

BOOTSTRAPPING ANALOGS OF THE ONE WAY MANOVA TEST

by

Hasthika S. Rupasinghe Arachchige Don
M.S., Southern Illinois University, 2013

A Dissertation

Submitted in Partial Fulfillment of the Requirements for the
Doctoral Degree

Department of Mathematics
in the Graduate School
Southern Illinois University Carbondale
July, 2017

AN ABSTRACT OF THE DISSERTATION OF

HASTHIKA S. RUPASINGHE ARACHCHIGE DON, for the Doctor of Philosophy degree in MATHEMATICS, presented on DATE OF DEFENSE, at Southern Illinois University Carbondale.

TITLE: BOOTSTRAPPING ANALOGS OF THE ONE WAY MANOVA TEST

MAJOR PROFESSOR: Dr. D. J. Olive

The classical one way MANOVA model is used to test whether the mean measurements are the same or differ across p groups, and assumes that the covariance matrix of each group is the same. This work suggests using the Olive (2017abc) bootstrap technique to develop analogs of one way MANOVA test. A large sample theory test has also been developed. The bootstrap tests can have considerable outlier resistance, and the tests do not need the population covariance matrices to be equal. The two sample Hotelling's T^2 test is the special case of the one way MANOVA model when $p = 2$.

TABLE OF CONTENTS

Abstract	i
List of Tables	iv
List of Figures	vi
1 Introduction	1
1.1 MANOVA	1
1.2 One Way MANOVA	3
1.3 Two Group case	4
1.3.1 Two Sample Hotelling's T^2 Test	4
2 Theory and Methods	7
2.1 Notation	7
2.1.1 Mahalanobis Distance	7
2.2 Prediction Region	7
2.3 Prediction region method	8
2.3.1 Testing $H_0 : \boldsymbol{\mu} = \mathbf{c}$ versus $H_1 : \boldsymbol{\mu} \neq \mathbf{c}$ Using the Prediction Region Method	9
2.4 A relationship between the one-way MANOVA test and the Hotelling Lawley trace test	9
2.5 Cell means model	16
3 Bootstrapping Analogs of the Two Sample Hotelling's T^2 Test	18
3.1 Applying the prediction region method to the two sample test	18
3.2 Real Data Example	19
4 Bootstrapping Analogs of the One Way MANOVA Test	21
4.1 An Alternative to the usual one way MANOVA	22
4.1.1 Test H_0 when $\hat{\boldsymbol{\Sigma}}\mathbf{w}$ is unknown or difficult to estimate.	26
4.2 Power comparison among the tests	27
4.3 Real data example	27

5	Simulations for Bootstrapping Analogs of the Two Sample Hotelling's T^2 Test	35
5.1	Simulation Setup	35
5.2	Simulation Output	36
5.2.1	Type I error rates simulation for clean data	36
5.2.2	Type I error rates simulation for contaminated data	45
5.2.3	Power Simulation	47
6	Simulations with three samples for Bootstrapping Analogs of the One-Way MANOVA Test	50
6.1	Simulation Setup	50
6.1.1	Simulations for type I error with clean data	51
6.1.2	Simulations for power with clean data	58
6.1.3	Simulations for type I error with contaminated data	63
7	Discussion	71
	Appendix	72
	References	90
	Vita	93

LIST OF TABLES

5.1	Coverages for clean multivariate normal data $p = 5$	37
5.2	Coverages for clean multivariate normal data $p = 15$	38
5.3	Coverages for clean $0.6N_p(\mathbf{0}, \mathbf{I}) + 0.4N_p(\mathbf{0}, 25\mathbf{I})$ data $p = 5$	39
5.4	Coverages for clean $0.6N_p(\mathbf{0}, \mathbf{I}) + 0.4N_p(\mathbf{0}, 25\mathbf{I})$ data $p = 15$	40
5.5	Coverages for clean multivariate t_4 data $p = 5$	41
5.6	Coverages for clean multivariate t_4 data $p = 15$	42
5.7	Coverages for clean lognormal data $p = 5$	43
5.8	Coverages for clean lognormal data $p = 15$	44
5.9	Coverages and cutoffs with outliers: $p = 4, n_1 = n_2 = 200, B = 200$	46
5.10	Coverages when H_0 is false for MVN data.	47
5.11	Coverages when H_0 is false for mixture data.	48
5.12	Coverages when H_0 is false for multivariate t_4 data.	48
5.13	Coverages when H_0 is false for lognormal data.	49
6.1	Type I error for clean MVN data with $\text{cov3I} = \text{F}$	52
6.2	Type I error for clean Mixture data with $\text{cov3I} = \text{F}$	53
6.3	Type I error for clean multivariate t data with $\text{cov3I} = \text{F}$	54
6.4	Type I error for clean lognormal data with $\text{cov3I} = \text{F}$	55
6.5	Type I error for clean MVN data with $\text{cov3I} = \text{T}$	56
6.6	Type I error for clean Mixture data with $\text{cov3I} = \text{T}$	56
6.7	Type I error for clean Multivariate t data with $\text{cov3I} = \text{T}$	57
6.8	Type I error for clean lognormal data with $\text{cov3I} = \text{T}$	57
6.9	Power for MVN data with $\delta_1 = 0.2$ and $\delta_3 = 0.5$	59
6.10	Power for Mixture data $\delta_1 = 0.2$ and $\delta_3 = 0.5$	60
6.11	Power for Multivariate t data $\delta_1 = 0.2$ and $\delta_3 = 0.5$	61
6.12	Power for lognormal data $\delta_1 = 0.2$ and $\delta_3 = 0.5$	62
6.13	Type I error with contaminated data: $m = 5, \gamma = 0.1$	64
6.14	Type I error with contaminated data: $m = 10, \gamma = 0.05$	65

6.15	Type I error with contaminated data: $m = 20, \gamma = 0.1$	66
6.16	Type I error with contaminated data: $m = 20, \gamma = 0.05$	67
6.17	Power curve for MVN data with $m = 5, \sigma_1 = 1, \sigma_2 = 1, \sigma_3 = 1, n_1 = 200, n_2 =$ 200 and $n_3 = 200$	68
6.18	Power curve for MVN data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 =$ 400 and $n_3 = 600$	68
6.19	Power curve for Mixture data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 =$ 400 and $n_3 = 600$	69
6.20	Power curve for Multivariate t_4 data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 =$ 200, $n_2 = 400$ and $n_3 = 600$	70

LIST OF FIGURES

3.1	DD plot for the age ≤ 50 group.	20
3.2	DD plot for the age > 50 group.	20
4.1	Power curve for clean MVN data with $m = 5, \sigma_1 = 1, \sigma_2 = 1, \sigma_3 = 1, n_1 = 200, n_2 = 200$ and $n_3 = 200$	28
4.2	Power curve for clean MVN data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 = 400$ and $n_3 = 600$	29
4.3	Power curve for clean Mixture data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 = 400$ and $n_3 = 600$	30
4.4	Power curve for clean multivariate t_4 data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 = 400$ and $n_3 = 600$	31
4.5	DD plots for Crime data	32
4.6	Side by side boxplots for Crime data	33
4.7	Scatterplot matrix for Crime data	34

CHAPTER 1

INTRODUCTION

1.1 MANOVA

Multivariate analysis of variance (MANOVA) is analogous to an ANOVA with more than one dependent variable. ANOVA tests for the difference in means between two or more groups, while MANOVA tests for the difference in two or more vectors of means.

Definition. The multivariate analysis of variance (MANOVA) model $\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i + \boldsymbol{\epsilon}_i$ for $i = 1, \dots, n$ has $m \geq 2$ response variables Y_1, \dots, Y_m and p predictor variables x_1, x_2, \dots, x_p . The i th case is $(\mathbf{x}_i^T, \mathbf{y}_i^T) = (x_{i1}, \dots, x_{ip}, Y_{i1}, \dots, Y_{im})$. If a constant $x_{i1} = 1$ is in the model, then x_{i1} could be omitted from the case.

For the MANOVA model predictors are indicator variables. Sometimes the trivial predictor $\mathbf{1}$ is also in the model.

The MANOVA model in matrix form is $\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{E}$ and has $E(\boldsymbol{\epsilon}_k) = \mathbf{0}$ and $Cov(\boldsymbol{\epsilon}_k) = \boldsymbol{\Sigma}_\epsilon = (\sigma_{ij})$ for $k = 1, \dots, n$. Also $E(\mathbf{e}_i) = \mathbf{0}$ while $Cov(\mathbf{e}_i, \mathbf{e}_j) = \sigma_{ij}\mathbf{I}_n$ for $i, j = 1, \dots, m$. Then \mathbf{B} and $\boldsymbol{\Sigma}_\epsilon$ are unknown matrices of parameters to be estimated.

$$\mathbf{Z} = \begin{pmatrix} Y_{1,1} & Y_{1,2} & \cdots & Y_{1,m} \\ Y_{2,1} & Y_{2,2} & \cdots & Y_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ Y_{n,1} & Y_{n,2} & \cdots & Y_{n,m} \end{pmatrix} = \begin{pmatrix} \mathbf{Y}_1 & \mathbf{Y}_2 & \cdots & \mathbf{Y}_m \end{pmatrix} = \begin{pmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_n^T \end{pmatrix}$$

The $n \times p$ matrix \mathbf{X} is not necessarily of full rank p , and

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,p} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,p} \end{pmatrix} = \begin{pmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_p \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix}$$

where often $\mathbf{v}_1 = \mathbf{1}$

The $p \times m$ matrix

$$\mathbf{B} = \begin{pmatrix} \beta_{1,1} & \beta_{1,2} & \cdots & \beta_{1,m} \\ \beta_{2,1} & \beta_{2,2} & \cdots & \beta_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{p,1} & \beta_{p,2} & \cdots & \beta_{p,m} \end{pmatrix} = \left(\boldsymbol{\beta}_1 \quad \boldsymbol{\beta}_2 \quad \cdots \quad \boldsymbol{\beta}_m \right).$$

The $n \times m$ matrix

$$\mathbf{E} = \begin{pmatrix} \epsilon_{1,1} & \epsilon_{1,2} & \cdots & \epsilon_{1,m} \\ \epsilon_{2,1} & \epsilon_{2,2} & \cdots & \epsilon_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ \epsilon_{n,1} & \epsilon_{n,2} & \cdots & \epsilon_{n,m} \end{pmatrix} = \left(\mathbf{e}_1 \quad \mathbf{e}_2 \quad \cdots \quad \mathbf{e}_m \right) = \begin{pmatrix} \boldsymbol{\epsilon}_1^T \\ \vdots \\ \boldsymbol{\epsilon}_n^T \end{pmatrix}.$$

Each response variable in a MANOVA model follows an ANOVA model $\mathbf{Y}_j = \mathbf{X}\boldsymbol{\beta}_j + \mathbf{e}_j$ for $j = 1, \dots, m$ where it is assumed that $E(\mathbf{e}_i) = \mathbf{0}$ and $Cov(\mathbf{e}_j) = \sigma_{jj}\mathbf{I}_n$.

MANOVA models are often fit by least squares. The least squares estimators $\hat{\mathbf{B}}$ of \mathbf{B} are

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{Z} = \left(\hat{\boldsymbol{\beta}}_1 \quad \hat{\boldsymbol{\beta}}_2 \quad \cdots \quad \hat{\boldsymbol{\beta}}_m \right)$$

where $(\mathbf{X}^T \mathbf{X})^{-}$ is a generalized inverse of $\mathbf{X}^T \mathbf{X}$. If \mathbf{X} has a full rank then $(\mathbf{X}^T \mathbf{X})^{-} = (\mathbf{X}^T \mathbf{X})^{-1}$ and $\hat{\mathbf{B}}$ is unique.

Definition. The predicted values or fitted values

$$\hat{\mathbf{Z}} = \mathbf{X}\hat{\mathbf{B}} = \left(\hat{\mathbf{Y}}_1 \quad \hat{\mathbf{Y}}_2 \quad \cdots \quad \hat{\mathbf{Y}}_m \right) = \begin{pmatrix} \hat{Y}_{1,1} & \hat{Y}_{1,2} & \cdots & \hat{Y}_{1,m} \\ \hat{Y}_{2,1} & \hat{Y}_{2,2} & \cdots & \hat{Y}_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{Y}_{n,1} & \hat{Y}_{n,2} & \cdots & \hat{Y}_{n,m} \end{pmatrix}.$$

The residuals $\hat{\mathbf{E}} = \mathbf{Z} - \hat{\mathbf{Z}} = \mathbf{Z} - \mathbf{X}\hat{\mathbf{B}}$.

Finally,

$$\hat{\Sigma}_{\boldsymbol{\epsilon}} = \frac{(\mathbf{Z} - \hat{\mathbf{Z}})^T (\mathbf{Z} - \hat{\mathbf{Z}})}{n - p} = \frac{\hat{\mathbf{E}}^T \hat{\mathbf{E}}}{n - p}.$$

1.2 ONE WAY MANOVA

Assume that there are independent random samples of size n_i from p different populations, or n_i cases are randomly assigned to p treatment groups. Let $n = \sum_{i=1}^p n_i$ be the total sample size. Also assume that m response variables $\mathbf{y}_{ij} = (Y_{ij1}, \dots, Y_{ijm})^T$ are measured for i th treatment group and j th case. Assume $E(\mathbf{y}_{ij}) = \boldsymbol{\mu}_i$ and $Cov(\mathbf{y}_{ij}) = \boldsymbol{\Sigma}_\epsilon$.

The one way MANOVA is used to test $H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \dots = \boldsymbol{\mu}_p$. Note that if $m = 1$ the one way MANOVA model becomes the one way ANOVA model. One might think that performing m ANOVA tests is sufficient to test the above hypotheses. But the separate ANOVA tests would not take the correlation between the m variables into account. On the other hand the MANOVA test will take the correlation into account.

Let $\bar{\mathbf{y}} = \sum_{i=1}^p \sum_{j=1}^{n_i} \mathbf{y}_{ij} / n$ be the overall mean. Let $\bar{\mathbf{y}}_i = \sum_{j=1}^{n_i} \mathbf{y}_{ij} / n_i$. Several $m \times m$ matrices will be useful. Let \mathbf{S}_i be the sample covariance matrix corresponding to the i th treatment group. Then the within sum of squares and cross products matrix is $\mathbf{W} = (n_1 - 1)\mathbf{S}_1 + \dots + (n_p - 1)\mathbf{S}_p = \sum_{i=1}^p \sum_{j=1}^{n_i} (\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)(\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)^T$. Then $\hat{\boldsymbol{\Sigma}}_\epsilon = \mathbf{W} / (n - p)$. The treatment or between sum of squares and cross products matrix is

$$\mathbf{B}_T = \sum_{i=1}^p n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}})(\bar{\mathbf{y}}_i - \bar{\mathbf{y}})^T.$$

The total corrected (for the mean) sum of squares and cross products matrix is $\mathbf{T} = \mathbf{B}_T + \mathbf{W} = \sum_{i=1}^p \sum_{j=1}^{n_i} (\mathbf{y}_{ij} - \bar{\mathbf{y}})(\mathbf{y}_{ij} - \bar{\mathbf{y}})^T$. Note that $\mathbf{S} = \mathbf{T} / (n - 1)$ is the usual sample covariance matrix of the \mathbf{y}_{ij} if it is assumed that all n of the \mathbf{y}_{ij} are iid so that the $\boldsymbol{\mu}_i \equiv \boldsymbol{\mu}$ for $i = 1, \dots, p$.

The one way MANOVA model is $\mathbf{y}_{ij} = \boldsymbol{\mu}_i + \boldsymbol{\epsilon}_{ij}$ where the $\boldsymbol{\epsilon}_{ij}$ are iid with $E(\boldsymbol{\epsilon}_{ij}) = \mathbf{0}$ and $Cov(\boldsymbol{\epsilon}_{ij}) = \boldsymbol{\Sigma}_\epsilon$. The summary one Way MANOVA table is shown bellow.

Source	matrix	df
Treatment or Between	\mathbf{B}_T	$p - 1$
Residual or Error or Within	\mathbf{W}	$n - p$
Total (Corrected)	\mathbf{T}	$n - 1$

There are three commonly used test statistics to test the above hypotheses. Namely,

1. Hotelling Lawley trace statistic: $U = tr(\mathbf{B}_T \mathbf{W}^{-1}) = tr(\mathbf{W}^{-1} \mathbf{B}_T)$.

2. Wilks' lambda: $\Lambda = \frac{|\mathbf{W}|}{|\mathbf{B}_T + \mathbf{W}|}$.
3. Pillai's trace statistic: $\mathbf{V} = \text{tr}(\mathbf{B}_T \mathbf{T}^{-1}) = \text{tr}(\mathbf{T}^{-1} \mathbf{B}_T)$.

If the $\mathbf{y}_{ij} - \boldsymbol{\mu}_j$ are iid with common covariance matrix $\boldsymbol{\Sigma}_\epsilon$, and if H_0 is true, then under regularity conditions Fujikoshi (2002) showed

1. $(n - m - p - 1)U \xrightarrow{D} \chi_{m(p-1)}^2$,
2. $-[n - 0.5(m + p - 2)]\log(\Lambda) \xrightarrow{D} \chi_{m(p-1)}^2$, and
3. $(n - 1)V \xrightarrow{D} \chi_{m(p-1)}^2$.

Note that the common covariance matrix assumption implies that each of the p treatment groups or populations has the same covariance matrix $\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}_\epsilon$ for $i = 1, \dots, p$, an extremely strong assumption. Kakizawa (2009) and Olive, Pelawa Watagoda, and Rupasinghe Arachchige Don (2015) show that similar results hold for the multivariate linear model. The common covariance matrix assumption, $\text{Cov}(\epsilon_k) = \boldsymbol{\Sigma}_\epsilon$ for $k = 1, \dots, n$, is often reasonable for the multivariate linear regression model.

1.3 TWO GROUP CASE

Suppose there are two independent random samples from two populations or groups. A common multivariate two sample test of hypotheses is $H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$ versus $H_1 : \boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2$ where $\boldsymbol{\mu}_i$ is a population location measure of the i th population for $i = 1, 2$. The two sample Hotelling's T^2 test is the classical method, and is a special case of the one way MANOVA model if the two populations are assumed to have the same population covariance matrix.

1.3.1 Two Sample Hotelling's T^2 Test

Suppose there are two independent random samples $\mathbf{x}_{1,1}, \dots, \mathbf{x}_{n_1,1}$ and $\mathbf{x}_{1,2}, \dots, \mathbf{x}_{n_2,2}$ from two populations or groups, and that it is desired to test $H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$ versus $H_1 : \boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2$ where the $\boldsymbol{\mu}_i$ are $p \times 1$ vectors. Assume that T_i satisfy a central limit type theorem $\sqrt{n}(T_i - \boldsymbol{\mu}_i) \xrightarrow{D} N_p(\mathbf{0}, \boldsymbol{\Sigma}_i)$ for $i = 1, 2$ where the $\boldsymbol{\Sigma}_i$ are positive definite.

To simplify large sample theory, assume $n_1 = kn_2$ for some positive real number k .

Let $\hat{\Sigma}_i$ be a consistent nonsingular estimator of Σ_i . Then

$$\begin{pmatrix} \sqrt{n_1} (T_1 - \boldsymbol{\mu}_1) \\ \sqrt{n_2} (T_2 - \boldsymbol{\mu}_2) \end{pmatrix} \xrightarrow{D} N_{2p} \left[\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \Sigma_2 \end{pmatrix} \right],$$

or

$$\begin{pmatrix} \sqrt{n_2} (T_1 - \boldsymbol{\mu}_1) \\ \sqrt{n_2} (T_2 - \boldsymbol{\mu}_2) \end{pmatrix} \xrightarrow{D} N_{2p} \left[\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \frac{\Sigma_1}{k} & \mathbf{0} \\ \mathbf{0} & \Sigma_2 \end{pmatrix} \right].$$

Hence

$$\sqrt{n_2} [(T_1 - T_2) - (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)] \xrightarrow{D} N_p \left(\mathbf{0}, \frac{\Sigma_1}{k} + \Sigma_2 \right).$$

Using $n\mathbf{B}^{-1} = \left(\frac{\mathbf{B}}{n}\right)^{-1}$ and $n_2k = n_1$, if $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$, then

$$\begin{aligned} n_2(T_1 - T_2)^T \left(\frac{\Sigma_1}{k} + \Sigma_2 \right)^{-1} (T_1 - T_2) &= \\ (T_1 - T_2)^T \left(\frac{\Sigma_1}{n_1} + \frac{\Sigma_2}{n_2} \right)^{-1} (T_1 - T_2) &\xrightarrow{D} \chi_p^2. \end{aligned}$$

Hence

$$T_0^2 = (T_1 - T_2)^T \left(\frac{\hat{\Sigma}_1}{n_1} + \frac{\hat{\Sigma}_2}{n_2} \right)^{-1} (T_1 - T_2) \xrightarrow{D} \chi_p^2. \quad (1.1)$$

Note that k drops out of the above result.

If the sequence of positive integers $d_n \rightarrow \infty$ and $Y_n \sim F_{p,d_n}$, then $Y_n \xrightarrow{D} \chi_p^2/p$. Using an F_{p,d_n} distribution instead of a χ_p^2 distribution is similar to using a t_{d_n} distribution instead of a standard normal $N(0,1)$ distribution for inference. Instead of rejecting H_0 when $T_0^2 > \chi_{p,1-\delta}^2$, reject H_0 when

$$T_0^2 > pF_{p,d_n,1-\delta} = \frac{pF_{p,d_n,1-\delta}}{\chi_{p,1-\delta}^2} \chi_{p,1-\delta}^2.$$

The term $\frac{pF_{p,d_n,1-\delta}}{\chi_{p,1-\delta}^2}$ can be regarded as a small sample correction factor that improves the test's performance for small samples. For example, use $d_n = \min(n_1 - p, n_2 - p)$. Here $P(Y_n \leq \chi_{p,\delta}^2) = \delta$ if Y_n has a χ_p^2 distribution, and $P(Y_n \leq F_{p,d_n,\delta}) = \delta$ if Y_n has an F_{p,d_n} distribution.

The two sample Hotelling's T^2 test is the classical method. If it is not assumed that the population covariance matrices are equal, then this test uses the sample mean and sample covariance matrix $T_i = \bar{\mathbf{x}}_i$ and $\hat{\Sigma}_i = \mathbf{S}_i$ applied to each sample. This test has considerable robustness to the assumption that both populations have a multivariate normal distribution and to the assumption that the populations have a common population covariance matrix Σ , but the test can be very poor if outliers are present.

Alternative statistics to the sample mean can be useful, but large sample tests of the form of (1.1) need practical consistent estimators $\hat{\Sigma}_i$ of the two asymptotic covariance matrices Σ_i .

Chapter 2 gives theory and methods for bootstrapping hypotheses tests and shows how to apply the bootstrap to test the hypothesis $H_0 : \boldsymbol{\mu} = \mathbf{c}$ versus $H_1 : \boldsymbol{\mu} \neq \mathbf{c}$. Chapter 3 suggests using the Olive (2017abc) bootstrap technique to develop analogs of the Hotelling's T^2 test that use a statistic T_i , such as the coordinatewise median, applied to the i th sample for $i = 1, 2$. These tests are useful if the asymptotic covariance matrix is unknown or difficult to estimate. The new tests can have considerable outlier resistance, and the tests do not need the population covariance matrices to be equal. Chapter 4 suggests using the Olive (2017abc) bootstrap technique to develop analogs of the one way MANOVA test. The new tests can have considerable outlier resistance, and the tests do not need the population covariance matrices to be equal. Chapters 5 and 6 give some simulations and examples.

CHAPTER 2

THEORY AND METHODS

2.1 NOTATION

2.1.1 Mahalanobis Distance

Let the $p \times 1$ column vector T be a multivariate location estimator, and let the $p \times p$ symmetric positive definite matrix \mathbf{C} be a dispersion estimator. Then the i th *squared sample Mahalanobis distance* is the scalar

$$D_i^2 = D_i^2(T, \mathbf{C}) = D_{\mathbf{x}_i}^2(T, \mathbf{C}) = (\mathbf{x}_i - T)^T \mathbf{C}^{-1} (\mathbf{x}_i - T) \quad (2.1)$$

for each observation \mathbf{x}_i . Notice that the Euclidean distance of \mathbf{x}_i from the estimate of center T is $D_i(T, \mathbf{I}_p)$ where \mathbf{I}_p is the $p \times p$ identity matrix. The classical Mahalanobis distance uses $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$, the sample mean and sample covariance matrix where

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \text{and} \quad \mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T. \quad (2.2)$$

2.2 PREDICTION REGION

A large sample $100(1 - \delta)\%$ prediction region is the hyperellipsoid

$$\{\mathbf{w} : D_{\mathbf{w}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(c)}^2\} = \{\mathbf{w} : D_{\mathbf{w}}(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(c)}\} \quad (2.3)$$

for appropriate c . Using $c = \lceil n(1 - \delta) \rceil$ covers about $100(1 - \delta)\%$ of the training data cases \mathbf{x}_i , but the prediction region will have coverage lower than the nominal coverage of $1 - \delta$ for moderate n . This result is not surprising since empirically statistical methods perform worse on test data. Increasing c will improve the coverage for moderate samples. Let $q_n = \min(1 - \delta + 0.05, 1 - \delta + p/n)$ for $\delta > 0.1$ and

$$q_n = \min(1 - \delta/2, 1 - \delta + 10\delta p/n), \quad \text{otherwise.} \quad (2.4)$$

If $1 - \delta < 0.999$ and $q_n < 1 - \delta + 0.001$, set $q_n = 1 - \delta$.

Let $D_{(U_n)}$ be the $100q_n$ th percentile of the D_i . Then the Olive (2013) large sample $100(1 - \delta)\%$ nonparametric prediction region for a future value \mathbf{x}_f given iid data $\mathbf{x}_1, \dots, \mathbf{x}_n$

is

$$\{\mathbf{w} : D_{\mathbf{w}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(U_n)}^2\}, \quad (2.5)$$

while the classical large sample $100(1 - \delta)\%$ prediction region is

$$\{\mathbf{w} : D_{\mathbf{w}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq \chi_{p,1-\delta}^2\}. \quad (2.6)$$

2.3 PREDICTION REGION METHOD

Olive (2017bc) shows that there is a useful relationship between prediction regions and confidence regions. Consider predicting a future $p \times 1$ test vector \mathbf{x}_f , given past training data $\mathbf{x}_1, \dots, \mathbf{x}_n$. A *large sample* $100(1 - \delta)\%$ *prediction region* is a set \mathcal{A}_n such that $P(\mathbf{x}_f \in \mathcal{A}_n) \rightarrow 1 - \delta$ while a large sample $100(1 - \delta)\%$ *confidence region* for a parameter $\boldsymbol{\mu}$ is a set \mathcal{A}_n such that $P(\boldsymbol{\mu} \in \mathcal{A}_n) \rightarrow 1 - \delta$ as $n \rightarrow \infty$. Consider testing $H_0 : \boldsymbol{\mu} = \mathbf{c}$ versus $H_1 : \boldsymbol{\mu} \neq \mathbf{c}$ where \mathbf{c} is a known $p \times 1$ vector.

The Olive (2017abc) prediction region method obtains a confidence region for $\boldsymbol{\mu}$ by applying the nonparametric prediction region (2.5) to the bootstrap sample T_1^*, \dots, T_B^* , and the theory for the method is sketched below. Let \bar{T}^* and \mathbf{S}_T^* be the sample mean and sample covariance matrix of the bootstrap sample. Following Bickel and Ren (2001), let the vector of parameters $\boldsymbol{\mu} = T(F)$, the statistic $T_n = T(F_n)$, and $T^* = T(F_n^*)$ where F is the cdf of iid $\mathbf{x}_1, \dots, \mathbf{x}_n$, F_n is the empirical cdf, and F_n^* is the empirical cdf of $\mathbf{x}_1^*, \dots, \mathbf{x}_n^*$, a sample from F_n using the nonparametric bootstrap. If $\sqrt{n}(F_n - F) \xrightarrow{D} \mathbf{z}_F$, a Gaussian random process, and if T is sufficiently smooth (Hadamard differentiable with a Hadamard derivative $\dot{T}(F)$), then $\sqrt{n}(T_n - \boldsymbol{\mu}) \xrightarrow{D} \mathbf{X}$ and $\sqrt{n}(T_i^* - \bar{T}^*) \xrightarrow{D} \mathbf{X}$ with $\mathbf{X} = \dot{T}(F)\mathbf{z}_F$. Olive (2017bc) uses these results to show that if $\mathbf{X} \sim N_p(\mathbf{0}, \boldsymbol{\Sigma}_T)$, then $\sqrt{n}(\bar{T}^* - T_n) \xrightarrow{D} \mathbf{0}$, $\sqrt{n}(\bar{T}^* - \boldsymbol{\mu}) \xrightarrow{D} \mathbf{X}$, and that the prediction region method large sample $100(1 - \delta)\%$ confidence region for $\boldsymbol{\mu}$ is

$$\{\mathbf{w} : (\mathbf{w} - \bar{T}^*)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - \bar{T}^*) \leq D_{(U_B)}^2\} = \{\mathbf{w} : D_{\mathbf{w}}^2(\bar{T}^*, \mathbf{S}_T^*) \leq D_{(U_B)}^2\} \quad (2.7)$$

where $D_{(U_B)}^2$ is computed from $D_i^2 = (T_i^* - \bar{T}^*)^T [\mathbf{S}_T^*]^{-1} (T_i^* - \bar{T}^*)$ for $i = 1, \dots, B$. Note that the corresponding test for $H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$ rejects H_0 if $(\bar{T}^* - \boldsymbol{\mu}_0)^T [\mathbf{S}_T^*]^{-1} (\bar{T}^* - \boldsymbol{\mu}_0) > D_{(U_B)}^2$. This procedure is basically the one sample Hotelling's T^2 test applied to the T_i^* using \mathbf{S}_T^* as the estimated covariance matrix and replacing the $\chi_{p,1-\delta}^2$ cutoff by $D_{(U_B)}^2$.

2.3.1 Testing $H_0 : \boldsymbol{\mu} = \mathbf{c}$ versus $H_1 : \boldsymbol{\mu} \neq \mathbf{c}$ Using the Prediction Region Method

The prediction region method for testing $H_0 : \boldsymbol{\mu} = \mathbf{c}$ versus $H_1 : \boldsymbol{\mu} \neq \mathbf{c}$ is simple. Let $\hat{\boldsymbol{\mu}}$ be a consistent estimator of $\boldsymbol{\mu}$ and make a bootstrap sample $\mathbf{w}_i = \hat{\boldsymbol{\mu}}_i^* - \mathbf{c}$ for $i = 1, \dots, B$. Make the nonparametric prediction region (2.7) for the \mathbf{w}_i and fail to reject H_0 if $\mathbf{0}$ is in the prediction region, reject H_0 otherwise.

The Bickel and Ren (2001) hypothesis testing method is equivalent to using confidence region (2.7) with \bar{T}^* replaced by T_n and U_B replaced by $\lceil B(1 - \delta) \rceil$. If region (2.7) or the Bickel and Ren (2001) region is a large sample $100(1 - \delta)\%$ confidence region, then so is the other region if $\sqrt{n}(\bar{T}^* - T_n) \xrightarrow{D} \mathbf{0}$. Hadamard differentiability and asymptotic normality are sufficient conditions for both regions to be large sample confidence regions if $\mathbf{S}_T^* \xrightarrow{D} \boldsymbol{\Sigma}_T$, but Bickel and Ren (2001) showed that their method can work when Hadamard differentiability fails.

The location model with means, medians, and trimmed means is one example where the Bickel and Ren (2001, p. 96) method works. Since the univariate sample mean, sample median, and sample trimmed mean are Hadamard differentiable and asymptotically normal, each coordinate satisfies $\sqrt{n}(T_{in} - \bar{T}_i^*) \xrightarrow{D} 0$ for $i = 1, \dots, p$. Hence $\sqrt{n}(T_n - \bar{T}^*) \xrightarrow{D} \mathbf{0}$, and (2.7) is a large sample $100(1 - \delta)\%$ confidence region if T_n is the coordinatewise sample mean, median, or trimmed mean.

Fréchet differentiability implies Hadamard differentiability, and many statistics are shown to be Hadamard differentiable in Bickel and Ren (2001), Clarke (1986, 2000), Fernholtz (1983), and Gill (1989). Also see Ren (1991) and Ren and Sen (1995).

2.4 A RELATIONSHIP BETWEEN THE ONE-WAY MANOVA TEST AND THE HOTELLING LAWLEY TRACE TEST

An alternative method for one way MANOVA is to use the model $\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{E}$ with \mathbf{X} , \mathbf{Z} and \mathbf{B} as follows.

Let

$$\mathbf{Y}_{ij} = \begin{pmatrix} Y_{ij1} \\ \vdots \\ Y_{ijm} \end{pmatrix} = \boldsymbol{\mu}_i + \mathbf{e}_{ij}, \quad E\mathbf{Y}_{ij} = \boldsymbol{\mu}_i = \begin{pmatrix} \mu_{ij1} \\ \vdots \\ \mu_{ijm} \end{pmatrix}$$

for $i = 1, \dots, p$ and $j = 1, \dots, n_i$

Then \mathbf{X} is a full rank where the i th column of \mathbf{X} is an indicator for group $i - 1$ for $i = 2, \dots, p$, and \mathbf{Z} are as follows.

$$\mathbf{Z} = \begin{pmatrix} \mathbf{Y}_{11}^T \\ \vdots \\ \mathbf{Y}_{1n_1}^T \\ \mathbf{Y}_{21}^T \\ \vdots \\ \mathbf{Y}_{2n_2}^T \\ \vdots \\ \mathbf{Y}_{p-1,1}^T \\ \vdots \\ \mathbf{Y}_{p-1,n_{p-1}}^T \\ \vdots \\ \mathbf{Y}_{p,1}^T \\ \vdots \\ \mathbf{Y}_{p,n_p}^T \end{pmatrix},$$

$$\mathbf{X} = \begin{pmatrix} 1 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 1 & 0 & \cdots & 0 \\ 1 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 0 & 0 & \cdots & 1 \\ 1 & 0 & 0 & \cdots & 1 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 0 & 0 & \cdots & 0 \end{pmatrix} \quad (2.8)$$

$$\mathbf{B} = \begin{pmatrix} \boldsymbol{\mu}_p^T \\ (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_p)^T \\ \vdots \\ (\boldsymbol{\mu}_{p-1} - \boldsymbol{\mu}_p)^T \end{pmatrix} \quad \text{and} \quad \mathbf{L} = \begin{pmatrix} \mathbf{0} & \mathbf{I}_{p-1} \end{pmatrix}. \quad \text{Note that } \mathbf{Y}_{ij}^T = \boldsymbol{\mu}_i^T + \mathbf{e}_{ij}^T.$$

Then

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} n & n_1 & n_2 & \cdots & n_{p-1} \\ n_1 & n_1 & 0 & \cdots & 0 \\ \vdots & & \ddots & & \vdots \\ n_{p-2} & 0 & \cdots & n_{p-2} & 0 \\ n_{p-1} & 0 & \cdots & 0 & n_{p-1} \end{pmatrix} \quad (2.9)$$

and

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{n_p} \begin{pmatrix} 1 & -1 & -1 & \cdots & -1 \\ -1 & 1 + \frac{n_p}{n_1} & 1 & \cdots & 1 \\ \vdots & & \ddots & & \vdots \\ -1 & 1 & \cdots & 1 + \frac{n_p}{n_{p-2}} & 1 \\ -1 & 1 & \cdots & 1 & 1 + \frac{n_p}{n_{p-1}} \end{pmatrix}. \quad (2.10)$$

Then the least square estimators $\hat{\mathbf{B}}$ of \mathbf{B} ,

$$\hat{\mathbf{B}} = \begin{pmatrix} \bar{\mathbf{y}}_p^T \\ (\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_p)^T \\ \vdots \\ (\bar{\mathbf{y}}_{p-1} - \bar{\mathbf{y}}_p)^T \end{pmatrix}, \quad \text{and} \quad \mathbf{L}\hat{\mathbf{B}} = \begin{pmatrix} (\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_p)^T \\ (\bar{\mathbf{y}}_2 - \bar{\mathbf{y}}_p)^T \\ \vdots \\ (\bar{\mathbf{y}}_{p-1} - \bar{\mathbf{y}}_p)^T \end{pmatrix}.$$

Then $\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T$ becomes

$$\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T = \frac{1}{n_p} \begin{pmatrix} 1 + \frac{n_p}{n_1} & 1 & 1 & \cdots & 1 \\ 1 & 1 + \frac{n_p}{n_2} & 1 & \cdots & 1 \\ \vdots & & \ddots & & \vdots \\ 1 & 1 & \cdots & 1 & 1 + \frac{n_p}{n_{p-1}} \end{pmatrix}. \quad (2.11)$$

It can be shown that the inverse of the above matrix is

$$[\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1} = \frac{1}{n} \begin{pmatrix} n_1(n - n_1) & -n_1 n_2 & -n_1 n_3 & \cdots & -n_1 n_{p-1} \\ -n_1 n_2 & n_2(n - n_2) & -n_2 n_3 & \cdots & -n_2 n_{p-1} \\ \vdots & & \ddots & & \vdots \\ -n_1 n_{p-1} & -n_2 n_{p-1} & \cdots & & n_{p-1}(n - n_{p-1}) \end{pmatrix}.$$

For convenience, write $[\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1}$ as follows

$$[\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1} = \frac{1}{n} \begin{pmatrix} -n_1^2 & -n_1 n_2 & -n_1 n_3 & \cdots & -n_1 n_{p-1} \\ -n_1 n_2 & -n_2^2 & -n_2 n_3 & \cdots & -n_2 n_{p-1} \\ \vdots & & \ddots & & \vdots \\ -n_1 n_{p-1} & -n_2 n_{p-1} & \cdots & & -n_{p-1}^2 \end{pmatrix} + \begin{pmatrix} n_1 & 0 & 0 & \cdots & 0 \\ 0 & n_2 & 0 & \cdots & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \cdots & 0 & n_{p-1} \end{pmatrix}.$$

Then,

$$\begin{aligned} & (\mathbf{L}\hat{\mathbf{B}})^T \left[\mathbf{L} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T \right]^{-1} (\mathbf{L}\hat{\mathbf{B}}) = \\ & -\frac{1}{n} \sum_{i=1}^{p-1} \sum_{j=1}^{p-1} n_i n_j (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)(\bar{\mathbf{y}}_j - \bar{\mathbf{y}}_p)^T + \sum_{i=1}^{p-1} n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)(\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)^T = \mathbf{H}. \end{aligned}$$

Let \mathbf{X} be as in (2.8). Then the multivariate linear regression (MREG) Hotelling Lawley test statistic for testing $H_0 : \mathbf{L}\mathbf{B} = \mathbf{0}$ versus $H_0 : \mathbf{L}\mathbf{B} \neq \mathbf{0}$ has

$$U = \text{tr}(\mathbf{W}^{-1}\mathbf{H}).$$

One way MANOVA is used to test $H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \dots = \boldsymbol{\mu}_p$. The Hotelling Lawley test statistic for testing for above hypotheses is

$$U = \text{tr}(\mathbf{W}^{-1}\mathbf{B}_T)$$

where

$$\mathbf{W} = (n-p)\hat{\boldsymbol{\Sigma}}\boldsymbol{\epsilon} \quad \text{and} \quad \mathbf{B}_T = \sum_{i=1}^p n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}})(\bar{\mathbf{y}}_i - \bar{\mathbf{y}})^T.$$

Theorem 2.1. *The one-way MANOVA test statistic and the Hotelling Lawley trace test statistic are the same for the design matrix as in (2.8).*

To show that the above two test statistics are equal it is sufficient to prove that $\mathbf{H} = \mathbf{B}_T$.

Proof. Special case I: $p = 2$ (Two group case)

Consider \mathbf{H} .

$$\mathbf{H} = -\frac{1}{n}n_1n_2(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T + n_1(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T. \text{ Since } n = n_1 + n_2,$$

$$\mathbf{H} = -\frac{1}{n}(nn_1 - n_1n_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T + n_1(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T$$

$$\mathbf{H} = -n_1(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T + \frac{n_1n_2}{n}(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T + n_1(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T$$

$$\mathbf{H} = \frac{n_1n_2}{n}(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T.$$

Now consider \mathbf{B}_T with $p = 2$.

Note that $\bar{\mathbf{y}} = (n_1\bar{\mathbf{y}}_1 + n_2\bar{\mathbf{y}}_2)/n$ and

$$\mathbf{B}_T = n_1(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}})(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}})^T + n_2(\bar{\mathbf{y}}_2 - \bar{\mathbf{y}})(\bar{\mathbf{y}}_2 - \bar{\mathbf{y}})^T$$

$$\mathbf{B}_T = \frac{n_1}{n^2}(n\bar{\mathbf{y}}_1 - n_1\bar{\mathbf{y}}_1 - n_2\bar{\mathbf{y}}_2)(n\bar{\mathbf{y}}_1 - n_1\bar{\mathbf{y}}_1 - n_2\bar{\mathbf{y}}_2)^T + \frac{n_2}{n^2}(n\bar{\mathbf{y}}_2 - n_1\bar{\mathbf{y}}_1 - n_2\bar{\mathbf{y}}_2)(n\bar{\mathbf{y}}_2 - n_1\bar{\mathbf{y}}_1 - n_2\bar{\mathbf{y}}_2)^T$$

$$\mathbf{B}_T = \frac{n_1 n_2}{n^2}(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T + \frac{n_1^2 n_2}{n^2}(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T$$

$$\mathbf{B}_T = \frac{n_1 n_2}{n}(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)(\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_2)^T.$$

Therefore $\mathbf{B}_T = \mathbf{H}$ when $p = 2$.

□

Proof. Special case II: $n_i = n_1 \forall i = 1, \dots, p$

$$\mathbf{H} = -\frac{1}{n} \sum_{i=1}^{p-1} \sum_{j=1}^{p-1} n_i n_j (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)(\bar{\mathbf{y}}_j - \bar{\mathbf{y}}_p)^T + \sum_{i=1}^{p-1} n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)(\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)^T.$$

Note that the i, j running from 1 through $p - 1$ and i, j running from 1 through p would yield the same \mathbf{H} . Therefore \mathbf{H} can be written as

$$\mathbf{H} = -\frac{1}{n} \sum_{i=1}^p \sum_{j=1}^p n_i n_j (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)(\bar{\mathbf{y}}_j - \bar{\mathbf{y}}_p)^T + \sum_{i=1}^p n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)(\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)^T.$$

Now consider the double sum in \mathbf{H} . Note that $n = n_1 p$ and

$$\begin{aligned} -\frac{1}{n} \sum_{i=1}^p \sum_{j=1}^p n_i n_j (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)(\bar{\mathbf{y}}_j - \bar{\mathbf{y}}_p)^T &= \frac{-n_1^2}{n_1 p} \sum_{i=1}^p \sum_{j=1}^p (\bar{\mathbf{y}}_i \bar{\mathbf{y}}_j^T - \bar{\mathbf{y}}_i \bar{\mathbf{y}}_p^T - \bar{\mathbf{y}}_p \bar{\mathbf{y}}_j^T + \bar{\mathbf{y}}_p \bar{\mathbf{y}}_p^T) \\ &= \frac{n_1}{p} \left[-\sum_{i=1}^p \sum_{j=1}^p (\bar{\mathbf{y}}_i \bar{\mathbf{y}}_j^T) + p \left(\sum_{i=1}^p \bar{\mathbf{y}}_i \right) \bar{\mathbf{y}}_p^T + p \bar{\mathbf{y}}_p \left(\sum_{j=1}^p \bar{\mathbf{y}}_j^T \right) - p^2 \bar{\mathbf{y}}_p \bar{\mathbf{y}}_p^T \right]. \end{aligned} \quad (2.12)$$

Now consider the rest of \mathbf{H} ,

$$n_1 \sum_{i=1}^p (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)(\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)^T = n_1 \sum_{i=1}^p \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^T - n_1 \left(\sum_{i=1}^p \bar{\mathbf{y}}_i \right) \bar{\mathbf{y}}_p^T - n_1 \bar{\mathbf{y}}_p \left(\sum_{i=1}^p \bar{\mathbf{y}}_i^T \right) + n_1 p \bar{\mathbf{y}}_p \bar{\mathbf{y}}_p^T. \quad (2.13)$$

Therefore by (2.12) and (2.13), it is clear that

$$\mathbf{H} = n_1 \sum_{i=1}^p \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^T - \frac{n_1}{p} \sum_{i=1}^p \sum_{j=1}^p \bar{\mathbf{y}}_i \bar{\mathbf{y}}_j^T. \quad (2.14)$$

Now consider

$$\mathbf{B}_T = n_1 \sum_{i=1}^p (\bar{\mathbf{y}}_i - \bar{\mathbf{y}})(\bar{\mathbf{y}}_i - \bar{\mathbf{y}})^T. \quad (2.15)$$

Let

$$\bar{\mathbf{Y}} = \begin{pmatrix} \bar{\mathbf{y}}_1^T \\ \bar{\mathbf{y}}_2^T \\ \vdots \\ \bar{\mathbf{y}}_p^T \end{pmatrix}. \quad \text{Then } \mathbf{B}_T = n_1 \left[\bar{\mathbf{Y}}^T \bar{\mathbf{Y}} - \frac{1}{p} \bar{\mathbf{Y}}^T \mathbf{1} \mathbf{1}^T \bar{\mathbf{Y}} \right].$$

Therefore, \mathbf{B}_T becomes

$$\mathbf{B}_T = n_1 \sum_{i=1}^p \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^T - \frac{n_1}{p} \sum_{i=1}^p \sum_{j=1}^p \bar{\mathbf{y}}_i \bar{\mathbf{y}}_j^T. \quad (2.16)$$

From (2.15) and (2.16) $\mathbf{B}_T = \mathbf{H}$.

□

Proof. General case:

$$\mathbf{H} = -\frac{1}{n} \sum_{i=1}^p \sum_{j=1}^p n_i n_j (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p) (\bar{\mathbf{y}}_j - \bar{\mathbf{y}}_p)^T + \sum_{i=1}^p n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p) (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)^T.$$

Now consider the double sum in \mathbf{H} .

$$\begin{aligned} & -\frac{1}{n} \sum_{i=1}^p \sum_{j=1}^p n_i n_j (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p) (\bar{\mathbf{y}}_j - \bar{\mathbf{y}}_p)^T = \\ & -\frac{1}{n} \sum_{i=1}^p \sum_{j=1}^p n_i n_j \bar{\mathbf{y}}_i \bar{\mathbf{y}}_j^T + \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^p n_i n_j \bar{\mathbf{y}}_i \bar{\mathbf{y}}_p^T + \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^p n_i n_j \bar{\mathbf{y}}_p \bar{\mathbf{y}}_j^T - \frac{1}{n} \bar{\mathbf{y}}_p \bar{\mathbf{y}}_p^T \sum_{i=1}^p \sum_{j=1}^p n_i n_j \end{aligned} \quad (2.17)$$

$$\begin{aligned} & -\frac{1}{n} \sum_{i=1}^p n_i \bar{\mathbf{y}}_i \sum_{j=1}^p n_j \bar{\mathbf{y}}_j^T + \frac{1}{n} \sum_{i=1}^p n_i \bar{\mathbf{y}}_i \sum_{j=1}^p n_j \bar{\mathbf{y}}_p^T + \frac{1}{n} \bar{\mathbf{y}}_p \sum_{i=1}^p n_i \sum_{j=1}^p n_j \bar{\mathbf{y}}_j^T - \frac{1}{n} \bar{\mathbf{y}}_p \bar{\mathbf{y}}_p^T n^2 \\ & -\frac{1}{n} n \bar{\mathbf{y}} n \bar{\mathbf{y}}^T + \frac{1}{n} \sum_{i=1}^p n_i \bar{\mathbf{y}}_i n \bar{\mathbf{y}}_p^T + \frac{1}{n} \bar{\mathbf{y}}_p n \sum_{j=1}^p n_j \bar{\mathbf{y}}_j^T - n \bar{\mathbf{y}}_p \bar{\mathbf{y}}_p^T \\ & -n \bar{\mathbf{y}} \bar{\mathbf{y}}^T + \sum_{i=1}^p n_i \bar{\mathbf{y}}_i \bar{\mathbf{y}}_p^T + \bar{\mathbf{y}}_p \sum_{j=1}^p n_j \bar{\mathbf{y}}_j^T - n \bar{\mathbf{y}}_p \bar{\mathbf{y}}_p^T. \end{aligned} \quad (2.18)$$

Now consider the rest of \mathbf{H} ,

$$\sum_{i=1}^p n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p) (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_p)^T = \sum_{i=1}^p n_i \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^T - \sum_{i=1}^p n_i \bar{\mathbf{y}}_i \bar{\mathbf{y}}_p^T - \bar{\mathbf{y}}_p \sum_{i=1}^p n_i \bar{\mathbf{y}}_i^T + n \bar{\mathbf{y}}_p \bar{\mathbf{y}}_p^T. \quad (2.19)$$

Therefore by (2.18) and (2.19)

$$\mathbf{H} = \sum_{i=1}^p n_i \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^T - n \bar{\mathbf{y}} \bar{\mathbf{y}}^T. \quad (2.20)$$

Now consider

$$\begin{aligned} \mathbf{B}_T &= \sum_{i=1}^p n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}) (\bar{\mathbf{y}}_i - \bar{\mathbf{y}})^T \\ \mathbf{B}_T &= \sum_{i=1}^p n_i \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^T - \sum_{i=1}^p n_i \bar{\mathbf{y}}_i \bar{\mathbf{y}}^T - \bar{\mathbf{y}} \sum_{i=1}^p n_i \bar{\mathbf{y}}_i^T + \bar{\mathbf{y}} \bar{\mathbf{y}}^T \sum_{i=1}^p n_i \\ \mathbf{B}_T &= \sum_{i=1}^p n_i \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^T - n \bar{\mathbf{y}} \bar{\mathbf{y}}^T - \bar{\mathbf{y}} n \bar{\mathbf{y}}^T + n \bar{\mathbf{y}} \bar{\mathbf{y}}^T \\ \mathbf{B}_T &= \sum_{i=1}^p n_i \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^T - n \bar{\mathbf{y}} \bar{\mathbf{y}}^T. \end{aligned} \quad (2.21)$$

(2.20) and (2.21) proves that $\mathbf{H} = \mathbf{B}_T$.

□

2.5 CELL MEANS MODEL

$$\mathbf{X} = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 1 & \cdots & 0 \\ 0 & 0 & \cdots & 1 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} \boldsymbol{\mu}_1^T \\ \vdots \\ \boldsymbol{\mu}_p^T \end{pmatrix} \quad \text{and} \quad \mathbf{L} = \begin{pmatrix} \mathbf{I}_{p-1} & -\mathbf{1} \end{pmatrix}$$

$$\hat{\mathbf{B}} = \begin{pmatrix} \bar{\mathbf{y}}_1^T \\ \vdots \\ \bar{\mathbf{y}}_p^T \end{pmatrix}, \quad \mathbf{L}\hat{\mathbf{B}} = \begin{pmatrix} (\bar{\mathbf{y}}_1 - \bar{\mathbf{y}}_p)^T \\ (\bar{\mathbf{y}}_2 - \bar{\mathbf{y}}_p)^T \\ \vdots \\ (\bar{\mathbf{y}}_{p-1} - \bar{\mathbf{y}}_p)^T \end{pmatrix}.$$

Then $\mathbf{X}^T \mathbf{X} = \text{diag}(n_1, \dots, n_{p-1})$ and $(\mathbf{X}^T \mathbf{X})^{-1} = \text{diag}\left(\frac{1}{n_1}, \dots, \frac{1}{n_{p-1}}\right)$.

Then $\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T$ becomes

$$\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T = \frac{1}{n_p} \begin{pmatrix} 1 + \frac{n_p}{n_1} & 1 & 1 & \cdots & 1 \\ 1 & 1 + \frac{n_p}{n_2} & 1 & \cdots & 1 \\ \vdots & & \ddots & & \vdots \\ 1 & 1 & \cdots & 1 & 1 + \frac{n_p}{n_{p-1}} \end{pmatrix}. \quad (2.22)$$

Corollary 2.2. *Theorem 2.1 does not depend on the full rank design matrix.*

Notice that the matrix equation (2.22) is the exactly same as (2.11). This is an indication that Theorem 2.1 does not depend on the design matrix.

CHAPTER 3

BOOTSTRAPPING ANALOGS OF THE TWO SAMPLE HOTELLING'S T^2 TEST

Suppose there are two independent random samples from two populations or groups. A common multivariate two sample test of hypotheses is $H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$ versus $H_1 : \boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2$ where $\boldsymbol{\mu}_i$ is a population location measure of the i th population for $i = 1, 2$. The two sample Hotelling's T^2 test is the classical method, and is a special case of the one way MANOVA model if the two populations are assumed to have the same population covariance matrix. This chapter suggests using the Olive (2017abc) bootstrap technique to develop analogs of Hotelling's T^2 test. The new tests can have considerable outlier resistance, and the tests do not need the population covariance matrices to be equal.

3.1 APPLYING THE PREDICTION REGION METHOD TO THE TWO SAMPLE TEST

The two sample test of $H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$ versus $H_1 : \boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2$ uses $\boldsymbol{\mu} = \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2 = \mathbf{c} = \mathbf{0}$ with $\mathbf{w}_i = T_{i1}^* - T_{i2}^*$ for $i = 1, \dots, B$. Make the prediction region (2.7) where $T_i^* = \mathbf{w}_i$. Fail to reject H_0 if $\mathbf{0}$ is in the prediction region, reject H_0 otherwise. A sample of size n_i is drawn with replacement from $\mathbf{x}_{1,i}, \dots, \mathbf{x}_{n_i,i}$ for $i = 1, 2$ to obtain the bootstrap sample.

For illustrative purposes, the simulation study will take T_i to be the coordinatewise median, the (Olive (2017b, ch. 4), Olive and Hawkins (2010), and Zhang, Olive, and Ye (2012)) RMVN estimator T_{RMVN} , the sample mean, and the 25% trimmed mean. The asymptotic covariance matrix of the coordinatewise median is difficult to estimate, while that of the RMVN estimator is unknown. The RMVN estimator has been shown to be \sqrt{n} consistent on a large class of elliptically contoured distributions, but has not yet been shown to be asymptotically normal. Hence the bootstrap "test" for the RMVN estimator should be used for exploratory purposes.

The RMVN estimator $(T_{RMVN}, \mathbf{C}_{RMVN})$ uses a concentration algorithm. Let $(T_{-1,j}, \mathbf{C}_{-1,j})$ be the j th start (initial estimator) and compute all n Mahalanobis distances $D_i(T_{-1,j}, \mathbf{C}_{-1,j})$. At the next iteration, the classical estimator $(T_{0,j}, \mathbf{C}_{0,j}) = (\bar{\mathbf{x}}_{0,j}, \mathbf{S}_{0,j})$

is computed from the $c_n \approx n/2$ cases corresponding to the smallest distances. This iteration can be continued for k *concentration steps* resulting in the sequence of estimators $(T_{-1,j}, \mathbf{C}_{-1,j}), (T_{0,j}, \mathbf{C}_{0,j}), \dots, (T_{k,j}, \mathbf{C}_{k,j})$. The result of the iteration $(T_{k,j}, \mathbf{C}_{k,j})$ is called the j th *attractor*. The algorithm estimator uses one of the attractors. The RMVN estimator uses the same two starts as the Olive (2004) MBA estimator: $(\bar{\mathbf{x}}, \mathbf{S})$ and $(MED(n), \mathbf{I}_p)$ where $MED(n)$ is the coordinatewise median. Then the location estimator T_{RMVN} can be used to test $H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$.

3.2 REAL DATA EXAMPLE

The Johnson (1996) STATLIB bodyfat data consists of 252 observations on 15 variables including the density determined from underwater weighing and the percent body fat measurement. Consider these two variables with two age groups: age ≤ 50 and age > 50 . The test with the RMVN estimator had $D_0 = 1.78$ while the test with the coordinatewise median had $D_0 = 1.35$. Both tests had cutoffs near 2.37 and fail to reject H_0 . The classical two sample Hotelling's T^2 test rejects H_0 with a test statistic of 4.74 and a p-value of 0.001.

The DD plots, shown in Figures 3.1 and 3.2, reveal five outliers. After deleting the outliers, the three tests all fail to reject H_0 . The RMVN test had $D_0 = 1.63$ with cutoff 2.25, the coordinatewise median test had $D_0 = 1.22$ with cutoff 2.38, and the classical test had test statistic 2.39 with a p-value of 0.09.

See the simulation set up and the simulation results in Chapter 5.

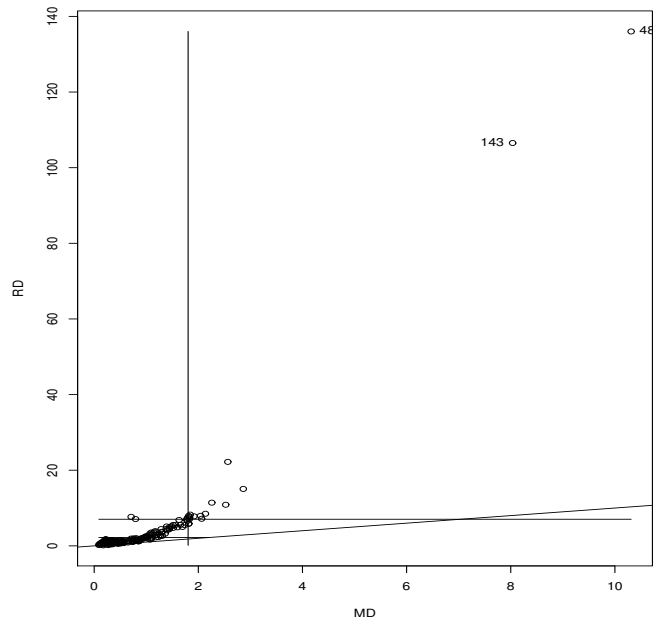


Figure 3.1. DD plot for the age ≤ 50 group.

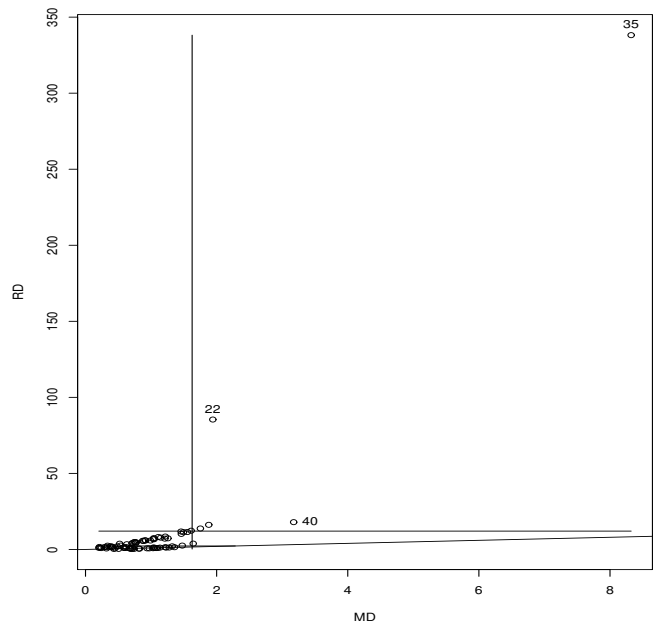


Figure 3.2. DD plot for the age > 50 group.

CHAPTER 4

BOOTSTRAPPING ANALOGS OF THE ONE WAY MANOVA TEST

The classical one way MANOVA model is used to test whether the mean measurements are the same or differ across p groups, and assumes that the covariance matrix of each group is the same. This chapter suggests using the Olive (2017abc) bootstrap technique to develop analogs of the one way MANOVA test. The new tests can have considerable outlier resistance, and the tests do not need the population covariance matrices to be equal.

The multivariate linear model

$$\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i + \boldsymbol{\epsilon}_i$$

for $i = 1, \dots, n$ has $m \geq 2$ response variables Y_1, \dots, Y_m and p predictor variables x_1, x_2, \dots, x_p . The i th case is $(\mathbf{x}_i^T, \mathbf{y}_i^T) = (x_{i1}, x_{i2}, \dots, x_{ip}, Y_{i1}, \dots, Y_{im})$. The model is written in matrix form as $\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{E}$ where the matrices are defined below. The model has $E(\boldsymbol{\epsilon}_k) = \mathbf{0}$ and $\text{Cov}(\boldsymbol{\epsilon}_k) = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} = (\sigma_{ij})$ for $k = 1, \dots, n$. Then the $p \times m$ coefficient matrix $\mathbf{B} = \begin{bmatrix} \boldsymbol{\beta}_1 & \boldsymbol{\beta}_2 & \dots & \boldsymbol{\beta}_m \end{bmatrix}$ and the $m \times m$ covariance matrix $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$ are to be estimated, and $E(\mathbf{Z}) = \mathbf{X}\mathbf{B}$ while $E(Y_{ij}) = \mathbf{x}_i^T \boldsymbol{\beta}_j$. The $\boldsymbol{\epsilon}_i$ are assumed to be iid. The univariate linear model corresponds to $m = 1$ response variable, and is written in matrix form as $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$. Subscripts are needed for the m univariate linear models $\mathbf{Y}_j = \mathbf{X}\boldsymbol{\beta}_j + \mathbf{e}_j$ for $j = 1, \dots, m$ where $E(\mathbf{e}_j) = \mathbf{0}$. For the multivariate linear model, $\text{Cov}(\mathbf{e}_i, \mathbf{e}_j) = \sigma_{ij} \mathbf{I}_n$ for $i, j = 1, \dots, m$ where \mathbf{I}_n is the $n \times n$ identity matrix.

The $n \times m$ matrix

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Y}_1 & \mathbf{Y}_2 & \dots & \mathbf{Y}_m \end{bmatrix} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_n^T \end{bmatrix}.$$

The $n \times p$ design matrix \mathbf{X} of predictor variables is not necessarily of full rank p , and

$$\mathbf{X} = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \dots & \mathbf{v}_p \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

where often $\mathbf{v}_1 = \mathbf{1}$.

The $p \times m$ matrix

$$\mathbf{B} = \begin{bmatrix} \boldsymbol{\beta}_1 & \boldsymbol{\beta}_2 & \dots & \boldsymbol{\beta}_m \end{bmatrix}.$$

The $n \times m$ matrix

$$\mathbf{E} = \begin{bmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \dots & \mathbf{e}_m \end{bmatrix} = \begin{bmatrix} \boldsymbol{\epsilon}_1^T \\ \vdots \\ \boldsymbol{\epsilon}_n^T \end{bmatrix}.$$

Considering the i th row of \mathbf{Z} , \mathbf{X} , and \mathbf{E} shows that $\mathbf{y}_i^T = \mathbf{x}_i^T \mathbf{B} + \boldsymbol{\epsilon}_i^T$.

The multivariate linear regression model and one way MANOVA model are special cases of the multivariate linear model, but using double subscripts will be useful for describing the one way MANOVA model. Suppose there are independent random samples of size n_i from p different populations (treatments), or n_i cases are randomly assigned to p treatment groups where $n = \sum_{i=1}^p n_i$. Assume that m response variables $\mathbf{y}_{ij} = (Y_{ij1}, \dots, Y_{ijm})^T$ are measured for the i th treatment group and the j th case (often an individual or thing) in the group. Hence $i = 1, \dots, p$ and $j = 1, \dots, n_i$. The Y_{ijk} follow different one way ANOVA models for $k = 1, \dots, m$. Assume $E(\mathbf{y}_{ij}) = \boldsymbol{\mu}_i$ and $\text{Cov}(\mathbf{y}_{ij}) = \boldsymbol{\Sigma}_\epsilon$. Hence the p treatments have different mean vectors $\boldsymbol{\mu}_i$, but common covariance matrix $\boldsymbol{\Sigma}_\epsilon$.

The one way MANOVA is used to test $H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \dots = \boldsymbol{\mu}_p$. Often $\boldsymbol{\mu}_i = \boldsymbol{\mu} + \boldsymbol{\tau}_i$, so H_0 becomes $H_0 : \boldsymbol{\tau}_1 = \dots = \boldsymbol{\tau}_p$. If $m = 1$, the one way MANOVA model is the one way ANOVA model. MANOVA is useful since it takes into account the correlations between the m response variables. The Hotelling's T^2 test that uses a common covariance matrix is a special case of the one way MANOVA model with $p = 2$.

4.1 AN ALTERNATIVE TO THE USUAL ONE WAY MANOVA

A useful one way MANOVA model is $\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{E}$ where \mathbf{X} is the full rank matrix where the first column of \mathbf{X} is $\mathbf{v}_1 = \mathbf{1}$ and the i th column \mathbf{v}_i of \mathbf{X} is an indicator for group $i - 1$ for $i = 2, \dots, p$. For example, $\mathbf{v}_3 = (\mathbf{0}^T, \mathbf{1}^T, \mathbf{0}^T, \dots, \mathbf{0}^T)^T$ where the p vectors in \mathbf{v}_3 have lengths n_1, n_2, \dots, n_p , respectively. Then $\hat{\beta}_{1k} = \bar{Y}_{p0k} = \hat{\mu}_{pk}$ for $k = 1, \dots, m$, and

$$\hat{\beta}_{ik} = \bar{Y}_{i-1,0k} - \bar{Y}_{p0k} = \hat{\mu}_{i-1,k} - \hat{\mu}_{pk}$$

for $k = 1, \dots, m$ and $i = 2, \dots, p$. Thus testing $H_0 : \boldsymbol{\mu}_1 = \dots = \boldsymbol{\mu}_p$ is equivalent to testing $H_0 : \mathbf{LB} = \mathbf{0}$ where $\mathbf{L} = [\mathbf{0} \ \mathbf{I}_{p-1}]$. Press (2005, p. 262) uses the above model. Then $\mathbf{y}_{ij} = \boldsymbol{\mu}_i + \boldsymbol{\epsilon}_{ij}$ and

$$\mathbf{B} = \begin{bmatrix} \boldsymbol{\mu}_p^T \\ (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_p)^T \\ (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_p)^T \\ \vdots \\ (\boldsymbol{\mu}_{p-2} - \boldsymbol{\mu}_p)^T \\ (\boldsymbol{\mu}_{p-1} - \boldsymbol{\mu}_p)^T \end{bmatrix}.$$

Then a test statistic for the one way Manova model is \mathbf{w} given by Equation (4.1) with $T_i = \hat{\boldsymbol{\mu}}_i = \bar{\mathbf{y}}_i$ where it is assumed that $\boldsymbol{\Sigma}_i \equiv \boldsymbol{\Sigma}\boldsymbol{\epsilon}$ for $i = 1, \dots, p$.

Large sample theory can be used to derive a better test that does not need the equal population covariance matrix assumption $\boldsymbol{\Sigma}_i \equiv \boldsymbol{\Sigma}\boldsymbol{\epsilon}$. To simplify the large sample theory, assume $n_i = \pi_i n$ where $0 < \pi_i < 1$ and $\sum_{i=1}^p \pi_i = 1$. Assume H_0 is true, and let $\boldsymbol{\mu}_i = \boldsymbol{\mu}$ for $i = 1, \dots, p$. Suppose the $\boldsymbol{\mu}_i = \boldsymbol{\mu}$ and $\sqrt{n_i}(T_i - \boldsymbol{\mu}) \xrightarrow{D} N_m(\mathbf{0}, \boldsymbol{\Sigma}_i)$, and $\sqrt{n}(T_i - \boldsymbol{\mu}) \xrightarrow{D} N_m\left(\mathbf{0}, \frac{\boldsymbol{\Sigma}_i}{\pi_i}\right)$. Let

$$\mathbf{w} = \begin{bmatrix} T_1 - T_p \\ T_2 - T_p \\ \vdots \\ T_{p-2} - T_p \\ T_{p-1} - T_p \end{bmatrix}. \quad (4.1)$$

Then $\sqrt{n}\mathbf{w} \xrightarrow{D} N_{m(p-1)}(\mathbf{0}, \boldsymbol{\Sigma}\mathbf{w})$ with $\boldsymbol{\Sigma}\mathbf{w} = (\boldsymbol{\Sigma}_{ij})$ where $\boldsymbol{\Sigma}_{ij} = \frac{\boldsymbol{\Sigma}_p}{\pi_p}$ for $i \neq j$, and $\boldsymbol{\Sigma}_{ii} = \frac{\boldsymbol{\Sigma}_i}{\pi_i} + \frac{\boldsymbol{\Sigma}_p}{\pi_p}$ for $i = j$. Hence

$$t_0 = n\mathbf{w}^T \hat{\boldsymbol{\Sigma}}\mathbf{w} = \mathbf{w}^T \left(\frac{\hat{\boldsymbol{\Sigma}}\mathbf{w}}{n} \right)^{-1} \mathbf{w} \xrightarrow{D} \chi_{m(p-1)}^2$$

as the $n_i \rightarrow \infty$ if H_0 is true. Here

$$\frac{\hat{\Sigma} \mathbf{w}}{n} = \begin{bmatrix} \frac{\hat{\Sigma}_1}{n_1} + \frac{\hat{\Sigma}_p}{n_p} & \frac{\hat{\Sigma}_p}{n_p} & \frac{\hat{\Sigma}_p}{n_p} & \cdots & \frac{\hat{\Sigma}_p}{n_p} \\ \frac{\hat{\Sigma}_p}{n_p} & \frac{\hat{\Sigma}_2}{n_2} + \frac{\hat{\Sigma}_p}{n_p} & \frac{\hat{\Sigma}_p}{n_p} & \cdots & \frac{\hat{\Sigma}_p}{n_p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\hat{\Sigma}_p}{n_p} & \frac{\hat{\Sigma}_p}{n_p} & \frac{\hat{\Sigma}_p}{n_p} & \cdots & \frac{\hat{\Sigma}_{p-1}}{n_{p-1}} + \frac{\hat{\Sigma}_p}{n_p} \end{bmatrix}$$

is a block matrix where the off diagonal block entries equal $\hat{\Sigma}_p/n_p$ and the i th diagonal block entry is $\frac{\hat{\Sigma}_i}{n_i} + \frac{\hat{\Sigma}_p}{n_p}$ for $i = 1, \dots, (p-1)$.

Reject H_0 if $t_0 > m(p-1)F_{m(p-1), d_n}(1-\alpha)$ where $d_n = \min(n_1, \dots, n_p)$. It may make sense to relabel the groups so that n_p is the largest n_i or $\hat{\Sigma}_p/n_p$ has the smallest generalized variance of the $\hat{\Sigma}_i/n_i$. This test may start to outperform the one way MANOVA test if $n \geq (m+p)^2$ and $n_i \geq 20m$ for $i = 1, \dots, p$.

Olive (2017b, ch. 10) has the above result where $T_i = \bar{\mathbf{y}}_i$ is the sample mean and $\hat{\Sigma}_i = \mathbf{S}_i$ is the sample covariance matrix of the i th group. Then Σ_i is the population covariance matrix of the i th group. The following theorem gives the general result.

Theorem 4.1. *If*

$$\begin{pmatrix} \sqrt{n_1} (T_1 - \boldsymbol{\mu}_1) \\ \vdots \\ \sqrt{n_p} (T_p - \boldsymbol{\mu}_p) \end{pmatrix} \xrightarrow{D} N_{mp} \left[\begin{pmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \Sigma_1 & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \Sigma_p \end{pmatrix} \right],$$

then under $H_0 : \boldsymbol{\mu}_1 = \cdots = \boldsymbol{\mu}_p$

$$\sqrt{n} \begin{pmatrix} T_1 - T_p \\ \vdots \\ T_{p-1} - T_p \end{pmatrix} \xrightarrow{D} N_{(m-1)p} \left[\begin{pmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \frac{\Sigma_1}{\pi_1} + \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \cdots & \frac{\Sigma_p}{\pi_p} \\ \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_2}{\pi_2} + \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \cdots & \frac{\Sigma_p}{\pi_p} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \cdots & \frac{\Sigma_{p-1}}{\pi_{p-1}} + \frac{\Sigma_p}{\pi_p} \end{pmatrix} \right].$$

Proof. To simplify large sample theory, assume $n_i = \pi_i n$ for some positive real π_i and $i = 1, \dots, p$. Let $\hat{\Sigma}_i$ be a consistent nonsingular estimator of Σ_i . Then

$$\begin{pmatrix} \sqrt{n_1} (T_1 - \boldsymbol{\mu}_1) \\ \vdots \\ \sqrt{n_p} (T_p - \boldsymbol{\mu}_p) \end{pmatrix} \xrightarrow{D} N_{mp} \left[\begin{pmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \Sigma_1 & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \Sigma_p \end{pmatrix} \right],$$

Under $H_0 : \boldsymbol{\mu}_1 = \cdots = \boldsymbol{\mu}_p = \boldsymbol{\mu}$

$$\begin{pmatrix} \sqrt{n} (T_1 - \boldsymbol{\mu}) \\ \vdots \\ \sqrt{n} (T_p - \boldsymbol{\mu}) \end{pmatrix} \xrightarrow{D} N_{mp} \left[\begin{pmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \frac{\Sigma_1}{\pi_1} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \frac{\Sigma_p}{\pi_p} \end{pmatrix} \right] = N_{mp}(\mathbf{0}, \Sigma).$$

Let A be a $(m-1)p \times mp$ matrix and

$$\mathbf{A} = \begin{pmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & -\mathbf{I} \\ \mathbf{0} & \mathbf{I} & \mathbf{0} & \cdots & \mathbf{0} & -\mathbf{I} \\ \mathbf{0} & \mathbf{0} & \mathbf{I} & \mathbf{0} & \cdots & -\mathbf{I} \\ \vdots & & & \ddots & & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{I} & -\mathbf{I} \end{pmatrix}.$$

Then,

$$\mathbf{A} \begin{pmatrix} \sqrt{n} (T_1 - \boldsymbol{\mu}) \\ \vdots \\ \sqrt{n} (T_p - \boldsymbol{\mu}) \end{pmatrix} = \sqrt{n} \begin{pmatrix} T_1 - T_p \\ \vdots \\ T_{p-1} - T_p \end{pmatrix}$$

and

$$\mathbf{A} \Sigma \mathbf{A}^T = \begin{pmatrix} \frac{\Sigma_1}{\pi_1} + \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \cdots & \frac{\Sigma_p}{\pi_p} \\ \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_2}{\pi_2} + \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \cdots & \frac{\Sigma_p}{\pi_p} \\ \vdots & & \ddots & & \vdots \\ \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \cdots & \frac{\Sigma_{p-1}}{k_{p-1}} + \frac{\Sigma_p}{\pi_p} \end{pmatrix}.$$

Therefore,

$$\sqrt{n} \begin{pmatrix} T_1 - T_p \\ \vdots \\ T_{p-1} - T_p \end{pmatrix} \xrightarrow{D} N_{(m-1)p} \left[\begin{pmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \frac{\Sigma_1}{\pi_1} + \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \dots & \frac{\Sigma_p}{\pi_p} \\ \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_2}{\pi_2} + \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \dots & \frac{\Sigma_p}{\pi_p} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \frac{\Sigma_p}{\pi_p} & \dots & \frac{\Sigma_{p-1}}{\pi_{p-1}} + \frac{\Sigma_p}{\pi_p} \end{pmatrix} \right]$$

□

If $\mathbf{T} = (T_1^T, T_2^T, \dots, T_p^T)^T$, $\boldsymbol{\theta} = (\boldsymbol{\mu}_1^T, \boldsymbol{\mu}_2^T, \dots, \boldsymbol{\mu}_p^T)^T$, \mathbf{c} is a constant vector, and \mathbf{A} is a full rank $r \times mp$ matrix with rank r , then a large sample test of the form $H_0 : \mathbf{A}\boldsymbol{\theta} = \mathbf{c}$ versus $H_1 : \mathbf{A}\boldsymbol{\theta} \neq \mathbf{c}$ uses

$$\mathbf{A}\sqrt{n}(\mathbf{T} - \boldsymbol{\theta}) \xrightarrow{D} N_r \left(\mathbf{0}, \mathbf{A} \text{diag} \left(\frac{\Sigma_1}{\pi_1}, \frac{\Sigma_2}{\pi_2}, \dots, \frac{\Sigma_p}{\pi_p} \right) \mathbf{A}^T \right).$$

When H_0 is true, the statistic

$$t_0 = [\mathbf{A}\mathbf{T} - \mathbf{c}]^T \left[\mathbf{A} \text{diag} \left(\frac{\hat{\Sigma}_1}{n_1}, \frac{\hat{\Sigma}_2}{n_2}, \dots, \frac{\hat{\Sigma}_p}{n_p} \right) \mathbf{A}^T \right]^{-1} [\mathbf{A}\mathbf{T} - \mathbf{c}] \xrightarrow{D} \chi_r^2.$$

The same statistic was used by Zhang and Liu (2013, p. 138) with $T_i = \bar{\mathbf{y}}_i$ and $\hat{\Sigma}_i = \mathbf{S}_i$.

4.1.1 Test H_0 when $\hat{\Sigma}_{\mathbf{w}}$ is unknown or difficult to estimate.

Since the common covariance matrix assumption $\text{Cov}(\boldsymbol{\epsilon}_k) = \Sigma_{\boldsymbol{\epsilon}}$ for $k = 1, \dots, p$ is extremely strong, using the prediction region method to test $H_0 : \mathbf{L}\mathbf{B} = \mathbf{0}$ may be a useful alternative. Take a sample of size n_k with replacement from the n_k cases for each group for $k = 1, 2, \dots, p$. Let $\hat{\mathbf{B}}_i^*$ be the i th bootstrap estimator of \mathbf{B} for $i = 1, \dots, B$. Let the $(p-1)m \times 1$ vector $\mathbf{w}_i = \text{vec}(\mathbf{L}\hat{\mathbf{B}}_i^*) = ((\hat{\boldsymbol{\mu}}_1^* - \hat{\boldsymbol{\mu}}_p^*)^T, \dots, (\hat{\boldsymbol{\mu}}_{p-1}^* - \hat{\boldsymbol{\mu}}_p^*)^T)^T$ for $i = 1, \dots, B$, where $\text{vec}(\mathbf{A})$ stacks columns of a matrix into a vector. For a robust test use $\mathbf{w}_i = ((T_1^* - T_p^*)^T, \dots, (T_{p-1}^* - T_p^*)^T)^T$ where T_k is a robust location estimator, such as the coordinatewise median or trimmed mean, applied to the cases in the k th treatment group. The prediction region method fails to reject H_0 if $\mathbf{0}$ is in the resulting confidence region. We likely need $n \geq 40mp$, $n \geq (m+p)^2$, and $n_i \geq 40m$.

4.2 POWER COMPARISON AMONG THE TESTS

Figures 4.1, 4.2 and 4.3 try to compare the powers among the tests mentioned above with the classical test. Here $\delta_1 \in \{0.00, 0.04, 0.08, 0.12, 0.16, 0.20, 0.24, 0.30, \dots\}$, $\delta_2 = 2 \times \delta_1$ and $\delta_3 = 3 \times \delta_1$. Group i has mean $\boldsymbol{\mu}_i = \delta_i \mathbf{1}$. When δ_1 increases, the distance between the mean vectors increases. Figure 4.1 shows the power curve for clean MVN data with a balanced design where the groups have the same covariance matrices while figure 4.2 shows clean MVN data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 = 400$ and $n_3 = 600$. Figure 4.3 uses a mixture distribution. Figure 4.4 is similar to 4.2 except it uses a multivariate t_4 distribution. See the actual simulation results in Chapter 6.

4.3 REAL DATA EXAMPLE

The North Carolina Crime data consists of 630 observations on 24 variables. This data set is available online at <https://vincentarelbundock.github.io/Rdatasets/datasets.html>. Region is a categorical variable with three categories: Central, West and Other with the number of observations 232, 146 and 245 respectively, and forms the three groups. This example uses “wsta” - weekly wage of state employees, “avgsen” - average sentence days, “prbarr” - ‘probability’ of arrest, “prbconv” - ‘probability’ of conviction and “taxpc” - tax revenue per capita as variables. The test with the coordinatewise median had $D_0 = 4.086$ with the cutoff of 4.32 and failed to reject H_0 . The classical one-way MANOVA test had a p-value of 0.001 and rejected the null hypothesis.

The DD plots in figure 4.5 reveal a few outliers. Furthermore the boxplots in figure 4.6 and the scatterplot matrix in figure 4.7 shows that the data are highly skewed. Hence the location measures other than the median likely do differ.

See the simulation set up and the simulation results in Chapter 6.

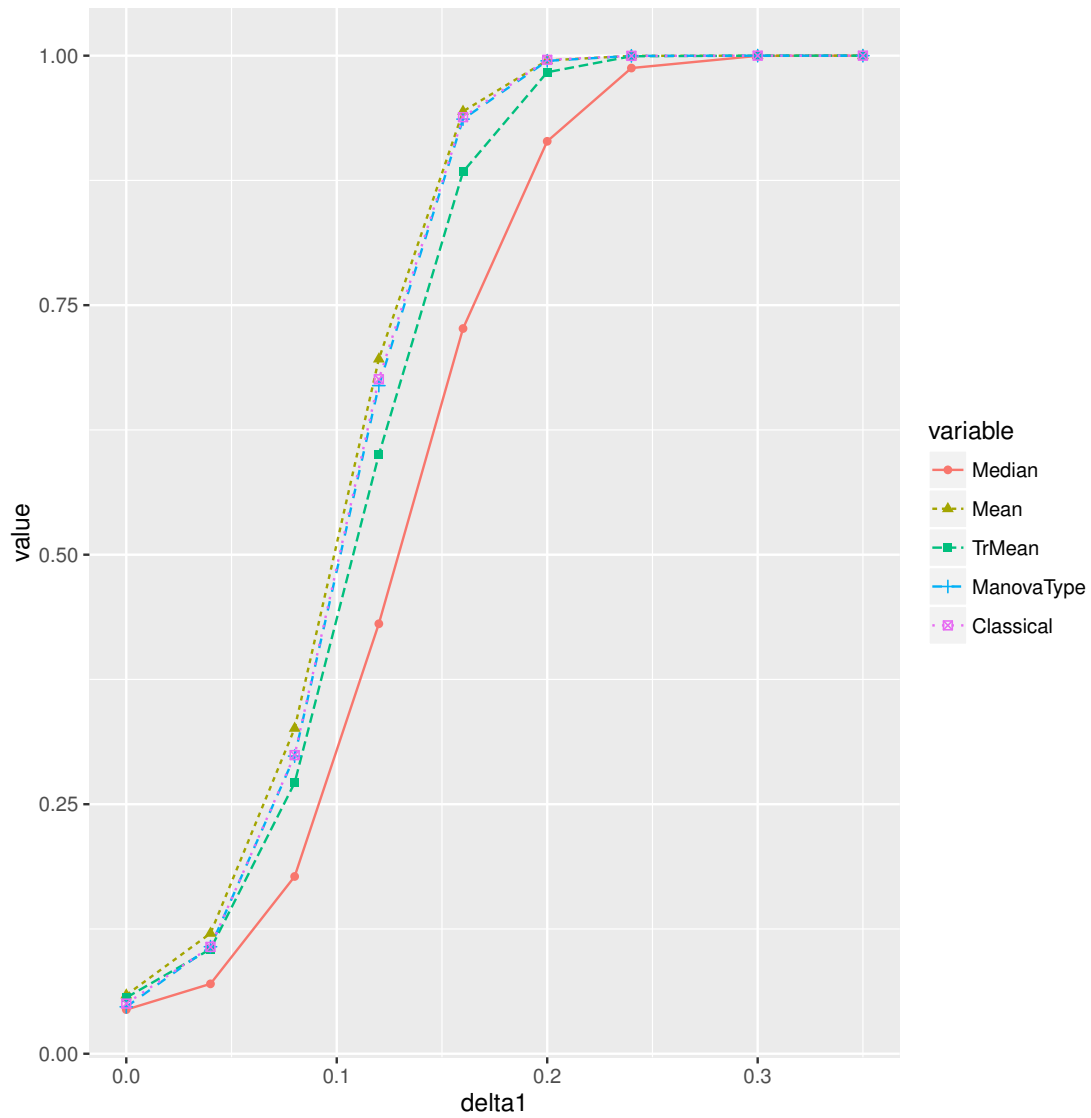


Figure 4.1. Power curve for clean MVN data with $m = 5, \sigma_1 = 1, \sigma_2 = 1, \sigma_3 = 1, n_1 = 200, n_2 = 200$ and $n_3 = 200$

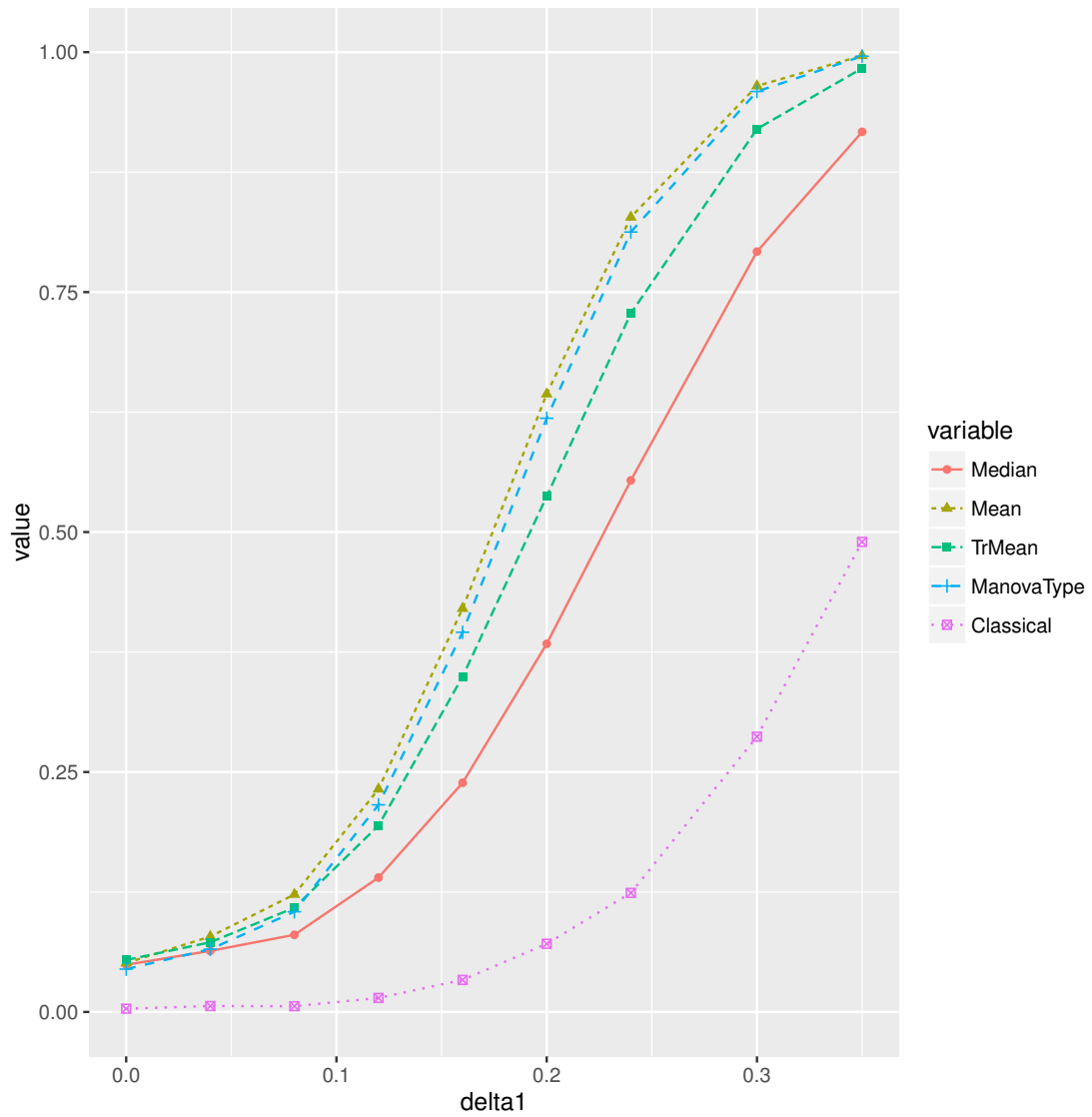


Figure 4.2. Power curve for clean MVN data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 = 400$ and $n_3 = 600$

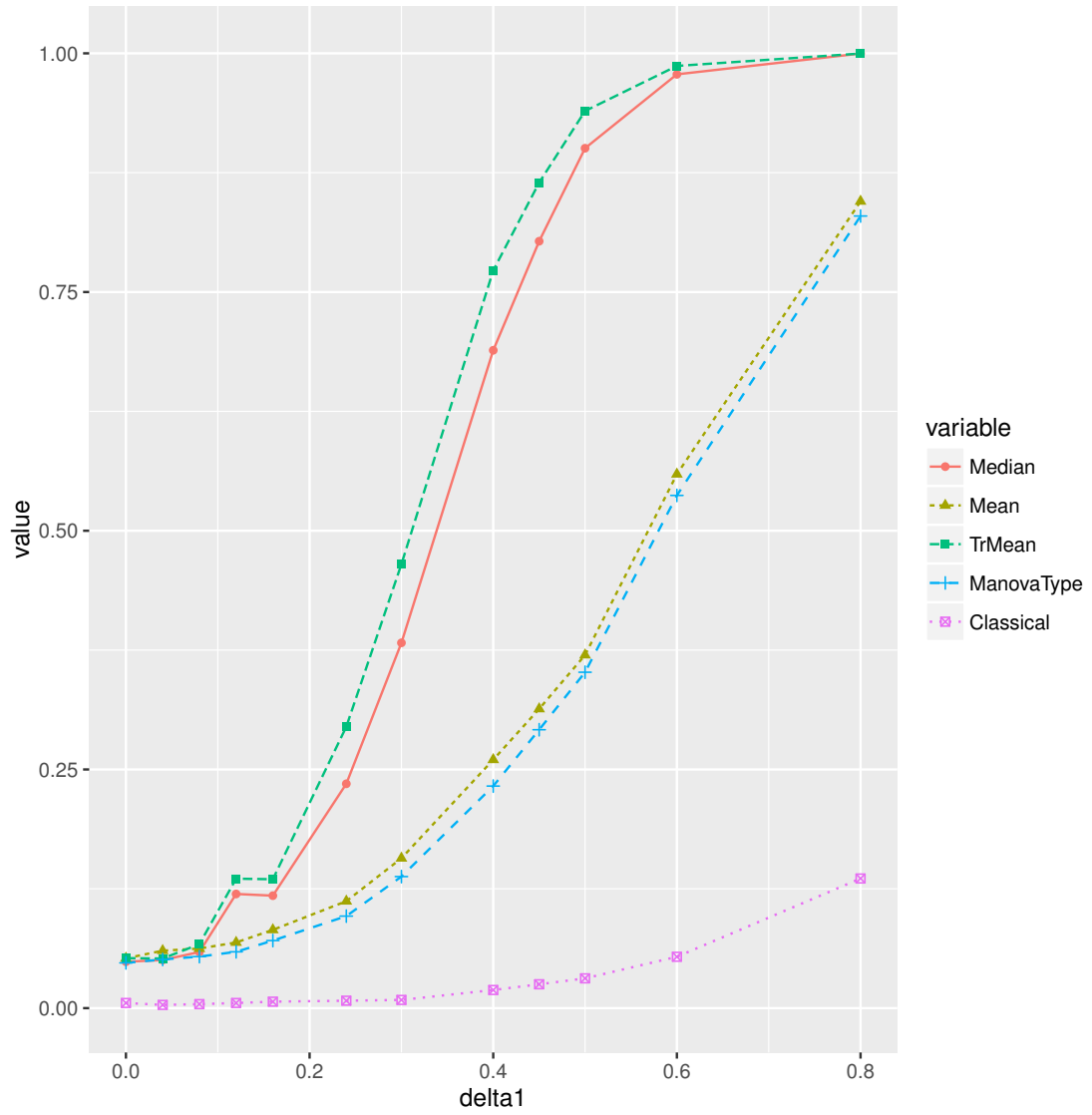


Figure 4.3. Power curve for clean Mixture data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 = 400$ and $n_3 = 600$

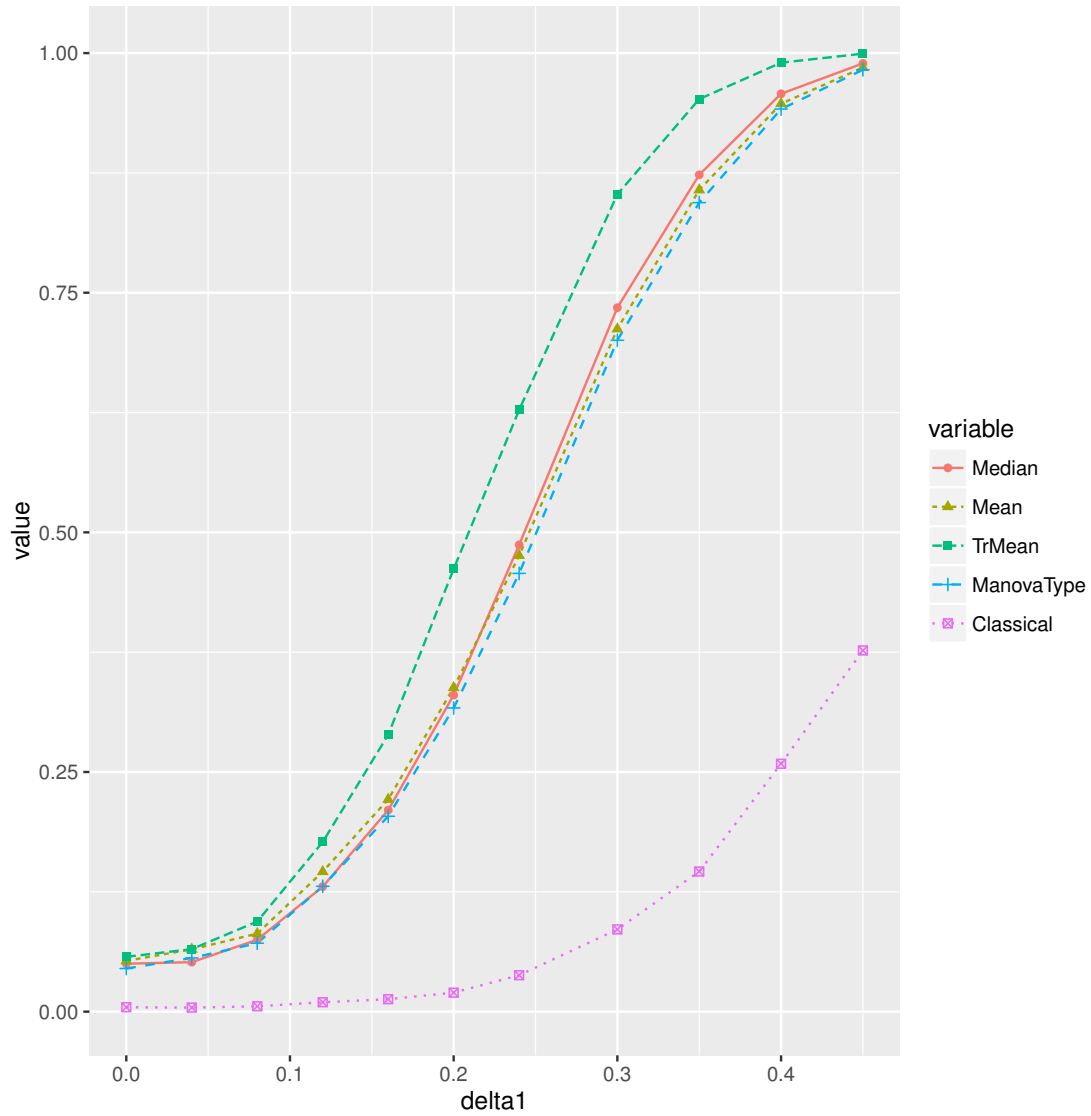


Figure 4.4. Power curve for clean multivariate t_4 data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 = 400$ and $n_3 = 600$

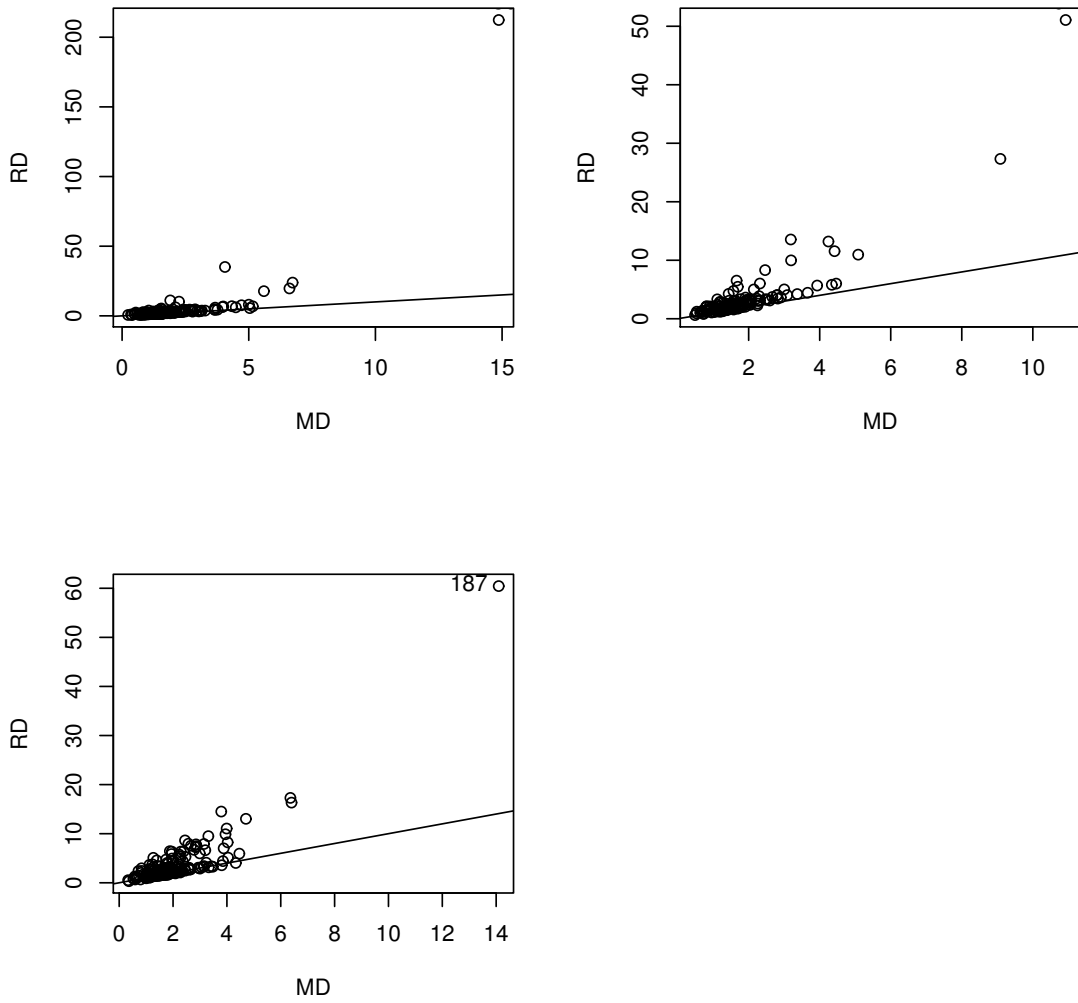


Figure 4.5. DD plots for Crime data

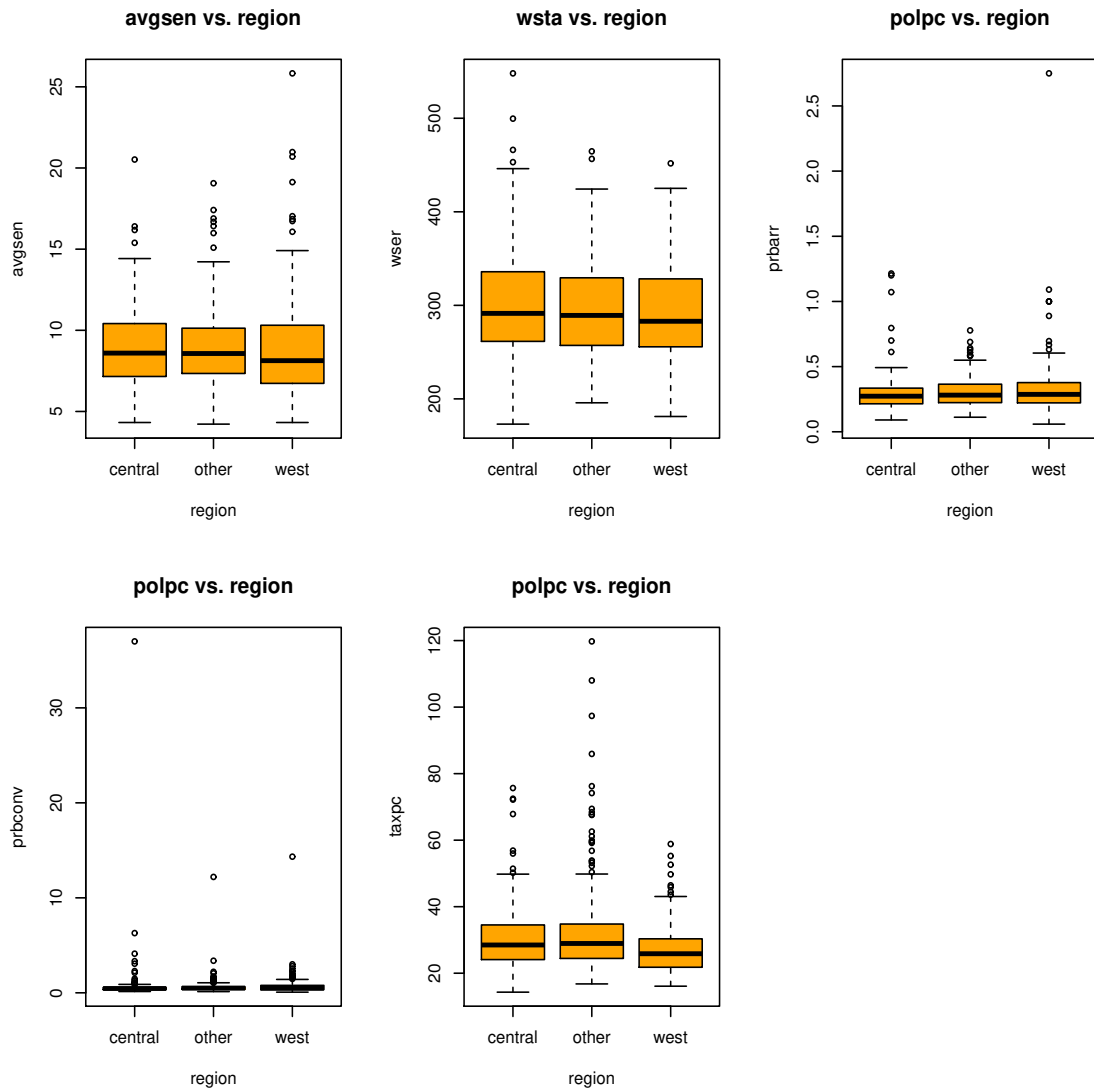


Figure 4.6. Side by side boxplots for Crime data

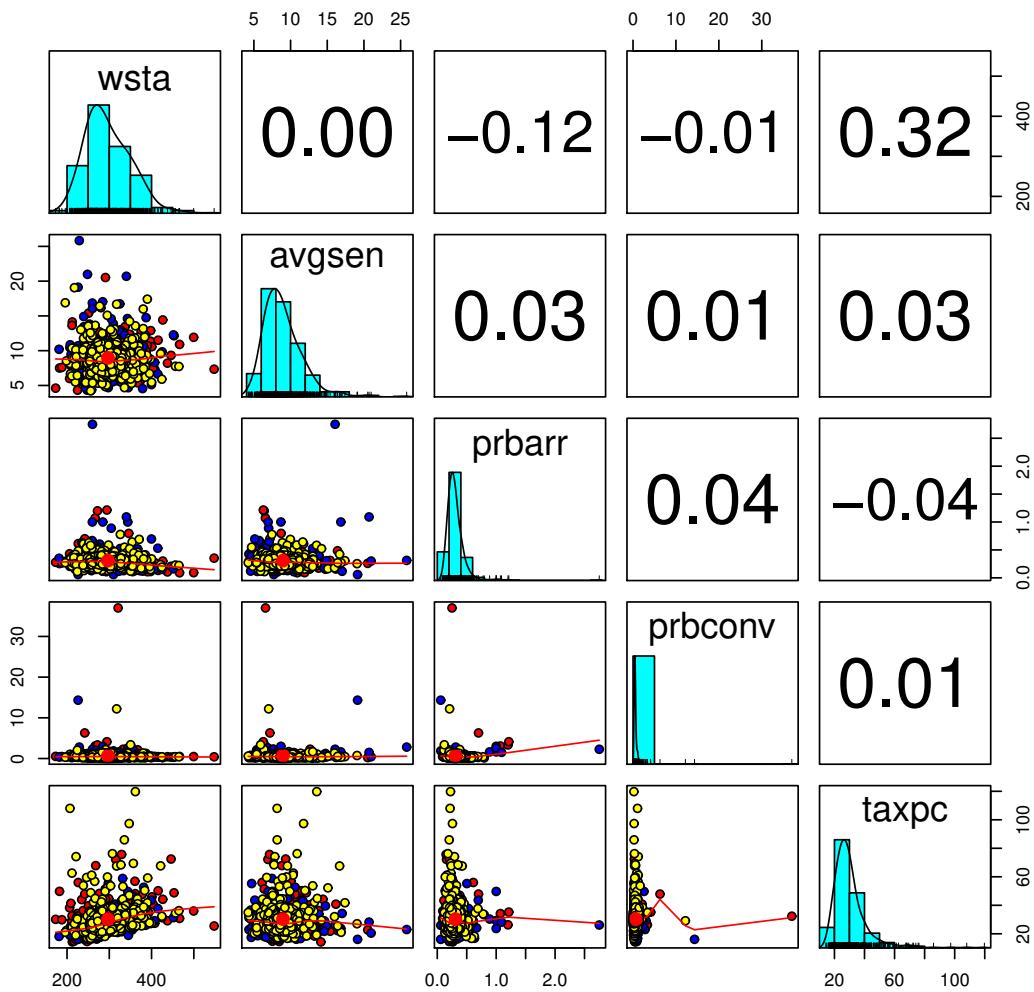


Figure 4.7. Scatterplot matrix for Crime data

CHAPTER 5
SIMULATIONS FOR BOOTSTRAPPING ANALOGS OF THE TWO
SAMPLE HOTELLING'S T^2 TEST

5.1 SIMULATION SETUP

The simulation used 5000 runs with B bootstrap samples. Olive (2017bc) suggests that the prediction region method can give good results when the number of bootstrap samples $B \geq 50p$ and if $n \geq 50p$, and the simulation used various values of B . See Rupasinghe Arachchige Don and Pelawa Watagoda (2017).

Four types of data distributions \mathbf{w}_i were considered that were identical for $i = 1, 2$. Then $\mathbf{x}_1 = \mathbf{A}\mathbf{w}_1 + \delta\mathbf{1}$ and $\mathbf{x}_2 = \sigma\mathbf{A}\mathbf{w}_2$ where $\mathbf{1} = (1, \dots, 1)^T$ is a vector of ones and $\mathbf{A} = \text{diag}(1, \sqrt{2}, \dots, \sqrt{p})$. The \mathbf{w}_i distributions were:

1. multivariate normal distribution $N_p(\mathbf{0}, \mathbf{I})$,
2. multivariate t distribution with 4 degrees of freedom,
3. mixture distribution $0.6N_p(\mathbf{0}, \mathbf{I}) + 0.4N_p(\mathbf{0}, 25\mathbf{I})$,
4. multivariate lognormal distribution shifted to have nonzero mean $\boldsymbol{\mu} = 0.649 \mathbf{1}$, but a population coordinatewise median of $\mathbf{0}$.

Note that $\text{Cov}(\mathbf{x}_2) = \sigma^2 \text{Cov}(\mathbf{x}_1)$, and for the first three distributions, $E(\mathbf{x}_i) = E(\mathbf{w}_i) = \mathbf{0}$ if $\delta = 0$.

Adding the same type and proportion of outliers to groups one and two often resulted in two distributions that were still similar. Hence outliers were added to the first group but not the second, making the covariance structures of the two groups quite different. The outlier proportion was $100\gamma\%$. Let $\mathbf{x}_1 = (x_{11}, \dots, x_{p1})^T$. The five outlier types for group 1 were:

1. type 1: a tight cluster at the major axis $(0, \dots, 0, z)^T$,
2. type 2: a tight cluster at the minor axis $(z, 0, \dots, 0)^T$,

3. type 3: a mean shift $N((z, \dots, z)^T, \text{diag}(1, \dots, p))$,
4. type 4: x_{1p} replaced by z , and
5. type 5: x_{11} replaced by z .

The quantity z determines how far the outliers are from the clean data.

Let the *coverage* be the proportion of times that H_0 is rejected. We want the *coverage* near 0.05 when H_0 is true and the coverage close to 1.0 for good power when H_0 is false. With 5000 runs, an observed *coverage* inside of (0.04, 0.06) suggests that the true *coverage* is close to the nominal 0.05 coverage when H_0 is true.

5.2 SIMULATION OUTPUT

5.2.1 Type I error rates simulation for clean data

Tables 5.1 - 5.6 were for clean elliptically contoured distributions (no outliers present), where H_0 is true and the different location estimators estimate $\boldsymbol{\mu} = \mathbf{0}$, the point of symmetry for the distribution. The chi-square cutoffs when $p = 5$ and $p = 15$ were 11.071 and 24.996, respectively. The *coverages* were often near the nominal value of 0.05, but the RMVN *coverages* were a bit low except for Table 5.4. The classical Hotelling's T^2 test does not use the bootstrap, and performed poorly when H_0 was true and both the sample sizes and the population covariance matrices were different.

For clean multivariate lognormal data, H_0 is true when $\sigma = 1$ (identical distributions for both groups), but H_0 is not true for the population mean when $\sigma = 2$. For $\sigma = 2$, the coordinatewise median had *coverages* near the nominal, while the sample mean had good power with *coverages* near 1. The RMVN coverage was a bit low when $\sigma = 1$ with power that was often less than that of the sample mean when $\sigma = 2$. See Table 5.7 and 5.8. The simulated cutoffs were quite similar to the chi-square cutoffs for Tables 5.1 through 5.8.

Table 5.1. *Coverages* for clean multivariate normal data $p = 5$

p	n_1	n_2	σ	B	Median	Mean	Tr.Mn	RMVN	Class
5	100	100	1	250	0.0418	0.0604	0.0546	0.0172	
				1000	0.0452	0.0678	0.0594	0.0198	
			2	250	0.0502	0.0706	0.0596	0.0258	
				1000	0.0470	0.0684	0.0638	0.0220	
	250	250	1	250	0.0470	0.0554	0.0568	0.0402	0.0560
				1000	0.0440	0.0606	0.0540	0.0414	
			2	250	0.0472	0.0550	0.0574	0.0422	0.0498
				1000	0.0420	0.0568	0.0538	0.0392	
100	200	1	250	0.0446	0.0670	0.0600	0.0228		
			1000	0.0434	0.0614	0.0582	0.0254		
		2	250	0.0488	0.0610	0.0568	0.0292		
			1000	0.0422	0.0518	0.0532	0.0234		
250	500	1	250	0.0490	0.0524	0.0496	0.0394	0.0552	
			1000	0.0462	0.0588	0.0584	0.0448		
		2	250	0.0460	0.0540	0.0524	0.0436	0.0070	
			1000	0.0470	0.0500	0.0534	0.0386		

Table 5.2. *Coverages* for clean multivariate normal data $p = 15$

p	n_1	n_2	σ	B	Median	Mean	Tr.Mn	RMVN	Class	
15	300	300	1	750	0.0454	0.0666	0.0622	0.0234		
				1000	0.0378	0.0578	0.0554	0.0256		
		750	750	1	750	0.0462	0.0626	0.0622	0.0466	0.0450
					1000	0.0390	0.0514	0.0470	0.0378	
	750	750	2	750	0.0484	0.0752	0.0674	0.0270		
				1000	0.0576	0.0730	0.0732	0.0296		
		1500	1500	2	750	0.0492	0.0598	0.0608	0.0464	0.0516
					1000	0.0474	0.0556	0.0568	0.0446	
300	600	600	1	750	0.0424	0.0650	0.0658	0.0286		
				1000	0.0440	0.0638	0.0592	0.0308		
		1500	1500	1	750	0.0466	0.0538	0.0550	0.0466	0.0480
					1000	0.0492	0.0556	0.0548	0.0444	
	1500	1500	2	750	0.0424	0.0538	0.0520	0.0454	0.0014	
				1000	0.0514	0.0532	0.0542	0.0426		

Table 5.3. *Coverages* for clean $0.6N_p(\mathbf{0}, \mathbf{I}) + 0.4N_p(\mathbf{0}, 25\mathbf{I})$ data $p = 5$

p	n_1	n_2	σ	B	Median	Mean	Tr.Mn	RMVN	Class
5	100	100	1	250	0.0294	0.0620	0.0388	0.0158	
				1000	0.0390	0.0544	0.0420	0.0130	
			2	250	0.0400	0.0606	0.0416	0.0184	
				1000	0.0422	0.0612	0.0386	0.0162	
	250	250	1	250	0.0420	0.0560	0.0480	0.0394	0.0462
				1000	0.0386	0.0532	0.0464	0.0336	
			2	250	0.0454	0.0550	0.0476	0.0416	0.0476
				1000	0.0370	0.0484	0.0400	0.0368	
100	200	1	250	0.0364	0.0546	0.0398	0.0190		
			1000	0.0344	0.0632	0.0394	0.0222		
		2	250	0.0372	0.0604	0.0462	0.0238		
			1000	0.0346	0.0616	0.0402	0.0228		
	250	500	1	250	0.0460	0.0542	0.0538	0.0416	0.0470
				1000	0.0368	0.0502	0.0416	0.0404	
			2	250	0.0480	0.0600	0.0474	0.0390	0.0060
				1000	0.0416	0.0598	0.0498	0.0416	

Table 5.4. *Coverages* for clean $0.6N_p(\mathbf{0}, \mathbf{I}) + 0.4N_p(\mathbf{0}, 25\mathbf{I})$ data $p = 15$

p	n_1	n_2	σ	B	Median	Mean	Tr.Mn	RMVN	Class
15	300	300	1	750	0.0414	0.0598	0.0490	0.0428	
				1000	0.0402	0.0592	0.0484	0.0414	
			2	750	0.0426	0.0620	0.0496	0.0502	
				1000	0.0414	0.0600	0.0448	0.0496	
	750	750	1	750	0.0434	0.0536	0.0540	0.0448	0.0496
				1000	0.0406	0.0598	0.0474	0.0396	
			2	750	0.0468	0.0626	0.0518	0.0456	0.0464
				1000	0.0456	0.0566	0.0490	0.0454	
300	600	1	750	0.0418	0.0582	0.0464	0.0474		
			1000	0.0430	0.0684	0.0514	0.0466		
		2	750	0.0394	0.0578	0.0466	0.0432		
			1000	0.0356	0.0606	0.0470	0.0422		
750	1500	1	750	0.0456	0.0584	0.0568	0.0488	0.0502	
			1000	0.0426	0.0550	0.0478	0.0438		
		2	750	0.0456	0.0576	0.0508	0.0442	0.0004	
			1000	0.0416	0.0572	0.0488	0.0510		

Table 5.5. *Coverages* for clean multivariate t_4 data $p = 5$

p	n_1	n_2	σ	B	Median	Mean	Tr.Mn	RMVN	Class
5	100	100	1	250	0.0478	0.0610	0.0592	0.0178	
				1000	0.0358	0.0548	0.0518	0.0164	
		2	250	0.0514	0.0608	0.0632	0.0238		
			1000	0.0444	0.0512	0.0558	0.0162		
	250	250	1	250	0.0442	0.0574	0.0570	0.0266	0.0456
				1000	0.0426	0.0570	0.0530	0.0282	
		2	250	0.0496	0.0618	0.0614	0.0328	0.0542	
			1000	0.0480	0.0558	0.0578	0.0292		
100	200	1	250	0.0432	0.0556	0.0576	0.0212		
			1000	0.0372	0.0552	0.0522	0.0200		
	2	250	0.0414	0.0586	0.0570	0.0232			
		1000	0.0446	0.0546	0.0568	0.0262			
250	500	1	250	0.0484	0.0512	0.0540	0.0346	0.0504	
			1000	0.0420	0.0488	0.0494	0.0310		
	2	250	0.0408	0.0580	0.0526	0.0348	0.0058		
		1000	0.0410	0.0492	0.0510	0.0348			

Table 5.6. *Coverages* for clean multivariate t_4 data $p = 15$

p	n_1	n_2	σ	B	Median	Mean	Tr.Mn	RMVN	Class
15	300	300	1	750	0.0392	0.0546	0.0562	0.0158	
				1000	0.0480	0.0590	0.0662	0.0140	
			2	750	0.0478	0.0572	0.0604	0.0134	
				1000	0.0512	0.0632	0.0640	0.0148	
	750	750	1	750	0.0470	0.0550	0.0562	0.0232	0.0414
				1000	0.0382	0.0526	0.0476	0.0228	
			2	750	0.0472	0.0572	0.0542	0.0248	0.0442
				1000	0.0502	0.0496	0.0556	0.0258	
	300	600	1	750	0.0448	0.0554	0.0598	0.0158	
				1000	0.0458	0.0602	0.0616	0.0184	
			2	750	0.0450	0.0564	0.0558	0.0178	
				1000	0.0400	0.0498	0.0546	0.0196	
750	1500	1	750	0.0482	0.0556	0.0528	0.0224	0.0446	
			1000	0.0464	0.0496	0.0528	0.0254		
		2	750	0.0442	0.0534	0.0502	0.0314	0.0016	
			1000	0.0452	0.0508	0.0554	0.0262		

Table 5.7. *Coverages* for clean lognormal data $p = 5$

p	n_1	n_2	σ	B	Median	Mean	Tr.Mn	RMVN	Class
5	100	100	1	250	0.0330	0.0462	0.0470	0.0096	
				1000	0.0290	0.0508	0.0390	0.0088	
			2	250	0.0442	0.6170	0.0600	0.0340	
				1000	0.0368	0.6200	0.0570	0.0352	
	250	250	1	250	0.0408	0.0460	0.0514	0.0274	0.0470
				1000	0.0388	0.0494	0.0474	0.0254	
			2	250	0.0436	0.9816	0.0858	0.1108	0.9968
				1000	0.0398	0.9846	0.0788	0.1168	
100	200	1	250	0.0346	0.0684	0.0492	0.0204		
			1000	0.0320	0.0548	0.0434	0.0138		
		2	250	0.0400	0.8898	0.0644	0.0692		
			1000	0.0428	0.8962	0.0700	0.0730		
250	500	1	250	0.0398	0.0540	0.0496	0.0316	0.0472	
			1000	0.0368	0.0588	0.0446	0.0292		
		2	250	0.0418	0.9998	0.1192	0.2492	0.9964	
			1000	0.0424	0.9994	0.1158	0.2520		

Table 5.8. *Coverages* for clean lognormal data $p = 15$

p	n_1	n_2	σ	B	Median	Mean	Tr.Mn	RMVN	Class
15	300	300	1	750	0.0326	0.0546	0.0478	0.0106	
				1000	0.0364	0.0558	0.0502	0.0120	
			2	750	0.0446	1.0000	0.1408	0.8530	
				1000	0.0474	1.0000	0.1450	0.8680	
	750	750	1	750	0.0402	0.0506	0.0480	0.0216	0.0502
				1000	0.0410	0.0444	0.0490	0.0238	
			2	750	0.0506	1.0000	0.3670	1.0000	1.0000
				1000	0.0510	1.0000	0.3748	1.0000	
300	600	1	750	0.0422	0.0684	0.0546	0.0188		
			1000	0.0406	0.0736	0.0560	0.0172		
			2	750	0.0396	1.0000	0.2344	0.9984	
				1000	0.0408	1.0000	0.2402	0.9990	
	750	1500	1	750	0.0420	0.0580	0.0514	0.0258	0.0514
				1000	0.0478	0.0558	0.0608	0.0284	
			2	750	0.0446	1.0000	0.6110	1.0000	1.0000
				1000	0.0464	1.0000	0.6256	1.0000	

5.2.2 Type I error rates simulation for contaminated data

Table 5.9 illustrates the simulated results where group 1 had outliers. The coordinate-wise median worked with a little higher type I error rate (around 0.08) than the nominal level of 0.05 for the mixture, multivariate t, and multivariate log normal distributions, but failed for the multivariate normal data when $\gamma = 0.4$. The sample mean (classical and bootstrap) and 25% trimmed mean failed to achieve the nominal level with any of the distributions used when H_0 was true for the clean data. The RMVN estimator worked with all four distributions with a better type I error rate compared to the other estimators. The chi-square cutoff was 9.488 since $p = 4$.

The coordinatewise median can achieve better coverages for smaller proportions of outliers with higher values of z (not shown in the tables), i.e. the outliers had to be far from the clean data compared to the RMVN estimator. The RMVN estimator can handle higher proportions of outliers as shown in the Table 5.9.

Table 5.9. *Coverages and cutoffs with outliers: $p = 4, n_1 = n_2 = 200, B = 200$*

Dist.	Otype	γ	z		Med	Mean	Tr.Me	RMVN	Class
MVN	1	0.4	10	Cov	0.6946	1.0000	1.0000	0.0330	1.0000
				cut	10.158	9.769	9.798	10.701	
	2	0.4	20	Cov	0.5232	1.0000	1.0000	0.0382	1.0000
				cut	9.836	9.776	9.809	9.268	
	3	0.4	20	Cov	0.8578	1.0000	1.0000	0.0402	1.0000
				cut	10.214	9.761	9.760	9.288	
	4	0.1	10	Cov	0.0980	0.8654	0.1450	0.0382	0.8684
				cut	9.898	9.771	9.777	9.851	
Mix	2	0.4	20	Cov	0.0828	1.0000	1.0000	0.0144	1.0000
				cut	10.542	9.788	9.878	11.300	
	5	0.1	10	Cov	0.0820	0.5306	0.1228	0.0184	0.5276
				cut	9.933	9.779	9.881	11.056	
MVT	1	0.4	10	Cov	0.0854	0.6700	0.1548	0.0204	1.0000
				cut	10.232	9.799	9.787	10.200	
	5	0.1	20	Cov	0.0864	1.0000	0.1418	0.0304	1.0000
				cut	9.924	9.795	9.795	9.830	
Log	3	0.4	20	Cov	0.0778	1.0000	1.0000	0.0162	1.0000
				cut	13.689	9.822	9.827	12.607	
	4	0.1	10	Cov	0.0842	0.3158	0.1482	0.0234	0.3044
				cut	10.013	9.875	9.872	10.416	

5.2.3 Power Simulation

In the power simulation, $\delta > 0$ was used. Hence for the first three distributions $\boldsymbol{\mu}_2 = \mathbf{0}$ and $\boldsymbol{\mu}_1 = \delta(1, \dots, 1)^T$. Then the Euclidean distance between the two means was $\sqrt{p}\delta$, where p is the number of parameters. Therefore the distance increases as p increase. The value of δ had to be fairly small so that the simulated power was not always 1. Also see Table 5.8 with $\sigma = 2$.

For Table 5.10, the sample mean (bootstrap and classical) had the best power while the sample median had the worst power. For Table 5.11, the RMVN estimator had the best power while the sample mean has the worst power. The trimmed mean had the best power for Table 5.12. For Table 5.13, the RMVN estimator had poor power when $p = 5$, $n = 250$, and $\sigma = 2$. No method was always best or worst.

Table 5.10. *Coverages* when H_0 is false for MVN data.

p	$n_1 = n_2$	σ	B	δ	Med	Mean	Tr.Me	RMVN	Class
5	250	1	250	0.35	0.9598	0.9990	0.9928	0.9942	0.9988
			1000	0.35	0.9684	0.9994	0.9970	0.9978	
	1000	2	250	0.35	0.5958	0.8442	0.7672	0.7604	0.8402
			1000	0.35	0.5832	0.8346	0.7438	0.7470	
15	750	1	750	0.15	0.7394	0.9552	0.9012	0.9268	0.9556
			1000	0.15	0.7474	0.9522	0.8984	0.9178	
	1000	2	750	0.15	0.3078	0.5318	0.4550	0.4468	0.5156
			1000	0.15	0.3118	0.5218	0.4430	0.4464	

Table 5.11. *Coverages* when H_0 is false for mixture data.

p	$n_1 = n_2$	σ	B	δ	Med	Mean	Tr.Me	RMVN	Class
5	250	1	250	0.45	0.8826	0.4062	0.9304	0.9938	0.4032
			1000	0.45	0.8858	0.4058	0.9338	0.9948	
		2	250	0.45	0.4458	0.1910	0.5222	0.7454	0.1642
			1000	0.45	0.4656	0.1890	0.5386	0.7626	
15	750	1	750	0.20	0.6204	0.2274	0.7148	0.9492	0.2114
			1000	0.20	0.6316	0.2228	0.7190	0.9494	
		2	750	0.20	0.2318	0.1154	0.2894	0.5034	0.1042
			1000	0.20	0.2438	0.1092	0.2916	0.4980	

Table 5.12. *Coverages* when H_0 is false for multivariate t_4 data.

p	$n_1 = n_2$	σ	B	δ	Med	Mean	Tr.Me	RMVN	Class
5	250	1	250	0.38	0.9642	0.9562	0.9916	0.9878	0.9548
			1000	0.38	0.9728	0.9572	0.9944	0.9880	
		2	250	0.38	0.5958	0.5960	0.7198	0.6488	0.6074
			1000	0.38	0.6188	0.6152	0.7490	0.6636	
15	750	1	750	0.20	0.9418	0.9270	0.9868	0.9714	0.9232
			1000	0.20	0.9422	0.9304	0.9860	0.9724	
		2	750	0.20	0.4934	0.4932	0.6422	0.5384	0.4754
			1000	0.20	0.4842	0.4916	0.6362	0.5252	

Table 5.13. *Coverages* when H_0 is false for lognormal data.

p	$n_1 = n_2$	σ	B	δ	Median	Mean	Tr.Me	RMVN	Class
5	250	1	250	0.45	0.9982	0.8256	0.9994	0.879	0.8208
			1000	0.45	0.9980	0.8324	0.9996	0.883	
		2	250	0.45	0.8210	0.4704	0.6488	0.0914	0.4630
			1000	0.45	0.8378	0.4646	0.6624	0.1038	
15	750	1	750	0.30	1.0000	0.9186	1.0000	0.8514	0.9120
			1000	0.30	1.0000	0.9178	1.0000	0.8544	
		2	750	0.30	0.9436	1.0000	0.5042	0.9438	1.0000
			1000	0.30	0.9484	1.0000	0.5022	0.9424	

CHAPTER 6

SIMULATIONS WITH THREE SAMPLES FOR BOOTSTRAPPING ANALOGS OF THE ONE-WAY MANOVA TEST

6.1 SIMULATION SETUP

The simulation used 5000 runs with B bootstrap samples and $p = 3$ groups. Olive (2017bc) suggests that the prediction region method can give good results when the number of bootstrap samples $B \geq 50m(p-1)$ and if $n \geq 50m(p-1)$, and the simulation used various values of B . The sample mean, coordinatewise median, and coordinatewise 25% trimmed mean were the statistics T used. The classical one way MANOVA Hotelling Lawley test statistic was also used.

Four types of data distributions \mathbf{w}_i were considered that were identical for $i = 1, 2$ and 3. Then $\mathbf{y}_1 = \mathbf{A}\mathbf{w}_1 + \delta_1\mathbf{1}$, $\mathbf{y}_2 = \sigma_2\mathbf{A}\mathbf{w}_2$, and $\mathbf{y}_3 = \sigma_3\mathbf{A}\mathbf{w}_3 + \delta_3\mathbf{1}$ or $\mathbf{y}_3 = \mathbf{w}_3$ where $\mathbf{1} = (1, \dots, 1)^T$ is a vector of ones and $\mathbf{A} = \text{diag}(1, \sqrt{2}, \dots, \sqrt{m})$. The \mathbf{w}_i distributions were the multivariate normal distribution $N_m(\mathbf{0}, \mathbf{I})$, the mixture distribution $0.6N_m(\mathbf{0}, \mathbf{I}) + 0.4N_m(\mathbf{0}, 25\mathbf{I})$, the multivariate t distribution with 4 degrees of freedom, and the multivariate lognormal distribution shifted to have nonzero mean $\boldsymbol{\mu} = 0.649 \mathbf{1}$, but a population coordinatewise median of $\mathbf{0}$. Note that $\text{Cov}(\mathbf{y}_2) = \sigma_2^2 \text{Cov}(\mathbf{y}_1)$, and for the first three distributions, $E(\mathbf{y}_i) = E(\mathbf{w}_i) = \mathbf{0}$ if $\delta_1 = \delta_3 = 0$. If $\mathbf{y}_3 = \mathbf{w}_3$ then $\text{Cov}(\mathbf{y}_3) = c\mathbf{I}_m$ for some constant $c > 0$. If $\mathbf{y}_3 = \sigma_3\mathbf{A}\mathbf{w}_3 + \delta_3\mathbf{1}$, then $\text{Cov}(\mathbf{y}_3) = \sigma_3^2 \text{Cov}(\mathbf{y}_1)$.

Adding the same type and proportion of outliers to all three groups often resulted in three distributions that were still similar. Hence outliers were added to the first group but not the second or third, making the covariance structures of the three groups quite different. The outlier proportion was $100\gamma\%$. Let $\mathbf{y}_1 = (y_{11}, \dots, y_{m1})^T$. The five outlier types for group 1 were type 1: a tight cluster at the major axis $(0, \dots, 0, z)^T$, type 2: a tight cluster at the minor axis $(z, 0, \dots, 0)^T$, type 3: $N((z, \dots, z)^T, \text{diag}(1, \dots, m))$, type 4: y_{m1} replaced by z , and type 5: y_{11} replaced by z . The quantity z determines how far the outliers are from the clean data.

Let the *coverage* be the proportion of times that H_0 is rejected. We want the *coverage*

near 0.05 when H_0 is true and the coverage close to 1.0 for good power when H_0 is false. With 5000 runs, an observed *coverage* inside of (0.04, 0.06) suggests that the true *coverage* is close to the nominal 0.05 coverage when H_0 is true.

6.1.1 Simulations for type I error with clean data

Tables 6.1 through 6.8 show simulation results for all for distributions with various covariance settings. We took $\delta_1 = \delta_3 = 0$, B : the number of bootstrap steps used also takes on different values throughout the simulation. Balanced and unbalanced designs have also been considered. For tables 6.5-6.8 $\sigma_2 = \sigma_3 = 1$. According to the tables, the new tests work well with all the distributions and with different covariance settings. The new tests could handle unbalanced designs as well. The classical test works well with the multivariate normal data and when the covariance matrices are the same, but the type I error is higher than the nominal level for different covariance settings. The classical test can be too conservative when the design is unbalanced. Having an unbalanced design and different covariance matrices is the worst case scenario for the classical test regardless of the data distribution.

Table 6.1. Type I error for clean MVN data with $\text{cov3I} = \text{F}$

m	n_1	n_2	n_3	B	σ_2	σ_3	Median	Mean	Tr.Mn	Class
5	200	200	200	400	1	1	0.0422	0.0562	0.0552	0.0460
				1000	1	1	0.0486	0.0602	0.0598	0.0510
				400	2	3	0.0506	0.0670	0.0606	0.0680
				1000	2	3	0.0482	0.0580	0.0590	0.0680
5	200	400	600	400	1	1	0.0506	0.0542	0.0598	0.0474
				1000	1	1	0.0492	0.0542	0.0554	0.0472
				400	2	3	0.0474	0.0580	0.0576	0.0066
				1000	2	3	0.0532	0.0626	0.0618	0.0074
10	400	400	400	800	1	1	0.0508	0.0724	0.0712	0.0558
				2000	1	1	0.0516	0.0652	0.0644	0.0526
				800	2	3	0.0562	0.0640	0.0686	0.0656
				2000	2	3	0.0554	0.0624	0.0630	0.0704
10	400	800	1200	800	1	1	0.0510	0.0594	0.0626	0.0456
				2000	1	1	0.0470	0.0578	0.0576	0.0494
				800	2	3	0.0468	0.0576	0.0572	0.0008
				2000	2	3	0.0440	0.0574	0.0534	0.0034
20	800	800	800	1600	1	1	0.0474	0.0724	0.0652	0.0496
				4000	1	1	0.0504	0.0662	0.0668	0.0494
				1600	2	3	0.0566	0.0728	0.0618	0.0772
				4000	2	3	0.0592	0.0644	0.0672	0.0638
20	800	1600	2400	1600	1	1	0.0562	0.0644	0.0648	0.0492
				4000	1	1	0.0504	0.0564	0.0618	0.0462
				1600	2	3	0.0530	0.0654	0.0650	0.0000
				4000	2	3	0.0472	0.0632	0.0620	0.0008

Table 6.2. Type I error for clean Mixture data with $\text{cov3I} = \text{F}$

m	n_1	n_2	n_3	B	σ_2	σ_3	Median	Mean	Tr.Mn	Class
5	200	200	200	400	1	1	0.0406	0.0526	0.0470	0.0446
				1000	1	1	0.0396	0.0600	0.0476	0.0492
				400	2	3	0.0448	0.0538	0.0484	0.0682
				1000	2	3	0.0418	0.0520	0.0412	0.0692
5	200	400	600	400	1	1	0.0440	0.0562	0.0438	0.0476
				1000	1	1	0.0376	0.0528	0.0486	0.0504
				400	2	3	0.0432	0.0518	0.0502	0.0082
				1000	2	3	0.0392	0.0528	0.0454	0.0060
10	400	400	400	800	1	1	0.0446	0.0604	0.0516	0.0498
				2000	1	1	0.0438	0.0592	0.0496	0.0502
				800	2	3	0.0454	0.0598	0.0478	0.0694
				2000	2	3	0.0460	0.0586	0.0468	0.0664
10	400	800	1200	800	1	1	0.0448	0.0590	0.0458	0.0494
				2000	1	1	0.0412	0.0590	0.0512	0.0532
				800	2	3	0.0490	0.0600	0.0528	0.0036
				2000	2	3	0.0444	0.0524	0.0464	0.0020
20	800	800	800	1600	1	1	0.0476	0.0628	0.0530	0.0472
				4000	1	1	0.0462	0.0606	0.0498	0.0490
				1600	2	3	0.0500	0.0680	0.0522	0.0738
				4000	2	3	0.0468	0.0676	0.0516	0.0720
20	800	1600	2400	1600	1	1	0.0522	0.0618	0.0560	0.0510
				4000	1	1	0.0480	0.0600	0.0504	0.0520
				1600	2	3	0.0488	0.0564	0.0566	0.0004
				4000	2	3	0.0432	0.0590	0.0476	0.0004

Table 6.3. Type I error for clean multivariate t data with cov3I = F

m	n_1	n_2	n_3	B	σ_2	σ_3	Median	Mean	Tr.Mn	Class
5	200	200	200	400	1	1	0.0420	0.0566	0.0524	0.0450
				1000	1	1	0.0388	0.0528	0.0542	0.0488
				400	2	3	0.0534	0.0596	0.0636	0.0666
				1000	2	3	0.0450	0.0564	0.0606	0.0650
5	200	400	600	400	1	1	0.0512	0.0532	0.0578	0.0478
				1000	1	1	0.0432	0.0596	0.0536	0.0526
				400	2	3	0.0458	0.0516	0.0556	0.0062
				1000	2	3	0.0464	0.0564	0.0578	0.0080
10	400	400	400	800	1	1	0.0448	0.0622	0.0588	0.0480
				2000	1	1	0.0490	0.0578	0.0608	0.0488
				800	2	3	0.0498	0.0622	0.0654	0.0680
				2000	2	3	0.0528	0.0576	0.0576	0.0652
10	400	800	1200	800	1	1	0.0470	0.0572	0.0530	0.0496
				2000	1	1	0.0448	0.0626	0.0544	0.0558
				800	2	3	0.0412	0.0508	0.0536	0.0030
				2000	2	3	0.0514	0.0582	0.0568	0.0026
20	800	800	800	1600	1	1	0.0474	0.0620	0.0576	0.0454
				4000	1	1	0.0532	0.0604	0.0634	0.0490
				1600	2	3	0.0534	0.0652	0.0640	0.0714
				4000	2	3	0.0556	0.0644	0.0652	0.0714
20	800	1600	2400	1600	1	1	0.0568	0.0630	0.0618	0.0504
				4000	1	1	0.0492	0.0584	0.0560	0.0532
				1600	2	3	0.0546	0.0570	0.0658	0.0008
				4000	2	3	0.0488	0.0544	0.0612	0.0004

Table 6.4. Type I error for clean lognormal data with cov3I = F

m	n_1	n_2	n_3	B	σ_2	σ_3	Median	Mean	Tr.Mn	Class
5	200	200	200	400	1	1	0.0368	0.0628	0.0478	0.0436
				1000	1	1	0.0402	0.0596	0.0486	0.0452
				400	2	3	0.0432	0.9996	0.1004	0.9994
				1000	2	3	0.0448	1.0000	0.0980	0.9996
5	200	400	600	400	1	1	0.0446	0.0768	0.0568	0.0476
				1000	1	1	0.0426	0.0724	0.0530	0.0530
				400	2	3	0.0428	1.0000	0.2068	1.0000
				1000	2	3	0.0428	1.0000	0.2002	1.0000
10	400	400	400	800	1	1	0.0450	0.0658	0.0622	0.0450
				2000	1	1	0.0472	0.0716	0.0542	0.0498
				800	2	3	0.0532	1.0000	0.2858	1.0000
				2000	2	3	0.0458	1.0000	0.2706	1.0000
10	400	800	1200	800	1	1	0.0434	0.0754	0.0542	0.0546
				2000	1	1	0.0502	0.0708	0.0526	0.0462
				800	2	3	0.0438	1.0000	0.6448	1.0000
				2000	2	3	0.0372	1.0000	0.6394	1.0000
20	800	800	800	1600	1	1	0.0482	0.0680	0.0580	0.0470
				4000	1	1	0.0412	0.0678	0.0582	0.0486
				1600	2	3	0.0530	1.0000	0.8714	1.0000
				4000	2	3	0.0516	1.0000	0.8622	1.0000
20	800	1600	2400	1600	1	1	0.0470	0.0756	0.0648	0.0532
				4000	1	1	0.0520	0.0684	0.0652	0.0464
				1600	2	3	0.0480	1.0000	0.9980	1.0000
				4000	2	3	0.0442	1.0000	0.9988	1.0000

Table 6.5. Type I error for clean MVN data with $\text{cov3I} = \text{T}$

m	n_1	n_2	n_3	B	Median	Mean	Tr.Mn	Class
5	200	200	200	400	0.0482	0.0682	0.0638	0.0650
				1000	0.0500	0.0684	0.0610	0.0592
5	200	400	600	400	0.0566	0.0604	0.0648	0.1354
				1000	0.0472	0.0526	0.0534	0.1278
10	400	400	400	800	0.0512	0.0636	0.0610	0.0604
				2000	0.0506	0.0608	0.0632	0.0584
10	400	800	1200	800	0.0570	0.0658	0.0642	0.2422
				2000	0.0536	0.0536	0.0536	0.2224
20	800	800	800	1600	0.0662	0.0740	0.0734	0.0638
				4000	0.0562	0.0668	0.0600	0.0604
20	800	1600	2400	1600	0.0566	0.0638	0.0628	0.4308
				4000	0.0560	0.0702	0.0658	0.4308

Table 6.6. Type I error for clean Mixture data with $\text{cov3I} = \text{T}$

m	n_1	n_2	n_3	B	Median	Mean	Tr.Mn	Class
5	200	200	200	400	0.0424	0.0614	0.0438	0.0570
				1000	0.0446	0.0618	0.0460	0.0550
5	200	400	600	400	0.0524	0.0572	0.0542	0.1284
				1000	0.0434	0.0540	0.0498	0.1270
10	400	400	400	800	0.0422	0.0620	0.0542	0.0598
				2000	0.0450	0.0638	0.0524	0.0642
10	400	800	1200	800	0.0468	0.0558	0.0484	0.2368
				2000	0.0522	0.0548	0.0518	0.2356
20	800	800	800	1600	0.0520	0.0620	0.0560	0.0596
				4000	0.0492	0.0662	0.0518	0.0680
20	800	1600	2400	1600	0.0568	0.0598	0.0546	0.4326
				4000	0.0536	0.0614	0.0524	0.4338

Table 6.7. Type I error for clean Multivariate t data with cov3I = T

m	n_1	n_2	n_3	B	Median	Mean	Tr.Mn	Class
5	200	200	200	400	0.0396	0.0568	0.0538	0.0554
				1000	0.0464	0.0618	0.0552	0.0582
5	200	400	600	400	0.0450	0.0568	0.0500	0.1228
				1000	0.0496	0.0588	0.0552	0.1326
10	400	400	400	800	0.0506	0.0608	0.0616	0.0582
				2000	0.0472	0.0572	0.0602	0.0612
10	400	800	1200	800	0.0542	0.0592	0.0570	0.2294
				2000	0.0492	0.0560	0.0524	0.2372
20	800	800	800	1600	0.0558	0.0684	0.0662	0.0642
				4000	0.0528	0.0608	0.0622	0.0610
20	800	1600	2400	1600	0.0574	0.0576	0.0654	0.4382
				4000	0.0602	0.0652	0.0636	0.4344

Table 6.8. Type I error for clean lognormal data with cov3I = T

m	n_1	n_2	n_3	B	Median	Mean	Tr.Mn	Class
5	200	200	200	400	0.0424	0.8744	0.0652	0.7208
				1000	0.0446	0.8790	0.0686	0.7220
5	200	400	600	400	0.0470	0.9950	0.0864	0.9980
				1000	0.0460	0.9976	0.0884	0.9990
10	400	400	400	800	0.0440	1.0000	0.2404	1.0000
				2000	0.0438	1.0000	0.2424	1.0000
10	400	800	1200	800	0.0524	1.0000	0.4256	1.0000
				2000	0.0520	1.0000	0.4384	1.0000
20	800	800	800	1600	0.0576	1.0000	0.9674	1.0000
				4000	0.0602	1.0000	0.9668	1.0000
20	800	1600	2400	1600	0.0588	1.0000	0.9994	1.0000
				4000	0.0504	1.0000	0.9996	1.0000

6.1.2 Simulations for power with clean data

Tables 6.9 through 6.12 show the power simulation results. In the table 6.9, the bootstrap test with the mean works similar to the classical test for multivariate normal data and produced the highest power. The test with the median had the worst power in this case. The classical test had the worst power in table 6.10 while the bootstrap test with the trimmed mean had the best power. In table 6.11, the bootstrap test with the trimmed mean was the best. The classical test was the worst except for the balanced design with the different covariance matrices case where the test with the median was the worst. In table 6.12, one of the three bootstrap tests was the best or the same (mainly the cases with power = 1) compared to the classical test.

Table 6.9. Power for MVN data with $\delta_1 = 0.2$ and $\delta_3 = 0.5$

m	n_1	n_2	n_3	B	σ_2	σ_3	Median	Mean	Tr.Mn	Class
5	200	200	200	400	1	1	0.9942	1.0000	1.0000	1.0000
				1000	1	1	0.9932	1.0000	0.9996	1.0000
				400	2	3	0.4198	0.6798	0.5902	0.8002
				1000	2	3	0.3056	0.5260	0.4394	0.6294
5	200	400	800	400	1	1	1.0000	1.0000	1.0000	1.0000
				1000	1	1	1.0000	1.0000	1.0000	1.0000
				400	2	3	0.8020	0.9738	0.9274	0.9026
				1000	2	3	0.8054	0.9648	0.9276	0.9004
10	400	400	400	800	1	1	1.0000	1.0000	1.0000	1.0000
				2000	1	1	1.0000	1.0000	1.0000	1.0000
				800	2	3	0.5860	0.8624	0.7628	0.9282
				2000	2	3	0.5942	0.8614	0.7740	0.9240
10	400	800	1200	800	1	1	1.0000	1.0000	1.0000	1.0000
				2000	1	1	1.0000	1.0000	1.0000	1.0000
				800	2	3	0.9792	0.9998	0.9984	0.9990
				2000	2	3	0.9794	0.9992	0.9970	0.9982
20	800	800	800	1600	1	1	1.0000	1.0000	1.0000	1.0000
				4000	1	1	1.0000	1.0000	1.0000	1.0000
				1600	2	3	0.9242	0.9970	0.9854	0.9990
				4000	2	3	0.9246	0.9968	0.9876	0.9992
20	800	1600	2400	1600	1	1	1.0000	1.0000	1.0000	1.0000
				4000	1	1	1.0000	1.0000	1.0000	1.0000
				1600	2	3	1.0000	1.0000	1.0000	1.0000
				4000	2	3	1.0000	1.0000	1.0000	1.0000

Table 6.10. Power for Mixture data $\delta_1 = 0.2$ and $\delta_3 = 0.5$

m	n_1	n_2	n_3	B	σ_2	σ_3	Median	Mean	Tr.Mn	Class
5	200	200	200	400	1	1	0.7742	0.3168	0.8548	0.2912
				1000	1	1	0.7954	0.3126	0.8654	0.2868
				400	2	3	0.1344	0.0876	0.1530	0.1040
				1000	2	3	0.1392	0.0824	0.1552	0.0992
5	200	400	800	400	1	1	0.9994	0.7694	1.0000	0.7646
				1000	1	1	0.9996	0.7598	1.0000	0.7550
				400	2	3	0.3950	0.1566	0.4742	0.0258
				1000	2	3	0.4038	0.1582	0.4776	0.0244
10	400	400	400	800	1	1	0.9928	0.6144	0.9976	0.5868
				2000	1	1	0.9910	0.6104	0.9968	0.5758
				800	2	3	0.2520	0.1166	0.3020	0.1500
				2000	2	3	0.2474	0.1222	0.2934	0.1568
10	400	800	1200	800	1	1	1.0000	0.9728	1.0000	0.9714
				2000	1	1	1.0000	0.9740	1.0000	0.9720
				800	2	3	0.6626	0.2240	0.7466	0.0574
				2000	2	3	0.6640	0.2218	0.7498	0.0496
20	800	800	800	1600	1	1	1.0000	0.9156	1.0000	0.9004
				4000	1	1	1.0000	0.9214	1.0000	0.9074
				1600	2	3	0.5058	0.1830	0.5920	0.2308
				4000	2	3	0.5030	0.1884	0.5914	0.2396
20	800	1600	2400	1600	1	1	1.0000	0.9998	1.0000	0.9998
				4000	1	1	1.0000	1.0000	1.0000	1.0000
				1600	2	3	0.9594	0.4310	0.9838	0.0830
				4000	2	3	0.9606	0.4208	0.9864	0.0758

Table 6.11. Power for Multivariate t data $\delta_1 = 0.2$ and $\delta_3 = 0.5$

m	n_1	n_2	n_3	B	σ_2	σ_3	Median	Mean	Tr.Mn	Class
5	200	200	200	400	1	1	0.9824	0.9790	0.9982	0.9770
				1000	1	1	0.9832	0.9806	0.9980	0.9778
				400	2	3	0.2564	0.2760	0.3450	0.3448
				1000	2	3	0.2622	0.2750	0.3538	0.3490
5	200	400	800	400	1	1	1.0000	1.0000	1.0000	1.0000
				1000	1	1	1.0000	1.0000	1.0000	1.0000
				400	2	3	0.7328	0.7132	0.8656	0.4634
				1000	2	3	0.7382	0.7210	0.8604	0.4620
10	400	400	400	800	1	1	1.0000	1.0000	1.0000	0.9998
				2000	1	1	1.0000	1.0000	1.0000	1.0000
				800	2	3	0.5048	0.5276	0.6646	0.6324
				2000	2	3	0.5154	0.5266	0.6750	0.6290
10	400	800	1200	800	1	1	1.0000	1.0000	1.0000	1.0000
				2000	1	1	1.0000	1.0000	1.0000	1.0000
				800	2	3	0.9576	0.9494	0.9908	0.8522
				2000	2	3	0.9574	0.9468	0.9892	0.8538
20	800	800	800	1600	1	1	1.0000	1.0000	1.0000	1.0000
				4000	1	1	1.0000	1.0000	1.0000	1.0000
				1600	2	3	0.8670	0.8698	0.9582	0.9288
				4000	2	3	0.8686	0.8660	0.9600	0.9272
20	800	1600	2400	1600	1	1	1.0000	1.0000	1.0000	1.0000
				4000	1	1	1.0000	1.0000	1.0000	1.0000
				1600	2	3	1.0000	0.9996	1.0000	0.9974
				4000	2	3	1.0000	0.9998	1.0000	0.9972

Table 6.12. Power for lognormal data $\delta_1 = 0.2$ and $\delta_3 = 0.5$

m	n_1	n_2	n_3	B	σ_2	σ_3	Median	Mean	Tr.Mn	Class
5	200	200	200	400	1	1	0.9928	0.7360	0.9974	0.6936
				1000	1	1	0.9954	0.7500	0.9986	0.7086
				400	2	3	0.2712	1.0000	0.5182	0.9998
				1000	2	3	0.2810	1.0000	0.5320	1.0000
5	200	400	800	400	1	1	1.0000	0.9934	1.0000	0.9940
				1000	1	1	1.0000	0.9942	1.0000	0.9944
				400	2	3	0.7906	1.0000	0.9922	1.0000
				1000	2	3	0.7956	1.0000	0.9940	1.0000
10	400	400	400	800	1	1	1.0000	0.9774	1.0000	0.9670
				2000	1	1	1.0000	0.9752	1.0000	0.9662
				800	2	3	0.5744	1.0000	0.9402	1.0000
				2000	2	3	0.5632	1.0000	0.9404	1.0000
10	400	800	1200	800	1	1	1.0000	1.0000	1.0000	1.0000
				2000	1	1	1.0000	1.0000	1.0000	1.0000
				800	2	3	0.9780	1.0000	1.0000	1.0000
				2000	2	3	0.9808	1.0000	1.0000	1.0000
20	800	800	800	1600	1	1	1.0000	1.0000	1.0000	1.0000
				4000	1	1	1.0000	0.9998	1.0000	0.9998
				1600	2	3	0.9182	1.0000	0.9998	1.0000
				4000	2	3	0.9212	1.0000	1.0000	1.0000
20	800	1600	2400	1600	1	1	1.0000	1.0000	1.0000	1.0000
				4000	1	1	1.0000	1.0000	1.0000	1.0000
				1600	2	3	1.0000	1.0000	1.0000	1.0000
				4000	2	3	1.0000	1.0000	1.0000	1.0000

6.1.3 Simulations for type I error with contaminated data

For the tables 6.13 through 6.16 the data was contaminated by using the five types of outliers mentioned earlier, and $\gamma = 10\%$ or 5% was used. In table 6.13 the test with the median works reasonably well (close to the nominal coverage) with 10% of outliers with all the distribution and for all the outlier types with the exception of outlier type 3. All the other tests including the classical test failed. Table 6.14 uses 10 parameters with 5% of outliers and results are similar to those in table 6.13. Tables 6.15 and 6.16 uses 20 parameters with 10% and 5% outliers respectively.

Tables 6.17 through 6.20 show the simulations used to get the power curves in section 4.2. It is clear from the graphs that classical test is worst when there is an unbalanced design with different covariance matrices. Also see figures 4.2-4.4. Figure 4.1 uses a balanced design with same covariance matrices. Then the bootstrap test with the mean, the large sample MANOVA type test and the classical test are best and all there of these test almost overlap. In figure 4.2, the bootstrap test with the mean does best while in 4.3 and 4.4 the bootstrap test with the trimmed mean does best.

Table 6.13. Type I error with contaminated data: $m = 5$, $\gamma = 0.1$

Dist.	m	$n_1 = n_2 = n_3$	B	outlier	Median	Mean	Tr.Mn	Class
1	5	200	1000	1	0.0638	0.8034	0.1572	0.9302
				2	0.0504	0.9826	0.1488	1.0000
				3	0.2720	0.9994	0.4024	1.0000
				4	0.0966	0.7862	0.1546	0.9236
				5	0.0840	0.9854	0.1268	1.0000
2	5	200	1000	1	0.0488	0.1832	0.1068	0.1812
				2	0.0376	0.4880	0.1042	0.5428
				3	0.1994	0.7502	0.2206	0.8978
				4	0.0858	0.1848	0.1080	0.1830
				5	0.0780	0.4688	0.0974	0.5400
3	5	200	1000	1	0.0598	0.6046	0.1554	0.7094
				2	0.0464	0.9356	0.1366	0.9946
				3	0.2590	0.9944	0.3882	1.0000
				4	0.0946	0.5824	0.1486	0.6926
				5	0.0828	0.9216	0.1270	0.9928
4	5	200	1000	1	0.0426	0.9880	0.1998	0.9624
				2	0.0416	0.9924	0.1396	0.9884
				3	0.1762	1.0000	0.3980	1.0000
				4	0.0708	0.9902	0.1892	0.9674
				5	0.0766	0.9950	0.1508	0.9932

Table 6.14. Type I error with contaminated data: $m = 10$, $\gamma = 0.05$

Dist.	m	$n_1 = n_2 = n_3$	B	outlier	Median	Mean	Tr.Mn	Class
1	10	800	1000	1	0.0274	0.9720	0.1318	0.9956
				2	0.0192	1.0000	0.1132	1.0000
				3	0.4710	1.0000	0.6760	1.0000
				4	0.0872	0.9660	0.1264	0.9974
				5	0.0756	1.0000	0.1084	1.0000
2	10	800	1000	1	0.0272	0.2586	0.1246	0.2734
				2	0.0152	0.9182	0.0992	0.9764
				3	0.4288	0.9976	0.5954	1.0000
				4	0.0852	0.2454	0.1102	0.2658
				5	0.0644	0.9032	0.0898	0.9750
3	10	800	1000	1	0.0252	0.8384	0.1266	0.9178
				2	0.0158	0.9996	0.1160	1.0000
				3	0.4596	1.0000	0.6544	1.0000
				4	0.0954	0.8254	0.1270	0.9158
				5	0.0804	1.0000	0.1054	1.0000
4	10	800	1000	1	0.0232	1.0000	0.6856	1.0000
				2	0.0146	1.0000	0.5036	1.0000
				3	0.4130	1.0000	0.9962	1.0000
				4	0.0806	1.0000	0.7362	1.0000
				5	0.0714	1.0000	0.6184	1.0000

Table 6.15. Type I error with contaminated data: $m = 20, \gamma = 0.1$

Dist.	m	$n_1 = n_2 = n_3$	B	outlier	Median	Mean	Tr.Mn	Class
1	20	800	2000	1	0.1522	0.8752	0.3994	0.8906
				2	0.0960	1.0000	0.2958	1.0000
				3	0.9986	1.0000	1.0000	1.0000
				4	0.2170	0.8612	0.3684	0.8800
				5	0.1740	1.0000	0.2890	1.0000
2	20	800	2000	1	0.1256	0.1496	0.3378	0.1342
				2	0.0988	0.9434	0.2560	0.9654
				3	0.9762	1.0000	0.9924	1.0000
				4	0.1994	0.1426	0.3028	0.1210
				5	0.1578	0.9438	0.2392	0.9664
3	20	800	2000	1	0.1438	0.5946	0.3928	0.5666
				2	0.0904	1.0000	0.2900	1.0000
				3	0.9970	1.0000	0.9998	1.0000
				4	0.2184	0.5888	0.3598	0.5716
				5	0.1546	1.0000	0.2804	1.0000
4	20	800	2000	1	0.2928	1.0000	0.8634	1.0000
				2	0.0892	1.0000	0.9490	1.0000
				3	0.9912	1.0000	1.0000	1.0000
				4	0.1756	1.0000	0.9986	1.0000
				5	0.1512	1.0000	0.9920	1.0000

Table 6.16. Type I error with contaminated data: $m = 20$, $\gamma = 0.05$

Dist.	m	$n_1 = n_2 = n_3$	B	outlier	Median	Mean	Tr.Mn	Class
1	20	800	1000	1	0.0046	0.7012	0.1330	0.7932
				2	0.0010	0.9974	0.1148	1.0000
				3	0.5950	0.9998	0.7496	1.0000
				4	0.0886	0.6820	0.1218	0.7758
				5	0.0878	0.9972	0.1138	1.0000
2	20	800	1000	1	0.0044	0.1468	0.1064	0.1338
				2	0.0014	0.7946	0.0960	0.9004
				3	0.5116	0.9776	0.6054	1.0000
				4	0.0800	0.1434	0.1038	0.1294
				5	0.0696	0.7718	0.0850	0.8908
3	20	800	1000	1	0.0032	0.4786	0.1346	0.5006
				2	0.0020	0.9908	0.1076	1.0000
				3	0.5762	0.9992	0.7190	1.0000
				4	0.0816	0.4524	0.1136	0.4848
				5	0.0784	0.9888	0.1020	1.0000
4	20	800	1000	1	0.0016	1.0000	0.9664	1.0000
				2	0.0014	1.0000	0.9216	1.0000
				3	0.4778	1.0000	1.0000	1.0000
				4	0.0822	1.0000	0.9892	1.0000
				5	0.0720	1.0000	0.9768	1.0000

Table 6.17. Power curve for MVN data with $m = 5, \sigma_1 = 1, \sigma_2 = 1, \sigma_3 = 1, n_1 = 200, n_2 = 200$ and $n_3 = 200$

δ_1	δ_2	δ_3	Med	Mean	Tr.Me	Clas	MaTyp
0	0	0	0.0442	0.0588	0.0562	0.0502	0.0470
0.04	0.08	0.12	0.0700	0.1206	0.1048	0.1070	0.1072
0.08	0.16	0.24	0.1776	0.326	0.2714	0.2992	0.2982
0.12	0.24	0.36	0.4308	0.696	0.6002	0.6758	0.6694
0.16	0.32	0.48	0.7266	0.9442	0.8836	0.9384	0.9364
0.20	0.40	0.60	0.9142	0.9952	0.9834	0.9958	0.9948
0.24	0.48	0.72	0.9876	0.9998	0.9994	0.9998	0.9998
0.30	0.60	0.90	0.9998	1	1	1	1
0.35	0.70	1.05	1	1	1	1	1

Table 6.18. Power curve for MVN data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 = 400$ and $n_3 = 600$

δ_1	δ_2	δ_3	Med	Mean	Tr.Me	Clas	MaTyp
0	0	0	0.0492	0.0510	0.0542	0.0034	0.0446
0.04	0.08	0.12	0.0636	0.0786	0.0726	0.0062	0.0654
0.08	0.16	0.24	0.0804	0.1224	0.1086	0.0060	0.1044
0.12	0.24	0.36	0.1400	0.2322	0.1942	0.0146	0.2158
0.16	0.32	0.48	0.2388	0.4204	0.3494	0.0334	0.3956
0.20	0.4	0.6	0.3836	0.6438	0.5372	0.0710	0.6186
0.24	0.48	0.72	0.5538	0.8280	0.7278	0.1240	0.8126
0.30	0.60	0.90	0.7922	0.9648	0.9200	0.2868	0.9590
0.35	0.70	1.05	0.9170	0.9964	0.9832	0.4898	0.9958

Table 6.19. Power curve for Mixture data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 = 400$ and $n_3 = 600$

δ_1	δ_2	δ_3	Med	Mean	Tr.Me	Clas	MaTyp
0	0	0	0.0486	0.0522	0.0524	0.0474	0.0054
0.04	0.08	0.12	0.0504	0.0602	0.0520	0.0510	0.0034
0.08	0.16	0.24	0.0588	0.0624	0.0670	0.0540	0.0042
0.12	0.24	0.36	0.1196	0.0688	0.1356	0.059	0.0054
0.16	0.32	0.48	0.1178	0.0820	0.1350	0.0708	0.0068
0.20	0.40	0.60	0.3846	0.6354	0.5432	0.6148	0.0644
0.24	0.48	0.72	0.2350	0.1120	0.2948	0.0964	0.0078
0.30	0.60	0.90	0.3826	0.1570	0.4650	0.1378	0.0086
0.40	0.80	1.20	0.689	0.2602	0.7726	0.2326	0.0190
0.45	0.90	1.35	0.8032	0.3134	0.8646	0.2916	0.0250
0.50	1.00	1.50	0.9006	0.3698	0.9396	0.3518	0.0312
0.60	1.20	1.80	0.9780	0.5592	0.9868	0.5370	0.0538
0.80	1.60	2.40	0.9998	0.8450	0.9998	0.8296	0.1358

Table 6.20. Power curve for Multivariate t_4 data with $m = 5, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 5, n_1 = 200, n_2 = 400$ and $n_3 = 600$

δ_1	δ_2	δ_3	Med	Mean	Tr.Me	Clas	MaTyp
0	0	0	0.0500	0.0530	0.0572	0.0046	0.0450
0.04	0.08	0.12	0.0516	0.0650	0.0652	0.0042	0.0560
0.08	0.16	0.24	0.0750	0.0814	0.0940	0.0056	0.0714
0.12	0.24	0.36	0.1304	0.1460	0.1770	0.0098	0.1308
0.16	0.32	0.48	0.2102	0.2214	0.2888	0.013	0.2038
0.20	0.40	0.60	0.3302	0.338	0.4618	0.0198	0.3168
0.24	0.48	0.72	0.487	0.4758	0.6278	0.0380	0.4572
0.30	0.60	0.90	0.7344	0.7122	0.8526	0.0858	0.7004
0.35	0.70	1.05	0.873	0.8572	0.9520	0.1462	0.8440
0.40	0.80	1.20	0.9576	0.9472	0.9900	0.2586	0.9416
0.45	0.90	1.35	0.9892	0.9848	0.9994	0.3768	0.9826

CHAPTER 7

DISCUSSION

The first part of this dissertation suggests a practical method to perform a multivariate two sample test when the asymptotic covariance matrix of the statistic T_i is difficult to estimate. Such tests may be useful when the data distribution is unknown or outliers are present. The method was illustrated with the coordinatewise median, sample mean, 25% trimmed mean, and RMVN estimators. All four estimators work well when the prediction region method was applied to the clean data, although care needs to be taken with the multivariate lognormal distribution where the four estimators T_i are estimating different parameters $\boldsymbol{\mu}_{T_i}$.

Both the sample mean and the 25% trimmed mean failed to achieve the nominal coverage when H_0 is true with the contaminated data. The coordinatewise median could handle up to 10% outliers, while the RMVN estimator could handle up to 40% outliers. Both estimators were robust to the equal covariance assumption.

This work also gives practical alternatives for the classical one way MANOVA test. One novel large sample test was given. This test does not use the unreasonably strong equal covariance assumption and worked well in the simulations. Another practical method to bootstrap analogs of the one way MANOVA test was also suggested. Such tests may be useful when the data distribution is unknown or outliers are present. The method was illustrated with the coordinatewise median, sample mean, 25% trimmed mean. All three estimators work well when the prediction region method was applied to the clean data. The classical test does not work well when the assumptions are not met. The classical test is worst in terms of type I error and power when working with an unbalanced design and the equal covariance assumption is violated.

The bootstrap test with the coordinatewise median is recommended when the outliers are present. To make the bootstrap test more efficient with outliers, a robust location estimator can be used (e.g. RMVN location estimator) as in the two sample case.

Konietschke, Bathke, Harrar, and Pauly (2015) suggest a method for bootstrapping

the MANOVA model, and Willems, Pison, Rousseeuw, and Van Aelst (2002) suggest a robust one sample Hotelling's T^2 type test. References for robust one way MANOVA tests are in Finch and French (2013), Todorov and Filzmoser (2010), Van Aelst and Willems (2011), and Wilcox (1995).

The *R* software was used in the simulation. See R Core Team (2016). Programs are in the Olive (2017b) collection of *R* functions *mpack.txt* available from (<http://lagrange.math.siu.edu/Olive/mpack.txt>). The function `hot2sim` was used to simulate the two sample tests of hypotheses, The Curran (2013) *R* package `Hotelling` was used to perform the classical 2 sample Hotelling's T^2 test. The function `manbtsim2` was used to simulate the bootstrapped one way MANOVA tests, `predreg` computes the confidence region given the bootstrap values from `rhot2boot`, and `manbtsim4` simulates the large sample one way MANOVA type test and to produce the power curves.

APPENDICES

R CODES FOR THE SIMULATIONS

APPENDIX I

The following is the R code for the simulations for Bootstrapping Analogs of the Two Sample Hotelling's T^2 Test.

```
source("http://lagrange.math.siu.edu/Olive/mpack.txt")
library(MASS)

hot2sim<-function(n = 100, n2 = 100, p = 2, csteps = 5, B= 100, gam = 0.4,
nruns = 100, xtype = 1, outliers = 0, z = 10, sig = 1, eps = 0.4,
dd=4, delta = 0){
# This R function simulates the two sample Hotelling's  $T^2$  test based on the
#bootstrap. Here x2 is clean, x = x1 may have outliers.
# Uses the coordinatewise mean, median, 25% trimmed mean, and RMVN location
estimators.
# Need  $p > 1$ . Want  $n > 20p$ ,  $n2 > 20p$ ,  $B > 20p$ .
# Multiply x by A where xtype = 1 for MVN  $N_p(0,I)$ ,
# 2, (with delta = eps) for  $(1 - \text{delta}) N_p(0,I) + \text{delta} N_p(0, 25 I)$ 
# 3 for multivariate  $t_d$  with  $d = \text{dd}$ 
# 4 for lognormal.
# outliers = 0 for no outliers and  $X \sim N(0, \text{diag}(1, \dots, p))$ ,
# 1 for outliers a tight cluster at major axis  $(0, \dots, 0, z)'$ 
# 2 for outliers a tight cluster at minor axis  $(z, 0, \dots, 0)'$ 
# 3 for outliers  $X \sim N((z, \dots, z)', \text{diag}(1, \dots, p))$ 
# 4 for outliers  $X[i, p] = z$ 
# 5 for outliers  $X[i, 1] = z$ 
# For the point mass outlier types 4 and 5, need gam much smaller than 0.4.
# Power can be estimated by increasing delta so  $\mu = \text{delta}(1, \dots, 1)$ 
# and  $\mu_2 = (0, \dots, 0)$ .
```

```

# Cov(x) = diag(1,2,...,p), Cov(x2) = sig^2 Cov(x) for clean data.
# For outliers=0, want hquant and rquant approx 1.
  A <- sqrt(diag(1:p))
  munot <- 0 * (1:p)
  mu <- delta * (1 + munot)
  val <- floor(gam * n)
    indx <- 1:n
    indx2 <- 1:n2
    medmus <- matrix(0,nrow=B,ncol=p)
    mnmus <- medmus
    tmnmus <- medmus
    rmvnmus <- medmus
    medcv <- 0
    mncv <- 0
    tmncv <- 0
    rmvncv <- 0
    chisqcut <- qchisq(p=0.95,df=p)
    cutoffs <- matrix(0,nrow=nruns,ncol=4)
    for(i in 1:nruns) {
#make data
      x <- matrix(rnorm(n * p), ncol = p, nrow = n)
      x2 <- matrix(rnorm(n2 * p), ncol = p, nrow = n2)
      if(xtype == 2) {
        zu <- runif(n)
        x[zu < eps, ] <- x[zu < eps, ] * 5
        zu <- runif(n2)
        x2[zu < eps, ] <- x2[zu < eps, ] * 5
      }
      if(xtype == 3) {
        zu <- sqrt(rchisq(n, dd)/dd)

```

```

x <- x/zu
      zu <- sqrt(rchisq(n2, dd)/dd)
x2 <- x2/zu
}
if(xtype == 4){ #Want pop coord med(x) = 0.
  x <- exp(x)
      x <- x - 1
      x2 <- exp(x2)
      x2 <- x2 - 1
  }
x <- x %**% A
  x2 <- x2 %**% A
  x2 <- sig * x2
if(outliers == 1) {
  x[1:val, ] <- matrix(rnorm(val * p, sd = 0.01), ncol
    = p, nrow = val)
  x[1:val, p] <- x[1:val, p] + z
}
if(outliers == 2) {
  x[1:val, ] <- matrix(rnorm(val * p, sd = 0.01), ncol
    = p, nrow = val)
  x[1:val, 1] <- x[1:val, 1] + z
}
if(outliers == 3) {
  tem <- z + 0 * 1:p
  x[1:val, ] <- x[1:val, ] + tem
}
if(outliers == 4) {
  x[1:val, p] <- z
}

```

```

if(outliers == 5) {
  x[1:val, 1] <- z
}
x <- mu + x
#clean x has mean  $\mu = \delta(1, \dots, 1)^T$ , x2 has mean  $(0, \dots, 0)^T$ 
#get bootstrapped Tx - Tx2 for various statistics T
for(j in 1:B){
  tem <- sample(indx,n,replace=T)
  medx <- apply(x[tem,],2,median)
  mnx <- apply(x[tem,],2,mean)
  tmnx <- apply(x[tem,],2,tmn)
  rmvnx <- covrmvn(x[tem,])$center
  tem <- sample(indx2,n2,replace=T)
  medx2 <- apply(x2[tem,],2,median)
  mnx2 <- apply(x2[tem,],2,mean)
  tmnx2 <- apply(x2[tem,],2,tmn)
  rmvnx2 <- covrmvn(x2[tem,])$center
  medmus[j,] <- medx-medx2
  mnmus[j,] <- mnx-mnx2
  tmnmus[j,] <- tmnx-tmnx2
  rmvnmus[j,] <- rmvnx-rmvnx2
}

outmed<-predreg(medmus)
medcv <- medcv + outmed$inr
outmn <- predreg(mnmus)
mncv <- mncv + outmn$inr
outtmn <- predreg(tmnmus)
tmncv <- tmncv + outtmn$inr
outrmvn <- predreg(rmvnmus)
rmvncv <- rmvncv + outrmvn$inr

```

```

        cutoffs[i,]<-
            c(outmed$cuplim,outmn$cuplim,outtmn$cuplim,outrmvn$cuplim)^2
    }
medcv <- 1 - medcv/nruns #prop of times Ho is rejected
    mncv <- 1 - mncv/nruns
    tmncv <- 1 - tmncv/nruns
    rmvncv <- 1 - rmvncv/nruns
    mncut <- apply(cutoffs,2,mean)
list(chisqcut = chisqcut, mncut=mncut, medcv = medcv,
      mncv = mncv,tmncv=tmncv,rmvncv=rmvncv)
}

tmn<-function(x, tp = 0.25)
{# computes 100tp% trimmed mean
    mean(x, trim = tp)
}

```

The following is the *R* code for the simulations for Bootstrapping Analogs of the One-Way MANOVA Test.

```

library(MASS)
source("http://lagrange.math.siu.edu/Olive/mpack.txt")

manbtsim2<-function(n = 100, n2 = 100, n3 = 100, m = 2, csteps = 5, B= 100,
gam = 0.4, nruns = 100, ytype = 1, outliers = 0, z = 10, sig2 = 1, sig3 = 1,
eps = 0.4, dd=4, delta = 0, delta3 = 0, cov3I = F){
# This R function simulates a one way Manova type bootstrap test for
# Ho:  $\mu_1 = \mu_2 = \mu_3$  where  $p = g = 3 =$  number of groups.
#Here y2 and y3 are clean,  $y = y_1$  may have outliers, and  $y$  is  $m$  by 1.
# Uses the coordinatewise mean, median, and 25% trimmed mean.
## Does not use the slow RMVN location estimators.
# Need  $m > 1$ . Want  $n > 20m$ ,  $n_2 > 20m$ ,  $n_3 > 20m$ ,  $B > 20m$ .
# Multiply  $y$  by  $A$  where  $ytype = 1$  for  $MVN Nm(0,I)$ ,
# 2 for  $(1 - \text{eps}) Nm(0,I) + \text{eps} Nm(0, 25 I)$ ,
# 3 for multivariate  $t_d$  with  $d = dd$ ,
# 4 for lognormal.
# outliers = 0 for no outliers and  $Y \sim N(0, \text{diag}(1, \dots, m))$ ,
# 1 for outliers a tight cluster at major axis  $(0, \dots, 0, z)'$ 
# 2 for outliers a tight cluster at minor axis  $(z, 0, \dots, 0)'$ 
# 3 for outliers  $Y \sim N((z, \dots, z)', \text{diag}(1, \dots, m))$ 
# 4 for outliers  $Y[i, m] = z$ 
# 5 for outliers  $Y[i, 1] = z$ 
# For the point mass outlier types 4 and 5, need  $\text{gam}$  much smaller than 0.4.
# Power can be estimated by increasing  $\text{delta}$  so  $\mu_1 = \text{delta}(1, \dots, 1)$ 
# and  $\mu_2 = (0, \dots, 0)$ ,  $\mu_3 = \text{delta}_3(1, \dots, 1)$ .
#  $\text{Cov}(y) = \text{diag}(1, 2, \dots, m)$ ,  $\text{Cov}(y_2) = \text{sig}^2 \text{Cov}(y)$  for clean data.
#  $\text{Cov}(y_3) = \text{sig}^3 \text{Cov}(y)$  for clean data if  $\text{cov3I} = F$ ,

```

```

# or Cov(y3) = cI_3 if cov3I = T.
# For outliers=0, want hquant and rquant approx 1.
A <- sqrt(diag(1:m))
      munot <- 0 * (1:m)
mu <- delta * (1 + munot)
      mu3 <- delta3 * (1 + munot)
val <- floor(gam * n)
      indx <- 1:n
      indx2 <- 1:n2
      indx3 <- 1:n3
      medmus <- matrix(0,nrow=B,ncol=2*m)
      mnmus <- medmus
      tmnmus <- medmus
      medcv <- 0
      mncv <- 0
      tmncv <- 0

crej <- 0
gp <- 0
grp <- 0
pval <- 0
out <- 0

      chisqcut <- qchisq(p=0.95,df=2*m)
      cutoffs <- matrix(0,nrow=nruns,ncol=3)
      for(i in 1:nruns) {
#make data
y <- matrix(rnorm(n * m), ncol = m, nrow = n)
      y2 <- matrix(rnorm(n2 * m), ncol = m, nrow = n2)
      y3 <- matrix(rnorm(n3 * m), ncol = m, nrow =n3)
if(ytype == 2) {

```

```

zu <- runif(n)
y[zu < eps, ] <- y[zu < eps, ] * 5
      zu <- runif(n2)
y2[zu < eps, ] <- y2[zu < eps, ] * 5
      zu <- runif(n3)
y3[zu < eps, ] <- y3[zu < eps, ] * 5
}
if(ytype == 3) {
zu <- sqrt(rchisq(n, dd)/dd)
y <- y/zu
      zu <- sqrt(rchisq(n2, dd)/dd)
y2 <- y2/zu
      zu <- sqrt(rchisq(n3, dd)/dd)
y3 <- y3/zu
}
if(ytype == 4){ #Want pop coord med(y) = 0.
y <- exp(y)
      y <- y - 1
      y2 <- exp(y2)
      y2 <- y2 - 1
      y3 <- exp(y3)
      y3 <- y3 - 1
}
y <- y %**% A
      y2 <- y2 %**% A
      y2 <- sig2 * y2
      if( cov3I != T){
      y3 <- y3 %**% A
      y3 <- sig3 * y3}
if(outliers == 1) {

```



```

        y[1:val, ] <- matrix(rnorm(val * m, sd = 0.01), ncol
= m, nrow = val)
y[1:val, m] <- y[1:val, m] + z
}
if(outliers == 2) {
y[1:val, ] <- matrix(rnorm(val * m, sd = 0.01), ncol
= m, nrow = val)
y[1:val, 1] <- y[1:val, 1] + z
}
if(outliers == 3) {
tem <- z + 0 * 1:m
y[1:val, ] <- y[1:val, ] + tem
}
if(outliers == 4) {
y[1:val, m] <- z
}
if(outliers == 5) {
y[1:val, 1] <- z
}
y <- mu + y

        y3 <- mu3 + y3
#clean y has mean  $\mu = \delta(1, \dots, 1)^T$ ,
#y2 has mean  $(0, \dots, 0)^T$  and y3 has mean  $\delta_3(1, \dots, 1)^T$ 
#get bootstrapped  $T_y - T_{y3}$ ,  $T_{y2} - T_{y3}$  for various statistics T
for(j in 1:B){
tem <- sample(indx,n,replace=T)
medy <- apply(y[tem,],2,median)
mny <- apply(y[tem,],2,mean)
tmny <- apply(y[tem,],2,tmn)
tem <- sample(indx2,n2,replace=T)

```

```

    medy2 <- apply(y2[tem,],2,median)
    mny2 <- apply(y2[tem,],2,mean)
    tmny2 <- apply(y2[tem,],2,tmn)
    tem <- sample(indx3,n3,replace=T)
    medy3 <- apply(y3[tem,],2,median)
    mny3 <- apply(y3[tem,],2,mean)
    tmny3 <- apply(y3[tem,],2,tmn)
    medmus[j,] <- c(medy-medy3,medy2-medy3)
    mnmus[j,] <- c(mny-mny3,mny2-mny3)
    tmnmus[j,] <- c(tmny-tmny3,tmny2-tmny3)
  }

  outmed<-predreg(medmus)
  medcv <- medcv + outmed$inr
  outmn <- predreg(mnmus)
  mncv <- mncv + outmn$inr
  outtmn <- predreg(tmnmus)
  tmncv <- tmncv + outtmn$inr
  cutoffs[i,]<-c(outmed$cuplim,outmn$cuplim,outtmn$cuplim)^2

#Get the classical test coverage
yall<-rbind(y,y2,y3)
yall<-as.matrix(yall)
  gp <- c(rep(1, n),rep(2, n2),rep(3, n3))
grp<-factor(gp)
out<-manova(yall~grp)

pval <- summary(out,test="Hotelling-Lawley")$stats[1,6]
  #pvalue for Hotelling-Lawley's test

if(pval < 0.05){

```

```
crej <- crej +1
}

}
medcv <- 1 - medcv/nruns #prop of times Ho is rejected
  mncv <- 1 - mncv/nruns
  tmncv <- 1 - tmncv/nruns
  mncut <- apply(cutoffs,2,mean)
ccv <- crej/nruns

list(chisqcut = chisqcut, mncut=mncut, medcv = medcv,
     mncv = mncv,tmncv=tmncv,ccv=ccv)
}
```

The following is the *R* code for the power curves for Bootstrapping Analogs of the One-Way MANOVA Test and the large sample test from the section 4.1.

```
manbtsim4 <- function(n1 = 100, n2 = 100, n3 = 100, m = 2, csteps = 5, B= 100,
gam = 0.4, nruns = 100, ytype = 1, outliers = 0, z = 10, sig2 = 1, sig3 = 1,
eps = 0.4, dd=4, delta1 = 0, delta2 = 0, delta3 = 0, cov3I = F){
#needs library(MASS)
# This R function simulates one way Manova type tests for
# Ho:  $\mu_1 = \mu_2 = \mu_3$  where  $p = g = 3 =$  number of groups.
#Can vary the mean vectors  $\mu_i$  better than in manbtsim2.
#Has one more test than manbtsim3.
#Here  $y_2$  and  $y_3$  are clean,  $y_1$  may have outliers, and  $y_1$  is  $m$  by 1.
# Uses the coordinatewise mean, median, and 25% trimmed mean.
#Also does the classical Hotelling Lawley one way MANOVA test.
#Partially written by Hasthika S. Rupasinghe Arachchige Don.
## Does not use the slow RMVN location estimators.
# Need  $m > 1$ . Want  $n > 20m$ ,  $n_2 > 20m$ ,  $n_3 > 20m$ ,  $B > 20m$ .
# Multiply  $y$  by  $A$  where  $ytype = 1$  for  $MVN N_m(0,I)$ ,
# 2 for  $(1 - \text{eps}) N_m(0,I) + \text{eps} N_m(0, 25 I)$ ,
# 3 for multivariate  $t_d$  with  $d = dd$ ,
# 4 for lognormal.
# outliers = 0 for no outliers and  $Y \sim N(0, \text{diag}(1, \dots, m))$ ,
# 1 for outliers a tight cluster at major axis  $(0, \dots, 0, z)'$ 
# 2 for outliers a tight cluster at minor axis  $(z, 0, \dots, 0)'$ 
# 3 for outliers  $Y \sim N((z, \dots, z)', \text{diag}(1, \dots, m))$ 
# 4 for outliers  $Y[i, m] = z$ 
# 5 for outliers  $Y[i, 1] = z$ 
# For the point mass outlier types 4 and 5, need  $\text{gam}$  much smaller than 0.4.
# Power can be estimated by using unequal  $\delta_i$  so  $\mu_1 = \delta_1(1, \dots, 1)$ 
```

```

# and mu2 = delta2(1, ..., 1), mu3 = delta3(1,...,1).
# Cov(y1) = diag(1,2,...,m), Cov(y2) = sig^2 Cov(y1) for clean data.
# Cov(y3) = sig^3 Cov(y1) for clean data if cov3I = F,
# or Cov(y3) = cI_3 if cov3I = T.
# For outliers=0, want hquant and rquant approx 1.
  A <- sqrt(diag(1:m))
  munot <- 0 * (1:m)
  mu1 <- delta1 * (1 + munot)
  mu2 <- delta2 * (1 + munot)
  mu3 <- delta3 * (1 + munot)
  val <- floor(gam * n1)
  indx1 <- 1:n1
  indx2 <- 1:n2
  indx3 <- 1:n3
  medmus <- matrix(0,nrow=B,ncol=2*m)
  mnmus <- medmus
  tmnmus <- medmus
  medcv <- 0
  mncv <- 0
  tmncv <- 0
  mantcov <- 0
  crej <- 0
  gp <- 0
  grp <- 0
  pval <- 0
  out <- 0

  chisqcut <- qchisq(p=0.95,df=2*m)
  cutoffs <- matrix(0,nrow=nruns,ncol=3)
  for(i in 1:nruns) {

```

```
#make data
```

```
y1 <- matrix(rnorm(n1 * m), ncol = m, nrow = n1)
y2 <- matrix(rnorm(n2 * m), ncol = m, nrow = n2)
y3 <- matrix(rnorm(n3 * m), ncol = m, nrow = n3)
if(ytype == 2) {
  zu <- runif(n)
  y1[zu < eps, ] <- y1[zu < eps, ] * 5
  zu <- runif(n2)
  y2[zu < eps, ] <- y2[zu < eps, ] * 5
  zu <- runif(n3)
  y3[zu < eps, ] <- y3[zu < eps, ] * 5
}
if(ytype == 3) {
  zu <- sqrt(rchisq(n1, dd)/dd)
  y1 <- y1/zu
  zu <- sqrt(rchisq(n2, dd)/dd)
  y2 <- y2/zu
  zu <- sqrt(rchisq(n3, dd)/dd)
  y3 <- y3/zu
}
if(ytype == 4){ #Want pop coord med(y) = 0.
  y1 <- exp(y1)
  y1 <- y1 - 1
  y2 <- exp(y2)
  y2 <- y2 - 1
  y3 <- exp(y3)
  y3 <- y3 - 1
}
y1 <- y1 %*% A
y2 <- y2 %*% A
```

```

y2 <- sig2 * y2
if( cov3I != T){
y3 <- y3 %**% A
y3 <- sig3 * y3}
if(outliers == 1) {
  y1[1:val, ] <- matrix(rnorm(val * m, sd = 0.01), ncol
                        = m, nrow = val)
  y1[1:val, m] <- y1[1:val, m] + z
}
if(outliers == 2) {
  y1[1:val, ] <- matrix(rnorm(val * m, sd = 0.01), ncol
                        = m, nrow = val)
  y1[1:val, 1] <- y1[1:val, 1] + z
}
if(outliers == 3) {
  tem <- z + 0 * 1:m
  y1[1:val, ] <- y1[1:val, ] + tem
}
if(outliers == 4) {
  y1[1:val, m] <- z
}
if(outliers == 5) {
  y1[1:val, 1] <- z
}
y1 <- mu1 + y1
y2 <- mu2 + y2
y3 <- mu3 + y3
#clean y1 has mean mu1 =delta1(1,...,1)^T,
#y2 has mean delta2(1,...,1)^T and y3 has mean delta3 (1,...,1)^T
#get bootstrapped Ty - Ty3, Ty2 - Ty3 for various statistics T

```

```

for(j in 1:B){
  tem <- sample(indx1,n1,replace=T)
  medy <- apply(y1[tem,],2,median)
  mny <- apply(y1[tem,],2,mean)
  tmny <- apply(y1[tem,],2,tmn)
  tem <- sample(indx2,n2,replace=T)
  medy2 <- apply(y2[tem,],2,median)
  mny2 <- apply(y2[tem,],2,mean)
  tmny2 <- apply(y2[tem,],2,tmn)
  tem <- sample(indx3,n3,replace=T)
  medy3 <- apply(y3[tem,],2,median)
  mny3 <- apply(y3[tem,],2,mean)
  tmny3 <- apply(y3[tem,],2,tmn)
  medmus[j,] <- c(medy-medy3,medy2-medy3)
  mnmus[j,] <- c(mny-mny3,mny2-mny3)
  tmnmus[j,] <- c(tmny-tmny3,tmny2-tmny3)
}

outmed<-predreg(medmus)
medcv <- medcv + outmed$inr
outmn <- predreg(mnmus)
mncv <- mncv + outmn$inr
outtmn <- predreg(tmnmus)
tmncv <- tmncv + outtmn$inr
cutoffs[i,]<-c(outmed$cuplim,outmn$cuplim,outtmn$cuplim)^2

#Get the classical test coverage
yall<-rbind(y1,y2,y3)
yall<-as.matrix(yall)
gp <- c(rep(1, n1),rep(2, n2),rep(3, n3))
grp<-factor(gp)

```



```

out<-manova(yall~grp)

pval <- summary(out,test="Hotelling-Lawley")$stats[1,6] #pvalue for
#Hotelling-Lawley's test

    if(pval < 0.05){
      crej <- crej +1
    }

#large sample Manova type test based on the sample mean
pval <- mantyp(y=yall,p=3,group=grp)$pval
if(pval < 0.05) mantcov <- mantcov + 1
}

medcv <- 1 - medcv/nruns #prop of times Ho is rejected
mncv <- 1 - mncv/nruns
tmncv <- 1 - tmncv/nruns
mncut <- apply(cutoffs,2,mean)
ccv <- crej/nruns
mantcov <- mantcov/nruns

list(chisqcut = chisqcut, mncut=mncut, medcv = medcv,
      mncv = mncv,tmncv=tmncv,ccv=ccv,mantcov=mantcov)
}

```

REFERENCES

- [1] Bickel, P. J., Ren, J. –J. (2001). The bootstrap in hypothesis testing. In: de Gunst, M., Klaassen, C., van der Vaart, A. Eds. *State of the Art in Probability and Statistics: Festschrift for William R. van Zwet*. The Institute of Mathematical Statistics, CA: Hayward, 91-112.
- [2] Clarke, B. R. (1986). Nonsmooth analysis and Fréchet differentiability of M functionals. *Probability Theory and Related Fields*, 73:137-209.
- [3] Clarke, B. R. (2000). A review of differentiability in relation to robustness with an application to seismic data analysis. *Proceedings of the Indian National Science Academy*, A 66:467-482.
- [4] Curran, J. M. (2013). Hotelling: Hotelling’s T-squared test and variants. *R Package version 1.0-2*, (<https://cran.r-project.org/package=Hotelling>).
- [5] Fernholtz, L. T. (1983). *von Mises Calculus for Statistical Functionals*. New York: Springer.
- [6] Finch, H., French, B. (2013). A “Monte Carlo” comparison of robust MANOVA test statistics. *Journal of Modern Applied Statistical Methods*, 12:35-81.
- [7] Fujikoshi, Y. (2002), Asymptotic Expansions for the Distributions of Multivariate Basic Statistics and One-Way MANOVA Tests Under Nonnormality, *Journal of Statistical Planning and Inference*, 108, 263-282.
- [8] Gill, R. D. (1989). Non- and semi-parametric maximum likelihood estimators and the von Mises method, part 1. *Scandinavian Journal of Statistics*, 16:97-128.
- [9] Johnson, R. W. (1996). Fitting percentage of body fat to simple body measurements. *Journal of Statistics Education* 4, online at (www.amstat.org/publications/jse/).
- [10] Kakizawa, Y. (2009), Third-Order Power Comparisons for a Class of Tests for Multivariate Linear Hypothesis Under General Distributions, *Journal of Multivariate Analysis*, 100, 473-496.
- [11] Konietzschke, F., Bathke, A. C., Harrar, S. W., Pauly, M. (2015). Parametric and non-parametric bootstrap methods for general MANOVA. *Journal of Multivariate Analysis*,

140:291-301.

- [12] Olive, D. J. (2004). A resistant estimator of multivariate location and dispersion. *Computational Statistics & Data Analysis*, 46:99-102.
- [13] Olive, D. J. (2013). Asymptotically optimal regression prediction intervals and prediction regions for multivariate data. *International Journal of Statistics and Probability*, 2:90-100.
- [14] Olive, D. J. (2017a). *Applications of hyperellipsoidal prediction regions*. Statistical Papers. To appear.
- [15] Olive, D. J. (2017b). *Robust Multivariate Analysis*. New York: Springer. To appear.
- [16] Olive, D. J. (2017c). Bootstrapping hypothesis tests and confidence regions. Unpublished Manuscript with the bootstrap material from Olive (2017b) at (<http://lagrange.math.siu.edu/Olive/ppvselboot.pdf>).
- [17] Olive, D. J., Hawkins, D. M. (2010). Robust multivariate location and dispersion. Unpublished Manuscript available from (<http://lagrange.math.siu.edu/Olive/pphbml.pdf>).
- [18] Olive, D.J., Pelawa Watagoda, L.C.R., and Rupasinghe Arachchige Don, H.S. (2015), Visualizing and Testing the Multivariate Linear Regression Model, *International Journal of Statistics and Probability*, 4, 126-137.
- [19] Press, S.J. (2005), *Applied Multivariate Analysis: Using Bayesian and Frequentist Methods of Inference*, 2nd ed., Dover, Mineola, NY.
- [20] R Core Team (2016). R: A Language and Environment for Statistical Computing. Vienna: R Foundation for Statistical Computing. (www.R-project.org).
- [21] Ren, J.-J. (1991), On Hadamard Differentiability of Extended Statistical Functional, *Journal of Multivariate Analysis*, 39, 30-43.
- [22] Ren, J.-J., and Sen, P.K. (1995), Hadamard Differentiability on $D[0,1]^p$, *Journal of Multivariate Analysis*, 55, 14-28.
- [23] Rupasinghe Arachchige Don, H.S., and Pelawa Watagoda, L.C.R. (2017), Bootstrapping Analogs of the Two Sample Hotelling's T^2 Test, *Preprint at* (<http://lagrange.math.siu.edu/Olive/stwosample.pdf>).
- [24] Todorov, V., Filzmoser, P. (2010). Robust statistics for the one-way MAVOVA. *Com-*

putational Statistics & Data Analysis, 54:37-48.

- [25] Van Aelst, S., Willems, G. (2011). Robust and efficient one-way MANOVA tests. *Journal of the American Statistical Association*, 106:706-718.
- [26] Wilcox, R. R. (1995). Simulation results on solutions to the multivariate “Behrens-Fisher” problem via trimmed means. *The Statistician*, 44:213-225.
- [27] Willems, G., Pison, G., Rousseeuw, P. J., Van Aelst, S. (2002). A robust “Hotelling” test. *Metrika*, 55:125-138.
- [28] Zhang, J.-T., and Liu, X. (2013), A Modified Bartlett Test for Heteroscedastic One-Way MANOVA, *Metrika*, 76, 135–152.
- [29] Zhang, J., Olive, D. J., Ye, P. (2012). Robust covariance matrix estimation with canonical correlation analysis. *International Journal of Statistics and Probability*, 1:119-136.

VITA

Graduate School
Southern Illinois University

Hasthika S. Rupasinghe Arachchige Don

Date of Birth: August 15, 1985

1457 West Lake Road, Murphysboro, Illinois 62966

hasthika@siu.edu (hasthika@appstate.edu)

Southern Illinois University at Carbondale
Master of Science, Mathematics, August 2013

Special Honors and Awards: Dissertation Research Assistantship Award 2017.
John M. H. Olmsted Award for Outstanding Teaching Performance in the Department of Mathematics at SIU 2014.

Research Paper Title:

Bootstrapping Analogs of the One Way Manova Test

Major Professor: Dr. David J. Olive

Publications:

1. "Visualizing and Testing the Multivariate Linear Regression Model", International Journal of Statistics and Probability, January 22, 2015, with David J Olive, Lasanthi Watagoda.
2. "Bootstrapping analogs of the Hotelling's T^2 test", Communications in Statistics Theory and Methods, submitted, with Lasanthi Watagoda.
3. "Bootstrapping Analogs of the One Way Manova Test", work in progress, with David J Olive.