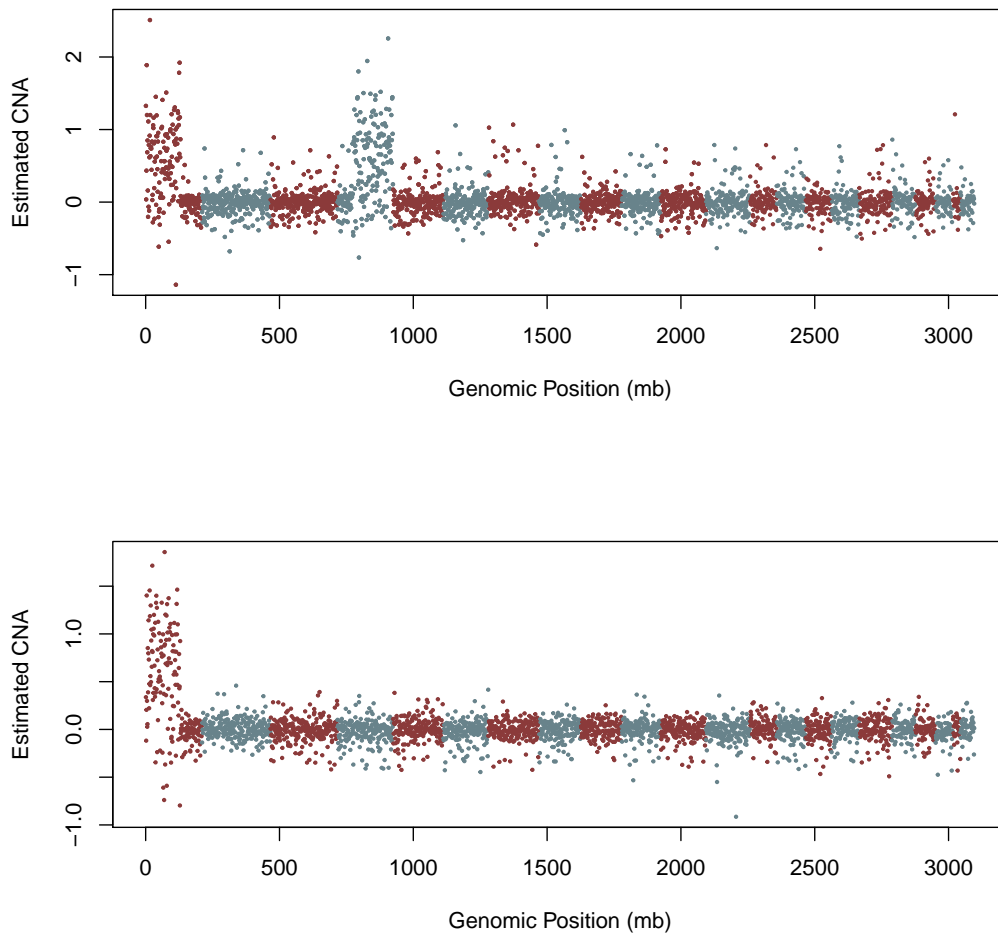


Figure 1.5: The estimated CNA for each probe along the genome for a sample from group 1 (top) and group 2 (bottom). The alternating colouring scheme indicates chromosomes 1–22.

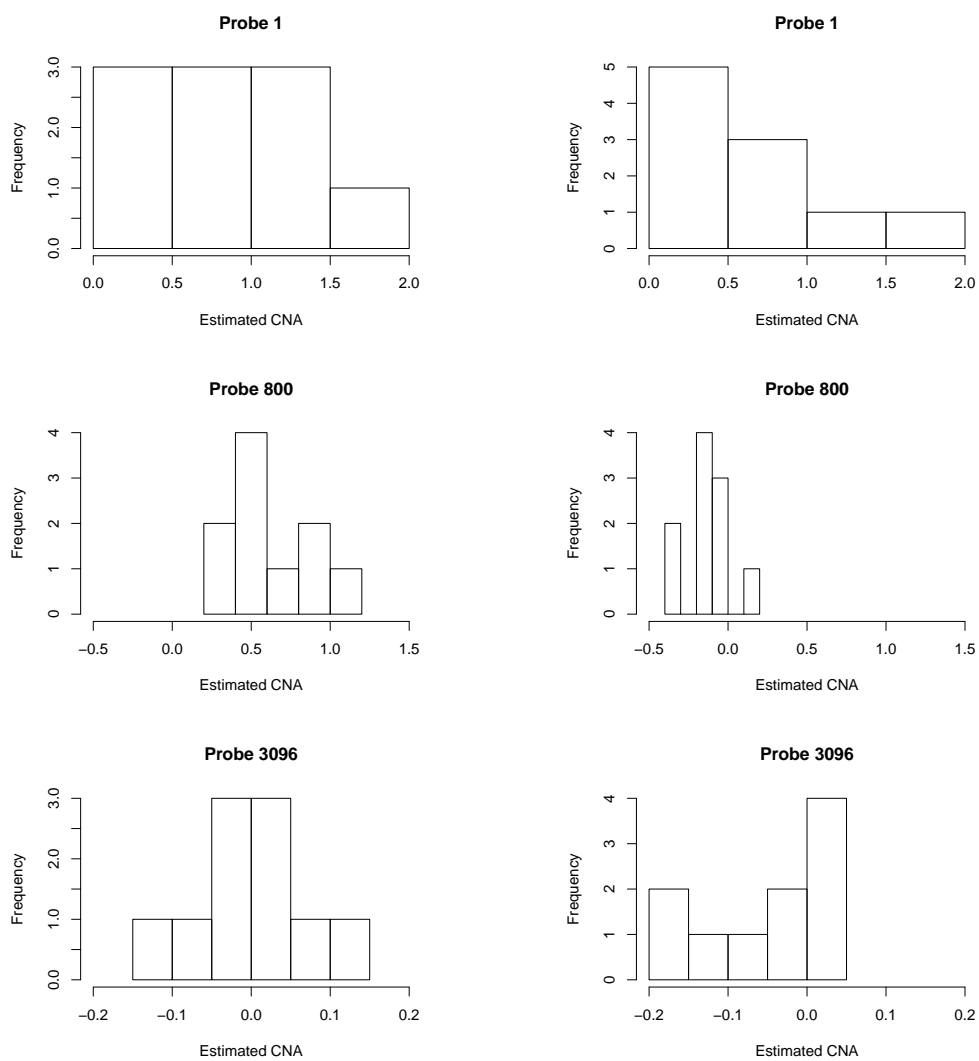


Figure 1.6: Histograms of the estimated CNA across samples from group 1 (left) and group 2 (right) of the sample data set. Probe 1 is taken from chromosome 1, probe 800 is taken from chromosome 4 and probe 3000 is taken from chromosome 22.

1. Introduction

Table 1.6 shows the number of windows which 1) pass the Shapiro-Wilks normality test, 2) have evidence of skewness and 3) have evidence of multi-modality for samples in group 1 and group 2 of the artificial data set respectively.

Evidence of	Group 1	Group 2
Normality	2694	2479
Skewness	584	729
Multi-Modality	2289	2413

Table 1.6: The number of windows which 1) pass the Shapiro-Wilks normality test, 2) have evidence of skewness and 3) have evidence of multi-modality for samples in group 1 and group 2 of the artificial data set respectively.

Table 1.6 shows that even though the majority of samples are passing the Shapiro-Wilks test, there are almost an equivalent amount which have evidence of multi-modality. A possible cause for this is the small number of observations in each sample; each sample has 10 observations. Note also however that in this case there are about three quarters of samples in each group that have evidence of multi-modality, which is not consistent with the lung cancer data set. This therefore shows that many features of a “true” genomic data set are missing from this artificial one. We will still continue to use this data set however to compare methods. Out of the windows which had evidence of multi-modality, Table 1.7 shows the quantity which have i number of peaks, $i = 1, \dots, 9$ for samples in group 1 and group 2 of the artificial data set respectively. To obtain the estimated number of peaks, a clustering algorithm was applied in R.

i	Group 1	Group 2
2	447	542
3	207	199
4	150	171
5	162	145
6	100	116
7	188	161
8	174	166
9	861	913

Table 1.7: The quantity of windows which have i number of peaks, $i = 1, \dots, 9$ for samples in group 1 and group 2 of the artificial data set respectively after applying the clustering algorithm in R.

1.7.2 KC Smart Methodology

In the sample data set provided in the KC Smart vignette, the copy number alterations are collected from a discrete number of probes along the genome. The use of a discrete number of probes is a limitation of the data collection method as CNA can be measured at continuous locations along the genome. Because of this, KC Smart first aims to “smooth” the data using a kernel convolution-based method (Parzen, 1962) to perform locally weighted regression (Atkeson et al., 1997). Locally weighted regression is applied because of the unequal spacing of the probes along the genome. The kernel smoothed estimate (KSE) or ‘KC score’ of the log₂ ratios at an arbitrary position x along the genome is given by the Nadaraya-Watson estimate

$$KSE(x) = \frac{\sum_{M_x} a_i \cdot g_i(x)}{\sum_{M_x} g_i(x)},$$

where a_i is the negative or positive log₂ values for probe i , $g_i(x)$ is the kernel function and M_x is the set of all probes contributing to the (KSE). The kernel function chosen by KC Smart is the flattop Gaussian kernel function, defined by

$$g_i(x) = I_{\{x \leq \mu_{i1}\}} \cdot e^{-\frac{(x-\mu_{i1})^2}{2\sigma^2}} + I_{\{x \leq \mu_{i2}\}} \cdot e^{-\frac{(x-\mu_{i2})^2}{2\sigma^2}} + I_{\{x \in [\mu_{i1}, \mu_{i2}]\}},$$

where the variables μ_{i1} and μ_{i2} represent the mapped genomic start and end points of the probe i and σ is the kernel width. The set M_x contains the probes that lie 4σ to the left and right of the sample point x . Boundary problems are corrected at the chromosome ends and the centromeres by mirroring the probes up to half of the kernel width from the boundary of the chromosomes.

For each sample, the KC score is calculated for each genomic position or probe, i , for positive and negative gains separately. The reason for this separation, as described by Klijn et al. (2008) is “since gains and losses are fundamentally different (only a few copies of a region can be lost, depending on the ploidy of the cell, but many copies can be gained)”. This has been done for the data shown in Figure 1.5 (top), and the kernel smooth estimate is shown in Figure 1.7 for positive and negative gains.

Then, for each genomic position j , the signal-to-noise ratio is calculated:

$$SNR(j) = \frac{\mu_{KC}^1(j) - \mu_{KC}^2(j)}{\sigma_{KC}^{1,2}(j) - f},$$

where $\mu_{KC}^1(j)$ and $\mu_{KC}^2(j)$ represent the mean KC score across all samples in group

1. Introduction

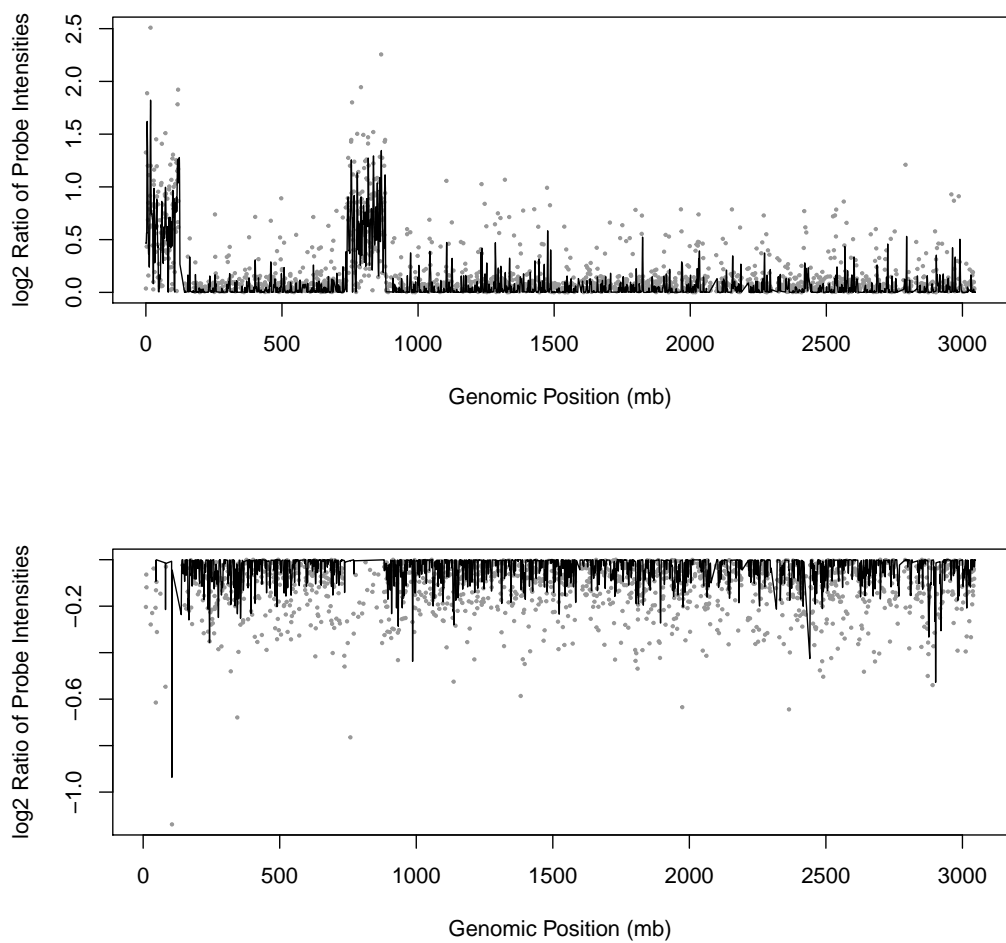


Figure 1.7: For one mouse or sample, the \log_2 ratios of probe intensities plotted against the genomic position of the probes for positive (top) and negative (bottom) gains (grey). The black lines represent the KC score calculated for arbitrary positions along the genome.

1 and group 2 respectively for genomic position j ; $\sigma_{KC}^{1,2}(j)$ represents the pooled variance of the KC scores across all samples for genomic position j and f represents the regularization factor equal to the 95th percentile of the pooled class standard deviation across all genomic positions. [de Ronde et al. \(2010\)](#) states that the “regularization factor prevents small variances from dominating the SNR statistic.”

To identify statistically significant regions of the genome, a permutation based method is adopted. Here, group labels are permuted and the signal-to-noise ratios calculated for each genomic position j . Then a significance threshold can be found depending on the users choice of false discovery rate. Any positions that are above this threshold are then determined to be significant.

1.7.3 Critical Assessment of KC Smart

KC Smart is currently the only known computer package available to compare groups of patients and locate genomic regions which differ significantly in CNA, and because of this it therefore has the advantage of performing unique analysis. However the method contains limitations for which we attempt to resolve.

Firstly, the use of a permutation based method for calculating the significance threshold means that this method will be computationally slow as the sample size and number of probes/windows increase - the method currently takes 3 seconds to analyse 20 samples each with 3096 probes. Another limitation of the KC Smart method is the use of the kernel smoother. Introducing a kernel smoother means introducing a smoothing parameter which could greatly affect the outcome of any analysis. For example for $\sigma = 10^6$ KC Smart is applied to the sample data and Figure 1.8 is produced. Note that the KC Smart package which produces Figure 1.8 does not allow the changing of the x and y axis labels. The y axis simply refers to the signal to noise ratio calculated for each genomic position.

KC Smart produces the table shown in Table 1.8 and displays the locations of significant regions. Here the chromosome is provided as well as the start and end positions in base pairs.

Chromosome	Start Position	End Position
4	48850001	169900001
4	170850001	190500001
6	20400001	21200001
13	39000001	39900001

Table 1.8: A tabular output of KC Smart showing the locations of the significant regions for a kernel width of $\sigma = 10^6$.

1. Introduction

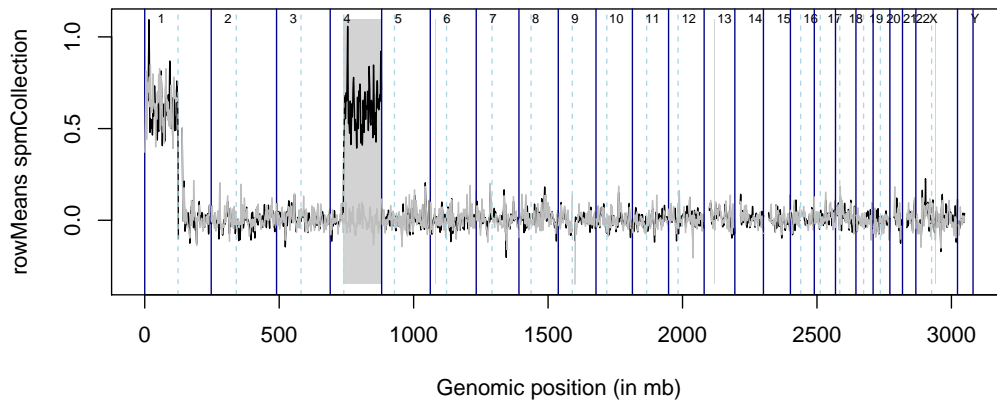


Figure 1.8: The output of the KC Smart analysis with $\sigma = 10^6$ on the sample data.

Now consider increasing the kernel width σ to 10^7 . For this kernel width, Figure 1.9 is produced. It is already clear from this graph that increasing the kernel width produces a much smoother KC score. Table 1.9 displays the locations of significant regions for this kernel width.

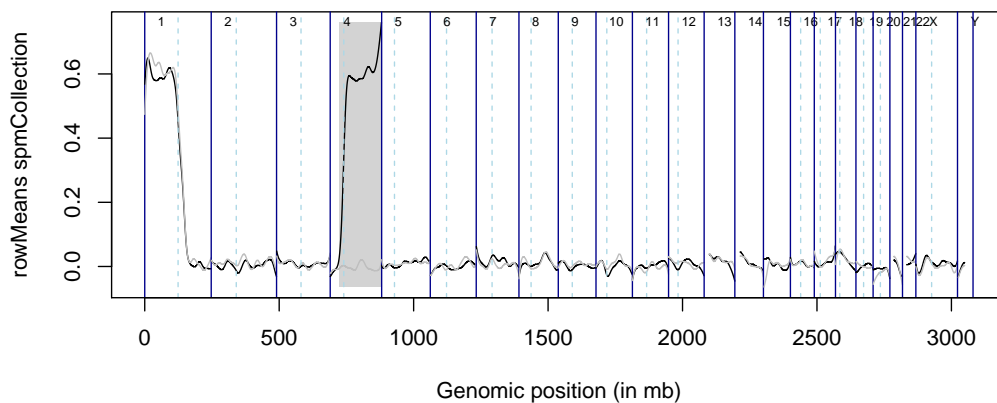


Figure 1.9: The output of the KC Smart analysis with $\sigma = 10^7$ on the sample data.

Chromosome	Start Position	End Position
4	33700001	190500001

Table 1.9: A tabular output of KC Smart showing the locations of the significant regions for a kernel width of $\sigma = 10^7$.

1.8 Hypothesis Testing for Identifying Genomic Regions of Interest

Clearly when the kernel width is changed, KC Smart produces very different outputs. So, which kernel width produces the correct significant regions? The answer to this question is not straightforward, and thus provides a limitation of this method.

We therefore attempt to create an alternative approach for identifying regions which differ significantly in CNA between groups of patients.

1.8 Hypothesis Testing for Identifying Genomic Regions of Interest

Recall that [Van De Wiel and Van Wieringen \(2007\)](#) and [Rueda and Diaz-Uriarte \(2009\)](#) perform hypothesis testing as a way of detecting regions of the genome which have significantly different CNA between groups of patients. In this section, we motivate the reason for using hypothesis testing in this way.

Identifying genomic regions of interest is equivalent to locating which windows have significantly different CNA between patients with one subtype of cancer compared to the other. Clearly, [Figure 1.1](#) from [Section 1.5.2](#) shows some differences and similarities between the estimated CNA in each window. Specifically there exists a large difference in estimated CNA in chromosome 3 between the two types of lung cancer. Many studies have been done on these subtypes of lung cancer and have found that a section of chromosome 3 experiences significant gains in patients with squamous carcinoma type lung cancer ([Björkqvist et al., 1998](#)), ([Wang et al., 2013](#)) and ([van Boerdonk et al., 2011](#)). Because of this, we would indeed expect to see a significant difference of CNA in chromosome 3 between each group of patients. However, for our case it is impossible to tell simply by looking at two patients whether the differences are global or whether they occur for these two patients alone.

Instead consider [Figure 1.10](#), the plots show the histograms of window 4170 with location at 132.75 Mbp in chromosome 3. Firstly note the negative CNA in the histogram across patients with adenocarcinoma type lung cancer. These negative values are a consequence of the normalisation process not performing correctly. Whilst this means there is a chance that any results we provide may be inaccurate, the methods we provide are unaffected by the inaccuracy of the normalisation step.

If indeed chromosome 3 is a significant region, then we would expect the two distributions in [Figure 1.10](#) to be significantly different. Indeed, it can be seen that this is the case. A simple two-sample t -test on these data produces a p -value

1. Introduction

of 1.42×10^{-9} , suggesting that mean of these two distributions are significantly different. We could apply a two-sample t -test on each individual region across the genome to determine the regions of the genome that differ in mean estimated CNA between the two types of lung cancer. However, Figure 1.10 shows that the distribution of CNA for patients with adenocarcinoma type lung cancer differs from the distribution of CNA for patients with squamous carcinoma type lung cancer in not only the mean, but also the variance, skewness and even multi-modality. Thus, we require a two-sample test that is sensitive at identifying any differences, including multi-modality, between the two samples.

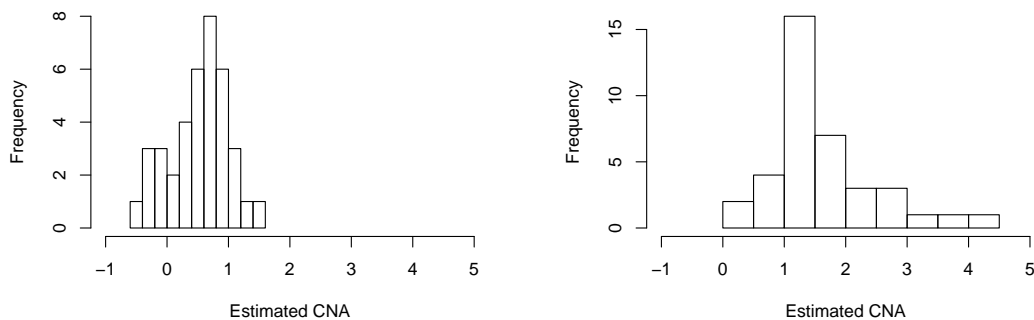


Figure 1.10: The estimated CNA in window 4170 across patients with adenocarcinoma type lung cancer (left) and squamous carcinoma type lung cancer (right). Both histograms have been plotted on the same x-axis scale to enable easy comparison.

Note that [Van De Wiel and Van Wieringen \(2007\)](#) and [Rueda and Diaz-Uriarte \(2009\)](#) use the Wilcoxon test and the absolute value of the difference in mean probability as a test statistic to test the difference in CNA between the two groups respectively. It is likely that these two tests are not sensitive enough to identify all differences between the groups. We aim to rectify this by obtaining a more suitable hypothesis test. For each region, we choose to test the null hypothesis that the distributions of estimated CNA for each type of lung cancer has a common cumulative distribution function. In a statistical sense, given X_1, \dots, X_n and Y_1, \dots, Y_m identically distributed random variables with distribution functions $F(x)$ and $G(y)$ respectively, the null hypothesis

$$H_0 : F(x) \equiv G(y) \tag{1.2}$$

is considered. By testing this null hypothesis, any difference in mean, variance,

skewness or multi-modality should be identified. Establishing a suitable two-sample statistic to test this null hypothesis is our first objective in this thesis.

One might question why hypothesis testing might be an approach worth researching to locate genomic regions which have a significant difference in CNA between two groups of patients. The advantage of such a method is that tests can be performed on regions which are only a couple of kilo base pairs long. In the studies done by Björkqvist et al. (1998), Wang et al. (2013) and van Boerdonk et al. (2011), large regions of significance were identified. If tests are performed which can only identify large regions, then perhaps smaller regions that have a significant difference in CNA between the two groups are being missed. By producing a method which uses hypothesis testing, we therefore hope to identify further smaller regions within the genome that have a significant difference in CNA between the groups.

1.9 Thesis Overview

Throughout this thesis we work towards creating an efficient and attractive method for identifying regions of interest between subtypes of cancer for the purpose of classifying patients on their subtype. We specifically work to identify regions of interest between patients with adenocarcinoma type lung cancer and squamous carcinoma type lung cancer. We establish a suitable test statistic to test the null hypothesis defined in Equation (1.2). We require the test statistic to be sensitive at identifying differences in the mean, variance, skewness and multi-modality between the two samples. We wish the application of the test on each window of the genome to be efficient and accurate.

We begin in Chapter 2 by investigating and critiquing current test statistics to be used to test the null hypothesis in Equation (1.2). We initially consider parametric tests before concluding that a non-parametric test may be more suitable. We investigate current well-known non-parametric tests as well as other less known tests. Out of the most popular non-parametric tests, we show that the Cramer test (Baringhaus and Franz, 2004) is the most preferable.

Properties of the Cramer test statistic are explored in Chapter 3 and we provide formulas for the exact calculation of the first four moments of the test statistic when the distribution of the data is unknown. We also provide results for the expectation and variance of the test statistic for various known common distributions of the data.

1. Introduction

To ensure the test is efficient, in Chapter 4 we look into alternative methods to estimate the p -value of the test. We discuss current resampling methods for obtaining the p -value including the method described by [Baringhaus and Franz \(2004\)](#). We then investigate empirical approximations of the null distribution using the method of moments by equating the parameters of the chosen approximate distribution to the formulae for the first four moments from Chapter 3. We consider two and three parameter approximations as well as transformations and the extreme value theorem.

In Chapter 5, we consider the application our method for identifying regions of interest. We firstly investigate the computational consideration of the method and discuss how we ensured fast calculations. We compare our method to calculate the p -value against other methods both in terms of speed and accuracy. We also compare the choice of test statistic to other well-known powerful statistics. Here we aim to show that the Cramer test was the most suitable choice of test statistic for our purposes. We also compare our method as a whole to KC Smart - the current method used by oncologists for identifying genomic regions of interest. Finally, we provide the results of applying our method to all windows in the lung cancer data set. We hope to prove that the results and conclusions using our method are similar to the conclusions obtained by [Belvedere et al. \(2012\)](#).

For these type of data, there usually exists high correlation between windows, especially windows which are adjacent or close to each other in the genome. To understand more about the correlation structures and show that high correlation exists in the lung cancer data set, Chapter 6 explores these correlation structures. Modelling the correlation structures is also investigated in this chapter which can be used for prediction purposes.

Given that we perform the Cramer test simultaneously on each of the correlated windows, we have a multiplicity problem that cannot be solved using only correction procedures like Bonferroni, which rely on the independence of all tests. Chapter 7 investigates methods to correct for multiplicity when the tests are dependent. We firstly implement an algorithm described by [Dudbridge and Gusnanto \(2008\)](#) to estimate the effective number of independent tests, which can be used in multiplicity correction procedures like Bonferroni. However, we find that methods to calculate the effective number of independent tests are slow and computationally expensive, thus we consider Fisher's combined probability test for testing the significance of a group of correlated p -values. For this, we extend the work done by [Brown \(1975\)](#) and [Kost and McDermott \(2002\)](#) to allow the method to be used when the distribution of the null distribution is unknown.

To conclude, Chapter 8 summarises our research and contributions. We also outline the potential for future work.

Chapter 2

Identifying a Two-Sample Test to Locate Genomic Regions of Interest

2.1 Introduction

The main purpose of our study is to locate the regions of the genome that differ significantly in CNA between groups or subtypes of cancer. We have seen in Section 1.8 that one method in which to do this is to apply a hypothesis test on each window or location of the genome independently. The more patients in the study will ensure a more accurate test per window or genomic location. The question however, is which hypothesis test is suitable for this task? [Wilting et al. \(2006\)](#), [Van De Wiel and Van Wieringen \(2007\)](#) and [Smeets et al. \(2006\)](#) choose to perform a hypothesis test using a rank based test, namely the Wilcoxon test, but we find that these tests are not powerful enough to pick up certain differences in the distribution. Therefore we aim to identify a hypothesis test which is sensitive at identifying the differences in the mean, variance, skewness and multi-modality. Such a test has not yet been identified and will ensure that any regions which differ significantly between groups of patients are identified. Using tests which fail to identify differences in certain aspects could potentially miss highly important regions of the genome which could help classify future patients on their subtype of cancer. Therefore it is vital that we choose a suitable hypothesis test.

In this chapter we perform a literature review on various choices for the hypothesis test and provide disadvantages of using each, before finding a suitable test which is sensitive at identifying all differences within the data. We provide the relevant derivations of the test statistics for each hypothesis test considered so

2. Identifying a Two-Sample Test to Locate Genomic Regions of Interest

the comparisons between each one can be clearly seen. In some of the hypothesis tests considered, the test statistics are very similar thus we believe it is important to see how each one compares to the others.

Given X_1, \dots, X_n and Y_1, \dots, Y_m identically distributed random variables with distribution functions $F(t)$ and $G(t)$ respectively, recall that the null hypothesis we wish to test is

$$H_0 : F(t) \equiv G(t), \quad (2.1)$$

versus the alternative hypothesis

$$H_1 : F(t) \not\equiv G(t). \quad (2.2)$$

2.2 Parametric Tests

Parametric tests rely strongly on a distributional assumption of the data or the central limit theorem (Heyde, 2014). For example, the t -test and F -test rely on the assumption that the summary statistics are normally distributed. Recall Figure 1.2 from Section 1.5.2. These histograms not only show that the distributions of CNA across patients per region for the lung cancer data set are not normally distributed, but that there exists multi-modality within the data. Because of this, it is reasonable to conclude that a parametric test would not be suitable in this case, which includes the use of a t -test and F -test. Even if one was to consider a transformation of the data, no transformations exist which would completely remove the multi-modality within the data. As well as this if such a transformation did exist, the important features of the data would be lost.

Balkin and Mallows (2001) considers an asymmetric, skew-adjusted two-sample t test which is considered as a potential choice in Appendix A. We omit this work from the main text as we have concluded that a parametric test would be unsuitable for our purposes.

2.3 Two Sample Tests Based on Empirical Cumulative Distribution Functions

The advantage of using a non-parametric test is that no distributional assumptions are placed on the data. As is seen in Figure 1.2, the shape of the distributions vary greatly between windows, with some distributions displaying a multimodal quality. The ability to perform a hypothesis test on the data without needing to place any distributional assumptions on the data will thus be ideal for our purposes.

2.3 Two Sample Tests Based on Empirical Cumulative Distribution Functions

Define $F_n(t)$ and $G_m(t)$ to be the empirical cumulative distribution functions of the samples respectively where

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq t}$$
$$G_m(t) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i \leq t}.$$

2.3.1 Kolmogorov-Smirnov

Perhaps one of the most well known non-parametric tests is the Kolmogorov-Smirnov two sample test (Kolmogorov, 1933), (Smirnov, 1939). The test statistic is defined as

$$D_{n,m} = \sup_t |F_n(t) - G_m(t)|$$

where “sup” is the supremum function. The test statistic calculates the maximum vertical distance between the two empirical cumulative distribution curves and the null hypothesis is rejected if the test statistic is larger than a critical value. Massey Jr (1951) provides a table of critical values. As an illustration, consider $X \sim N(0,1)$ and $Y \sim N(1,1)$ and sample $n = m = 100$ observations from each distribution. The empirical cumulative distribution functions of the two samples are shown in Figure 2.1.

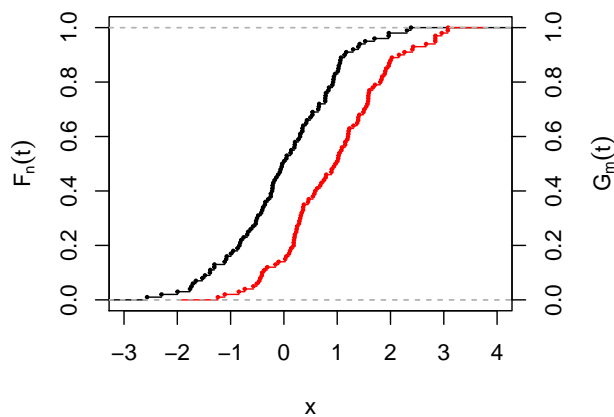


Figure 2.1: The empirical cumulative distribution functions of two samples drawn from distributions X and Y such that $X \sim N(0,1)$ and $Y \sim N(1,1)$ and $n = m = 100$.

The maximum distance between the two ECDF curves shown in Figure 2.1 is calculated to be 0.37 and after performing the Kolmogorov-Smirnov test on the

2. Identifying a Two-Sample Test to Locate Genomic Regions of Interest

two samples, H_0 is rejected at 5% level with a p -value of 2.27×10^{-6} . Now consider $X \sim N(0,2)$ and $Y \sim N(0,1)$ and sample $n = m = 100$ observations from each distribution. The empirical cumulative distribution functions of the two samples are shown in Figure 2.2.

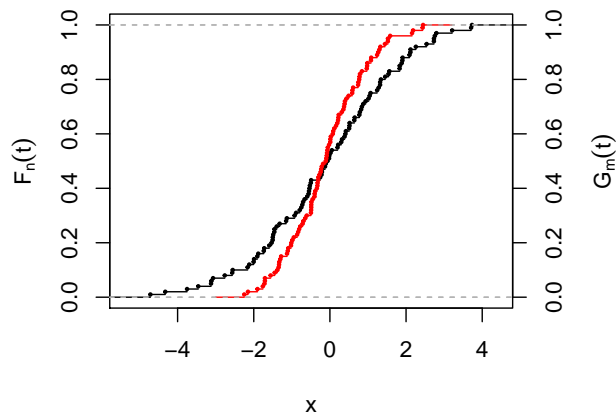


Figure 2.2: The empirical cumulative distribution functions of two samples drawn from distributions X and Y such that $X \sim N(0,2)$ and $Y \sim N(0,1)$ and $n = m = 100$.

The maximum distance between the two ECDF curves shown in Figure 2.2 is 0.17 and, after performing the Kolmogorov-Smirnov test on the two samples, H_0 is not rejected at 5% level with a p -value of 0.11. Clearly, the Kolmogorov-Smirnov is more sensitive at picking up differences when there exists a change in mean compared to a change in variance. Not only this, but the Kolmogorov-Smirnov test performs poorly when the sample size for n and m is small. Consider again $X \sim N(0,1)$ and $Y \sim N(1,1)$ and sample $n = m = 40$ observations from each distribution. The empirical cumulative distribution functions of the two samples are shown in Figure 2.3.

The maximum distance between the two ECDF curves shown in Figure 2.3 is 0.33 and after performing the Kolmogorov-Smirnov test on the two samples, the p -value is 0.029. Thus at 5% level H_0 is rejected, but at 1% level H_0 is not rejected. Compare the p -value obtained here to the p -value obtained when $n = m = 100$, clearly when $n = m$ is larger the null hypothesis is rejected at a smaller significance level.

It is important that the hypothesis test chosen for comparing the distribution of copy number alterations is sensitive enough to identify differences in the mean, variance and skewness etc including when sample sizes are small. If the

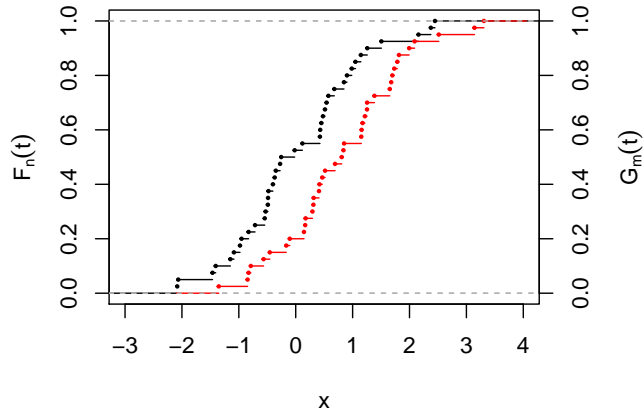


Figure 2.3: The empirical cumulative distribution functions of two samples drawn from distributions X and Y such that $X \sim N(0,1)$ and $Y \sim N(1,1)$ and $n = m = 40$.

Kolmogorov-Smirnov test is unable to identify differences in the variance successfully, there may be some windows which are wrongly concluded to be insignificant. Not only this, but for a sample size of 40, the Kolmogorov-Smirnov test has smaller sensitivity when identifying a difference in the mean. In the lung cancer data set, distributions of samples with only $n = m = 38$ observations will be compared, thus clearly the Kolmogorov-Smirnov test is an unsuitable test for this data.

2.3.2 Cramer-von Mises and Anderson-Darling

The general Cramer-von Mises one-sample test statistic is defined by

$$W^2 = n \int_{-\infty}^{\infty} \lambda(F(t))(F_n(t) - F(t))^2 dF(t), \quad (2.3)$$

where $\lambda(F(t))$ is a chosen weight function. The fundamental ideas behind this test were developed by [Cramér \(1928\)](#), [Von Mises \(1931\)](#) and [Smirnov \(1936\)](#). By setting $\lambda(F(t)) = 1$, W^2 becomes the one-sample Cramer-von Mises (CvM) test statistic

$$w^2 = n \int_{-\infty}^{\infty} (F_n(t) - F(t))^2 dF(t), \quad (2.4)$$

for which [Csorgo and Faraway \(1996\)](#) describe the exact and asymptotic properties. One such property, which can be shown by results in [Smirnov \(1936\)](#) and [Götze \(1979\)](#), is that $\lim_{n \rightarrow \infty} F_n(t) = F(t)$. This property therefore claims that for large n , the cumulative distribution function can be approximated by the empirical

2. Identifying a Two-Sample Test to Locate Genomic Regions of Interest

cumulative distribution function.

Using the fact that $F_n(t)$ can be considered as a binomial random variable, [Anderson and Darling \(1954\)](#) show that $E[F_n(t)] = F(t)$ and $\text{Var}[F_n(t)] = F(t)(1 - F(t))$. Because of this, [Anderson and Darling \(1954\)](#) propose to modify the statistic w^2 by setting the weight $\lambda(F(t)) = (F(t)(1 - F(t)))^{-1}$. In this formulation, they want to “equalize the sampling error over the entire range of t by weighting the deviation by the reciprocal of the variance”. [Anderson and Darling \(1954\)](#) state that another consequence of choosing $\lambda(F(t)) = (F(t)(1 - F(t)))^{-1}$ is that the test statistic weights the tails of the distribution more than its centre. This can be considered an advantage of this test because it is more sensitive at identifying differences in the tails of the distributions where small but important deviations can occur. By choosing this weight, Equation (2.3) becomes the one-sample Anderson-Darling test statistic,

$$A_n^2 = n \int_{-\infty}^{\infty} \frac{(F_n(t) - F(t))^2}{F(t)(1 - F(t))} dF(t). \quad (2.5)$$

[Anderson \(1962\)](#) subsequently adapts the one-sample Cramer-von Mises test (2.4) into a two-sample version

$$T = \frac{nm}{n+m} \int_{-\infty}^{\infty} \lambda(H_{n+m}(t))(F_n(t) - G_m(t))^2 dH_{n+m}(t), \quad (2.6)$$

where $H_{n+m}(t)$ is the ECDF of the pooled data. For $n, m \leq 7$, [Anderson \(1962\)](#) obtains the sample distribution of T . Following that, [Burr \(1964\)](#) obtains the sample distribution for sample pairs such that $n, m \geq 4$ and $n + m \leq 17$. The asymptotics of this statistic are also considered in [Lehmann \(1951\)](#), [Rosenblatt et al. \(1952\)](#), [Fisz \(1960\)](#) and [Darling \(1957\)](#).

[Pettitt \(1976\)](#) also modifies (2.5) into a two sample version

$$A_{nm}^2 = \frac{nm}{n+m} \int_{-\infty}^{\infty} \frac{(F_n(t) - G_m(t))^2}{H_{n+m}(t)(1 - H_{n+m}(t))} dH_{n+m}(t). \quad (2.7)$$

For the Anderson-Darling (AD) test statistic (2.7), [Pettitt \(1976\)](#) calculates its asymptotic distribution. However, [Pettitt \(1976\)](#) also indicates that finding moments higher than the mean explicitly is “impossible”, which is a disadvantage of the Anderson-Darling test. [Baumgartner et al. \(1998\)](#) also considers a test statistic similar to the Anderson-Darling test statistic and compares the power of the test to the Kolmogorov-Smirnov test, the Cramer-von Mises test and the Wilcoxon test.

2.3 Two Sample Tests Based on Empirical Cumulative Distribution Functions

Whilst it can be shown that the Cramer-von Mises and Anderson-Darling tests perform well at identifying differences between distributions which differ in mean, variance and skewness etc. the two tests are less effective when dealing with multi-modal data. Consider sampling $n = m = 100$ observations from X and Y such that $X \sim N(1, 0.25)$ and Y which follows a mixture of normals distribution with probability density function $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$. The distribution of Y was chosen to reflect the shape of the distribution of estimated CNA across groups of patients per region of the genome. Here we have chosen the distribution to have a larger peak centred at 1 with a smaller peak centred at 3. This is perhaps a more extreme case of the data we expect to see, but as can be seen in Figure 1.2 from Section 1.5.2 observing CNA of more than 2 away from the most common CNA is possible. Therefore it is vital that any hypothesis test we choose can adequately deal with data of this type.

Figure 2.4 shows the histograms of the two samples from X and Y . The number of breaks was chosen to be about 10 and the two histograms are plotted on the same scale for easy comparison. The empirical cumulative distribution functions of the two samples are shown in Figure 2.5.

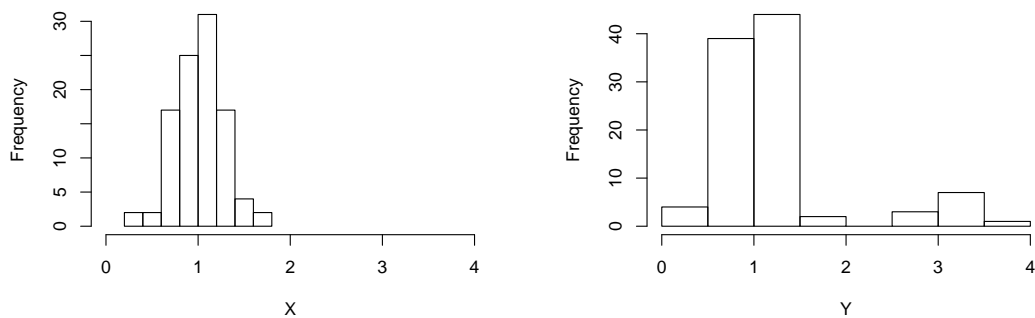


Figure 2.4: The histograms of the samples X (left) and Y (right).

Performing the CvM test on the data gives a p -value of 0.502 and performing the AD test gives a p -value of 0.231. Thus both tests conclude that the two distributions of X and Y are equivalent despite the graphs in Figures 2.4 and 2.5 clearly showing large differences between the distributions. A lack of sensitivity when dealing with multi-modal data is a major disadvantage of the CVM and AD tests.

2. Identifying a Two-Sample Test to Locate Genomic Regions of Interest

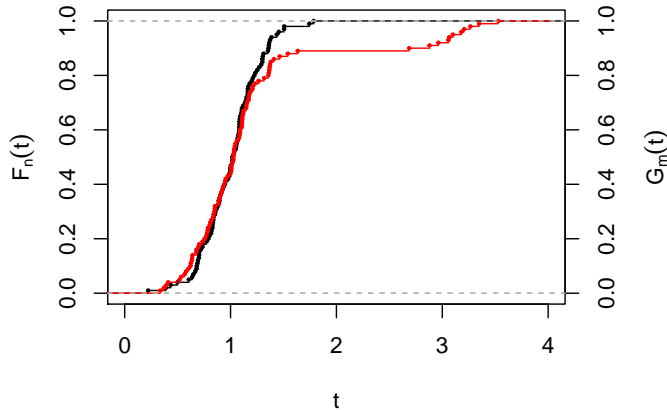


Figure 2.5: The empirical cumulative distribution functions of two samples drawn from distributions X and Y such that $X \sim N(1, 0.25)$ and Y which follows a mixture of normals distribution with probability density function $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$. and $n = m = 100$.

2.3.3 Cramer Test

[Baringhaus and Franz \(2004\)](#) presents the following test in one dimension;

$$T_{n,m} = \frac{nm}{n+m} \int_{-\infty}^{\infty} (F_n(t) - G_m(t))^2 dt, \quad (2.8)$$

which can be considered a modification of the two sample Cramer-von Mises test (2.6) and is thus named the Cramer test for this reason. [Baringhaus and Franz \(2004\)](#) states that the limiting distribution of $T_{n,m}$ as $n, m \rightarrow \infty$ is

$$\int_{-\infty}^{\infty} B^2(H(t)) dt, \quad (2.9)$$

where $B(u)$, $0 \leq u \leq 1$ is the classical Brownian bridge. Note that the limiting distribution depends on $H(t)$, which consequently means the test is not completely non-parametric. To calculate a p -value, [Baringhaus and Franz \(2004\)](#) suggests using a traditional bootstrap or permutation estimate of the distribution of $T_{n,m}$ or the limiting distribution defined in Equation (2.9). Note that using either a bootstrap approach or the permutation method for calculating the p -value can be slow for large n . In our application, this is a practical limitation and is considered further in Chapter 4.

Recall that the two sample Cramer-von Mises test (2.6) uses a weight function $\lambda(u)$ for $0 \leq u \leq 1$. Consider instead a weight function $\lambda(t)$ for $-\infty \leq t \leq \infty$ and

2.3 Two Sample Tests Based on Empirical Cumulative Distribution Functions

let $\lambda(t) = \frac{1}{h(t)}$, where $h(t)$ is the common probability density function. With this modified weight function, the Cramer-von Mises test statistic (2.6) becomes

$$\frac{nm}{n+m} \int_{-\infty}^{\infty} \frac{1}{h(t)} (F_n(t) - G_m(t))^2 dH_{n+m}(t). \quad (2.10)$$

As

$$\frac{dH(t)}{dt} = h(t),$$

Equation (2.10) is equivalent to Equation (2.8). Recall that the weight function for the Anderson-Darling test is of the form

$$\lambda(u) = \frac{1}{u(1-u)}$$

for $0 \leq u \leq 1$. With this choice of weight function it is easy to see that $\lambda(u) \rightarrow \infty$ when $u = 0$ or $u = 1$. Thus for $u = H(t)$, more weight will be given when $H(t)$ is close to 0 or 1, which occurs at the tails of the pooled data.

Using the weight function $\lambda(t) = \frac{1}{h(t)}$ will put more weight on the test statistic when $h(t) \approx 0$. In particular, $h(t) \approx 0$ will not only occur at the tails of the distribution but could also occur more often if the data is multi-modal. Thus we can argue that the Cramer test statistic will be more sensitive at detecting differences when the data is multi-modal.

To test this, consider once again $X \sim N(1, 0.25)$ and Y which follows a mixture of normals distribution with probability density function $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$. Table 2.1 shows the results of a simulation designed to show the sensitivity of the Cramer test against the Cramer-von Mises and Anderson-Darling tests when dealing with multi-modal data with $n = m = 100$. The table shows the number of times each test statistic rejects the null hypothesis that the distributions of X and Y are equivalent at the 5% significance level out of 100 simulated datasets.

Test	Number of rejected H_0
Cramer-von Mises Test	34
Anderson-Darling Test	52
Cramer Test	100

Table 2.1: Number of rejections of H_0 out of 100 simulated datasets in which the Cramer-von Mises test, the Anderson-Darling test and the Cramer test is performed on 100 observations from X and Y such that $X \sim N(1, 0.25)$ and Y which follows a mixture of normals distribution with probability density function $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$.

Clearly, Table 2.1 shows that the Cramer-von Mises test fails to detect any

2. Identifying a Two-Sample Test to Locate Genomic Regions of Interest

significant differences between the distributions of X and Y about a third of the time. The Anderson-Darling test performs better than the Cramer-von Mises test, but still fails to detect any significant differences about a half of the time. The Cramer test is able to identify significant differences 100% of the time. This therefore shows that if the data exhibits multi-modality, like for the lung cancer data set, the Cramer test is the most suitable.

The integration variable plays a major role in why the Cramer test is more powerful in the case of multi-modal data. To see why, consider the following function

$$\eta(t) = \frac{(F_n(t) - G_m(t))^2}{H_{n+m}(t)(1 - H_{n+m}(t))},$$

which is the integrand of the Anderson-Darling test statistic. To show how the integration variable affects the test statistics, Figure 2.6 displays two plots, namely $\eta(t)$ plotted against t' where t' is the linearly transformed version of t so that $t' \in [0, 1]$, and $\eta(t)$ plotted against $H_{n+m}(t)$.

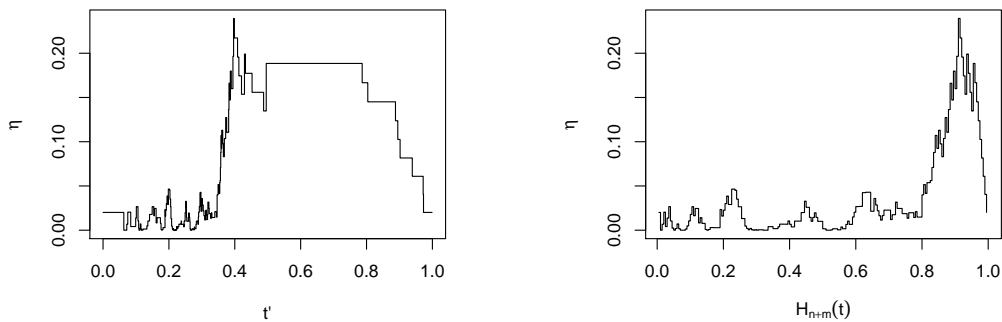


Figure 2.6: The graphs of $\eta(t)$ plotted against t' where t' is the linearly transformed version of t so that $t' \in [0, 1]$ (left), and $\eta(t)$ plotted against $H_{n+m}(t)$ (right).

The area under the curve in Figure 2.6 (right) is the test statistic for the Anderson-Darling test. It is clear that the area under the curve in Figure 2.6 (left) is larger than the area under the curve in the right graph. This shows that the test statistic will be larger when integrating with respect to t and thus suggests that the Cramer test is more likely to detect significant differences between the distributions of X and Y .

Note that the Anderson-Darling test statistic is divided by $H_{n+m}(t)(1 - H_{n+m}(t))$ and $T_{n,m}$ is not. Whilst it may be preferable to divide the integrand of $T_{n,m}$ by $H_{n+m}(t)(1 - H_{n+m}(t))$, this sacrifices the ability to calculate exact formulas for the moments of the test statistic.

2.4 Comparison to Current Literature

As has already been mentioned [Wilting et al. \(2006\)](#), [Van De Wiel and Van Wieringen \(2007\)](#) and [Smeets et al. \(2006\)](#) choose to perform a hypothesis test using a rank based test, namely the Wilcoxon test [Wilcoxon \(1945\)](#). Like the t -test, the Wilcoxon test is very powerful in locating differences in the mean but fails to identify differences when they occur in the variance or multi-modality. As before, due to the complex nature of the data, we require a test which can locate differences based on the mean, variance, skewness and multi-modality so whilst performing the Wilcoxon test on two groups CNA data will identify some genomic regions which have significantly different CNA between the groups many regions will be overlooked. We therefore hope that the Cramer test will solve this problem.

2.5 Further Two Sample Tests

[Fernández et al. \(2008\)](#) describes a two-sample test based on the empirical characteristic functions. The test statistic is defined as

$$D_{n,m} = \int_0^1 |C_n(t) - C_m(t)|^2 dB(t) \quad (2.11)$$

where $C_n(t)$ and $C_m(t)$ are the empirical characteristic functions defined by

$$C_n(t) = \frac{1}{n} \sum_{j=1}^n e^{itX_j}; \quad C_m(t) = \frac{1}{m} \sum_{j=1}^m e^{itY_j}.$$

In Equation (2.11), $B(t)$ is a chosen distribution function on \mathbb{R}^d . The distribution function $B(t)$ can be chosen to maximise the power of the test. To calculate the p -value [Fernández et al. \(2008\)](#) uses a bootstrap approximation to approximate the null distribution. [Fernández et al. \(2008\)](#) shows that the test is powerful when observing a change in scale and location. [Jiménez-Gamero et al. \(2017\)](#) considers the test statistic defined by Equation 2.11 and provides fast algorithms using a weighted bootstrap approach to calculate a p -value in 0.01 seconds for dimension $d = 1$ and $n = m = 20$. A massive advantage of this test over the Cramer-von Mises and Anderson-Darling is that it can be easily generalised to any dimension. However, as the Cramer test can also be generalised to any dimension and avoids exponentials, we find it more preferable.

Many other methods to solve the two-sample problem have been studied in the literature. For example, [Gretton et al. \(2007\)](#) suggests using a kernel approach

2. Identifying a Two-Sample Test to Locate Genomic Regions of Interest

for the two-sample problem and [Cao and Van Keilegom \(2006\)](#) uses an empirical likelihood ratio test using kernel density estimates. Aside from the Wilcoxon test already mentioned, other rank based tests are also well studied in the literature, e.g. [Mann and Whitney \(1947\)](#) suggests the Mann-Whitney test. [Curry et al. \(2018\)](#) also introduce a new rank-based Cramer-von Mises type test in which the power of their test is compared to the Wilcoxon test, the Cramer test (2.8) and the test described in [Fernández et al. \(2008\)](#) (2.11).

2.6 Discussion

This chapter investigates current parametric and non-parametric tests for the purpose of locating genomic regions where the distribution of CNA is significantly different between subtypes of cancer. After exploring the shape of the data it was clear that parametric tests are unsuitable due to the incorrect normality assumption. [Balkin and Malloys \(2001\)](#) try to address this issue by introducing an asymmetric, skew-adjusted t -test. However, when taking the variance and skewness of both distributions into account, we have found that when $n = m$, the test is generalised to the standard Welch's t -test. Aside from the incorrect normality assumption, another disadvantage of most parametric tests is they are only sensitive in identifying a single difference between the distributions, i.e. in the mean or variance. Whilst the Cramer-von Mises test and the Anderson-Darling test sufficiently deal with this issue, in the case of multi-modality the two tests are less effective. As it has been identified that multi-modality is a common feature of the lung cancer data set, having a test which can deal with this feature is of importance. The Cramer test is shown to be more sensitive when the data is multi-modal, the reason for this is the choice of weight function and thus the integration variable. We conclude, therefore, that a suitable choice of hypothesis test is the one-dimensional Cramer test. However, as we have identified, a limitation of this test is the computational burden in calculating the p -value. We investigate this in more detail in Chapter 4.

Chapter 3

Properties of the Cramer Test

3.1 Introduction

In Chapter 2, we showed that multi-modality was a common feature of the lung cancer data set. We also showed that when the data is multi-modal, the Cramer test (Baringhaus and Franz, 2004) is more sensitive in identifying differences compared to the Cramer-von Mises and Anderson-Darling tests.

In this chapter, we calculate and present formulas to calculate the moments of the Cramer test statistic up to the fourth moment which has not been previously calculated. Additionally, we provide the first two moments for a selection of distributions. We also consider properties of the Cramer test statistic when a transformation of the data is performed.

3.2 Cramer Test Statistic

Consider two samples of data, x_1, \dots, x_n and y_1, \dots, y_m , where we can assume without loss of generality that $n > m$, from distributions f and g respectively. If we wish to test the null hypothesis $H_0 : F \equiv G$, we can use the Cramer test with test statistic $T_{n,m}$ defined by

$$T_{n,m} = \frac{nm}{n+m} \int_{-\infty}^{\infty} (F_n(t) - G_m(t))^2 dt. \quad (3.1)$$

3.3 Moments

Exact formulas for the moments of the Cramer test statistic can be calculated. Proofs are given either in the text or appendices D, E and F. In the following

3. Properties of the Cramer Test

equations the cumulative distribution function $H(t)$ is used to denote the distribution when H_0 is true. As in most cases $H(t)$ is unknown, we will later use the empirical cumulative distribution function $H_{n+m}(t)$ to approximate $H(t)$. Recall a property of the Cramer-von Mises test statistic, w^2 (2.4), mentioned by Csorgo and Faraway (1996), which states that, for large n , the cumulative distribution function can be approximated by the empirical cumulative distribution function. Thus for large $n + m$, H_{n+m} provides an accurate approximation of $H(t)$.

3.3.1 Expectation

The expectation of $T_{n,m}$ under the null hypothesis is given by

$$\mathbb{E}[T_{n,m}|F = G] = \int_{-\infty}^{\infty} H(t)(1 - H(t)) dt. \quad (3.2)$$

Proof By substituting the empirical cumulative distribution functions $F_n(t)$ and $G_m(t)$ with the relevant summations, the Cramer test statistic becomes

$$\begin{aligned} T_{n,m} &= \frac{nm}{n+m} \int_{-\infty}^{\infty} \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}} - \frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m \mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{y_j \leq t\}} \\ &\quad + \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \mathbb{1}_{\{y_i \leq t\}} \mathbb{1}_{\{y_j \leq t\}} dt. \end{aligned}$$

Taking the expectation of $T_{n,m}$ and using Fubini's theorem (Fubini, 1907) which gives the conditions under which integrals can be switched, gives

$$\begin{aligned} \mathbb{E}[T_{n,m}] &= \frac{nm}{n+m} \int_{-\infty}^{\infty} \left(\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}}] \right. \\ &\quad \left. - \frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{y_j \leq t\}}] + \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \mathbb{E} [\mathbb{1}_{\{y_i \leq t\}} \mathbb{1}_{\{y_j \leq t\}}] \right) dt \\ &= \frac{nm}{n+m} \int_{-\infty}^{\infty} \left(\frac{n-1}{n} \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}] \mathbb{E} [\mathbb{1}_{\{x_j \leq t\}}] + \frac{1}{n} \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}] \right. \\ &\quad \left. - 2 \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}] \mathbb{E} [\mathbb{1}_{\{y_j \leq t\}}] + \frac{m-1}{m} \mathbb{E} [\mathbb{1}_{\{y_i \leq t\}}] \mathbb{E} [\mathbb{1}_{\{y_j \leq t\}}] \right. \\ &\quad \left. + \frac{1}{m} \mathbb{E} [\mathbb{1}_{\{y_i \leq t\}}] \right) dt. \end{aligned}$$

Under the null hypothesis $E[\mathbb{1}_{\{x_i \leq t\}}] = E[\mathbb{1}_{\{y_i \leq t\}}] = H(t)$, thus

$$E[T_{n,m}] = \int_{-\infty}^{\infty} (H(t) - H(t)^2) dt, \quad (3.3)$$

□

Note that

$$E\left[\frac{nm}{n+m} (F_n(t) - G_m(t))^2\right] = H(t)(1 - H(t)),$$

is the divisor used in the Anderson-Darling test statistic, which confirms that the expectation of the Anderson-Darling test statistic is one (Pettitt, 1976). It is also of interest to note that the expectation does not depend on n and m .

3.3.2 Variance

The variance of $T_{n,m}$ under the null hypothesis for $n > m$ is given by

$$\begin{aligned} \text{Var}[T_{n,m}|F = G] = & \frac{2}{\mathcal{V}} \int_{-\infty}^{\infty} \int_{-\infty}^t H(s) \left(1 + 2(\mathcal{V} - 2)H(s) - 3H(t) \right. \\ & \left. - 2(2\mathcal{V} - 5)H(s)H(t) + 2H(t)^2 + 2(\mathcal{V} - 3)H(s)H(t)^2 \right) ds dt \end{aligned} \quad (3.4)$$

where $\mathcal{V} = \frac{nm(n+m)^2}{n^3+m^3}$. The proof of this is located in Appendix D. When $n = m$, i.e. $\mathcal{V} = 2n$, Equation (3.4) becomes

$$\begin{aligned} \text{Var}[T_{n,m}|F = G] = & \frac{1}{n} \int_{-\infty}^{\infty} \int_{-\infty}^t H(s) \left(1 + 4(n - 1)H(s) - 3H(t) \right. \\ & \left. - 2(4n - 5)H(s)H(t) + 2H(t)^2 + 2(2n - 3)H(s)H(t)^2 \right) ds dt. \end{aligned}$$

As $n = m \rightarrow \infty$ the variance can be approximated by

$$\text{Var}[T_{n,m}|F = G] \rightarrow 4 \int_{-\infty}^{\infty} \int_{-\infty}^t H(s)^2 (H(t) - 1)^2 ds dt. \quad (3.5)$$

3. Properties of the Cramer Test

3.3.3 Skewness

The third non-centralized moment of $T_{n,m}$ for $n > m$ is given by

$$\begin{aligned}
\mathbb{E}[T_{n,m}^3 | F = G] &= \frac{6}{\mathcal{G}} \int_{-\infty}^{\infty} \int_{-\infty}^s \int_{-\infty}^t H(t) \left(1 + 2 (\mathcal{G}(7(m^2 + n^2) - 10nm) - 8) H(t) \right. \\
&\quad + 2 (\mathcal{G}(5(m^2 + n^2) - 7nm) - 6) H(s) + (\mathcal{G}(m^2 + n^2 - nm) - 3) H(r) \\
&\quad + 5 (\mathcal{G}(2mn(m + n) - 19(m^2 + n^2) + 25mn) + 18) H(s)H(t) \\
&\quad + (\mathcal{G}(2mn(m + n) - 19(m^2 + n^2) + 25mn) + 18) H(s)^2 \\
&\quad - (\mathcal{G}(m^2 + n^2 - nm) - 2) H(r)^2 \\
&\quad + 2 (\mathcal{G}(mn(m + n) - 19(m^2 + n^2) + 26mn) + 20) H(r)H(t) \\
&\quad + (\mathcal{G}(mn(m + n) - 27(m^2 + n^2) + 37mn) + 30) H(r)H(s) \\
&\quad - 4 (\mathcal{G}(3mn(m + n) - 26(m^2 + n^2) + 34mn) + 24) H(s)^2 H(t) \\
&\quad - 5 (\mathcal{G}(5mn(m + n) - 45(m^2 + n^2) + 59mn) + 42) H(r)H(s)H(t) \\
&\quad - 2 (\mathcal{G}(mn(m + n) - 12(m^2 + n^2) + 16mn) + 12) H(r)^2 H(t) \\
&\quad - (\mathcal{G}(5mn(m + n) - 45(m^2 + n^2) + 59mn) + 42) H(r)H(s)^2 \\
&\quad - (\mathcal{G}(mn(m + n) - 17(m^2 + n^2) + 23mn) + 18) H(r)^2 H(s) \\
&\quad + 9 (\mathcal{G}(3mn(m + n) - 26(m^2 + n^2) + 34mn) + 24) H(r)H(s)^2 H(t) \\
&\quad + 5 (\mathcal{G}(3mn(m + n) - 26(m^2 + n^2) + 34mn) + 24) H(r)^2 H(s)H(t) \\
&\quad + (\mathcal{G}(3mn(m + n) - 26(m^2 + n^2) + 34mn) + 24) H(r)^2 H(s)^2 \\
&\quad \left. - 5 (\mathcal{G}(3mn(m + n) - 26(m^2 + n^2) + 34mn) + 24) H(r)^2 H(s)^2 H(t) \right) \\
&\quad dt ds dr
\end{aligned} \tag{3.6}$$

where $\mathcal{G} = \frac{nm(n+m)^2}{n^5+m^5}$. The proof of this is located in Appendix E. The skewness,

γ_T , can then be calculated as

$$\gamma_T = \frac{\mathbb{E}[T_{n,m}^3 | F = G] - 3\mathbb{E}[T_{n,m} | F = G]\text{Var}[T_{n,m} | F = G] - \mathbb{E}[T_{n,m} | F = G]^3}{\text{Var}[T_{n,m} | F = G]^{(\frac{3}{2})}}. \tag{3.7}$$

When $n = m$, Equation (3.6) simplifies to

$$\begin{aligned} \mathbb{E}[T_{n,m}^3 | F = G] &= \frac{3}{2n^2} \int_{-\infty}^{\infty} \int_{-\infty}^r \int_{-\infty}^s H(t) \left(1 + 16(n-1)H(t) + 12(n-1)H(s) \right. \\ &\quad + (2n-3)H(r) + 10(n-1)(4n-9)H(s)H(t) + 2(n-1)(4n-9)H(s)^2 \\ &\quad - 2(n-1)H(r)^2 + 8(n-1)(n-5)H(r)H(t) + 2(n-1)(2n-15)H(r)H(s) \\ &\quad - 48(n-1)(n-2)H(s)^2H(t) - 10(n-1)(10n-21)H(r)H(s)H(t) \\ &\quad - 8(n-1)(n-3)H(r)^2H(t) - 2(n-1)(10n-21)H(r)H(s)^2 \\ &\quad - 2(n-1)(2n-9)H(r)^2H(s) + 108(n-1)(n-2)H(r)H(s)^2H(t) \\ &\quad + 60(n-1)(n-2)H(r)^2H(s)H(t) + 12(n-1)(n-2)H(r)^2H(s)^2 \\ &\quad \left. - 60(n-1)(n-2)H(r)^2H(s)^2H(t) \right) dt ds dr. \end{aligned}$$

As $n = m \rightarrow \infty$ the third non-centralised moment can be approximated by

$$\begin{aligned} \mathbb{E}[T_{n,m}^3 | F = G] &\rightarrow 6 \int_{-\infty}^{\infty} \int_{-\infty}^r \int_{-\infty}^s H(t) \left(10H(s)H(t) + 2H(s)^2 + 2H(r)H(t) \right. \\ &\quad + H(r)H(s) - 12H(s)^2H(t) - 25H(r)H(s)H(t) - 2H(r)^2H(t) \\ &\quad - 30H(r)H(s)^2 - 6H(r)^2H(s) + 27H(r)H(s)^2H(t) + 15H(r)^2H(s)H(t) \\ &\quad \left. + 3H(r)^2H(s)^2 - 15H(r)^2H(s)^2H(t) \right) dr ds dt \end{aligned}$$

3.3.4 Kurtosis

The fourth non-centralized moment of $T_{n,m}$ along with the proof is given in Appendix F. When $n = m$, the fourth non-centralized moment is given by

$$\begin{aligned} \mathbb{E}[T_{n,m}^4 | F = G] &= \frac{3}{n^3} \int_{-\infty}^{\infty} \int_{-\infty}^s \int_{-\infty}^t \int_{-\infty}^u H(v) \left(1 + 64(n-1)H(v) + 48(n-1)H(u) \right. \\ &\quad + 12(n-1)H(t) + (2n-3)H(s) + 10(n-1)(40n-81)H(u)H(v) \\ &\quad + 2(n-1)(40n-81)H(u)^2 + 2(n-1)(4n-9)H(t)^2 - 2(n-1)H(s)^2 \\ &\quad + 24(n-1)(7n-15)H(t)H(v) + 32(n-1)(n-5)H(s)H(v) \\ &\quad + 2(n-1)(62n-135)H(t)H(u) + 24(n-1)(n-5)H(s)H(u) \\ &\quad + 2(n-1)(2n-15)H(s)H(t) + 192(n-1)(n-2)(n-4)H(u)^2H(v) \\ &\quad + 80(n-1)(n-2)(5n-21)H(t)H(u)H(v) \\ &\quad + 16(n-1)(n-2)(5n-21)H(t)H(u)^2 \\ &\quad \left. + 16(n-1)(n-2)(n-9)H(t)^2H(u) + 32(n-1)(n-2)(n-6)H(t)^2H(v) \right) \end{aligned}$$

3. Properties of the Cramer Test

$$\begin{aligned}
& + 10(n-1)(8n^2 - 110n + 189)H(s)H(u)H(v) - 32(n-1)(n-3)H(s)^2H(v) \\
& + 2(n-1)(8n^2 - 110n + 189)H(s)H(u)^2 - 24(n-1)(n-3)H(s)^2H(u) \\
& - 2(n-1)(10n - 21)H(s)H(t)^2 - 2(n-1)(2n - 9)H(s)^2H(t) \\
& + 8(n-1)(2n^2 - 55n + 105)H(s)H(t)H(v) \\
& + 2(n-1)(4n^2 - 160n + 315)H(s)H(t)H(u) \\
& - 108(n-1)(n-2)(8n - 25)H(t)H(u)^2H(v) \\
& - 60(n-1)(n-2)(8n - 25)H(t)^2H(u)H(v) \\
& - 12(n-1)(n-2)(8n - 25)H(t)^2H(u)^2 \\
& - 20(n-1)(n-2)(4n - 27)H(s)^2H(u)H(v) \\
& - 4(n-1)(n-2)(4n - 27)H(s)^2H(u)^2 \\
& + 12(n-1)(n-2)H(s)^2H(t)^2 \\
& - 20(n-1)(n-2)(50n - 189)H(s)H(t)H(u)H(v) \\
& - 96(n-1)(n-2)(5n - 18)H(s)H(u)^2H(v) \\
& - 16(n-1)(n-2)(5n - 27)H(s)H(t)^2H(v) \\
& - 16(n-1)(n-2)(n - 15)H(s)^2H(t)H(v) \\
& - 4(n-1)(n-2)(50n - 189)H(s)H(t)H(u)^2 \\
& - 4(n-1)(n-2)(10n - 81)H(s)H(t)^2H(u) \\
& - 4(n-1)(n-2)(2n - 45)H(s)^2H(t)H(u) \\
& + 720(n-1)(n-2)(n-3)H(t)^2H(u)^2H(v) \\
& + 108(n-1)(n-2)(18n - 55)H(s)H(t)H(u)^2H(v) \\
& + 96(n-1)(n-2)(3n - 10)H(s)^2H(u)^2H(v) \\
& + 60(n-1)(n-2)(18n - 55)H(s)H(t)^2H(u)H(v) \\
& + 300(n-1)(n-2)(2n - 7)H(s)^2H(t)H(u)H(v) \\
& + 48(n-1)(n-2)(n - 5)H(s)^2H(t)^2H(v) \\
& + 12(n-1)(n-2)(18n - 55)H(s)H(t)^2H(u)^2 \\
& + 60(n-1)(n-2)(2n - 7)H(s)^2H(t)H(u)^2 \\
& + 12(n-1)(n-2)(2n - 15)H(s)^2H(t)^2H(u) \\
& - 1560(n-1)(n-2)(n-3)H(s)H(t)^2H(u)^2H(v) \\
& - 1080(n-1)(n-2)(n-3)H(s)^2H(t)H(u)^2H(v) \\
& - 600(n-1)(n-2)(n-3)H(s)^2H(t)^2H(u)H(v) \\
& - 120(n-1)(n-2)(n-3)H(s)^2H(t)^2H(u)^2 \\
& + 840(n-1)(n-2)(n-3)H(s)^2H(t)^2H(u)^2H(v) \Big) ds dt du dv. \tag{3.8}
\end{aligned}$$

The kurtosis, κ , can then be calculated by

$$\begin{aligned} \kappa = \frac{1}{\text{Var}[T_{n,m}|F = G]^2} & \left(\text{E}[T_{n,m}^4|F = G] - 4\text{E}[T_{n,m}^3|F = G]\text{E}[T_{n,m}|F = G] \right. \\ & \left. + 6\text{Var}[T_{n,m}|F = G]^2\text{E}[T_{n,m}|F = G]^2 + 3\text{E}[T_{n,m}|F = G]^4 \right). \end{aligned} \quad (3.9)$$

As $n = m \rightarrow \infty$ the fourth non-centralized moment is given by

$$\begin{aligned} \text{E}[T_{n,m}^4|F = G] \rightarrow & 24 \int_{-\infty}^{\infty} \int_{-\infty}^s \int_{-\infty}^t \int_{-\infty}^u H(v) \left(24H(u)^2H(v) + 50H(t)H(u)H(v) \right. \\ & + 10H(t)H(u)^2 + 2H(t)^2H(u) + 4H(t)^2H(v) + 10H(s)H(u)H(v) \\ & + 2H(s)H(u)^2 + 2H(s)H(t)H(v) + H(s)H(t)H(u) \\ & - 108H(t)H(u)^2H(v) - 60H(t)^2H(u)H(v) - 12H(t)^2H(u)^2 \\ & - 10H(s)^2H(u)H(v) - 2H(s)^2H(u)^2 - 125H(s)H(t)H(u)H(v) \\ & - 60H(s)H(u)^2H(v) - 10H(s)H(t)^2H(v) - 2H(s)^2H(t)H(v) \\ & - 25H(s)H(t)H(u)^2 - 5H(s)H(t)^2H(u) - H(s)^2H(t)H(u) \\ & + 90H(t)^2H(u)^2H(v) + 243H(s)H(t)^2H(u)H(v) + 36H(s)^2H(u)^2H(v) \\ & + 135H(s)H(t)^2H(u)H(v) + 75H(s)^2H(t)H(u)H(v) + 6H(s)^2H(t)^2H(v) \\ & + 27H(s)H(t)^2H(u)^2 + 15H(s)^2H(t)H(u)^2 + 3H(s)^2H(t)^2H(u) \\ & - 195H(s)H(t)^2H(u)^2H(v) - 135H(s)^2H(t)H(u)^2H(v) \\ & - 75H(s)^2H(t)^2H(u)H(v) - 15H(s)^2H(t)^2H(u)^2 \\ & \left. + 105H(s)^2H(t)^2H(u)^2H(v) \right) ds dt du dv. \end{aligned}$$

3.4 Transformation of data

A sense of the magnitude of the Cramer test statistic may be useful in day-to-day practice. When the Cramer test statistic is calculated, we ought to be able to sense whether the magnitude of the test statistic is meaningful or not from hypothesis testing purposes. For this purpose, the pooled data needs to be standardised and so we first need to consider the effect of standardisation of data on the Cramer test statistic $T_{n,m}$.

Consider applying a linear transformation of the pooled data Z , such that $Z' = aZ + b$, where $a > 0$ and b are known constants. Since

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{z_i \leq t\}} = F_n(t),$$

3. Properties of the Cramer Test

then for the transformed data Z' , we can write

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{z'_i \leq t\}} &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{az_i + b \leq t\}} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{z_i \leq \frac{t-b}{a}\}} \\ &= F_n \left(\frac{t-b}{a} \right). \end{aligned}$$

The test statistic for the transformed data is given by

$$T'_{n,m} = \frac{nm}{n+m} \int_{-\infty}^{\infty} \left(F_n \left(\frac{t-b}{a} \right) - G_m \left(\frac{t-b}{a} \right) \right)^2 dt.$$

Let $s = \frac{t-b}{a}$, then a change of variables can be performed to give

$$\begin{aligned} T'_{n,m} &= \frac{anm}{n+m} \int_{-\infty}^{\infty} (F_n(s) - G_m(s))^2 ds \\ &= a \cdot T_{n,m}. \end{aligned} \tag{3.10}$$

This result therefore implies that

$$E[T'_{n,m}] = aE[T_{n,m}] \tag{3.11}$$

$$\text{Var}[T'_{n,m}] = a^2 \text{Var}[T_{n,m}]. \tag{3.12}$$

It can be easily shown that a change in location or scale will not affect the skewness or kurtosis of the test statistic. Hence, a change in location of the data does not affect the expectation, variance, skewness or kurtosis of the test statistic, but scaling affects the expectation and variance.

Note that Equation 3.10 shows that a linear transformation of the data will scale the test statistic by the same factor as the data. Then, if $t_{n,m}$ and $t'_{n,m}$ represents the test statistics for Z and Z' under H_0 respectively and $T_{n,m}$ and $T'_{n,m}$ are the observed test statistics for Z and Z' respectively, then

$$\begin{aligned} \Pr\{t'_{n,m} \geq T'_{n,m}\} &= \Pr\{at_{n,m} \geq aT_{n,m}\} \\ &= \Pr\{t_{n,m} \geq T_{n,m}\}. \end{aligned}$$

Similarly, we can show that $\Pr\{t'_{n,m} \leq T'_{n,m}\} = \Pr\{t_{n,m} \leq T_{n,m}\}$. Thus we can conclude that a linear transformation of the data will not affect the p -value of the test.

3.5 Expectation and Variance for Known Distribution

Assume that the null hypothesis is true, i.e. $F \equiv G \equiv H$ where H represents the true distribution of F and G under the null hypothesis. Also assume for simplicity that $n = m$. Consider we have data z_1, \dots, z_{2n} with cumulative distribution function H , which is the pool of two samples x_1, \dots, x_n and y_1, \dots, y_n . As seen in Section 3.4, scaling the pooled data Z affects the expectation and variance of the test statistic. Thus, it is of interest to see how the expectation and variance changes depending on the distribution of Z . For certain distributions of Z , we can obtain exact results for the expectation and variance.

Proposition For a random variable Z following various underlying distributions and $n = m$, the expectation and variance of the test statistic $T_{n,n}$ respectively follow the format

$$\begin{aligned} \mathbb{E}[T_{n,n}] &= \tilde{a} \cdot \text{sd}(Z), \\ \text{Var}[T_{n,n}] &= \left(\tilde{b} + \frac{\tilde{c}}{n} \right) \cdot \text{Var}(Z), \end{aligned} \tag{3.13}$$

where \tilde{a} , \tilde{b} , and \tilde{c} are real numbers and n is half the number of observations sampled from Z .

Proof Consider data sampled from a random variable Z with known cumulative distribution $H(t)$. The data can be standardised to form a new random variable Z' using the formula

$$Z' = \frac{Z - \mathbb{E}[Z]}{\text{sd}(Z)},$$

which can be rearranged to give

$$Z = \text{sd}(Z) \cdot Z' + \mathbb{E}[Z]. \tag{3.14}$$

Note that Equation (3.14) is of the same format as the linear transformation described in Section 3.4 with $a = \text{sd}(Z)$ and $b = \mathbb{E}[Z]$.

Now, define the test statistic calculated from the standardised data, Z' , to be $T'_{n,n}$ and the test statistic calculated from the original data, Z , to be $T_{n,n}$. From Equation (3.2), the expectation of $T_{n,n}$ is constant with respect to n . Equation (3.4) has terms which are constant in n as well as terms in $\frac{1}{n}$. Thus the expectation

3. Properties of the Cramer Test

and variance of $T'_{n,n}$, when calculated using $H(t)$, will be of the form

$$\begin{aligned} \mathbb{E}[T'_{n,n}] &= \tilde{a}, \\ \text{Var}[T'_{n,n}] &= \left(\tilde{b} + \frac{\tilde{c}}{n} \right), \end{aligned}$$

where \tilde{a} , \tilde{b} , and \tilde{c} are constants which depend on H and n is half the number of observations sampled from Z . We know from Section 3.4 that

$$\begin{aligned} \mathbb{E}[T'_{n,n}] &= a\mathbb{E}[T_{n,m}], \\ \text{Var}[T'_{n,n}] &= a^2\text{Var}[T_{n,m}], \end{aligned}$$

and here $a = \text{sd}(Z)$, thus

$$\begin{aligned} \mathbb{E}[T_{n,n}] &= \tilde{a} \cdot \text{sd}(Z), \\ \text{Var}[T_{n,n}] &= \left(\tilde{b} + \frac{\tilde{c}}{n} \right) \cdot \text{Var}(Z). \end{aligned}$$

Clearly, as $n \rightarrow \infty$ the value of \tilde{c} will be obsolete, as is evident from Equation (3.5). □

3.5.1 Example - Continuous Uniform Distribution

Suppose $Z \sim \text{Unif}(0,1)$, then $H(t) = t$ and $0 \leq t \leq 1$. Thus, Equation (3.2) becomes

$$\begin{aligned} \mathbb{E}[T_{n,n}] &= \int_0^1 H(t)(1 - H(t)) dt \\ &= \int_0^1 t(1 - t) dt \\ &= \left[\frac{t^2}{2} - \frac{t^3}{3} \right]_0^1 = \frac{1}{6}. \end{aligned}$$

As $\text{sd}(Z) = \frac{1}{2\sqrt{3}}$, this implies that $\tilde{a} = \frac{1}{\sqrt{3}}$.

Similarly, Equation (3.4) becomes

$$\begin{aligned}
 \text{Var}[T_{n,n}] &= \frac{1}{n} \int_{-\infty}^{\infty} \int_{-\infty}^t H(s) \left(1 + 4(n-1)H(s) - 3H(t) - 2(4n-5)H(s)H(t) \right. \\
 &\quad \left. + 2H(t)^2 + 2(2n-3)H(s)H(t)^2 \right) ds dt \\
 &= \frac{1}{n} \int_0^1 \int_0^t s \left(1 + 4(n-1)s - 3t - 2(4n-5)st + 2t^2 \right. \\
 &\quad \left. + 2(2n-3)st^2 \right) ds dt \\
 &= \frac{1}{n} \int_0^1 \left[\frac{s^2}{2} + \frac{4(n-1)s^3}{3} - \frac{3s^2t}{2} - \frac{2(4n-5)s^3t}{3} + s^2t^2 \right. \\
 &\quad \left. + \frac{2(2n-3)s^3t^2}{3} \right]_0^t dt \\
 &= \frac{1}{45} - \frac{1}{120n}.
 \end{aligned}$$

This implies that $\tilde{b} = \frac{4}{15}$ and $\tilde{c} = -\frac{1}{10}$. □

3.5.2 Example - Standard Normal Distribution

Suppose $Z \sim N(0,1)$, results for a non-standard normal distribution can be used using Equations (3.11) and (3.12). Here, we will use the standard notation of $\phi(t)$ and $\Phi(t)$ to represent the density function and cumulative distribution function of the random variable Z respectively. Thus Equation (3.2) becomes

$$\begin{aligned}
 \text{E}[T_{n,n}] &= \int_{-\infty}^{\infty} H(t)(1-H(t)) dt \\
 &= \int_{-\infty}^{\infty} \Phi(t)(1-\Phi(t)) dt
 \end{aligned}$$

To calculate this integral we use Proposition 1 with $\mu = 0$ and $\sigma = 1$ in Appendix C. Thus

$$\begin{aligned}
 \text{E}[T_{n,n}] &= \left[t\Phi(t) + \phi(t) - t\Phi(t)^2 - 2\phi(t)\Phi(t) + \frac{1}{\sqrt{\pi}}\Phi(\sqrt{2}t) \right]_{-\infty}^{\infty} \\
 &= \frac{1}{\sqrt{\pi}}.
 \end{aligned}$$

3. Properties of the Cramer Test

Thus $\tilde{a} = \frac{1}{\sqrt{\pi}}$. □

For $\text{Var}[T_{n,n}]$, it was difficult to find a solution analytically due to requiring the integration of $\Phi(t)^3$ and $\Phi(t)^4$, thus a numerical approach was applied. It is known that $\text{Var}[T_{n,n}] = \tilde{b} + \frac{\tilde{c}}{n}$, and as $n \rightarrow \infty$, $\text{Var}[T_{n,n}] = \tilde{b}$. Thus Equation (3.5) is calculated numerically to estimate $\tilde{b} = 0.20$. Now, using Proposition 2 in Appendix C, $T_{1,1} = \frac{|x_1 - y_1|}{2}$ and

$$\begin{aligned} \text{Var}[T_{1,1}] &= \text{E}[T_{1,1}^2] - \text{E}[T_{1,1}]^2 \\ &= \text{E} \left[\left(\frac{|X_1 - Y_1|}{2} \right)^2 \right] - \frac{1}{\pi} \\ &= \frac{1}{4} \text{E}[X_1^2] + \text{E}[Y_1^2] - 2\text{E}[X_1]\text{E}[Y_1] - \frac{1}{\pi} \\ &= \frac{1}{2} - \frac{1}{\pi}. \end{aligned}$$

Thus when $n = 1$, $\tilde{b} + \tilde{c} = \frac{1}{2} - \frac{1}{\pi}$, hence $\tilde{c} = \frac{1}{2} - \frac{1}{\pi} - \tilde{b}$. As numerically, $\tilde{b} = 0.2$, this implies that $\tilde{c} = 0.3 - \frac{1}{\pi}$.

3.5.3 Results for Other Unimodal Distributions Z

Table 3.1 provides the values of \tilde{a} , \tilde{b} and \tilde{c} for a selection of distributions for Z . From the results in Table 3.1, we can sense the magnitude of the test statistic by standardising the data Z to have standard deviation one. Note that the Bernoulli distribution is an outlier as \tilde{a} , \tilde{b} and \tilde{c} all depend on the parameter p . If $\text{sd}(Z) = 1$, then $\text{E}[T_{n,n}]$ is no larger than 0.58 and $\text{Var}[T_{n,n}]$ is no larger than 0.5 as $n \rightarrow \infty$. It was shown in Section 3.4 that scaling the data does not affect the p -value, thus performing this standardisation will not affect the results of the test. If the calculated test statistic (after standardisation of the data) is much larger than $0.58 + 1.96\sqrt{0.5} = 1.97$ then the null hypothesis can be immediately rejected. It is suggested, however, that the reader should calculate the p -value to ensure accuracy when the test statistic falls in the range of (1, 3) for standardised data. When the test statistic is less than one or more than three, the reader can immediately decide to not reject or to reject the null hypothesis at the 5% significance level, respectively, without the need to calculate the p -value.

Note that in some cases in Table 3.1, the value of \tilde{c} is negative and in some cases positive. This means that for increasing n , for some underlying distributions, $\text{Var}[T_{n,n}]$ will increase and for other underlying distributions, decrease. The cause for this is unclear, as no pattern can be seen between the underlying distributions with a negative \tilde{c} and those with a positive \tilde{c} .

Distribution of Z	\tilde{a}	\tilde{b}	\tilde{c}
Uniform(a, b)	$\frac{1}{\sqrt{3}}$	$\frac{4}{15}$	$-\frac{1}{10}$
Exponential(λ)	$\frac{1}{2}$	$\frac{1}{6}$	$\frac{1}{12}$
Laplace(μ, λ)	$\frac{3\sqrt{2}}{8}$	$\frac{7}{48}$	$\frac{7}{96}$
Bernoulli(p)	$\sqrt{p(1-p)}$	$2p(1-p)$	$-\frac{6p(1-p)+1}{2}$
Poisson(λ)	0.52 *	0.25 *	0.23 *
Normal(μ, σ)	$\frac{1}{\sqrt{\pi}}$	0.20 *	$0.3 - \frac{1}{\pi}$ *

Table 3.1: The values of \tilde{a} , \tilde{b} , and \tilde{c} in Eq. (3.13) for different underlying distributions of Z , in the expectation and variance of the test statistic $T_{n,m}$ under the null hypothesis. The symbol * indicates the value is obtained numerically.

3.5.4 Results for Multi-Modal Distributions Z

Recall that multi-modality is a common feature of the lung cancer data set. Let X and Y be distributed as mixture distributions which follow $p_1N(1, 0.25) + p_2N(1 + d, 0.25)$, thus the distribution of Z will also follow this setup. The values for \tilde{a} , \tilde{b} and \tilde{c} are calculated numerically and are shown in Table 3.2.

p_1	p_2	d	\tilde{a}	\tilde{b}	\tilde{c}
$\frac{7}{8}$	$\frac{1}{8}$	2	0.44 *	0.18 *	0.30 *
$\frac{1}{2}$	$\frac{1}{2}$	2	0.54 *	0.37 *	-0.62 *
$\frac{1}{8}$	$\frac{1}{8}$	4	0.38 *	0.19 *	-0.39 *
$\frac{1}{2}$	$\frac{1}{2}$	4	0.52 *	0.43 *	-0.54 *

Table 3.2: The values of \tilde{a} , \tilde{b} , and \tilde{c} in Eq. (3.13) for different multi-modal underlying distributions of Z , in the expectation and variance of the test statistic $T_{n,m}$ under the null hypothesis. The symbol * indicates the value is obtained numerically.

3.6 Discussion

The first four moments of the Cramer test statistic have been obtained. We have provided formulas to calculate the exact moments of the distribution when the underlying distribution function H is known or unknown. When the distribution function, H , is unknown the moments can be calculated using the empirical cumulative distribution functions, thus ensuring the test remains distribution free. Csorgo and Faraway (1996) states that we can use $H_{n+m}(t)$ to approximate $H(t)$ for large $n + m$. The test is also invariant to a linear transformation, suggesting that data can be standardised before performing the test without affecting the p -value. This then provides a way of rejecting the null hypothesis without the need for estimating a p -value.

Chapter 4

A Faster Approach to Estimate the p -value of the Cramer Test

4.1 Introduction

When performing a statistical significance test, once a test statistic is calculated the next step is to obtain a p -value. In Chapter 3, we obtained formulas for calculating the first four moments of the Cramer test (Baringhaus and Franz, 2004) when the distribution of the data is unknown. These formulas can be used to fit known distributions to a sample of test statistics by estimating the parameters through method of moments. If a distribution can be found which suitably describes a sample of test statistics, then the parameters can be estimated and thus a p -value obtained.

In this chapter, we present current reliable resampling methods including the method used by Baringhaus and Franz (2004) to estimate the null distribution - that is the distribution of the test statistics under the null hypothesis - and thus find a p -value. We then explore choices of distributions $\pi(\boldsymbol{\theta})$, where $\boldsymbol{\theta}$ represents the parameters of the distribution, to approximate the null distribution.

The purpose of this exploration is to find a method that will enable a faster calculation of a p -value compared to the resampling approaches. Recall that we wish to locate genomic regions for which the difference in CNA is significantly different between subtypes of cancer. To do this, we will perform the two-sample test on each genomic location. The number of genomic locations can be much greater than 10,000, thus ensuring a fast calculation of each p -value will ensure the overall analysis is not too slow. As well as this, in genomic studies, we require a very small p -value for a genome wide significance ($< 10^{-6}$). To obtain that level of p -value, the number of resampling that needs to be done for a single window

4. A Faster Approach to Estimate the p -value of the Cramer Test

is massive, thereby slowing down the calculation even further. For test statistics t_1, \dots, t_n , we will define the density function of the null distribution to be $\nu(t)$ and the distribution function $\mathcal{N}(t)$.

4.2 Resampling Approaches

4.2.1 Permutation Test

A current and reliable method for estimating the p -value is the permutation test (Fisher, 1935). In order to use this method, one must compare the test statistic $T_{n,m}$ calculated from the dataset to a distribution of values $t_i, i = 1, \dots, d$ obtained under H_0 . To calculate t_i , one must first permute the labels of the dataset and then calculate the test statistic of the new dataset, the test statistic is then defined as t_i for the i th permutation. Usually, to get an accurate estimation of the p -value, the value of d is chosen to be large.

This permutation test can be reasonably fast for small n , however an issue arises when n becomes large, namely the process to calculate a single p -value is computationally costly as well as taking a large amount of time. Note that whilst the lung cancer data set only has 38 observations in each sample, as technology gets better the sample sizes will increase. Thus to ensure the test can be used on larger datasets we aim to find a faster approach to calculate the p -value.

4.2.2 Bootstrapping the Limiting Distribution

Recall that Baringhaus and Franz (2004) states that the limiting distribution for $T_{n,m}$ as $n, m \rightarrow \infty$ is

$$\int_{-\infty}^{\infty} B^2(H(t))dt,$$

where $B(u), 0 \leq u \leq 1$ is the classical Brownian bridge. To calculate a p -value Baringhaus and Franz (2004) suggests to bootstrap (Efron, 1992) the limiting distribution, i.e. by estimating the test statistic by

$$\hat{T}_{n,m} = \int_{-\infty}^{\infty} B^2(H_{n,m}(t))dt. \quad (4.1)$$

If z_1, \dots, z_{n+m} are the pooled ordered data of x_1, \dots, x_n and y_1, \dots, y_m , Equation (4.1) can be calculated using

$$\sum_{i=1}^{n+m-1} (z_{i+1} - z_i) \left[B\left(\frac{i}{n}\right) \right]^2. \quad (4.2)$$

To bootstrap, samples are drawn from z_1, \dots, z_{n+m} with replacement and $\hat{T}_{n,m}$ calculated using Equation (4.2). This is repeated a large number of times.

Whilst this method is faster than the permutation test, it can take a long time when n and m get larger. Also, to obtain accurate estimates for the null distribution, the number of times you resample needs to be large and will thus slow down the calculation. Again, if the aim is to ensure a fast calculation of the p -value for larger sample sizes, an alternative approach is required.

4.3 Empirical Approximations

If the null distribution can be approximated by a known distribution $\pi(\boldsymbol{\theta})$ for some parameters $\boldsymbol{\theta}$, which depend on $H(t)$, then the calculations to obtain a p -value will become less computationally costly. As formulas for the first four moments are obtained, estimating the parameters $\boldsymbol{\theta}$ should be possible via the method of moments.

The aim here is to find a two parameter distribution, $\pi(\boldsymbol{\theta})$, which can approximate the null distribution. To estimate the null distribution, consider simulating n observations from random variables X and Y with some distribution functions $F \equiv G$ and obtain a sampled test statistic t_1 . If this process is repeated k times we will therefore have k test statistics, t_1, \dots, t_k , for which the empirical cumulative distribution (ECD) function can be plotted. This ECD function - provided k is large enough - will be a good estimate of the null distribution. We can thus use the ECD function as a tool for finding a distribution which closely matches the null distribution.

4.3.1 ECD Function of $T_{n,m}$ for Various Distributions of X and Y

The shape of the null distribution is investigated for three distributional forms for $H(t)$. We choose three distributional forms which reflect the shape of the lung cancer data set. We therefore choose a unimodal normal distribution, a multi-modal distribution with two peaks and a multi-modal distribution with three peaks.

For the unimodal normal distribution, consider $X \sim N(0,1)$, $Y \sim N(0,1)$, $n = m = 10000$ and $k = 10000$. Figure 4.1 shows the histogram (left) and the ECD function (right) of the 10000 test statistics.

Next, consider X and Y which both follow mixture distributions defined by $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$ for $n = m = 10000$ and $k = 10000$. Also, for the

4. A Faster Approach to Estimate the p -value of the Cramer Test

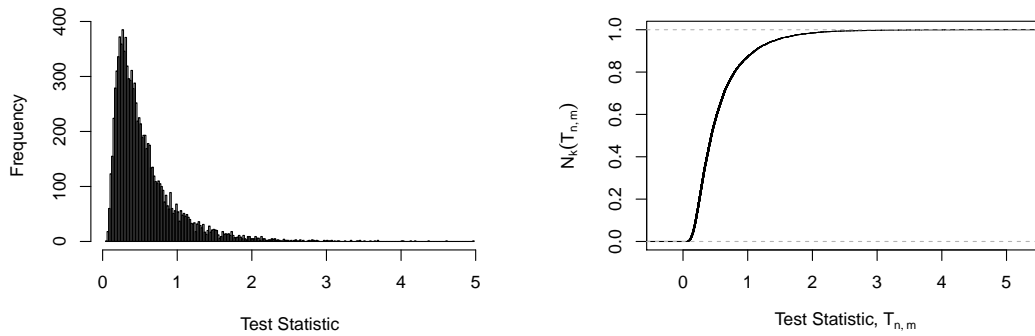


Figure 4.1: The histogram (left) and the ECD function (right) of the 10000 test statistics where $X \sim N(0,1)$, $Y \sim N(0,1)$, $n = m = 10000$ and $k = 10000$.

purpose of comparison, consider standardising X and Y . Figure 4.2 shows the histogram (left) and the ECD function (right) of the 10000 test statistics.

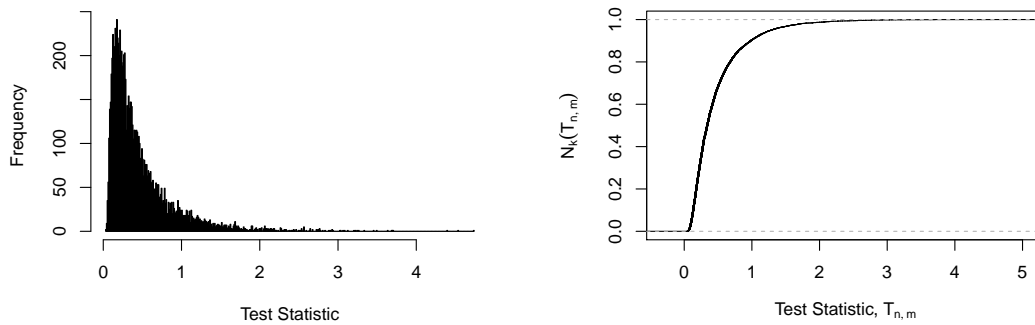


Figure 4.2: The histogram (left) and the ECD function (right) of the 10000 test statistics where X and Y are both mixture distributions defined by $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$ for $n = m = 10000$ and $k = 10000$.

Finally, consider X and Y which both follow mixture distributions defined by $\frac{5}{9}N(1, 0.25) + \frac{3}{9}N(3, 0.25) + \frac{1}{9}N(5, 0.25)$ for $n = m = 10000$ and $k = 10000$. Also, for the purpose of comparison, consider standardising X and Y . Figure 4.3 shows the histogram (left) and the ECD function (right) of the 10000 test statistics.

To compare the ECDF for each choice of distribution, Figure 4.4 shows all three ECDF curves plotted on a single graph.

It is clear by comparing from Figure 4.4 that the shape of the null distribution changes depending on the distributions of X and Y . Thus $\nu(z)$ and $\mathcal{N}(z)$ will depend on F and G , or equivalently H under the null hypothesis. As H_{n+m} is used to approximate H in application, this means that for each variable a different null distribution will be used.

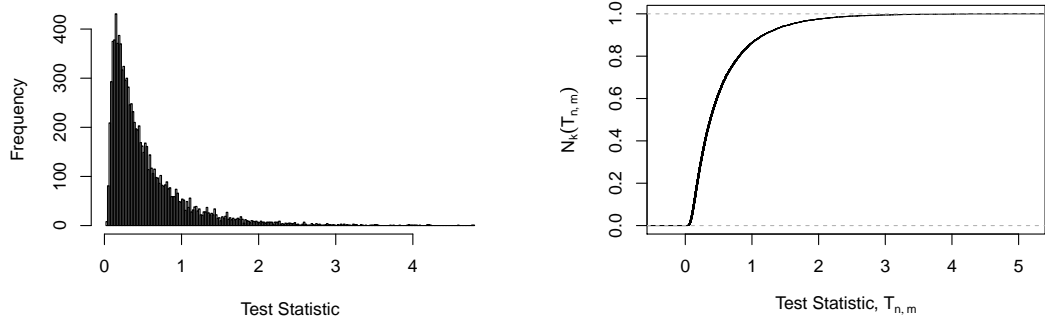


Figure 4.3: The histogram (left) and the ECD function (right) of the 10000 test statistics where X and Y are both mixture distributions defined by $\frac{5}{9}N(1, 0.25) + \frac{3}{9}N(3, 0.25) + \frac{1}{9}N(5, 0.25)$ for $n = m = 10000$ and $k = 10000$.

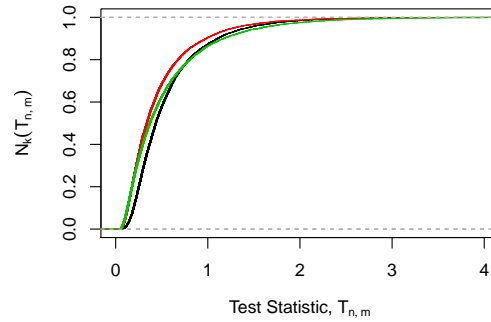


Figure 4.4: The ECDF curves of test statistics when X and Y are sampled from each distributional form.

4. A Faster Approach to Estimate the p -value of the Cramer Test

Note also that when $X \sim N(0,1)$ and $Y \sim N(0,1)$ and $n = m = 1$, it can be shown that $T_{1,1} = \frac{|x_1 - y_1|}{2}$ has a scaled chi-distribution, because of this, the first choice of $\pi(\boldsymbol{\theta})$ is the scaled chi-square distribution.

In the following sections, we now consider the non-standardised data as it is more reflective of the lung cancer data set.

4.3.2 Generalised Pareto Distribution

In Section B.2, the Generalised Pareto distribution was used to estimate the null distribution specifically in the right tail. The parameters σ and ξ were calculated using the information provided by the data and the parameter μ was chosen to be the threshold λ . The main reason for not using this method was due to the estimation of \bar{t}_λ and s_λ , the mean and variance of the test statistics above the threshold λ .

The probability density function of the generalised Pareto distribution is defined in Equation (B.3) and the mean, variance and skewness of the distribution can be calculated using the following three equations;

$$\begin{aligned} E[X] &= \mu + \frac{\sigma}{1 - \xi} & \xi < 1; \\ \text{Var}[X] &= \frac{\sigma^2}{(1 - \xi)^2(1 - 2\xi)} & \xi < \frac{1}{2}; \\ \gamma_X &= \frac{2(1 + \xi)\sqrt{1 - 2\xi}}{1 - 3\xi} & \xi < \frac{1}{3}. \end{aligned}$$

By calculating $E[T_{n,m}]$, $\text{Var}[T_{n,m}]$ and γ_T using the formulae in Section 3.3, we can use method of moments to estimate the parameters μ , σ and ξ . By doing this, we can fit a GPD to the entire null distribution, not just the right tail. Consider the test statistics calculated when $X \sim N(0,1)$, $Y \sim N(0,1)$, the parameters are estimated to be $\mu = 0.121$, $\sigma = 0.473$, and $\xi = -0.025$. Figure 4.5 shows the QQ-plot comparing the percentiles of the 10,000 test statistics and the percentiles of the fitted GPD. Figure 4.5 shows that all points lie in a straight line with a slight exception in the left tail. However as we are more interested in the fit in the right tail, this is not cause for concern.

Next consider the test statistics when X and Y both follow mixture distributions defined by $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$, the parameters are calculated to be $\mu = 0.052$, $\sigma = 0.260$ and $\xi = 0.071$. Figure 4.6 shows the QQ-plot comparing the percentiles of the 10,000 test statistics and the percentiles of the fitted GPD. Figure 4.6 shows that again, the fit in the right tail is reasonably accurate.

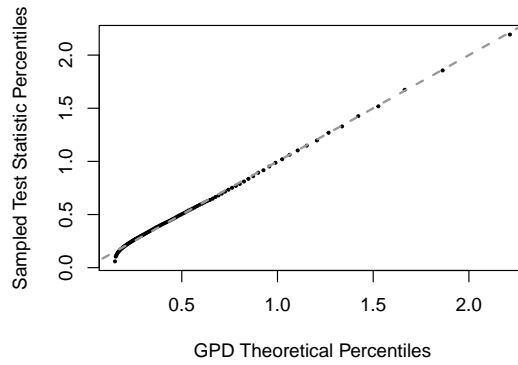


Figure 4.5: The percentiles of the 10,000 sampled test statistics when X and Y are distribution as $N(0, 1)$ plotted against the percentiles of the fitted GPD with $\mu = 0.121$, $\sigma = 0.473$, and $\xi = -0.025$.

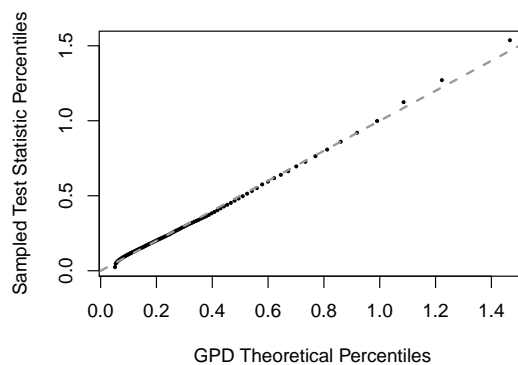


Figure 4.6: The percentiles of the 10,000 sampled test statistics when X and Y both follow mixture distributions defined by $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$ plotted against the percentiles of the GPD with $\mu = 0.052$, $\sigma = 0.260$ and $\xi = 0.071$.

4. A Faster Approach to Estimate the p -value of the Cramer Test

Finally consider the test statistics when X and Y both follow mixture distributions defined by $\frac{5}{9}N(1, 0.25) + \frac{3}{9}N(3, 0.25) + \frac{1}{9}N(5, 0.25)$, the parameters are calculated to be $\mu = 0.083$, $\sigma = 0.628$ and $\xi = 0.083$. Figure 4.7 shows the QQ-plot comparing the percentiles of the 10,000 test statistics and the percentiles of the fitted GPD. Figure 4.7 shows that the generalised Pareto distribution is an accurate representation of the null distribution for data distributed as a Poisson distribution.

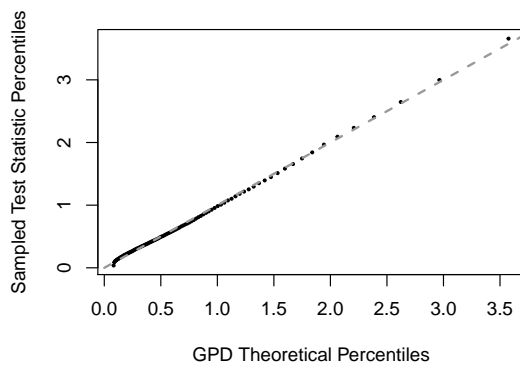


Figure 4.7: The percentiles of the 10,000 sampled test statistics plotted against the percentiles of the fitted GPD with $\mu = 0.083$, $\sigma = 0.628$ and $\xi = 0.083$.

Other three-parameter distributions such as the three-parameter gamma and the three-parameter log-normal distributions were also investigated but the approximation was less satisfactory compared to that of the GPD.

4.3.3 Measuring the Accuracy in the Right Tail

Recall that to measure the accuracy of the fit in the right tail of the distribution we can calculate

$$RT_{acc} = \int_{0.9}^1 (p - \mathcal{N}_k(\mathcal{N}^{-1}(p)))^2 dp$$

For the test statistics calculated when $X \sim N(0, 1)$, $Y \sim N(0, 1)$, $RT_{acc} = 5.87 \times 10^{-8}$, for the test statistics calculated when X and Y both follow mixture distributions defined by $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$, $RT_{acc} = 3.03 \times 10^{-7}$ and for the test statistics calculated when X and Y both follow mixture distributions defined by $\frac{5}{9}N(1, 0.25) + \frac{3}{9}N(3, 0.25) + \frac{1}{9}N(5, 0.25)$, $RT_{acc} = 2.15 \times 10^{-7}$. Clearly, the GPD is a very good fit to the null distribution when the underlying distributions are either normally distributed, or multi-modally distributed.

4.3.4 Comparing Empirical and Theoretical Quantiles

Let X and Y be distributed as mixture distributions which follows $p_1N(1, 0.25) + p_2N(1 + d, 0.25) + p_3N(1 + 2d, 0.25)$. To show that the GPD is satisfactory choice of $\pi(\boldsymbol{\theta})$ to estimate the null hypothesis, for different choices of p_1, p_2, p_3 and d , sampled with $n = m = 100$, of data, Table 4.1 shows

1. the 95th quantile of the empirical cumulative distribution function, denoted $H_k^{-1}(0.95)$, for $k = 10,000$,
2. the 95th quantile of the cumulative adjusted GPD, denoted $H^{-1}(0.95)$, and
3. the probability of obtaining a value from the fitted adjusted GPD larger than the 95th quantile of the empirical cumulative distribution function, denoted p_k , for $k = 10,000$.

p_1	p_2	p_3	d	$H_{10000}^{-1}(0.95)$ (empirical)	$H^{-1}(0.95)$ (GPD)	$p^{0.95}$
1	0	0	0	1.432	1.450	0.052
$\frac{7}{8}$	$\frac{1}{8}$	0	2	0.905	0.906	0.050
$\frac{1}{2}$	$\frac{1}{2}$	0	2	1.788	1.760	0.048
$\frac{7}{8}$	$\frac{1}{8}$	0	4	1.686	1.681	0.050
$\frac{5}{9}$	$\frac{3}{9}$	$\frac{1}{9}$	2	2.171	2.166	0.050

Table 4.1: For X and Y distributed as mixture distributions which follow $p_1N(1, 0.25) + p_2N(1 + d, 0.25) + p_3N(1 + 2d, 0.25)$ with $n = m = 100$, the 95th percentile of the empirical cumulative distribution function for $k = 10,000$, the 95th percentile of the fitted GPD, and the probability p_k of obtaining a value from the fitted GPD larger than the 95th quantile of the empirical cumulative distribution function. ϕ denotes the normal probability density function.

For different distributions, sampled with $n = m = 100$, Table 4.2 shows

1. the 99.5th quantile of the empirical cumulative distribution function, denoted $H_k^{-1}(0.995)$, for $k = 10,000$,
2. the 99.5th quantile of the cumulative adjusted GPD, denoted $H^{-1}(0.995)$, and
3. the probability of obtaining a value from the fitted adjusted GPD larger than the 99.5th quantile of the empirical cumulative distribution function, denoted p_k , for $k = 10,000$.

4. A Faster Approach to Estimate the p -value of the Cramer Test

p_1	p_2	p_3	d	$H_{10000}^{-1}(0.95)$ (empirical)	$H^{-1}(0.95)$ (GPD)	$p^{0.95}$
1	0	0	0	2.635	2.636	0.0050
$\frac{7}{8}$	$\frac{1}{8}$	0	2	1.638	1.648	0.0052
$\frac{1}{2}$	$\frac{1}{2}$	0	2	3.487	3.486	0.0050
$\frac{7}{8}$	$\frac{1}{8}$	0	4	3.412	3.398	0.0049
$\frac{5}{9}$	$\frac{3}{9}$	$\frac{1}{9}$	2	4.198	4.195	0.0050

Table 4.2: For X and Y distributed as mixture distributions which follow $p_1N(1, 0.25) + p_2N(1 + d, 0.25) + p_3N(1 + 2d, 0.25)$ with $n = m = 100$, the 99.5th percentile of the empirical cumulative distribution function for $k = 10,000$, the 99.5th percentile of the fitted GPD, and the probability p_k of obtaining a value from the fitted GPD larger than the 99.5th quantile of the empirical cumulative distribution function, for various distributions. ϕ denotes the normal probability density function.

Table 4.1 and 4.2 shows that the null distribution of the proposed test statistic $T_{n,m}$ can be well approximated by the GPD for different distributions of data at different significance levels for $n = m = 100$. For the lung cancer data set, $n = m = 38$ hence it is worth showing that the null distribution of the proposed test statistic $T_{n,m}$ can be well approximated by the GPD for a smaller sample size. Thus Tables 4.1 and 4.2 are repeated for $n = m = 30$.

For different distributions, sampled with $n = m = 30$, Table 4.3 shows

1. the 95th quantile of the empirical cumulative distribution function, denoted $H_k^{-1}(0.95)$, for $k = 10,000$,
2. the 95th quantile of the cumulative adjusted GPD, denoted $H^{-1}(0.95)$, and
3. the probability of obtaining a value from the fitted adjusted GPD larger than the 95th quantile of the empirical cumulative distribution function, denoted p_k , for $k = 10,000$.

For different distributions, sampled with $n = m = 30$, Table 4.4 shows

1. the 99.5th quantile of the empirical cumulative distribution function, denoted $H_k^{-1}(0.995)$, for $k = 10,000$,
2. the 99.5th quantile of the cumulative adjusted GPD, denoted $H^{-1}(0.995)$, and
3. the probability of obtaining a value from the fitted adjusted GPD larger than the 99.5th quantile of the empirical cumulative distribution function, denoted p_k , for $k = 10,000$.

p_1	p_2	p_3	d	$H_{10000}^{-1}(0.95)$ (empirical)	$H^{-1}(0.95)$ (GPD)	$p^{0.95}$
1	0	0	0	1.460	1.443	0.048
$\frac{7}{8}$	$\frac{1}{8}$	0	2	0.886	0.882	0.049
$\frac{1}{2}$	$\frac{1}{2}$	0	2	1.728	1.691	0.047
$\frac{7}{8}$	$\frac{1}{8}$	0	4	1.654	1.639	0.049
$\frac{5}{9}$	$\frac{3}{9}$	$\frac{1}{9}$	2	2.089	2.077	0.049

Table 4.3: For X and Y distributed as mixture distributions which follow $p_1N(1, 0.25) + p_2N(1 + d, 0.25) + p_3N(1 + 2d, 0.25)$ with $n = m = 30$, the 95th percentile of the empirical cumulative distribution function for $k = 10,000$, the 95th percentile of the fitted GPD, and the probability p_k of obtaining a value from the fitted GPD larger than the 95th quantile of the empirical cumulative distribution function, for various distributions. ϕ denotes the normal probability density function.

p_1	p_2	p_3	d	$H_{10000}^{-1}(0.95)$ (empirical)	$H^{-1}(0.95)$ (GPD)	$p^{0.95}$
1	0	0	0	2.617	2.608	0.0049
$\frac{7}{8}$	$\frac{1}{8}$	0	2	1.667	1.650	0.0048
$\frac{1}{2}$	$\frac{1}{2}$	0	2	3.318	3.315	0.0050
$\frac{7}{8}$	$\frac{1}{8}$	0	4	3.261	3.267	0.0050
$\frac{5}{9}$	$\frac{3}{9}$	$\frac{1}{9}$	2	3.983	3.922	0.0047

Table 4.4: For X and Y distributed as mixture distributions which follow $p_1N(1, 0.25) + p_2N(1 + d, 0.25) + p_3N(1 + 2d, 0.25)$ with $n = m = 30$, the 99.5th percentile of the empirical cumulative distribution function for $k = 10000$, the 99.5th percentile of the fitted GPD, and the probability p_k of obtaining a value from the fitted GPD larger than the 99.5th quantile of the empirical cumulative distribution function, for various distributions. ϕ denotes the normal probability density function.

Tables 4.3 and 4.4 show that the null distribution of the proposed test statistic $T_{n,m}$ can be well approximated by the GPD for different distributions of data at different significance levels for $n = m = 30$.

4.4 Alternative Approaches

In Appendix B we discuss a few other methods which were attempted in order to obtain a suitable null distribution approximation. We firstly discuss matching a two-parameter distribution but found none to be suitable enough for our purposes. We also consider a transformation of variables as well as the extreme value theorem. All methods either didn't provide an accurate enough approximation or in the case of the extreme value theorem was found to be difficult to implement in practice.

4. A Faster Approach to Estimate the p -value of the Cramer Test

It was mainly due to the extreme value theorem that the GPD was considered and ultimately chosen.

In Bayesian statistics, one could estimate an unknown distribution using Markov chain Monte Carlo (MCMC) or approximate Bayesian computation (ABC) algorithms. In our case, we cannot assume that the distribution of the data X and Y is constant across all regions of the genome and as was seen in Section 4.3.1 the distribution of test statistics depends on the distributions for X and Y . Therefore to implement these techniques they would need to be performed for every test statistic calculated to obtain its null distribution. As we have already mentioned we want the technique to be as fast as is reasonably possible, thus repeating this step over all genomic regions may not be feasible. For argument purposes, let's say implementing the algorithm on every genomic region is feasible, for the case of MCMC, the question of what proposal distribution we use however is an issue as it is required to be proportional to the posterior distribution - for which a reasonable choice is completely unknown. We therefore still face the issue of finding which distribution closely fits the null distribution for any data X and Y . Once this has been done of course MCMC and ABC could be used to estimate the parameters of this distribution. However, estimating the parameters using the method of moments seemed at the time the fastest approach. Note that other Bayesian techniques exist to estimate unknown distributions and could therefore be used to estimate the null distribution for each genomic region. However, it is highly unlikely that a method exists which can estimate the distribution without knowing its general form from a single value or test statistic without the need for resampling. As we wish to avoid resampling techniques, we are content that our method can be considered a suitable solution.

4.5 Discussion

The permutation test and the bootstrap method used by [Baringhaus and Franz \(2004\)](#) is firstly described. It is noted that both these methods can be computationally expensive for large n and m . Also as the test is needed to be repeated for each window of the genome (17,000 times for the lung cancer data set), ensuring the test runs as quickly as possible is important. Thus methods for calculating the p -value which does not involve resampling is considered.

Whilst a two-parameter distribution was found to be an unsuitable choice for $\pi(\theta)$, methods such as transforming the test statistic and applying the extreme

value theorem were investigated. It was found that these methods were also found to be unsuitable.

In our application using the extreme value theorem is unsuitable as the expectation and variance of the test statistic after a certain threshold λ is required. Whilst formulas for the expectation and variance of the test statistic has been developed, obtaining the formulas given the data is greater than λ , proves to be challenging.

Finally, three-parameter distributions were considered for $\pi(\boldsymbol{\theta})$. The generalised Pareto distribution was found to closely fit the null distribution when the underlying distribution follows any form. The distribution can easily be estimated by calculating the parameters of the distribution using the expectation, variance and skewness formulas in Chapter 3, Section 3.3. Hence, once the distribution is estimated, a p -value can be obtained.

One could also consider distributions which rely on more than three parameters. As we have calculated the fourth moment of the test statistic, this could be a reasonably exercise. However, we understand that by increasing the number of parameters, the complexity of the distribution increases. Thus we preferred to consider a distribution which sufficiently describes the null distribution without being too complicated.

It should be noted that perhaps no 2 or 3 parameter distribution is a suitable fit for the entire null distribution. We present in this chapter enough evidence to assume that the generalised Pareto distribution accurately fits the right tail of the null distribution. Of course in this case we do not claim that the GPD is a suitable fit for the whole distribution. The simulations we have chosen in this case are aimed to reflect the shape of the lung cancer data set as much as possible. Therefore as we have shown that the GPD is a suitable fit in the chosen simulation scenarios we can assume that the GPD will be a suitable fit for the lung cancer data set.

Chapter 5

Application of Two Sample Test

5.1 Introduction

We found, in Chapter 4, a method to calculate the p -value that does not involve resampling and instead uses the method of moments to estimate parameters of a fitted generalised Pareto distribution. The idea behind creating such a method was due to the speed of the resampling approaches. Consider the lung cancer data set, in this case over 17,000 simultaneous Cramer tests are required to be performed. If for example, each test was to take a second to run and we assume that the user does not have access to parallel processing, this means it will take just under 5 hours to calculate over 17,000 simultaneous p -values, which is reasonable. However, as next generation sequencing technologies are improving, it is becoming possible to obtain more accurate read counts for smaller genomic regions, thereby causing the total number of genomic regions or windows to increase. In particular, we could be expected to locate genomic regions of interest by comparing over 200,000 genomic regions. In this case, if each test was expected to take a second to run, then it will take just under 2 and a half days to calculate 200,000 p -values.

In this chapter we focus on the application of the Cramer test and in particular the use of our method to calculate the p -value. We begin by considering the computational time to calculate the p -value and investigate and implement methods for speeding up the computational calculation, thereby ensuring a fast user friendly test when required to perform hundreds of thousands of simultaneous tests. To this end, we create a new R package which implements these methods. We also investigate the false positive rate when calculating the p -value using our approach from Chapter 4. We do this to ensure a properly controlled false positive rate which will provide the user with confidence that the correct level of false positives are being identified. Finally, we apply the Cramer test on all genomic locations of

5. Application of Two Sample Test

the lung cancer data set to identify regions of significance. We will then compare our results to other methods to see whether the results obtained from our method are consistent with other results.

5.2 Computational Considerations

5.2.1 Test Statistic

In order to calculate the test statistic of the Cramer test a single integral needs to be performed, namely

$$T_{n,m} = \frac{nm}{n+m} \int_{-\infty}^{\infty} (F_n(t) - G_m(t))^2 dt.$$

The computer programme R can easily compute this integral and on a standard MacBook pro this takes approximately 0.017 seconds to run for $n = m = 100$. However when applying this test to over 17,000 genomic regions simultaneously, the time increases to approximately 5 hours. This in itself is manageable, however only the test statistic has been obtained and the calculation of the p -value is yet to be performed. Thus any time which can be saved at this stage will be vital to ensure that the total computational time is minimal.

Consider sampling $x_1 < \dots < x_n$ and $y_1 < \dots < y_m$, $n = m = 100$, from a standard normal distribution and let

$$\mathcal{J}(t) = (F_n(t) - G_m(t))^2$$

be the integrand of the test statistic.

Proposition The integrand of the test statistic $\mathcal{J}(t) = 0$ when $t < \min(x_1, y_1)$ or $t > \max(x_n, y_m)$.

Proof For $t < x_1$, $F_n(t) = 0$ and for $t < y_1$, $G_m(t) = 0$, thus $t < \min(x_1, y_1)$ implies $\mathcal{J}(t) = 0$. For $t > x_n$, $F_n(t) = 1$ and for $t > y_m$, $G_m(t) = 1$, thus $t > \max(x_n, y_m)$ implies $\mathcal{J}(t) = 0$. \square

Thus using the results of the proposition, we can instead calculate

$$T_{n,m} = \frac{nm}{n+m} \int_{\min(x_1, y_1)}^{\max(x_n, y_m)} (F_n(t) - G_m(t))^2 dt.$$

This has already reduced the time required to compute the integral to 0.013 seconds, however the calculation can be further reduced in time when considering the

integral as a sum. Because $F_n(t)$ and $G_m(t)$ are step functions, $J(t)$ will also be a step function. Also, because $J(t) > 0 \forall t$, the integral can easily be reduced to a summation. To see this, consider Figure 5.1 which shows the function $J(t)$ plotted against $t \in [-3, 3]$.

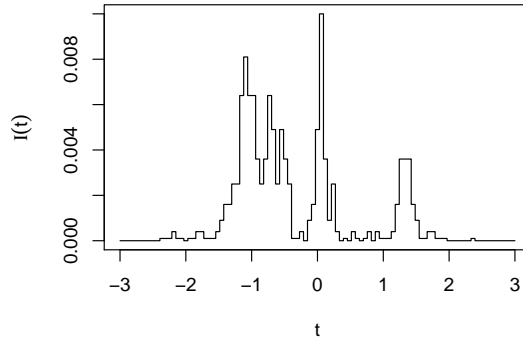


Figure 5.1: The integrand of the test statistic $J(t)$ plotted against $t \in [-3, 3]$.

For the sampled data, $\min(x_1, y_1) = -2.79$ and $\max(x_n, y_m) = 2.36$. Figure 5.1 clearly shows that for any value of t outside the range $[-2.79, 2.36]$, $J(t) = 0$ as well as showing that $J(t)$ is indeed a step function.

Let $z_1 < \dots < z_{n+m}$ be the sorted pooled sample of the two groups of data. Then the integral can be reduced to the following summation;

$$T_{n,m} = \frac{nm}{n+m} \sum_{i=1}^{n+m-1} (F_n(z_i) - G_m(z_i))^2 \cdot (z_{i+1} - z_i),$$

which can be easily coded into R. When we code the test statistic in this way, the computational time reduces to less than 0.001 seconds.

5.2.2 Moments

To calculate the p -value, three further integrals need to be calculated, namely the mean, variance and third moment of the test statistic (see Section 3.3). The formula for the mean of the test statistic involves a single integral, the variance, a double integral, and the third moment, a triple integral. Due to the nature of the integrals, they will take much longer to calculate than the test statistic. We can, however, use the same techniques to speed up these calculations as we did for the test statistic whilst also using a trick which is adopted in the trapezoidal rule.

5. Application of Two Sample Test

Firstly let

$$ET_{\mathcal{L}} = H_{n+m}(z_{i+\mathcal{L}})(1 - H_{n+m}(z_{i+\mathcal{L}})) \quad (5.1)$$

$$\begin{aligned} VT_{\mathcal{L}} = & H_{n+m}(z_{j+\mathcal{L}}) \left(1 + 2(\mathcal{V} - 2) H_{n+m}(z_{j+\mathcal{L}}) - 3H_{n+m}(z_{i+\mathcal{L}}) \right. \\ & - 2(2\mathcal{V} - 5) H_{n+m}(z_{j+\mathcal{L}})H_{n+m}(z_{i+\mathcal{L}}) + 2H_{n+m}(z_{i+\mathcal{L}})^2 \\ & \left. + 2(\mathcal{V} - 3) H_{n+m}(z_{j+\mathcal{L}})H_{n+m}(z_{i+\mathcal{L}})^2 \right) \end{aligned} \quad (5.2)$$

$$\begin{aligned} GT_{\mathcal{L}} = & H(z_{k+\mathcal{L}}) \left(1 + 2(\mathcal{G}(7(m^2 + n^2) - 10nm) - 8) H(z_{k+\mathcal{L}}) \right. \\ & + 2(\mathcal{G}(5(m^2 + n^2) - 7nm) - 6) H(z_{j+\mathcal{L}}) \\ & + (\mathcal{G}(m^2 + n^2 - nm) - 3) H(z_{i+\mathcal{L}}) \\ & + 5(\mathcal{G}(2mn(m + n) - 19(m^2 + n^2) + 25mn) + 18) H(z_{j+\mathcal{L}})H(z_{k+\mathcal{L}}) \\ & + (\mathcal{G}(2mn(m + n) - 19(m^2 + n^2) + 25mn) + 18) H(z_{j+\mathcal{L}})^2 \\ & - (\mathcal{G}(m^2 + n^2 - nm) - 2) H(z_{i+\mathcal{L}})^2 \\ & + 2(\mathcal{G}(mn(m + n) - 19(m^2 + n^2) + 26mn) + 20) H(z_{i+\mathcal{L}})H(z_{k+\mathcal{L}}) \\ & + (\mathcal{G}(mn(m + n) - 27(m^2 + n^2) + 37mn) + 30) H(z_{i+\mathcal{L}})H(z_{j+\mathcal{L}}) \\ & - 4(\mathcal{G}(3mn(m + n) - 26(m^2 + n^2) + 34mn) + 24) H(z_{j+\mathcal{L}})^2H(z_{k+\mathcal{L}}) \\ & - 5(\mathcal{G}(5mn(m + n) - 45(m^2 + n^2) + 59mn) \\ & \quad \left. + 42)H(z_{i+\mathcal{L}})H(z_{j+\mathcal{L}})H(z_{k+\mathcal{L}}) \right. \\ & - 2(\mathcal{G}(mn(m + n) - 12(m^2 + n^2) + 16mn) + 12) H(z_{i+\mathcal{L}})^2H(z_{k+\mathcal{L}}) \\ & - (\mathcal{G}(5mn(m + n) - 45(m^2 + n^2) + 59mn) + 42) H(z_{i+\mathcal{L}})H(z_{j+\mathcal{L}})^2 \\ & - (\mathcal{G}(mn(m + n) - 17(m^2 + n^2) + 23mn) + 18) H(z_{i+\mathcal{L}})^2H(z_{j+\mathcal{L}}) \\ & + 9(\mathcal{G}(3mn(m + n) - 26(m^2 + n^2) + 34mn) \\ & \quad \left. + 24)H(z_{i+\mathcal{L}})H(z_{j+\mathcal{L}})^2H(z_{k+\mathcal{L}}) \right. \\ & + 5(\mathcal{G}(3mn(m + n) - 26(m^2 + n^2) + 34mn) \\ & \quad \left. + 24)H(z_{i+\mathcal{L}})^2H(z_{j+\mathcal{L}})H(z_{k+\mathcal{L}}) \right. \\ & + (\mathcal{G}(3mn(m + n) - 26(m^2 + n^2) + 34mn) + 24) H(z_{i+\mathcal{L}})^2H(z_{j+\mathcal{L}})^2 \\ & - 5(\mathcal{G}(3mn(m + n) - 26(m^2 + n^2) + 34mn) \\ & \quad \left. + 24)H(z_{i+\mathcal{L}})^2H(z_{j+\mathcal{L}})^2H(z_{k+\mathcal{L}}) \right) \end{aligned} \quad (5.3)$$

where $\mathcal{V} = \frac{nm(n+m)^2}{n^3+m^3}$ and $\mathcal{G} = \frac{nm(n+m)^2}{n^5+m^5}$. Then the integrals for the mean, variance and third moments of the test statistic reduces to

$$\mathbb{E}[T_{n,m}] = \sum_{i=1}^{n+m-1} \frac{(ET_1 + ET_0)}{2} \cdot (z_{i+1} - z_i) \quad (5.4)$$

$$\text{Var}[T_{n,m}] = \sum_{i=1}^{n+m-1} \sum_{j=1}^i \frac{(VT_1 + VT_0)}{2} \cdot (z_{i+1} - z_i) \cdot (z_{j+1} - z_j) \quad (5.5)$$

$$\begin{aligned} \mathbb{E}[T_{n,m}^3] = & \sum_{i=1}^{n+m-1} \sum_{j=1}^i \sum_{k=1}^j \frac{(GT_1 + GT_0)}{2} \cdot (z_{i+1} - z_i) \cdot (z_{j+1} - z_j) \\ & \cdot (z_{k+1} - z_k). \end{aligned} \quad (5.6)$$

Now, whilst R is very quick at calculating the sum in Equation (5.4), it is very slow at dealing with Equations (5.5) and (5.6). The slowness comes from the double and triple nested summations in Equations (5.5) and (5.6). Thus to ensure faster calculations, all three equations have been coded in C++. Figure 5.2 shows the speed in seconds for calculating the mean, variance and third moment of the test statistic for $n = m \in [2, 250]$ when the equations are coded in C++.

Figure 5.2 shows that after coding in C++, the calculation for the mean and arguably the variance also, remains fast. However the third moment calculation is very slow for larger samples. Currently the grid points chosen in the summation to approximate the integral is the pooled data points, thus as $n + m$ increases, the calculations will get exponentially slower. However, the calculations can be made quicker by choosing fewer equally spaced grid points along the range $[z_1, z_{n+m}]$. It should be noted that whilst choosing fewer grid points will speed up the calculation, accuracy is sacrificed.

5.2.3 Reducing Number of Grid Points

Consider the equally spaced grid points in the range $[z_1, z_{n+m}]$ to be $\zeta_1, \dots, \zeta_{N_g}$. In Equations (5.1) to (5.6), we replace the z_1, \dots, z_{n+m} with $\zeta_1, \dots, \zeta_{N_g}$ and sum from $i = 1$ to $i = N_g$. The question which remains is what value of N_g will ensure the calculations are fast whilst still maintaining accuracy. To investigate the accuracy, we calculate the moments $\mathbb{E}[T_{n,m}]_z$, $\text{Var}[T_{n,m}]_z$ and $\mathbb{E}[T_{n,m}^3]_z$ using the data points as the grid points, we also calculate the moments $\mathbb{E}[T_{n,m}]_{N_g}$, $\text{Var}[T_{n,m}]_{N_g}$ and $\mathbb{E}[T_{n,m}^3]_{N_g}$ using N_g equally spaced grid points. The ratios of the moments calculated using N_g equally spaced grid points over the moments calculated using the data points as the grid points are then calculated. By calculating these ratios for $N_g \in [2, 250]$, a suitable value of N_g can be found which not only ensures a fast

5. Application of Two Sample Test

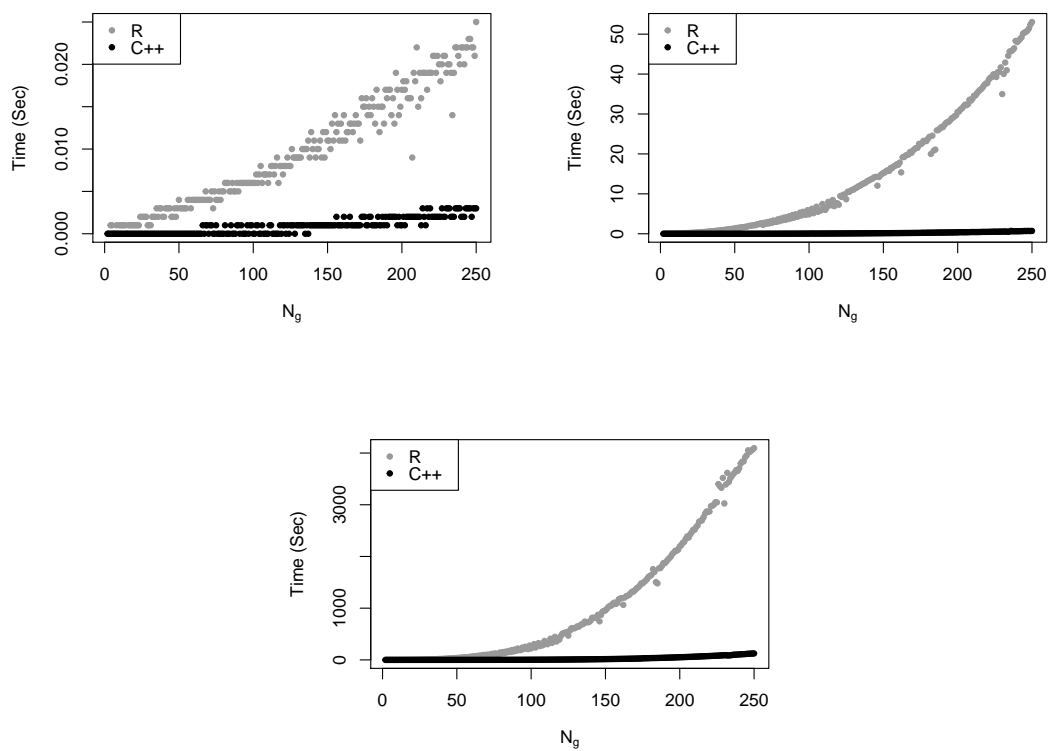


Figure 5.2: The time in seconds for calculating the mean (top left), variance (top right) and third moment (bottom) of the test statistic when $n = m \in [2, 250]$ using R and C++.

calculation but also an accurate estimate. The plots in Figure 5.3 show the ratios $\frac{E[T_{n,m}]_{N_g}}{E[T_{n,m}]_z}$ (top left), $\frac{\text{Var}[T_{n,m}]_{N_g}}{\text{Var}[T_{n,m}]_z}$ (top right) and $\frac{E[T_{n,m}^3]_{N_g}}{E[T_{n,m}^3]_z}$ (bottom) plotted against $N_g \in [2, 250]$.

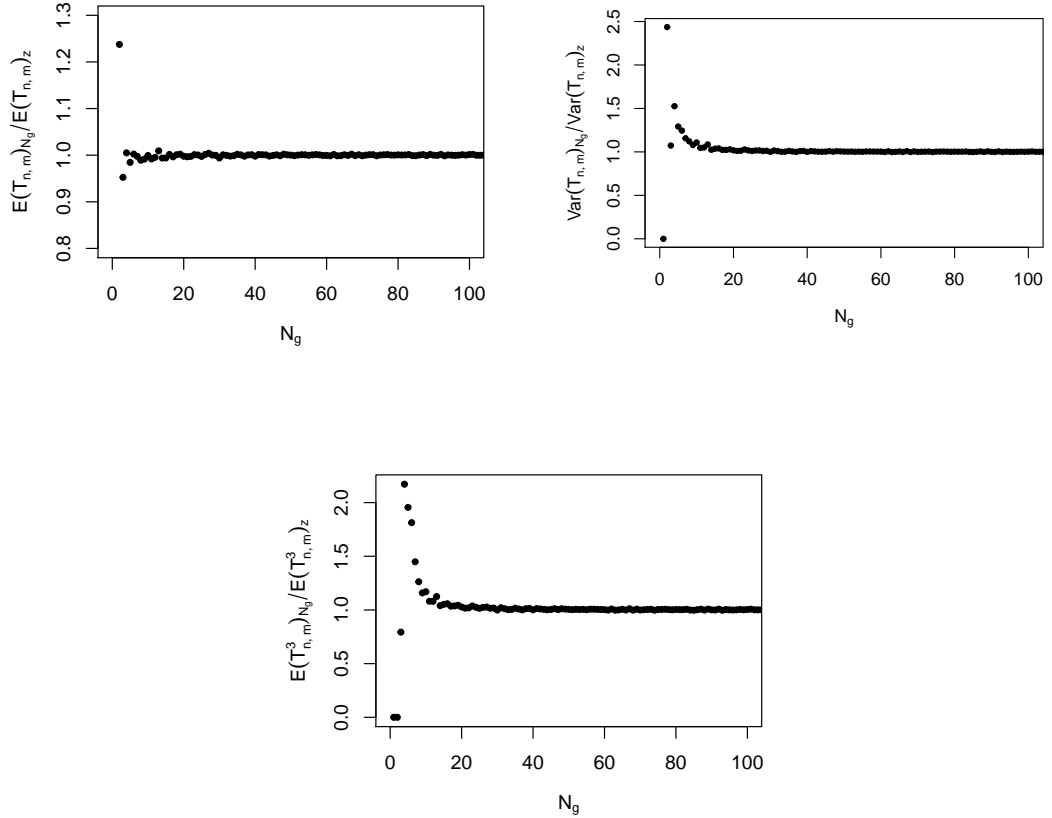


Figure 5.3: The ratios $\frac{E[T_{n,m}]_{N_g}}{E[T_{n,m}]_z}$ (top left), $\frac{\text{Var}[T_{n,m}]_{N_g}}{\text{Var}[T_{n,m}]_z}$ (top right) and $\frac{E[T_{n,m}^3]_{N_g}}{E[T_{n,m}^3]_z}$ (bottom) plotted against $N_g \in [2, 250]$.

Figure 5.3 shows that for $N_g > 20$, all three moments can be accurately estimated and Figure 5.2 shows that for $N_g = 20$ the calculation speed is close to 0. Thus we can conclude that $N_g = 20$ is a reasonable recommendation for the number of equally spaced grid points. Note however that $N_g = 20$ is only a suggestion and the ultimate decision of the number of grid points remains with the user.

5.3 Atest - an R Package

For the purpose of calculating the p -value for the Cramer test using our approach, namely fitting a GPD using the method of moments, we created an R package

5. Application of Two Sample Test

called `A.test`. The main function in the package, `A.test` will perform the one-dimensional two-sample Cramer test on two given samples of data and calculate the p -value using the GPD method.

5.3.1 The Functions

All the functions included in the package are listed below. For each function, its purpose, the arguments and the output are described.

`A.test(Data1, Data2, P.Value = T, GridPoints=50)` This is the main function in the package. It takes as input two samples, namely `Data1` and `Data2`, each of length n and m and will produce the one-dimensional two-sample Cramer test statistic (`Statistic`) and the p -value using the GPD method (`P.Value`). If `P.Value = F` the p -value will not be calculated thus ensuring a faster calculation of the test statistic alone. The argument `GridPoints` controls the number of grid points used in the calculation of the moments. A smaller number of grid points will mean a faster calculation but with less accuracy. Alternatively, a larger number of grid points will mean a slower calculation with more accuracy. The recommended number of grid points is shown to be 20.

`TestStatExpectation(sample, Regions)` This function calculates the expectation of the test statistic using Equation 5.1. It is used in `A.test` to calculate the p -value and is coded using C++ to enable a faster calculation. The argument `sample` is `Data1` and `Data2` pooled and `Regions` is the grid points. The grid points are obtained by taking a sequence of length `GridPoints` between the minimum value of `sample` and the maximum value of `sample`. Note that the lengths of `Data1` and `Data2` are not provided here as the expectation of the test statistic does not depend on n and m .

`TestStatVariance(n1, n2, sample, Regions)` This function calculates the variance of the test statistic using Equation 5.2. It is used in `A.test` to calculate the p -value and is coded using C++ to enable a faster calculation. The argument `sample` is `Data1` and `Data2` pooled and `Regions` is the grid points. The grid points are obtained by taking a sequence of length `GridPoints` between the minimum value of `sample` and the maximum value of `sample`. Here, `n1` and `n2` represents the length of `Data1` and `Data2` respectively.

`TestStatMoment3(n1, n2, sample, Regions)` This function calculates the third moment of the test statistic using Equation 5.3. It is used in `A.test` to calculate the p -value and is coded using C++ to enable a faster calculation. The argument `sample` is `Data1` and `Data2` pooled and `Regions` is the grid points. The grid points are obtained by taking a sequence of length `GridPoints` between the minimum value of `sample` and the maximum value of `sample`. Here, `n1` and `n2` represents the length of `Data1` and `Data2` respectively.

`Calc.Xi(xi, skewness)` This function is optimised in `A.test` for the purpose of estimating the generalised Pareto distribution parameter ξ . It's arguments are ξ , which is to be optimised, and `skewness` which is the skewness of the test statistic calculated using `TestStatMoment3`.

`Calc.Sigma(sigma, xi, variance)` This function is optimised in `A.test` for the purpose of estimating the generalised Pareto distribution parameter σ . It's arguments are σ , which is to be optimised, ξ , which is obtained from optimising `Calc.Xi` and `variance` which is the variance of the test statistic calculated using `TestStatVariance`.

`Calc.Mu(mu, xi, sigma, mean)` This function is optimised in `A.test` for the purpose of estimating the generalised Pareto distribution parameter μ . It's arguments are μ , which is to be optimised, σ , which is obtained from optimising `Calc.Sigma`, ξ , which is obtained from optimising `Calc.Xi` and `mean` which is the expectation of the test statistic calculated using `TestStatExpectation`.

5.3.2 An Example

Consider sampling $n = m = 100$ values from $X \sim N(0, 1)$ and $Y \sim N(5, 10)$. We can calculate the Cramer test statistic and the p -value by inputting the following code into R;

```
sample1 = rnorm(100,0,1)
sample2 = rnorm(100,5,10)
A.test(sample1, sample2, P.Value = T, GridPoints = 50)
```

which gives the following output;

```
$Statistic
[1] 141.9328
```

```
$P.Value
[1] 2.87028e-10
```

In this case, the Cramer test concludes that the samples from distributions X and Y are not equivalent, which was to be expected.

5.4 Comparing GPD Method to Bootstrap and Permutation Approach

One of the main reasons for investigating other methods of obtaining a p -value was to ensure a faster calculation which doesn't involve resampling. As we have found in Section 4.3.2, the GPD can accurately estimate the null distribution and parameters can be obtained using the formulas for the moments (Section 3.3). Now computational considered has been taken into account it is then natural to compare the results obtained through the GPD method, the bootstrap approach to approximate the limiting distribution which is considered in [Baringhaus and Franz \(2004\)](#) and the permutation method. We compare the speed of all three methods as well as the accuracy against the permutation approach.

5.4.1 Speed

Consider sampling X and Y from a mixture of normals distribution with probability density function $\frac{3}{4}N(0, 0.5) + \frac{1}{4}N(2, 0.5)$. For each simulation setting we take a value of $n = m \in [2, 200]$ and calculate a p -value using all three methods and record the time taken to calculate each p -value. We repeat this 5 times for each value of $n = m$ and take the average time over the 5 replications. For the permutation and bootstrap approach, the number of replications is 1,000 and for the GPD method, $N_g = 50$. In Figure 5.4 we record the amount of time required if 200,000 simultaneous hypothesis tests are performed.

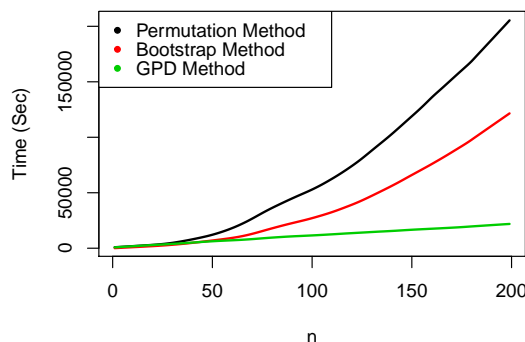


Figure 5.4: The average speed over 5 replications if 200,000 simultaneous hypothesis tests are performed using the permutation method, the bootstrap approach and the GPD method with $N_g = 50$ to calculate the p -value. In this scenario two samples are drawn from a mixture of normals distribution with probability density function $\frac{3}{4}N(0, 0.5) + \frac{1}{4}N(2, 0.5)$ and $n = m \in [2, 200]$.

5.4 Comparing GPD Method to Bootstrap and Permutation Approach

Clearly, for $n = m > 40$ the GPD method is faster than both the bootstrap and permutation approach. Also for $n = m = 200$, the GPD method still takes approximately 0.1 seconds to compute a single p -value, which corresponds to a total time of just under 6 hours to calculate 200,000. If we were to calculate 200,000 p -values using the bootstrap approach it would take approximately 33 hours and approximately 57 hours for the permutation approach. Thus clearly there exists a large difference in time between the GPD method and the resampling approaches when the number of simultaneous hypothesis tests required is large. Therefore, it can be argued that to ensure a more user friendly test - especially if the user does not have access to parallel processing - the GPD method is the most preferable.

5.4.2 Accuracy

It is one thing to have a faster method for calculating the p -value, but to ensure that the GPD method is just as accurate as the other two methods is important. Therefore, to test the accuracy of all three methods, sample X and Y 100 times from a mixture of normals distribution with probability density function $\frac{3}{4}N(0, 0.5) + \frac{1}{4}N(2, 0.5)$ and $n = m = 50$. For each simulation, the p -value for each method is obtained and recorded. We are only concerned with the accuracy of the p -values when they are less than 0.10, as it is more imperative that the p -values are accurate in the right tail. Thus when we calculate the p -values we only store p -values which are less than 0.10 for all methods.

Here, the number of replications for the permutation and bootstrap approach is 10,000 and again $N_g = 50$ for the GPD method. We increase the number of replications for the permutation and bootstrap approach to ensure a more accurate p -value. Figure 5.5 (top left) shows the ratio of the p -values which are less than 0.10 calculated using the GPD method over the permutation approach, Figure 5.5 (top right) shows the ratio of the p -values which are less than 0.10 calculated using the GPD method over the bootstrap approach and finally Figure 5.5 (bottom) shows the ratio of the p -values which are less than 0.10 calculated using the bootstrap approach over the permutation approach.

As the points are randomly scattered around the line $y = 1$ in Figure 5.5 (top left), it can be concluded that the p -values calculated using the GPD method are just as accurate as the p -values calculated using the permutation approach. Similar conclusions can be made from Figure 5.5 (top right) and (bottom). The mean ratios across the 100 simulations are calculated to be 1.005, 0.999 and 1.008 respectively. Note also that if the p -values were corrected using multiple testing then the ratios between the p -values calculated using different methods would not be affected. Thus the graphs in Figure 5.5 provide a good representation of the

5. Application of Two Sample Test

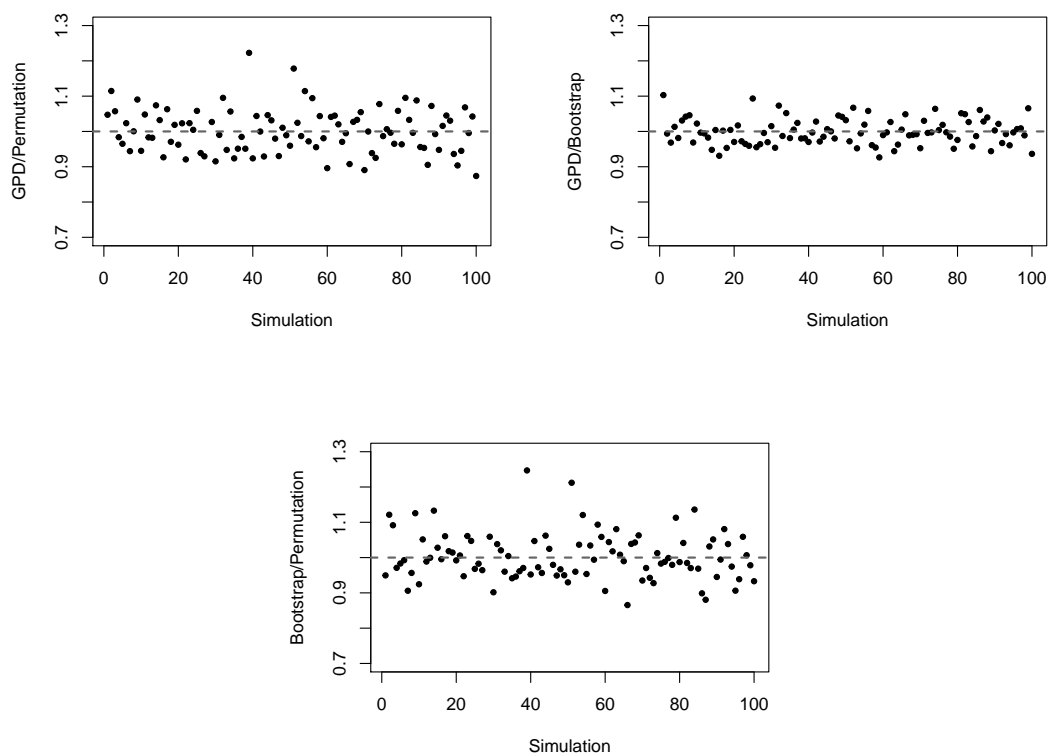


Figure 5.5: Top Left: The ratio of the p -values which are less than 0.10 calculated using the GPD method over the permutation approach. Top Right: The ratio of the p -values which are less than 0.10 calculated using the GPD method over the bootstrap approach. Bottom: The ratio of the p -values which are less than 0.10 calculated using the bootstrap approach over the permutation approach. Here, the sample size is $n = m = 50$, the number of replicates for the permutation and bootstrap approach is 10,000 and for the GPD method $N_g = 50$.

accuracy of the p -values regardless of their magnitude.

5.5 Simulation Study

5.5.1 Type-I error control

To show that the Cramer test has a proper false positive rate (FPR) control using our method to calculate the p -value, four simulations have been done under the null hypothesis. In each simulation, two samples of 100 observations are drawn from a skewed normal distribution $SN(\mu, \sigma, \alpha)$, where μ , σ and α represent the location, scale and shape parameters respectively. A skewed normal distribution was chosen here as we could control not only the mean and variance of the distribution but also the skewness, which was important because it better reflected the shape of the estimated CNA in each window and subtype of cancer. In each simulated dataset, we perform the Cramer test and calculate the corresponding p -value with $N_g = 50$ by estimating the null distribution using the GPD. We repeat this 100,000 times and calculate the proportion of p -values which are less than 0.05.

Each simulation considers the false positive rate of the Cramer test when some parameters vary while the other parameters are fixed. The first simulation considers the FPR when μ varies in the range $[0,1]$ with $\sigma = 1$ and $\alpha = 0$, the second one when σ varies in the range $[0.1,1]$ with $\mu = 0$ and $\alpha = 0$, the third one when α varies in the range $[0,1]$ with $\mu = 0$ and $\sigma = 1$, and the final one when both σ and α vary in the range $[0.1,1]$ with $\sigma = \alpha$ and $\mu = 0$.

For the purpose of comparison, we also calculate the FPR for other tests: (two-sample) t -test, Cramer-von Mises test, Anderson-Darling test, F -test, and Kolmogorov-Smirnov test. The FPR figures for those tests in each simulation are presented in Figure 5.6.

Figure 5.6 indicates that the different tests in the simulation manage to control FPR properly, except for the Kolmogorov-Smirnov test which exhibits lower FPR than the other tests. It can be shown that, when the sample size increases to be much larger than $n = m = 100$, the false positive rate for the Kolmogorov-Smirnov test converges to 0.05. This therefore implies that, for a small sample size, the Kolmogorov-Smirnov test is more likely to fail to reject windows which are significant.

Consider now sampling X and Y from a multi-modal mixture distribution which follows $p_1N(1, 1) + p_2N(1 + d, 1)$. We now perform simulations to ensure that the FPR is properly controlled when the data is multi-modal. As multi-modality is a common feature of the lung cancer data set, it is vital that the FPR is properly controlled in this case. The first simulation samples 100 observations

5. Application of Two Sample Test

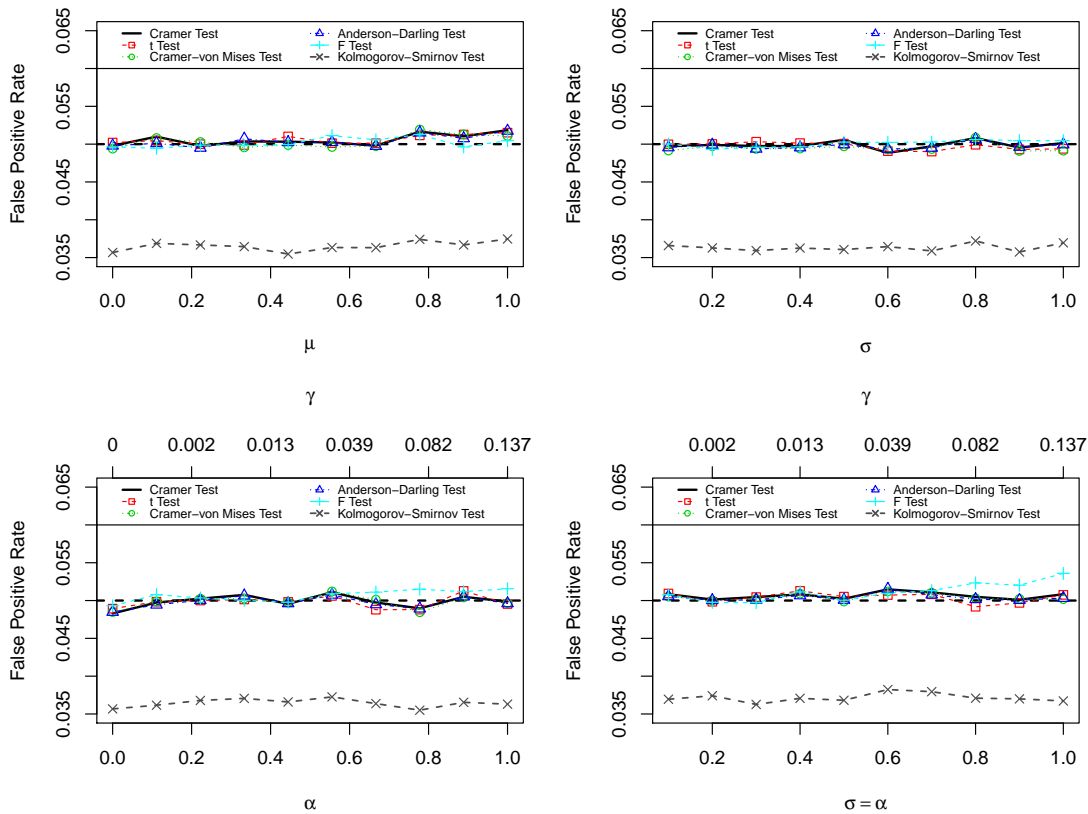


Figure 5.6: False positive rates for the Cramer test, (two-sample) t -test, Cramer-von Mises test, Anderson Darling test, F -test, and Kolmogorov-Smirnov test at different simulation settings: varying μ (top left panel), σ (top right panel), α (bottom left panel), and both α and σ with $\alpha = \sigma$ (bottom right panel), from skew-normal distribution (see Section A.0.1) $SN(\mu, \sigma, \alpha)$. In the bottom row figures, the values of α are within the interval $[0, 1]$ (bottom horizontal axis), which have been accompanied by the corresponding values of skewness γ (top horizontal axis).

from X and Y where $p_1 = \frac{7}{8}$, $p_2 = \frac{1}{8}$ and d varies in the range $[1,11]$. The second simulation samples 100 observations from X and Y where $p_1 = 1 - p_2$, p_2 varies in the range $[0.01,0.5]$, and $d = 2$. In each simulated dataset, we perform the Cramer test and calculate the corresponding p -value with $N_g = 50$ by estimating the null distribution using the GPD. We repeat this 1,000 times and calculate the proportion of p -values which are less than 0.05.

For the purpose of comparison, we also calculate the FPR for other tests: (two-sample) t -test, Cramer-von Mises test, Anderson-Darling test, F -test, and Kolmogorov-Smirnov test. The FPR figures for those tests in each simulation are presented in Figure 5.7.

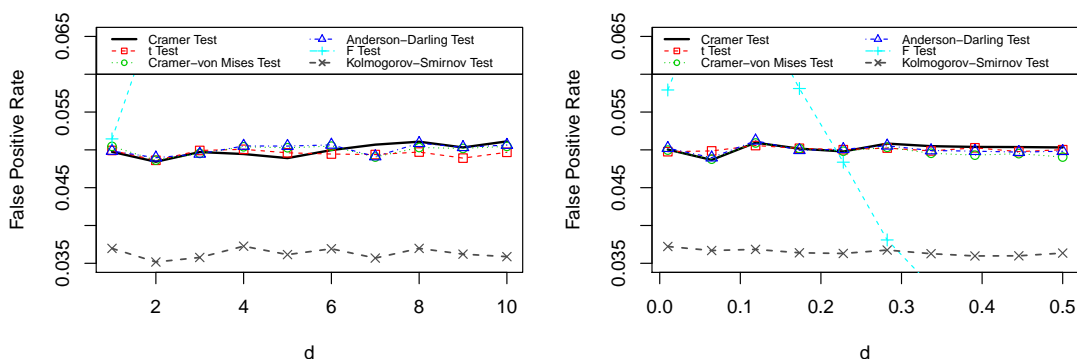


Figure 5.7: False positive rates for the Cramer test, (two-sample) t -test, Cramer-von Mises test, Anderson Darling test, F -test, and Kolmogorov-Smirnov test at different simulation settings: varying d (left panel) and p_1 (right panel) from a multi-modal mixture distribution which follows $p_1N(1, 1) + p_2N(1 + d, 1)$.

Figure 5.7 shows that the false positive rate is still controlled even when the distributions for X and Y are still multi-modal. This means that we can be confident that there aren't too many or indeed too few significant results being obtained using the GPD to calculate the p -value. One interesting feature of the graphs in Figure 5.7 is that the F -test does not properly control the false positive rate when the data is multi-modal. This suggests that if the F -test were to be performed on the lung cancer data set, there may be too many regions being incorrectly identified as significant.

5.5.2 Sensitivity

To investigate the power of the Cramer test using our method for calculating the p -value, four simulations have been performed under the alternative hypothesis. Specifically, in each simulated dataset, a sample of 100 observations are drawn

5. Application of Two Sample Test

from a skewed normal distribution $SN(\mu, \sigma, \alpha)$, and in the second sample, another 100 observations are drawn from a $SN(0, 1, 0)$ distribution. In each simulated dataset, we perform the Cramer test and calculate the corresponding p -value with $N_g = 50$. We repeat this 100,000 times and calculate the proportion of p -values which are less than 0.05.

The first simulation considers the power when μ varies in the range $[0,1]$ with $\sigma = 1$ and $\alpha = 0$, the second one when σ varies in the range $[0.1,1]$ with $\mu = 0$ and $\alpha = 0$, the third one when α varies in the range $[0,1]$ with $\mu = 0$ and $\sigma = 1$, and the final one when both σ and α vary in the range $[0.1,1]$ with $\sigma = \alpha$ and $\mu = 0$. The choice of these settings were to investigate the power in situations in which (1) only the mean differs, (2) only the variance differs, (3) only the skewness differs and (4) only the variance and skewness differs between the two distributions. Particularly in the setting of (1) and (2), the Cramer test will be compared to the t -test and F -test, which are the “gold standard” tests for these scenarios.

As in the previous simulation, we also calculate the sensitivity for other tests for comparison to the Cramer test and these are presented in Figure 5.8. It can be seen that the Cramer test has a good sensitivity to detect differences in mean, variance, skewness, and joint skewness and variance between two samples. As expected, (two sample) t -test is powerful to detect differences in mean, but not the other parameters. Similarly, the F -test is powerful to detect differences in the variance (right column of Figure 5.8), but not the mean nor skewness. The Cramer test has the same sensitivity with either the Anderson-Darling test or Cramer-von Mises test. As expected, the Cramer test, AD test and CvM test is less powerful than the F -test when only the variance differs. There are some other situations in which the Cramer test is superior to the Anderson-Darling and Cramer-von Mises tests, and more superior than the F -test. We consider two additional simulations to highlight this.

Firstly, consider now a simulation in which 100 observations have been drawn from a $N(1, 1)$ distribution in the first sample, and 100 observations from a multi-modal mixture distribution which follows $(1 - p_1)N(1, 1) + p_1N(1 + d, 1)$. The sensitivity of the different tests are then calculated from 1,000 simulated datasets in two cases: (1) d varies in the range $[0,9]$ and p_1 is fixed at $\frac{1}{8}$, and (2) p_1 varies in the range $[0.01,0.5]$, d is fixed at 2. Figure 5.9 shows the results of the simulation.

Figure 5.9 shows that the Cramer test has better sensitivity than both the Anderson-Darling test and the Cramer-Von Mises test. It is clear however that the F -test performs better and the t -test perform just as well as the Cramer test in the simulation. This is because the simulation setting in both cases inevitably give a difference in mean and variance between the two samples.

Secondly, consider now a simulation in which 100 observations are drawn from

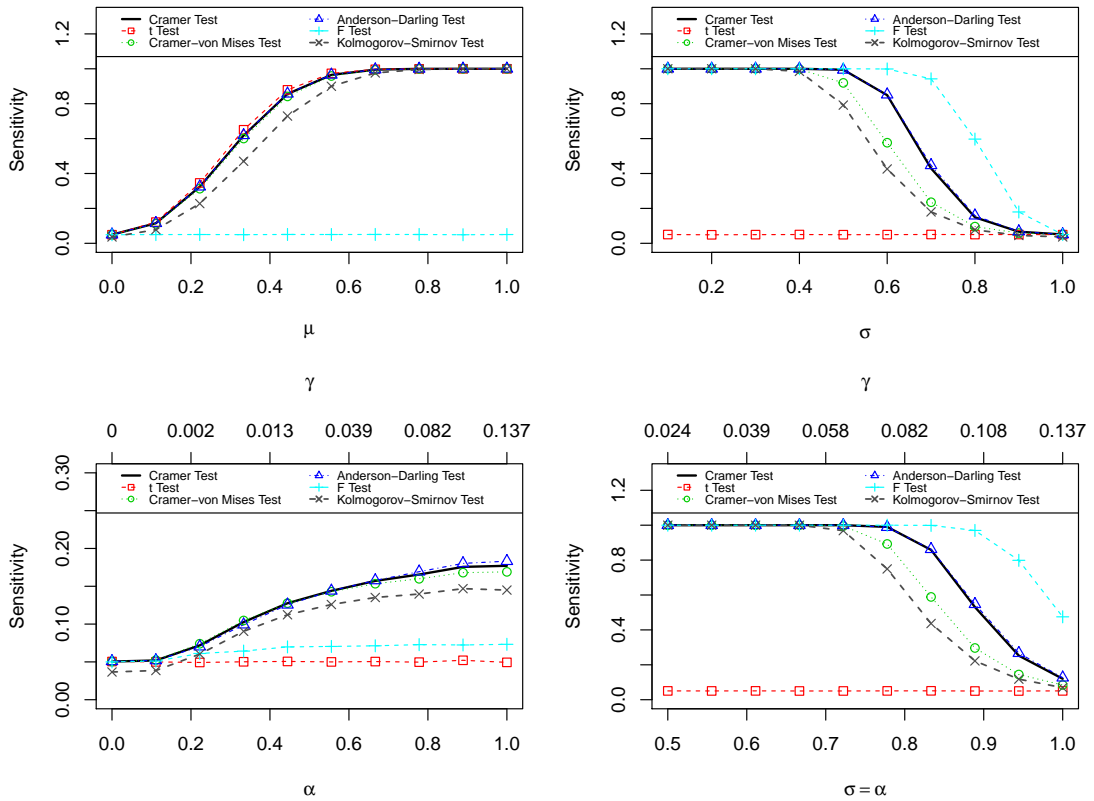


Figure 5.8: Sensitivity for the Cramer test, (two-sample) t -test, Cramer-von Mises test, Anderson Darling test, F -test, and Kolmogorov-Smirnov test at different simulation settings: varying μ (top left panel), σ (top right panel), α (bottom left panel), and both α and σ with $\alpha = \sigma$ (bottom right panel), from skew-normal distribution $SN(\mu, \sigma, \alpha)$ in the first sample. In the second sample, the observations are drawn from $SN(0, 1, 0)$. In the bottom row figures, the values of α are within the interval $[0, 1]$ (bottom horizontal axis), which have been accompanied by the corresponding values of skewness γ (top horizontal axis).

5. Application of Two Sample Test

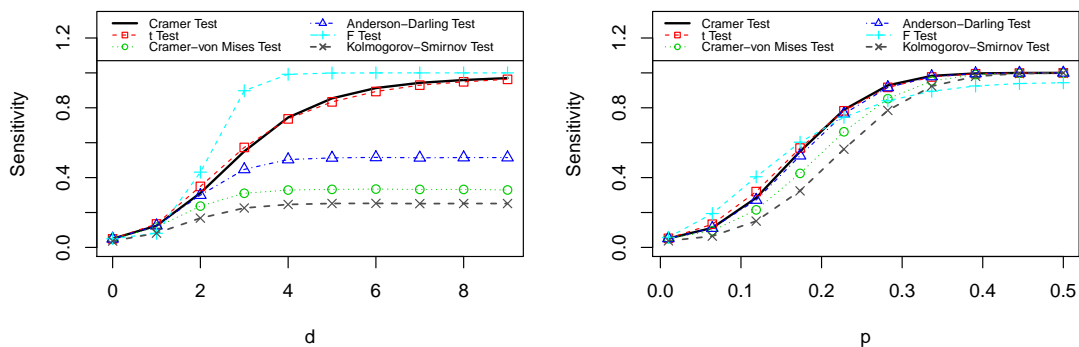


Figure 5.9: Sensitivity for the Cramer test, (two-sample) t -test, Cramer-von Mises test, Anderson Darling test, F -test, and Kolmogorov-Smirnov test, where 100 observations are drawn from $N(1, 1)$ in the first sample, and 100 observations from a multi-modal mixture distribution which follows $(1 - p_1)N(1, 1) + p_1N(1 + d, 1)$ in the second sample. The left panel is the setting where d varies in the range $[0, 9]$ and p_1 is fixed at $\frac{1}{8}$. The right panel is the setting where p_1 varies in the range $[0.01, 0.5]$, d is fixed at 2.

a $N(0, 10)$ distribution in the first sample, and 100 observations are drawn from a $N(-d, \sigma)^\pi \cdot N(d, \sigma)^{1-\pi}$ distribution where $\pi \sim \text{Bernoulli}(\frac{1}{2})$, $d \in [7, 10]$ and $\sigma = \sqrt{100 - d^2}$ in the second sample. We now have the situation where the mean and variance between the two samples do not differ. Figure 5.10 shows the results of this simulation.

Figure 5.10 indicates that the Cramer test has a good sensitivity across different values of d , along with the Anderson-Darling test, Cramer-von Mises test, and the Kolmogorov-Smirnov test. In this setting, the t -test and F -test cannot recover sensitivity because the simulation setting implies that there is no mean nor variance difference between the two samples.

Clearly, the Cramer test is highly competitive against other similar two-sample parametric and non-parametric tests when the p -value is calculated with $N_g = 50$. We now compare our two-sample method for locating genomic regions of interest to KC Smart.

5.6 Genomic Results

We now apply the Cramer test to each and every genomic region in the lung cancer data to test the null hypothesis that the distribution of CNA are equal in both pathological subtypes. As a comparison, we also consider the (two-sample) t -test, F -test, Kolmogorov-Smirnov test, Anderson-Darling test, and Cramer-von Mises

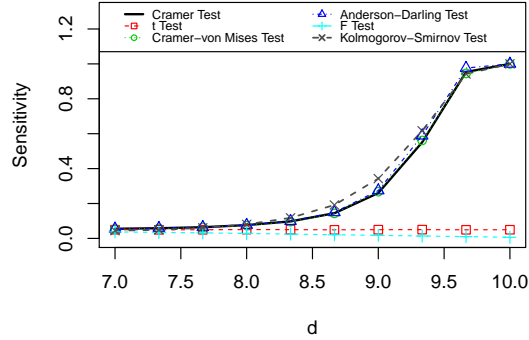


Figure 5.10: Sensitivity for the Cramer test, (two-sample) t -test, Cramer-von Mises test, Anderson Darling test, F -test, and Kolmogorov-Smirnov test, where 100 observations are drawn from $N(0, 10)$ in the first sample, and 100 observations from $N(-d, \sigma)^\pi \cdot N(d, \sigma)^{1-\pi}$ distribution where $\pi \sim \text{Bernoulli}(\frac{1}{2})$, $d \sim [7, 10)$ and $\sigma = \sqrt{100 - d^2}$, in the second sample. The sensitivity figures are plotted as a function of d .

test with their respective null hypotheses. The number of genomic regions which have p -value less than 5% are presented in Table 5.1 for each test. The table indicates that the Cramer test is able to detect more significant regions than each of KS, AD, and CvM tests. Note that there are approximately 92.6% (KS), 94.5% (AD), and 93.4% (CvM) of the regions which are in common with the Cramer test. This indicates that the Cramer test is able to capture almost all significant regions identified by those tests, and some more. Given that the false positive rates were shown to be controlled for the Cramer test in scenarios where X and Y are multi-modal, we can be confident that the extra regions being identified as significant by the Cramer test are not simply false positives and are instead significant regions which are being missed by the other tests. The F -test is clearly able to identify more significant regions than any other test, however only has 69.7% of the regions which are in common with the Cramer test. It could be argued that the F -test is identifying too many regions to be significant. We have shown in Figure 5.7 that when the data is multi-modal the F -test does not properly control the false positive rate.

5.6.1 Results of Cramer Test

Recall that previous studies (Björkqvist et al. (1998), Wang et al. (2013) and van Boerdonk et al. (2011)) already identified a large gain in chromosome 3 for patients with squamous carcinoma type lung cancer. Therefore we can use this knowledge to test the accuracy of our test. If indeed a region of chromosome 3 is found to

5. Application of Two Sample Test

	Cramer test	t -test	F -test	KS test	AD test	CvM test
Cramer test	7,981	6,464	5,567	6,298	7,435	7,061
t -test		6,923	4,793	5,116	6,119	5,777
F -test			10,909	4,474	5,285	4,937
KS test				6,800	6,480	6,481
AD test					7,868	7,342
CvM test						7,558

Table 5.1: The number of genomic regions out of 17,613 with (unadjusted) p -values less than 0.05 under the Cramer test, the t -test, the F -test, the KS test, the AD test and the CvM test in our lung cancer dataset. The (i, j) th entry indicates the number of significant windows in both the i th and j th test.

have a significant difference in CNA between the two groups of patients then this is evidence that our method is correctly identifying regions which differ in CNA.

Figure 5.11 presents the p -values of individual genomic regions across the genome. Some of the p -values exceed the Bonferroni corrected significance threshold, and the corresponding genomic regions are considered significant. In our lung cancer data set, we have a total of 669 significant regions. The large significant region spans windows 4045-4603 within chromosome 3 in the genome. This therefore confirms that the Cramer test is correctly identifying regions which differ in CNA.

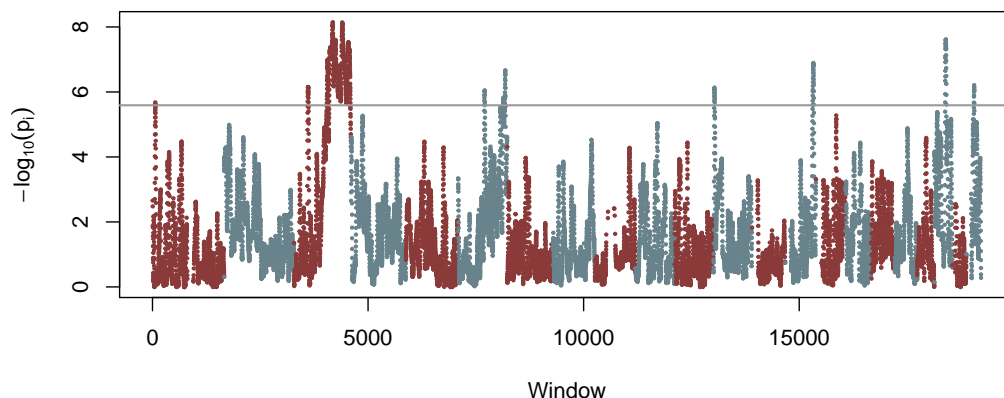


Figure 5.11: The p -values of individual genomic regions across the genome using the Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey line represents the Bonferroni corrected significance threshold for which 669 regions pass this threshold. The alternating colouring scheme indicates chromosomes 1–22 from the left. Sex chromosomes are excluded from the analysis.

Now, as can be also seen from Figure 5.11 there exists regions in other chromo-

somes which were also identified as significant by the Cramer test. Many of these regions have not been previously identified which could be an indication that the Cramer test is much more sensitive at identifying differences in smaller regions of the genome - a massive advantage of our method. Because we have confirmed in Section 5.5.1 that the false positive rate is being properly controlled we can assume that these extra regions being identified as significant are not primarily false positives and therefore conclude that these regions are in fact regions which differ in CNA between the two groups. Perhaps this new insight into the data could be the key to better classification of future patients, or indeed better treatments for those patients.

We have, of course, only applied our method to the lung cancer data set for which many studies have been done. As our method proves to be highly apt in locating regions of the genome - no matter how big or small - which displays a significant difference in CNA between groups of patients, it could also be applied to compare clinical groups of patients with different cancers for which no previous studies have been done and in the future we hope to do just that.

5.6.2 Results of KC Smart

We now apply KC Smart to the lung cancer data set to compare the results to applying the Cramer test to each region. We hope that KC Smart identifies similar regions of interest to confirm that the Cramer test is performing correctly. Note that we have already identified advantages of our method over KC Smart, so even if it appears that both tests are equally as capable we can still argue that our method is preferable. Figure 5.12 presents the results after applying KC Smart to the lung cancer data.

By comparing the output of the KC Smart analysis in Figure 5.12 to the output of the Cramer test in Figure 5.11 it is clear that both methods are able to identify similar significant regions. For example, both methods identify regions of significance in chromosome 3, 6, 12, 14, 20 and 22. This provides further evidence that the Cramer test is once again identifying significant regions correctly. Table 5.2 shows the number of genomic windows which are identified as significant and not significant for both the Cramer test and KC Smart. For the Cramer test, windows are identified as significant if they are larger than the Bonferroni corrected significance threshold.

Table 5.2 shows that KC Smart is able to identify further significant regions compared to the Cramer test. It is unclear whether the extra significant regions identified by KC Smart are in fact false positives or whether these regions should be correctly identified and are being missed by the Cramer test. To see whether

5. Application of Two Sample Test

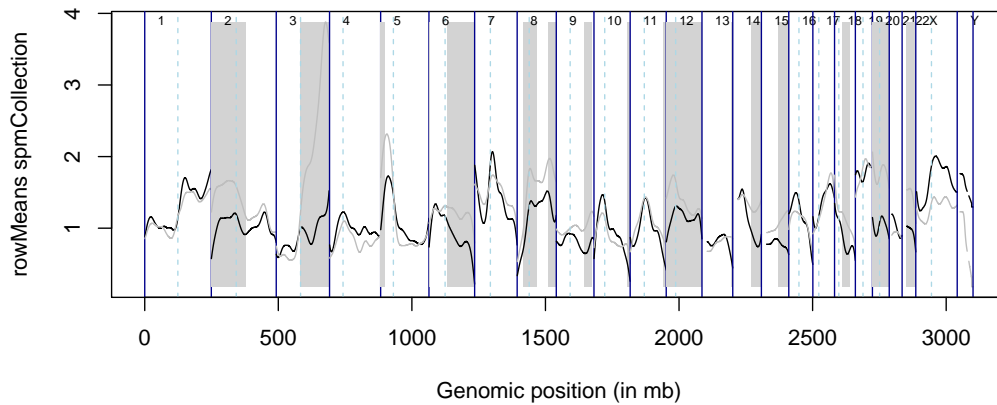


Figure 5.12: The output after applying KC Smart to the lung cancer data set. The same mirror locations were used as the artificial data set from the KC Smart vignette.

		Cramer Test	
		Significant	Not Significant
KC Smart	Significant	645	2981
	Not Significant	24	13963

Table 5.2: The number of significant genomic regions after Bonferroni correction out of 17,613 when using the Cramer test and KC Smart.

KC Smart consistently identifies further regions which differ significantly between two groups of patients compared to the Cramer test, we apply both methods to the artificial data from the KC Smart vignette.

5.7 Application to KC Smart Data

Recall Figures 1.8 and 1.9 which display the output of applying KC Smart to the artificial data in the KC Smart vignette. Both Figures clearly identified a section of chromosome 4 as significant. We now apply the Cramer test to each region of this data and make comparisons between the results. Figure 5.13 presents the p -values of individual regions across the genome for this artificial data set.

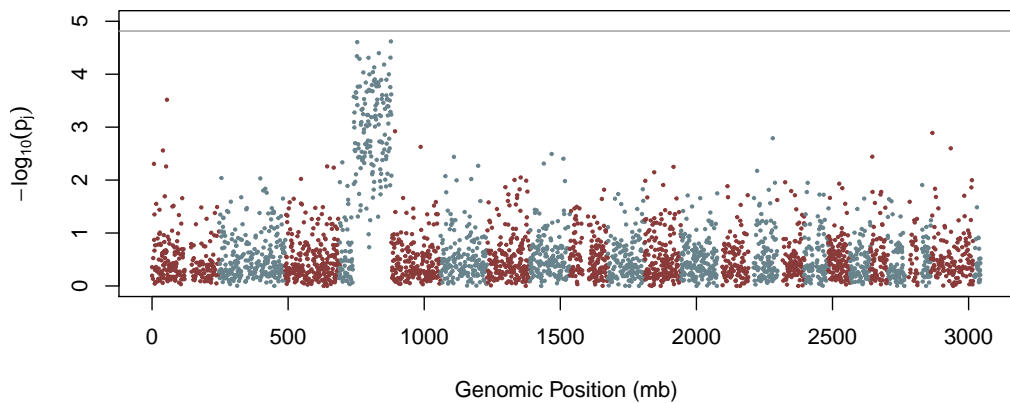


Figure 5.13: The p -values across regions in the genome using the Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey line represents the Bonferroni corrected significance threshold. The alternating colouring scheme indicates chromosomes 1–22. Sex chromosomes are excluded from the analysis.

Clearly, Figure 5.13 shows that no regions are significant as no points lie above the Bonferroni corrected significance threshold. This could be providing further evidence which suggests that KC Smart is overestimating the significance of some regions and in fact the significance of chromosome 4 is a false positive. Further work is needed here to investigate the false positive rate of KC Smart.

Let's now assume that the results of KC Smart is correct meaning that the Cramer test could be missing vital significant regions. One possible explanation of this is due to the multiplicity correction we perform on the results. As we are performing 17,613 simultaneous results, a multiplicity correction is required. Here we choose to perform the Bonferroni correction. Most multiplicity corrections assume that all tests performed are independent and the Bonferroni correction

5. Application of Two Sample Test

divides the significance threshold of a single test by the number of simultaneous independent tests performed. If indeed our tests are not independent then it could be that the Bonferroni significance threshold is too high and many regions are therefore being considered insignificant. We investigate this further by taking any potential correlation into account prior to performing our test.

5.7.1 Segmenting the Data

If two probes are highly correlated, then only a single test could be required to determine the significance of both probes. If “segments” of correlated probes are identified across the genome and only a single test performed on each then not only will noise be removed from the data, but the number of “independent” tests performed will be lower. If the number of “independent” tests performed is lower, the significance threshold will also be lower; therefore, more regions could be considered significant. To investigate whether correlation has an effect on the outcome of the Cramer test, consider segmenting the artificial data using circular binary segmentation (CBS) (Olshen et al., 2004) (Section 1.6.2). CBS has been performed on the artificial data set in the KC Smart vignette. Figure 5.14 shows the results of segmenting the first sample.

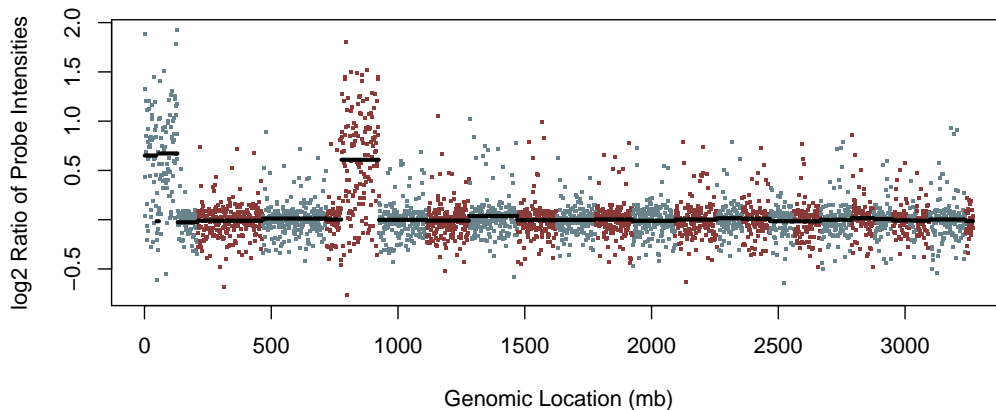


Figure 5.14: The first sequence in the artificial KC Smart data. The black line represents the segment means. The alternating colouring scheme indicates chromosomes 1–22, X, Y.

Each sample, i , has been split into different segments with genomic location start points $S_1^i, \dots, S_{N_i}^i \in \mathcal{S}_i$, and end points $E_1^i, \dots, E_{N_i}^i \in \mathcal{E}_i$ with N_i varying across samples. To be able to apply the Cramer test, all samples need to be segmented in the same way. To do this, two new sets \mathcal{S}_M and \mathcal{E}_M has been created

such that $\mathcal{S}_M = \mathcal{S}_1 \cup \dots \cup \mathcal{S}_{20}$ and $\mathcal{E}_M = \mathcal{E}_1 \cup \dots \cup \mathcal{E}_{20}$. The number M will be the total number of unique genomic location start points and end points across the samples. Now, let $S_j \in \mathcal{S}_M$ be such that $S_j \notin \mathcal{S}_i$ and $E_j \in \mathcal{E}_M$ be such that $E_j \notin \mathcal{E}_i$. Thus, there must exist an $S^* \in \mathcal{S}_i$ and a $E^* \in \mathcal{E}_i$ such that $S^* < E_j < S_j < E^*$. Hence, to add S_j and E_j to the sets \mathcal{S}_i and \mathcal{E}_i , we would split the segment with start and end point S^* and E^* respectively and create two new segments with the same segment mean. The first segment will have start point S^* and end point E_j , the second segment will have start point S_j and end point E^* . Elements from \mathcal{S}_M and \mathcal{E}_M will be added to the \mathcal{S}_i and \mathcal{E}_i for all i until all sets have the same segments.

Once all sets have the same segments, the Cramer test can be applied to compare the segment means between the two groups for each segment. Figure 5.15 shows the results after applying the Cramer test to each segment.

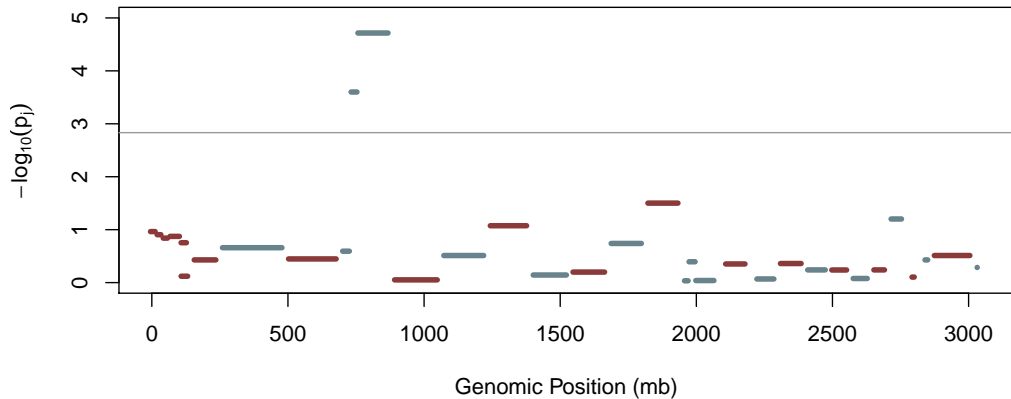


Figure 5.15: The p -values of each segment using the Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey line represents the Bonferroni corrected significance threshold. The alternating colouring scheme indicates chromosomes 1–22, X, Y.

Figure 5.15 shows that after segmenting the data and applying the Cramer test, chromosome 4 is now considered to be significant. This is because less tests have been performed so the significance threshold is lower. Also, the significance of the segments in chromosome 4 is higher because noise has been removed from the data. We can therefore conclude that the correlation between probes play an important part in the significance of genomic regions.

5.8 Discussion

In order to maintain a fast and accurate calculation for the test statistic and p -value, computational time and accuracy was considered. By changing the integration into summations, coding in C++ and reducing the number of grid points, the calculation remained fast without sacrificing too much accuracy. In order to perform analysis using our method, we have created an R packages called *Atest*, which is available from <http://www1.maths.leeds.ac.uk/~arief/R/Atest>. The package performs the Cramer test on two samples of data and calculates the p -value using our method of fitting a generalised Pareto distribution to estimate the null distribution.

The simulation study showed that the false positive rate for the Cramer test is properly controlled as well as identifying scenarios in which the Cramer test is preferable to other parametric and non-parametric tests. By comparing the Cramer test to the Cramer-von Mises and Anderson-Darling tests, further justifications to prefer the Cramer test in application were made.

After performing the Cramer test on the lung cancer data, chromosome 3 was identified as a significant genomic region. Prior to our study, it was known that chromosome 3 is a significant region for this data set, showing that our method can correctly identify significant regions. The results of our method was also compared to performing KC Smart on the data and showed that all significant regions identified by the Cramer test were also identified by KC Smart. However, KC Smart was also able to identify further regions of interest that requires further research to determine whether KC Smart is identifying non significant regions as significant.

The Cramer test was also compared to KC Smart. After applying both KC Smart and the Cramer test to the artificial data supplied in the KC Smart vignette, it showed that KC Smart identified chromosome 4 as a significant region whereas the Cramer test did not. We show, through use of the circular binary segmentation (CBS) technique, that a possible cause for the potential mis-identification of chromosome 4 is because of the high correlation between probes.

Because we have discovered that there could be some significant regions of the genome being missed by the Cramer test due to the correlation within the data. It is important to understand firstly whether high correlation exists within the lung cancer data set and secondly whether methods exist to account for such correlation. This is the main focus of our research in Chapters 6 and 7.

Chapter 6

Examination of Correlation Structure in the Data

6.1 Introduction

In large data sets such as the lung cancer data set, there exists many different types of correlation. For example, correlation could occur between variables or windows as well as between observations or patients. It is important to understand the correlation structures when performing any kind of analysis. In previous chapters, we have applied the Cramer test independently on each window, thereby assuming complete independence between windows. This assumption may not be true; therefore, this chapter explores the correlation structures which can exist within data sets and specifically the lung cancer data set. We start by discussing ways in which the correlation structures could be modelled and then investigate and provide examples of the different kinds of correlation structures which exist within the lung cancer data set.

6.1.1 Notation

In the sections which follow we use the following notation to represent the correlation within the data.

- The notation ρ^{pa} and ρ^{ps} represents the correlation in CNA between patients with adenocarcinoma type lung cancer only and squamous carcinoma type lung cancer only, respectively.
- The notation ρ^p represents the correlation between patients with any type of lung cancer.

6. Examination of Correlation Structure in the Data

- The notation ρ^{wa} and ρ^{ws} represents the correlation in CNA between the windows of the genome only for data from patients with adenocarcinoma type lung cancer and squamous carcinoma type lung cancer respectively.
- The notation ρ^w represents the correlation in CNA between the windows of the genome when data from both groups of patients are considered.
- The notation ρ^{pval} represents the correlation in p -values between the windows of the genome.

6.2 Modelling the Correlation Structure Between Variables

When performing a hypothesis test on p variables simultaneously, understanding the correlation structure of the variables is important. It would therefore be advantageous if we could model the correlation structure of the variables. The models can be used for predicting correlation between variables, simulating a correlation structure similar to that of the lung cancer data set as well as giving us insight into the structure itself. Note that hidden Markov models which were mentioned in Section 1.6.2 could be used to model the correlation structure within the data - this is done in [Rabiner \(1989\)](#). We choose however to focus on simpler non-bayesian models as modelling the correlation structure is not one of the main focuses in this thesis. Indeed, further work could be done in this area to identify the best method for modelling such correlation structures.

6.2.1 Autocorrelation

In time series analysis, autocorrelation is the Pearson correlation coefficient calculated between two times ([Box et al., 2015](#)). Let X_t be a realisation of a random process X evaluated at times $t \in \mathbb{N}$ with mean μ and variance σ^2 independent of time t . The autocorrelation function with lag τ is then

$$R(\tau) = \frac{\mathbb{E}[(X_t - \mu)(X_{t+\tau} - \mu)]}{\sigma^2}.$$

We can calculate the autocorrelation at lag τ for the lung cancer data set by considering t to be windows of the genome thereby calculating spatial autocorrelation instead of time autocorrelation. For a single patient with adenocarcinoma type lung cancer (top) and squamous carcinoma type lung cancer (bottom), [Figure 6.1](#) shows the autocorrelation function evaluated at $\tau \in [1, 100]$. Note that a lag

6.2 Modelling the Correlation Structure Between Variables

of $\tau = 1$ may not correspond to adjacent windows due to some windows being removed because of lack of data.

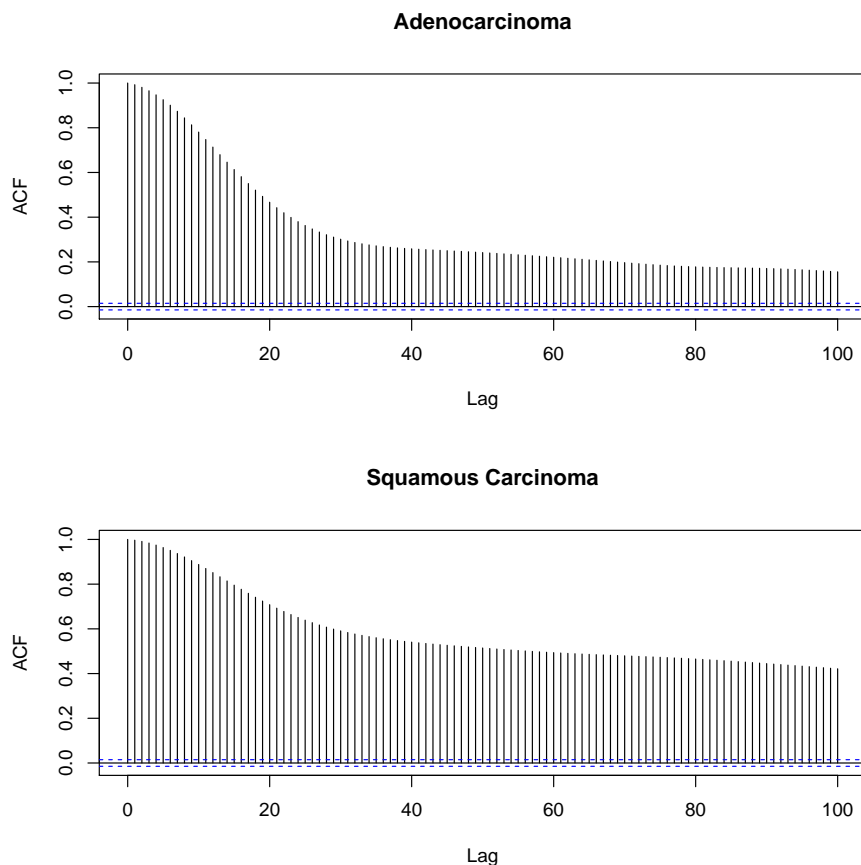


Figure 6.1: The autocorrelation function evaluated at $\tau \in [1, 100]$ for a single patient with adenocarcinoma type lung cancer (top) and squamous carcinoma type lung cancer (bottom).

Figure 6.1 shows that high correlation persists until about lag 20 for the patients with either adenocarcinoma or squamous carcinoma type lung cancer. Note however that the autocorrelation decay for patients with squamous carcinoma type lung cancer is much slower. This strong persistence in correlation suggests that the number of effective independent windows in the data set is much less than the total number of windows (17,613). This will therefore affect the location of the significance threshold when determining the significant regions.

6.2.2 Autoregressive Models

Looking at the autocorrelation functions of two patients with each type of lung cancer, it could be plausible to model the data as an autoregressive model of order q (AR(q)). Consider again X_t to be a realisation of a process X at time t . The

6. Examination of Correlation Structure in the Data

definition of an AR(q) process is

$$X_t = \sum_{i=1}^q \alpha_i X_{t-i} + \epsilon_t,$$

where α is the parameter of the model and ϵ_t is white noise. The autocorrelation function of an AR(1) process with $n = 17,613$ (the total number of windows in the genome), $E[X_t] = 0$ and $\alpha = 0.99$ is plotted in Figure 6.2.

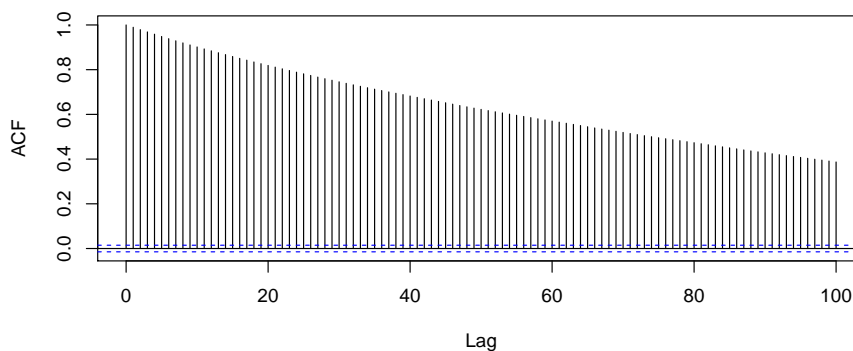


Figure 6.2: The autocorrelation function of an AR(1) process with $n = 17,613$, $E[X_t] = 0$ and $\alpha = 0.99$.

The level of decay in the autocorrelation function in Figure 6.2 is similar to that of the patient with squamous carcinoma type lung cancer. However the shape of the decay appears more linear in Figure 6.2 compared to the autocorrelation functions in Figure 6.1. This suggests that whilst an AR process may be suitable, the order q may be larger than 1. To see whether this is the case, for the two patients whose autocorrelation functions are plotted in Figure 6.1, AR(q) models are fitted. To fit the models and obtain a suitable order q , the computer package R was used. The AR models fitted to each patient with adenocarcinoma and squamous carcinoma type lung cancer were an AR(22) model and a AR(23) model respectively. Realisations of the autocorrelation functions for the fitted models with $n = 17,613$ are shown in Figure 6.3.

By comparing the autocorrelation functions of the AR(22) and AR(23) models in Figure 6.3 to the observed autocorrelation functions in Figure 6.1, it is unclear whether the fitted models are an accurate representation of the data. Up to lag 20 both models appear to describe the correlation structure well, however after lag 20 it becomes less accurate.

Whilst AR(q) processes could be a good way of accurately describing the correlation structure for small lag, the more parameters the model has, the more

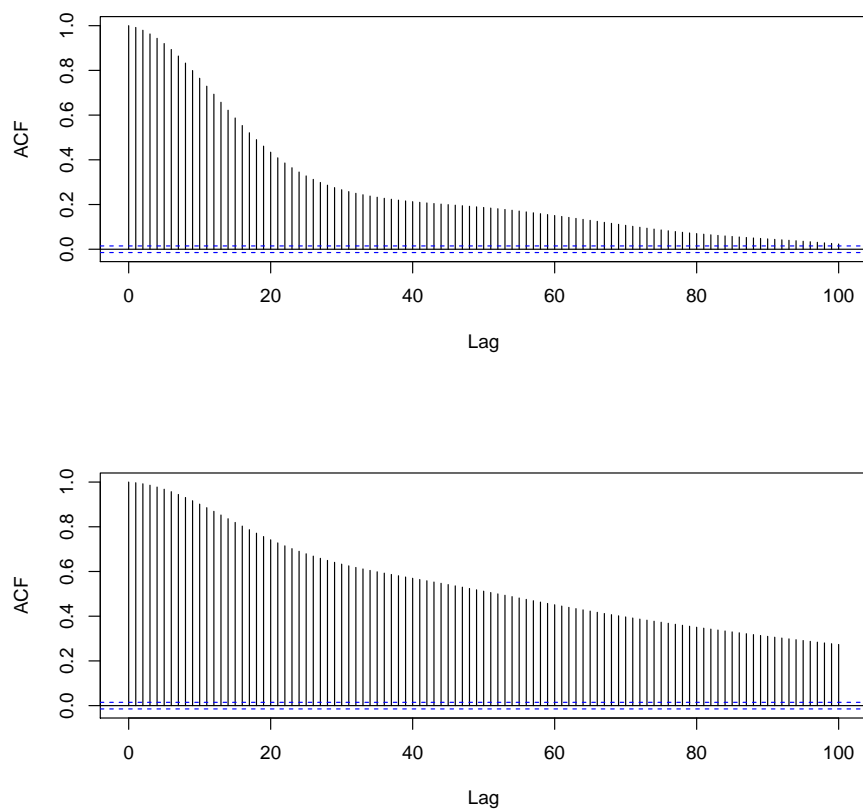


Figure 6.3: Realisation of an AR(22) (top) and AR(23) (bottom) model with parameters fitted using the CNA of two patients with adenocarcinoma and squamous carcinoma type lung cancer respectively.

6. Examination of Correlation Structure in the Data

complicated it becomes. Thus performing analyses with these models could become too complicated to deal with. Therefore, we instead consider a more simpler model for describing the correlation structures.

6.2.3 Multivariate Normal Distribution

Consider instead an AR(1) process with $\alpha = r$. This model is equivalent to a multivariate normal distribution with the following setup: Let $X \sim MVN(\vec{\mu}, \Sigma_x)$ with the mean vector and covariance matrix defined as

$$\vec{\mu} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix},$$

$$\Sigma_x = \begin{pmatrix} 1 & r & r^2 & \dots & r^p \\ r & 1 & & & \\ r^2 & & \ddots & & \\ \vdots & & & & \\ r^p & r^{(p-1)} & r^{(p-2)} & \dots & 1 \end{pmatrix}.$$

As correlation is not affected by the mean, we can therefore consider the mean to be zero across all windows. We acknowledge that the mean will not be zero in application, but for modelling the correlation structure the choice of mean is not important. Note that the process is stationary as the correlation $|r| < 1$.

Consider taking a sample of 10 consecutive windows from the genome. We can calculate the correlation matrix for these 10 windows for patients with adenocarcinoma and squamous carcinoma type lung cancer and define the matrices to be S^a and S^s respectively. Let $S_{i,j}^a$ be the (i, j) th element of the matrix S^a , then $\frac{1}{9} \sum_{i=1}^9 S_{i,i+1}^a$ is the mean correlation between the adjacent windows. We will therefore choose r to be equal to the mean sample correlation between adjacent windows and construct the matrix Σ_a with $p = 10$. We can similarly construct the matrix Σ_s by choosing $r = \frac{1}{9} \sum_{i=1}^9 S_{i,i+1}^s$, the mean correlation between adjacent windows for patients with squamous carcinoma type lung cancer. Then calculate the following matrices;

$$M_a = S^a - \Sigma_a + I_{10}$$

$$M_s = S^s - \Sigma_s + I_{10},$$

where I_{10} is the identity matrix with dimension 10. To test whether the matrix S^a is similar to Σ_a and equivalently whether S^s is similar to Σ_s , we can use a method

described by [Steiger \(1980\)](#) to statistically test the null hypothesis that M_a and M_s are equal to the identity matrix.

The first set of 10 consecutive windows after removing some windows in the lung cancer data set is windows 9 to 18 in chromosome 1, which corresponds to location 1350–2700 kbp. For these windows, the test is performed on M_a and M_s and p -values are calculated to be equal to 1 for each test. Therefore, the null hypothesis is not rejected and it is concluded that each of the two matrices are similar to the identity matrix. Thus for a small portion of the genome using a multivariate normal construction seems plausible. Testing this construction on the entire genome is currently unachievable due to the computational strain of calculating the correlation matrix for 17,613 windows.

In the next few sections, we now explore the different kinds of correlation which exists within the lung cancer data set and provide examples for each. We no longer assume that the data is stationary.

6.3 Correlation between Patients, ρ^{pa} , ρ^{ps} and ρ^p

Recall that observing correlation between patients may be an indication of potential clinical groupings among patients. Heat maps can be used to identify potential groupings by observing which patients have the strongest correlations. Here we consider correlations between patients in the lung cancer data set. We would expect to observe the correlation of patients within the same clinical groups - which in this case are patients with adenocarcinoma and patients with squamous carcinoma - to be high. We firstly consider the correlation between patients within the same group before considering the correlation between patients across both groups.

For the random variables X_{ij} and Y_{ij} , $i = 1, \dots, n$, $j = 1, \dots, p$, the correlation between patients in the same group across all windows can either correspond to the correlation between X_k and X_l (correlation between patients k and l with adenocarcinoma type lung cancer across all windows) or the correlation between Y_k and Y_l (correlation between patients k and l with squamous carcinoma type lung cancer across all windows). Thus the covariance between patients with adenocarcinoma type lung cancer is

$$\text{Cov}(X_k, X_l) = \frac{1}{p-1} \sum_{j=1}^p (x_{kj} - \bar{x}_k)(x_{lj} - \bar{x}_l).$$

Similarly for patients with squamous carcinoma type lung cancer. The correlation

6. Examination of Correlation Structure in the Data

can be calculated by

$$\rho_{kl}^{pa} = \frac{\text{Cov}(X_{k\cdot}, X_{l\cdot})}{\text{sd}(X_{k\cdot})\text{sd}(X_{l\cdot})}$$

and the correlation between patients with squamous carcinoma type lung cancer is

$$\rho_{kl}^{ps} = \frac{\text{Cov}(Y_{k\cdot}, Y_{l\cdot})}{\text{sd}(Y_{k\cdot})\text{sd}(Y_{l\cdot})}.$$

The correlation between patients in different groups across all windows corresponds to the correlation between $X_{k\cdot}$ and $Y_{l\cdot}$ (correlation between patient k with adenocarcinoma type lung cancer and patient l with squamous carcinoma type lung cancer). Thus the covariance between patients with different types of lung cancer is

$$\text{Cov}(X_{k\cdot}, Y_{l\cdot}) = \frac{1}{p-1} \sum_{j=1}^p (x_{kj} - \bar{x}_k)(y_{lj} - \bar{y}_l),$$

and the correlation can be calculated by

$$\rho_{kl}^p = \frac{\text{Cov}(X_{k\cdot}, Y_{l\cdot})}{\text{sd}(X_{k\cdot})\text{sd}(Y_{l\cdot})}.$$

6.3.1 Lung Cancer Data Set Application

Figure 6.4 shows the estimated CNA across windows for two patients with adenocarcinoma type lung cancer (left) and two patients with squamous carcinoma type lung cancer (right).

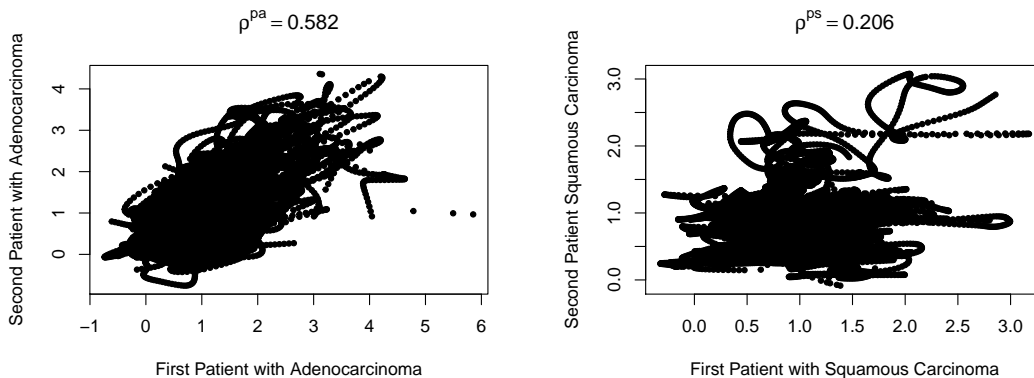


Figure 6.4: Estimated CNA across windows for two patients with adenocarcinoma type lung cancer (left) and two patients with squamous carcinoma type lung cancer (right).

Figure 6.4 (right) shows that the correlation between two patients with squamous carcinoma type lung cancer is quite low, whilst the correlation between two patients with adenocarcinoma type lung cancer is reasonably high.

The estimated CNA across windows between patient 1 with adenocarcinoma type lung cancer and patient 1 with squamous carcinoma type lung cancer is plotted in Figure 6.5. Note that the correlation between these two patients is small.

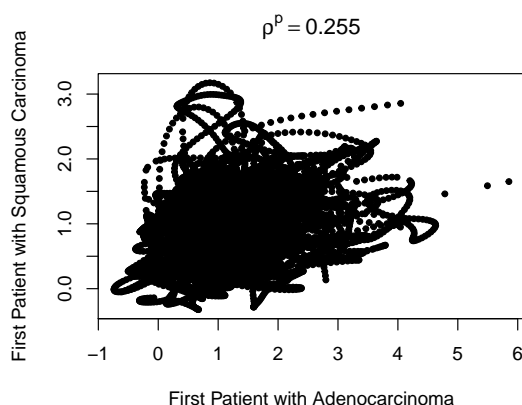


Figure 6.5: Estimated CNA across windows for patient 1 with adenocarcinoma type lung cancer and patient 1 with squamous carcinoma type lung cancer.

Here we have considered the correlation between two patients with either the same type of lung cancer or differing types. To get a clearer view of the correlation structure between patients, the correlation between all patients – with either the same type of cancer or different – is now explored.

6.3.2 Correlation Across all Patients

Between each patient either with adenocarcinoma or squamous carcinoma type lung cancer, the correlations ρ^{pa} , ρ^{ps} and ρ^p are calculated. Figure 6.6 shows a heat map of the correlations ρ^{pa} , ρ^{ps} and ρ^p across all patients.

In Figure 6.6, it can be seen that between patients with the same type of lung cancer there exists a strong positive correlation, whilst the correlation between patients with different types of lung cancers is low. By producing a heat map, the separation of patients with adenocarcinoma type lung cancer and patients with squamous carcinoma type lung cancer can be clearly seen and shows just how well plotting the heat maps of correlation can distinguish between clinical groups. There is a potential for further analysis to be done here on the correlations between patients to identify further clinical groups.

6. Examination of Correlation Structure in the Data

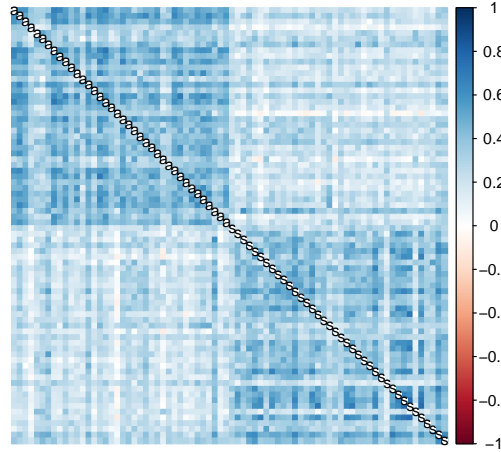


Figure 6.6: Heat map showing the value of ρ^{pa} , ρ^{ps} and ρ^p calculated for each pair of patients with either adenocarcinoma or squamous carcinoma type lung cancer. The letters on the diagonal represent the subtype of cancer, i.e. “a” represents a patient with adenocarcinoma type lung cancer and “s” represents a patient with squamous carcinoma type lung cancer.

6.4 Correlation between Windows, ρ^{wa} , ρ^{ws} and ρ^w

We have already explored the correlation between the samples or patients of our data set but probably the most important type of correlation within the data set is the correlation between variables or windows. As the regions of the genome are continuous, by taking measurements at discrete points, the correlation of the CNA between two spatially adjacent measurements will be high. In this section, we explore the correlation in the CNA between the windows of the genome and show how strong the correlation is.

For random variables X_{ij} and Y_{ij} , $i = 1, \dots, n$, $j = 1, \dots, p$, the correlation between windows for patients in the same group can either correspond to the correlation between X_k and X_l (correlation between windows k and l for patients with adenocarcinoma type lung cancer) or the correlation between Y_k and Y_l (correlation between windows k and l for patients with squamous carcinoma type lung cancer). Thus for patients with adenocarcinoma type lung cancer, the covariance between windows is

$$\text{Cov}(X_k, X_l) = \frac{1}{n-1} \sum_{i=1}^n (x_{ik} - \bar{x}_k)(x_{il} - \bar{x}_l).$$

Similarly for patients with squamous carcinoma type lung cancer. The correlation

can be calculated by

$$\rho_{kl}^{wa} = \frac{\text{Cov}(X_{.k}, X_{.l})}{\text{sd}(X_{.k})\text{sd}(X_{.l})}$$

and the correlation between windows for patients with squamous carcinoma type lung cancer is

$$\rho_{kl}^{ws} = \frac{\text{Cov}(Y_{.k}, Y_{.l})}{\text{sd}(Y_{.k})\text{sd}(Y_{.l})}.$$

6.4.1 Lung Cancer Data Set Application

For patients with adenocarcinoma type lung cancer, Figure 6.7 shows four plots investigating the correlation between windows. To see how strong the correlation between two adjacent windows is, Figure 6.7 top left shows the estimated CNA for window 6 (chromosome 1, 750kbp to 900kbp) against window 7 (chromosome 1, 901kbp, 1050kbp). Whilst it is expected that two adjacent windows in the same chromosome is highly correlated, it might not be the case for two windows in different chromosomes, thus Figure 6.7 top right shows the estimated CNA for window 1662 (chromosome 1, 249.15Mbp to 249.3Mbp) plotted against window 1663 (chromosome 2, 0 to 150kbp). Next, the estimated CNA for window 6 (chromosome 1, 750kbp to 900kbp) is plotted against window 1662 (chromosome 1, 249.3Mbp to 249.45Mbp) in Figure 6.7 bottom left to investigate the strength of the correlation between the start and end of a chromosome. Finally, given that the centromere occurs at location 123.4 Mbp in chromosome 1, the estimated CNA for window 810 (chromosome 1, 121.5Mbp to 121.65Mbp) is plotted against window 951 (chromosome 1, 142.50Mbp to 142.65Mbp) in Figure 6.7 bottom right to see the strength of the correlation just before and after the centromere.

Figure 6.7 shows that two adjacent windows in the same chromosome are very highly correlated, whilst two windows in different chromosomes are not. Also it seems that the further away the two windows are from each other the smaller the correlation, which was to be expected. Finally, the centromere doesn't effect the correlation between windows as the correlation is very high for two windows either side of the centromere. This is not surprising as the genomic distance between the two windows is small, therefore we would expect to see a higher correlation.

To see whether the same conclusions can be drawn from patients with squamous carcinoma type lung cancer, Figure 6.8 shows four plots investigating the correlation between the same pairs of windows.

Indeed, we can draw the same conclusions regarding the strength of correlations between certain pairs of windows from Figure 6.8 as we did from Figure 6.7. Whilst we have explored the correlation between specific windows, the correlation structure between two chromosomes will become clearer if we consider the

6. Examination of Correlation Structure in the Data

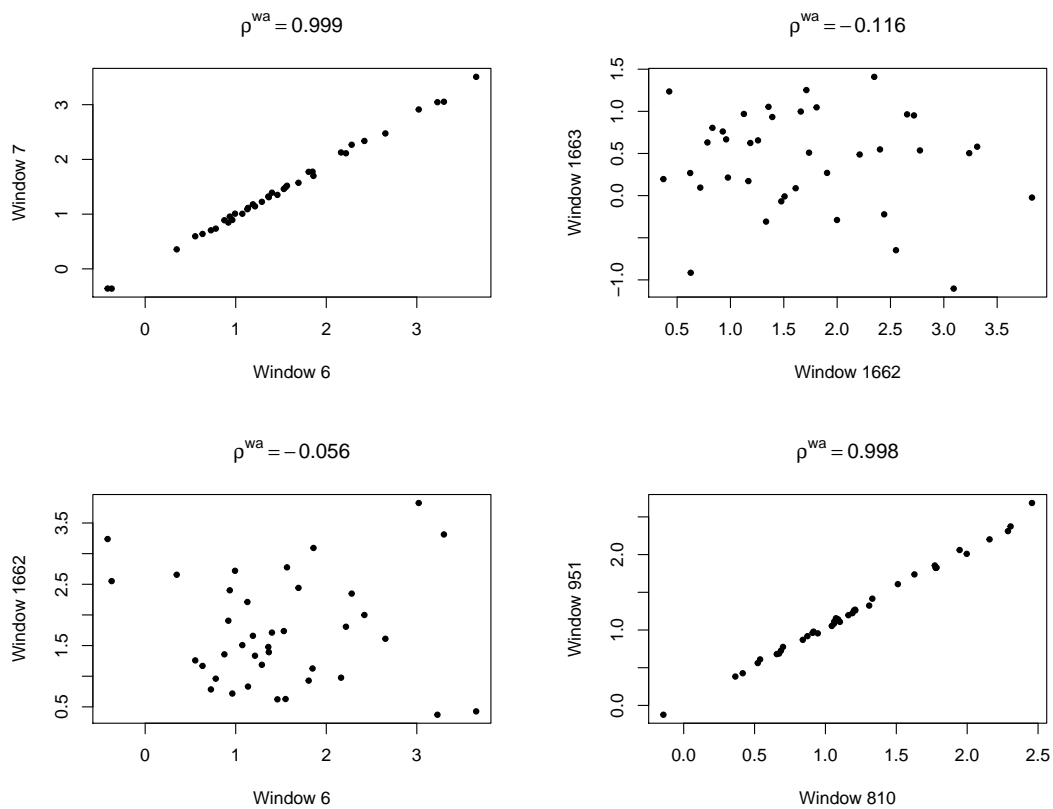


Figure 6.7: Estimated CNA plotted for patients with adenocarcinoma type lung cancer for window 6 against window 7 (top left), window 1662 against window 1663 (top right), window 6 against window 1663 (bottom left) and window 810 against window 951 (bottom right).

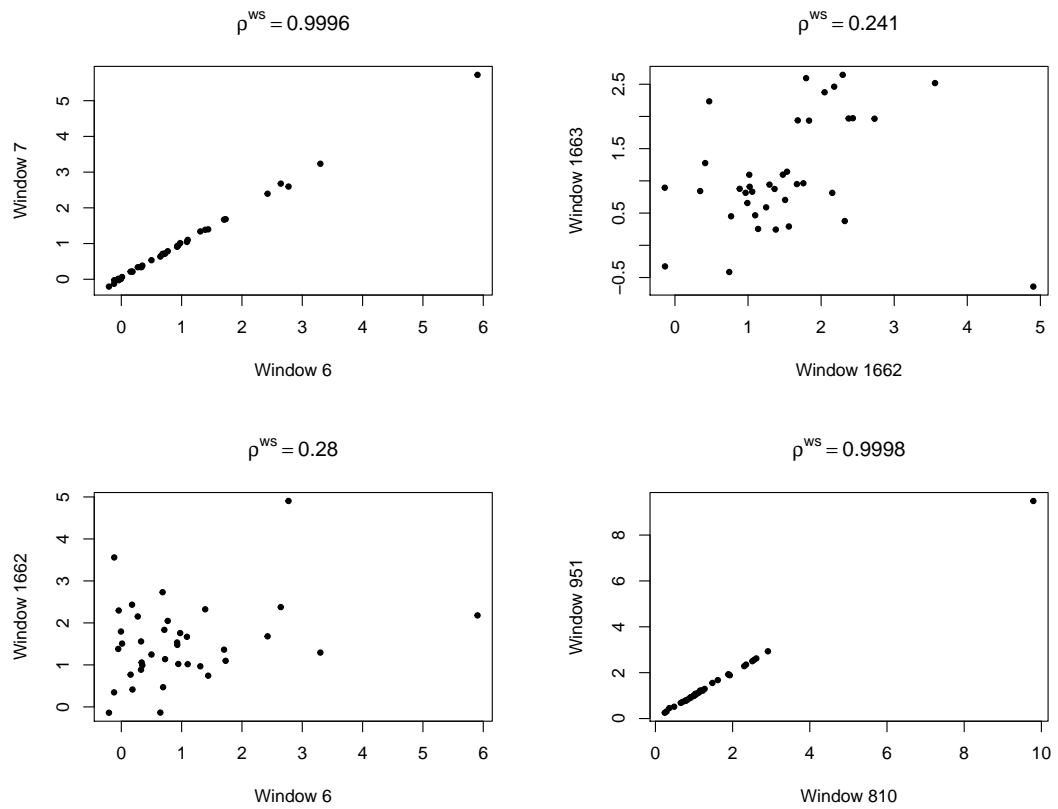


Figure 6.8: Estimated CNA for patients with squamous carcinoma type lung cancer plotted for window 6 against window 7 (top left), window 1662 against window 1663 (top right), window 6 against window 1663 (bottom left) and window 810 against window 951 (bottom right).

6. Examination of Correlation Structure in the Data

correlation between all the windows in two chromosomes.

6.4.2 Correlation Between all Windows in Two Chromosomes

As chromosomes 21 and 22 have fewer windows compared to other chromosomes, the computational time of calculating the correlation for each pair of windows is smaller than for any other pair of chromosomes. Because of this, we will use these chromosomes to display the correlation structure within chromosomes. Between each pair of windows in chromosomes 21 and 22, ρ^{wa} and ρ^{ws} has been calculated. Figure 6.9 shows a heat map of ρ^{wa} across the two chromosomes and Figure 6.10 shows the heat map of ρ^{ws} across the two chromosomes. Note that adjacent pixels may not always represent adjacent windows in the genome due to the removal of some windows.

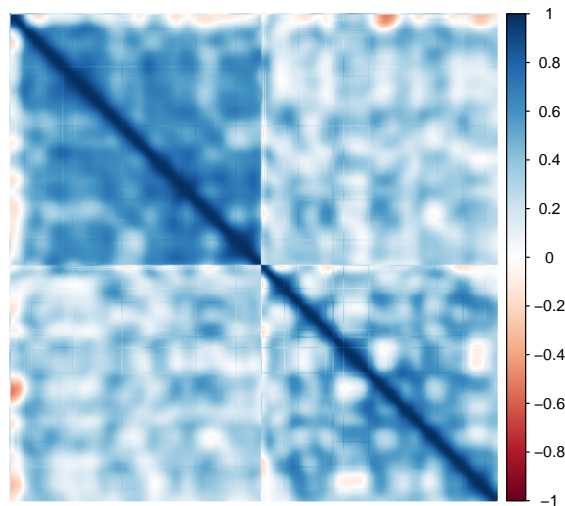


Figure 6.9: Heat map showing the value of ρ^{wa} calculated for each pair of windows in Chromosome 21 and 22.

Each heat map in Figure 6.9 and Figure 6.10 is split into four quadrants which shows the separation between chromosome 21 and 22. Strong positive correlation can also be seen around the diagonal line which is expected as the closer two windows are to each other in the genome the more highly correlated they will be. In Figure 6.10 strong positive correlation seems to remain persistent as the windows get further apart within each chromosome, whereas this feature only occurs in chromosome 21 in Figure 6.9. This suggests that high correlation persists until a larger distance between two windows for patients with squamous carcinoma compared to patients with adenocarcinoma. This is consistent with the conclusions

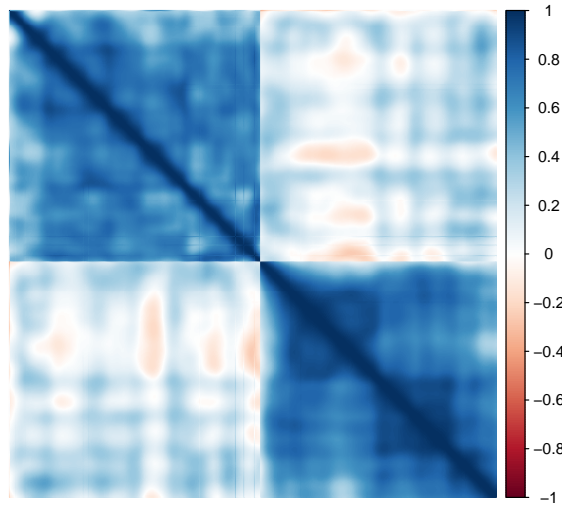


Figure 6.10: Heat map showing the value of ρ^{ws} calculated for each pair of windows in Chromosome 21 and 22.

drawn from the autocorrelation function from Section 6.2.1.

As we have shown, high correlation exists between windows. Because of this, it is likely that after performing the Cramer test on each window, the p -values will also be correlated. We now explore the correlation between the p -values.

6.5 Correlation between p -values, ρ^{pval}

Consider a set of p random variables P_1, P_2, \dots, P_p , where P_i represents all possible p -values after applying the Cramer test on variable or window i . For each P_i consider sampling n p -values, p_{i1}, \dots, p_{in} . The covariance between P_k and P_l is then

$$\text{Cov}(P_k, P_l) = \frac{1}{n-1} \sum_{j=1}^n (p_{kj} - \bar{p}_k)(p_{lj} - \bar{p}_l).$$

Correlations ρ^{pval} can then be obtained. It is natural to assume that a larger correlation between the data will cause a larger correlation between the p -values. However, the relationship between ρ^{wa} , ρ^{ws} and ρ^{pval} is complicated, for which we will now show.

6.5.1 Finding a Relationship between ρ^{wa} , ρ^{ws} and ρ^{pval}

It is impossible in our application to calculate ρ^{pval} directly. This is because for each variable only a single p -value will be obtained, and thus a sample correlation cannot be obtained. However, the sample correlations ρ^{wa} , ρ^{ws} can be calculated

6. Examination of Correlation Structure in the Data

from the data. Thus in order to be able to estimate these correlations, the relationship between ρ^{wa} , ρ^{ws} and ρ^{pval} is explored. Consider $X \sim MVN(\vec{\mu}, \Sigma_x)$ and $Y \sim MVN(\vec{\mu}, \Sigma_y)$ where $\vec{\mu} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$, $\Sigma_x = \begin{pmatrix} 1 & \rho^{wa} \\ \rho^{wa} & 1 \end{pmatrix}$ and $\Sigma_y = \begin{pmatrix} 1 & \rho^{ws} \\ \rho^{ws} & 1 \end{pmatrix}$ and $n = 500$. We can consider the variables of the multivariate normal distributions to be the windows of the genome. The Cramer test can be applied to $X_{.1}$ and $Y_{.1}$ to produce $p_{1,1}$ and can be applied to $X_{.2}$ and $Y_{.2}$ to produce $p_{1,2}$. This can be repeated 50,000 times to obtain $p_{1,1}, p_{2,1}, \dots, p_{50000,1}$ and $p_{1,2}, p_{2,2}, \dots, p_{50000,2}$ for which the correlation between the two sets of p -values, ρ^{pval} , can be calculated. This can then be done for each $\rho^{wa} \in [-1, 1]$ and $\rho^{ws} \in [-1, 1]$. Table 6.1, 6.2 and 6.3 displays some of the results of this simulation as a look-up table to obtain ρ^{pval} given ρ^{wa} and ρ^{ws} . Tables 6.1 and 6.2 display almost identical results suggesting that the correlation between p -values is not affected by the sign of ρ^{wa} and ρ^{ws} as long as both are the same sign.

$\rho^{wa} \backslash \rho^{ws}$	-1	-0.96	-0.92	-0.88	-0.84	-0.80	-0.76	-0.72
-1	1	0.898	0.821	0.758	0.704	0.654	0.610	0.571
-0.96		0.812	0.743	0.686	0.636	0.591	0.551	0.515
-0.92			0.682	0.631	0.583	0.543	0.506	0.469
-0.88				0.583	0.537	0.499	0.464	0.431
-0.84					0.497	0.461	0.428	0.399
-0.80						0.426	0.399	0.369
-0.76							0.369	0.338
-0.72								0.317

Table 6.1: The value of ρ^{pval} for each $\rho^{wa} \in [-1, -0.7]$ and $\rho^{ws} \in [-1, -0.7]$ when X and Y are simulated from a multivariate normal distribution.

$\rho^{wa} \backslash \rho^{ws}$	0.72	0.76	0.80	0.84	0.88	0.92	0.96	1
0.72	0.318	0.343	0.369	0.398	0.432	0.472	0.515	0.570
0.76		0.368	0.396	0.430	0.465	0.504	0.552	0.612
0.80			0.426	0.463	0.500	0.543	0.591	0.653
0.84				0.496	0.540	0.582	0.637	0.703
0.88					0.581	0.630	0.686	0.759
0.92						0.683	0.744	0.823
0.96							0.812	0.898
1								1

Table 6.2: The value of ρ^{pval} for each $\rho^{wa} \in [0.7, 1]$ and $\rho^{ws} \in [0.7, 1]$ when X and Y are simulated from a multivariate normal distribution.

$\rho^{wa} \backslash \rho^{ws}$	-1	-0.96	-0.92	-0.88	-0.84	-0.80	-0.76	-0.72
0.72	0.108	0.082	0.064	0.054	0.045	0.039	0.031	0.028
0.76	0.108	0.089	0.072	0.059	0.053	0.040	0.039	0.032
0.80	0.113	0.091	0.075	0.063	0.055	0.048	0.046	0.040
0.84	0.122	0.098	0.082	0.071	0.061	0.055	0.049	0.047
0.88	0.126	0.102	0.090	0.076	0.071	0.064	0.058	0.055
0.92	0.137	0.119	0.099	0.088	0.082	0.073	0.068	0.066
0.96	0.155	0.131	0.119	0.103	0.097	0.089	0.087	0.080
1	0.181	0.155	0.139	0.129	0.121	0.113	0.110	0.104

Table 6.3: The value of ρ^{pval} for each $\rho^{wa} \in [0.7, 1]$ and $\rho^{ws} \in [0.7, 1]$ when X and Y are simulated from a multivariate normal distribution.

It is clear from all tables that the relationship between ρ^{wa} , ρ^{ws} and ρ^{pval} is not straightforward. Regression techniques, including fitting generalised linear models and adding higher polynomial terms in the linear model, were used to find a formula for the relationship which can be used for predictions, however finding a satisfactory formula proved unsuccessful. The look-up tables we have created however provide a quick way of estimating the correlation between p -values given the correlations between data for each subtype.

We will now repeat these simulations to ensure robustness of the results. We will now attempt to simulate correlated random variables X and Y where X and Y follow the multi-modal mixture distribution $\frac{7}{8}N(1, 1) + \frac{1}{8}N(3, 1)$. We choose this multi-modal distribution as it is representative of the windows in the lung cancer data set which are multi-modal. We simulate this data using the following algorithm;

1. Simulate $n = 500$ observations from a multivariate normal distribution with mean vector μ and covariance matrix Σ as before.
2. Apply the univariate normal cumulative distribution function to derive probabilities for each variable.
3. Apply the inverse cumulative distribution function for the mixture distribution of X and Y to simulate draws from the distribution.

We therefore repeat the previous set up to obtain the correlation between the two sets of p -values and obtain ρ^{pval} for different values of ρ^{wa} and ρ^{ws} . Tables 6.4, 6.5 and 6.6 shows the results of this simulation.

As you can see from Tables 6.4, 6.5 and 6.6, the results obtained are not too dissimilar from that of Tables 6.1, 6.2 and 6.3 suggesting that the results are fairly robust. The mean difference between Tables 6.4 and 6.1 is 0.0025, the mean

6. Examination of Correlation Structure in the Data

$\rho^{wa} \backslash \rho^{ws}$	-1	-0.96	-0.92	-0.88	-0.84	-0.80	-0.76	-0.72
-1	0.987	0.890	0.818	0.758	0.703	0.660	0.616	0.575
-0.96		0.810	0.743	0.687	0.640	0.599	0.558	0.511
-0.92			0.683	0.634	0.586	0.544	0.515	0.471
-0.88				0.584	0.544	0.510	0.473	0.438
-0.84					0.498	0.461	0.428	0.405
-0.80						0.431	0.400	0.368
-0.76							0.373	0.347
-0.72								0.322

Table 6.4: The value of ρ^{pval} for each $\rho^{wa} \in [-1, -0.7]$ and $\rho^{ws} \in [-1, -0.7]$ when X and Y are simulated from correlated mixture distributions.

$\rho^{wa} \backslash \rho^{ws}$	0.72	0.76	0.80	0.84	0.88	0.92	0.96	1
0.72	0.322	0.350	0.376	0.409	0.440	0.480	0.525	0.575
0.76		0.374	0.406	0.435	0.472	0.507	0.555	0.614
0.80			0.437	0.471	0.511	0.555	0.597	0.657
0.84				0.503	0.547	0.592	0.645	0.707
0.88					0.585	0.639	0.687	0.745
0.92						0.693	0.750	0.829
0.96							0.815	0.903
1								1

Table 6.5: The value of ρ^{pval} for each $\rho^{wa} \in [0.7, 1]$ and $\rho^{ws} \in [0.7, 1]$ when X and Y are simulated from correlated mixture distributions.

$\rho^{wa} \backslash \rho^{ws}$	-1	-0.96	-0.92	-0.88	-0.84	-0.80	-0.76	-0.72
0.72	0.106	0.084	0.073	0.060	0.047	0.042	0.036	0.032
0.76	0.113	0.084	0.066	0.060	0.058	0.045	0.039	0.034
0.80	0.112	0.101	0.083	0.060	0.054	0.051	0.049	0.036
0.84	0.120	0.094	0.075	0.069	0.059	0.056	0.055	0.043
0.88	0.140	0.111	0.095	0.070	0.068	0.068	0.060	0.053
0.92	0.141	0.123	0.100	0.099	0.085	0.076	0.073	0.063
0.96	0.161	0.130	0.117	0.103	0.100	0.093	0.082	0.079
1	0.184	0.165	0.146	0.139	0.128	0.118	0.114	0.112

Table 6.6: The value of ρ^{pval} for each $\rho^{wa} \in [0.7, 1]$ and $\rho^{ws} \in [0.7, 1]$ when X and Y are simulated from correlated mixture distributions.

difference between 6.5 and 6.2 is 0.0059 and the mean difference between 6.6 and 6.3 is 0.0102. Thus, there does exist a small increase in the correlation between p -values when the distributions of X and Y are now multi-modal. However, because this difference is very small, it will not make much difference in further calculations

performed later in Chapter 7.

6.6 Discussion

In this chapter, we explored the various correlation structures within the lung cancer data set. The main observation was that there exists high positive correlation between adjacent windows in the same chromosome, and this high correlation persists up to at least lag 20 for patients with either type of lung cancer. Positive correlation also exists between patients with the same lung cancer type.

When modelling the correlation structure, it was found that a suitable model was a multivariate normal distribution with a 0 mean vector and a covariance matrix with element (i, j) equal to $r^{|i-j|}$. This model was then used to help find the relationship between the correlation between windows for each subtype, ρ^{wa} and ρ^{ws} , and ρ^{pval} , the correlation between the p -values. To this end, we were able to create a look-up table to determine ρ^{pval} for given ρ^{wa} and ρ^{ws} .

Chapter 7

Multiple Testing for Dependent p -values

The multiple testing problem occurs when one considers performing a hypothesis test on m independent variables simultaneously (Miller Jr, 1981). Each of the m tests are performed with a significance level α . This means that if m tests are performed when the null hypothesis is true, $\alpha \cdot m$ tests will be considered significant. In other words, there will be $\alpha \cdot m$ Type I errors. Now, the Family Wise Error Rate (FWER) is defined as the probability of making at least one Type I error in the family of simultaneous hypothesis tests. If each test is performed at the 5% significance level, the FWER is $(1 - \alpha)^m > \alpha$. To control the FWER, procedures like the Bonferroni procedure (Benjamini and Hochberg, 1995) and Sidak correction (Šidák, 1967) are used.

In our application, we apply the Cramer test to 17,613 windows of the genome simultaneously, thereby causing a multiple testing problem. The difference in our case however, is that the tests we perform are not independent (as seen in Chapter 6), thus the true number of independent tests performed, m , is unknown. As the FWER corrections rely on knowing m , the number of independent tests, we cannot apply the usual FWER corrections. An issue we face when assuming that all tests are independent is that the Bonferroni corrected significance threshold, or indeed any corrected significance threshold, may be much higher than it should be. We suspect that for the lung cancer data set, the total number of tests $p \gg m$, the effective number of independent tests performed.

In this chapter, we outline some methods to deal with multiple testing for dependent p -values. We begin by discussing a method known for approximating the multiplicity burden of the data, which will therefore help calculate the effective number of independent p -values. We apply this method to the lung cancer and artificial data set in the KC Smart vignette in order to adjust the results obtained

7. Multiple Testing for Dependent p -values

in Chapter 5 and test its effectiveness. Finally, we discuss Fisher’s combined probability test for dependent p -values. We initially describe the methodology of [Brown \(1975\)](#) but extend it for when the test statistics have an unknown distribution. This method will allow us to “group” together the p -values obtained from performing the Cramer test on the lung cancer data set to perform a test which incorporates the correlations between the p -values.

7.1 Multiplicity Burden

Multiplicity is a common problem in almost all genomic testing ([Manly et al., 2004](#)) and our case is no different. Consider performing two hypothesis tests on two highly correlated variables to obtain p -values, p_1 and p_2 . The total number of tests performed here is $p = 2$. However if p_1 is significant then most likely so will p_2 , thus to determine the significance of both tests, only a single test is required. Or in other words, the number of effective independent tests is approximately $m = 1$. Once the number of effective independent tests are identified, correction procedures like the Bonferroni correction can be applied. The issue is how to approximate m .

7.1.1 Estimating m

[Dudbridge and Gusnanto \(2008\)](#) discusses methods for obtaining an estimate for the effective number of significant tests m . They compare a method proposed by [Patterson et al. \(2006\)](#) to standard permutation methods, and found that the permutation methods provided a more accurate estimate for m . Not only this but permutation methods preserve the correlation structure within the sample ([Dudbridge and Gusnanto, 2008](#)). The permutation method to find m for the lung cancer data set is implemented as follows. Randomly assign half the samples as adenocarcinoma type lung cancer and the other half as squamous carcinoma type lung cancer. The Cramer test is applied to all windows and p -values obtained. The smallest p -value is recorded and the procedure is repeated 10,000 times to obtain 10,000 minimum p -values.

The Sidak correction, by [Šidák \(1967\)](#), calculates the multiplicity corrected significance level α_1 by solving

$$\alpha = 1 - (1 - \alpha_1)^m \tag{7.1}$$

for α_1 . Here α_1 represents the significance level for one test, α represents the significance threshold for all tests and m is the number of independent tests. The α

is calculated by considering the probability that at least one of the m independent tests are significant, which is equivalent to 1 minus the probability that none of the m independent tests are significant. The probability that none of the m independent tests are significant is

$$\Pr\{\min(p\text{-value}) > \alpha_1\} = (1 - \alpha_1)^m,$$

thus Equation (7.1) can be rewritten as

$$\Pr\{\min(p\text{-value}) \leq \alpha_1\} = 1 - (1 - \alpha_1)^m. \quad (7.2)$$

Now, Equation (7.2) is equivalent to the cumulative distribution function for a Beta distribution with parameters 1 and m . Hence, in order to estimate m , one can fit a Beta(1, m) distribution to the minimum p -values obtained after each permutation replicate. To estimate m , the formula

$$m = \frac{1 - \bar{p}_0}{\bar{p}_0}$$

is used where \bar{p}_0 represents the mean of the simulated minimum p -values.

7.1.2 Example

Consider simulating two samples of length $n = 100$ from a multivariate normal distribution with $p = 50$ variables. The mean vectors for each sample are equal to

$$\vec{\mu}_1 = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \quad \vec{\mu}_2 = \begin{pmatrix} 0.6 \\ \vdots \\ 0.6 \end{pmatrix}, \quad (7.3)$$

respectively and let the elements $S_{(i,j)}$ of the covariance matrix be defined as

$$S_{(i,j)} = \begin{cases} 1 & i = j \\ 0.999 & i > j, j - i > j \pmod{5} - 6 \\ 0 & i > j, j - i \leq j \pmod{5} - 6 \end{cases} \quad (7.4)$$

where $S_{(j,i)} = S_{(i,j)}$.

For simplicity, a two-sample t -test can be performed on each variable and the p -values obtained. As 50 hypothesis tests are being carried out simultaneously, a 5% Bonferroni corrected significance threshold is calculated as $\frac{0.05}{50} = 0.01$. Figure 7.1 shows a plot of the p -values along with the 5% Bonferroni corrected significance threshold with $m = 50$.

7. Multiple Testing for Dependent p -values

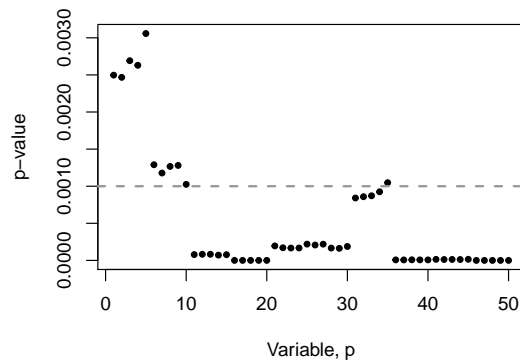


Figure 7.1: The p -values for each variable p when comparing two samples from a multivariate normal distribution with mean vectors defined in Equation (7.3) and covariance matrix defined by Equation (7.4). The horizontal grey line represents the 5% Bonferroni corrected significance threshold with $m = 50$.

Note that only 39 out of the 50 p -values are considered significant when comparing them to the 5% Bonferroni corrected significance threshold with $m = 50$. It can be seen from Figure 7.1 that the p -values are grouped together in blocks of five, this is due to the simulated correlation structure of the data. As the data are blocked in this way, instead of calculating fifty simultaneous hypothesis tests, only ten are needed; one for each block of five. This is because, if a p -value is significant, then other p -values in the same highly correlated block will also most likely be significant. Note however that this is not the case for the seventh block from the left as a single p -value in the block is not significant whereas the rest are.

We use the method described by [Dudbridge and Gusnanto \(2008\)](#) to estimate the number of effective independent tests. We permute the group labels 10,000 times, perform hypothesis tests on each variable and calculate the smallest p -value for each permutation. Figure 7.2 shows the histogram of smallest p -values along with the fitted $\text{Beta}(1, m)$ distribution.

Here $m = 10.68$, which is very close to the true number of independent tests, 10. When we recalculate the 5% Bonferroni corrected significance threshold with $m = 10$, the threshold is now $\frac{0.05}{10} = 0.005$. Figure 7.3 shows the plot of the p -values along with the 5% Bonferroni corrected significance threshold with $m = 10$.

Now Figure 7.3 shows that 50 out of the 50 hypothesis tests are significant at the 5% level. These results are what we would expect to see given the nature of the simulated data. This shows that by estimating the number of effective independent tests, we could indeed correct for the correlation within the data. We therefore attempt to apply this method to the lung cancer data set to investigate whether further regions are identified as significant.

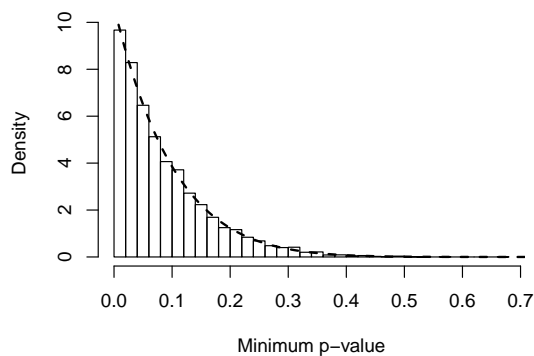


Figure 7.2: The histogram of the minimum p -values calculated for each 10,000 permutations. The dashed black line represents the fitted Beta(1,10.68) distribution.

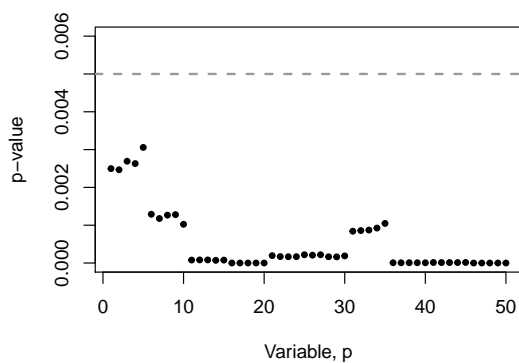


Figure 7.3: The p -values for each variable p when comparing two samples from a multivariate normal distribution with mean vectors defined in Equation (7.3) and covariance matrix defined by Equation (7.4). The horizontal grey line represents the 5% Bonferroni corrected significance threshold with $m = 10$.

7.1.3 Lung Cancer Data Set

Because each chromosome is biologically independent to all others, the effective number of independent tests can be calculated for each chromosome of the lung cancer data set. For each chromosome, the method proposed by [Dudbridge and Gusnanto \(2008\)](#) is used and a $\text{Beta}(1,m)$ distribution is fitted to the 10,000 minimum p -values. Figure 7.4 shows the histogram of the 10,000 minimum p -values when the method is performed on chromosome 1 and the ECDF plotted with the fitted $\text{Beta}(1,m)$ cumulative distribution curve.

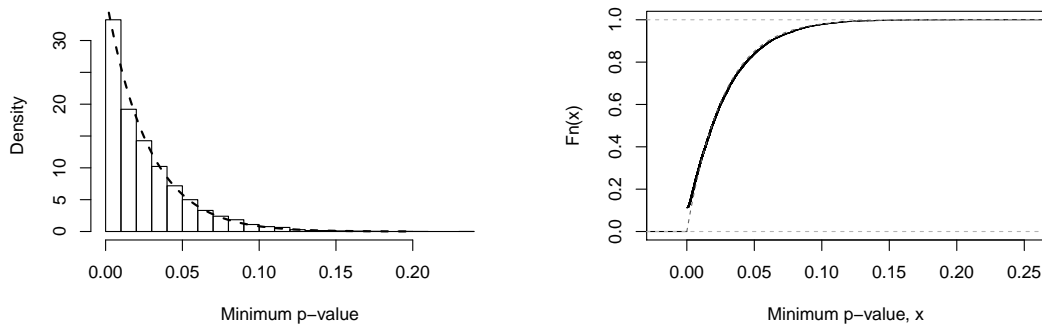


Figure 7.4: The histogram (left) and the ECDF (right) of the 10,000 minimum p -values when the method proposed by [Dudbridge and Gusnanto \(2008\)](#) is used on chromosome 1. The dashed grey line represents the fitted $\text{Beta}(1,m)$ distribution.

Table 7.1 displays the percentage of effective number of independent tests when [Dudbridge and Gusnanto \(2008\)](#)'s method is performed on each chromosome. We can perform the method on the data as a whole, but it is more informative to estimate m for each chromosome as we already know of their independence.

Table 7.1 shows that the percentage of effective number of independent tests is very small and in most cases is less than 4% for each chromosome. This would suggest a massive difference in the location of the Bonferroni corrected significance thresholds for each chromosome. The threshold for each chromosome is calculated as follows; let m_i , $i = 1, \dots, 22$, be the effective number of independent windows within each chromosome obtained from Table 7.1, thus m_1 will be the effective number of independent windows in chromosome 1 etc. Then the significance threshold for chromosome i is calculated as

$$\frac{0.05}{m_i}.$$

A modified plot with a lower significance threshold for each chromosome is shown in Figure 7.5. Figure 7.5 shows that with a lower significance threshold per chro-

Chromosome	m_i	Percentage of Effective Number of Independent Tests
1	37	2.5%
2	27	1.7%
3	34	2.6%
4	31	2.5%
5	28	2.4%
6	27	2.4%
7	19	1.9%
8	18	1.9%
9	19	2.6%
10	20	2.3%
11	24	2.8%
12	23	2.6%
13	17	2.7%
14	18	3.1%
15	19	3.7%
16	21	4.1%
17	19	3.7%
18	16	3.3%
19	17	4.5%
20	12	3.1%
21	14	5.7%
22	11	4.8%

Table 7.1: The number of effective independent tests m , and the percentage of effective number of independent tests when [Dudbridge and Gusnanto \(2008\)](#)'s method is performed on each chromosome.

mosome, more windows are now considered to be significant. In fact, there are now 3,165 significant windows compared to only 669 when $m = 17$, 613 was used in the Bonferroni correction. Recall that KC Smart identified 3,626 regions that display a significant difference of CNA between the groups of patients. After modifying the multiplicity correction to account for the correlation within the data, the Cramer test now identifies a similar amount, but KC Smart still identifies a further 500 significant regions. Again, as previously mentioned in [Section 5.7](#) this could be due to KC Smart overestimating the significance in some regions and further work is required to understand whether this is indeed the case.

7.1.4 KC Smart Data Set

The effective number of independent tests can be calculated for the artificial data set in the KC Smart vignette. The method proposed by [Dudbridge and Gusnanto \(2008\)](#) is used and a $\text{Beta}(1, m)$ distribution is fitted to the 10,000 minimum p -values. [Figure 7.6](#) shows the histogram of the 10,000 minimum p -values and the

7. Multiple Testing for Dependent p -values

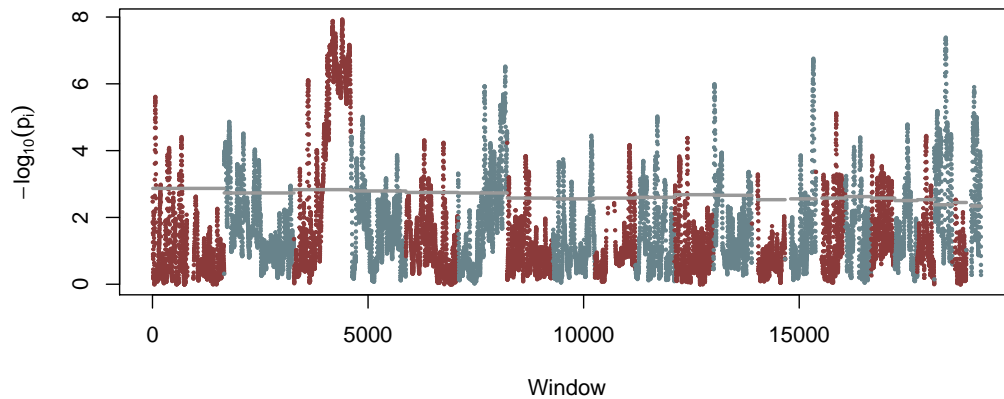


Figure 7.5: The p -values of individual genomic regions across the genome using the Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey lines represent the Bonferroni corrected significance threshold for each chromosome taking into account the percentage of effective independent tests from Table 7.1. The alternating colouring scheme indicates different chromosomes, starting with chromosome 1, 2, \dots , 22 from the left. Sex chromosomes are excluded from the analysis.

ECDF plotted with the fitted $\text{Beta}(1, m)$ cumulative distribution curve.

Here, $m = 551$, meaning that out of 3,268 regions only 551 are effectively independent. Note however that the fitted $\text{Beta}(1, m)$ distribution is not an accurate fit to the empirical distribution function of the minimum p -values suggesting that an effective number of independent tests might not exist (Dudbridge and Koeleman, 2004). This could be due to the artificial nature of the data and the simulated correlation structure may not be representative of real data. If we however assume that the effective number of independent tests is $m = 551$, then the significance threshold will be much lower than the threshold used in Figure 5.13. A modified plot with a lower significance threshold is shown in Figure 7.7.

Figure 7.7 now suggests that some probes in chromosome 4 are significant, however there are still many probes in chromosome 4 which are not identified as significant. This could suggest - assuming that the results from the KC Smart analysis are correct - that an effective number of independent tests does not sufficiently describe the correlation structure (Dudbridge and Koeleman, 2004), i.e. the correlation structure is not sufficiently captured through the permutation of the group labels. Alternatively, this could again be further evidence supporting the fact that KC Smart is over estimating the significance of chromosome 4.

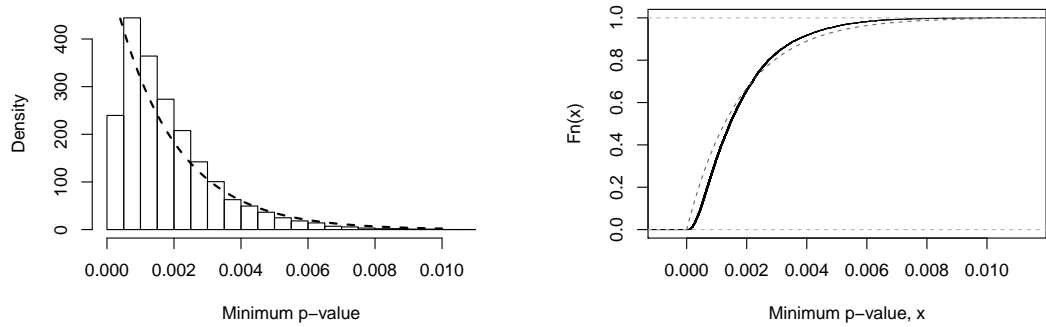


Figure 7.6: The histogram (left) and the ECDF (right) of the 10,000 minimum p -values. The dashed grey line represents the fitted Beta(1, m) distribution.

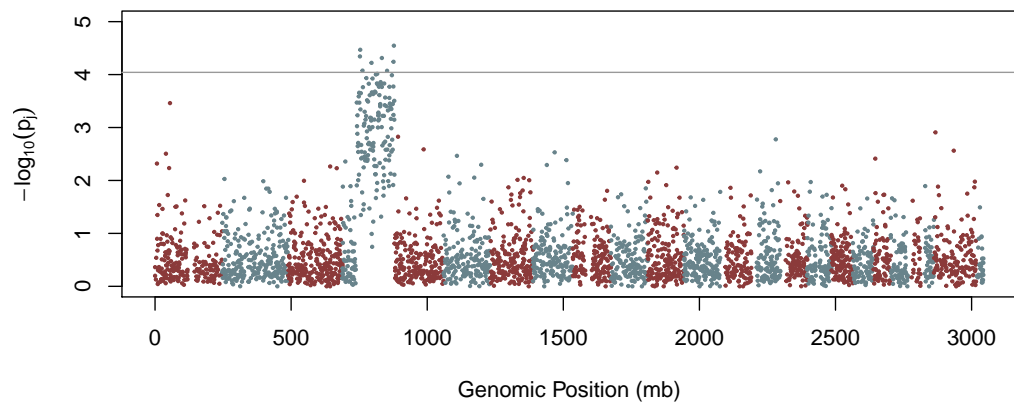


Figure 7.7: The p -values of individual genomic regions across the genome using the Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey line represents the Bonferroni corrected significance threshold with $m = 551$ effective independent tests. The alternating colouring scheme indicates different chromosomes.

7. Multiple Testing for Dependent p -values

7.1.5 Estimating Multiplicity Burden

In application, in order to estimate the multiplicity burden one could use the method described by [Dudbridge and Gusnanto \(2008\)](#). However, as the method involves arduous permutation this means that obtaining a suitable significance threshold is computationally slow. One could however, estimate the multiplicity burden if the mean correlation between adjacent windows is known. A simulation has been performed to estimate the multiplicity burden for two samples of data simulated from a multivariate normal distribution with mean vector and covariance matrix,

$$\vec{\mu} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix},$$

$$\Sigma_i = \begin{pmatrix} 1 & r_i & r_i^2 & \cdots & r_i^p \\ r_i & 1 & & & \\ r_i^2 & & \ddots & & \\ \vdots & & & & \\ r_i^p & r_i^{(p-1)} & r_i^{(p-2)} & \cdots & 1 \end{pmatrix} \quad i = 1, 2.$$

Here, r_1 will equal the mean of ρ^{wa} across all pairs of variables for the first sample and r_2 is equal to the mean of ρ^{ws} across all pairs of variables for the second sample. For high correlation to exist in the data and affect the multiplicity burden, it is expected that the correlation between adjacent variables will be higher than 0.7. Thus for $n = 100$, $p = 1000$, $r_1 \in [0.7, 1]$ and $r_2 \in [0.7, 1]$, two samples have been simulated from the multivariate normal distribution and the minimum p -value across all variables obtained. This was repeated 1000 times to obtain 1000 minimum p -values for each $r_1 \in [0.7, 1]$ and $r_2 \in [0.7, 1]$. Table 7.2 shows the estimated multiplicity burden for each $r_1 \in [0.7, 1]$ and $r_2 \in [0.7, 1]$.

$r_1 \backslash r_2$	0.7	0.75	0.8	0.85	0.9	0.95	1
0.7	481.2	450.5	447.7	425.0	424.0	361.0	89.1
0.75		451.5	425.6	407.7	389.7	347.9	94.7
0.8			404.1	386.3	356.5	318.5	88.0
0.85				359.3	342.0	295.2	86.0
0.9					293.0	255.8	69.1
0.95						203.3	44.8
1							1.1

Table 7.2: The estimated multiplicity burden for each $r_1 \in [0.7, 1]$ and $r_2 \in [0.7, 1]$ for $p = 1000$ variables and $n = 100$ observations.

If r_1 and r_2 are known, the multiplicity burden can be estimated from the table without the need for any arduous calculations. However, not all datasets will have $p = 1000$ variables, so separate tables would need to be created for all p . This, therefore, is an unrealistic method for estimating the multiplicity burden. We therefore explore another approach to multiple testing for dependent p -values, namely Fisher's combined probability test for dependent p -values.

7.2 Fisher's Combined Probability Test for Dependent p -values

Fisher et al. (1925) introduced a method to combine the p -values calculated from independent hypothesis tests which are tested under the same null hypothesis. Fisher et al. (1925) states that although individually some tests are identified as not significant, when tested in a group with other tests, it could then be identified as significant. Therefore, Fisher et al. (1925) proposes a single test to perform on groups of p -values to determine the significance of the group. This method has been called Fisher's combined probability test.

7.2.1 Fisher's Combined Probability Test for Independent p -values

Consider a set of p independent hypotheses tested to give p -values p_1, \dots, p_p . Now consider taking a subset of length k of the p -values to obtain a smaller set p_1^*, \dots, p_k^* . To test whether the subset of tests is significant as a group, firstly calculate

$$Z^2 = -2 \sum_{i=1}^k \log(p_i^*).$$

Now, as each p_i^* follows a uniform distribution in $[0,1]$, each $-\log(p_i^*)$ will follow an exponential distribution with parameter 1. Thus, $-2 \log(p_i^*)$ will be distributed as a χ^2 distribution with 2 degrees of freedom. Then taking the sum of k χ^2 distributions yield another χ^2 distribution with $2k$ degrees of freedom, thus

$$Z^2 \sim \chi_{2k}^2.$$

Assume as an example that m independent tests have been performed and the p -values for all the tests are equal to 0.04. Performing Fisher's combined probability test will yield a test statistic of $6.44 \cdot m$ which for any m will have a combined p -value of ≈ 0.04 . Thus as expected, Fisher's combined probability test

7. Multiple Testing for Dependent p -values

will conclude that all tests are significant. However, for $\alpha = 0.05$, the Bonferroni correction yields a significance threshold of $\frac{0.05}{m}$. Thus for $m > 1$, all tests are deemed insignificant. This shows that Fisher's combined probability test and the Bonferroni correction will provide opposite results for this example. However, it cannot be stated that one method is "wrong" as each method has its justifications. We will therefore compare the results after applying the Bonferroni correction using the effective number of independent tests to applying Fisher's combined probability test on dependent p -values. However, we first need to develop the theory behind the use of Fisher's combined probability test on dependent p -values.

7.2.2 Adapting Fisher's Method for Dependent p -values

When the test statistics can be modelled using a multivariate normal distribution with known covariance matrix Σ , [Brown \(1975\)](#) suggests an adaptation to Fisher's combined probability test. [Kost and McDermott \(2002\)](#) extends this adaptation for when the covariance matrix Σ is known up to a scalar quantity. It was found, using the Cramer test, that the null distribution was not a normal distribution, thus we adapt the method for the case where the distribution of the test statistics is unknown.

The adaptation used by both [Brown \(1975\)](#) and [Kost and McDermott \(2002\)](#) states that when all k null hypotheses are true, Z^2 has a scaled χ^2 distribution, i.e.

$$Z^2 \sim c\chi_f^2.$$

To estimate the values of c and f the method of moments can be used by equating

$$\begin{aligned} f &= \frac{2\mathbb{E}[Z^2]^2}{\text{Var}[Z^2]} \\ c &= \frac{\text{Var}[Z^2]}{2\mathbb{E}[Z^2]}. \end{aligned}$$

Again as each p -value, p_i^* , is distributed uniformly in the region $[0,1]$, $-\log(p_i^*)$ will be exponentially distributed with 1 degree of freedom. Thus $\mathbb{E}[-\log(p_i^*)] = 1$ and $\text{Var}[-\log(p_i^*)] = 1$. Hence,

$$\begin{aligned} \mathbb{E}[Z^2] &= \mathbb{E}\left[-2 \sum_{i=1}^k \log(p_i^*)\right] \\ &= 2 \sum_{i=1}^k \mathbb{E}[-\log(p_i^*)] \\ &= 2k, \end{aligned}$$

and

$$\begin{aligned}
 \text{Var}[Z^2] &= \text{Var} \left[-2 \sum_{i=1}^k \log(p_i^*) \right] \\
 &= 4 \sum_{i=1}^k \text{Var}[-\log(p_i^*)] + 2 \sum_{i < j} \text{Cov}(-2 \log(p_i^*), -2 \log(p_j^*)) \\
 &= 4k + 2 \sum_{i < j} \text{Cov}(-2 \log(p_i^*), -2 \log(p_j^*)). \tag{7.5}
 \end{aligned}$$

Thus in order to estimate the variance of Z^2 , the covariance between two p -values is required. To calculate the covariance, we obtain

$$\begin{aligned}
 \text{Cov}(-2 \log(p_i^*), -2 \log(p_j^*)) &= \text{E}[4 \log(p_i^*) \log(p_j^*)] - \text{E}[-2 \log(p_i^*)] \text{E}[-2 \log(p_j^*)] \\
 &= \text{E}[4 \log(p_i^*) \log(p_j^*)] - 4 \text{E}[\log(p_i^*)] \text{E}[\log(p_j^*)] \\
 &= 4(\text{E}[\log(p_i^*) \log(p_j^*)] - 1)
 \end{aligned}$$

and

$$\text{E}[\log(p_i^*) \log(p_j^*)] = \int_0^1 \int_0^1 \log p_i^* \log p_j^* f_{P_i^*, P_j^*}(p_i^*, p_j^*) dp_i^* dp_j^*. \tag{7.6}$$

Hence, the joint expectation of $\log(p_i^*)$ and $\log(p_j^*)$ and thus the covariance between p -values can be estimated by calculating the integral in Equation (7.6).

To estimate the joint expectation and therefore the covariance, [Brown \(1975\)](#) lets $W_i = -2 \log(p_i^*) = -2 \log \Phi(Z_i)$, where Φ represents the cumulative distribution function of the standard normal distribution and $Z_i = \Phi^{-1}(p_i^*)$. Therefore [Brown \(1975\)](#) obtains the joint density of W_i and W_j using transformation of variables. Similarly, [Kost and McDermott \(2002\)](#) assumes that the $W_i = -2 \log(p_i^*) = -2 \log \mathcal{J}(T_i)$ where \mathcal{J} represents the cumulative distribution function of the t distribution with v degrees of freedom. [Brown \(1975\)](#) then evaluates the joint expectation using Gaussian quadrature ([Krylov, 1962](#)) and [Kost and McDermott \(2002\)](#) uses numerical integration techniques. When the distribution of the test statistics is unknown, it is impossible to determine the joint distribution by replacing the p -values by a function of their corresponding test statistics. Thus instead we evaluate the joint expectation by considering the joint distribution of p_i^* and p_j^* .

7.2.3 Calculating the Joint Expectation

When all null hypotheses are true, the p_i^* , $i = 1, \dots, k$, are uniformly distributed. Thus $f_{P_i^*, P_j^*}(p_i^*, p_j^*)$ is a joint distribution such that the marginal distributions are correlated uniform distributions. [Ferguson \(1995\)](#) describes a method for defining a joint distribution with correlated uniform distributed marginals, and is imple-

7. Multiple Testing for Dependent p -values

mented by Demirtas (2014). For the purpose of simplicity we define such a joint distribution to be a bivariate uniform distribution, but note that the definition is defined here loosely and simply refers to any joint distribution which has uniformly distributed marginal distributions which are correlated.

Definition Suppose that G is absolutely continuous with density $g(u)$ for $0 \leq u \leq 1$. Then

$$f_{X,Y}(x, y) = \frac{1}{2}[g(|x - y|) + g(1 - |1 - x - y|)], \quad (7.7)$$

is the probability density function of a bivariate uniform distribution where marginally X and Y follow a uniform distribution, $0 \leq x \leq 1$, $0 \leq y \leq 1$. \square

Recall that ρ^{pv} defines the correlation between p -values, which in our case will also be the correlation between the bivariate uniform random variables. Ferguson (1995) proved that

$$\rho^{pv} = 1 - 6E[U^2] + 4E[U^3]$$

and by choosing $U \sim \text{Beta}(a, 1)$ (Demirtas, 2014), i.e. $g(u) = au^{(a-1)}$, it can be shown that

$$\rho^{pv} = \frac{(1-a)(6+a)}{(2+a)(3+a)}. \quad (7.8)$$

Rearranging Equation (7.8) gives

$$a = -\frac{5}{2} + \frac{1}{2}\sqrt{\frac{\rho^{pv} + 49}{\rho^{pv} + 1}}. \quad (7.9)$$

Figure 7.8 shows a one to one relationship between ρ^{pv} and a for $\rho^{pv} \in [-1, 1]$. Thus one can define a bivariate uniform distribution with any correlation, ρ^{pv} , between -1 and 1 by using Equation (7.9) to determine a suitable choice of a . The correlation ρ^{pv} can be estimated by calculating ρ^{wa} and ρ^{ws} for the data and using Tables 6.1, 6.2 and 6.3.

The definition for the bivariate uniform distribution defined in Equation (7.7) along with the choice of $g(u) = au^{(a-1)}$ can be substituted into Equation (7.6) to give

$$E[\log(p_i^*) \log(p_j^*)] = \frac{a}{2} \int_0^1 \int_0^1 \log p_i^* \log p_j^* [(|p_i^* - p_j^*|)^{(a-1)} + (1 - |1 - p_i^* - p_j^*|)^{(a-1)}] dp_i^* dp_j^*. \quad (7.10)$$

To calculate the integral in Equation (7.10) the range of integration is split into four areas:

1. **The case: $p_i^* \leq p_j^*$ and $p_i^* + p_j^* \leq 1$**

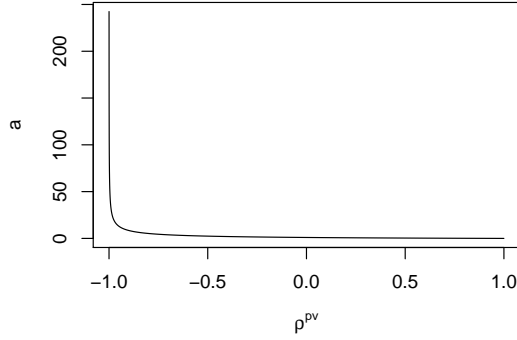


Figure 7.8: The relationship between ρ^{pv} and a for $\rho^{pv} \in [-1, 1]$.

Here, the bivariate uniform distribution defined in Equation (7.7) with $g(u) = au^{(a-1)}$ is equal to

$$f_{P_i^*, P_j^*}(p_i^*, p_j^*) = \frac{a}{2}[(p_j^* - p_i^*)^{(a-1)} + (p_i^* + p_j^*)^{(a-1)}]$$

and the integral in Equation (7.6) becomes

$$\frac{a}{2} \int_0^1 \int_0^1 \log p_i^* \log p_j^* [(p_j^* - p_i^*)^{(a-1)} + (p_i^* + p_j^*)^{(a-1)}] dp_i^* dp_j^*. \quad (7.11)$$

Now let $s = p_i^* - p_j^*$ and $t = p_i^* + p_j^*$. Given the range restrictions for p_i^* and p_j^* , this implies that $0 \leq t \leq 1$ and $-t \leq s \leq 0$. Thus after performing a change of variables, (7.11) becomes

$$\frac{a}{4} \int_0^1 \int_{-t}^0 \log \frac{1}{2}(s+t) \log \frac{1}{2}(t-s) [(-s)^{(a-1)} + t^{(a-1)}] ds dt. \quad (7.12)$$

2. **The case: $p_i^* \geq p_j^*$ and $p_i^* + p_j^* \leq 1$**

Here, the bivariate uniform distribution defined in Equation (7.7) with $g(u) = au^{(a-1)}$ is equal to

$$f_{P_i^*, P_j^*}(p_i^*, p_j^*) = \frac{a}{2}[(p_i^* - p_j^*)^{(a-1)} + (p_i^* + p_j^*)^{(a-1)}]$$

and the integral in Equation (7.6) becomes

$$\frac{a}{2} \int_0^1 \int_0^1 \log p_i^* \log p_j^* [(p_i^* - p_j^*)^{(a-1)} + (p_i^* + p_j^*)^{(a-1)}] dp_i^* dp_j^*. \quad (7.13)$$

Now let $s = p_i^* - p_j^*$ and $t = p_i^* + p_j^*$. Given the range restrictions for p_i^*

7. Multiple Testing for Dependent p -values

and p_j^* , this implies that $0 \leq t \leq 1$ and $0 \leq s \leq t$. Thus after performing a change of variables, Equation (7.13) becomes

$$\frac{a}{4} \int_0^1 \int_0^{-t} \log \frac{1}{2}(s+t) \log \frac{1}{2}(t-s)[s^{(a-1)} + t^{(a-1)}] ds dt. \quad (7.14)$$

3. The case: $\mathbf{p}_i^* \leq \mathbf{p}_j^*$ and $\mathbf{p}_i^* + \mathbf{p}_j^* \geq 1$

Here, the bivariate uniform distribution defined in Equation (7.7) with $g(u) = au^{(a-1)}$ is equal to

$$f_{P_i^*, P_j^*}(p_i^*, p_j^*) = \frac{a}{2} [(p_j^* - p_i^*)^{(a-1)} + (2 - p_i^* - p_j^*)^{(a-1)}]$$

and the integral in Equation (7.6) becomes

$$\frac{a}{2} \int_0^1 \int_0^1 \log p_i^* \log p_j^* [(p_j^* - p_i^*)^{(a-1)} + (2 - p_i^* - p_j^*)^{(a-1)}] dp_i^* dp_j^*. \quad (7.15)$$

Now let $s = p_i^* - p_j^*$ and $t = p_i^* + p_j^*$. Given the range restrictions for p_i^* and p_j^* , this implies that $1 \leq t \leq 2$ and $t - 2 \leq s \leq 0$. Thus after performing a change of variables, Equation (7.15) becomes

$$\frac{a}{4} \int_1^2 \int_{t-2}^0 \log \frac{1}{2}(s+t) \log \frac{1}{2}(t-s)[(-s)^{(a-1)} + (2-t)^{(a-1)}] ds dt. \quad (7.16)$$

4. The case: $\mathbf{p}_i^* \geq \mathbf{p}_j^*$ and $\mathbf{p}_i^* + \mathbf{p}_j^* \geq 1$

Here, the bivariate uniform distribution defined in Equation (7.7) with $g(u) = au^{(a-1)}$ is equal to

$$f_{P_i^*, P_j^*}(p_i^*, p_j^*) = \frac{a}{2} [(p_i^* - p_j^*)^{(a-1)} + (2 - p_i^* - p_j^*)^{(a-1)}]$$

and the integral in Equation (7.6) becomes

$$\frac{a}{2} \int_0^1 \int_0^1 \log p_i^* \log p_j^* [(p_i^* - p_j^*)^{(a-1)} + (2 - p_i^* - p_j^*)^{(a-1)}] dp_i^* dp_j^*. \quad (7.17)$$

Now let $s = p_i^* - p_j^*$ and $t = p_i^* + p_j^*$. Given the range restrictions for p_i^* and p_j^* , this implies that $1 \leq t \leq 2$ and $0 \leq s \leq 2 - t$. Thus after performing a change of variables, Equation (7.17) becomes

$$\frac{a}{4} \int_1^2 \int_0^{2-t} \log \frac{1}{2}(s+t) \log \frac{1}{2}(t-s)[s^{(a-1)} + (2-t)^{(a-1)}] ds dt. \quad (7.18)$$

Due to symmetry Equations (7.11) and (7.13) are equivalent and so are Equations

7.2 Fisher's Combined Probability Test for Dependent p -values

(7.15) and (7.17). Equations (7.11), (7.13), (7.15) and (7.17) can be combined to obtain the joint expectation $E[\log(p_i^*) \log(p_j^*)]$. Thus Equation (7.6) is equal to

$$E[\log(p_i^*) \log(p_j^*)] = a \left[\int_0^1 \int_{-t}^0 \log \frac{1}{2}(s+t) \log \frac{1}{2}(t-s)[(-s)^{(a-1)} + t^{(a-1)}] ds dt + \int_1^2 \int_{t-2}^0 \log \frac{1}{2}(s+t) \log \frac{1}{2}(t-s)[(-s)^{(a-1)} + (2-t)^{(a-1)}] ds dt \right],$$

and so the joint expectation can be substituted into

$$\text{Cov}(-2 \log(p_i^*), -2 \log(p_j^*)) = 4(E[\log(p_i^*) \log(p_j^*)] - 1)$$

to obtain the covariance between $-2 \log(p_i^*)$ and $-2 \log(p_j^*)$, thus enabling an estimate to be obtained for the variance of Z^2 .

Now that an estimate for both the mean and variance can be found, the method of moments can be used to find estimates for the parameters c and f of the scaled χ^2 distribution and hence a p -value can be obtained.

7.2.4 Using Fisher's Combined Probability Test on Dependent p -values

Consider a set of p p -values p_1, \dots, p_p , and let the correlation between two p -values p_i and p_j be high when $|i - j| = 1$ and decay when $|i - j|$ increases. Estimate the value of $k = |i - j|$ in which the correlation between p -values is no longer large. There now exists two ways of proceeding, the choice of which is down to the user and each comes with advantages and disadvantages.

The first method is to consider a sliding block $B_i, i = 1, \dots, p - k + 1$ of k consecutive p -values. Apply Fisher's combined probability test to each block of p -values where $B_i = \{p_i, \dots, p_{k+i-1}\}$. This method has advantages in that every possible block of consecutive k p -values are tested using Fisher's combined probability test. However a big disadvantage is that subsequent adjacent p -values obtained after applying Fisher's combined probability test will remain highly correlated.

An alternative approach is to consider non-overlapping blocks $B_i, i = 1, \dots, \frac{p}{k}$ of k consecutive p -values. An obvious disadvantage to this method is the requirement of $\frac{p}{k} \in \mathbb{N}$. Apply Fisher's combined probability test to each block of p -values where $B_i = \{p_{k+i-1}, \dots, p_{ik}\}$. The advantages of this method is that the correlation between Fisher's combined probability test p -values should no longer be high, however many sets of k consecutive p -values are not tested, and thus important information may be lost.

7. Multiple Testing for Dependent p -values

7.2.5 Example

Consider the same two samples from a multivariate normal distribution with mean vectors defined in Equation (7.3) and covariance matrix defined in Equation (7.4) from Section 7.1.2. The value of $|i - j|$ in which p_i and p_j are no longer highly correlated is 5, thus we will consider blocks of size 5.

Firstly consider a sliding block of 5 p -values. We can perform Fisher's combined probability test on each block $B_i, i = 1, \dots, 46$. As 46 simultaneous Fisher's combined probability tests are being performed, the 5% Bonferroni corrected significance threshold is calculated for $m = 46$. Note that as all blocks remain highly correlated, the number of effective independent tests will be smaller than 46. Figure 7.9 shows the p -values obtained after performing Fisher's combined probability test on each block $B_i, i = 1, \dots, 46$.

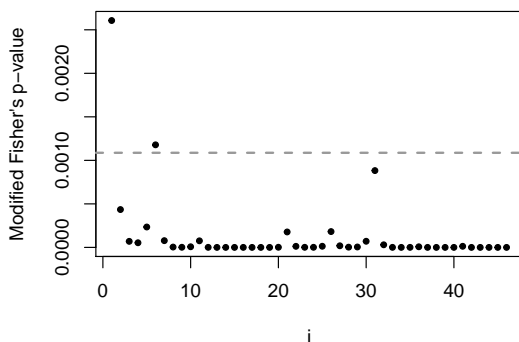


Figure 7.9: Using a sliding block of 5 p -values, the Fisher's combined probability test p -values are plotted when the test is performed on each block $B_i, i = 1, \dots, 46$, of highly correlated p -values. The horizontal grey dashed line represents the 5% Bonferroni corrected significance threshold.

Figure 7.9 shows that only blocks B_1 and B_6 are considered as not significant, however as all the p -values in B_1 and B_6 appear in other significant blocks, we could conclude that all p -values are significant. This conclusion is consistent with the results from Section 7.1.2. We could alternatively conclude that all p -values in blocks B_1 and B_6 are not significant, however as there still exists high correlation between the modified Fisher's method p -values, the number of effective independent tests is still about 10, and thus the Bonferroni corrected significance threshold is too high. When the Bonferroni corrected significance threshold is calculated for $m = 10$ all blocks are significant. In application however, it is impossible to know the number of effective independent tests without using estimation procedures like from Section 7.1.1. Therefore this method should be used with caution.

7.2 Fisher's Combined Probability Test for Dependent p -values

Next consider non-overlapping blocks of 5 p -values. We can perform Fisher's combined probability test on each block B_i , $i = 1, \dots, 10$, of highly correlated p -values. As ten simultaneous Fisher's combined probability tests are being performed, the 5% Bonferroni corrected significance threshold is calculated for $m = 10$. Figure 7.10 shows the p -values obtained after performing Fisher's combined probability test on each block B_i , $i = 1, \dots, 10$, of highly correlated p -values.

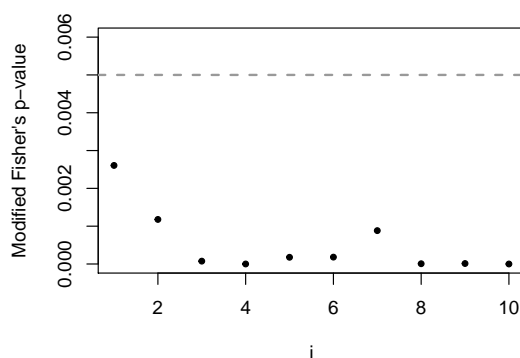


Figure 7.10: Using non-overlapping blocks of 5 p -values, the Fisher's combined probability test p -values are plotted using when the test is performed on each block B_i , $i = 1, \dots, 10$, of highly correlated p -values. The horizontal grey dashed line represents the 5% Bonferroni corrected significance threshold.

Here, Figure 7.10 shows that all blocks are significant. We therefore conclude that all p -values are identified as significant. Using non-overlapping blocks of 5 p -values works well here, however in application the data might not be correlated in perfect blocks, thus using a sliding block might be more suitable as all groups of adjacent p -values are considered. Note that the conclusions drawn after applying Fisher's combined probability test using a sliding block system and non-overlapping block system are the same. As well as this, both methods manage to identify that, for each variable, the distributions of each sample are not the same. We have therefore shown that Fisher's combined probability test may be a suitable post-hoc test to perform to identify further significant regions by accounting for correlation within the data.

7.2.6 Lung Cancer Data Set

We found in Section 6.2.1 that high correlation is present until at least lag 20 for patients with each subtype of lung cancer. Because of this, we consider a sliding block of 20 p -values within each chromosome and perform Fisher's combined probability test on each block. By using a sliding window on each chromosome,

7. Multiple Testing for Dependent p -values

we ensure that a block of p -values will not overlap more than one chromosome. Figure 7.11 shows the $-\log_{10}(p_i)$ plotted against the block of 20 p -values, where p_i represents the p -values obtained after performing Fisher's combined probability test on each block.

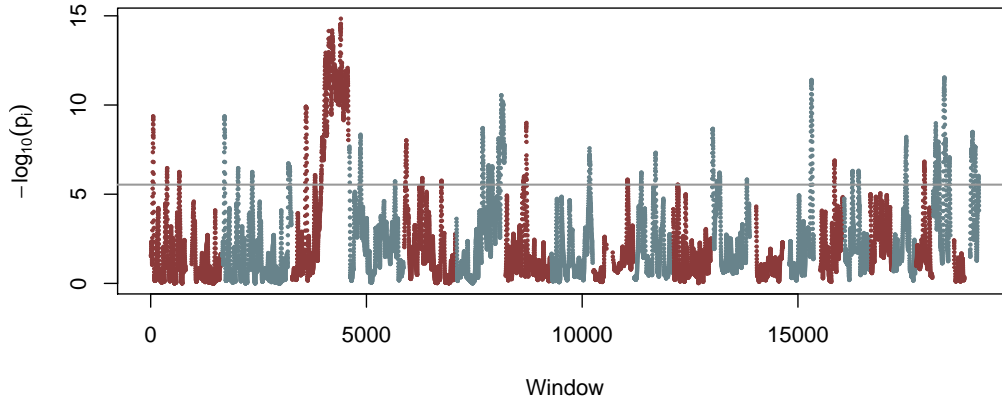


Figure 7.11: The p -values of each sliding block of 20 Cramer test p -values using the adjusted Fisher's method. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The horizontal grey line represents the Bonferroni corrected significance threshold with $m = 17173$. The alternating colouring scheme indicates different chromosomes.

After performing Fisher's combined probability test on sliding blocks of 20 p -values for each chromosome, the number of significant blocks is 1,725. This corresponds to a total of 2,580 significant windows when we count the number of unique significant windows across all the blocks. Table 7.3 shows the number of windows which were identified as significant by (1) both Fisher's combined probability test and using the Bonferroni correction (with the estimated number of effective independent tests), (2) Fisher's combined probability test only, (3) using the Bonferroni correction (with the estimated number of effective independent tests) only and (4) neither Fisher's combined probability test or using the Bonferroni correction (with the estimated number of effective independent tests).

		Fisher's Combined Probability Test	
		Significant	Not Significant
Bonferroni Correction	Significant	2,236	929
	Not Significant	344	15,377

Table 7.3: The number of significant windows (after Bonferroni correction for Fisher's combined probability test) when using Fisher's combined probability test and using the Bonferroni correction (with the estimated number of effective independent tests).

7.2 Fisher’s Combined Probability Test for Dependent p -values

Table 7.3 shows that both methods were able to identify a further 2,236 significant windows more than when no adjustment for correlation is taken into account. However, there were just over 900 windows identified as significant using the Bonferroni correction (with the estimated number of effective independent tests) for which Fisher’s combined probability test identified as insignificant. One potential cause is that high correlation still exists between the Fisher’s combined probability tests, thus the significance threshold may be too high. One solution to this would be to perform Fisher’s combined probability test and then calculate the effective number of independent tests to adjust the Bonferroni correction. This would however have the disadvantage we tried to avoid which is that of being computationally slow.

To show which windows are identified as significant by each method, Figure 7.12 shows a plot of the p -values using the univariate Cramer test and the GPD to calculate the p -value for each window of the genome. The p -values which are identified as significant by each method are plotted on the graph using a different symbol. It can be seen from Figure 7.12 that if the significance level was corrected for the estimated number of effective independent tests, Fisher’s combined probability test may identify more windows in common with just using the Bonferroni correction (with the estimated number of effective independent tests) on the Cramer test p -values.

7.2.7 KC Smart Data Set

For the artificial KC Smart data set we will perform Fishers combined probability test on each chromosome to test the significance of each one. Given that the data was created so there exists a difference in CNA between the two groups in chromosome 4, we would expect to see this and only this chromosome as significant. For this example, we use non-overlapping blocks B_i , $i = 1, \dots, 22$. Each block B_i will contain the p -values of each region in chromosome i . Table 7.4 shows the Fishers combined probability test p -values after performing the test on each block B_i .

Figure 7.13 shows a plot of the Fishers probability test p -values for each chromosome. The horizontal grey line represents the Bonferroni corrected significance threshold. In this case, any p -values below the threshold is considered significant.

It can be seen from Figure 7.13 that the only chromosome which is significant is chromosome 4. This is consistent with the results obtained after applying KC Smart and therefore suggests that chromosome 4 is indeed significant after taking the correlations into account.

After calculating the effective number of independent tests for this dataset

7. Multiple Testing for Dependent p -values

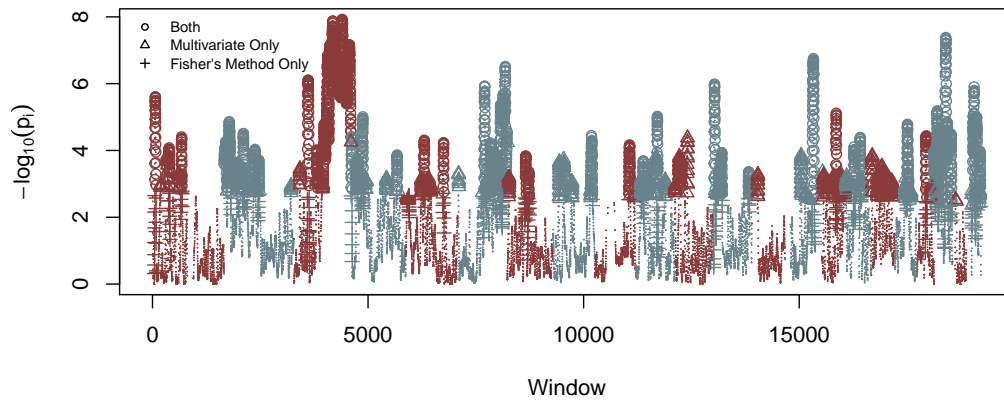


Figure 7.12: The p -values using the univariate Cramer test and the GPD to calculate the p -value for each window of the genome. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. Larger circles are plotted at windows which are identified as significant by both Fisher's combined probability test and using the Bonferroni correction (with the estimated number of effective independent tests), triangles are plotted at windows which are identified as significant by using the Bonferroni correction (with the estimated number of effective independent tests) only and plus signs are plotted at windows which are identified as significant by Fisher's combined probability test only. The alternating colouring scheme indicates different chromosomes.

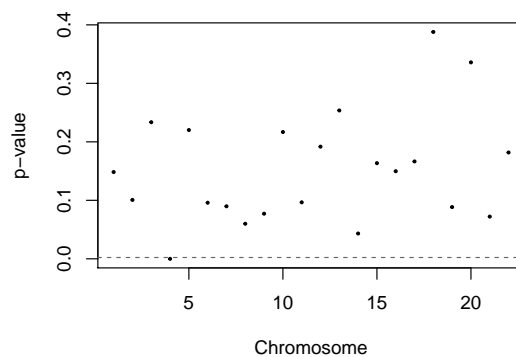


Figure 7.13: The p -values of applying Fisher's combined probability test to each block B_i of p -values for each chromosome.

7.2 Fisher's Combined Probability Test for Dependent p -values

Chromosome	p -value
1	0.148
2	0.101
3	0.234
4	0
5	0.220
6	0.096
7	0.090
8	0.060
9	0.077
10	0.217
11	0.097
12	0.192
13	0.254
14	0.043
15	0.164
16	0.150
17	0.166
18	0.388
19	0.088
20	0.336
21	0.072
22	0.182

Table 7.4: The Fishers combined probability test p -values when applied to blocks of p -values B_i from each chromosome i .

7. Multiple Testing for Dependent p -values

in Section 7.1.4 and correcting the significance threshold using the Bonferroni correction, the probes 717–924 were not considered significant. This therefore provides an example for which Fisher’s combined probability test is able to identify regions of significance which other methods do not.

7.2.8 Using the Multivariate Version of the Cramer Test

Baringhaus and Franz (2004) not only introduces the univariate Cramer test but also extends it to a multivariate case. The multivariate test statistic is defined by

$$T_{m,n} = \gamma_d \frac{nm}{n+m} \int_{S^{d-1}} \int_{-\infty}^{\infty} [F_n^a(t) - G_m^a(t)]^2 dt d\mu(a),$$

where

$$\gamma_d = \frac{\sqrt{\pi}(d-1)\Gamma\left(\frac{d-1}{2}\right)}{2\Gamma\left(\frac{d}{2}\right)},$$

μ is the uniform distribution on $S^{d-1} = \{x \in \mathbb{R}^d : \|x\| = 1\}$, the surface of the unit sphere in \mathbb{R}^d and d is the dimension of the data.

This version of the test is therefore comparable with using the univariate version of the test followed by Fisher’s combined probability test. In the multivariate Cramer test, p -values are calculated by bootstrapping the limiting distribution. Thus to directly compare the two methods, a sliding block of 20 windows is tested using the multivariate Cramer test. Figure 7.14 shows the $-\log_{10}(p_i)$ plotted against the block of 20 windows, where p_i represents the p -values obtained after performing the multivariate Cramer test.

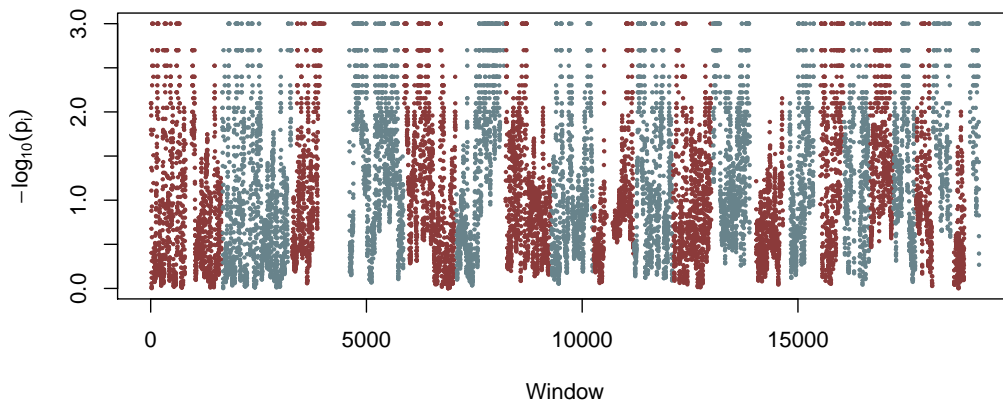


Figure 7.14: The p -values of each sliding block of 20 windows using the multivariate Cramer test. The scale of the vertical axis is $-\log_{10}(\text{p-value})$. The alternating colouring scheme indicates different chromosomes.

The Bonferroni significance threshold with $m = 17,173$ is $-\log_{10}(\frac{0.05}{17613}) = 5.55$, thus it can be seen from Figure 7.14 that only blocks of 20 windows with a p -value of 0, where $-\log_{10}(0) = \infty$, will be larger than this threshold. The number of blocks of windows with a p -value of 0 is 2,297. This corresponds to a total of 4,714 significant windows when we consider each window separately. Table 7.5 shows the number of separate windows which were identified as significant by (1) both Fisher’s combined probability test and the multivariate Cramer test, (2) Fisher’s combined probability test only, (3) the multivariate Cramer test only and (4) neither Fisher’s combined probability test or the multivariate Cramer test.

		Fisher’s Combined Probability Test	
		Significant	Not Significant
MV Cramer Test	Significant	2,312	2,402
	Not Significant	268	15,301

Table 7.5: The number of significant windows (after Bonferroni correction for Fisher’s combined probability test) when using Fisher’s combined probability test and the multivariate version of the Cramer test.

Table 7.5 shows that the multivariate version of the Cramer test identified 2,402 significant windows for which Fisher’s combined probability test identified as insignificant. Again this shows that Fisher’s combined probability test may be misidentifying significant windows because of a misplaced significance threshold. Therefore we again suggest to correct the Bonferroni corrected significance threshold by estimating the effective number of independent Fisher’s combined probability tests.

Figure 7.15 shows a plot of the p -values using the univariate Cramer test and the GPD to calculate the p -value for each window of the genome. The p -values which are identified as significant by each method are plotted on the graph using a different symbol. Note that some of the windows for which only the multivariate Cramer test identified as significant has a very small $-\log_{10}(p\text{-value})$ in Figure 7.15. It is unclear whether these windows are being correctly identified as significant. Therefore further research is needed to discover whether the results from the multivariate Cramer test is “correct”. For example, the false positive rate can be examined for the multivariate Cramer test to ensure it is properly controlled. We leave this as future work.

7.3 Segmentation Methods

Recall the segmentation techniques used to analyse CNA data per sequence in Section 1.6.2. Typically these techniques are performed prior to any further analysis

7. Multiple Testing for Dependent p -values

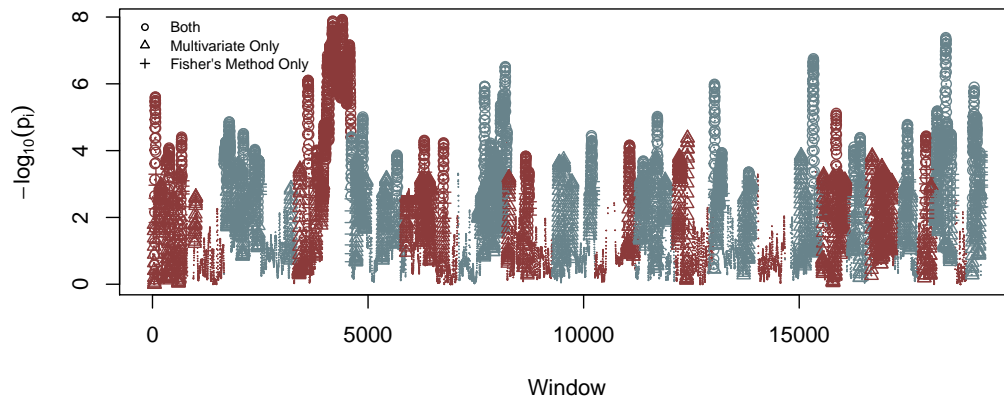


Figure 7.15: The p -values using the univariate Cramer test and the GPD to calculate the p -value for each window of the genome. The scale of the vertical axis is $-\log_{10}(p\text{-value})$. Larger circles are plotted at windows which are identified as significant by both the multivariate Cramer test and Fisher's combined probability test, triangles are plotted at windows which are identified as significant by the multivariate Cramer test only and plus signs are plotted at windows which are identified as significant by Fisher's combined probability test only. The alternating colouring scheme indicates different chromosomes.

being carried out and the lung cancer data set is not different. For our analysis we perform the Cramer test on the segmented CNA data and obtain correlated test statistics and p -values for each genomic window. In this chapter, we investigated various techniques which aims to solve the multiple testing problem for dependent p -values. However, an alternative approach would be to use the segmentation techniques described in Section 1.6.2 to segment the p -values or equivalently the test statistics.

As many segmentation techniques consider the correlation in the data, this seems like a suitable approach. There is however still a problem of choosing a suitable significance threshold. If, for example, the CBS technique was used to segment the test statistics or p -values, a suitable significance threshold correction could be to divide the chosen significance level by the total number of segments - similar to the Bonferroni correction. Figure 7.16 shows the results after applying the CBS segmentation technique to the p -values of the Cramer test applied to each window of the lung cancer data set. The total number of segments obtained after applying CBS to the p -values is 2389. The total number of significant regions is 113 which equates to a total number of 901 significant windows. This once again provides an alternative number of significant windows compared to other techniques described in this chapter and further analysis is needed to determine

which method provides the best results.

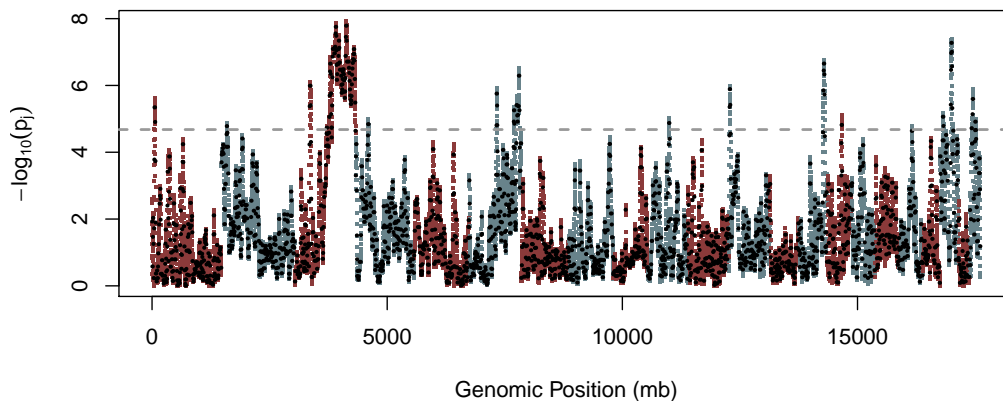


Figure 7.16: The results after applying the CBS segmentation technique to the p -values of the Cramer test applied to each window of the lung cancer data set. The scale of the vertical axis is $-\log_{10}(p\text{-value})$. The black lines represent the segments obtained after applying CBS. The horizontal grey line represents the significance threshold after dividing 0.05 by 2389 - the total number of segments. The alternating colouring scheme indicates different chromosomes.

Recall in Section 1.6.2 the use of hidden Markov models as a method for analysing CNA data per sequence by segmenting the data into various states. These methods could also be used as a way of dealing with dependent p -values. By considering only two states, namely significant or not significant, a hidden Markov model could be fitted to the p -values obtained from the test. It is therefore an alternative approach at identifying which genomic regions are significant by using the results from applying the Cramer test whilst also considering the correlation within the data. However, the challenge would be to identify the transition probabilities and emission probabilities as the transition probabilities would rely on the correlation between adjacent p -values and the emission probabilities would rely on the desired false discovery rate. This method would provide an alternative approach at dealing with multiple testing for dependent p -values and we consider this to be future work.

7.4 Discussion

In this chapter, we explore two main methods for dealing with dependent p -values in the multiple testing problem. The first method is described by [Dudbridge and Gusnanto \(2008\)](#) which enables the user to estimate m - the number of effective

7. Multiple Testing for Dependent p -values

independent tests. This method was shown to work well on a toy data set and was also applied to both the lung cancer data set for each chromosome and the artificial data set from the KC Smart vignette. The major disadvantage of this method is the speed of computation. Another disadvantage is that if the correlation structure cannot be described by a number of effective independent tests, this method may not work sufficiently.

We then looked to adapt Fisher’s combined probability test for dependent p -values when the distribution of the test statistics is unknown. To do this, we considered the joint distribution of the p -values instead of the joint distribution of the test statistics in the calculation of $\text{Var}[Z^2]$. The method was applied to a toy data set to show how it works in the case of a sliding block of k p -values and non-overlapping blocks of k p -values. We also applied the method to the lung cancer data set and the artificial data set from the KC Smart vignette. After applying Fisher’s combined probability test to the lung cancer data set, a further 2,580 windows were found to be significant. The Bonferroni correction (with the estimated effective number of independent Cramer tests) identified 929 significant windows that Fisher’s combined probability test did not. We suggest that the cause for this is due to the correlation between the Fisher’s combined probability tests which are performed simultaneously. We therefore suggest that the effective number of independent Fisher’s combined probability tests are calculated to adjust the significance threshold accordingly. This however, will make the computation much slower. When Fisher’s combined probability test was performed on chromosome 4 in the artificial dataset in the KC Smart vignette, the probes were identified as significant, which is consistent with the results obtained after applying KC Smart.

We also compared Fisher’s combined probability test to the multivariate version of the Cramer test and again found that the multivariate Cramer test identified more significant windows. We suggest however that further research is required to understand whether the windows which are identified as significant by the multivariate Cramer test is “correct”. This will therefore enable a fairer comparison between the methods.

It should be noted that other methods for combining p -values exist which are similar to Fisher’s combined probability test. Namely [Tippett et al. \(1931\)](#) suggests using $Z^2 = \min(p_i^*)$ and [Liptak \(1958\)](#) considers $Z^2 = \sum_{i=1}^k \Phi^{-1}(1 - p_i^*)$ for which [Hartung \(1999\)](#) develops a dependent p -value version with the assumption that the test statistics are normally distributed. It would be of interest to compare the results from applying these methods to the results from applying Fisher’s combined probability test and observe whether similar regions are being identified as significant. We consider this as future work.

The fact that applying Fisher’s combined probability test to the p -values iden-

tifies further genomic regions which display a significant difference of CNA between groups of patients suggests that correlation plays an important role in identifying such genomic regions. It is therefore of interest to perform further research into other techniques which can correct for the correlation in the data. As mentioned in Section 1.2.2 there are multiple ways for which correlation can be incorporated into the test. Further research could be done in all three areas and the methods compared to fully understand which approach is the best based on accuracy of results as well as speed.

Chapter 8

Discussion

In this thesis, we developed an alternative approach for identifying genomic regions of interest using Copy Number Alterations (CNA) to determine tumour subtypes. We developed our method by comparing the CNA between patients with two subtypes of lung cancer, namely adenocarcinoma and squamous carcinoma type lung cancer. Firstly, we implement and perform the Cramer test on each window and obtain the p -values by fitting a Generalised Pareto Distribution (GPD) to the null distribution. As the p -values are not all independent, we then suggest to perform Fisher's combined probability test on blocks of p -values to obtain larger regions of significance when the correlation between p -values is considered.

To find a suitable two-sample test to compare the distribution of estimated copy number alterations between two subtypes of cancer, Chapter 2 investigated well-known parametric and non-parametric tests. We show that our data displays evidence of multi-modality as well as skewness. Thus we required a test statistic which can not only deal with this type of data, but also be sensitive at identifying differences in not only the mean and variance, but also the skewness and multi-modality. Because of this, we immediately determined that the two-sample t -test and F -test would be unsuitable. We discuss the skew-adjusted t -test, which adjusts for skewness, as a potential choice but show that under certain conditions it is equivalent to Welch's t -test. As the test statistic needs to be more flexible, we focus on non-parametric tests. For testing our null hypothesis, $H_0 : F(x) \equiv G(y)$, a Cramer-von Mises type test can be implemented. We show by simulation that the Cramer test, which is a modification of the Cramer-von Mises test, is more sensitive at identifying differences between multi-modal data than the Cramer-von Mises test and Anderson-Darling test. Whilst other tests could be considered, we choose to adopt the Cramer test to identify genomic regions of significance.

We obtain the first four moments of the Cramer test statistic in Chapter 3 when the distribution function $H(t)$ is unknown. For this purpose, we can use the

8. Discussion

empirical cumulative distribution function $H_{n+m}(t)$ to approximate $H(t)$, thereby ensuring the test remains distribution free. For various known forms of $H(t)$, we also provide the expectation and variance of the test statistic. We show that the Cramer test is invariant to a linear transformation and thus prove that data can be standardised without affecting the p -value. Because the data can be standardised without affecting the p -value, this provides a method of rejecting the null hypothesis without the need for calculating a p -value.

Chapter 4 explores methods for calculating the p -value without the need for resampling. We first discuss resampling techniques, i.e. the permutation approach and the bootstrap approach for which [Baringhaus and Franz \(2004\)](#) adopts. We show that these methods are too slow and computationally expensive when we require the test to be repeated on tens of thousands of variables. We therefore investigate an empirical approach for estimating the p -value. By finding a suitable empirical approximation to the null distribution, the parameters can be obtained by using the method of moments. We show that two-parameter approximations, transformations of the test statistic and the extreme value theorem are all viable methods but do not provide a good approximation to the p -value in this case. Instead we show that the three-parameter Generalised Pareto Distribution (GPD) is a suitable empirical approximation. We compare our method to calculate the p -value to the resampling techniques in Chapter 5.

We started Chapter 5 by considering the computational burden of the test. We successfully modified the computational calculation to provide fast and accurate results. We did this by firstly considering the trapezium rule as an approximation for the integrals as well as using C++ for the calculations. We compared the speed and accuracy of our method for calculating the p -value to the permutation test and the bootstrap approach for which [Baringhaus and Franz \(2004\)](#) adopts. The results of this comparison showed that our method was much faster and just as accurate for p -values less than 0.10 as the two other methods. To locate genomic regions of interest, the test is required to be performed on over 10,000 variables simultaneously. Hence, using the GPD to estimate the null distribution and obtain a p -value seems the most appropriate method in this case.

We prove, through a simulation study in Chapter 5, that the Cramer test was a good choice of two-sample test to locate genomic regions of interest. We show that it is able to correctly control the false positive rate under different conditions whilst also being more sensitive than the Cramer-von Mises and Anderson-Darling test at identifying differences between two samples of multi-modal data. In situations where there is only a difference of mean, the Cramer test was just as sensitive as other tests including the t -test. For a difference in the variances, the F -test was still superior, but the Cramer test was no less sensitive than the Cramer-von Mises

and Anderson-Darling tests.

When comparing our two-sample test approach for identifying genomic regions of interest to KC Smart in Chapter 5, we note that KC Smart is able to identify the same significant regions as our method and more. We cannot be certain however that the regions which KC Smart identifies as significant are correctly identified. To examine this, we would need to determine whether the false positive rate is properly controlled. We leave this as future work. After segmenting the data and performing the Cramer test on each segment, we were able to identify similar significant regions as KC Smart. This indicated that the correlation between tests or p -values was affecting the number of variables which were significant. We look into this further in Chapters 6 and 7.

We end Chapter 5 by applying our method to the lung cancer dataset. We find that a large region in chromosome 3 was identified as significant. As [Belvedere et al. \(2012\)](#) states that chromosome 3 is a significant region to determine the subtype of lung cancer, we are confident that the Cramer test is correctly identifying significant regions. We also find the Cramer test identifies further significant regions which have not been previously identified. This information could be of interest to Oncologists to further help classify patients on their subtype of lung cancer.

In Chapter 6 we explore the correlation structures of the lung cancer data set. We were able to model the correlation structure adequately and find that a multivariate normal distribution is the most suitable model without being too complicated. We show that very high correlation exists between adjacent windows and the correlation remains high up until at least lag 20 for each type of lung cancer. We investigate the correlation between p -values and produce a look-up table which will give the value of this correlation when the correlation between variables is known. This table was used later on in Chapter 7.

As our method performs the Cramer test on all windows simultaneously, we end up with a multiplicity problem. Not only this, but as high correlation exists between windows, correcting for multiplicity becomes harder as the number of independent tests is unknown. In Chapter 7 we discuss multiplicity correction techniques when the p -values are not independent. We firstly discuss a method used by [Dudbridge and Gusnanto \(2008\)](#) which estimates the number of independent tests using permutation to fit a Beta(1, m) distribution to the “minimum p -values”. Once m has been estimated, the Bonferroni correction using m can be done in the usual way on the significance threshold. We applied this method to the lung cancer dataset and found that $m \ll p$, the total number of simultaneous tests performed, for each chromosome. Because of this, many more windows were identified as significant. We find however, that whilst this method is effective, it

8. Discussion

is computationally slow.

We then consider Fisher’s combined probability test for dependent p -values that tests the overall significance on groups of p -values. We extend work done by [Brown \(1975\)](#) and [Kost and McDermott \(2002\)](#) so Fisher’s combined probability test can be used when the distribution of test statistics is unknown. We demonstrate this method by applying it to some toy datasets. Then we apply Fisher’s combined probability test to the lung cancer data set and identify further significant regions. We compare the results of applying Fisher’s combined probability test to using the Bonferroni correction (with the effective number of independent Cramer tests) and the multivariate Cramer test. We conclude that Fisher’s combined probability test may be misidentifying some regions of significance due to the correlation between the tests and therefore an incorrectly placed significance threshold. Despite this however, Fisher’s combined probability test identified over 2,000 further significant windows in common with using the Bonferroni correction (with the effective number of independent Cramer tests) and the multivariate Cramer test. Further research is also needed to determine whether the significant windows identified by the multivariate Cramer test is “correct”.

We present in this thesis the results of applying a new method to identify genomic regions of interest between subtypes of lung cancer. Not many methods currently exist for the purpose of comparing multiple groups of patients. There are many examples in the literature where statistical testing was performed of regions of the genome ([Wilting et al. \(2006\)](#), [Van De Wiel and Van Wieringen \(2007\)](#) and [Smeets et al. \(2006\)](#)). We believe that the method we have created not only applies a suitable choice of test which is sensitive at identifying any differences within the data, but also accounts for the correlation within the data. Whilst it seems that KC Smart is also a capable method, it still has its disadvantages. The method relies heavily on a choice of smoothing parameter which could greatly affect the number of significant regions, something which our method avoids. We also discover in Chapters 5 and 7 that KC Smart could be overestimating the significance of some of the regions therefore leading to more significant regions being identified by this method. As we have already mentioned, further work is required to investigate the false positive rate of this method. If it is indeed the case that KC Smart is overestimating the significance and identifying too many false positives, this could also be a major disadvantage of this method. Because of all this, we believe that the method we have created provides many advantages over existing tools.

We believe that further work is required to give a more concrete justification in the use of Fisher’s combined probability test. For example, to determine whether the false positive rate is properly controlled for Fisher’s combined probability test

and the multivariate Cramer test to ensure that the windows which were identified by each test as significant were correctly identified. If the false positive rate is incorrectly controlled by Fisher's combined probability test, then the cause for this needs to be identified and corrected. If the false positive rate is incorrectly controlled for the multivariate version of the Cramer test, then depending on the cause, this could provide a justification for the use of Fisher's combined probability test.

Currently our methods only work when comparing two clinical groups of patients. We believe that further work is required in order to facilitate the comparison between more than two clinical groups. One-way ANOVA type methods currently exist which extend the Kolmogorov-Smirnov test to compare multiple groups, eg. the Kruskal-Wallis test (Kruskal and Wallis, 1952). First steps would involve research into whether this type of extension can be done for the Cramer test. To our knowledge, no literature exists which looks into comparing two or more clinical groups of patients, thus these types of methods would be a welcome addition to the literature.

A major extension to this work is the use of our research to classify new patients on their subtype of cancer. Using the methods described in this thesis, genomic markers are identified which display a significant difference in the estimated CNA between subtypes of cancer. This information could be used to formulate classification trees, however incorporating a large amount of information into the classification trees would involve extensive research to ensure the classification process is accurate as well as fast. To test accuracy, the classification techniques will be used on data sets in which the patients subtypes are known.

Appendix A

Alternative Choices of Hypothesis Tests

This appendix contains research performed which provided the stepping stones needed to reach the ultimate goal but does not contribute to the thesis overall.

A.0.1 Skew-Normal Distribution

It could be argued that the data shown in Figure 1.2 follows more closely to a skew-normal distribution (Azzalini, 1985) rather than a normal distribution. Let $X \sim \text{SN}(\mu, \sigma, \alpha)$, then the probability density function of X is defined as

$$f(x) = \frac{2}{\sigma} \phi\left(\frac{x - \mu}{\sigma}\right) \Phi\left(\alpha \left(\frac{x - \mu}{\sigma}\right)\right)$$

where $\phi(\cdot)$ and $\Phi(\cdot)$ represents the density and distribution function of the standard normal distribution respectively. The mean, variance and skewness of the distribution respectively are given by

$$\begin{aligned} E[X] &= \mu + \sigma \delta \sqrt{\frac{\pi}{2}}, \\ \text{Var}[X] &= \sigma^2 \left(1 - \frac{2\delta^2}{\pi}\right) \\ \gamma &= \frac{4 - \pi}{2} \frac{(\delta \sqrt{2/\pi})^3}{(1 - 2\delta^2/\pi)^{3/2}} \end{aligned}$$

where

$$\delta = \frac{\alpha}{\sqrt{1 + \alpha^2}}.$$

Using the method of moments, a skew-normal distribution can be fitted to each data set in Figure 1.2. Figure A.1 shows the histograms of the data for windows

A. Alternative Choices of Hypothesis Tests

448, 2023, and 9546 separated for patients with adenocarcinoma type lung cancer and patients with squamous carcinoma type lung cancer. The histograms are plotted with a solid black line representing the probability density function of the fitted skew-normal distribution. For each histogram, the number of bins was chosen to be around 10 and for comparison purposes, the histograms for each type of lung cancer is plotted on the same scale. Figure A.1 shows that a skew-normal distribution is a better representation of the data compared to a normal distribution.

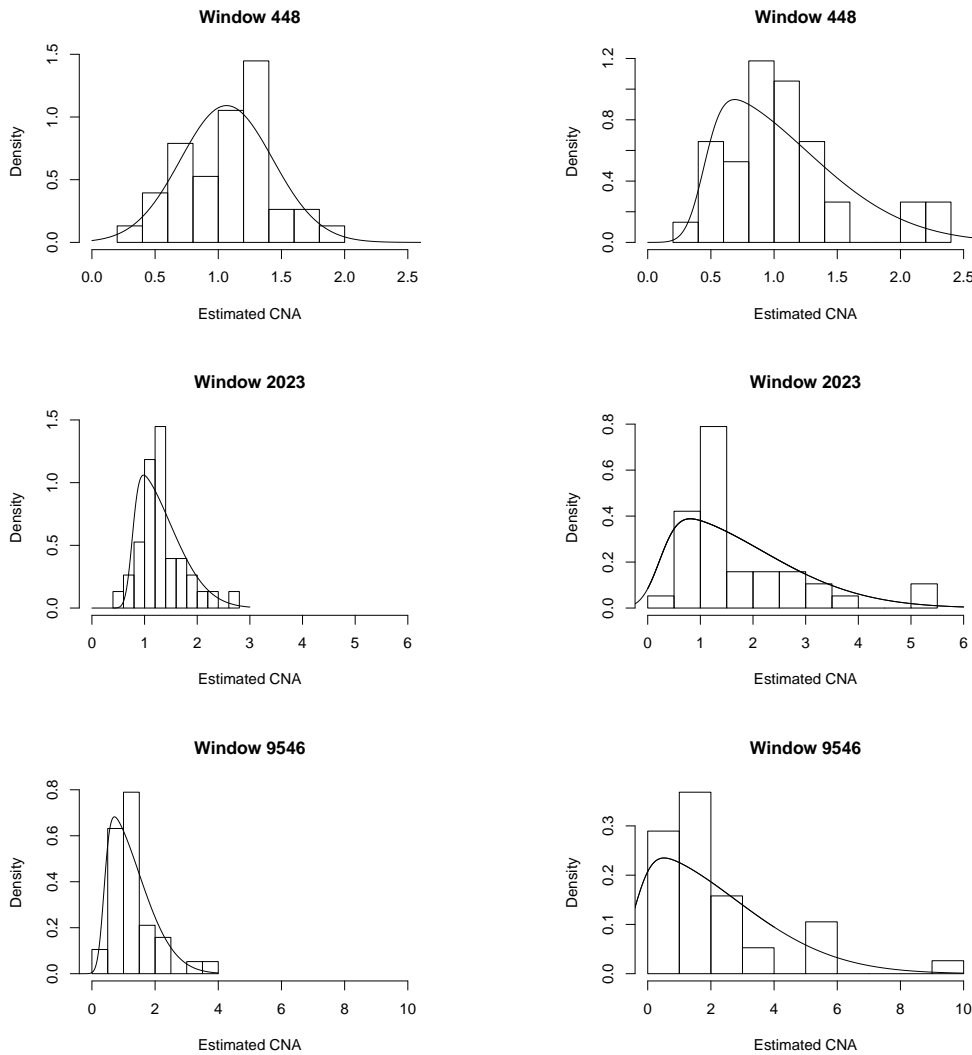


Figure A.1: Density scaled histograms of the estimated CNA in various windows across patients with adenocarcinoma type lung cancer (left) and squamous carcinoma type lung cancer (right). Each window is located at a specific position in a specific chromosome: Window 448 (67.05 – 67.2 Mbp, chromosome 1), Window 2023 (53.85 – 54 Mbp, chromosome 2) and Window 9546 (38.1 – 38.25 Mbp, chromosome 8). The solid black line represents the probability density function of the fitted skew-normal distribution.

A.1 Skew-Adjusted t test

As a skewed distribution is a closer fit to the data compared to a normal distribution, a test statistic which is adjusted for skewness is considered.

A.1.1 Asymmetric, Skew-Adjusted t -test

The test statistic of a two sample t -test is

$$t = \frac{\bar{x} - \bar{y}}{S}. \quad (\text{A.1})$$

For $S^2 = \frac{s_X^2}{n} + \frac{s_Y^2}{m}$, t is considered to be Welch's t -test where s_X^2 and s_Y^2 are the sample variances for random variables X_i and Y_j , $i = 1, \dots, n$, $j = 1, \dots, m$ respectively. Note that the two-sample t -test is commonly used to test the null hypothesis $H_0^* : \mu_X \equiv \mu_Y$ and whilst a rejection of this null hypothesis also implies a rejection of the null hypothesis in Equation (2.1) it is not always the case that if H_0^* is true, H_0 is also true.

[Balkin and Mallows \(2001\)](#) considers an asymmetric, skew-adjusted two-sample t test which slackens the assumption of normality whilst also assuming that the variance of Y_j is larger than the variance of X_i . For the second assumption, they instead use $S^2 = s_X^2 \left(\frac{1}{n} + \frac{1}{m}\right)$ and adjust Equation (A.1) for the skewness of X_i with this choice of S . The asymmetric skew-adjusted two sample t -test proposed by [Balkin and Mallows \(2001\)](#) uses the adjustment

$$t_{\text{adj}} = t_{\text{asymmetric}} + \frac{g}{6} \frac{n + 2m}{\sqrt{nm(n + m)}} \left(t_{\text{asymmetric}}^2 + \frac{m - n}{m + 2n} \right),$$

where $t_{\text{asymmetric}}$ is the t -test statistic defined in (A.1) with $S = s_X^2 \left(\frac{1}{n} + \frac{1}{m}\right)$ and g is the estimate for the third standardized moment, γ_X ;

$$g = \frac{1}{ns_X^3} \sum_{i=1}^n (x_i - \bar{x})^3,$$

and

$$\gamma_X = \frac{\text{E}[(X_i - \mu_X)^3]}{\text{Var}[X_i]^{3/2}}.$$

Note that this asymmetric version of the two sample t -test only considers the skewness of X_i and not Y_j . If both random variables are highly skewed then only considering the skewness of X_i can be considered a disadvantage of this method.

A. Alternative Choices of Hypothesis Tests

A.1.2 Using Welch's t -test

Note that the adjusted t -test described in Section A.1.1 considers the variance and skewness of only one of the random variables - in this case X_i . This means that information regarding the other random variable is not being considered, this is a disadvantage of Balkin and Mallows (2001) test. We have, however, obtained a skew-adjusted Welch's t -test, i.e. using $S^2 = \frac{s_X^2}{n} + \frac{s_Y^2}{m}$, by following the steps used by Balkin and Mallows (2001). We modify Welch's t -test using the Cornish-Fisher expansion (Cornish and Fisher, 1938); given some X_i , the Cornish-Fisher expansion of X_i is

$$X_i = CF_{X_i}(Z) = \mu + \sigma Z + \frac{\mu_3}{6\sigma^2} + \dots$$

where μ , σ and μ_3 are the mean, variance and centralised third moment of X_i respectively.

Proposition The Cornish-Fisher expansion of $\bar{X} - \bar{Y}$ under the null hypothesis $H_0 : F(x) \equiv G(y)$ is

$$CF_{\bar{X}-\bar{Y}}(Z) \approx \sigma_X \left[\sqrt{\frac{1}{n} + \frac{1}{m}} Z + \frac{\gamma_X}{6} \left(\frac{1}{n} - \frac{1}{m} \right) (Z^2 - 1) \right].$$

Proof Under the null hypothesis, we know that $\mu_Y = \mu_X$, $\sigma_Y = \sigma_X$ and $\gamma_Y = \gamma_X$. Now the mean and variance of $\bar{X} - \bar{Y}$ is

$$\mu = E[\bar{X} - \bar{Y}] = \mu_X - \mu_X = 0,$$

and

$$\sigma^2 = \text{Var}[\bar{X} - \bar{Y}] = \text{Var}[\bar{X}] + \text{Var}[\bar{Y}] = \sigma_X^2 \left(\frac{1}{n} + \frac{1}{m} \right),$$

respectively. The centralised third moment of $\bar{X} - \bar{Y}$ is

$$\begin{aligned} \mu_3 &= E[(\bar{X} - \bar{Y})^3] \\ &= E[\bar{X}^3] - 3E[\bar{X}^2]E[\bar{Y}] + 3E[\bar{X}]E[\bar{Y}^2] - E[\bar{Y}^3]. \end{aligned}$$

Now,

$$\begin{aligned}
 E[\bar{X}^3] &= E \left[\frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n X_i X_j X_k \right] \\
 &= \frac{1}{n^3} (nE[X_i^3] + 3n(n-1)E[X_i^2]E[X_j] + n(n-1)(n-2)E[X_i]E[X_j]E[X_k]) \\
 &= \frac{1}{n^2} E[X_i^3] + \frac{3(n-1)}{n^2} \mu_X (\sigma_X^2 + \mu_X^2) + \frac{(n-1)(n-2)}{n^2} \mu_X^3. \tag{A.2}
 \end{aligned}$$

The skewness of X_i can be expressed as

$$\gamma_X = \frac{E[(X_i - E[X_i])^3]}{\sigma_X^3},$$

and by expanding brackets and using the linearity property of expectations,

$$\begin{aligned}
 E[(X_i - E[X_i])^3] &= E[X_i^3] - 3E[X_i^2]\mu_X + 3E[X_i]\mu_X^2 - \mu^3 \\
 &= E[X_i^3] - 3\mu_X(\sigma_X^2 + \mu_X^2) + 2\mu^3.
 \end{aligned}$$

Thus we have

$$E[X_i^3] = \gamma_X \sigma_X^3 + 3\mu_X(\sigma_X^2 + \mu_X^2) - 2\mu_X^3,$$

and substituting this into Equation (A.2) gives

$$\begin{aligned}
 E[\bar{X}^3] &= \frac{1}{n^2} \gamma_X \sigma_X^3 + \frac{3\mu_X}{n^2} (\sigma_X^2 + \mu_X^2) - \frac{2}{n^2} \mu_X^3 + \frac{3(n-1)}{n^2} \mu_X (\sigma_X^2 + \mu_X^2) \\
 &\quad + \frac{(n-1)(n-2)}{n^2} \mu_X^3 \\
 &= \frac{1}{n^2} \gamma_X \sigma_X^3 + \frac{3\mu_X}{n} (\sigma_X^2 + \mu_X^2) + \frac{(n-3)}{n} \mu_X^3.
 \end{aligned}$$

Similarly,

$$E[\bar{Y}^3] = \frac{1}{m^2} \gamma_X \sigma_X^3 + \frac{3\mu_X}{m} (\sigma_X^2 + \mu_X^2) + \frac{(m-3)}{m} \mu_X^3.$$

It can also be shown that

$$\begin{aligned}
 E[\bar{X}^2]E[\bar{Y}] &= \mu_X \left(\frac{\sigma_X}{n} + \mu_X^2 \right) \\
 E[\bar{X}]E[\bar{Y}^2] &= \mu_X \left(\frac{\sigma_X}{m} + \mu_X^2 \right).
 \end{aligned}$$

A. Alternative Choices of Hypothesis Tests

Then

$$\begin{aligned} \mathbb{E}[(\bar{X} - \bar{Y})^3] &= \gamma_X \sigma_X^3 \left(\frac{1}{n^2} - \frac{1}{m^2} \right) + 3\mu_X (\sigma_X^2 + \mu_X) \left(\frac{1}{n} - \frac{1}{m} \right) \\ &\quad - 3\mu_X^3 \left(\frac{1}{n} - \frac{1}{m} \right) - 3\mu_X \sigma_X^2 \left(\frac{1}{n} - \frac{1}{m} \right) \\ &= \gamma_X \sigma_X^3 \left(\frac{1}{n^2} - \frac{1}{m^2} \right). \end{aligned}$$

Thus the Cornish-Fisher expansion of $\bar{X} - \bar{Y}$ is

$$\begin{aligned} CF_{\bar{X}-\bar{Y}}(Z) &\approx \sigma_X \sqrt{\frac{1}{n} + \frac{1}{m}} Z + \frac{\gamma_X \sigma_X^3 \left(\frac{1}{n^2} - \frac{1}{m^2} \right)}{6\sigma_X^2 \left(\frac{1}{n} + \frac{1}{m} \right)} (Z^2 - 1) \\ &= \sigma_X \sqrt{\frac{1}{n} + \frac{1}{m}} Z + \frac{\gamma_X}{6} \sigma_X \left(\frac{1}{n} - \frac{1}{m} \right) (Z^2 - 1) \\ &= \sigma_X \left[\sqrt{\frac{1}{n} + \frac{1}{m}} Z + \frac{\gamma_X}{6} \left(\frac{1}{n} - \frac{1}{m} \right) (Z^2 - 1) \right]. \end{aligned}$$

□

Proposition The Cornish-Fisher expansion of $S^2 = \frac{s_X^2}{n} + \frac{s_Y^2}{m}$ under the null hypothesis is approximated by

$$CF_{S^2}(W) \approx \sigma_X^2 \left(\frac{1}{n} + \frac{1}{m} \right) \left[1 + \sqrt{\frac{n^3 + m^3}{nm(n+m)^2}} \sqrt{\gamma_2} W \right]$$

where

$$\sqrt{\gamma_2} = \frac{\mu_4^X - \sigma_X^4}{\sigma_X^4}.$$

Proof The mean of S^2 is

$$\begin{aligned} \mathbb{E}[S^2] &= \mathbb{E} \left[\frac{s_X^2}{n} + \frac{s_Y^2}{m} \right] \\ &= \frac{1}{n} \mathbb{E}[s_X^2] + \frac{1}{m} \mathbb{E}[s_Y^2] \\ &= \frac{1}{n} \sigma_X^2 + \frac{1}{m} \sigma_Y^2 \\ &= \sigma_X^2 \left(\frac{1}{n} + \frac{1}{m} \right), \end{aligned}$$

as $\sigma_X^2 = \sigma_Y^2$ under the null hypothesis. The variance of s_X^2 is proved by [Cho and Cho \(2009\)](#) and is

$$\begin{aligned}\text{Var}[s_X^2] &= \frac{1}{n} \left(\mu_4^X - \frac{n-3}{n-1} \sigma_X^4 \right) \\ &= \frac{1}{n} (\mu_4^X - \sigma_X^4) + \mathcal{O}(n^{-2})\end{aligned}$$

where μ_4^X is the fourth centralised moment. Thus,

$$\begin{aligned}\text{Var}[S^2] &= \text{Var} \left[\frac{s_X^2}{n} + \frac{s_Y^2}{m} \right] \\ &= \frac{1}{n^2} \text{Var}[s_X^2] + \frac{1}{m^2} \text{Var}[s_Y^2] \\ &= \frac{1}{n^3} (\mu_4^X - \sigma_X^4) + \frac{1}{m^3} (\mu_4^Y - \sigma_Y^4) \\ &= (\mu_4^X - \sigma_X^4) \left(\frac{1}{n^3} + \frac{1}{m^3} \right),\end{aligned}$$

as $\mu_4^X = \mu_4^Y$ under the null hypothesis. Hence, the Cornish-Fisher expansion of S^2 is approximated by

$$CF_{S^2}(W) \approx \sigma_X^2 \left(\frac{1}{n} + \frac{1}{m} \right) \left[1 + \sqrt{\frac{n^3 + m^3}{nm(n+m)^2}} \sqrt{\gamma_2} W \right].$$

□

Now, assume that the skew-adjusted t -test is of the form

$$t_{\text{adj}} = t + \lambda + \nu t^2. \tag{A.3}$$

By noting that we can write t as the ratio of $CF_{\bar{X}-\bar{Y}}$ and CF_{S^2} , substituting the Cornish-Fisher expansions into the adjusted t -test statistic in Equation (A.3), we get

$$\begin{aligned}t_{\text{adj}} &= \left[Z + \frac{\gamma}{6} \frac{m-n}{\sqrt{nm(n+m)}} (Z^2 - 1) \right] \times \left(1 + \sqrt{\frac{m^3 + n^3}{nm(n+m)^2}} \sqrt{\gamma_2} W \right)^{-\frac{1}{2}} \\ &\quad + \lambda + \nu \left[Z^2 \left(1 + \sqrt{\frac{m^3 + n^3}{nm(n+m)^2}} \sqrt{\gamma_2} W \right)^{-1} \right] + \mathcal{O}(n^{-\frac{1}{2}}).\end{aligned} \tag{A.4}$$

A. Alternative Choices of Hypothesis Tests

Using the Taylor series expansion, Equation (A.4) becomes

$$t_{\text{adj}} = Z + Z^2 \left[\frac{\gamma_X}{6} \frac{m-n}{\sqrt{mn(n+m)}} + \nu \right] - \frac{\gamma_X}{6} \frac{m-n}{\sqrt{mn(n+m)}} + \lambda - \frac{1}{2} \sqrt{\frac{m^3+n^3}{nm(n+m)^2}} \sqrt{\gamma_2} W Z.$$

Now, it can be shown that

$$\begin{aligned} \text{Cov}(\bar{X} - \bar{Y}, s_X^2) &= \frac{\mu_3^X}{n} \\ \text{Cov}(\bar{X} - \bar{Y}, s_Y^2) &= -\frac{\mu_3^Y}{m}. \end{aligned}$$

Thus

$$\begin{aligned} \text{Cov}\left(\bar{X} - \bar{Y}, \frac{s_X^2}{n} + \frac{s_Y^2}{m}\right) &= \frac{\mu_3^X}{n^2} - \frac{\mu_3^Y}{m^2} \\ &= \mu_3^X \left(\frac{1}{n^2} - \frac{1}{m^2} \right) \end{aligned}$$

as $\mu_3^X = \mu_3^Y$ under the null hypothesis. Then using the Cornish-Fisher expansions,

$$\begin{aligned} \mu_3^X \left(\frac{1}{n^2} - \frac{1}{m^2} \right) &= \text{Cov}\left(\sigma_X \sqrt{\frac{1}{n} + \frac{1}{m}} Z, \sigma^2 \sqrt{\frac{1}{m^3} + \frac{1}{n^3}} \sqrt{\gamma_2} W\right) \\ &= \sigma_X^3 \frac{\sqrt{(n+m)(n^3+m^3)}}{(nm)^2} \sqrt{\gamma_2} \text{Cov}(Z, W). \end{aligned}$$

Thus

$$\text{Cov}(Z, W) = \frac{\gamma_X}{\sqrt{\gamma_2}} (m-n) \sqrt{\frac{n+m}{n^3+m^3}}.$$

Now, let $W = \frac{\gamma_X}{\sqrt{\gamma_2}} (m-n) \sqrt{\frac{n+m}{n^3+m^3}} Z + Z^*$ where Z and Z^* are both independent normally distributed random variables. Then we have

$$\begin{aligned} t_{\text{adj}} &= Z + Z^2 \left[\frac{\gamma_X}{6} \frac{m-n}{\sqrt{mn(n+m)}} + \nu - \frac{\gamma_X}{2} \frac{m-n}{\sqrt{mn(n+m)}} \right] - \frac{\gamma_X}{6} \frac{m-n}{\sqrt{mn(n+m)}} \\ &\quad + \lambda - \frac{1}{2} \sqrt{\frac{m^3+n^3}{nm(n+m)^2}} \sqrt{\gamma_2} Z Z^*. \end{aligned}$$

Now we solve for λ and ν so as to remove the constant term and the term in Z^2

respectively. Thus

$$\begin{aligned}\lambda &= \frac{\gamma_X}{6} \frac{m-n}{\sqrt{mn(n+m)}} \\ \nu &= \frac{\gamma_X}{3} \frac{m-n}{\sqrt{mn(n+m)}}.\end{aligned}$$

Hence, the skew-adjusted t -test statistic is

$$t_{\text{adj}} = t + \frac{g}{6} \frac{m-n}{\sqrt{mn(n+m)}} (2t^2 + 1).$$

Note that when the sample size is the same, i.e. $n = m$ the adjusted t -test will be equal to Welch's t -test. Thus this adjustment is only applicable when the sample sizes are unequal. For the lung cancer data set, $n = m = 36$, thus applying the skew-adjusted t -test will be equivalent to applying Welch's t -test which, as we have already stated, is not an optimal choice of test for our data.

Appendix B

Alternative Methods to Obtain a Suitable Null Distribution

In this appendix, we discuss other methods which were considered when finding a suitable null distribution.

B.0.1 The (Scaled) Chi-Square Distribution

Recall that $T_{n,m}$ is distributed as a scaled Chi-Square distribution if $T_{n,m} \sim c\chi_f^2$. To estimate the parameters of the scaled Chi-Square distribution, we use method of moments as we can easily estimate $E[T_{n,m}]$ and $\text{Var}[T_{n,m}]$ from the sampled test statistics and in an application setting we can use the formulas obtained in Chapter 3, Section 3.3 using H_{n+m} in place of H . The formulas to calculate the parameters of the scaled Chi-Square distribution are

$$f = \frac{2E[T_{n,m}]^2}{\text{Var}[T_{n,m}]}, \quad c = \frac{\text{Var}[T_{n,m}]}{2E[T_{n,m}]}.$$

When X and Y both follow a $N(0, 1)$ distribution, the estimates for c and f are $c = 0.174$ and $f = 3.353$. Figure B.1 shows the QQ-plot comparing the percentiles of the 10000 test statistics and the percentiles of the fitted scaled Chi-Square distribution. As the points in the QQ-plot do not lie on a straight line, it is clear that the scaled Chi-Square distribution is not a suitable fit to the null distribution. We therefore do not test to see whether the Chi-Square distribution will work when X and Y follow a multi-modal distribution with 2 or 3 peaks and try a different choice of distribution for $\nu(z)$.

B. Alternative Methods to Obtain a Suitable Null Distribution

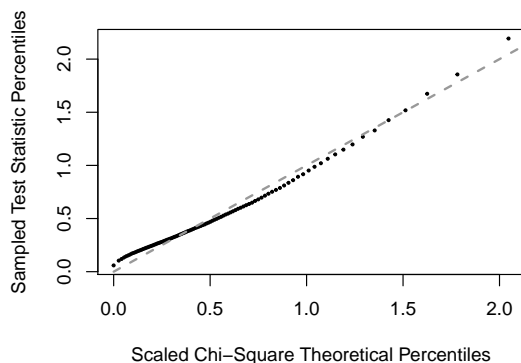


Figure B.1: The percentiles of the 10000 sampled test statistics plotted against the percentiles of the fitted scaled Chi-Square distribution with $c = 0.174$ and $f = 3.353$.

B.0.2 The Gamma Distribution

Figure 4.1 suggests that a Gamma distribution could be suitable. Thus we shall next consider the Gamma distribution as an estimate of the null distribution. Again, method of moments can be used to estimate the parameters α and β . The formulas to calculate the parameters of the Gamma distribution are

$$\alpha = \frac{E[T_{n,m}]^2}{\text{Var}[T_{n,m}]}, \quad \beta = \frac{E[T_{n,m}]}{\text{Var}[T_{n,m}]}.$$

When X and Y both follow a $N(0,1)$ distribution, the estimates for α and β are $\alpha = 1.68$ and $\beta = 2.88$. Figure B.2 shows the QQ-plot comparing the percentiles of the 10000 test statistics and the percentiles of the fitted Gamma distribution. As the points in the QQ-plot do not lie on a straight line, it is clear that the Gamma distribution is not a suitable fit to the null distribution. We therefore do not test to see whether the Gamma distribution will work when X and Y follow a multi-modal distribution with 2 or 3 peaks and try a different choice of distribution for $\nu(z)$.

B.0.3 The Log-Normal Distribution

Next, we consider a log-normal distribution as an estimate to the null distribution. Again, method of moments can be used to estimate the parameters μ and σ . The

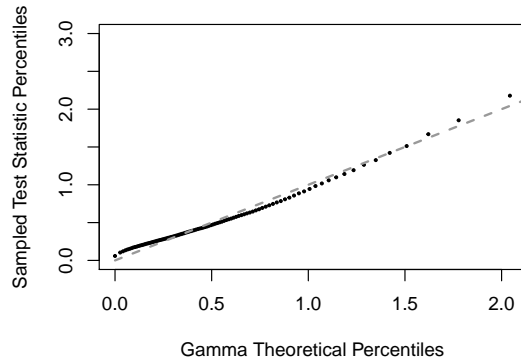


Figure B.2: The percentiles of the 10000 sampled test statistics plotted against the percentiles of the fitted Gamma distribution with $\alpha = 1.63$ and $\beta = 2.93$.

formulas to calculate the parameters of the log-normal distribution are

$$\mu = \ln \left(\frac{E[T_{n,m}]}{\sqrt{1 + \frac{\text{Var}[T_{n,m}]}{E[T_{n,m}]^2}}} \right) \quad \sigma^2 = \ln \left(1 + \frac{\text{Var}[T_{n,m}]}{E[T_{n,m}]^2} \right).$$

When X and Y both follow a $N(0, 1)$ distribution, the estimates for μ and σ are $\mu = -0.82$ and $\sigma = 0.69$. Figure B.3 shows the QQ-plot comparing the percentiles of the 10000 test statistics and the percentiles of the fitted log-normal distribution. As the points in the QQ-plot lie close to a straight line, this indicates that the log-normal distribution could be a reasonable fit for the null distribution.

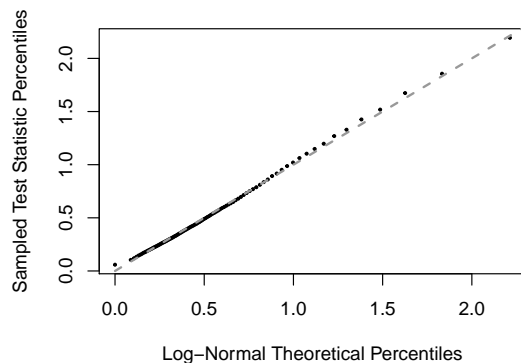


Figure B.3: The percentiles of the 10000 sampled test statistics plotted against the percentiles of the fitted log-normal distribution with $\mu = -0.82$ and $\sigma = 0.69$.

When X and Y both follow a $\frac{7}{8}N(1, 0.25) + \frac{1}{8}N(3, 0.25)$ mixture distribution, the estimates for μ and σ are $\mu = -1.47$ and $\sigma = 0.78$. Figure B.4 shows the

B. Alternative Methods to Obtain a Suitable Null Distribution

QQ-plot comparing the percentiles of the 10000 test statistics and the percentiles of the fitted log-normal distribution. It is now clear that when X and Y are multi-modal, the log-normal distribution is not an accurate enough fit for the null distribution. As multi-modality is a common feature in the lung cancer dataset, it is important that the estimate of the null distribution is accurate enough when X and Y is multi-modal.

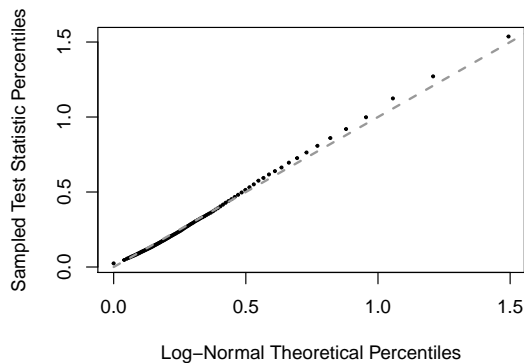


Figure B.4: The percentiles of the 10000 sampled test statistics plotted against the percentiles of the fitted log-normal distribution with $\mu = -1.47$ and $\sigma = 0.78$.

B.0.4 Other Two-Parameter Distributions

Whilst the Chi-Square, Gamma and log-normal distributions are the only distributions presented in this chapter, other two-parameter distributions were also considered as a choice for the null distribution. Many suitable two-parameter distributions such as the two-parameter Pareto, Beta-Prime, F, and Weibull distributions were considered. The parameters were estimated using method of moments and the fitted cumulative distribution curves were compared to the empirical cumulative distribution curve when the underlying distributions were normally and multi-modally distributed. In almost all cases, the fitted CDFs were not a close match to the ECDF when the underlying distributions were normally distributed. In the case of the Beta-Prime distribution, the fitted CDF was a close match for the ECDF when the underlying distributions were normally distributed. However, in all cases, the fitted CDFs were not a close match to the ECDF when the underlying distributions were multi-modally distributed. Thus it became clear that a two-parameter distribution would not be a suitable choice for the null distribution.

B.1 Transformation of $T_{n,m}$

As no two-parameter distributions were found to be suitable for estimating the null distribution, it may be the case that after transforming the data, a two parameter distribution will be a reasonable fit for the transformed data. We consider a transformation of the form $T_{n,m}^\psi$, $\psi \in \mathbb{R}$. It is important however, that the parameters of the null distribution for the transformed test statistic can still be calculated given only the mean and variance of the untransformed test statistic. Because of this, the choice of two parameter distributions for the null distribution of the transformed test statistics is limited. Consider $U = T_{n,m}^\psi$ where t_1, \dots, t_n and u_1, \dots, u_n are sampled data from random variables $T_{n,m}$ and U respectively.

B.1.1 Finding the Optimal ψ .

If U follows a two parameter distribution, there must exist a ψ such that $T_{n,m}^\psi$ also follows the distribution. When performing a hypothesis test, the significance level of the test is usually less than 10%. Thus, it is more important for the null distribution approximation to be accurate in the right tail as opposed to the left tail. When finding an optimal ψ only the fit in the right tail will be optimised for $\psi > 0$ and the fit in the left tail will be optimised for $\psi < 0$. More specifically, we will consider optimising the fit of the distribution between the 90th and 100th percentiles when $\psi > 0$ and between the 0th and 10th percentiles when $\psi < 0$.

If $U = T_{n,m}^\psi$ and the distribution of U is known, then

$$RT_{acc} = \int_{0.9}^1 (p - \mathcal{N}_k(\mathcal{N}^{-1}(p)))^2 dp$$

will give a measure for how suitable U is for the right tail of the null distribution and

$$LT_{acc} = \int_0^{0.1} (p - \mathcal{N}_k(\mathcal{N}^{-1}(p)))^2 dp$$

will give a measure for how suitable U is for the left tail of the null distribution. To find the value of ψ for which $T_{n,m}^\psi$ follows the same distribution as U , RT_{acc} and LT_{acc} are minimised over $\psi \in (0, 1]$ and $\psi \in [-1, 0)$ respectively.

B.1.2 Transforming $T_{n,m}$ to a Log-Normal Distribution

Suppose $U \sim \text{LN}(\mu, \sigma)$, then

$$\nu(u) = \frac{1}{u\sigma\sqrt{2\pi}} e^{-\frac{(\log u - \mu)^2}{2\sigma^2}}.$$

B. Alternative Methods to Obtain a Suitable Null Distribution

The probability density function of $T_{n,m}$ can be calculated using

$$\nu(t) = \left| \frac{du}{dt} \right| \nu(u).$$

As

$$\begin{aligned} \left| \frac{du}{dt} \right| &= |\psi t^{\psi-1}| \\ &= |\psi| t^{\psi-1}, \end{aligned}$$

the probability density function for $T_{n,m}$ is thus

$$\begin{aligned} \nu(t) &= |\psi| t^{\psi-1} \frac{1}{t^\psi \sigma \sqrt{2\pi}} e^{-\frac{(\log t^\psi - \mu)^2}{2\sigma^2}} \\ &= \frac{1}{t(\sigma/|\psi|)\sqrt{2\pi}} e^{-\frac{(\log t - (\mu/|\psi|))^2}{2(\sigma/|\psi|)^2}}. \end{aligned}$$

Thus if $U \sim \text{LN}(\mu, \sigma)$, then $T_{n,m} \sim \text{LN}(\frac{\mu}{|\psi|}, \frac{\sigma}{|\psi|})$. As seen in Section 4.3, the null distribution of the untransformed test statistics did not follow a Log-Normal distribution, thus it can be concluded that the transformed test statistics will also not follow a Log-Normal distribution.

B.1.3 Transforming $T_{n,m}$ to a Gamma Distribution

Suppose $U \sim \text{Gamma}(\alpha, \beta)$, then

$$\nu(u) = \frac{\beta^\alpha u^{\alpha-1} e^{-u\beta}}{\Gamma(\alpha)}.$$

The probability density function of $T_{n,m}$ can therefore be calculated using

$$\nu(t) = \left| \frac{du}{dt} \right| \nu(u).$$

As

$$\begin{aligned} \left| \frac{du}{dt} \right| &= |\psi t^{\psi-1}| \\ &= |\psi| t^{\psi-1}, \end{aligned}$$

the probability density function for $T_{n,m}$ becomes

$$\begin{aligned}\nu(t) &= |\psi| t^{\psi-1} \frac{\beta^\alpha u^{\alpha-1} e^{-u\beta}}{\Gamma(\alpha)} \\ &= |\psi| t^{\psi-1} \frac{\beta^\alpha t^{\psi(\alpha-1)} e^{-t^\psi \beta}}{\Gamma(\alpha)} \\ &= \frac{|\psi| \beta^\alpha t^{\psi\alpha-1} e^{-t^\psi \beta}}{\Gamma(\alpha)}.\end{aligned}$$

The moments of this distribution can be calculated by

$$\begin{aligned}\mathbb{E}[T_{n,m}^d] &= \int_0^\infty t^d \frac{|\psi| \beta^\alpha t^{\psi\alpha-1} e^{-t^\psi \beta}}{\Gamma(\alpha)} dt \\ &= \int_0^\infty \frac{|\psi| \beta^\alpha t^{d+\psi\alpha-1} e^{-t^\psi \beta}}{\Gamma(\alpha)} dt \\ &= \frac{|\psi| \beta^\alpha}{\Gamma(\alpha)} \int_0^\infty t^{d+\psi\alpha-1} e^{-\left(\frac{\beta^{-\frac{1}{\psi}}}{t}\right)^{-\psi}} dt.\end{aligned}\tag{B.1}$$

Consider now the generalised inverse gamma distribution, with probability density function

$$f(x; \delta, \epsilon, \zeta) = \frac{\zeta \epsilon^{\delta\zeta}}{\Gamma(\delta)} x^{-\delta\zeta-1} e^{-\left(\frac{\epsilon}{x}\right)^\zeta}.$$

Note that the integrand in equation (B.1) is proportional to the probability density function for the generalised inverse gamma distribution with parameters

$$\delta = \alpha + \frac{d}{\psi}; \quad \epsilon = \beta^{-\frac{1}{\psi}}; \quad \zeta = -\psi.$$

Thus

$$\begin{aligned}\mathbb{E}[T_{n,m}^d] &= -\frac{|\psi| \beta^\alpha \Gamma(\alpha + \frac{d}{\psi})}{\psi \beta^{\alpha + \frac{d}{\psi}} \Gamma(\alpha)} \int_0^\infty -\frac{\psi \beta^{\alpha + \frac{d}{\psi}}}{\Gamma(\alpha + \frac{d}{\psi})} t^{d+\psi\alpha-1} e^{-\left(\frac{\beta^{-\frac{1}{\psi}}}{t}\right)^{-\psi}} dt \\ &= -\frac{|\psi| \beta^\alpha \Gamma(\alpha + \frac{d}{\psi})}{\psi \beta^{\alpha + \frac{d}{\psi}} \Gamma(\alpha)} \\ &= -\frac{|\psi| \beta^{-\frac{d}{\psi}} \Gamma(\alpha + \frac{d}{\psi})}{\psi \Gamma(\alpha)}.\end{aligned}\tag{B.2}$$

We can therefore solve a system of two equations when $d = 1$ and $d = 2$ to calculate the parameters of α and β using the formulas for $\mathbb{E}[T_{n,m}]$ and $\text{Var}[T_{n,m}]$.

If U is assumed to follow a Gamma distribution with parameters α and β , then the value of ψ can be found such that $T_{n,m}^\psi$ also follows the Gamma distribution

B. Alternative Methods to Obtain a Suitable Null Distribution

with parameters α and β . Consider the test statistics when X and Y followed a $N(0, 1)$ distribution. For $\psi \in (0, 1]$, RT_{acc} has been calculated and the results plotted in Figure B.5. The plot shows that for $\psi = 0.81$, $T_{n,m}^\psi$ is equal to the distribution of U .

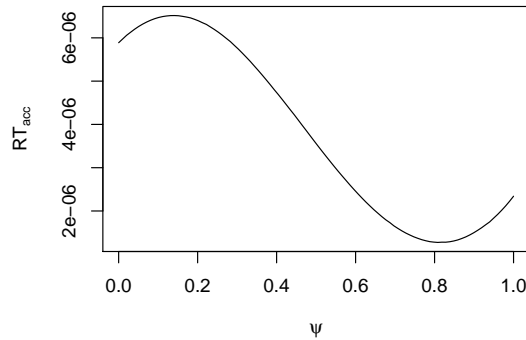


Figure B.5: The values of RT_{acc} plotted against $\psi \in (0, 1]$.

For $\psi \in [-1, 0)$, LT_{acc} has been calculated and the results plotted in Figure B.6. The plot shows that for $\psi = -0.4$, $T_{n,m}^\psi$ is equal to the distribution of U .

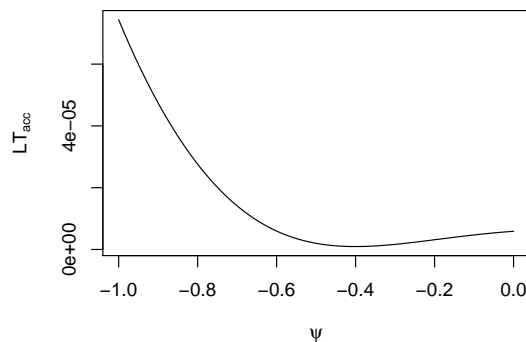


Figure B.6: The values of LT_{acc} plotted against $\psi \in [-1, 0)$.

Now that the values of ψ for which $T_{n,m}^\psi$ is equal to the distribution of U has been found, it should be tested to see if transformations of the form $U_1 = T_{n,m}^{0.81}$ and $U_2 = T_{n,m}^{-0.4}$ indeed follow a Gamma distribution. Thus Figure B.7 shows the QQ-plots comparing the percentiles of the fitted Gamma distributions against U_1 (left) and U_2 (right).

Clearly, Figure B.7 shows that the value of $\psi = -0.4$ provides a much closer fit to a Gamma distribution in the left tail than the value of $\psi = 0.81$ in the right

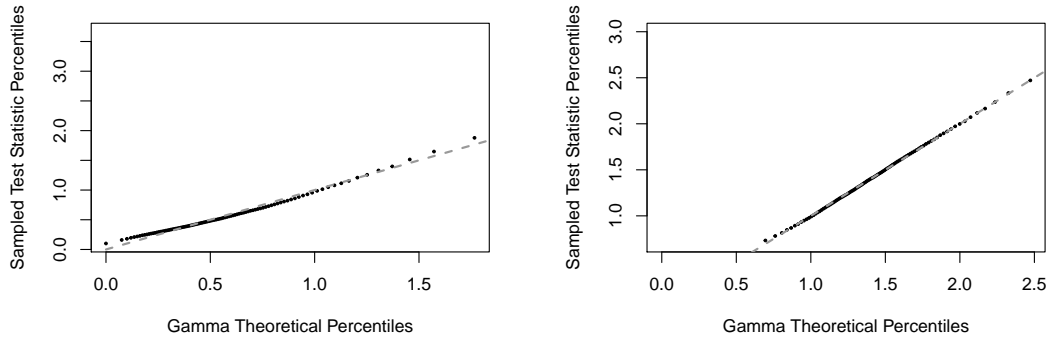


Figure B.7: The QQ-plots comparing the percentiles of the fitted Gamma distributions against U_1 (left) and U_2 (right).

tail. To test the fit in the left tail for $\psi = -0.4$, a Kolmogorov-Smirnov goodness of fit test can be performed. The KS test is useful here simply to determine how well a Gamma distribution fits to the random variable U . Here we have sampled test statistics t_1, \dots, t_{10000} where $X \sim N(0,1)$, $Y \sim N(0,1)$, $n = m = 10000$ and $k = 10000$ and after transforming we have u_1, \dots, u_{10000} . Consider taking the 1000 smallest values from the sample u_1, \dots, u_{10000} , to give a smaller sample u_1^*, \dots, u_{1000}^* . Now sample 10000 values from $V \sim \text{Gamma}(\alpha, \beta)$ where α and β are estimated from the data u_1, \dots, u_{10000} . Again, take the 1000 smallest values from the sample to obtain v_1^*, \dots, v_{1000}^* . The Kolmogorov-Smirnov test on the two samples u^* and v^* is performed and a p -value calculated. When this process is repeated 1000 times, 676 of those times gave a p -value less than 0.05, suggesting that U might not closely fit a Gamma distribution. If this is the case, then obtaining an accurate p -value using this method is unlikely.

B.2 Extreme Value Theorem

Given that so far an attempt to fit a two-parameter distribution to the null distribution has been unsuccessful, it may be worth fitting a distribution to only the tail of the null distribution. Recall that it is more important to have an accurate fit between the 90th and 100th percentiles for calculating a p -value. Considering this, it was natural to consider the extreme value theorem and in particular the Pickands-Balkema-de Haan Theorem (Balkema and De Haan, 1974), (Pickands III, 1975). The theorem states that under certain maximum domain of attraction conditions, the limiting distribution of the excesses above a certain threshold λ is a three-parameter Generalised Pareto Distribution (GPD) with parameters μ , σ and

B. Alternative Methods to Obtain a Suitable Null Distribution

ξ . The density function of the Generalised Pareto distribution is

$$f(t; \mu, \sigma, \xi) = \frac{1}{\sigma} \left(1 + \frac{\xi(t - \mu)}{\sigma} \right)^{\left(-\frac{1}{\xi}-1\right)}, \quad (\text{B.3})$$

for $t \geq \mu$ when $\xi \geq 0$, and $\mu \leq t \leq \mu - \frac{\sigma}{\xi}$ when $\xi < 0$. The parameters μ , σ and ξ can be estimated using the following formulas;

$$\begin{aligned} \mu &= \lambda \\ \sigma &= \frac{1}{2}(\bar{t}_\lambda - \lambda) \left(\frac{(\bar{t}_\lambda - \lambda)^2}{s_u^2} + 1 \right) \\ \xi &= \frac{1}{2} \left(\frac{(\bar{t}_\lambda - \lambda)^2}{s_u^2} - 1 \right), \end{aligned}$$

where \bar{t}_λ and s_λ represent the mean and standard deviation of the excesses above the threshold λ .

The cumulative distribution function of the excesses above the threshold λ is defined to be

$$\mathcal{N}_\lambda(t) = \Pr\{T_{n,m} - \lambda \leq t | T_{n,m} > \lambda\}.$$

Now, the cumulative distribution function for $T_{n,m}$ can be rewritten as

$$\begin{aligned} \mathcal{N}(t) &= \Pr\{T_{n,m} \leq t\} = \Pr\{T_{n,m} \leq \lambda\} + \Pr\{\lambda < T_{n,m} \leq t\} \\ &= \Pr\{T_{n,m} \leq \lambda\} + \Pr\{T_{n,m} \leq t \cap T_{n,m} > \lambda\} \\ &= \Pr\{T_{n,m} \leq \lambda\} + \Pr\{T_{n,m} \leq t | T_{n,m} > \lambda\} \Pr\{T_{n,m} > \lambda\} \\ &= \Pr\{T_{n,m} \leq \lambda\} + \Pr\{T_{n,m} \leq t | T_{n,m} > \lambda\} (1 - \Pr\{T_{n,m} \leq \lambda\}). \end{aligned}$$

and for $t \geq \lambda$,

$$\begin{aligned} \mathcal{N}(t) &= \Pr\{T_{n,m} \leq \lambda\} + \mathcal{N}_\lambda(t - \lambda)(1 - \Pr\{T_{n,m} \leq \lambda\}) \\ &= \mathcal{N}(\lambda) + \mathcal{N}_\lambda(t - \lambda)(1 - \mathcal{N}(\lambda)). \end{aligned}$$

It is seen in Proposition 3 in Appendix C that $\hat{\mathcal{N}}(t)$ also follows a Generalised Pareto distribution with parameters

$$\begin{aligned} \tilde{\xi} &= \xi \\ \tilde{\sigma} &= \sigma(1 - \mathcal{N}(\lambda))^{\tilde{\xi}} \\ \tilde{\lambda} &= \lambda - \frac{\tilde{\sigma}}{\tilde{\xi}} \left((1 - \mathcal{N}(\lambda))^{-\tilde{\xi}} - 1 \right). \end{aligned}$$

Thus, once a threshold λ is chosen, and \bar{t}_λ and s_λ is known, the fitted Generalised

Pareto distribution for the tail can be estimated. As has been mentioned previously, for p -value estimation, a suitable choice of λ is the 90th percentile of the test statistics. In our application it is very easy to calculate the expectation and variance of $T_{n,m}$, however obtaining estimates for the mean and variance of the test statistic for $T_{n,m} \geq \lambda$ is difficult. Calculating the moments for $T_{n,m} \geq \lambda$ directly proved to be impossible and no strong relationship between $E[T_{n,m}]$, $\text{Var}[T_{n,m}]$ and $E[T_{n,m} \geq \lambda]$, $\text{Var}[T_{n,m} \geq \lambda]$ was found. Thus using the Pickands-Balkema-de Haan Theorem to estimate the p -value was not successful.

Appendix C

Miscellaneous Propositions

This appendix contains various propositions not included in the main text.

Proposition 1 Let $\phi(x)$ and $\Phi(x)$ be the density function and distribution function of a random variable X where $X \sim N(\mu, \sigma)$. Then the following statements are true;

$$\int \Phi(x) dx = (x - \mu)\Phi(x) + \sigma^2\phi(x) + C \quad (\text{C.1})$$

$$\int \Phi(x)^2 dx = (x - \mu)\Phi(x)^2 + 2\sigma^2\Phi(x)\phi(x) - \frac{\sigma}{\sqrt{\pi}}\Phi\left(\sqrt{2}x + \mu(1 - \sqrt{2})\right) + C. \quad (\text{C.2})$$

Proof First we prove (C.1). Using integration by parts we can show that

$$\int \Phi(x) dx = x\Phi(x) - \int x\phi(x) dx + C.$$

Now,

$$\begin{aligned} \int x\phi(x) dx &= \int x \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\} dx \\ &= \int (x + \mu) \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2\sigma^2}\right\} dx \\ &= \int x \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2\sigma^2}\right\} dx + \mu \int \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2\sigma^2}\right\} dx \\ &= \frac{1}{\sigma\sqrt{2\pi}} \int x \exp\left\{-\frac{x^2}{2\sigma^2}\right\} dx + \mu\Phi(x) \end{aligned}$$

C. Miscellaneous Propositions

$$\begin{aligned}
 &= -\sigma^2 \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2\sigma^2}\right\} + \mu\Phi(x) + C \\
 &= -\sigma^2\phi(x) + \mu\Phi(x) + C.
 \end{aligned}$$

Thus

$$\begin{aligned}
 \int \Phi(x) dx &= x\Phi(x) + \sigma^2\phi(x) - \mu\Phi(x) + C \\
 &= (x - \mu)\Phi(x) + \sigma^2\phi(x) + C.
 \end{aligned}$$

Next we prove (C.2). Using integration by parts we can show that

$$\begin{aligned}
 \int \Phi(x)^2 dx &= x\Phi(x)^2 - \int 2x\phi(x)\Phi(x) dx + C \\
 &= x\Phi(x)^2 + 2\sigma^2\phi(x)\Phi(x) - 2\mu\Phi(x)^2 + \int -2\sigma^2\phi(x)^2 \\
 &\quad + 2\mu\phi(x)\Phi(x) dx + C.
 \end{aligned}$$

Now,

$$\begin{aligned}
 \int \phi(x)^2 dx &= \int \frac{1}{2\pi\sigma^2} \exp\left\{-\frac{(x-\mu)^2}{\sigma^2}\right\} dx \\
 &= \frac{1}{\sigma\sqrt{2\pi}} \int \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{\sigma^2}\right\} dx.
 \end{aligned}$$

Let $y = \sqrt{2}x + \mu(1 - \sqrt{2})$, then

$$\begin{aligned}
 \int \phi(x)^2 dx &= \frac{1}{2\sigma\sqrt{\pi}} \int \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(y-\mu)^2}{2\sigma^2}\right\} dy \\
 &= \frac{1}{2\sigma\sqrt{\pi}} \Phi\left(\sqrt{2}x + \mu(1 - \sqrt{2})\right).
 \end{aligned}$$

Also, using integration by parts,

$$\int \phi(x)\Phi(x) dx = \Phi(x)^2 - \int \phi(x)\Phi(x) dx + C$$

which implies that

$$\int \phi(x)\Phi(x) dx = \frac{1}{2}\Phi(x)^2 + C.$$

Thus

$$\begin{aligned}
\int \Phi(x)^2 dx &= x\Phi(x)^2 + 2\sigma^2\phi(x)\Phi(x) - 2\mu\Phi(x)^2 - \frac{\sigma}{\sqrt{\pi}}\Phi\left(\sqrt{2}x - \mu(1 - \sqrt{2})\right) \\
&\quad + \mu\Phi(x)^2 + C \\
&= (x - \mu)\Phi(x)^2 + 2\sigma^2\phi(x)\Phi(x) - \frac{\sigma}{\sqrt{\pi}}\Phi\left(\sqrt{2}x - \mu(1 - \sqrt{2})\right) + C.
\end{aligned}$$

□

Proposition 2 Consider two samples of data x_1, \dots, x_n and y_1, \dots, y_m and let $F_n(t)$ and $G_m(t)$ represent the empirical cumulative distribution functions for the two samples respectively. Let

$$T_{n,m} = \frac{nm}{n+m} \int_{-\infty}^{\infty} (F_n(t) - G_m(t))^2 dt.$$

Then for $n = m = 1$,

$$T_{1,1} = \frac{|x_1 - y_1|}{2}.$$

Proof If $n = m = 1$, then we have a single observation in each sample, namely x_1 and y_1 . Firstly consider $x_1 \leq y_1$. Then

$$\begin{aligned}
T_{1,1} &= \frac{1}{2} \int_{-\infty}^{\infty} (F_1(t) - G_1(t))^2 dt \\
&= \frac{1}{2} \left(\int_{-\infty}^{x_1} (F_1(t) - G_1(t))^2 dt + \int_{x_1}^{y_1} (F_1(t) - G_1(t))^2 dt + \int_{y_1}^{\infty} (F_1(t) - G_1(t))^2 dt \right) \\
&= \frac{1}{2} \left(0 + \int_{x_1}^{y_1} 1 dt + 0 \right) \\
&= \frac{y_1 - x_1}{2}.
\end{aligned}$$

Now let $x_1 > y_1$, then

$$\begin{aligned}
T_{1,1} &= \frac{1}{2} \int_{-\infty}^{\infty} (F_1(t) - G_1(t))^2 dt \\
&= \frac{1}{2} \left(\int_{-\infty}^{y_1} (F_1(t) - G_1(t))^2 dt + \int_{y_1}^{x_1} (F_1(t) - G_1(t))^2 dt + \int_{x_1}^{\infty} (F_1(t) - G_1(t))^2 dt \right) \\
&= \frac{1}{2} \left(0 + \int_{y_1}^{x_1} 1 dt + 0 \right) \\
&= \frac{x_1 - y_1}{2}.
\end{aligned}$$

C. Miscellaneous Propositions

Hence, for $n = m = 1$,

$$T_{1,1} = \frac{|x_1 - y_1|}{2}.$$

□

Proposition 3 Let $\mathcal{N}_\lambda(t - \lambda)$ be the cumulative distribution function of a two-parameter Generalised Pareto distribution with parameters σ and ξ . Then for $t \geq \lambda$, $\mathcal{N}(t) = \mathcal{N}(\lambda) + \mathcal{N}_\lambda(t - \lambda)(1 - \mathcal{N}(\lambda))$, is a three-parameter Generalised Pareto distribution with parameters

$$\begin{aligned}\tilde{\xi} &= \xi \\ \tilde{\sigma} &= \sigma(1 - \mathcal{N}(\lambda))^{\tilde{\xi}} \\ \tilde{\lambda} &= \lambda - \frac{\tilde{\sigma}}{\tilde{\xi}} \left((1 - \mathcal{N}(\lambda))^{-\tilde{\xi}} - 1 \right).\end{aligned}$$

Proof It is known that

$$\mathcal{N}_\lambda(t - \lambda) = 1 - \left(1 + \frac{\xi(t - \lambda)}{\sigma} \right)^{-\frac{1}{\xi}}.$$

Thus for $t \geq \lambda$,

$$\begin{aligned}\mathcal{N}(t) &= \mathcal{N}(\lambda) + \mathcal{N}_\lambda(t - \lambda)(1 - \mathcal{N}(\lambda)) \\ &= \mathcal{N}(\lambda) + (1 - \mathcal{N}(\lambda)) \left[1 - \left(1 + \frac{\xi(t - \lambda)}{\sigma} \right)^{-\frac{1}{\xi}} \right] \\ &= 1 - (1 - \mathcal{N}(\lambda)) \left(1 + \frac{\xi(t - \lambda)}{\sigma} \right)^{-\frac{1}{\xi}} \\ &= 1 - \left(\frac{1}{(1 - \mathcal{N}(\lambda))^\xi} + \frac{\xi(t - \lambda)}{\sigma(1 - \mathcal{N}(\lambda))^\xi} \right)^{-\frac{1}{\xi}} \\ &= 1 - \left(1 + \frac{\xi(t - \lambda) - \sigma((1 - \mathcal{N}(\lambda))^\xi - 1)}{\sigma(1 - \mathcal{N}(\lambda))^\xi} \right)^{-\frac{1}{\xi}} \\ &= 1 - \left(1 + \frac{\xi \left(t - \left(\lambda + \frac{\sigma}{\xi} ((1 - \mathcal{N}(\lambda))^\xi - 1) \right) \right)}{\sigma(1 - \mathcal{N}(\lambda))^\xi} \right)^{-\frac{1}{\xi}} \\ &= 1 - \left(1 + \frac{\tilde{\xi}(t - \tilde{\lambda})}{\tilde{\sigma}} \right)^{-\frac{1}{\xi}}\end{aligned}$$

where

$$\begin{aligned}\tilde{\xi} &= \xi \\ \tilde{\sigma} &= \sigma(1 - \mathcal{N}(\lambda))^{\tilde{\xi}} \\ \tilde{\lambda} &= \lambda - \frac{\tilde{\sigma}}{\tilde{\xi}} \left((1 - \mathcal{N}(\lambda))^{-\tilde{\xi}} - 1 \right).\end{aligned}$$

□

Appendix D

Variance of $T_{n,m}$ Proof

In this appendix we provide the proof for calculating the variance of the test statistic, $\text{Var}[T_{n,m}]$.

The variance of the test statistic $T_{n,m}$ where $n > m$ follows the following formula;

$$\begin{aligned} \text{Var}[T_{n,m}|F = G] = \frac{2}{\mathcal{V}} \int_{-\infty}^{\infty} \int_{-\infty}^t H(s) & \left(1 + 2(\mathcal{V} - 2)H(s) - 3H(t) \right. \\ & \left. - 2(2\mathcal{V} - 5)H(s)H(t) + 2H(t)^2 + 2(\mathcal{V} - 3)H(s)H(t)^2 \right) ds dt \end{aligned}$$

where $\mathcal{V} = \frac{nm(n+m)^2}{n^3+m^3}$.

Proof Consider the test statistic

$$T_{n,m} = \frac{nm}{n+m} \int_{-\infty}^{\infty} (\hat{F}(t) - \hat{G}(t))^2 dt \quad (\text{D.1})$$

where

$$\hat{F}(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_i \leq t\}}$$

and

$$\hat{G}(t) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{y_i \leq t\}}$$

are the empirical cumulative distribution functions for the distributions f and g , respectively. From the standard theory,

$$\text{Var}[T_{n,m}] = \text{E}[T_{n,m}^2] - \{\text{E}[T_{n,m}]\}^2. \quad (\text{D.2})$$

D. Variance of $T_{n,m}$ Proof

We have already identified that

$$\mathbb{E}[T_{n,m}] = \int_{-\infty}^{\infty} (H(t) - H(t)^2) dt,$$

thus we require the expression for $\mathbb{E}[T_{n,m}^2]$. Now by squaring the test statistic and using Fubini's theorem (Fubini, 1907), we get

$$\begin{aligned} \mathbb{E}[T_{n,m}^2] &= \frac{n^2 m^2}{(n+m)^2} \mathbb{E} \left[\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\hat{F}(t) - \hat{G}(t))^2 (\hat{F}(s) - \hat{G}(s))^2 ds dt \right] \\ &= \frac{n^2 m^2}{(n+m)^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mathbb{E} \left[(\hat{F}(t) - \hat{G}(t))^2 (\hat{F}(s) - \hat{G}(s))^2 \right] ds dt \quad (\text{D.3}) \end{aligned}$$

Substituting the expressions for $\hat{F}(t)$ and $\hat{G}(t)$ into equation (D.3) and expanding gives

$$\begin{aligned} \mathbb{E}[T_{n,m}^2] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mathbb{E} \left[\frac{1}{n^4} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}} \mathbb{1}_{\{x_l \leq s\}} \right. \\ &\quad - \frac{2}{n^3 m} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^m \mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}} \mathbb{1}_{\{y_l \leq s\}} \\ &\quad + \frac{1}{n^2 m^2} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^m \sum_{l=1}^m \mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{y_k \leq s\}} \mathbb{1}_{\{y_l \leq s\}} \\ &\quad - \frac{2}{n^3 m} \sum_{i=1}^n \sum_{j=1}^m \sum_{k=1}^n \sum_{l=1}^n \mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{y_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}} \mathbb{1}_{\{x_l \leq s\}} \\ &\quad + \frac{4}{n^2 m^2} \sum_{i=1}^n \sum_{j=1}^m \sum_{k=1}^n \sum_{l=1}^m \mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{y_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}} \mathbb{1}_{\{y_l \leq s\}} \\ &\quad - \frac{2}{nm^3} \sum_{i=1}^n \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{y_j \leq t\}} \mathbb{1}_{\{y_k \leq s\}} \mathbb{1}_{\{y_l \leq s\}} \\ &\quad + \frac{1}{n^2 m^2} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^n \sum_{l=1}^n \mathbb{1}_{\{y_i \leq t\}} \mathbb{1}_{\{y_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}} \mathbb{1}_{\{x_l \leq s\}} \\ &\quad - \frac{2}{nm^3} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^n \sum_{l=1}^m \mathbb{1}_{\{y_i \leq t\}} \mathbb{1}_{\{y_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}} \mathbb{1}_{\{y_l \leq s\}} \\ &\quad \left. + \frac{1}{m^4} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \mathbb{1}_{\{y_i \leq t\}} \mathbb{1}_{\{y_j \leq t\}} \mathbb{1}_{\{y_k \leq s\}} \mathbb{1}_{\{y_l \leq s\}} \right] ds dt. \quad (\text{D.4}) \end{aligned}$$

The expectation of each summation in equation (D.4) can be calculated separately by firstly considering all the cases for i, j, k, l . For example, Table D.1 displays the different cases for i, j, k, l and the number of times each case occurs

for the first summation in equation (D.4). Each summation can then be split into each case for i, j, k, l and using the properties that under the null hypothesis $F \equiv G \equiv H$

$$\mathbb{E} [\mathbb{1}_{Z_i \leq \mathcal{A}}] = \Pr\{Z_i \leq \mathcal{A}\} = H(\mathcal{A}) \quad (\text{D.5})$$

where $Z = \{X_1, \dots, X_n, Y_1, \dots, Y_m\}$, each expectation can be calculated.

Case	Number of Occurrences
$i = j = k = l$	n
$i = j = k, i \neq l$	$2n(n - 1)$
$i = k = l, i \neq j$	$2n(n - 1)$
$i = j, k = l, i \neq k$	$n(n - 1)$
$i = k, j = l, i \neq j$	$2n(n - 1)$
$i = j, i \neq k \neq l$	$n(n - 1)(n - 2)$
$k = l, i \neq j \neq k$	$n(n - 1)(n - 2)$
$i = k, i \neq j \neq l$	$4n(n - 1)(n - 2)$
$i \neq j \neq k \neq l$	$n(n - 1)(n - 2)(n - 3)$

Table D.1: The number of occurrences for each case of i, j, k, l for the first summation in equation (D.4). Here i and j are interchangeable, similarly, k and l are interchangeable.

Using Table D.1, the first summation in equation (D.4) can be expanded;

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{n^4} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}} \mathbb{1}_{\{x_l \leq s\}}] \, ds \, dt = \\ & \frac{1}{n^4} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left\{ n \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}^2 \mathbb{1}_{\{x_i \leq s\}}^2] + 2n(n - 1) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}^2 \mathbb{1}_{\{x_i \leq s\}} \mathbb{1}_{\{x_j \leq s\}}] \right. \\ & + 2n(n - 1) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{x_i \leq s\}}^2] + n(n - 1) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}^2 \mathbb{1}_{\{x_j \leq s\}}^2] \\ & + 2n(n - 1) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_i \leq s\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{x_i \leq s\}}] \\ & + n(n - 1)(n - 2) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}^2 \mathbb{1}_{\{x_j \leq s\}} \mathbb{1}_{\{x_k \leq s\}}] \\ & + n(n - 1)(n - 2) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}}^2] \\ & + 4n(n - 1)(n - 2) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_i \leq s\}} \mathbb{1}_{\{x_j \leq s\}} \mathbb{1}_{\{x_k \leq t\}}] \\ & \left. + n(n - 1)(n - 2)(n - 3) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}} \mathbb{1}_{\{x_l \leq s\}}] \right\} \, ds \, dt. \quad (\text{D.6}) \end{aligned}$$

By using the property in equation (D.5), under the null hypothesis, equation (D.6)

D. Variance of $T_{n,m}$ Proof

is equivalent to

$$\begin{aligned} & \frac{1}{n^4} \int_{-\infty}^{\infty} \left\{ \int_{-\infty}^t \left\{ nH(s) + 4n(n-1)H(s)^2 + 3n(n-1)H(s)H(t) \right. \right. \\ & \quad + 5n(n-1)(n-2)H(s)^2H(t) + n(n-1)(n-2)H(s)H(t)^2 \\ & \quad \left. \left. + n(n-1)(n-2)(n-3)H(s)^2H(t)^2 \right\} ds \right. \\ & + \int_t^{\infty} \left\{ nH(t) + 4n(n-1)H(t)^2 + 3n(n-1)H(s)H(t) \right. \\ & \quad + n(n-1)(n-2)H(s)^2H(t) + 5n(n-1)(n-2)H(s)H(t)^2 \\ & \quad \left. \left. + n(n-1)(n-2)(n-3)H(s)^2H(t)^2 \right\} ds \right\} dt. \end{aligned}$$

Note that care needs to be taken for other summations where data X and Y are involved. In the case of the second summation in equation (D.4), Table D.2 displays the different cases for i, j, k, l and the number of times each case occurs. Recall that $n > m$.

Case	Number of Occurrences
$i = j = k = l$	m
$i = j = k, i \neq l$	$m(n-1)$
$i = j = l, i \neq k$	$m(n-1)$
$i = k = l, i \neq j$	$2m(n-1)$
$i = j, k = l, i \neq k$	$m(n-1)$
$i = k, j = l, i \neq j$	$2m(n-1)$
$i = j, i \neq k \neq l$	$m(n-1)(n-2)$
$i = k, i \neq j \neq l$	$2m(n-1)(n-2)$
$i = l, i \neq j \neq k$	$3m(n-1)(n-2)$
$i \neq j \neq k \neq l$	$m(n-1)(n-2)(n-3)$

Table D.2: The number of occurrences for each case of i, j, k, l for the second summation in equation (D.4). Here only i and j are interchangeable.

Using Table D.2, the second summation in equation (D.4) can be expanded;

$$\begin{aligned}
& \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} -\frac{2}{n^3 m} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^m \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}} \mathbb{1}_{\{y_l \leq s\}}] \, ds \, dt = \\
& -\frac{2}{n^3 m} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left\{ m \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}^2 \mathbb{1}_{\{x_i \leq s\}} \mathbb{1}_{\{y_i \leq s\}}] \right. \\
& + m(n-1) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}^2 \mathbb{1}_{\{x_i \leq s\}} \mathbb{1}_{\{y_j \leq s\}}] \\
& + m(n-1) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}^2 \mathbb{1}_{\{y_i \leq s\}} \mathbb{1}_{\{x_j \leq s\}}] \\
& + 2m(n-1) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_i \leq s\}} \mathbb{1}_{\{y_i \leq s\}} \mathbb{1}_{\{x_j \leq t\}}] \\
& + m(n-1) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}^2 \mathbb{1}_{\{x_j \leq s\}} \mathbb{1}_{\{y_j \leq s\}}] \\
& + 2m(n-1) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_i \leq s\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{y_j \leq s\}}] \\
& + m(n-1)(n-2) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}}^2 \mathbb{1}_{\{x_j \leq s\}} \mathbb{1}_{\{y_k \leq s\}}] \\
& + 2m(n-1)(n-2) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_i \leq s\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{y_k \leq s\}}] \\
& + 3m(n-1)(n-2) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{y_i \leq s\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}}] \\
& \left. + m(n-1)(n-2)(n-3) \mathbb{E} [\mathbb{1}_{\{x_i \leq t\}} \mathbb{1}_{\{x_j \leq t\}} \mathbb{1}_{\{x_k \leq s\}} \mathbb{1}_{\{y_l \leq s\}}] \right\} \, ds \, dt. \quad (\text{D.7})
\end{aligned}$$

By using the property in equation (D.5), under the null hypothesis, equation (D.7) is equivalent to

$$\begin{aligned}
& -\frac{2}{n^3 m} \int_{-\infty}^{\infty} \left\{ \int_{-\infty}^t \left\{ mnH(s)^2 + 6m(n-1)H(s)^2H(t) \right. \right. \\
& \quad + 3m(n-1)(n-2)H(s)^2H(t) + 3m(n-1)(n-2)H(s)^2H(t)^2 \\
& \quad \left. \left. + m(n-1)(n-2)(n-3)H(s)^2H(t)^2 \right\} \, ds \right. \\
& \quad + \int_t^{\infty} \left\{ nmH(t)H(s) + 2m(n-1)H(s)^2H(t) \right. \\
& \quad + m(n-1)(n-2)H(s)^2H(t) + 4m(n-1)H(s)H(t)^2 \\
& \quad + 2m(n-1)(n-2)H(t)^2H(s) + 3m(n-1)(n-2)H(s)^2H(t)^2 \\
& \quad \left. \left. + m(n-1)(n-2)(n-3)H(s)^2H(t)^2 \right\} \, ds \right\} \, dt.
\end{aligned}$$

Once we repeat this process for each summation in equation (D.4) and collect

D. Variance of $T_{n,m}$ Proof

together all the terms, we get

$$\begin{aligned} \mathbb{E}[T_{n,m}^2] = & \int_{-\infty}^{\infty} \left\{ \int_{-\infty}^t \left\{ \frac{1}{\mathcal{V}} H(s) + 2 \left(1 - \frac{2}{\mathcal{V}}\right) H(s)^2 + \left(1 - \frac{3}{\mathcal{V}}\right) H(s)H(t) \right. \right. \\ & - 5 \left(1 - \frac{2}{\mathcal{V}}\right) H(s)^2 H(t) - \left(1 - \frac{2}{\mathcal{V}}\right) H(s)H(t)^2 \\ & \left. \left. + 3 \left(1 - \frac{2}{\mathcal{V}}\right) H(s)^2 H(t)^2 \right\} ds \right. \end{aligned} \quad (\text{D.8})$$

$$\begin{aligned} & \left. + \int_t^{\infty} \left\{ \frac{1}{\mathcal{V}} H(t) + 2 \left(1 - \frac{2}{\mathcal{V}}\right) H(t)^2 + \left(1 - \frac{3}{\mathcal{V}}\right) H(s)H(t) \right. \right. \\ & - 5 \left(1 - \frac{2}{\mathcal{V}}\right) H(s)H(t)^2 - \left(1 - \frac{2}{\mathcal{V}}\right) H(s)^2 H(t) \\ & \left. \left. + 3 \left(1 - \frac{2}{\mathcal{V}}\right) H(s)^2 H(t)^2 \right\} ds \right\} dt. \end{aligned} \quad (\text{D.9})$$

Note that for if the order of integration is swapped and the variables s and t permuted in the second integral of equation (D.9), it becomes equivalent to the first integral. Thus

$$\begin{aligned} \mathbb{E}[T_{n,m}^2] = & \int_{-\infty}^{\infty} \left\{ \int_{-\infty}^t \left\{ \frac{2}{\mathcal{V}} H(s) + 4 \left(1 - \frac{2}{\mathcal{V}}\right) H(s)^2 + 2 \left(1 - \frac{3}{\mathcal{V}}\right) H(s)H(t) \right. \right. \\ & - 10 \left(1 - \frac{2}{\mathcal{V}}\right) H(s)^2 H(t) - 2 \left(1 - \frac{2}{\mathcal{V}}\right) H(s)H(t)^2 \\ & \left. \left. + 6 \left(1 - \frac{2}{\mathcal{V}}\right) H(s)^2 H(t)^2 \right\} ds \right\} dt \end{aligned} \quad (\text{D.10})$$

Now, as

$$\mathbb{E}[T_{n,m}] = \int_{-\infty}^{\infty} (H(t) - H(t)^2) dt,$$

we find that

$$\begin{aligned} \{\mathbb{E}[T_{n,m}]\}^2 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} H(t)(1 - H(t))H(s)(1 - H(s)) ds dt \\ &= 2 \int_{-\infty}^{\infty} \int_{-\infty}^t H(s)H(t)(1 - H(s))(1 - H(t)) ds dt \\ &= 2 \int_{-\infty}^{\infty} \int_{-\infty}^t (H(s)H(t) - H(s)H(t)^2 - H(s)^2 H(t) \\ & \quad + H(s)^2 H(t)^2) ds dt. \end{aligned} \quad (\text{D.11})$$

Finally, using equation (D.2) we obtain

$$\begin{aligned} \text{Var}[T_{n,m}|F = G] &= \int_{-\infty}^{\infty} \int_{-\infty}^t \frac{2}{\mathcal{V}} H(s) + 4 \left(1 - \frac{2}{\mathcal{V}}\right) H(s)^2 - \frac{6}{\mathcal{V}} H(s)H(t) \\ &\quad - 4 \left(2 - \frac{5}{\mathcal{V}}\right) H(s)^2 H(t) + \frac{4}{\mathcal{V}} H(s)H(t)^2 + 4 \left(1 - \frac{3}{\mathcal{V}}\right) H(s)^2 H(t)^2 \, ds \, dt, \end{aligned}$$

which can be factorised to give

$$\begin{aligned} \text{Var}[T_{n,m}|F = G] &= \frac{2}{\mathcal{V}} \int_{-\infty}^{\infty} \int_{-\infty}^t H(s) \left(1 + 2(\mathcal{V} - 2) H(s) - 3H(t) \right. \\ &\quad \left. - 2(2\mathcal{V} - 5) H(s)H(t) + 2H(t)^2 + 2(\mathcal{V} - 3) H(s)H(t)^2 \right) \, ds \, dt \end{aligned}$$

where $\mathcal{V} = \frac{nm(n+m)^2}{n^3+m^3}$.

□

Appendix E

Third Moment of $T_{n,m}$ Proof

In this appendix we provide the proof for calculating the third moment of the test statistic, $E[T_{n,m}^3]$. Note that this proof is very similar to the proof for $\text{Var}[T_{n,m}]$ in Appendix D and we will refer to it throughout.

The third non-centralised moment of $T_{n,m}$ for $n > m$ is given by

$$\begin{aligned}
 E[T_{n,m}^3|F = G] = & \frac{6}{\mathcal{G}nm(n+m)} \int_{-\infty}^{\infty} \int_{-\infty}^s \int_{-\infty}^t H(t) \left(1 \right. \\
 & + 2 \left(\mathcal{G}(7(m^2 + n^2) - 10nm) - 8 \right) H(t) + 2 \left(\mathcal{G}(5(m^2 + n^2) - 7nm) - 6 \right) H(s) \\
 & + \left(\mathcal{G}(m^2 + n^2 - nm) - 3 \right) H(r) \\
 & + 5 \left(\mathcal{G}(2mn(m+n) - 19(m^2 + n^2) + 25mn) + 18 \right) H(s)H(t) \\
 & + \left(\mathcal{G}(2mn(m+n) - 19(m^2 + n^2) + 25mn) + 18 \right) H(s)^2 \\
 & - \left(\mathcal{G}(m^2 + n^2 - nm) - 2 \right) H(r)^2 \\
 & + 2 \left(\mathcal{G}(mn(m+n) - 19(m^2 + n^2) + 26mn) + 20 \right) H(r)H(t) \\
 & + \left(\mathcal{G}(mn(m+n) - 27(m^2 + n^2) + 37mn) + 30 \right) H(r)H(s) \\
 & - 4 \left(\mathcal{G}(3mn(m+n) - 26(m^2 + n^2) + 34mn) + 24 \right) H(s)^2H(t) \\
 & - 5 \left(\mathcal{G}(5mn(m+n) - 45(m^2 + n^2) + 59mn) + 42 \right) H(r)H(s)H(t) \\
 & - 2 \left(\mathcal{G}(mn(m+n) - 12(m^2 + n^2) + 16mn) + 12 \right) H(r)^2H(t) \\
 & - \left(\mathcal{G}(5mn(m+n) - 45(m^2 + n^2) + 59mn) + 42 \right) H(r)H(s)^2 \\
 & - \left(\mathcal{G}(mn(m+n) - 17(m^2 + n^2) + 23mn) + 18 \right) H(r)^2H(s) \\
 & + 9 \left(\mathcal{G}(3mn(m+n) - 26(m^2 + n^2) + 34mn) + 24 \right) H(r)H(s)^2H(t) \\
 & + 5 \left(\mathcal{G}(3mn(m+n) - 26(m^2 + n^2) + 34mn) + 24 \right) H(r)^2H(s)H(t) \\
 & + \left(\mathcal{G}(3mn(m+n) - 26(m^2 + n^2) + 34mn) + 24 \right) H(r)^2H(s)^2 \\
 & \left. - 5 \left(\mathcal{G}(3mn(m+n) - 26(m^2 + n^2) + 34mn) + 24 \right) H(r)^2H(s)^2H(t) \right) \\
 & dt ds dr
 \end{aligned} \tag{E.1}$$

E. Third Moment of $T_{n,m}$ Proof

where $\mathcal{G} = \frac{nm(n+m)^2}{n^5+m^5}$.

Proof Consider the test statistic

$$T_{n,m} = \frac{nm}{n+m} \int_{-\infty}^{\infty} (\hat{F}(t) - \hat{G}(t))^2 dt.$$

Now by cubing the test statistic and using Fubini's theorem (Fubini, 1907), we obtain

$$\begin{aligned} \mathbb{E}[T_{n,m}^3] &= \mathbb{E} \left[\frac{n^3 m^3}{(n+m)^3} \iiint_{\mathbb{R}^3} (\hat{F}(t) - \hat{G}(t))^2 (\hat{F}(s) - \hat{G}(s))^2 ds dt \right] \\ &= \frac{n^3 m^3}{(n+m)^3} \iiint_{\mathbb{R}^3} \mathbb{E} \left[(\hat{F}(t) - \hat{G}(t))^2 (\hat{F}(s) - \hat{G}(s))^2 \right] ds dt \quad (\text{E.2}) \end{aligned}$$

Substituting the expressions for $\hat{F}(t)$ and $\hat{G}(t)$ into equation (E.2) and expanding gives

$$\begin{aligned} \mathbb{E}[T_{n,m}^3] &= \frac{n^3 m^3}{(n+m)^3} \iiint_{\mathbb{R}^3} \mathbb{E} \left[\left(\hat{F}(s) - \hat{G}(s) \right)^2 \left(\hat{F}(t) - \hat{G}(t) \right)^2 \right. \\ &\quad \left. \cdot \left(\hat{F}(u) - \hat{G}(u) \right)^2 \right] dt ds dr \\ &= \frac{n^3 m^3}{(n+m)^3} \iiint_{\mathbb{R}^3} \mathbb{E} \left[\frac{1}{n^6} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^n \sum_{p=1}^n (\mathbb{1}_{\{x_i \geq r\}} \mathbb{1}_{\{x_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \right. \\ &\quad \left. \cdot \mathbb{1}_{\{x_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{x_p \geq t\}} \right) \\ &\quad - \frac{2}{n^5 m} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^n \sum_{p=1}^m (\mathbb{1}_{\{x_i \geq r\}} \mathbb{1}_{\{x_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \\ &\quad \left. \cdot \mathbb{1}_{\{x_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \right) \\ &\quad + \frac{1}{n^4 m^2} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^m \sum_{o=1}^m \sum_{p=1}^m (\mathbb{1}_{\{x_i \geq r\}} \mathbb{1}_{\{x_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \\ &\quad \left. \cdot \mathbb{1}_{\{x_l \geq s\}} \mathbb{1}_{\{y_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \right) \\ &\quad - \frac{2}{n^5 m} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^m \sum_{l=1}^n \sum_{o=1}^n \sum_{p=1}^n (\mathbb{1}_{\{x_i \geq r\}} \mathbb{1}_{\{x_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \\ &\quad \left. \cdot \mathbb{1}_{\{y_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{x_p \geq t\}} \right) \\ &\quad + \frac{4}{n^4 m^2} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^m \sum_{p=1}^m (\mathbb{1}_{\{x_i \geq r\}} \mathbb{1}_{\{x_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \\ &\quad \left. \cdot \mathbb{1}_{\{y_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \right) \end{aligned}$$

E. Third Moment of $T_{n,m}$ Proof

$$\begin{aligned}
& + \frac{4}{n^2 m^4} \sum_{i=1}^n \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^n \sum_{p=1}^m \left(\mathbb{1}_{\{x_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{y_k \geq s\}} \right. \\
& \quad \cdot \mathbb{1}_{\{y_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \left. \right) \\
& - \frac{2}{n m^5} \sum_{i=1}^n \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^m \sum_{p=1}^m \left(\mathbb{1}_{\{x_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{y_k \geq s\}} \right. \\
& \quad \cdot \mathbb{1}_{\{y_l \geq s\}} \mathbb{1}_{\{y_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \left. \right) \\
& + \frac{1}{n^4 m^2} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^n \sum_{p=1}^n \left(\mathbb{1}_{\{y_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \right. \\
& \quad \cdot \mathbb{1}_{\{x_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{x_p \geq t\}} \left. \right) \\
& - \frac{2}{n^3 m^3} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^n \sum_{p=1}^m \left(\mathbb{1}_{\{y_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \right. \\
& \quad \cdot \mathbb{1}_{\{x_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \left. \right) \\
& + \frac{1}{n^2 m^4} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^m \sum_{p=1}^m \left(\mathbb{1}_{\{y_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \right. \\
& \quad \cdot \mathbb{1}_{\{x_l \geq s\}} \mathbb{1}_{\{y_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \left. \right) \\
& - \frac{2}{n^3 m^3} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^n \sum_{l=1}^m \sum_{o=1}^n \sum_{p=1}^n \left(\mathbb{1}_{\{y_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \right. \\
& \quad \cdot \mathbb{1}_{\{y_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{x_p \geq t\}} \left. \right) \\
& + \frac{4}{n^2 m^4} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^n \sum_{l=1}^m \sum_{o=1}^n \sum_{p=1}^m \left(\mathbb{1}_{\{y_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \right. \\
& \quad \cdot \mathbb{1}_{\{y_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \left. \right) \\
& - \frac{2}{n m^5} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^n \sum_{l=1}^m \sum_{o=1}^m \sum_{p=1}^m \left(\mathbb{1}_{\{y_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \right. \\
& \quad \cdot \mathbb{1}_{\{y_l \geq s\}} \mathbb{1}_{\{y_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \left. \right) \\
& + \frac{1}{n^2 m^4} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^n \sum_{p=1}^n \left(\mathbb{1}_{\{y_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{y_k \geq s\}} \right. \\
& \quad \cdot \mathbb{1}_{\{y_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{x_p \geq t\}} \left. \right) \\
& - \frac{2}{n m^5} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^n \sum_{p=1}^m \left(\mathbb{1}_{\{y_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{y_k \geq s\}} \right. \\
& \quad \cdot \mathbb{1}_{\{y_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \left. \right)
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{m^6} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^m \sum_{p=1}^m \left(\mathbb{1}_{\{y_i \geq r\}} \mathbb{1}_{\{y_j \geq r\}} \mathbb{1}_{\{y_k \geq s\}} \right. \\
& \quad \left. \cdot \mathbb{1}_{\{y_l \geq s\}} \mathbb{1}_{\{y_o \geq t\}} \mathbb{1}_{\{y_p \geq t\}} \right) \Bigg]. \tag{E.3}
\end{aligned}$$

It is important to note that we can split the triple integral into the sum of six triple integrals for the following six cases;

1. $s \geq t \geq u$,
2. $s \geq u \geq t$,
3. $t \geq s \geq u$,
4. $t \geq u \geq s$,
5. $u \geq s \geq t$,
6. $u \geq t \geq s$.

We also know that due to symmetry, all six triple integrals will be equal. Thus we will only consider the first case and multiply the result by six.

The expectation of each summation in equation (E.3) can be calculated separately by firstly considering all the cases for i, j, k, l, o, p (see for example Table D.1 and Table D.2). Each summation can then be split into each case for i, j, k, l, o, p and using the properties that under the null hypothesis $F \equiv G \equiv H$

$$\mathbb{E} [\mathbb{1}_{Z_i \leq \mathcal{A}}] = \Pr\{Z_i \leq \mathcal{A}\} = H(\mathcal{A}) \tag{E.4}$$

where $Z = \{X_1, \dots, X_n, Y_1, \dots, Y_m\}$, each expectation can be calculated.

For example, the first summation in equation (E.3) can be expanded;

$$\begin{aligned}
& \frac{1}{n^6} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^n \sum_{p=1}^n \mathbb{E} \left[\mathbb{1}_{\{x_i \geq r\}} \mathbb{1}_{\{x_j \geq r\}} \mathbb{1}_{\{x_k \geq s\}} \mathbb{1}_{\{x_l \geq s\}} \mathbb{1}_{\{x_o \geq t\}} \mathbb{1}_{\{x_p \geq t\}} \right] \\
& = \frac{1}{n^6} \left\{ n \mathbb{E} \left[\mathbb{1}_{\{x_i \geq r\}}^2 \mathbb{1}_{\{x_i \geq s\}}^2 \mathbb{1}_{\{x_i \geq t\}}^2 \right] \right. \\
& \quad + 2n(n-1) \mathbb{E} \left[\mathbb{1}_{\{x_i \geq r\}}^2 \mathbb{1}_{\{x_i \geq s\}}^2 \mathbb{1}_{\{x_i \geq t\}} \mathbb{1}_{\{x_j \geq t\}} \right] \\
& \quad + 2n(n-1) \mathbb{E} \left[\mathbb{1}_{\{x_i \geq r\}}^2 \mathbb{1}_{\{x_i \geq s\}} \mathbb{1}_{\{x_j \geq s\}} \mathbb{1}_{\{x_i \geq t\}}^2 \right] \\
& \quad + 2n(n-1) \mathbb{E} \left[\mathbb{1}_{\{x_i \geq r\}} \mathbb{1}_{\{x_j \geq r\}} \mathbb{1}_{\{x_i \geq s\}}^2 \mathbb{1}_{\{x_i \geq t\}}^2 \right] \\
& \quad + n(n-1) \mathbb{E} \left[\mathbb{1}_{\{x_i \geq r\}}^2 \mathbb{1}_{\{x_i \geq s\}}^2 \mathbb{1}_{\{x_j \geq t\}}^2 \right] \\
& \quad + n(n-1) \mathbb{E} \left[\mathbb{1}_{\{x_i \geq r\}}^2 \mathbb{1}_{\{x_j \geq s\}}^2 \mathbb{1}_{\{x_i \geq t\}}^2 \right] \\
& \quad \left. + n(n-1) \mathbb{E} \left[\mathbb{1}_{\{x_j \geq r\}}^2 \mathbb{1}_{\{x_i \geq s\}}^2 \mathbb{1}_{\{x_i \geq t\}}^2 \right] \right\}
\end{aligned}$$

E. Third Moment of $T_{n,m}$ Proof

the first summation in equation (E.3) is equivalent to

$$\begin{aligned}
& \frac{1}{n^6} \int_{-\infty}^{\infty} \int_{-\infty}^s \int_{-\infty}^t nH(u) + 16n(n-1)H(u)^2 + 12n(n-1)H(t)H(u) \\
& + 3n(n-1)H(s)H(u) + 45n(n-1)(n-2)H(t)H(u)^2 \\
& + 9n(n-1)(n-2)H(t)^2H(u) + n(n-1)(n-2)H(s)^2H(u) \\
& + 20n(n-1)(n-2)H(s)H(u)^2 + 15n(n-1)(n-2)H(s)H(t)H(u) \\
& + 16n(n-1)(n-2)(n-3)H(t)^2H(u)^2 \\
& + 35n(n-1)(n-2)(n-3)H(s)H(t)H(u)^2 \\
& + 4n(n-1)(n-2)(n-3)H(s)^2H(u)^2 \\
& + 7n(n-1)(n-2)(n-3)H(s)H(t)^2H(u) \\
& + 3n(n-1)(n-2)(n-3)H(s)^2H(t)H(u) \\
& + 9n(n-1)(n-2)(n-3)(n-4)H(s)H(t)^2H(u)^2 \\
& + 5n(n-1)(n-2)(n-3)(n-4)H(s)^2H(t)H(u)^2 \\
& + n(n-1)(n-2)(n-3)(n-4)H(s)^2H(t)^2H(u) \\
& + n(n-1)(n-2)(n-3)(n-4)(n-5)H(s)^2H(t)^2H(u)^2 dt ds dr.
\end{aligned}$$

As for the proof of $\text{Var}[T_{n,m}]$, this process can be repeated for each summation in equation (E.3) with care being taken where both data X and Y are involved. Once we repeat this process for each summation in equation (E.3), multiply by 6 and collect together all the terms, we get

$$\begin{aligned}
\mathbb{E}[T_{n,m}^3 | F = G] &= \frac{6}{\mathcal{G}nm(n+m)} \int_{-\infty}^{\infty} \int_{-\infty}^s \int_{-\infty}^t H(t) \left(1 \right. \\
& + 2 (\mathcal{G}(7(m^2 + n^2) - 10nm) - 8) H(t) + 2 (\mathcal{G}(5(m^2 + n^2) - 7nm) - 6) H(s) \\
& + (\mathcal{G}(m^2 + n^2 - nm) - 3) H(r) \\
& + 5 (\mathcal{G}(2mn(m+n) - 19(m^2 + n^2) + 25mn) + 18) H(s)H(t) \\
& + (\mathcal{G}(2mn(m+n) - 19(m^2 + n^2) + 25mn) + 18) H(s)^2 \\
& - (\mathcal{G}(m^2 + n^2 - nm) - 2) H(r)^2 \\
& + 2 (\mathcal{G}(mn(m+n) - 19(m^2 + n^2) + 26mn) + 20) H(r)H(t) \\
& + (\mathcal{G}(mn(m+n) - 27(m^2 + n^2) + 37mn) + 30) H(r)H(s) \\
& - 4 (\mathcal{G}(3mn(m+n) - 26(m^2 + n^2) + 34mn) + 24) H(s)^2H(t) \\
& - 5 (\mathcal{G}(5mn(m+n) - 45(m^2 + n^2) + 59mn) + 42) H(r)H(s)H(t) \\
& - 2 (\mathcal{G}(mn(m+n) - 12(m^2 + n^2) + 16mn) + 12) H(r)^2H(t) \\
& \left. - (\mathcal{G}(5mn(m+n) - 45(m^2 + n^2) + 59mn) + 42) H(r)H(s)^2 \right)
\end{aligned}$$

$$\begin{aligned}
& - (\mathcal{G}(mn(m+n) - 17(m^2 + n^2) + 23mn) + 18) H(r)^2 H(s) \\
& + 9 (\mathcal{G}(3mn(m+n) - 26(m^2 + n^2) + 34mn) + 24) H(r)H(s)^2 H(t) \\
& + 5 (\mathcal{G}(3mn(m+n) - 26(m^2 + n^2) + 34mn) + 24) H(r)^2 H(s)H(t) \\
& + (\mathcal{G}(3mn(m+n) - 26(m^2 + n^2) + 34mn) + 24) H(r)^2 H(s)^2 \\
& - 5 (\mathcal{G}(3mn(m+n) - 26(m^2 + n^2) + 34mn) + 24) H(r)^2 H(s)^2 H(t) \\
& dt ds dr
\end{aligned}$$

where $\mathcal{G} = \frac{nm(n+m)^2}{n^5+m^5}$.

□

Appendix F

Fourth Moment of $T_{n,m}$ Proof

In this appendix we provide the proof for calculating the fourth moment of the test statistic, $E[T_{n,m}^4]$. Note that this proof is very similar to the proof for $\text{Var}[T_{n,m}]$ in Appendix D and will refer to it throughout.

The fourth non-centralised moment of $T_{n,m}$ for $n > m$ is given by

$$\begin{aligned} E[T_{n,m}^4] = & \frac{24}{\mathcal{K}(nm(n+m))^2} \int_{-\infty}^{\infty} \int_{-\infty}^s \int_{-\infty}^t \int_{-\infty}^u H(v) \left(1 \right. \\ & + 2 \left(\mathcal{K}(31(m^4 + n^4) - 56mn(m^2 + n^2) + 66m^2n^2) - 32 \right) H(v) \\ & + 2 \left(\mathcal{K}(23(m^4 + n^4) - 41mn(m^2 + n^2) + 48m^2n^2) - 24 \right) H(u) \\ & + \left(\mathcal{K}(10(m^4 + n^4) - 15mn(m^2 + n^2) + 16m^2n^2) - 12 \right) H(t) \\ & + \left(\mathcal{K}((m^4 + n^4) - mn(m^2 + n^2) + m^2n^2) - 3 \right) H(s) \\ & + 5 \left(\mathcal{K}(50mn(m^3 + n^3) - 30m^2n^2(m+n) - 211(m^4 + n^4) \right. \\ & \quad \left. + 341mn(m^2 + n^2) - 381m^2n^2) + 162 \right) H(u)H(v) \\ & + \left(\mathcal{K}(50mn(m^3 + n^3) - 30m^2n^2(m+n) - 211(m^4 + n^4) \right. \\ & \quad \left. + 341mn(m^2 + n^2) - 381m^2n^2) + 162 \right) H(u)^2 \\ & + \left(\mathcal{K}(2mn(m^3 + n^3) - 19(m^4 + n^4) \right. \\ & \quad \left. + 29mn(m^2 + n^2) - 33m^2n^2) + 18 \right) H(t)^2 \\ & - \left(\mathcal{K}((m^4 + n^4) - mn(m^2 + n^2) + m^2n^2) - 2 \right) H(s)^2 \\ & + 2 \left(\mathcal{K}(48mn(m^3 + n^3) - 27m^2n^2(m+n) - 226(m^4 + n^4) \right. \\ & \quad \left. + 366mn(m^2 + n^2) - 412m^2n^2) + 180 \right) H(t)H(v) \\ & + 2 \left(\mathcal{K}(7mn(m^3 + n^3) - 3m^2n^2(m+n) - 85(m^4 + n^4) \right. \\ & \quad \left. + 146mn(m^2 + n^2) - 170m^2n^2) + 80 \right) H(s)H(v) \\ & + \left(\mathcal{K}(67mn(m^3 + n^3) - 36m^2n^2(m+n) - 333(m^4 + n^4) \right. \\ & \quad \left. + 538mn(m^2 + n^2) - 607m^2n^2) + 270 \right) H(t)H(u) \end{aligned}$$

F. Fourth Moment of $T_{n,m}$ Proof

$$\begin{aligned}
& + 2(\mathcal{K}(5mn(m^3 + n^3) - 2m^2n^2(m + n) - 63(m^4 + n^4) \\
& \quad + 107mn(m^2 + n^2) - 124m^2n^2) + 60)H(s)H(u) \\
& + (\mathcal{K}(mn(m^3 + n^3) - 27(m^4 + n^4) \\
& \quad + 40mn(m^2 + n^2) - 43m^2n^2) + 30)H(s)H(t) \\
& + 4(\mathcal{K}(6m^2n^2 - 844mn)(m^2 + n^2) - 165mn(m^3 + n^3) + 81m^2n^2(m + n) \\
& \quad + 542(m^4 + n^4) + 12m^3n^3 + 940m^2n^2) - 384)H(u)^2H(v) \\
& + 10(\mathcal{K}(5m^2n^2 - 736mn)(m^2 + n^2) - 142mn(m^3 + n^3) + 70m^2n^2(m + n) \\
& \quad + 472(m^4 + n^4) + 10m^3n^3 + 820m^2n^2) - 336)H(t)H(u)H(v) \\
& + 2(\mathcal{K}(2m^2n^2 - 412mn)(m^2 + n^2) - 73mn(m^3 + n^3) + 37m^2n^2(m + n) \\
& \quad + 262(m^4 + n^4) + 4m^3n^3 + 460m^2n^2) - 192)H(t)^2H(v) \\
& + 2(\mathcal{K}(5m^2n^2 - 736mn)(m^2 + n^2) - 142mn(m^3 + n^3) + 70m^2n^2(m + n) \\
& \quad + 472(m^4 + n^4) + 10m^3n^3 + 820m^2n^2) - 336)H(t)H(u)^2 \\
& + 2(\mathcal{K}(m^2n^2 - 304mn)(m^2 + n^2) - 50mn(m^3 + n^3) + 26m^2n^2(m + n) \\
& \quad + 192(m^4 + n^4) + 2m^3n^3 + 340m^2n^2) - 336)H(t)^2H(u) \\
& + 5(\mathcal{K}(2m^2n^2 - 805mn)(m^2 + n^2) - 131mn(m^3 + n^3) + 72m^2n^2(m + n) \\
& \quad + 505(m^4 + n^4) + 4m^3n^3 + 899m^2n^2) - 378)H(s)H(u)H(v) \\
& - 2(\mathcal{K}(7mn(m^3 + n^3) - 3m^2n^2(m + n) - 54(m^4 + n^4) \\
& \quad + 90mn(m^2 + n^2) - 104m^2n^2) - 48)H(s)^2H(v) \\
& + (\mathcal{K}(2m^2n^2 - 805mn)(m^2 + n^2) - 131mn(m^3 + n^3) + 72m^2n^2(m + n) \\
& \quad + 505(m^4 + n^4) + 4m^3n^3 + 899m^2n^2) - 378)H(s)H(u)^2 \\
& - 2(\mathcal{K}(5mn(m^3 + n^3) - 2m^2n^2(m + n) - 40(m^4 + n^4) \\
& \quad + 66mn(m^2 + n^2) - 76m^2n^2) + 36)H(s)^2H(u) \\
& - (\mathcal{K}(5mn(m^3 + n^3) - 45(m^4 + n^4) \\
& \quad + 69mn(m^2 + n^2) - 79m^2n^2) + 42)H(s)H(t)^2 \\
& - (\mathcal{K}(mn(m^3 + n^3) - 17(m^4 + n^4) \\
& \quad + 25mn(m^2 + n^2) - 27m^2n^2) + 18)H(s)^2H(t) \\
& + 2(\mathcal{K}(m^2n^2 - 866mn)(m^2 + n^2) - 125mn(m^3 + n^3) + 68m^2n^2(m + n) \\
& \quad + 540(m^4 + n^4) + 2m^3n^3 + 972m^2n^2) - 420)H(s)H(t)H(v) \\
& + (\mathcal{K}(m^2n^2 - 1274mn)(m^2 + n^2) - 174mn(m^3 + n^3) + 92m^2n^2(m + n) \\
& \quad + 795(m^4 + n^4) + 2m^3n^3 + 1433m^2n^2) - 630)H(s)H(t)H(u) \\
& - 9(\mathcal{K}(12m^2n^2 - 1336mn)(m^2 + n^2) - 279mn(m^3 + n^3) + 132m^2n^2(m + n) \\
& \quad + 866(m^4 + n^4) + 24m^3n^3 + 1486m^2n^2) - 600)H(t)H(u)^2H(v)
\end{aligned}$$

$$\begin{aligned}
& - 5(\mathcal{K}(12m^2n^2 - 1336mn)(m^2 + n^2) - 279mn(m^3 + n^3) + 132m^2n^2(m + n) \\
& \quad + 866(m^4 + n^4) + 24m^3n^3 + 1486m^2n^2) - 600)H(t)^2H(u)H(v) \\
& - (\mathcal{K}(12m^2n^2 - 1336mn)(m^2 + n^2) - 279mn(m^3 + n^3) + 132m^2n^2(m + n) \\
& \quad + 866(m^4 + n^4) + 24m^3n^3 + 1486m^2n^2) - 600)H(t)^2H(u)^2 \\
& - 5(\mathcal{K}(2m^2n^2 - 464mn)(m^2 + n^2) - 81mn(m^3 + n^3) + 42m^2n^2(m + n) \\
& \quad + 294(m^4 + n^4) + 4m^3n^3 + 518m^2n^2) - 216)H(s)^2H(u)H(v) \\
& - (\mathcal{K}(2m^2n^2 - 464mn)(m^2 + n^2) - 81mn(m^3 + n^3) + 42m^2n^2(m + n) \\
& \quad + 294(m^4 + n^4) + 4m^3n^3 + 518m^2n^2) - 216)H(s)^2H(u)^2 \\
& + (\mathcal{K}(3mn(m^3 + n^3) - 26(m^4 + n^4) \\
& \quad + 40mn(m^2 + n^2) - 46m^2n^2) + 24)H(s)^2H(t)^2 \\
& - 5(\mathcal{K}(25m^2n^2 - 3328mn)(m^2 + n^2) - 659mn(m^3 + n^3) + 320m^2n^2(m + n) \\
& \quad + 2142(m^4 + n^4) + 50m^3n^3 + 3706m^2n^2) - 1512)H(s)H(t)H(u)H(v) \\
& - 4(\mathcal{K}(15m^2n^2 - 1908mn)(m^2 + n^2) - 383mn(m^3 + n^3) + 185m^2n^2(m + n) \\
& \quad + 2130(m^4 + n^4) + 30m^3n^3 + 2124m^2n^2) - 864)H(s)H(u)^2H(v) \\
& - 2(\mathcal{K}(5m^2n^2 - 932mn)(m^2 + n^2) - 169mn(m^3 + n^3) + 85m^2n^2(m + n) \\
& \quad + 594(m^4 + n^4) + 10m^3n^3 + 1040m^2n^2) - 432)H(s)H(t)^2H(v) \\
& - 2(\mathcal{K}(m^2n^2 - 500mn)(m^2 + n^2) - 77mn(m^3 + n^3) + 41m^2n^2(m + n) \\
& \quad + 314(m^4 + n^4) + 2m^3n^3 + 560m^2n^2) - 240)H(s)^2H(t)H(v) \\
& - (\mathcal{K}(25m^2n^2 - 3328mn)(m^2 + n^2) - 659mn(m^3 + n^3) + 320m^2n^2(m + n) \\
& \quad + 2142(m^4 + n^4) + 50m^3n^3 + 3706m^2n^2) - 1512)H(s)H(t)H(u)^2 \\
& - (\mathcal{K}(5m^2n^2 - 1376mn)(m^2 + n^2) - 231mn(m^3 + n^3) + 120m^2n^2(m + n) \\
& \quad + 870(m^4 + n^4) + 10m^3n^3 + 1538m^2n^2) - 648)H(s)H(t)^2H(u) \\
& - (\mathcal{K}(m^2n^2 - 736mn)(m^2 + n^2) - 107mn(m^3 + n^3) + 56m^2n^2(m + n) \\
& \quad + 462(m^4 + n^4) + 2m^3n^3 + 826m^2n^2) - 360)H(s)^2H(t)H(u) \\
& + 6(\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720)H(t)^2H(u)^2H(v) \\
& + 9(\mathcal{K}(27m^2n^2 - 2499mn)(m^2 + n^2) - 619mn(m^3 + n^3) + 292m^2n^2(m + n) \\
& \quad + 1910(m^4 + n^4) + 54m^3n^3 + 3274m^2n^2) - 1320)H(s)H(t)H(u)^2H(v) \\
& + 4(\mathcal{K}(9m^2n^2 - 1064mn)(m^2 + n^2) - 218mn(m^3 + n^3) + 104m^2n^2(m + n) \\
& \quad + 688(m^4 + n^4) + 18m^3n^3 + 1184m^2n^2) - 480)H(s)^2H(u)^2H(v) \\
& + 5(\mathcal{K}(27m^2n^2 - 2499mn)(m^2 + n^2) - 619mn(m^3 + n^3) + 292m^2n^2(m + n) \\
& \quad + 1910(m^4 + n^4) + 54m^3n^3 + 3274m^2n^2) - 1320)H(s)H(t)^2H(u)H(v)
\end{aligned}$$

F. Fourth Moment of $T_{n,m}$ Proof

$$\begin{aligned}
& + 5(\mathcal{K}(15m^2n^2 - 1856mn)(m^2 + n^2) - 375mn(m^3 + n^3) + 180m^2n^2(m + n) \\
& \quad + 1198(m^4 + n^4) + 30m^3n^3 + 2066m^2n^2) - 840)H(s)^2H(t)H(u)H(v) \\
& + 2(\mathcal{K}(3m^2n^2 - 520mn)(m^2 + n^2) - 96mn(m^3 + n^3) + 48m^2n^2(m + n) \\
& \quad + 332(m^4 + n^4) + 6m^3n^3 + 580m^2n^2) - 240)H(s)^2H(t)^2H(v) \\
& + (\mathcal{K}(27m^2n^2 - 2499mn)(m^2 + n^2) - 619mn(m^3 + n^3) + 292m^2n^2(m + n) \\
& \quad + 1910(m^4 + n^4) + 54m^3n^3 + 3274m^2n^2) - 1320)H(s)H(t)^2H(u)^2 \\
& + (\mathcal{K}(15m^2n^2 - 1856mn)(m^2 + n^2) - 375mn(m^3 + n^3) + 180m^2n^2(m + n) \\
& \quad + 1198(m^4 + n^4) + 30m^3n^3 + 2066m^2n^2) - 840)H(s)^2H(t)H(u)^2 \\
& + (\mathcal{K}(3m^2n^2 - 768mn)(m^2 + n^2) - 131mn(m^3 + n^3) + 68m^2n^2(m + n) \\
& \quad + 486(m^4 + n^4) + 6m^3n^3 + 858m^2n^2) - 360)H(s)^2H(t)^2H(u) \\
& - 13(\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720)H(s)H(t)^2H(u)^2H(v) \\
& - 9(\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720)H(s)^2H(t)H(u)^2H(v) \\
& - 5(\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720)H(s)^2H(t)^2H(u)H(v) \\
& - (\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720)H(s)^2H(t)^2H(u)^2 \\
& + 7(\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720)H(s)^2H(t)^2H(u)^2H(v) \\
& dv du dt ds \tag{F.1}
\end{aligned}$$

where $\mathcal{K} = \frac{nm(n+m)^2}{n^7+m^7}$.

Proof Consider the test statistic

$$T_{n,m} = \frac{nm}{n+m} \int_{-\infty}^{\infty} (\hat{F}(t) - \hat{G}(t))^2 ds.$$

Now by taking the fourth power of the test statistic and using Fubini's theorem (Fubini, 1907), we obtain

$$\begin{aligned} \mathbb{E}[T_{n,m}^4] &= \mathbb{E} \left[\frac{n^4 m^4}{(n+m)^4} \iiint \int_{\mathbb{R}^4} (\hat{F}(s) - \hat{G}(s))^2 (\hat{F}(t) - \hat{G}(t))^2 \right. \\ &\quad \cdot (\hat{F}(u) - \hat{G}(u))^2 (\hat{F}(v) - \hat{G}(v))^2 dv du dt ds \left. \right] \\ &= \frac{n^4 m^4}{(n+m)^4} \iiint \int_{\mathbb{R}^4} \mathbb{E} \left[(\hat{F}(s) - \hat{G}(s))^2 (\hat{F}(t) - \hat{G}(t))^2 \right. \\ &\quad \cdot (\hat{F}(u) - \hat{G}(u))^2 (\hat{F}(v) - \hat{G}(v))^2 \left. \right] dv du dt ds \end{aligned} \quad (\text{F.2})$$

Substituting the expressions for $\hat{F}(t)$ and $\hat{G}(t)$ into equation (F.2) and expanding gives

$$\begin{aligned} \mathbb{E}[T_{n,m}^4] &= \frac{n^4 m^4}{(n+m)^4} \iiint \int_{\mathbb{R}^4} \mathbb{E} \left[\frac{1}{n^8} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^n \sum_{p=1}^n \sum_{q=1}^n \sum_{r=1}^n \left(\mathbb{1}_{\{x_i \geq s\}} \mathbb{1}_{\{x_j \geq s\}} \right. \right. \\ &\quad \cdot \mathbb{1}_{\{x_k \geq t\}} \mathbb{1}_{\{x_l \geq t\}} \mathbb{1}_{\{x_o \geq u\}} \mathbb{1}_{\{x_p \geq u\}} \mathbb{1}_{\{x_q \geq v\}} \mathbb{1}_{\{x_r \geq v\}} \\ &\quad - \frac{2}{n^7 m} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^n \sum_{p=1}^n \sum_{q=1}^m \sum_{r=1}^m \left(\mathbb{1}_{\{x_i \geq s\}} \mathbb{1}_{\{x_j \geq s\}} \mathbb{1}_{\{x_k \geq t\}} \right. \\ &\quad \cdot \mathbb{1}_{\{x_l \geq t\}} \mathbb{1}_{\{x_o \geq u\}} \mathbb{1}_{\{x_p \geq u\}} \mathbb{1}_{\{x_q \geq v\}} \mathbb{1}_{\{y_r \geq v\}} \left. \right) \\ &\quad + \frac{1}{n^6 m^2} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^n \sum_{p=1}^m \sum_{q=1}^m \sum_{r=1}^m \left(\mathbb{1}_{\{x_i \geq s\}} \mathbb{1}_{\{x_j \geq s\}} \mathbb{1}_{\{x_k \geq t\}} \right. \\ &\quad \cdot \mathbb{1}_{\{x_l \geq t\}} \mathbb{1}_{\{x_o \geq u\}} \mathbb{1}_{\{x_p \geq u\}} \mathbb{1}_{\{y_q \geq v\}} \mathbb{1}_{\{y_r \geq v\}} \left. \right) \\ &\quad - \frac{2}{n^7 m} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^m \sum_{p=1}^n \sum_{q=1}^n \sum_{r=1}^n \left(\mathbb{1}_{\{x_i \geq s\}} \mathbb{1}_{\{x_j \geq s\}} \mathbb{1}_{\{x_k \geq t\}} \right. \\ &\quad \cdot \mathbb{1}_{\{x_l \geq t\}} \mathbb{1}_{\{x_o \geq u\}} \mathbb{1}_{\{y_p \geq u\}} \mathbb{1}_{\{x_q \geq v\}} \mathbb{1}_{\{x_r \geq v\}} \left. \right) \\ &\quad + \frac{4}{n^6 m^2} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \sum_{o=1}^m \sum_{p=1}^m \sum_{q=1}^m \sum_{r=1}^m \left(\mathbb{1}_{\{x_i \geq s\}} \mathbb{1}_{\{x_j \geq s\}} \mathbb{1}_{\{x_k \geq t\}} \right. \\ &\quad \cdot \mathbb{1}_{\{x_l \geq t\}} \mathbb{1}_{\{x_o \geq u\}} \mathbb{1}_{\{y_p \geq u\}} \mathbb{1}_{\{x_q \geq v\}} \mathbb{1}_{\{y_r \geq v\}} \left. \right) \end{aligned}$$

$$\begin{aligned}
& - \frac{2}{n^3 m^5} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^n \sum_{p=1}^m \sum_{q=1}^n \sum_{r=1}^n \left(\mathbb{1}_{\{y_i \geq s\}} \mathbb{1}_{\{y_j \geq s\}} \mathbb{1}_{\{y_k \geq t\}} \right. \\
& \quad \left. \cdot \mathbb{1}_{\{y_l \geq t\}} \mathbb{1}_{\{x_o \geq u\}} \mathbb{1}_{\{y_p \geq u\}} \mathbb{1}_{\{x_q \geq v\}} \mathbb{1}_{\{x_r \geq v\}} \right) \\
& + \frac{4}{n^2 m^6} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^n \sum_{p=1}^m \sum_{q=1}^n \sum_{r=1}^m \left(\mathbb{1}_{\{y_i \geq s\}} \mathbb{1}_{\{y_j \geq s\}} \mathbb{1}_{\{y_k \geq t\}} \right. \\
& \quad \left. \cdot \mathbb{1}_{\{y_l \geq t\}} \mathbb{1}_{\{x_o \geq u\}} \mathbb{1}_{\{y_p \geq u\}} \mathbb{1}_{\{x_q \geq v\}} \mathbb{1}_{\{y_r \geq v\}} \right) \\
& - \frac{2}{n m^7} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^n \sum_{p=1}^m \sum_{q=1}^m \sum_{r=1}^m \left(\mathbb{1}_{\{y_i \geq s\}} \mathbb{1}_{\{y_j \geq s\}} \mathbb{1}_{\{y_k \geq t\}} \right. \\
& \quad \left. \cdot \mathbb{1}_{\{y_l \geq t\}} \mathbb{1}_{\{x_o \geq u\}} \mathbb{1}_{\{y_p \geq u\}} \mathbb{1}_{\{y_q \geq v\}} \mathbb{1}_{\{y_r \geq v\}} \right) \\
& + \frac{1}{n^2 m^6} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^m \sum_{p=1}^m \sum_{q=1}^n \sum_{r=1}^n \left(\mathbb{1}_{\{y_i \geq s\}} \mathbb{1}_{\{y_j \geq s\}} \mathbb{1}_{\{y_k \geq t\}} \right. \\
& \quad \left. \cdot \mathbb{1}_{\{y_l \geq t\}} \mathbb{1}_{\{y_o \geq u\}} \mathbb{1}_{\{y_p \geq u\}} \mathbb{1}_{\{x_q \geq v\}} \mathbb{1}_{\{x_r \geq v\}} \right) \\
& - \frac{2}{n m^7} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^m \sum_{p=1}^m \sum_{q=1}^n \sum_{r=1}^m \left(\mathbb{1}_{\{y_i \geq s\}} \mathbb{1}_{\{y_j \geq s\}} \mathbb{1}_{\{y_k \geq t\}} \right. \\
& \quad \left. \cdot \mathbb{1}_{\{y_l \geq t\}} \mathbb{1}_{\{y_o \geq u\}} \mathbb{1}_{\{y_p \geq u\}} \mathbb{1}_{\{x_q \geq v\}} \mathbb{1}_{\{y_r \geq v\}} \right) \\
& + \frac{1}{m^8} \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{l=1}^m \sum_{o=1}^m \sum_{p=1}^m \sum_{q=1}^m \sum_{r=1}^m \left(\mathbb{1}_{\{y_i \geq s\}} \mathbb{1}_{\{y_j \geq s\}} \mathbb{1}_{\{y_k \geq t\}} \right. \\
& \quad \left. \cdot \mathbb{1}_{\{y_l \geq t\}} \mathbb{1}_{\{y_o \geq u\}} \mathbb{1}_{\{y_p \geq u\}} \mathbb{1}_{\{y_q \geq v\}} \mathbb{1}_{\{y_r \geq v\}} \right) \Big] dv du dt ds \tag{F.3}
\end{aligned}$$

It is important to note that we can split the quadruple integral into the sum of 24 quadruple integrals for the following twenty four cases;

1. $s \leq t \leq u \leq v$,
2. $s \leq t \leq v \leq u$,
3. $s \leq u \leq t \leq v$,
4. $s \leq u \leq v \leq t$,
5. $s \leq v \leq t \leq u$,
6. $s \leq v \leq u \leq t$,
7. $t \leq s \leq u \leq v$,
8. $t \leq s \leq v \leq u$,

F. Fourth Moment of $T_{n,m}$ Proof

9. $t \leq u \leq t \leq v$,
10. $t \leq u \leq v \leq t$,
11. $t \leq v \leq t \leq u$,
12. $t \leq v \leq u \leq t$,
13. $u \leq s \leq t \leq v$,
14. $u \leq s \leq v \leq t$,
15. $u \leq t \leq s \leq v$,
16. $u \leq t \leq v \leq s$,
17. $u \leq v \leq s \leq t$,
18. $u \leq v \leq t \leq s$,
19. $v \leq s \leq t \leq u$,
20. $v \leq s \leq u \leq t$,
21. $v \leq t \leq s \leq u$,
22. $v \leq t \leq u \leq s$,
23. $v \leq u \leq s \leq t$,
24. $v \leq u \leq t \leq s$.

We also know that due to symmetry, all twenty four quadruple integrals will be equal. Thus we will only consider the first case and multiply the result by twenty four.

The expectation of each summation in equation (F.3) can be calculated separately by firstly considering all the cases for i, j, k, l, o, p, q, r (see for example Table D.1 and Table D.2). Each summation can then be split into each case for i, j, k, l, o, p, q, r and using the properties that under the null hypothesis $F \equiv G \equiv H$

$$\mathbb{E}[\mathbb{1}_{Z_i \leq \mathcal{A}}] = \Pr\{Z_i \leq \mathcal{A}\} = H(\mathcal{A}) \quad (\text{F.4})$$

where $Z = \{X_1, \dots, X_n, Y_1, \dots, Y_m\}$, each expectation can be calculated.

F. Fourth Moment of $T_{n,m}$ Proof

$$\begin{aligned}
& + n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)E\left[\mathbb{1}_{\{x_j \geq s\}}\mathbb{1}_{\{x_k \geq s\}}\mathbb{1}_{\{x_l \geq t\}}\mathbb{1}_{\{x_o \geq t\}}\right. \\
& \quad \left. \cdot \mathbb{1}_{\{x_i \geq u\}}^2\mathbb{1}_{\{x_p \geq v\}}\mathbb{1}_{\{x_q \geq v\}}\right] \\
& + n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)E\left[\mathbb{1}_{\{x_j \geq s\}}\mathbb{1}_{\{x_k \geq s\}}\mathbb{1}_{\{x_l \geq t\}}\mathbb{1}_{\{x_o \geq t\}}\right. \\
& \quad \left. \cdot \mathbb{1}_{\{x_p \geq u\}}\mathbb{1}_{\{x_q \geq u\}}\mathbb{1}_{\{x_i \geq v\}}^2\right] \\
& + 4n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)E\left[\mathbb{1}_{\{x_i \geq s\}}\mathbb{1}_{\{x_j \geq s\}}\mathbb{1}_{\{x_l \geq t\}}\mathbb{1}_{\{x_k \geq t\}}\right. \\
& \quad \left. \cdot \mathbb{1}_{\{x_l \geq u\}}\mathbb{1}_{\{x_o \geq u\}}\mathbb{1}_{\{x_p \geq v\}}\mathbb{1}_{\{x_q \geq v\}}\right] \\
& + 2n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)E\left[\mathbb{1}_{\{x_i \geq s\}}\mathbb{1}_{\{x_j \geq s\}}\mathbb{1}_{\{x_l \geq t\}}\mathbb{1}_{\{x_o \geq t\}}\right. \\
& \quad \left. \cdot \mathbb{1}_{\{x_i \geq u\}}\mathbb{1}_{\{x_k \geq u\}}\mathbb{1}_{\{x_p \geq v\}}\mathbb{1}_{\{x_q \geq v\}}\right] \\
& + 2n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)E\left[\mathbb{1}_{\{x_i \geq s\}}\mathbb{1}_{\{x_j \geq s\}}\mathbb{1}_{\{x_l \geq t\}}\mathbb{1}_{\{x_o \geq t\}}\right. \\
& \quad \left. \cdot \mathbb{1}_{\{x_i \geq u\}}\mathbb{1}_{\{x_k \geq u\}}\mathbb{1}_{\{x_p \geq v\}}\mathbb{1}_{\{x_q \geq v\}}\right] \\
& + 4n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)E\left[\mathbb{1}_{\{x_i \geq s\}}\mathbb{1}_{\{x_j \geq s\}}\mathbb{1}_{\{x_l \geq t\}}\mathbb{1}_{\{x_o \geq t\}}\right. \\
& \quad \left. \cdot \mathbb{1}_{\{x_p \geq u\}}\mathbb{1}_{\{x_q \geq u\}}\mathbb{1}_{\{x_i \geq v\}}\mathbb{1}_{\{x_k \geq v\}}\right] \\
& + 4n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)E\left[\mathbb{1}_{\{x_i \geq s\}}\mathbb{1}_{\{x_o \geq s\}}\mathbb{1}_{\{x_l \geq t\}}\mathbb{1}_{\{x_j \geq t\}}\right. \\
& \quad \left. \cdot \mathbb{1}_{\{x_i \geq u\}}\mathbb{1}_{\{x_k \geq u\}}\mathbb{1}_{\{x_p \geq v\}}\mathbb{1}_{\{x_q \geq v\}}\right] \\
& + 4n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)E\left[\mathbb{1}_{\{x_l \geq s\}}\mathbb{1}_{\{x_o \geq s\}}\mathbb{1}_{\{x_i \geq t\}}\mathbb{1}_{\{x_j \geq t\}}\right. \\
& \quad \left. \cdot \mathbb{1}_{\{x_p \geq u\}}\mathbb{1}_{\{x_q \geq u\}}\mathbb{1}_{\{x_i \geq v\}}\mathbb{1}_{\{x_k \geq v\}}\right] \\
& + 4n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)E\left[\mathbb{1}_{\{x_l \geq s\}}\mathbb{1}_{\{x_o \geq s\}}\mathbb{1}_{\{x_p \geq t\}}\mathbb{1}_{\{x_q \geq t\}}\right. \\
& \quad \left. \cdot \mathbb{1}_{\{x_i \geq u\}}\mathbb{1}_{\{x_j \geq u\}}\mathbb{1}_{\{x_i \geq v\}}\mathbb{1}_{\{x_k \geq v\}}\right] \\
& + n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)(n-7)E\left[\mathbb{1}_{\{x_i \geq s\}}\mathbb{1}_{\{x_j \geq s\}}\mathbb{1}_{\{x_k \geq t\}}\right. \\
& \quad \left. \cdot \mathbb{1}_{\{x_l \geq t\}}\mathbb{1}_{\{x_o \geq u\}}\mathbb{1}_{\{x_p \geq u\}}\mathbb{1}_{\{x_q \geq v\}}\mathbb{1}_{\{x_r \geq v\}}\right] \Bigg\}.
\end{aligned}$$

By using the property in equation (F.4) and simplifying, under the null hypothesis, the first summation in equation (F.3) is equivalent to

$$\begin{aligned}
& \frac{1}{n^8} \int_{-\infty}^{\infty} \int_{-\infty}^s \int_{-\infty}^t \int_{-\infty}^u \left\{ nH(v) + 64n(n-1)H(v)^2 + 48n(n-1)H(u)H(v) \right. \\
& \quad + 12n(n-1)H(t)H(v) + 405n(n-1)(n-2)H(u)H(v)^2 \\
& \quad + 81n(n-1)(n-2)H(u)^2H(v) + n(n-1)(n-2)H(s)^2H(v) \\
& \quad + 180n(n-1)(n-2)H(t)H(v)^2 + 80n(n-1)(n-2)H(s)H(v)^2 \\
& \quad + 135n(n-1)(n-2)H(t)H(u)H(v) + 60n(n-1)(n-2)H(s)H(u)H(v) \\
& \quad + 15n(n-1)(n-2)H(s)H(t)H(v) + 256n(n-1)(n-2)(n-3)H(u)^2H(v)^2 \\
& \quad + 560n(n-1)(n-2)(n-3)H(t)H(u)H(v)^2 \\
& \quad \left. + 64n(n-1)(n-2)(n-3)H(t)^2H(v)^2 \right\}
\end{aligned}$$

$$\begin{aligned}
& + 112n(n-1)(n-2)(n-3)H(t)H(u)^2H(v) \\
& + 48n(n-1)(n-2)(n-3)H(t)^2H(u)H(v) \\
& + 315n(n-1)(n-2)(n-3)H(s)H(u)H(v)^2 \\
& + 16n(n-1)(n-2)(n-3)H(s)^2H(v)^2 \\
& + 63n(n-1)(n-2)(n-3)H(s)H(u)^2H(v) \\
& + 12n(n-1)(n-2)(n-3)H(s)^2H(u)H(v) \\
& + 7n(n-1)(n-2)(n-3)H(s)H(t)^2H(v) \\
& + 3n(n-1)(n-2)(n-3)H(s)^2H(t)H(v) \\
& + 140n(n-1)(n-2)(n-3)H(s)H(t)H(v)^2 \\
& + 105n(n-1)(n-2)(n-3)H(s)H(t)H(u)H(v) \\
& + 225n(n-1)(n-2)(n-3)(n-4)H(t)H(u)^2H(v)^2 \\
& + 125n(n-1)(n-2)(n-3)(n-4)H(t)^2H(u)H(v)^2 \\
& + 25n(n-1)(n-2)(n-3)(n-4)H(t)^2H(u)^2H(v) \\
& + 45n(n-1)(n-2)(n-3)(n-4)H(s)^2H(u)H(v)^2 \\
& + 9n(n-1)(n-2)(n-3)(n-4)H(s)^2H(u)^2H(v) \\
& + n(n-1)(n-2)(n-3)(n-4)H(s)^2H(t)^2H(v) \\
& + 315n(n-1)(n-2)(n-3)(n-4)H(s)H(t)H(u)H(v)^2 \\
& + 144n(n-1)(n-2)(n-3)(n-4)H(s)H(u)^2H(v)^2 \\
& + 36n(n-1)(n-2)(n-3)(n-4)H(s)H(t)^2H(v)^2 \\
& + 20n(n-1)(n-2)(n-3)(n-4)H(s)^2H(t)H(v)^2 \\
& + 63n(n-1)(n-2)(n-3)(n-4)H(s)H(t)H(u)^2H(v) \\
& + 27n(n-1)(n-2)(n-3)(n-4)H(s)H(t)^2H(u)H(v) \\
& + 15n(n-1)(n-2)(n-3)(n-4)H(s)^2H(t)H(u)H(v) \\
& + 36n(n-1)(n-2)(n-3)(n-4)(n-5)H(t)^2H(u)^2H(v)^2 \\
& + 99n(n-1)(n-2)(n-3)(n-4)(n-5)H(s)H(t)H(u)^2H(v)^2 \\
& + 16n(n-1)(n-2)(n-3)(n-4)(n-5)H(s)^2H(u)^2H(v)^2 \\
& + 55n(n-1)(n-2)(n-3)(n-4)(n-5)H(s)H(t)^2H(u)H(v)^2 \\
& + 35n(n-1)(n-2)(n-3)(n-4)(n-5)H(s)^2H(t)H(u)H(v)^2 \\
& + 4n(n-1)(n-2)(n-3)(n-4)(n-5)H(s)^2H(t)^2H(v)^2 \\
& + 11n(n-1)(n-2)(n-3)(n-4)(n-5)H(s)H(t)^2H(u)^2H(v)^2 \\
& + 7n(n-1)(n-2)(n-3)(n-4)(n-5)H(s)^2H(t)H(u)^2H(v) \\
& + 3n(n-1)(n-2)(n-3)(n-4)(n-5)H(s)^2H(t)^2H(u)H(v) \\
& + 13n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)H(s)H(t)^2H(u)^2H(v)^2
\end{aligned}$$

F. Fourth Moment of $T_{n,m}$ Proof

$$\begin{aligned}
& + 9n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)H(s)^2H(t)H(u)^2H(v)^2 \\
& + 5n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)H(s)^2H(t)^2H(u)H(v)^2 \\
& + n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)H(s)^2H(t)^2H(u)^2H(v) \\
& + n(n-1)(n-2)(n-3)(n-4)(n-5)(n-6)(n-7)H(s)^2H(t)^2H(u)^2H(v)^2 \} \\
& dv du dt ds.
\end{aligned}$$

As for the proof of $\text{Var}[T_{n,m}]$, this process can be repeated for each summation in equation (F.3) with care being taken where both data X and Y are involved. Once we repeat this process for each summation in equation (F.3), multiply by 24 and collect together all the terms, we get

$$\begin{aligned}
\mathbb{E}[A^4] = & \frac{24}{\mathcal{K}(nm(n+m))^2} \int_{-\infty}^{\infty} \int_{-\infty}^s \int_{-\infty}^t \int_{-\infty}^u H(v) \left(1 \right. \\
& + 2 \left(\mathcal{K}(31(m^4 + n^4) - 56mn(m^2 + n^2) + 66m^2n^2) - 32 \right) H(v) \\
& + 2 \left(\mathcal{K}(23(m^4 + n^4) - 41mn(m^2 + n^2) + 48m^2n^2) - 24 \right) H(u) \\
& + \left(\mathcal{K}(10(m^4 + n^4) - 15mn(m^2 + n^2) + 16m^2n^2) - 12 \right) H(t) \\
& + \left(\mathcal{K}((m^4 + n^4) - mn(m^2 + n^2) + m^2n^2) - 3 \right) H(s) \\
& + 5 \left(\mathcal{K}(50mn(m^3 + n^3) - 30m^2n^2(m+n) - 211(m^4 + n^4) \right. \\
& \quad \left. + 341mn(m^2 + n^2) - 381m^2n^2) + 162 \right) H(u)H(v) \\
& + \left(\mathcal{K}(50mn(m^3 + n^3) - 30m^2n^2(m+n) - 211(m^4 + n^4) \right. \\
& \quad \left. + 341mn(m^2 + n^2) - 381m^2n^2) + 162 \right) H(u)^2 \\
& + \left(\mathcal{K}(2mn(m^3 + n^3) - 19(m^4 + n^4) \right. \\
& \quad \left. + 29mn(m^2 + n^2) - 33m^2n^2) + 18 \right) H(t)^2 \\
& - \left(\mathcal{K}((m^4 + n^4) - mn(m^2 + n^2) + m^2n^2) - 2 \right) H(s)^2 \\
& + 2 \left(\mathcal{K}(48mn(m^3 + n^3) - 27m^2n^2(m+n) - 226(m^4 + n^4) \right. \\
& \quad \left. + 366mn(m^2 + n^2) - 412m^2n^2) + 180 \right) H(t)H(v) \\
& + 2 \left(\mathcal{K}(7mn(m^3 + n^3) - 3m^2n^2(m+n) - 85(m^4 + n^4) \right. \\
& \quad \left. + 146mn(m^2 + n^2) - 170m^2n^2) + 80 \right) H(s)H(v) \\
& + \left(\mathcal{K}(67mn(m^3 + n^3) - 36m^2n^2(m+n) - 333(m^4 + n^4) \right. \\
& \quad \left. + 538mn(m^2 + n^2) - 607m^2n^2) + 270 \right) H(t)H(u) \\
& + 2 \left(\mathcal{K}(5mn(m^3 + n^3) - 2m^2n^2(m+n) - 63(m^4 + n^4) \right. \\
& \quad \left. + 107mn(m^2 + n^2) - 124m^2n^2) + 60 \right) H(s)H(u) \\
& + \left(\mathcal{K}(mn(m^3 + n^3) - 27(m^4 + n^4) \right. \\
& \quad \left. + 40mn(m^2 + n^2) - 43m^2n^2) + 30 \right) H(s)H(t)
\end{aligned}$$

$$\begin{aligned}
& + 4(\mathcal{K}(6m^2n^2 - 844mn)(m^2 + n^2) - 165mn(m^3 + n^3) + 81m^2n^2(m + n) \\
& \quad + 542(m^4 + n^4) + 12m^3n^3 + 940m^2n^2) - 384)H(u)^2H(v) \\
& + 10(\mathcal{K}(5m^2n^2 - 736mn)(m^2 + n^2) - 142mn(m^3 + n^3) + 70m^2n^2(m + n) \\
& \quad + 472(m^4 + n^4) + 10m^3n^3 + 820m^2n^2) - 336)H(t)H(u)H(v) \\
& + 2(\mathcal{K}(2m^2n^2 - 412mn)(m^2 + n^2) - 73mn(m^3 + n^3) + 37m^2n^2(m + n) \\
& \quad + 262(m^4 + n^4) + 4m^3n^3 + 460m^2n^2) - 192)H(t)^2H(v) \\
& + 2(\mathcal{K}(5m^2n^2 - 736mn)(m^2 + n^2) - 142mn(m^3 + n^3) + 70m^2n^2(m + n) \\
& \quad + 472(m^4 + n^4) + 10m^3n^3 + 820m^2n^2) - 336)H(t)H(u)^2 \\
& + 2(\mathcal{K}(m^2n^2 - 304mn)(m^2 + n^2) - 50mn(m^3 + n^3) + 26m^2n^2(m + n) \\
& \quad + 192(m^4 + n^4) + 2m^3n^3 + 340m^2n^2) - 336)H(t)^2H(u) \\
& + 5(\mathcal{K}(2m^2n^2 - 805mn)(m^2 + n^2) - 131mn(m^3 + n^3) + 72m^2n^2(m + n) \\
& \quad + 505(m^4 + n^4) + 4m^3n^3 + 899m^2n^2) - 378)H(s)H(u)H(v) \\
& - 2(\mathcal{K}(7mn(m^3 + n^3) - 3m^2n^2(m + n) - 54(m^4 + n^4) \\
& \quad + 90mn(m^2 + n^2) - 104m^2n^2) - 48)H(s)^2H(v) \\
& + (\mathcal{K}(2m^2n^2 - 805mn)(m^2 + n^2) - 131mn(m^3 + n^3) + 72m^2n^2(m + n) \\
& \quad + 505(m^4 + n^4) + 4m^3n^3 + 899m^2n^2) - 378)H(s)H(u)^2 \\
& - 2(\mathcal{K}(5mn(m^3 + n^3) - 2m^2n^2(m + n) - 40(m^4 + n^4) \\
& \quad + 66mn(m^2 + n^2) - 76m^2n^2) + 36)H(s)^2H(u) \\
& - (\mathcal{K}(5mn(m^3 + n^3) - 45(m^4 + n^4) \\
& \quad + 69mn(m^2 + n^2) - 79m^2n^2) + 42)H(s)H(t)^2 \\
& - (\mathcal{K}(mn(m^3 + n^3) - 17(m^4 + n^4) \\
& \quad + 25mn(m^2 + n^2) - 27m^2n^2) + 18)H(s)^2H(t) \\
& + 2(\mathcal{K}(m^2n^2 - 866mn)(m^2 + n^2) - 125mn(m^3 + n^3) + 68m^2n^2(m + n) \\
& \quad + 540(m^4 + n^4) + 2m^3n^3 + 972m^2n^2) - 420)H(s)H(t)H(v) \\
& + (\mathcal{K}(m^2n^2 - 1274mn)(m^2 + n^2) - 174mn(m^3 + n^3) + 92m^2n^2(m + n) \\
& \quad + 795(m^4 + n^4) + 2m^3n^3 + 1433m^2n^2) - 630)H(s)H(t)H(u) \\
& - 9(\mathcal{K}(12m^2n^2 - 1336mn)(m^2 + n^2) - 279mn(m^3 + n^3) + 132m^2n^2(m + n) \\
& \quad + 866(m^4 + n^4) + 24m^3n^3 + 1486m^2n^2) - 600)H(t)H(u)^2H(v) \\
& - 5(\mathcal{K}(12m^2n^2 - 1336mn)(m^2 + n^2) - 279mn(m^3 + n^3) + 132m^2n^2(m + n) \\
& \quad + 866(m^4 + n^4) + 24m^3n^3 + 1486m^2n^2) - 600)H(t)^2H(u)H(v) \\
& - (\mathcal{K}(12m^2n^2 - 1336mn)(m^2 + n^2) - 279mn(m^3 + n^3) + 132m^2n^2(m + n) \\
& \quad + 866(m^4 + n^4) + 24m^3n^3 + 1486m^2n^2) - 600)H(t)^2H(u)^2
\end{aligned}$$

F. Fourth Moment of $T_{n,m}$ Proof

$$\begin{aligned}
& - 5(\mathcal{K}(2m^2n^2 - 464mn)(m^2 + n^2) - 81mn(m^3 + n^3) + 42m^2n^2(m + n) \\
& \quad + 294(m^4 + n^4) + 4m^3n^3 + 518m^2n^2) - 216)H(s)^2H(u)H(v) \\
& - (\mathcal{K}(2m^2n^2 - 464mn)(m^2 + n^2) - 81mn(m^3 + n^3) + 42m^2n^2(m + n) \\
& \quad + 294(m^4 + n^4) + 4m^3n^3 + 518m^2n^2) - 216)H(s)^2H(u)^2 \\
& + (\mathcal{K}(3mn(m^3 + n^3) - 26(m^4 + n^4) \\
& \quad + 40mn(m^2 + n^2) - 46m^2n^2) + 24)H(s)^2H(t)^2 \\
& - 5(\mathcal{K}(25m^2n^2 - 3328mn)(m^2 + n^2) - 659mn(m^3 + n^3) + 320m^2n^2(m + n) \\
& \quad + 2142(m^4 + n^4) + 50m^3n^3 + 3706m^2n^2) - 1512)H(s)H(t)H(u)H(v) \\
& - 4(\mathcal{K}(15m^2n^2 - 1908mn)(m^2 + n^2) - 383mn(m^3 + n^3) + 185m^2n^2(m + n) \\
& \quad + 2130(m^4 + n^4) + 30m^3n^3 + 2124m^2n^2) - 864)H(s)H(u)^2H(v) \\
& - 2(\mathcal{K}(5m^2n^2 - 932mn)(m^2 + n^2) - 169mn(m^3 + n^3) + 85m^2n^2(m + n) \\
& \quad + 594(m^4 + n^4) + 10m^3n^3 + 1040m^2n^2) - 432)H(s)H(t)^2H(v) \\
& - 2(\mathcal{K}(m^2n^2 - 500mn)(m^2 + n^2) - 77mn(m^3 + n^3) + 41m^2n^2(m + n) \\
& \quad + 314(m^4 + n^4) + 2m^3n^3 + 560m^2n^2) - 240)H(s)^2H(t)H(v) \\
& - (\mathcal{K}(25m^2n^2 - 3328mn)(m^2 + n^2) - 659mn(m^3 + n^3) + 320m^2n^2(m + n) \\
& \quad + 2142(m^4 + n^4) + 50m^3n^3 + 3706m^2n^2) - 1512)H(s)H(t)H(u)^2 \\
& - (\mathcal{K}(5m^2n^2 - 1376mn)(m^2 + n^2) - 231mn(m^3 + n^3) + 120m^2n^2(m + n) \\
& \quad + 870(m^4 + n^4) + 10m^3n^3 + 1538m^2n^2) - 648)H(s)H(t)^2H(u) \\
& - (\mathcal{K}(m^2n^2 - 736mn)(m^2 + n^2) - 107mn(m^3 + n^3) + 56m^2n^2(m + n) \\
& \quad + 462(m^4 + n^4) + 2m^3n^3 + 826m^2n^2) - 360)H(s)^2H(t)H(u) \\
& + 6(\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720)H(t)^2H(u)^2H(v) \\
& + 9(\mathcal{K}(27m^2n^2 - 2499mn)(m^2 + n^2) - 619mn(m^3 + n^3) + 292m^2n^2(m + n) \\
& \quad + 1910(m^4 + n^4) + 54m^3n^3 + 3274m^2n^2) - 1320)H(s)H(t)H(u)^2H(v) \\
& + 4(\mathcal{K}(9m^2n^2 - 1064mn)(m^2 + n^2) - 218mn(m^3 + n^3) + 104m^2n^2(m + n) \\
& \quad + 688(m^4 + n^4) + 18m^3n^3 + 1184m^2n^2) - 480)H(s)^2H(u)^2H(v) \\
& + 5(\mathcal{K}(27m^2n^2 - 2499mn)(m^2 + n^2) - 619mn(m^3 + n^3) + 292m^2n^2(m + n) \\
& \quad + 1910(m^4 + n^4) + 54m^3n^3 + 3274m^2n^2) - 1320)H(s)H(t)^2H(u)H(v) \\
& + 5(\mathcal{K}(15m^2n^2 - 1856mn)(m^2 + n^2) - 375mn(m^3 + n^3) + 180m^2n^2(m + n) \\
& \quad + 1198(m^4 + n^4) + 30m^3n^3 + 2066m^2n^2) - 840)H(s)^2H(t)H(u)H(v) \\
& + 2(\mathcal{K}(3m^2n^2 - 520mn)(m^2 + n^2) - 96mn(m^3 + n^3) + 48m^2n^2(m + n) \\
& \quad + 332(m^4 + n^4) + 6m^3n^3 + 580m^2n^2) - 240)H(s)^2H(t)^2H(v)
\end{aligned}$$

$$\begin{aligned}
& + (\mathcal{K}(27m^2n^2 - 2499mn)(m^2 + n^2) - 619mn(m^3 + n^3) + 292m^2n^2(m + n) \\
& \quad + 1910(m^4 + n^4) + 54m^3n^3 + 3274m^2n^2) - 1320) H(s)H(t)^2H(u)^2 \\
& + (\mathcal{K}(15m^2n^2 - 1856mn)(m^2 + n^2) - 375mn(m^3 + n^3) + 180m^2n^2(m + n) \\
& \quad + 1198(m^4 + n^4) + 30m^3n^3 + 2066m^2n^2) - 840) H(s)^2H(t)H(u)^2 \\
& + (\mathcal{K}(3m^2n^2 - 768mn)(m^2 + n^2) - 131mn(m^3 + n^3) + 68m^2n^2(m + n) \\
& \quad + 486(m^4 + n^4) + 6m^3n^3 + 858m^2n^2) - 360) H(s)^2H(t)^2H(u) \\
& - 13(\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720) H(s)H(t)^2H(u)^2H(v) \\
& - 9(\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720) H(s)^2H(t)H(u)^2H(v) \\
& - 5(\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720) H(s)^2H(t)^2H(u)H(v) \\
& - (\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720) H(s)^2H(t)^2H(u)^2 \\
& + 7(\mathcal{K}(15m^2n^2 - 1608mn)(m^2 + n^2) - 340mn(m^3 + n^3) + 160m^2n^2(m + n) \\
& \quad + 1044(m^4 + n^4) + 30m^3n^3 + 1788m^2n^2) - 720) H(s)^2H(t)^2H(u)^2H(v)
\end{aligned}$$

$dv du dt ds$

where $\mathcal{K} = \frac{nm(n+m)^2}{n^7+m^7}$.

□

References

- Abyzov, A., Urban, A. E., Snyder, M., and Gerstein, M. (2011). CNVnator: An approach to discover, genotype and characterize typical and atypical CNVs from family and population genome sequencing. *Genome research*, pages gr-114876. [13](#)
- Aguirre, A. J., Brennan, C., Bailey, G., Sinha, R., Feng, B., Leo, C., Zhang, Y., Zhang, J., Gans, J. D., Bardeesy, N., et al. (2004). High-resolution characterization of the pancreatic adenocarcinoma genome. *Proceedings of the National Academy of Sciences*, 101(24):9067–9072.
- Anderson, T. W. (1962). On the distribution of the two-sample Cramer-von Mises criterion. *The Annals of Mathematical Statistics*, 33:1148–1159. [40](#)
- Anderson, T. W. and Darling, D. A. (1954). A test of goodness of fit. *Journal of the American Statistical Association*, 49(268):765–769. [40](#)
- Atkeson, C. G., Moore, A. W., and Schaal, S. (1997). Locally weighted learning for control. In *Lazy learning*, pages 75–113. Springer. [25](#)
- Azzalini, A. (1985). A class of distributions which includes the normal ones. *Scandinavian journal of statistics*, pages 171–178. [157](#)
- Balkema, A. A. and De Haan, L. (1974). Residual life time at great age. *The Annals of probability*, pages 792–804. [175](#)
- Balkin, S. D. and Mallows, C. L. (2001). An adjusted, asymmetric two-sample t test. *The American Statistician*, 55(3):203–206. [36](#), [46](#), [159](#), [160](#)
- Baringhaus, L. and Franz, C. (2004). On a new multivariate two-sample test. *Journal of multivariate analysis*, 88(1):190–206. [31](#), [32](#), [42](#), [47](#), [61](#), [62](#), [72](#), [84](#), [144](#), [152](#)
- Baumgartner, W., Weiß, P., and Schindler, H. (1998). A nonparametric test for the general two-sample problem. *Biometrics*, pages 1129–1135. [40](#)

REFERENCES

- Behjati, S. and Tarpey, P. S. (2013). What is next generation sequencing? *Archives of Disease in Childhood-Education and Practice*, 98(6):236–238. [6](#)
- Belvedere, O., Berri, S., Chalkley, R., Conway, C., Barbone, F., Pisa, F., MacLennan, K., Daly, C., Alsop, M., Morgan, J., et al. (2012). A computational index derived from whole-genome copy number analysis is a novel tool for prognosis in early stage lung squamous cell carcinoma. *Genomics*, 99(1):18–24. [7](#), [32](#), [153](#)
- Ben-Yaacov, E. and Eldar, Y. C. (2008). A fast and flexible method for the segmentation of acgh data. *Bioinformatics*, 24(16):i139–i145. [16](#)
- Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the royal statistical society. Series B (Methodological)*, pages 289–300. [121](#)
- Beroukhi, R., Getz, G., Nghiemphu, L., Barretina, J., Hsueh, T., Linhart, D., Vivanco, I., Lee, J. C., Huang, J. H., Alexander, S., et al. (2007). Assessing the significance of chromosomal aberrations in cancer: methodology and application to glioma. *Proceedings of the National Academy of Sciences*, 104(50):20007–20012. [19](#)
- Beroukhi, R., Mermel, C. H., Porter, D., Wei, G., Raychaudhuri, S., Donovan, J., Barretina, J., Boehm, J. S., Dobson, J., Urashima, M., et al. (2010). The landscape of somatic copy-number alteration across human cancers. *Nature*, 463(7283):899. [1](#)
- Björkqvist, A.-M., Husgafvel-Pursiainen, K., Anttila, S., Karjalainen, A., Tammlahti, L., Mattson, K., Vainio, H., and Knuutila, S. (1998). DNA gains in 3q occur frequently in squamous cell carcinoma of the lung, but not in adenocarcinoma. *Genes, Chromosomes and Cancer*, 22(1):79–82. [29](#), [31](#), [93](#)
- Boeva, V., Zinovyev, A., Bleakley, K., Vert, J.-P., Janoueix-Lerosey, I., Delattre, O., and Barillot, E. (2010). Control-free calling of copy number alterations in deep-sequencing data using GC-content normalization. *Bioinformatics*, 27(2):268–269. [13](#)
- Box, G. E., Jenkins, G. M., Reinsel, G. C., and Ljung, G. M. (2015). *Time series analysis: forecasting and control*. John Wiley & Sons. [102](#)
- Brown, M. B. (1975). 400: A method for combining non-independent, one-sided tests of significance. *Biometrics*, pages 987–992. [32](#), [122](#), [132](#), [133](#), [154](#)
- Burr, E. (1964). Small-sample distributions of the two-sample Cramer-von Mises' W^2 and Watson's U^2 . *The Annals of Mathematical Statistics*, 35:1091–1098. [40](#)

-
- Cao, R. and Van Keilegom, I. (2006). Empirical likelihood tests for two-sample problems via nonparametric density estimation. *Canadian Journal of Statistics*, 34(1):61–77. [46](#)
- Cho, E. and Cho, M. J. (2009). Variance of sample variance with replacement. *International Journal of Pure and Applied Mathematics*, 52:43–47. [163](#)
- Choi, W., Ochoa, A., McConkey, D. J., Aine, M., Höglund, M., Kim, W. Y., Real, F. X., Kiltie, A. E., Milsom, I., Dyrskjøt, L., et al. (2017). Genetic alterations in the molecular subtypes of bladder cancer: illustration in the cancer genome atlas dataset. *European urology*, 72(3):354–365. [1](#)
- Choo-Wosoba, H., Albert, P. S., and Zhu, B. (2018). hsegHMM: hidden Markov model-based allele-specific copy number alteration analysis accounting for hypersegmentation. *BMC bioinformatics*, 19(1):424. [17](#)
- Colella, S., Yau, C., Taylor, J. M., Mirza, G., Butler, H., Clouston, P., Bassett, A. S., Seller, A., Holmes, C. C., and Ragoussis, J. (2007). QuantiSNP: an Objective Bayes Hidden-Markov Model to detect and accurately map copy number variation using SNP genotyping data. *Nucleic acids research*, 35(6):2013–2025. [16](#)
- Cornish, E. A. and Fisher, R. A. (1938). Moments and cumulants in the specification of distributions. *Revue de l'Institut international de Statistique*, pages 307–320. [160](#)
- Cramér, H. (1928). On the composition of elementary errors: First paper: Mathematical deductions. *Scandinavian Actuarial Journal*, 1928(1):13–74. [39](#)
- Csorgo, S. and Faraway, J. J. (1996). The exact and asymptotic distributions of Cramér-von Mises statistics. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 221–234. [39](#), [48](#), [59](#)
- Curry, J., Dang, X., and Sang, H. (2018). A rank-based Cramer-von-Mises-type test for two samples. *arXiv preprint arXiv:1802.06332*. [46](#)
- Darling, D. A. (1957). The Kolmogorov-Smirnov, Cramer-von Mises tests. *The Annals of Mathematical Statistics*, 28(4):823–838. [40](#)
- De Ronde, J. and Klijn, C. (2013). KC-SMART Vignette. [20](#)
- de Ronde, J., Klijn, C., Velds, A., de Ronde, M. J., and aCGH, M. (2019). Package KCsmart. [20](#)

REFERENCES

- de Ronde, J. J., Klijn, C., Velds, A., Holstege, H., Reinders, M. J., Jonkers, J., and Wessels, L. F. (2010). KC-SMARTR: An R package for detection of statistically significant aberrations in multi-experiment aCGH data. *BMC research notes*, 3(1):298. [19](#), [27](#)
- Demirtas, H. (2014). Generating bivariate uniform data with a full range of correlations and connections to bivariate binary data. *Communications in Statistics-Theory and Methods*, 43(17):3574–3579. [134](#)
- Diskin, S. J., Eck, T., Greshock, J., Mosse, Y. P., Naylor, T., Stoeckert, C. J., Weber, B. L., Maris, J. M., and Grant, G. R. (2006). STAC: A method for testing the significance of DNA copy number aberrations across multiple array-CGH experiments. *Genome research*, 16(9):1149–1158. [19](#)
- Donoho, D. L. and Johnstone, I. M. (1995). Adapting to unknown smoothness via wavelet shrinkage. *Journal of the american statistical association*, 90(432):1200–1224. [18](#)
- Dudbridge, F. and Gusnanto, A. (2008). Estimation of significance thresholds for genomewide association scans. *Genetic epidemiology*, 32(3):227–234. [xviii](#), [xxv](#), [32](#), [122](#), [124](#), [126](#), [127](#), [130](#), [147](#), [153](#)
- Dudbridge, F. and Koeleman, B. P. (2004). Efficient computation of significance levels for multiple associations in large studies of correlated data, including genomewide association studies. *The American Journal of Human Genetics*, 75(3):424–435. [128](#)
- Efron, B. (1992). Bootstrap methods: another look at the jackknife. In *Breakthroughs in statistics*, pages 569–593. Springer. [62](#)
- Engler, D. A., Mohapatra, G., Louis, D. N., and Betensky, R. A. (2006). A pseudolikelihood approach for simultaneous analysis of array comparative genomic hybridizations. *Biostatistics*, 7(3):399–421. [17](#)
- Ferguson, T. S. (1995). A class of symmetric bivariate uniform distributions. *Statistical Papers*, 36(1):31. [133](#), [134](#)
- Fernández, V. A., Gamero, M. J., and García, J. M. (2008). A test for the two-sample problem based on empirical characteristic functions. *Computational statistics & data analysis*, 52(7):3730–3748. [45](#), [46](#)
- Fisher, R. A. (1935). *The design of experiments*. Oliver & Boyd. [62](#)

- Fisher, R. A. et al. (1925). *Statistical methods for research workers*. Edinburgh: Oliver and Boyd. 131
- Fisz, M. (1960). On a result by M. Rosenblatt concerning the von Mises-Smirnov test. *The Annals of Mathematical Statistics*, pages 427–429. 40
- Fridlyand, J., Snijders, A. M., Pinkel, D., Albertson, D. G., and Jain, A. N. (2004). Hidden Markov models approach to the analysis of array CGH data. *Journal of multivariate analysis*, 90(1):132–153. 16
- Fubini, G. (1907). Sugli integrali multipli. *Rend. Acc. Naz. Lincei*, 16:608–614. 48, 186, 194, 207
- Furey, T. S. and Haussler, D. (2003). Integration of the cytogenetic map with the draft human genome sequence. *Human Molecular Genetics*, 12(9):1037–1044. 14
- Genome Reference Consortium (2017). Human Genome Assembly GRCh38.p12. <https://www.ncbi.nlm.nih.gov/grc/human/data>. Accessed: 2018-11-30. 7
- Götze, F. (1979). Asymptotic expansions for bivariate von Mises functionals. *Probability Theory and Related Fields*, 50(3):333–355. 39
- Govindarajan, R., Duraiyan, J., Kaliyappan, K., and Palanisamy, M. (2012). Microarray and its applications. *Journal of pharmacy & bioallied sciences*, 4(Suppl 2):S310. 5
- Gretton, A., Borgwardt, K. M., Rasch, M., Schölkopf, B., and Smola, A. J. (2007). A kernel method for the two-sample-problem. In *Advances in neural information processing systems*, pages 513–520. 45
- Gusnanto, A., Taylor, C. C., Nafisah, I., Wood, H. M., Rabbitts, P., and Berri, S. (2014). Estimating optimal window size for analysis of low-coverage next-generation sequence data. *Bioinformatics*, 30(13):1823–1829. 6
- Gusnanto, A., Tcherveniakov, P., Shuweihi, F., Samman, M., Rabbitts, P., and Wood, H. M. (2015). Stratifying tumour subtypes based on copy number alteration profiles using next-generation sequence data. *Bioinformatics*, 31(16):2713–2720. 1, 7
- Gusnanto, A., Wood, H. M., Pawitan, Y., Rabbitts, P., and Berri, S. (2011). Correcting for cancer genome size and tumour cell content enables better estimation of copy number alterations from next-generation sequence data. *Bioinformatics*, 28(1):40–47. 6, 13, 14

REFERENCES

- Hartung, J. (1999). A note on combining dependent tests of significance. *Biometrical Journal: Journal of Mathematical Methods in Biosciences*, 41(7):849–855. [148](#)
- Heyde, C. (2014). Central limit theorem. *Wiley StatsRef: Statistics Reference Online*. [36](#)
- Hsu, L., Self, S. G., Grove, D., Randolph, T., Wang, K., Delrow, J. J., Loo, L., and Porter, P. (2005). Denoising array-based comparative genomic hybridization data using wavelets. *Biostatistics*, 6(2):211–226. [17](#), [18](#)
- Huang, J., Gusnanto, A., O’Sullivan, K., Staaf, J., Borg, Å., and Pawitan, Y. (2007). Robust smooth segmentation approach for array cgh data analysis. *Bioinformatics*, 23(18):2463–2469. [17](#), [19](#), [20](#)
- Hurd, P. J. and Nelson, C. J. (2009). Advantages of next-generation sequencing versus the microarray in epigenetic research. *Briefings in Functional Genomics and Proteomics*, 8(3):174–183. [6](#)
- Jianqing, F. and Gijbels, I. (1996). Local polynomial modelling and its applications. *Monographs on Statistics and Applied Probability*. Chapman & Hall/CRC. [18](#)
- Jiménez-Gamero, M.-D., Alba-Fernández, M., Jodrá, P., and Barranco-Chamorro, I. (2017). Fast tests for the two-sample problem based on the empirical characteristic function. *Mathematics and Computers in Simulation*, 137:390–410. [45](#)
- Kallioniemi, A., Kallioniemi, O.-P., Sudar, D., Rutovitz, D., Gray, J. W., Waldman, F., and Pinkel, D. (1992). Comparative genomic hybridization for molecular cytogenetic analysis of solid tumors. *Science*, 258(5083):818–821. [4](#), [5](#)
- Klijn, C., Holstege, H., de Ridder, J., Liu, X., Reinders, M., Jonkers, J., and Wessels, L. (2008). Identification of cancer genes using a statistical framework for multi-experiment analysis of non-discretized array CGH data. *Nucleic acids research*, 36(2):e13–e13. [25](#)
- Kolmogorov, A. N. (1933). *Sulla determinazione empirica di una legge di distribuzione*, volume 4. Giornale dell’Istituto Italiano degli Attuari. [37](#)
- Kost, J. T. and McDermott, M. P. (2002). Combining dependent p-values. *Statistics & Probability Letters*, 60(2):183–190. [32](#), [132](#), [133](#), [154](#)

-
- Kruskal, W. H. and Wallis, W. A. (1952). Use of ranks in one-criterion variance analysis. *Journal of the American statistical Association*, 47(260):583–621. [155](#)
- Krylov, V. (1962). Approximate calculation of integrals. *New York*, pages 100–111. [133](#)
- Lachman, H. M., Pedrosa, E., Petruolo, O. A., Cockerham, M., Papolos, A., Novak, T., Papolos, D. F., and Stopkova, P. (2007). Increase in GSK3 β gene copy number variation in bipolar disorder. *American Journal of Medical Genetics Part B: Neuropsychiatric Genetics*, 144(3):259–265. [4](#)
- Lai, W. R., Johnson, M. D., Kucherlapati, R., and Park, P. J. (2005). Comparative analysis of algorithms for identifying amplifications and deletions in array cgh data. *Bioinformatics*, 21(19):3763–3770. [16](#)
- Lehmann, E. L. (1951). Consistency and unbiasedness of certain nonparametric tests. *The Annals of Mathematical Statistics*, pages 165–179. [40](#)
- Li, W., Lee, A., and Gregersen, P. K. (2009). Copy-number-variation and copy-number-alteration region detection by cumulative plots. *BMC bioinformatics*, 10(1):S67. [17](#)
- Lin, S., Wang, C., Zarei, S., Bell, D. A., Kerr, S. E., Runger, G. C., and Kocher, J.-P. A. (2018). A data science approach for the classification of low-grade and high-grade ovarian serous carcinomas. *BMC genomics*, 19(1):841. [2](#)
- Liptak, T. (1958). On the combination of independent tests. *Magyar Tud Akad Mat Kutato Int Kozl*, 3:171–197. [148](#)
- Liu, Z., Li, A., Schulz, V., Chen, M., and Tuck, D. (2010). MixHMM: inferring copy number variation and allelic imbalance using SNP arrays and tumor samples mixed with stromal cells. *PloS one*, 5(6):e10909. [17](#)
- Loo, L. W., Wang, Y., Flynn, E. M., Lund, M. J., Bowles, E. J. A., Buist, D. S., Liff, J. M., Flagg, E. W., Coates, R. J., Eley, J. W., et al. (2011). Genome-wide copy number alterations in subtypes of invasive breast cancers in young white and African American women. *Breast cancer research and treatment*, 127(1):297–308. [1](#)
- Lu, X., Zhang, Q., Wang, Y., Zhang, L., Zhao, H., Chen, C., Wang, Y., Liu, S., Lu, T., Wang, F., et al. (2018). Molecular classification and subtype-specific characterization of skin cutaneous melanoma by aggregating multiple genomic platform data. *Journal of Cancer Research and Clinical Oncology*, pages 1–13. [2](#)

REFERENCES

- Malekpour, S. A., Pezeshk, H., and Sadeghi, M. (2018). MSeq-CNV: accurate detection of Copy Number Variation from Sequencing of Multiple samples. *Scientific reports*, 8(1):4009. [19](#)
- Manly, K. F., Nettleton, D., and Hwang, J. G. (2004). Genomics, prior probability, and statistical tests of multiple hypotheses. *Genome Research*, 14(6):997–1001. [122](#)
- Mann, H. B. and Whitney, D. R. (1947). On a test of whether one of two random variables is stochastically larger than the other. *The Annals of Mathematical Statistics*, 18:50–60. [46](#)
- Massey Jr, F. J. (1951). The Kolmogorov-Smirnov test for goodness of fit. *Journal of the American statistical Association*, 46(253):68–78. [37](#)
- Miller Jr, R. G. (1981). *Simultaneous statistical inference 2nd Edition*. Springer-Verlag. [121](#)
- Muzzey, D., Evans, E. A., and Lieber, C. (2015). Understanding the basics of NGS: from mechanism to variant calling. *Current genetic medicine reports*, 3(4):158–165. [6](#)
- Nguyen, N., Huang, H., Orintara, S., and Wang, Y. (2007). Denoising of array-based dna copy number data using the dual-tree complex wavelet transform. In *2007 IEEE 7th International Symposium on BioInformatics and BioEngineering*, pages 137–144. IEEE. [18](#)
- Olshen, A. B., Venkatraman, E., Lucito, R., and Wigler, M. (2004). Circular binary segmentation for the analysis of array-based DNA copy number data. *Biostatistics*, 5(4):557–572. [14](#), [98](#)
- Ostroverkhova, N., Nazarenko, S., and Cheremnykh, A. (2002). Comparative genomic hybridization as a new method for detection of genomic imbalance. *Russian Journal of Genetics*, 38(2):95–104. [5](#)
- Parzen, E. (1962). On estimation of a probability density function and mode. *The annals of mathematical statistics*, 33(3):1065–1076. [25](#)
- Patterson, N., Price, A. L., and Reich, D. (2006). Population structure and eigenanalysis. *PLoS genetics*, 2(12):e190. [122](#)
- Percival, D. B. and Walden, A. T. (2006). *Wavelet methods for time series analysis*, volume 4. Cambridge university press. [18](#)

-
- Pettitt, A. N. (1976). A two-sample Anderson-Darling rank statistic. *Biometrika*, 63(1):161–168. [40](#), [49](#)
- Picard, F., Robin, S., Lavielle, M., Vaisse, C., and Daudin, J.-J. (2005). A statistical approach for array cgh data analysis. *BMC bioinformatics*, 6(1):27. [17](#)
- Pickands III, J. (1975). Statistical inference using extreme order statistics. *the Annals of Statistics*, pages 119–131. [175](#)
- Pinkel, D., Segraves, R., Sudar, D., Clark, S., Poole, I., Kowbel, D., Collins, C., Kuo, W.-L., Chen, C., Zhai, Y., et al. (1998). High resolution analysis of DNA copy number variation using comparative genomic hybridization to microarrays. *Nature genetics*, 20(2):207. [5](#)
- Pollack, J. R., Sørlie, T., Perou, C. M., Rees, C. A., Jeffrey, S. S., Lonning, P. E., Tibshirani, R., Botstein, D., Børresen-Dale, A.-L., and Brown, P. O. (2002). Microarray analysis reveals a major direct role of DNA copy number alteration in the transcriptional program of human breast tumors. *Proceedings of the National Academy of Sciences*, 99(20):12963–12968. [4](#)
- Rabiner, L. R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286. [102](#)
- Rosenblatt, M. et al. (1952). Limit theorems associated with variants of the von Mises statistic. *The Annals of Mathematical Statistics*, 23(4):617–623. [40](#)
- Rueda, O. M. and Diaz-Uriarte, R. (2009). Detection of recurrent copy number alterations in the genome: taking among-subject heterogeneity seriously. *Bmc Bioinformatics*, 10(1):308. [19](#), [29](#), [30](#)
- Rueda, O. M. and Diaz-Uriarte, R. (2010). Finding recurrent copy number alteration regions: a review of methods. *Current Bioinformatics*, 5(1):1–17. [18](#)
- Schildkraut, C. L., Marmur, J., and Doty, P. (1961). The formation of hybrid DNA molecules and their use in studies of DNA homologies. *Journal of molecular biology*, 3(5):595–IN16. [4](#)
- Schuster, S. C. (2007). Next-generation sequencing transforms today’s biology. *Nature methods*, 5(1):16. [5](#)
- Shah, S. P., Lam, W. L., Ng, R. T., and Murphy, K. P. (2007). Modeling recurrent DNA copy number alterations in array CGH data. *Bioinformatics*, 23(13):i450–i458. [19](#)

REFERENCES

- Šidák, Z. (1967). Rectangular confidence regions for the means of multivariate normal distributions. *Journal of the American Statistical Association*, 62(318):626–633. [121](#), [122](#)
- Singleton, A., Farrer, M., Johnson, J., Singleton, A., Hague, S., Kachergus, J., Hulihan, M., Peuralinna, T., Dutra, A., Nussbaum, R., et al. (2003). α -Synuclein locus triplication causes Parkinson’s disease. *Science*, 302(5646):841–841. [4](#)
- Smeets, S. J., Braakhuis, B. J., Abbas, S., Snijders, P. J., Ylstra, B., van de Wiel, M. A., Meijer, G. A., Leemans, C. R., and Brakenhoff, R. H. (2006). Genome-wide DNA copy number alterations in head and neck squamous cell carcinomas with or without oncogene-expressing human papillomavirus. *Oncogene*, 25(17):2558. [19](#), [20](#), [35](#), [45](#), [154](#)
- Smirnov, N. V. (1936). Sur la distribution de w^2 . *Comp. Rend. Acad. Sci*, 202:449–452. [39](#)
- Smirnov, N. V. (1939). Estimate of deviation between empirical distribution functions in two independent samples. *Bulletin Moscow University*, 2(2):3–16. [37](#)
- Smith, J. C. and Sheltzer, J. M. (2018). Systematic identification of mutations and copy number alterations associated with cancer patient prognosis. *Elife*, 7:e39217. [1](#)
- Steiger, J. H. (1980). Testing pattern hypotheses on correlation matrices: Alternative statistics and some empirical results. *Multivariate Behavioral Research*, 15(3):335–352. [107](#)
- Sun, W., Wright, F. A., Tang, Z., Nordgard, S. H., Loo, P. V., Yu, T., Kristensen, V. N., and Perou, C. M. (2009). Integrated study of copy number states and genotype calls using high-density SNP arrays. *Nucleic acids research*, 37(16):5365–5377. [16](#), [17](#)
- Tang, Y.-C. and Amon, A. (2013). Gene copy-number alterations: a cost-benefit analysis. *Cell*, 152(3):394–405. [4](#)
- Theisen, A. (2008). Microarray-based comparative genomic hybridization (aCGH). *Nature Education*, 1(1):45. [4](#), [5](#)
- Tippett, L. H. C. et al. (1931). The Methods of Statistics. an introduction mainly for workers in the biological sciences. *The Methods of Statistics. An introduction mainly for workers in the biological sciences*. [148](#)

- van Boerdonk, R. A., Sutedja, T. G., Snijders, P. J., Reinen, E., Wilting, S. M., van de Wiel, M. A., Thunnissen, F. B., Duin, S., Kooi, C., Ylstra, B., et al. (2011). DNA copy number alterations in endobronchial squamous metaplastic lesions predict lung cancer. *American journal of respiratory and critical care medicine*, 184(8):948–956. [29](#), [31](#), [93](#)
- Van De Wiel, M. A. and Van Wieringen, W. N. (2007). CGHregions: dimension reduction for array CGH data with minimal information loss. *Cancer informatics*, 3:117693510700300031. [19](#), [20](#), [29](#), [30](#), [35](#), [45](#), [154](#)
- Van Loo, P., Nordgard, S. H., Lingjærde, O. C., Russnes, H. G., Rye, I. H., Sun, W., Weigman, V. J., Marynen, P., Zetterberg, A., Naume, B., et al. (2010). Allele-specific copy number analysis of tumors. *Proceedings of the National Academy of Sciences*, 107(39):16910–16915. [17](#)
- Von Mises, R. (1931). *Wahrscheinlichkeitsrechnung und ihre anwendung in der statistik und theorestischen physik*. Franz Deuticke. [39](#)
- Wang, H., Liang, L., Fang, J.-Y., and Xu, J. (2016). Somatic gene copy number alterations in colorectal cancer: new quest for cancer drivers and biomarkers. *Oncogene*, 35(16):2011. [1](#)
- Wang, H., Yan, W., Zhang, S., Gu, Y., Wang, Y., Wei, Y., Liu, H., Wang, F., Wu, Q., and Zhang, Y. (2017). Survival differences of CIMP subtypes integrated with CNA information in human breast cancer. *Oncotarget*, 8(30):48807. [2](#)
- Wang, J., Qian, J., Hoeksema, M. D., Zou, Y., Espinosa, A. V., Rahman, S. J., Zhang, B., and Massion, P. P. (2013). Integrative genomics analysis identifies candidate drivers at 3q26-29 amplicon in squamous cell carcinoma of the lung. *Clinical Cancer Research*, 19(20):5580–5590. [29](#), [31](#), [93](#)
- Wang, K., Li, M., Hadley, D., Liu, R., Glessner, J., Grant, S. F., Hakonarson, H., and Bucan, M. (2007). PennCNV: an integrated hidden Markov model designed for high-resolution copy number variation detection in whole-genome SNP genotyping data. *Genome research*, 17(11):1665–1674. [16](#)
- Wang, Y. and Wang, S. (2007). A novel stationary wavelet denoising algorithm for array-based dna copy number data. *International journal of bioinformatics research and applications*, 3(2):206–222. [18](#)
- Weiss, M., Hermsen, M., Meijer, G., Van Grieken, N., Baak, J., Kuipers, E., and Van Diest, P. (1999). Comparative genomic hybridisation. *Molecular Pathology*, 52(5):243. [4](#)

REFERENCES

- Wilcoxon, F. (1945). Individual comparisons by ranking methods. *Biometrics bulletin*, 1(6):80–83. [45](#)
- Wilting, S. M., Snijders, P. J., Meijer, G. A., Ylstra, B., van den IJssel, P. R., Snijders, A. M., Albertson, D. G., Coffa, J., Schouten, J. P., van de Wiel, M. A., et al. (2006). Increased gene copy numbers at chromosome 20q are frequent in both squamous cell carcinomas and adenocarcinomas of the cervix. *The Journal of Pathology: A Journal of the Pathological Society of Great Britain and Ireland*, 209(2):220–230. [19](#), [20](#), [35](#), [45](#), [154](#)
- Wolf, M., Mousses, S., Hautaniemi, S., Karhu, R., Huusko, P., Allinen, M., Elkahloun, A., Monni, O., Chen, Y., Kallioniemi, A., et al. (2004). High-resolution analysis of gene copy number alterations in human prostate cancer using CGH on cDNA microarrays: impact of copy number on gene expression. *Neoplasia*, 6(3):240–247. [4](#)
- Wood, H. M., Belvedere, O., Conway, C., Daly, C., Chalkley, R., Bickerdike, M., McKinley, C., Egan, P., Ross, L., Hayward, B., et al. (2010). Using next-generation sequencing for high resolution multiplex analysis of copy number variation from nanogram quantities of DNA from formalin-fixed paraffin-embedded specimens. *Nucleic acids research*, 38(14):e151–e151. [6](#)
- Wu, H.-T., Hajirasouliha, I., and Raphael, B. J. (2014). Detecting independent and recurrent copy number aberrations using interval graphs. *Bioinformatics*, 30(12):i195–i203. [3](#)
- Xie, C. and Tammi, M. T. (2009). CNV-seq, a new method to detect copy number variation using high-throughput sequencing. *BMC bioinformatics*, 10(1):80. [6](#), [13](#)
- Yoon, S., Xuan, Z., Makarov, V., Ye, K., and Sebat, J. (2009). Sensitive and accurate detection of copy number variants using read depth of coverage. *Genome research*. [6](#)