ABSTRACT

A Power Contrast of Tests for Homogeneity of Covariance Matrices in a
High-Dimensional Setting

Ben J. Barnard, Ph.D.

Chairperson: Dean M. Young, Ph.D.

Multivariate statistical analyses, such as linear discriminant analysis, MANOVA, and profile analysis, have a covariance-matrix homogeneity assumption. Until recently, homogeneity testing of covariance matrices was limited to the well-posed problem, where the number of observations is much larger than the data dimension. Linear dimension reduction has many applications in classification and regression but has been used very little in hypothesis testing for equal covariance matrices. In this manuscript, we first contrast the powers of five current tests for homogeneity of covariance matrices under a high-dimensional setting for two population covariance matrices using Monte Carlo simulations. We then derive a linear dimension reduction method specifically constructed for testing homogeneity of high-dimensional covariance matrices. We also explore the effect of our proposed linear dimension reduction for two or more covariance matrices on the power of four tests for homogeneity of covariance matrices under a high-dimensional setting for two- and three-population covariance matrices. We determine that our proposed linear dimension reduction method, when applied to the original data before using an appropriate test, can yield a substantial increase in power.

A Power Contrast of Tests for Homogeneity of Covariance Matrices in a
High-Dimensional Setting

by

Ben J. Barnard, B.S., M.S., M.S.Ed.

A Dissertation

Approved by the Department of Statistical Science

_____

James D. Stamey, Ph.D., Chairperson

Submitted to the Graduate Faculty of
Baylor University in Partial Fulfillment of the
Requirements for the Degree
of
Doctor of Philosophy

Approved by the Dissertation Committee

_____
Dean M. Young, Ph.D., Chairperson

_____
James D. Stamey, Ph.D.

_____
David J. Kahle, Ph.D.

_____
Dennis A. Johnston, Ph.D.

_____
Darryn S. Willoughby, Ph.D.

Accepted by the Graduate School
December 2018

_____
J. Larry Lyon, Ph.D., Dean

TABLE OF CONTENTS

LIST OF FIGURES

LIST OF TABLES

# CHAPTER ONE

## Introduction

In multivariate statistical analysis, if the data dimension $p$ is greater than the sample size $n$, then the corresponding sample covariance matrix is singular. We, therefore, cannot perform classical multivariate statistical methods because we are unable to invert the sample covariance matrix. The *high-dimensional* statistical inference problem concerns scenarios when the data dimension exceeds the sample size. High-dimensional covariance-matrix testing has become necessary because of the increasing ability to collect, store, and analyze high-dimensional data. The high-dimensional problem has been explored in the literature for covariance matrices in articles such as Schott (2007), Ledoit and Wolf (2002), and Srivastava (2005).

In this dissertation, we are interested in the problem of contrasting the power of tests for homogeneity of population covariance matrices in the high-dimensional setting (*HPCHDS*). In Chapter Two we contrast power curves of five two-population *HPCHDS* tests. We restrict our results to differences in the hyper-volume of two and three covariance matrices. We also perform randomization tests using these five *HPCHDS* tests on real data. Essentially, we determine that a test proposed by Srivastava et al. (2014) yields the largest omnibus power in our power simulations.

In Chapter Three we derive and apply a new linear dimension reduction (*LDR*) method for two covariance matrices and examine the effect on the power for the two-covariance matrices hypothesis testing problem for *HPCHDS*. More specifically, using five covariance matrix structures previously utilized in the *HPCHDS* literature, we apply our new *LDR* method and contrast the subsequent powers of four tests for *HPCHDS*. We also perform a randomization test on real data using our *LDR* method in combination with the four *HPCHDS* test statistics . We conclude that

while the test proposed by Chaipitak and Chongcharoen (2013) does not uniformly yield the largest power, it is the best omnibus *HPCHDS* test in the sample-size and high-dimensional scenarios examined here.

Last, in Chapter Four we derive a new *LDR* matrix for $(k > 2)$ high-dimensional population covariance matrices. We then devise a *LDR* matrix for $(k > 2)$ high-dimensional sample covariance matrices and examine the effect of *LDR* on the powers of four *HPCHDS* tests. Again, in our simulations, we use covariance matrix structures previously utilized in the *HPCHDS* literature to apply *LDR* and then contrast estimated powers of four high-dimensional tests for homogeneity of covariance matrices. Also, we perform randomization tests on a real dataset and contrast the effect of using linear dimension reduction versus the effect of using no dimension reduction on the powers on the four *HPCHDS* tests. In summary, we determine that the application of our proposed *LDR* method prior to the use of an appropriate *HPCHDS* test can yield a substantial power increase.

CHAPTER TWO

A Power Contrast of Five Tests for Homogeneity of Population Covariance Matrices
in a High-Dimensional Setting

*Abstract*

We compare and contrast the powers of five tests for testing for homogeneity of
two population covariance matrices in a high-dimensional setting for various sample-
size and parameter configuration scenarios. To examine the powers of the five tests, we
conduct relatively large Monte Carlo simulations for five population covariance-matrix
structures. Also, we perform permutation tests on a real dataset and compare the
results. We determine that a test proposed by Srivastava et al. (2014) generally yields
the largest power of the five competing tests for the covariance-matrix structures and
parameter configurations considered here, as well as for the leukemia dataset of Golub
et al. (1999). Furthermore, we determine that tests proposed by Schott (2007) and
Srivastava and Yanagihara (2010) yield relatively poor power.

## 2.1 Introduction

Data where the dimension is greater than the sample size, which is commonly
referred to as high-dimensional data, is an increasingly prominent complication in
present-day statistical applications. Many high-dimensional data applications arise in
disciplines such as genomics, portfolio analysis, and functional data imaging. Gener-
ally, one cannot use conventional multivariate analysis procedures to analyze these
data configurations because many of these statistical methods require the sample size
to be greater than the data dimension.

Here, we are interested in the relative efficacy of five recently-proposed hypoth-
esis tests for homogeneity of population covariance matrices in a high-dimensional
scenario (*HPCHDS*). That is, here we consider power contrasts for two population

*HPCHDS* hypothesis tests. When we conduct this *HPCHDS* hypothesis test, the respective null and alternative hypotheses are $H_0 : \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2$ and $H_A : \boldsymbol{\Sigma}_1 \neq \boldsymbol{\Sigma}_2$. Although multiple *HPCHDS* tests have been proposed, little work on contrasting the powers of these tests has been performed. Here, we compare and contrast the powers of five *HPCHDS* tests via Monte Carlo simulations for five population covariance-matrix structures.

We contrast the estimated powers of these *HPCHDS* tests as a function of the sample size and the data dimension for each covariance structure. Also, using these five *HPCHDS* tests on a real dataset, we contrast the results via permutation tests. Overall, we determine that an *HPCHDS* test proposed by Srivastava et al. (2014) generally yields the largest power of the five competing tests for the covariance-matrix structures and parameter configurations considered here. The test from Srivastava et al. (2014) also performs well on the Golub dataset.

In recent years, multiple *HPCHDS* tests have been proposed. Many such tests have been motivated by the work of Ledoit and Wolf (2004), who proposed a high-dimensional squared Frobenius norm (*HDSFN*) as a criterion for *HPCHDS* testing. For $\boldsymbol{\Sigma}_i \in \mathbb{R}_p^{\geq}, i = 1, 2$, we have that

$$HDSFN := \frac{1}{p} \operatorname{tr}\left(\boldsymbol{\Sigma}_1^2\right) + \frac{1}{p} \operatorname{tr}\left(\boldsymbol{\Sigma}_2^2\right) - \frac{2}{p} \operatorname{tr}\left(\boldsymbol{\Sigma}_1\boldsymbol{\Sigma}_2\right). \tag{2.1}$$

The divisor $p$ in (2.1) is typically omitted in the classical squared Frobenius norm for two real matrices but yields several desirable properties as described in Ledoit and Wolf (2004).

The general consensus in the *HPCHDS* literature is that the first *HPCHDS* test was introduced by Schott (2007). Other *HPCHDS* tests based on the *HDSFN* include those from Srivastava (2007), Srivastava et al. (2014), and Ahmad (2017). Ishii et al. (2016) have proposed an *HPCHDS* test that is based on the first principal component and the corresponding loadings, and employs noise-reduced eigenvalue estimators. Also, Srivastava and Yanagihara (2010) and Chaipitak and Chongcharoen

4

(2013) have proposed tests that use estimated summands of the *HDSFN*. In contrast, the methods of Chen et al. (2010) and Peng et al. (2016) have used banded estimators and transformations to increase the power of their tests.

We have structured the remainder of this paper as follows. In Section 2.2, we introduce notation and define estimators for the population means and population covariance matrices. In Section 2.3, we describe five recently published two-population *HPCHDS* tests. We then outline our Monte Carlo simulations for contrasting their powers in Section 2.4 and present our power simulation results for contrasting the powers in Section 2.5. We next apply the five *HPCHDS* tests to a gene-expression dataset in Section 2.6 and compare and contrast the test characteristics. Finally, we conclude with a brief discussion of our overall power results in Section 2.7.

<div align="center">

*2.2   Notation and Preliminaries*

</div>

We use the notation $\mathbb{R}_{m \times n}$ to represent the vector space of all $m \times n$ matrices over the real field $\mathbb{R}$. Also, we let $\mathbb{R}_n^S$ represent the set of all $n \times n$ symmetric matrices with entries from $\mathbb{R}$. In addition, we let $\mathbb{R}_n^{\geq}$ be the cone of all $n \times n$ symmetric positive semi-definite matrices over $\mathbb{R}$ and let $\mathbb{R}_n^{>}$ denote the interior of the cone $\mathbb{R}_n^{\geq}$. Moreover, we let $\mathcal{C}(\mathbf{A})$ represent the column space of $\mathbf{A} \in \mathbb{R}_{m \times n}$, and we let $\mathbf{I}_n$ signify the $n$-dimensional identity matrix. We define the transpose, trace, and rank of $\mathbf{A}$ by $\mathbf{A}^T$, $\mathrm{tr}(\mathbf{A})$, and $\mathrm{rank}(\mathbf{A})$, respectively. In addition, we use the notation $MN_{pn}(\mathbf{M}, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$ to denote a matrix-normal distribution, where $\mathbf{M} \in \mathbb{R}_{p \times n}$ is the mean matrix, $\boldsymbol{\Sigma} \in \mathbb{R}_p^{>}$, and $\otimes$ signifies the Kronecker product. Also, let $\mathbf{X}_i := \begin{bmatrix} \mathbf{x}_{i1} \vdots \mathbf{x}_{i2} \vdots \cdots \vdots \mathbf{x}_{in_i} \end{bmatrix} \in \mathbb{R}_{p \times n_i}$, represent a data matrix of $n_i$ observations sampled from the $i^{th}$ population, so that $\mathbf{X}_i \sim MN_{pn_i}(\mathbf{M}_i, \mathbf{I}_{n_i} \otimes \boldsymbol{\Sigma}_i)$ with $\mathbf{M}_i \in \mathbb{R}_{p \times n_i}$, $\boldsymbol{\Sigma}_i \in \mathbb{R}_p^{>}$, and $i = 1, 2$. Therefore, $\mathbf{x}_{ij} \sim N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ are independent $p$-dimensional random vectors for $i = 1, 2$, and $j = 1, 2, ..., n_i$. In addition, we use the horizontally concatenated matrix $\mathbf{X} := \begin{bmatrix} \mathbf{X}_1 \vdots \mathbf{X}_2 \end{bmatrix} \in \mathbb{R}_{p \times N}$ to represent the complete data matrix, where $N := \sum_{i=1}^{2} n_i$.

<div align="center">

5

</div>

For $i = 1, 2$, we use the $i^{th}$ sample mean

$$\overline{\mathbf{x}}_i := \frac{1}{n_i} \mathbf{X}_i \mathbf{j}$$

to estimate the $i^{th}$ population mean $\mu_i$, where $\mathbf{j}$ is an $n_i \times 1$ vector of ones. Also, we use the $i^{th}$ sample covariance matrix

$$\mathbf{S}_i := \frac{1}{n_i - 1} \mathbf{X}_i \left( \mathbf{I}_{n_i} - \frac{1}{n_i} \mathbf{J}_{n_i} \right) \mathbf{X}_i^T, \tag{2.2}$$

to estimate the $i^{th}$ population covariance matrix, where $\mathbf{J}_{n_i}$ is an $n_i \times n_i$ matrix of ones. Additionally, let $\mathbf{V}_i := (n_i - 1)\mathbf{S}_i$ represent the Gram matrix of (2.2) and, finally, let

$$\mathbf{S} := \frac{1}{n} \sum_{i=1}^{k} \mathbf{V}_i$$

denote the pooled sample covariance, where $n := N - 2$.

### 2.3  Five Tests for the Homogeneity of Two Population Covariance Matrices Performed Under a High-Dimensional Setting

We now describe five *HPCHDS* tests for two population covariance matrices. Namely, we consider the tests derived in Schott (2007), Srivastava and Yanagihara (2010), Chaipitak and Chongcharoen (2013), Srivastava et al. (2014), and Ahmad (2017). We also state their asymptotic distributions under the null hypothesis for these tests.

#### 2.3.1  Schott (2007)

Schott (2007) has proposed a test for *HPCHDS* based on *HDSFN* given in (2.1). Because researchers routinely examine thousands of gene expressions with sample sizes typically less than $N = 100$, Schott (2007) used a DNA microarray-data example to motivate his work. For notational consistency throughout this paper, we do not use the original notation of the test from Schott (2007). Srivastava and Yanagihara (2010)

have expressed the test statistic proposed in Schott (2007) as

$$Q_{Sc} := \frac{\hat{a}_{21} + \hat{a}_{22} - \frac{2}{p}\operatorname{tr}(\mathbf{S}_1\mathbf{S}_2)}{\sqrt{\widehat{\operatorname{Var}(q_{Sc})}}}, \tag{2.3}$$

where $\hat{a}_{21}$ and $\hat{a}_{22}$ are given in (A.2), $q_{Sc} := \hat{a}_{21} + \hat{a}_{22} - \frac{2}{p}\operatorname{tr}(\mathbf{S}_1\mathbf{S}_2)$ estimates the sum of squared elements of $[\mathbf{\Sigma}_2 - \mathbf{\Sigma}_1]$, and

$$\widehat{\operatorname{Var}(q_{Sc})} = 4\hat{a}_2^2 \left(\frac{1}{n_1 - 1} + \frac{1}{n_2 - 1}\right)^2.$$

We refer to the *HPCHDS* test using the test statistic $Q_{Sc}$ by $T_{Sc}$. Schott (2007) has shown that $\widehat{\operatorname{Var}(q_{Sc})} \xrightarrow{P} \operatorname{Var}(q_{Sc})$ if H$_0$ is true and also that $Q_{Sc} \overset{\cdot}{\sim} N(0,1)$ under H$_0$ as $(p, n_1, n_2) \to \infty$. However, because of a dearth of competing *HPCHDS* tests at the time of publication, this paper includes only a contrast between the power of $T_{Sc}$ and the power of the likelihood ratio test for $n_i > p$, $i = 1, 2$.

### 2.3.2 *Srivastava and Yanagihara (2010)*

Next, Srivastava and Yanagihara (2010) have proposed a test based on an estimator of $\left[\operatorname{tr}(\mathbf{\Sigma}_i^2) / \operatorname{tr}(\mathbf{\Sigma}_i)^2\right]$, $i = 1, 2$. That is, the test presented by Srivastava and Yanagihara (2010) has compared the sum of squared elements to the sum of the squared eigenvalues via a ratio. Their *HPCHDS* test statistic is

$$Q_{S10} := \frac{\hat{a}_{21}/\hat{a}_{11}^2 - \hat{a}_{22}/\hat{a}_{12}^2}{\sqrt{\widehat{\operatorname{Var}(q_{S10})}}},$$

where $\hat{a}_{21}$ and $\hat{a}_{22}$ are given in (A.2), $\hat{a}_{11}$ and $\hat{a}_{12}$ are given in (A.1), $q_{S10} := \hat{a}_{21}/\hat{a}_{11}^2 - \hat{a}_{22}/\hat{a}_{12}^2$, and

$$\widehat{\operatorname{Var}(q_{S10})} = \left(\frac{4\hat{a}_2^2}{\hat{a}_1^4} + \frac{2}{p}\left(\frac{\hat{a}_2^3}{\hat{a}_1^6} - \frac{2\hat{a}_2\hat{a}_3}{\hat{a}_1^5} + \frac{\hat{a}_4}{\hat{a}_1^4}\right)\right)\left(\frac{1}{(n_1 - 1)} + \frac{1}{(n_2 - 1)}\right),$$

where $\hat{a}_2$ is given in (A.3), $\hat{a}_3$ is defined in (A.4), and $\hat{a}_4$ is defined in (A.5). We refer to the *HPCHDS* test using the test statistic $Q_{S10}$ by $T_{S10}$.

Srivastava and Yanagihara (2010) have proven that $\widehat{\operatorname{Var}(q_{S10})} \xrightarrow{P} \operatorname{Var}(q_{S10})$. Also, provided H$_0$ holds, they have shown that $Q_{S10} \overset{\cdot}{\sim} N(0,1)$ as $(p, n_1, n_2) \to \infty$.

7

Furthermore, they have provided simulated power contrast of the power of $T_{S10}$ and the power of $T_{Sc}$.

### 2.3.3 Chaipitak and Chongcharoen (2013)

Chaipitak and Chongcharoen (2013) have developed a test based on an estimator of the ratio $\text{tr}\left(\mathbf{\Sigma}_i^2\right) / \text{tr}\left(\mathbf{\Sigma}_j^2\right)$, where $i, j = 1, 2$, and $i \neq j$. Their test statistic is

$$Q_C := \frac{\hat{a}_{21}/\hat{a}_{22} - 1}{\sqrt{\widehat{\text{Var}(q_C)}}},$$

where $\hat{a}_{21}$ and $\hat{a}_{22}$ are given in (A.2), $q_C := \hat{a}_{21}/\hat{a}_{22} - 1$, and

$$\widehat{\text{Var}(q_C)} = 4\left\{\frac{2\hat{a}_4^*}{p\hat{a}_2^2}\left(\frac{1}{n_1 - 1} + \frac{1}{n_2 - 1}\right) + \left(\frac{1}{(n_1 - 1)^2} + \frac{1}{(n_2 - 1)^2}\right)\right\},$$

where $\hat{a}_2$ is given in (A.3) and $\hat{a}_4^*$ is defined in (A.6). We refer to the *HPCHDS* test using the test statistic $Q_C$ by $T_C$.

Chaipitak and Chongcharoen (2013) have shown that $\widehat{\text{Var}(q_C)} \xrightarrow{P} \text{Var}(q_C)$, assuming $H_0$ is true, and have also shown that $Q_C \overset{\cdot}{\sim} N(0, 1)$ under $H_0$ as $(p, n_1, n_2) \to \infty$. Additionally, they have compared and contrasted powers of $T_C$, $T_{Sc}$, and $T_{S10}$ for four different covariance matrix structures for the two-population covariance matrices case.

### 2.3.4 Srivastava et al. (2014)

Srivastava et al. (2014) have improved upon the test (2.3) by replacing $\hat{a}_{2i}$ and $\hat{a}_2$ with the unbiased, consistent estimators $\hat{a}_{2si}$ and $\hat{a}_{2s}$, respectively. For $i = 1, 2$, they let

$$\hat{a}_{2si} := \frac{(n_i - 2)(n_i - 1)\text{tr}\left(\mathbf{V}_i^2\right) - n(n - k)\text{tr}\left(\mathbf{D}_i^2\right) + \text{tr}\left(\mathbf{V}_i\right)^2}{pn_i(n_i - 1)(n_i - 2)(n_i - 3)}$$

and

$$\hat{a}_{2s} := \frac{1}{N}\left(\sum_{i=1}^{k} n_i\hat{a}_{2i}\right),$$

where $\mathbf{D}_i := \mathrm{diag}\left(\mathbf{b}_{i1}^T\mathbf{b}_{i1}, \ldots, \mathbf{b}_{in_i}^T\mathbf{b}_{in_i}\right)$, $\mathbf{b}_{ij} = (\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)$, $i = 1, 2$, and $j = 1, \ldots, n_i$. We simplify their test statistic to

$$Q_{S14} := \frac{\hat{a}_{21_{Sr}} + \hat{a}_{22_{Sr}} - \frac{2}{p}\,\mathrm{tr}\,(\mathbf{S}_1\mathbf{S}_2)}{\sqrt{\widehat{\mathrm{Var}(q_{S14})}}},$$

where $q_{S14} := \hat{a}_{21_{Sr}} + \hat{a}_{22_{Sr}} - \frac{2}{p}\,\mathrm{tr}\,(\mathbf{S}_1\mathbf{S}_2)$ and

$$\widehat{\mathrm{Var}(q_{S14})} = 4\hat{a}_{2s}^2\left(\frac{1}{(n_1 - 1)^2} + \frac{1}{(n_2 - 1)^2}\right).$$

We refer to the *HPCHDS* test using the test statistic $Q_{S14}$ by $T_{S14}$.

Srivastava et al. (2014) have also shown that $\widehat{\mathrm{Var}(q_{S14})} \xrightarrow{P} \mathrm{Var}\,(q_{S14})$, assuming $\mathrm{H}_0$ is true, and that $Q_{S14} \stackrel{.}{\sim} N(0, 1)$ under $\mathrm{H}_0$ as $(p, n_1, n_2) \to \infty$. Additionally, they have shown that the power of $T_{S14}$ is larger than the power of $T_{Sc}$ for one type of covariance matrix structure.

*2.3.5  Ahmad (2017)*

Ahmad (2017) has proposed one of the newest tests in the *HPCHDS* literature that is based on the squared Frobenius norm without the $p$-divisor. We give the test statistic from Ahmad (2017) as

$$Q_A := \frac{E_1 + E_2 - 2E_{12}}{\sqrt{\widehat{Var(q_A)}}},$$

where $E_{12} := \mathrm{tr}\,(\mathbf{S}_1\mathbf{S}_2)$ and $q_A := E_1 + E_2 - 2E_{12}$. For $i = 1, 2$,

$$E_i := \frac{(n_i - 1)}{n_i\,(n_i - 2)\,(n_i - 3)}\Bigg\{ (n_i - 1)\,(n_i - 2)\,\mathrm{tr}\,(\mathbf{S}_i)^2 + \left[\,\mathrm{tr}\,(\mathbf{S}_i)\,\right]^2$$
$$- \frac{n_i}{(n_i - 1)}\sum_{j=1}^{n_i}(\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)^T(\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)(\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)^T(\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)\Bigg\}$$

and

$$\widehat{\mathrm{Var}(q_A)} = 4\left[\,\mathrm{tr}\,(\mathbf{S}^2)\,\right]^2\left(\frac{1}{n_1} + \frac{1}{n_2}\right)^2.$$

Ahmad (2017) has shown that $\mathbf{E}_{12} \xrightarrow{P} \mathrm{tr}\,(\mathbf{\Sigma}_1\mathbf{\Sigma}_2)$ and $\mathbf{E}_i \xrightarrow{P} \mathrm{tr}\,(\mathbf{\Sigma}_i^2)$, $i = 1, 2$, and has further demonstrated that $\widehat{\mathrm{Var}(q_A)} \xrightarrow{P} \mathrm{Var}\,(q_A)$, assuming $\mathrm{H}_0$ is true. Additionally,

under $H_0$, he has shown that $Q_A \overset{\cdot}{\sim} N(0,1)$ as $(p, n_1, n_2) \to \infty$. However, the power of $T_A$ has not been contrasted with the power of any other *HPCHDS* tests.

## 2.4  Monte Carlo Power Simulation Description

We now describe the Monte Carlo simulation designs we used to contrast $POW(T_{Sc})$, $POW(T_{S14})$, $POW(T_A)$, $POW(T_{S10})$, and $POW(T_C)$, where $POW(T_{(*)})$ represents the estimated power of the test $T_{(*)}$.

### 2.4.1  Simulation Covariance Structures

The covariance matrix structures used in our Monte Carlo simulations were selected from the *HPCHDS* tests literature. We have compared test powers across five different covariance matrix structures with balanced group sample sizes of 10.

First, we used the constant-times-identity covariance-matrix structure. For our simulation, the parameters for the null and alternative hypotheses are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{I}_p$$

and

$$H_A : \mathbf{\Sigma}_1 = \mathbf{I}_p \text{ and } \mathbf{\Sigma}_2 = \sigma^2 \mathbf{I}_p.$$

Second, we utilized the compound-symmetric covariance matrix class. For our simulation, the null and alternative hypotheses for testing for homogeneity of compound-symmetric population covariance matrices are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \sigma_1^2 \mathbf{I}_p + \sigma_2^2 \mathbf{J}_p$$

and

$$H_A : \mathbf{\Sigma}_1 = \sigma_1^2 \mathbf{I}_p + \sigma_2^2 \mathbf{J}_p \text{ and } \mathbf{\Sigma}_2 = \sigma_{1A}^2 \mathbf{I}_p + \sigma_{2A}^2 \mathbf{J}_p.$$

Third, we used the autoregressive covariance-matrix structure. Here, the null and alternative hypotheses for the homogeneity of autoregressive covariance matrices

are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{U}_0$$

and

$$H_A : \mathbf{\Sigma}_1 = \mathbf{U}_0 \text{ and } \mathbf{\Sigma}_2 = \mathbf{U}_1,$$

where $\mathbf{U}_0 = \sigma_{ij} = 0.1^{|i-j|}$ and $\mathbf{U}_1 = \sigma_{ij} = 0.3^{|i-j|}$ and where $1 \leq i, j \leq p$.

Fourth, we used the heterogeneous autoregressive covariance-matrix structure. These heterogeneous autoregressive covariance-matrix structures are similar to those in Srivastava et al. (2014), and were created as follows. First, let $\sigma_l := 1 + (-1)^{l+1} Q_l / 2$, where $Q_l \sim Unif(0,1)$ and $l = 1, 2, \ldots, p$. Then, the null and alternative hypotheses for the homogeneity of heterogeneous autoregressive covariance matrices are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \sigma_i \sigma_j 0.1^{|i-j|^{\frac{1}{10}}}$$

and

$$H_A : \mathbf{\Sigma}_1 = \sigma_i \sigma_j 0.1^{|i-j|^{\frac{1}{10}}} \text{ and } \mathbf{\Sigma}_2 = \sigma_i \sigma_j 0.3^{|i-j|^{\frac{1}{10}}},$$

respectively, where $1 \leq i, j \leq p$.

Last, we examined an unstructured covariance-matrix configuration that has no discernible pattern, which we modeled as

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{U}_1$$

and

$$H_A : \mathbf{\Sigma}_1 = \mathbf{U}_1 \text{ and } \mathbf{\Sigma}_2 = \mathbf{U}_2.$$

Here, the parameters for the null and alternative hypotheses are

$$\mathbf{U}_1 = \sigma_{ij} = \sigma_{ji} := \begin{cases} (-1)^{i+j} \left( \frac{0.10}{j} \right), & i < j \\ 1, & i = j \end{cases}$$

and

$$\mathbf{U}_2 = \sigma_{ij} = \sigma_{ji} := \begin{cases} (-1)^{i+j} \left( \frac{0.05}{j} \right), & i < j \\ 1, & i = j. \end{cases}$$

### 2.4.2  Simulated Critical Values

First, we used Monte Carlo simulations to check the accuracy of the suggested asymptotic critical values and corresponding significance levels of the *HPCHDS* tests $T_A$, $T_{Sc}$, $T_{S10}$, $T_C$, and $T_{S14}$. The results of these simulations are given in Appendix A. Because many of the asymptotic significance levels were so different from $\alpha = 0.05$, we simulated critical values for each of the five test for various combinations of data dimension and class sample size. More specifically, we generated 10,000 independent multivariate normal datasets of $n_i$ observations, $i = 1, 2$, from a $N_p(\mathbf{0}, \mathbf{\Sigma})$ distribution, where $\mathbf{\Sigma}$ is the common population covariance matrix if $H_0$ is true. We next determined the appropriate simulated critical values for each test by calculating

$$SCV_{1-\alpha/2} := \inf \left\{ x \in \mathbb{R} : 1 - \alpha/2 \leq \widehat{F}_{T_{(*)}}(x) \right\},$$

and

$$SCV_{\alpha/2} := \sup \left\{ x \in \mathbb{R} : \alpha/2 \geq \widehat{F}_{T_{(*)}}(x) \right\},$$

where $\widehat{F}_{T_{(*)}}(x)$ is the empirical *CDF* for the test $T_{(*)}$ and $\alpha = 0.05$.

### 2.4.3  Powers of the HPCHDS tests

To estimate the power of each *HPCHDS* test for a given $p$ and $n_i, i = 1, 2$, we generated 10,000 independent multivariate normal datasets from each of the populations modeled as $N_p(\mathbf{0}, \mathbf{\Sigma}_i)$, $i = 1, 2$, where $\mathbf{\Sigma}_i$ are the population covariance matrices under the alternative hypothesis. Using the alternative-hypothesis parameters, we calculated the test values $T_{A,j}$, $T_{Sc,j}$, $T_{S10,j}$, $T_{C,j}$, $T_{S14,j}$ for each of the $j$ simulated datasets, $1 \leq j \leq 10,000$. We then calculated the estimated powers for each test $T_{(*)}$

12

and for each common sample size and data dimension using

$$POW(T_{(*)}) := \frac{\sum_{j=1}^{10,000} I\left[T_{(*),j} \in RR(T_{(*)})\right]}{10,000},$$

where $RR(T_{(*)})$ is the rejection region for the test $T_{(*)}$ and $I[\cdot]$ is the indicator function. We performed the power simulations in parallel using `R` and the `covTestR` package.

### 2.4.4  Simulation Design Summary

Tables [2.1, A.6– A.9] display the simulated powers for of the *HPCHDS* tests. We performed Monte Carlo simulations to contrast the powers of the *HPCHDS* tests for common sample sizes of $n_i \in \{5, 10, 15, 20\}, i = 1, 2$, and data dimensions of $p \in \{20, 40, 80, 160\}$. In Figures $[2.1 - 2.5]$, we have displayed Monte Carlo simulation power-curve comparisons plotted against common sample sizes $n_i \in \{5, 10, 15, 20\}$, $i = 1, 2$, with $p = 160$. For Figures $[A.1 - A.5]$, we have fixed the common sample size at $n_i = 10$ and have conducted simulations for $p \in \{11, 12, \ldots, 160\}$. We computed the simulations in parallel using `R` and the `covTestR` package.

### 2.5  Simulations Contrasting the Powers of Five Tests for Homogeneity of Two Population Covariance Matrices in a High-Dimensional Setting

Here, we present power curve figures for $p = 160$ and common sample sizes ranging from $n_i = 5$ to $n_i = 40, i = 1, 2$. We also include power curves for fixed sample size $n_i = 10$ and $p = 2, 3, \ldots, 160$. The power curves were fitted using generalized linear models with b-splines.

### 2.5.1  A Power-Simulation Summary Table

Table 2.1 shows $POW(T_A), POW(T_C), POW(T_{Sc}), POW(T_{S10})$, and $POW(T_{S14})$ when we simulated the five *HPCHDS* tests for two heterogeneous autoregressive covariance-matrix structure described in Subsection 2.4.1. The columns of Table 2.1 correspond to the tests $T_A$, $T_C$, $T_{Sc}$, $T_{S10}$, and $T_{S14}$, and the rows are the combina-

tions of data dimension and sample size in ascending order. We see that $POW(T_{S14})$ yielded the largest power for each combination of sample size and data dimension considered here and that $POW(T_C)$ and $POW(T_{Sc})$ tended to yield substantially inferior powers for most data dimensions and sample sizes.

Table 2.1: A table contrasting $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing $HPCHDS$ on two heterogeneous autoregressive population covariance structures.

| $p$ | $n_1 = n_2$ | $POW(T_A)$ | $POW(T_C)$ | $POW(T_{Sc})$ | $POW(T_{S10})$ | $POW(T_{S14})$ |
|-----|------------|-----------|-----------|--------------|---------------|---------------|
| 20 | 5 | 0.12 | 0.06 | 0.06 | 0.06 | 0.13 |
| | 10 | 0.21 | 0.13 | 0.14 | 0.22 | 0.25 |
| | 15 | 0.40 | 0.26 | 0.20 | 0.30 | 0.45 |
| 40 | 5 | 0.21 | 0.07 | 0.08 | 0.07 | 0.21 |
| | 10 | 0.38 | 0.17 | 0.19 | 0.35 | 0.46 |
| | 15 | 0.45 | 0.23 | 0.38 | 0.50 | 0.48 |
| | 20 | 0.58 | 0.30 | 0.43 | 0.71 | 0.78 |
| 80 | 5 | 0.24 | 0.11 | 0.10 | 0.08 | 0.29 |
| | 10 | 0.52 | 0.25 | 0.27 | 0.47 | 0.62 |
| | 15 | 0.64 | 0.33 | 0.34 | 0.69 | 0.78 |
| | 20 | 0.80 | 0.42 | 0.58 | 0.86 | 0.93 |
| 160 | 5 | 0.34 | 0.18 | 0.14 | 0.12 | 0.40 |
| | 10 | 0.67 | 0.42 | 0.38 | 0.64 | 0.78 |
| | 15 | 0.84 | 0.57 | 0.56 | 0.83 | 0.92 |
| | 20 | 0.92 | 0.66 | 0.73 | 0.94 | 0.97 |

### 2.5.2 Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for Constant-Times-Identity Covariance Matrix Structures

Figure 2.1 presents curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for two population covariance matrices with a constant-times-identity covariance matrix structure with parameters $\mathbf{\Sigma}_1 = \mathbf{I}_p$ and $\mathbf{\Sigma}_2 = (1.5)\mathbf{I}_p$. The data dimension was $p = 160$, and common sample sizes were $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$. For all common sample sizes considered here $POW(T_{S14}) = 1.0$ uniformly. In addition, $POW(T_A)$ and $POW(T_C)$ yielded similar power curves with $POW(T_A) > POW(T_C)$. Both $POW(T_A)$ and $POW(T_C)$ attained a power of 1.0 for $n_1 = n_2 \geq 17$, and $POW(T_{S10})$ was essentially negligible regardless of the value of $n_i$, $i = 1, 2$. Finally,

$POW(T_{Sc})$ attained mild gains as $n_i$ increased and attained a maximum power value of 0.30 at $n_i = 40, i = 1, 2$.



Figure 2.1: Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing $HPCHDS$ with two constant-times-identity population co-variance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{I}_p$ and $\mathbf{\Sigma}_2 = (1.5)\mathbf{I}_p$. The common sample sizes were $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$, with $p = 160$.

### 2.5.3 Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for Two Compound-Symmetric Covariance Matrix Structures

Figure 2.2 displays the curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing for $HPCHDS$ having two compound-symmetric covariance matrix structures with parameters $\mathbf{\Sigma}_1 = (0.99)\mathbf{I}_p + (0.01)\mathbf{J}_p$ and $\mathbf{\Sigma}_2 = (0.95)\mathbf{I}_p + (0.05)\mathbf{J}_p$. The data dimension was $p = 160$, and the common sample size was $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$. Here, $POW(T_{S14})$ yielded the largest power value of 0.90

when $n_i = 40, i = 1, 2$. Also, $POW(T_{Sc})$ and $POW(T_C)$ produced similar curves with $POW(T_C)$ being slightly larger than $POW(T_{Sc})$ for all considered values of $n_i, i = 1, 2$. In addition, $POW(T_A)$ yielded the largest power for smaller common sample sizes but was overtaken by $POW(T_{S14})$ near $n_1 = n_2 = 10$. Moreover, $POW(T_{S10})$ yielded the smallest power value for $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$.



Figure 2.2: Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for $HPCHDS$ with two compound-symmetric population covariance matrices with parameters $\mathbf{\Sigma}_1 = (0.99)\mathbf{I}_p + (0.01)\mathbf{J}_p$ and $\mathbf{\Sigma}_2 = (0.95)\mathbf{I}_p + (0.05)\mathbf{J}_p$. The common sample sizes were $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$, with $p = 160$.

### 2.5.4 Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for Two Autoregressive Covariance Matrix Structures

Figure 2.3 presents the power curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for two autoregressive population covariance struc-

16

tures with parameters $\mathbf{\Sigma}_1 = 0.1^{|i-j|}$ and $\mathbf{\Sigma}_2 = 0.3^{|i-j|}$, $p = 160$, and common sample $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$. As $n_i$ increased, $POW(T_{S14})$ attained the largest power values, as shown in Figure 2.3. Also, $POW(T_{Sc})$ gained little power as $n_i$ increased, and $POW(T_{S10})$ did not show increased power until $n_i = 25$, $i = 1, 2$. In addition, $POW(T_A)$ produced the largest power for small common sample sizes, but $POW(T_{S14})$ yielded superior power for all $n_i \geq 15, i = 1, 2$.



Figure 2.3: Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing *HPCHDS* with two autoregressive population covariance matrix structures with parameters $\mathbf{\Sigma}_1 = 0.1^{|i-j|}$ and $\mathbf{\Sigma}_2 = 0.3^{|i-j|}$. The common sample sizes were $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$, with $p = 160$.

Figure 2.4 presents power curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$,
and $POW(T_{S14})$, for testing $HPCHDS$ for two population covariance matrices with
heterogeneous autoregressive covariance matrices. Here, $p = 160$, and the common
sample size was $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$. We observed that $POW(T_{S14})$ was supe-
rior for all $n_i$ and attained a maximum $POW(T_{S14}) = 1.0$ at $n_i = 30$. Also, $POW(T_A)$
and $POW(T_{S10})$ produced comparable curves attaining maximum powers of 1.0 for
$n_i \geq 30, i = 1, 2$. In contrast, $POW(T_{Sc})$ and $POW(T_C)$ yielded the two smallest
curves for all $n_i$ considered here.



Figure 2.4: Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and
$POW(T_{S14})$ for testing the $HPCHDS$ with two heterogeneous autoregressive covari-
ance matrix structures. The common sample sizes were $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$, for
$p = 160$.

### 2.5.6 Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for Two Unstructured Population Covariance Matrices

Figure 2.5 displays power curves for the five *HPCHDS* tests for two unstructured covariance matrices. Here, $p = 160$, and the common sample size was $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$. For all values of $n_i$, $POW(T_{S14})$ was superior and attained a maximum $POW(T_{S14}) = 1.0$ at $n_i = 32$. Also, $POW(T_C)$ and $POW(T_{Sc})$ curves were similar and were uniformly least powerful.



Figure 2.5: Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing the *HPCHDS* with two unstructured population covariance matrix structures. The common sample sizes were $n_i \in \{5, 6, ..., 40\}$, $i = 1, 2$, for $p = 160$.

## 2.6 A Real-Data Example for Testing the Homogeneity of Two High-Dimensional Population Covariance Matrices Contrasting the Tests $T_A$, $T_C$, $T_{Sc}$, $T_{S10}$, and $T_{S14}$

In this section, we evaluate the differences in the performance of the *HPCHDS* tests $T_A, T_C, T_{Sc}, T_{S10}$, and $T_{S14}$ using real data. The data we use originated in Golub et al. (1999). However, we utilized the dataset from Ramey (2016). The Golub dataset contains 47 patients with acute lymphoblastic leukemia, which can be further separated into luekemia from T-cell and B-cell lymphocytes, and 25 patients with acute myeloid leukemia.

Because we could not rely on the accuracy of suggested asymptotic critical values, we performed permutation tests for each *HPCHDS* statistic considered here and, thus, assumed the observations were exchangeable. We first determined a critical value for each statistic and then compared each of the original empirical test scores with the corresponding critical values using the permutation test procedure. Table 2.2 presents the results of our five *HPCHDS* tests using permutation tests. For the tests $T_A$, $T_{S10}$, and $T_{S14}$, we rejected the null hypothesis that the population covariance matrix for the subjects with acute lymphoblastic leukemia was homogeneous to the covariance matrix for subjects with acute myeloid leukemia. In contrast, we failed to reject the hypothesis of homogeneous population covariance matrices when utilizing the tests $T_C$ and $T_{Sc}$. This power-contrast result is a similar conclusion to that obtained in the simulated power studies in the previous section.

Table 2.2: Test-results summary table of the *HPCHDS* tests $T_A$, $T_C$, $T_{Sc}$, $T_{S10}$, and $T_{S14}$ applied to the Golub et al. (1999) data set.

| Test | Lower Crit. Val. | Upper Crit. Val. | test Val. | $p$-value | Decision |
|------|-----------------|------------------|-----------|-----------|----------|
| $T_C$ | -4.198 | 4.557 | 4.258 | 0.056 | FTR $H_0$ |
| $T_{Sc}$ | -9.648 | 10.611 | 9.230 | 0.051 | FTR $H_0$ |
| $T_{S10}$ | -4.449 | 5.004 | 6.563 | 0.011 | Reject $H_0$ |
| $T_{S14}$ | -5.237 | 4.782 | 5.441 | 0.024 | Reject $H_0$ |
| $T_A$ | -1.736 | 1.273 | 1.312 | 0.045 | Reject $H_0$ |

## 2.7 Discussion

In summary, we have compared and contrasted the powers for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for two population covariance matrices. We have examined the powers for these five $HPCHDS$ tests for five different covariance structures via Monte Carlo simulations and have shown that, except for very small sample sizes, the test proposed by Srivastava et al. (2014) was generally the most powerful of the competing $HPCHDS$ tests considered here. Also, we have shown that $POW(T_C)$ and $POW(T_{Sc})$ were inferior to $POW(T_{S14})$ for all considered sample sizes and data dimensions considered here.

Finally, we contrasted the characteristics of the $HPCHDS$ tests $T_A$, $T_C$, $T_{Sc}$, $T_{S10}$, and $T_{S14}$ for testing two population covariance matrices using permutation tests on the Golub dataset. For the tests $T_A$, $T_{S10}$, and $T_{S14}$, we rejected the null hypothesis of covariance homogeneity at the $\alpha = 0.05$ level. These test results mirrored the conclusions determined from our Monte Carlo power simulations. Thus, we concluded that $T_{S14}$ yielded superior power to the $HPCHDS$ tests $T_C$ and $T_{Sc}$ and is marginally more powerful than $T_A$ and $T_{S10}$.

CHAPTER THREE

Linear Dimension Reduction for Power Improvement of Tests for Homogeneity of
Two Population Covariance Matrices in a High-Dimensional Setting

ABSTRACT

We develop a linear dimension reduction ($LDR$) technique to improve the power of tests for the homogeneity of two population covariance matrices in a high-dimensional scenario ($HPCHDS$). Using Monte Carlo simulations, we contrast the powers of four $HPCHDS$ tests calculated with reduced-dimension data from our $LDR$ method with the powers of these four tests calculated from the original data. We also perform a permutation tests using real data to contrast the no-$LDR$ and post-$LDR$ test characteristics. Our proposed $LDR$ technique yields substantial power increases for certain $HPCHDS$ tests considered here. We conclude that the test of Chaipitak and Chongcharoen (2013), when calculated with reduced-dimensional data using our $LDR$ method, yields the best power for most of the considered population covariance structures and on the a high-dimensional real dataset from Alon et al. (1999).

### 3.1   Introduction

In many scientific disciplines, including biomedical imaging, magnetic resonance imaging, tomography, and financial portfolio analysis, one may collect data where the data dimension is greater than the group sample size. We label this type of data as "high-dimensional data." For a fixed sample size, increasing the data dimension increases the estimator variability, thus making statistical inference more inaccurate. Also, if the data dimension is greater than the group sample size, then the corresponding sample covariance matrix is singular and, therefore, one cannot conduct classical multivariate statistical analysis.

In this paper, we derive and apply an *LDR* matrix for two sample covariance matrices to reduce the number of estimated parameters, thus often yielding more powerful hypothesis tests for *HPCHDS*. In particular, we compare the powers of four *HPCHDS* tests calculated after applying our proposed *LDR* method on the original data to the powers of these four tests calculated from the original data. We employ Monte Carlo simulations with five covariance structures previously utilized in the *HPCHDS* literature. We restrict the differences in the covariance matrices to the differences in the hyper-volume as measured by the determinant while ignoring differences in the eigenvector orientation. We also contrast the four *HPCHDS* tests calculated with no-*LDR* data against the four tests calculated with post-*LDR* data using a permutation test procedure on a real high-dimensional dataset. For both the simulations and the real dataset, we determine that the proposed LDR method, when used with appropriate tests, can yield a considerable increase in power.

The current consensus in the *HPCHDS* literature is that Schott (2007) was the first to introduce a *HPCHDS* test which has been based on a high-dimensional squared Frobenius norm (*HDSFN*) for two symmetric nonegative-definite matrices given by Ledoit and Wolf (2004). This norm is

$$HDSFN := \frac{1}{p}\operatorname{tr}\left(\boldsymbol{\Sigma}_1^2\right) + \frac{1}{p}\operatorname{tr}\left(\boldsymbol{\Sigma}_2^2\right) - \frac{2}{p}\operatorname{tr}\left(\boldsymbol{\Sigma}_1\boldsymbol{\Sigma}_2\right). \tag{3.1}$$

The inclusion of the divisor $p$ in (3.1) yields several desirable properties that one can find in Ledoit and Wolf (2004). Some *HPCHDS* tests motivated by summands of (3.1) include those of Srivastava (2007), Srivastava et al. (2014), and Ahmad (2017).

Also, Ishii et al. (2016) have proposed a test based on the first principal components from each sample covariance matrix and the corresponding loadings composed of noise-reduced estimators. Srivastava and Yanagihara (2010) and Chaipitak and Chongcharoen (2013) have proposed *HPCHDS* tests that use ratios of the summands of (3.1). In contrast, the methods of Chen et al. (2010) and Peng et al. (2016) use banded estimators and transformations to increase the power of their tests.

The remainder of the paper is organized as follows. In Section 3.2 we define notation used throughout the remainder of the paper. We then present four two population *HPCHDS* tests in Section 3.3. Next, in Section 3.4, we describe a new *LDR* method that for two-population *HPCHDS* tests to ostensibly increase their power. We then describe our simulation design and the covariance matrix structures used for the simulations in Section 3.5. In Section 3.6, we present our power-contrast results from the Monte Carlo simulations for the tests from Chaipitak and Chongcharoen (2013), Ahmad (2017), Srivastava and Yanagihara (2010), and Ishii et al. (2016), calculated from the original data with post-*LDR* data. We then apply these *HPCHDS* tests using no-*LDR* and post-*LDR* data to a real high-dimensional dataset in Section 3.7 and contrast the corresponding p-values and hypothesis decision results. Finally, we offer some concluding remarks in Section 3.8.

### 3.2   Notation

We use the notation $\mathbb{R}_{m \times n}$ and $\mathbb{R}_{n \times n}$ to represent the vector space of all $m \times n$ and $n \times n$ matrices over the real field $\mathbb{R}$, respectively, and we let the symbol $\mathbb{R}_{n \times n}^{S}$ represent all $n \times n$ symmetric matrices of real numbers. In addition, we use the symbol $\mathbb{R}_{n}^{\geq}$ to represent the cone of all symmetric nonnegative-definite matrices in $\mathbb{R}_{n \times n}$ and the symbol $\mathbb{R}_{n}^{>}$ to represent the interior of the cone of all symmetric positive-definite matrices in $\mathbb{R}_{n \times n}$. We also use $\mathcal{C}(\mathbf{A})$ to represent the column space of $\mathbf{A} \in \mathbb{R}_{m \times n}$.

We let $\mathbf{I}_n \in \mathbb{R}_{n \times n}$ signify the $n \times n$ identity matrix. For $\mathbf{A} \in \mathbb{R}_{m \times n}$, we define the transpose of $\mathbf{A}$ and the Moore-Penrose pseudoinverse of $\mathbf{A}$ by $\mathbf{A}^T$ and $\mathbf{A}^+$, respectively, and we note that $\mathbf{A}\mathbf{A}^+$ is the orthogonal projection onto $\mathcal{C}(\mathbf{A})$. We denote the trace of a matrix $\mathbf{A} \in \mathbb{R}_{n \times n}$ by $\text{tr}(\mathbf{A})$ and the rank of $\mathbf{A} \in \mathbb{R}_{m \times n}$ is denoted by rank($\mathbf{A}$). Also, we use $SVD(\mathbf{A})$ to represent the singular value decomposition of $\mathbf{A}$.

In addition, we use the notation $MN_{pn}(\mathbf{M}, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$ to denote a matrix-normal distribution where $\mathbf{M} \in \mathbb{R}_{p \times n}$ is the mean matrix and $\boldsymbol{\Sigma} \in \mathbb{R}_{p}^{>}$. Also, let $\mathbf{X}_i :=$

$\left[ \mathbf{x}_{i1} \vdots \mathbf{x}_{i2} \vdots \cdots \vdots \mathbf{x}_{in_i} \right] \in \mathbb{R}_{p \times n_i}$, represent a data matrix sampled from the $i^{th}$ population so that $\mathbf{X}_i \sim MN_{pn}(\mathbf{M}_i, \mathbf{I}_n \otimes \boldsymbol{\Sigma}_i)$ with $\mathbf{M}_i \in \mathbb{R}_{p \times n}$, $\boldsymbol{\Sigma}_i \in \mathbb{R}_p^>$, and $\otimes$ is the Kronecker product. Therefore, $\mathbf{x}_{ij} \sim N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ are independent $p$-dimensional multivariate normal random vectors for $i = 1, 2$, and $j = 1, 2, \ldots, n_i$. We also use the horizontally concatenated matrix $\mathbf{X} := \left[ \mathbf{X}_1 \vdots \mathbf{X}_2 \right]$ to represent the complete data matrix. Additionally, we use the notation $POW(T_{(R*)})$ to indicate the estimated power of the test $T_{(R*)}$, and $DPOW\left(T_{(R*)}, T_{(*)}\right) := [POW(T_{(R*)}) - POW(T_{(*)})]$, where $T_{(R*)}$ represents a test calculated with post-LDR data and $T_{(*)}$ represents a test calculated from the original data.

For $i = 1, 2$, we use the estimator of $\boldsymbol{\mu}_i$,

$$\overline{\mathbf{x}}_i := \frac{1}{n_i} \mathbf{X}_i \mathbf{j}$$

and the estimator of $\boldsymbol{\Sigma}_i$,

$$\mathbf{S}_i := \frac{1}{n_i - 1} \mathbf{X}_i \left( \mathbf{I}_{n_i} - \frac{1}{n_i} \mathbf{J}_{n_i} \right) \mathbf{X}_i^T,$$

where $\mathbf{J}_{n_i} \in \mathbb{R}_{n_i \times n_i}$ is a matrix of ones and $\mathbf{j} \in \mathbb{R}_{n_i \times 1}$ is a vector of ones. Then, let

$$\mathbf{V}_i := (n_i - 1)\mathbf{S}_i$$

be the gram matrix of the $i^{th}$ sample covariance matrix $\mathbf{S}_i$, $i = 1, 2$, and let

$$\mathbf{S} := \frac{\mathbf{V}_1 + \mathbf{V}_2}{n_1 + n_2 - 2}$$

be the pooled sample covariance matrix.

## 3.3 Four Hypothesis Tests for the Homogeneity of Two Covariance Matrices in the High-Dimensional Setting

Using a Monte Carlo simulation, we first compared the powers of an assortment of *HPCHDS* tests and found four tests that performed relatively well. These tests have been proposed by Srivastava and Yanagihara (2010), Ishii et al. (2016), Chaipitak and Chongcharoen (2013), and Ahmad (2017).

### 3.3.1 Srivastava and Yanagihara (2010)

Srivastava and Yanagihara (2010) have proposed the test statistic

$$Q_{S10} := \sum_{i=1}^{k} \frac{(\hat{a}_{21} - \hat{a}_{22})}{\widehat{\mathrm{Var}(q_{S10})}},$$

where $\hat{a}_{21}$ and $\hat{a}_{22}$ are defined in (B.1), $q_{S10} := (\hat{a}_{21} - \hat{a}_{22})$, and

$$\widehat{\mathrm{Var}(q_{S10})} := \frac{4\hat{a}_2}{p} \left( \frac{p\hat{a}_2 + 2(n_1 - 1)}{(n_1 - 1)^2} + \frac{p\hat{a}_2 + 2(n_2 - 1)}{(n_2 - 1)^2} \right),$$

where $\hat{a}_2$ is defined in (B.2). We refer to the *HPCHDS* test performed using the test statistic $Q_A$ by $T_A$.

Srivastava and Yanagihara (2010) have proven that $\widehat{\mathrm{Var}(q_{S10})} \xrightarrow{P} \mathrm{Var}(q_{S10})$ and, assuming $\mathrm{H}_0$ is true, have also shown that $Q_{S10} \stackrel{\cdot}{\sim} N(0,1)$ as $(p, n_1, n_2) \to \infty$. The power of $T_{S10}$ was first contrasted to powers of other *HPCHDS* tests in Srivastava and Yanagihara (2010).

### 3.3.2 Ishii et al. (2016)

Ishii et al. (2016) have proposed an *HPCHDS* test using ratios of the largest eigenvalues and the corresponding eigenvectors of the two sample covariance matrices $\mathbf{S}_i, i = 1, 2$. Ishii et al. (2016) have proposed the *HPCHDS* test statistic

$$Q_I := \widetilde{\lambda}_* \widetilde{h}_* \widetilde{\gamma}_*, \tag{3.2}$$

where

$$\widetilde{\lambda}_* := \frac{\max\left(\widetilde{\lambda}_{11}, \widetilde{\lambda}_{21}\right)}{\min\left(\widetilde{\lambda}_{11}, \widetilde{\lambda}_{21}\right)}$$

is the ratio of the larger of the two noise-reduced eigenvalues to the smaller of the two noise-reduced eigenvalues and

$$\widetilde{\lambda}_{i1} := \widehat{\lambda}_{i1} - \frac{\mathrm{tr}(\mathbf{S}_i) - \widehat{\lambda}_{i1}}{n_i - 2}, \quad i = 1, 2,$$

are the noise-reduced eigenvalues of $\mathbf{S}_i, i = 1, 2$, respectively. The term $\widetilde{h}_*$ is the ratio of the noise-reduced first eigenvectors,

$$\widetilde{h}_* := \max\left( \left|\widetilde{\mathbf{h}}_1^T \widetilde{\mathbf{h}}_2\right|, \left|\widetilde{\mathbf{h}}_1^T \widetilde{\mathbf{h}}_2\right|^{-1} \right),$$

26

where

$$\widetilde{\mathbf{h}}_{i1} := \left\{(n-1)\widetilde{\lambda}_{i1}\right\}^{-1/2} \left(\mathbf{X}_i - \overline{\mathbf{X}}_i\right) \hat{\mathbf{u}}_{i1},$$

are the first noise-reduced principal component direction vector for group $i$ and $\hat{\mathbf{u}}_{i1}$ is the first unit eigenvector of $i, i = 1, 2$. The final component of the test (3.2) is

$$\widetilde{\gamma}_* := \max\left(\frac{\widetilde{\kappa}_1}{\widetilde{\kappa}_2}, \frac{\widetilde{\kappa}_2}{\widetilde{\kappa}_1}\right),$$

where

$$\widetilde{\kappa}_i := \operatorname{tr}\left(\mathbf{S}_i\right) - \widetilde{\lambda}_{i1}.$$

Ishii et al. (2016) have shown that under $H_0$, $Q_I \overset{.}{\sim} F_{n_1-1, n_2-1}$ as $p \to \infty$. We denote the *HPCHDS* test conducted using the test statistic $Q_I$ by $T_I$.

The test (3.2) compares only the information contained in the first principal component and the loading of the covariance matrices minus some noise. Also, Ishii et al. (2016) did not compare the power of their *HPCHDS* test with the powers of competing tests.

### 3.3.3 *Chaipitak and Chongcharoen* (2013)

Chaipitak and Chongcharoen (2013) have developed a test based on an estimator of the ratio $\left[\operatorname{tr}\left(\mathbf{\Sigma}_i^2\right) / \operatorname{tr}\left(\mathbf{\Sigma}_j^2\right)\right]$. We write their test statistic as

$$Q_C := \frac{\hat{a}_{21}/\hat{a}_{22} - 1}{\sqrt{\widehat{\operatorname{Var}(q_C)}}},$$

where $\hat{a}_{21}$ and $\hat{a}_{22}$ are defined in (B.1), $q_C := \hat{a}_{21}/\hat{a}_{22} - 1$, and

$$\widehat{\operatorname{Var}(q_C)} = 4\left\{\frac{2\hat{a}_4^*}{p\hat{a}_2^2}\left(\frac{1}{n_1 - 1} + \frac{1}{n_2 - 1}\right) + \left(\frac{1}{(n_1 - 1)^2} + \frac{1}{(n_2 - 1)^2}\right)\right\},$$

where $\hat{a}_2$ is given in (B.2) and $\hat{a}_4^*$ is defined in (B.3). We refer to the *HPCHDS* test performed using the test statistic $Q_C$ by $T_C$.

In addition, Chaipitak and Chongcharoen (2013) have shown that $\widehat{\operatorname{Var}(q_C)} \overset{P}{\to}$ $\operatorname{Var}(q_C)$, assuming $H_0$ is true, and that $Q_C \overset{.}{\sim} N(0, 1)$ under $H_0$ as $(p, n_1, n_2) \to$

27

$\infty$. Also, they have contrasted the power of $T_C$ with the powers of *HPCHDS* tests from Schott (2007), and Srivastava and Yanagihara (2010) on four covariance matrix structures for the two-covariance-matrix *HPCHDS* case.

### *3.3.4   Ahmad* (2017)

Ahmad (2017) has proposed the *HPCHDS* test statistic

$$Q_A := \frac{E_1 + E_2 - 2E_{12}}{\sqrt{\widehat{\mathrm{Var}(q_A)}}},$$

where $q_A := E_1 + E_2 - 2E_{12}$, $E_{12} := \mathrm{tr}\,(\boldsymbol{S}_1 \boldsymbol{S}_2)$,

$$E_i := \frac{(n_i - 1)}{n_i\,(n_i - 2)\,(n_i - 3)} \left\{ (n_i - 1)\,(n_i - 2)\,\mathrm{tr}\,(\mathbf{S}_i^2) + \left[\,\mathrm{tr}\,(\mathbf{S}_i)\,\right]^2 \right.$$
$$\left. - \frac{n_i}{(n_i - 1)} \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)^T (\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)(\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)^T (\mathbf{x}_{ij} - \overline{\mathbf{x}}_i) \right\},$$

$i = 1, 2$, and

$$\widehat{\mathrm{Var}(q_A)} = 4\left[\,\mathrm{tr}\,(\mathbf{S}^2)\,\right]^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)^2.$$

We refer to the *HPCHDS* test conducted using the test statistic $Q_A$ by $T_A$.

Ahmad (2017) has proven that $\widehat{\mathrm{Var}(q_A)} \xrightarrow{P} \mathrm{Var}\,(q_A)$, assuming $\mathrm{H}_0$ is true, and has shown that $Q_A \stackrel{.}{\sim} N(0, 1)$ under $\mathrm{H}_0$ as $(p, n_1, n_2) \to \infty$. However, he did not contrast the power of $T_A$ with the power of other *HPCHDS* tests.

### *3.4   Linear Dimension Reduction for Improved Power of Tests of Covariance-Matrix Homogeneity for Two Population Covariance Matrices Under a High-Dimensional Setting*

Below, we first prove a separability theorem for a new *LDR* method for two population covariance matrices. Our new *LDR* model matrix for two high-dimensional covariance matrices is based on a property proposed by Peters et al. (1978), who defined the concept of a linear sufficient matrix for reducing the dimension of two multivariate normal population covariance matrices. We also propose a new *LDR* method for two sample covariance matrices. Then, we employ our *LDR* matrix for two

sample covariance matrices to increase the power of two population *HPCHDS* tests. Using the singular value decomposition (*SVD*) of Eckart and Young (1936), we derive the *LDR* matrix for two sample covariance matrices that allows us to retain most of the distinguishing information in $[\mathbf{S}_2 - \mathbf{S}_1]$ and, therefore, distinguishing information for $[\mathbf{\Sigma}_2 - \mathbf{\Sigma}_1]$.

To prove our *LDR* model theorem, we present the symmetrized Kullback-Leibler separability measure for two positive-definite matrices as

$$SKL(\mathbf{\Sigma}_1, \mathbf{\Sigma}_2) := \left[\log \frac{|\mathbf{\Sigma}_2|}{|\mathbf{\Sigma}_1|} + \text{tr}(\mathbf{\Sigma}_2^{-1}\mathbf{\Sigma}_1)\right] + \left[\log \frac{|\mathbf{\Sigma}_1|}{|\mathbf{\Sigma}_2|} + \text{tr}(\mathbf{\Sigma}_1^{-1}\mathbf{\Sigma}_2)\right]. \qquad (3.3)$$

We now derive our new *LDR* model matrix for two population covariance matrices and demonstrate its ability to preserve information concerning $[\mathbf{\Sigma}_1 - \mathbf{\Sigma}_2]$ in a reduced dimension.

**Theorem.** *Suppose we have two multivariate normal populations with covariance matrices* $\mathbf{\Sigma}_i \in \mathbb{R}_p^>$, $i = 1, 2$, *and let*

$$\mathbf{H} := [\mathbf{\Sigma}_2 - \mathbf{\Sigma}_1]. \qquad (3.4)$$

*Next, let* $SVD(\mathbf{H}) = \mathbf{F}\mathbf{\Lambda}\mathbf{G} \in \mathbb{R}_{p \times p}$, *where* $\mathbf{F} \in \mathbb{R}_{p \times r}$ *and* $\text{rank}(\mathbf{F}) = \text{rank}(\mathbf{H}) = r < p$. *Also, let the separability measure* $SKL(\mathbf{\Sigma}_1, \mathbf{\Sigma}_2)$ *be defined in* (3.3). *Then,*

$$SKL(\mathbf{\Sigma}_1, \mathbf{\Sigma}_2) = SKL(\mathbf{F}^+\mathbf{\Sigma}_1\mathbf{F}^{+T}, \mathbf{F}^+\mathbf{\Sigma}_2\mathbf{F}^{+T}).$$

*Proof.* The proof follows from Lemmas B.2.2 and B.2.3 in Appendix B.

Next, let

$$\widehat{\mathbf{H}} := [\mathbf{S}_2 - \mathbf{S}_1], \qquad (3.5)$$

where $\text{rank}(\mathbf{S}_i) < p, i = 1, 2$, be an estimator of $\mathbf{H}$, given in (3.4). Because $\text{rank}(\widehat{\mathbf{H}}) = n_1 + n_2 - \dim\left[\mathcal{C}(\mathbf{S}_1) \cap \mathcal{C}(\mathbf{S}_2)\right]$ and $\mathbf{S}_i \in \mathbb{R}_p^\geq$, $i = 1, 2$, one cannot directly apply the theorem to obtain the *LDR* matrix $\mathbf{F}^+ \in \mathbb{R}_{q \times p}$ when we desire a $q < \text{rank}(\widehat{\mathbf{H}})$. Moreover, we often wish to obtain a low-dimensional representation of the original

data with dimension $q$, where $1 \leq q \ll \text{rank}(\widehat{\mathbf{H}}) \ll p$. However, motivated by the theorem, we can construct an *LDR* matrix that preserves much of the original $p$-dimensional information in the original data for the distinguishing aspects of the two covariance matrices $\mathbf{S}_i \in \mathbb{R}_p^{\geq}$, $i = 1, 2$, by using the *SVD*.

Let $SVD(\widehat{\mathbf{H}}) = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$ and let $\mathbf{U}^{(q)}$ denote the first $q$ eigenvectors of $\widehat{\mathbf{H}}$ corresponding to the $q < p$ largest singular values of $\widehat{\mathbf{H}}$. Next, let $\widehat{\mathbf{F}}_{p \times q} := \mathbf{U}^{(q)}$. Then, because $\left[\widehat{\mathbf{F}}_{p \times q}\right]^T \in \mathbb{R}_{q \times p}$ and $\left[\widehat{\mathbf{F}}_{p \times q}\right]^+ \in \mathbb{R}_{q \times p}$ span the same subspace, we use $\left[\widehat{\mathbf{F}}_{p \times q}\right]^T$ as the *LDR* matrix for reducing the original feature dimension to the reduced dimension $q$, where $1 \leq q < \text{rank}(\widehat{\mathbf{H}})$, while preserving much of the separability information between the estimated covariance matrices $\mathbf{S}_i \in \mathbb{R}_p^{\geq}, i = 1, 2$. Therefore, we believe that mapping the high-dimensional data matrix onto $\mathcal{C}([\widehat{\mathbf{F}}_{p \times q}]^T)$ will enhance our ability to detect differences in the individual population covariance matrices $\mathbf{\Sigma}_i \in \mathbb{R}_p^{>}, i = 1, 2$, because of the decreased number of parameters that must be estimated.

### 3.5  Monte Carlo Simulation Design

#### 3.5.1  Simulation Covariance Structures

The five covariance matrix structures used in our Monte Carlo simulations were selected from the *HPCHDS* literature. We have compared test powers across five different covariance matrix structures. The group sample sizes used in the simulations were $n_1 = n_2$.

First, we used the constant-times-identity covariance matrix structure. For our simulation, the parameters for the null and alternative hypotheses are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{I}_p$$

and

$$H_A : \mathbf{\Sigma}_1 = \mathbf{I}_p \text{ and } \mathbf{\Sigma}_2 = \sigma^2 \mathbf{I}_p.$$

Second, we use the compound-symmetric covariance matrix class. For our simulation, the null and alternative hypothesis parameters are

$$H_0 : \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \sigma_1^2 \mathbf{I}_p + \sigma_2^2 \mathbf{J}_p$$

and

$$H_A : \boldsymbol{\Sigma}_1 = \sigma_1^2 \mathbf{I}_p + \sigma_2^2 \mathbf{J}_p \text{ and } \boldsymbol{\Sigma}_2 = \sigma_{1A}^2 \mathbf{I}_p + \sigma_{2A}^2 \mathbf{J}_p,$$

respectively, where $\mathbf{J}_p \in \mathbb{R}_{p \times p}$ is a matrix of ones.

Third, we use the autoregressive covariance matrix structure. For our simulation, the parameters for the null and alternative hypotheses are

$$H_0 : \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \mathbf{U}_0$$

and

$$H_A : \boldsymbol{\Sigma}_1 = \mathbf{U}_0 \text{ and } \boldsymbol{\Sigma}_2 = \mathbf{U}_1,$$

where $\mathbf{U}_0 = \sigma_{ij} = 0.1^{|i-j|}$, $\mathbf{U}_1 = \sigma_{ij} = 0.3^{|i-j|}$, and $1 \leq i, j \leq p$.

Fourth, we use the heterogeneous autoregressive covariance matrix structure. For our simulation, we use the heterogeneous autoregressive covariance matrix structures similar to those in Srivastava et al. (2014), which we create as follows. First, let $\sigma_l := 1 + (-1)^{l+1} Q_l / 2$, where $Q_l \sim Unif(0,1)$ and $l = 1, 2, \ldots, p$. Then, the parameters for the null and alternative hypotheses are

$$H_0 : \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \sigma_i \sigma_j 0.1^{|i-j|^{\frac{1}{10}}}$$

and

$$H_A : \boldsymbol{\Sigma}_1 = \sigma_i \sigma_j 0.1^{|i-j|^{\frac{1}{10}}} \text{ and } \boldsymbol{\Sigma}_2 = \sigma_i \sigma_j 0.3^{|i-j|^{\frac{1}{10}}},$$

respectively, where $1 \leq i, j \leq p$.

Last, we examine an unstructured covariance matrix structure that has no structure or pattern, which we model as

$$\mathbf{U}_2 = \sigma_{ij} = \sigma_{ji} := \begin{cases} (-1)^{i+j} \left( \frac{0.10}{j} \right), & i < j \\ 1, & i = j \end{cases}$$

and

$$\mathbf{U}_3 = \sigma_{ij} = \sigma_{ji} := \begin{cases} (-1)^{i+j} \left( \frac{0.05}{j} \right), & i < j \\ 1, & i = j. \end{cases}$$

Then, for testing *HPCHDS* with unstructured covariance matrices, we note that the parameters for the null and alternative hypotheses are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{U}_2$$

and

$$H_A : \mathbf{\Sigma}_1 = \mathbf{U}_2 \text{ and } \mathbf{\Sigma}_2 = \mathbf{U}_3.$$

### 3.5.2   *Monte Carlo Power Simulation Description*

We now describe the simulation design we use to contrast the powers for the two population post-*LDR HPCHDS* tests $T_{RA}$, $T_{RCB}$, $T_{RSc}$, and $T_{RS10}$ and the two-population, no-*LDR HPCHDS* tests $T_A$, $T_{CB}$, $T_{Sc}$, and $T_{S10}$. Using the statistical programming software R, we generated 10,000 independent multivariate normal vectors from $N_p(\mathbf{0}, \mathbf{\Sigma})$, where $\mathbf{\Sigma} \in \mathbb{R}_p^{>}$ is the common covariance matrix under $H_0$ for group $i = 1, 2$. For each sample dataset $\mathbf{X}^{(j)}$, $j = 1, \ldots, 10,000$, we estimated $\mathbf{M}$ by $\widehat{\mathbf{M}}^{(j)}$, as in (3.5), and calculated $SVD(\widehat{\mathbf{M}}^{(j)})$ to extract our *LDR* matrix $\left[ \widehat{\mathbf{F}}_{p \times q}^{(j)} \right]^T$. Next, we reduced the dimension of the sample data by mapping the full-dimensional data matrix $\mathbf{X}^{(j)}$ onto $\mathcal{C}\left( \left[ \widehat{\mathbf{F}}_{p \times q}^{(j)} \right]^T \right)$ so that $\mathbf{Y}^{(j)} = \left[ \widehat{\mathbf{F}}_{p \times q}^{(j)} \right]^T \mathbf{X}^{(j)}$, where $q$ is the targeted reduced-data dimension. We then calculated the empirical test value $T_{(*),j}$ from the

32

reduced data vectors $\mathbf{y}_{ij}$, $i = 1, 2$, and $j = 1, 2, \ldots, n_i$, and determined the simulated critical values ($SCVs$) for $T_{RA}$, $T_{RC}$, and $T_{RS10}$ using

$$SCV_{1-\alpha/2} := \inf \left\{ x \in \mathbb{R} : 1 - \alpha/2 \leq \widehat{F}_{T_{(R*)}}(x) \right\} \tag{3.6}$$

and

$$SCV_{\alpha/2} := \sup \left\{ x \in \mathbb{R} : \alpha/2 \geq \widehat{F}_{T_{(R*)}}(x) \right\}, \tag{3.7}$$

where $\widehat{F}_{T_{(R*)}}(x)$ is the empirical distribution function of the test $T_{(R*)}$. For $T_{RI}$ the $SCV$ was calculated as

$$SCV_{1-\alpha} := \inf \left\{ x \in \mathbb{R} : 1 - \alpha \leq \widehat{F}_{T_{(RI)}}(x) \right\}. \tag{3.8}$$

We used the significance level $\alpha = 0.05$ as our decision criterion to perform each individual $HPCHDS$ hypothesis test. Additionally, we determined the critical values for the tests $T_A$, $T_{CB}$, $T_{Sc}$, and $T_{S10}$ via the methods described in (3.6), (3.7), and (3.8), but using the original-dimension datasets $\mathbf{X}^{(j)}$.

To determine the power for each test, given $q$ and $n_i, i = 1, 2$, we generated 10,000 independent multivariate normal random vectors from each of the $N_p(\mathbf{0}, \boldsymbol{\Sigma}_i)$ populations, where $\boldsymbol{\Sigma}_i \in \mathbb{R}_p^{>}$, $i = 1, 2$, are the covariance matrices assuming $\mathrm{H_A}$. For each complete sample dataset $\mathbf{X}^{(j)} \in \mathbb{R}_{p \times N}$, where $N := \sum_{i=1}^{2} n_i$, we extracted the $LDR$ matrix $\left[ \widehat{\mathbf{F}}_{p \times q}^{(j)} \right]^T$ and then reduced the data dimension by mapping the full-dimensional data matrix $\mathbf{X}^{(j)}$ onto $\mathcal{C}\left( [\widehat{\mathbf{F}}_{p \times q}^{(j)}]^T \right)$ so that $\mathbf{Y}^{(j)} = \left[ \widehat{\mathbf{F}}_{p \times q}^{(j)} \right]^T \mathbf{X}^{(j)}$. We then calculated the test $T_{(R*),j}$ for each $j$, where $1 \leq j \leq 10,000$, using each reduced-data matrix $\mathbf{Y}^{(j)}$ and estimated the power by

$$POW(T_{(R*)}) = \frac{\sum_{j=1}^{10,000} I[T_{(R*),j} \in RR(T_{(R*)}]}{10,000}, \tag{3.9}$$

where $RR(T_{(R*)})$ is the rejection region for the test $T_{(R*)}$ and $I[\cdot]$ is the indicator function. We applied a method similar to that given in (3.9) but using the original unreduced datasets $\mathbf{X}^{(j)}$ to estimate $POW(T_A)$, $POW(T_C)$, $POW(T_I)$, and $POW(T_{S10})$.

We performed power simulations for the common sample sizes $n_1 = n_2 \in$ $5, 10, 15, 20$ and the full-data dimensions $p \in 20, 40, 80, 160$ and summarized the results for $POW(T_{RA})$, $POW(T_{RC})$, $POW(T_{RI})$, $POW(T_{RS10})$ and for $POW(T_A)$, $POW(T_C)$, $POW(T_I)$, and $POW(T_{S10})$ in 3.1. Also, in the five figures shown below, we displayed curves for $POW(T_{RA})$, $POW(T_{RC})$, $POW(T_{RI})$, and $POW(T_{RS10})$ plotted against the reduced dimension $q$ for $q \in \{1, 2, ..., 159\}$. We also provided power-difference curves plotted versus the reduced dimension $q$ for $q \in \{1, 2, ..., 159\}$ $DPOW(T_{RA}, T_A)$, $DPOW(T_{RC}, T_C)$, $DPOW(T_{RI}, T_I)$, and $DPOW(T_{RS10}, T_{S10})$ , where $DPOW\left(T_{(R*)}, T_{(*)}\right) := [POW(T_{(R*)}) - POW(T_{(*)})]$ for the tests $T_{(R*)}$ and $T_{(*)}$. The original-data dimension was $p = 160$ and the sample size was $n_1 = n_2 = 10$. The simulations were run in parallel using `R` and the `covTestR` package.

### 3.6    Monte Carlo Power Simulation Results

In this section we present our simulated power-contrast results for the $HPCHDS$ tests $T_A$, $T_C$, $T_I$, and $T_{S10}$, calculated with no $LDR$, and the $HPCHDS$ tests $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$, calculated after we applied $LDR$ to the sample observations from each group. We display the power-simulation results in Table 3.1 below and in Figures 3.1– Figures 3.5. We used generalized linear models with b-splines to fit the power curves and the power-difference curves.

### 3.6.1    A Simulation-Summary Table for POW(T_{RA}), POW(T_{RC}), POW(T_{RI}), and POW(T_{RS10}) for Two Autoregressive Covariance Matrices

Table 3.1 shows the post-$LDR$ powers $POW(T_{RA})$, $POW(T_{RC})$, $POW(T_{RI})$, and $POW(T_{RS10})$ for two autoregressive covariance matrix structures with parameters $\mathbf{\Sigma}_1 = 0.1^{|i-j|}$ and $\mathbf{\Sigma}_2 = 0.3^{|i-j|}$, $q \in \{5, 10, 15, 20\}$ when $n_i = 5$ and for $q \in \{10, 20, 30, 40\}$ when $n_i = 10$. We also report the no-$LDR$ powers $POW(T_A)$, $POW(T_C)$, $POW(T_I)$, and $POW(T_{S10})$ for configuration scenarios with $p \in \{80, 160\}$ and $n_i = \{5, 10\}$, $i = 1, 2$.

Table 3.1: A table contrasting POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) with POW($T_A$), POW($T_C$), POW($T_I$), and POW($T_{S10}$) when testing for *HPCHDS* for two autoregressive population covariance matrices with parameters $\mathbf{\Sigma}_1 = 0.1^{|i-j|}$ and $\mathbf{\Sigma}_2 = 0.3^{|i-j|}$.

| $p$ | $n_1 = n_2$ | $q$ | POW($T_{RA}$) | POW($T_{RC}$) | POW($T_{RI}$) | POW($T_{RS10}$) |
|---|---|---|---|---|---|---|
| 80 | 5 | 5 | 0.00 | 1.00 | 0.02 | 0.00 |
| | | 10 | 0.00 | 0.93 | 0.00 | 0.00 |
| | | 15 | 0.01 | 0.78 | 0.35 | 0.00 |
| | | 20 | 0.01 | 0.66 | 0.29 | 0.00 |
| | 10 | 10 | 0.00 | 1.00 | 0.00 | 1.00 |
| | | 20 | 0.00 | 0.98 | 0.01 | 0.36 |
| | | 30 | 0.00 | 0.83 | 0.12 | 0.16 |
| | | 40 | 0.00 | 0.55 | 0.10 | 0.05 |
| 160 | 5 | 5 | 0.00 | 1.00 | 0.03 | 0.99 |
| | | 10 | 0.00 | 0.99 | 0.00 | 0.00 |
| | | 15 | 0.00 | 0.97 | 0.40 | 0.00 |
| | | 20 | 0.00 | 0.94 | 0.40 | 0.00 |
| | 10 | 10 | 0.00 | 1.00 | 0.00 | 0.9 |
| | | 20 | 0.00 | 1.00 | 0.01 | 0.52 |
| | | 30 | 0.00 | 1.00 | 0.19 | 0.20 |
| | | 40 | 0.00 | 0.99 | 0.25 | 0.13 |

| $p$ | $n_1 = n_2$ | $p$ | POW($T_A$) | POW($T_C$) | POW($T_I$) | POW($T_{S10}$) |
|---|---|---|---|---|---|---|
| 80 | 5 | 80 | 0.10 | 0.03 | 0.04 | 0.06 |
| | 10 | 80 | 0.16 | 0.02 | 0.06 | 0.05 |
| 160 | 5 | 160 | 0.10 | 0.02 | 0.05 | 0.06 |
| | 10 | 160 | 0.14 | 0.02 | 0.05 | 0.06 |

As shown in Table 3.1, $POW(T_{RC})$ was increased significantly over $POW(T_C)$ because of the application of our *LDR* method prior to the *HPCHDS* hypothesis test. However, $POW(T_{RI})$ was less than $POW(T_I)$ and $POW(T_{RA})$ was considerably less than $POW(T_A)$ because of the use of *LDR* on the full-dimensional data prior to the hypothesis test. We also found circumstances in which the use of *LDR* increased $POW(T_{RS10})$ over $POW(T_{S10})$. Thus, we observed that a reduction in the data dimension prior to calculating the empirical test value by using our *LDR* method can yield a considerable increase in test power or can cause a relatively large decrease in power, depending on the hypothesis test used.

### 3.6.2  Power Curves and Power-Difference Curves for $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ for Two Constant-Times-Identity Covariance Structures

In Figure 3.1, we present plots for the post-*LDR* power curves for $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ and for the difference in the post-*LDR* power curves $DPOW(T_{R*}, T_*)$. The powers for the *HPCHDS* tests when $p = 160$ were $POW(T_A) = 0.86$, $POW(T_C) = 0.13$, $POW(T_I) = 0.11$, and $POW(T_{S10}) = 0.01$. The *HPCHDS* hypothesis tests were performed for two constant-times-identity covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{I}_p$ and $\mathbf{\Sigma}_2 = (1.5)\mathbf{I}_p$. Comparing the two groups of plots, we see similar trends in power and power-difference curves. Three of the power curves have maximum power at or near the common sample size $n_i = q = 10$, $i = 1, 2$.

We see that $POW(T_{RC})$ was relatively small for $q > 75$ but increased as $q$ decreased, attaining a maximum power of 1.0 at $q = 20$. Also, $POW(T_{RS10})$ yielded discernible power for $q < 25$ and peaked at $q = 10$ with $POW(T_{RS10}) = 0.95$. In addition, $POW(T_{RI}) = 0.40$ was the maximum power for $T_{RI}$, which occurred at $q = 19$, while $POW(T_{RA})$ decreased as $q$ decreased. However, $DPOW(T_{RA}, T_A) < 0$ for all $q$, where $q \in \{1, 2, ..., 159\}$. Also, $DPOW(T_{RS10}, T_{S10})$ and $DPOW(T_{RC}, T_C)$ peaked to nearly 1.0 and 0.90 at $q = 10$, respectively. Thus, applying our *LDR* method for two covariance matrices prior to calculating the *HPCHDS* tests $T_{RS10}$ and $T_{RC}$ yielded a significant increase in the maximum power. The $DPOW(T_{RI}, T_I)$ curve peaked at $q = 13$ with a power increase of 0.27. Also, $POW(T_A)$ had the largest power among the reported no-*LDR HPCHDS* tests. The fact that $DPOW(T_{RA}, T_A) < 0$ for $q$, where $q \in \{1, 2, ..., 159\}$, as seen in the power-difference plot, suggested that we lost a substantial amount of information concerning the difference between covariance matrices because of the use of *LDR*.

Figure 3.1: Reduced-dimension power curves and power-difference curves for POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) for two constant-times-identity covariance matrices with parameters $\boldsymbol{\Sigma}_1 = \mathbf{I}_p$ and $\boldsymbol{\Sigma}_2 = (1.5)\mathbf{I}_p$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 10, i = 1, 2, 3$.

### 3.6.3 Power Curves and Power-Difference Curves for $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ for Two Compound-Symmetric Covariance Matrices

In Figure 3.2, we present plots for the power curves constructed after application of our *LDR* method for the tests $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ and plots for the difference in the post-*LDR* power curves from the original-dimension powers. The *HPCHDS* hypothesis tests were performed for two compound symmetric covariance matrices with parameters $\boldsymbol{\Sigma}_1 = (0.99)\mathbf{I}_p + (0.01)\mathbf{J}_p$ and $\boldsymbol{\Sigma}_2 = (0.95)\mathbf{I}_p + (0.05)\mathbf{J}_p$. The no-*LDR* powers were $POW(T_A) = 0.30$, $POW(T_C) = 0.02$, $POW(T_I) = 0.05$, and

$POW(T_{S10}) = 0.02$ at $p = 160$. Three power curves have maximums at or near the common sample size $n_i = q = 10, i = 1, 2$, which is denoted by the vertical line in Figure 3.2.



Figure 3.2: Reduced-dimension power curves and power-difference curves for $POW(T_{RA})$, $POW(T_{RC})$, $POW(T_{RI})$, and $POW(T_{RS10})$ for two compound symmetric population covariance matrices with parameters $\mathbf{\Sigma}_1 = (0.99)\mathbf{I}_p + (0.01)\mathbf{J}_p$ and $\mathbf{\Sigma}_2 = (0.95)\mathbf{I}_p + (0.05)\mathbf{J}_p$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 10, i = 1, 2, 3$.

The curve for $POW(T_{RS10})$ peaked to 0.87 at $q = 10$, which is also the common sample size. Additionally, $POW(T_{RC}) \approx 1.0$ for $q \in \{10, 11, ..., 35\}$ and consistently yielded the largest power when $q \leq 112$. In addition, $POW(T_{RI})$ was maximized at $q = 27$ with $POW(T_{RI}) = 0.37$, which was relatively small. Again, $POW(T_A)$ was

the largest power among the four no-$LDR$ $HPCHDS$ tests. However, $POW\left(T_{RA}\right)$ decreased as $q$ decreased and never yielded a power value greater than $POW\left(T_A\right)$.

The curve for $POW\left(T_{RC}\right)$ considerably increased as $q$ was reduced and peaked at $q = 20$ with $DPOW\left(T_{RC}, T_C\right) \approx 1.0$. Also, $DPOW\left(T_{RS10}, T_{S10}\right)$ peaked in power increase at $q = 10$ with a difference of 0.87. We see that $POW\left(T_{RI}\right)$ peaked at 0.37 when $q = 19$. In the power-difference plot, $DPOW\left(T_{RA}, T_A\right) < 0$ for $q \in \{1, 2, ..., 159\}$, and, therefore, we again found that the application of $LDR$ to the original data decreased the power for $T_A$ for all considered values of $q$.

### 3.6.4  Power Curves and Power-Difference Curves for $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ for Two Autoregressive Covariance Structures

In Figure 3.3, we present plots for the post-$LDR$ power curves for tests $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ and plots for the difference in the reduced-dimension power curves from the original-dimension powers $POW(T_A) = 0.14$, $POW(T_C) = 0.02$, $POW(T_I) = 0.05$, and $POW(T_{S10}) = 0.06$ for $p = 160$. The hypothesis tests for $HPCHDS$ were performed for two autoregressive covariance matrices with parameters $\boldsymbol{\Sigma}_1 = 0.1^{|i-j|}$ and $\boldsymbol{\Sigma}_2 = 0.3^{|i-j|}$. We first see that two power curves have maximums at the common sample size denoted by the solid vertical line. The maximum for $POW(T_{RS10})$ occurred at $q = 16$. Also, for this autoregressive covariance matrix structure, $POW\left(T_{RC}\right) = 1$ was the maximum power for each of the reduced data dimensions $q \in \{9, 10, ..., 30\}$ and was the most powerful post-$LDR$ $HPCHDS$ test examined here.

Also, $POW\left(T_{RS10}\right)$ had almost no discernible power increase for $q > 50$. However, $POW\left(T_{RS10}\right)$ was maximized at $q = 10$ where $POW\left(T_{RS10}\right) \approx 0.86$. In addition, $POW\left(T_{RI}\right) = 0.40$ was the maximum power at $q = 19$, and $POW\left(T_{RC}\right)$ produced the largest power among all post-LDR tests when $q$ was reduced to $q \in \{9, 10, ..., 30\}$. Once again, the plot of $POW(T_{RA})$ monotonically decreased as $q$ was reduced.

Figure 3.3: Reduced-dimension power curves and power-difference curves for $POW(T_{RA})$, $POW(T_{RC})$, $POW(T_{RI})$, and $POW(T_{RS10})$ for two autoregressive covariance matrices with parameters $\mathbf{\Sigma}_1 = 0.1^{|i-j|}$ and $\mathbf{\Sigma}_2 = 0.3^{|i-j|}$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 10, i = 1, 2, 3$.

We also see that $DPOW(T_{RS10}, T_{S10})$ had a maximum power increase of 0.88 at $q = 11$ and that $DPOW(T_{RC}, T_C) \approx 1.0$ for $q < 50$. IN addition, $POW(T_I)$ was maximized at $q = 25$, where $DPOW(T_{RI}, T_I) \approx 0.38$. Clearly, $DPOW(T_{RC}, T_C)$ had the largest positive power difference for all $q$. Additionally, $DPOW(T_{RA}, T_A) < 0$ for all $q$, where $q \in \{1, 2, ..., 159\}$.

### 3.6.5 Power Curves and Power-Curve Differences for the tests $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ for Two Heterogeneous Autoregressive Covariance Matrices

In Figure 3.4, we present plots for the reduced-dimension power curves for tests $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$. We also show the power difference in the reduced-dimension power curves from the original-dimension powers, which were $POW(T_A) = 0.67$, $POW(T_C) = 0.01$, $POW(T_I) = 0.07$, and $POW(T_{S10}) = 0.52$ at $p = 160$. The hypothesis tests for $HPCHDS$ were performed for two heterogeneous autoregressive covariance matrices. Comparing the power and power-difference plots, we see similar trends in the power and power-difference curves. We also see that two power curves have maximum power at or near the common sample size $n_i = q = 10, i = 1, 2$, which is denoted by the solid vertical red line. For the heterogeneous autoregressive covariance matrix structure, $POW(T_A)$ had the largest power for tests calculated with the unreduced data.

We see that $POW(T_{RC})$ increased as $q$ was reduced and peaked at $q = 10$ with $POW(T_{RC}) = 0.90$ and that $POW(T_{RS10})$ decreased as $q$ was initially reduced but then increased in power for $q \in \{9, 10, ..., 20\}$. Also, $POW(T_{RS10})$ attained its maximum of 0.75 at $q = 9$. Additionally, $POW(T_{RI})$ peaked at 0.40 for $q = 19$. Not surprisingly, we found that $POW(T_A)$ produced the largest power among the considered no-$LDR$ $HPCHDS$ tests.

In the power-difference plots, $DPOW(T_{RC}, T_C)$ attained a maximum power increase of 0.89 at $q = 10$, whereas $DPOW(T_{RA}, T_A) < 0$ for $q \in \{1, 2, ..., 159\}$ once again. Also, $DPOW(T_{RS10}, T_{S10})$ yielded a maximum power increase of 0.24 over the no-$LDR$ power of $T_{S10}$ when $q = 9$, and $DPOW(T_{RI}, T_I) = 0.36$ was the maximum power increase which occurred at $q = 20$.

Figure 3.4: Reduced-dimension power curves and power-difference curves for $POW(T_{RA})$, $POW(T_{RC})$, $POW(T_{RI})$, and $POW(T_{RS10})$ for two heterogeneous autoregressive covariance matrices. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 10, i = 1, 2, 3$.

### 3.6.6 Power Curves and Power-Curve Differences for the tests $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ for Two Unstructured Covariance Matrices

In Figure 3.5 we display the power curves for $POW(T_{RA})$, $POW(T_{RC})$, $POW(T_{RI})$, and $POW(T_{RS10})$ for two unstructured population covariance matrices plotted against $q$ for $q \in \{1, 2, ..., 159\}$. The curves for $POW(T_{RI})$ was hardly discernible for $q > 25$, and the curves for $POW(T_{RA})$ and $POW(T_{RS10})$ were similar. We see that $POW(T_{RA}) = 0.27$ was the maximum power and it occurred at $q = 159$ and $POW(T_{RS10}) \approx 0.14$ was the maximum power that was attained at $q = 159$. Also,

$POW\left(T_{RC}\right)$ had a small peak in power of 0.07 at $q = 7$. From Figure 3.5 we see that the only practically significant increase in power was $DPOW\left(T_{RI}, T_I\right) = 0.38$.



Figure 3.5: Reduced-dimension power curves and power-difference curves for $POW(T_{RA})$, $POW(T_{CR})$, $POW(T_{RI})$, and $POW(T_{RS10})$ for two unstructured co-variance matrices. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 10, i = 1, 2, 3$.

### 3.6.7 Selection of the Reduced Dimension $q$ for LDR

We have shown that $LDR$ can substantially increase the power of certain tests for $HPCHDS$ provided an appropriate $q < p$ is chosen. However, the practitioner testing for $HPCHDS$ cannot feasibly check every possible value of $q$.

A common approach would be to use one of many existing methods to determine the essential rank of $\widehat{\mathbf{M}}$ and, thus, an appropriate $q$. However, we emphasize that an

optimal value of $q$ that yields maximum power depends on the chosen test, the type of population covariance matrices being tested, and the group sample sizes. Therefore, we do not recommend using only a rank-estimation method such as those found in Cook and Forzani (2009) and Rohde and Tsybakov (2011).

In Figures 3.1 – 3.5 we have shown the common sample size $n_i$ with a solid vertical line. For the tests that showed improvement in power, most obtained maximum power at or near the common sample size $n_i = q, i = 1, 2$. Though some tests plateaued in power increase for $q > n_i, i = 1, 2$, and have multiple optimal values of $q$, we recommend reducing the data dimension to $q = \min_{i=1,2} n_i$ when the group sample sizes are approximately equal. This reduced data dimension criteria will not always yield the largest possible increase in power but should supply the researcher with a relatively good increase in power from the original-data dimension.

### 3.7 A Contrast of the Original-Data tests $T_A$, $T_C$, $T_I$, and $T_{S10}$ and the Reduced-Data tests $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ on Real Data

In this section, we examine the efficacy of our proposed $LDR$ technique for the conservation of separability of the covariance matrices on a real high-dimensional dataset from Alon et al. (1999), as curated by Ramey (2016). The Alon dataset contains gene expression levels for 40 tumor and 22 normal colon tissues of 6500 human genes from an Affymetrix oligonucleotide array. Alon et al. (1999) recommended using the 2000 genes with the highest minimal intensity across the 62 samples.

To test for $HPCHDS$ for two population covariance matrices, we performed permutation tests for each test considered in this paper to determine critical values for all tests considered here both with and without first applying $LDR$. We created 1000 permutations of the Alon dataset and the post-$LDR$ dataset by randomly assigning the 62 samples into the two classes. Using the permutation test distributions, we then calculated the permuted critical values for the test statistics, $T_A$, $T_I$, $T_{S10}$, and $T_C$

and for the test statistics, $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ at the $\alpha = 0.05$ level. Table 3.2 presents the results of our permutation tests.

Table 3.2: Comparison and contrast of characteristics of the *HPCHDS* tests $T_A$, $T_C$, $T_I$, and $T_{S10}$ calculated with the original data and the *HPCHDS* tests $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$ calculated after the application of *LDR* to the Alon dataset.

| Full-Dimension test | $T_C$ | $T_I$ | $T_{S10}$ | $T_A$ |
| --- | --- | --- | --- | --- |
| Full-Dimension Lower Crit. Val. | -9.111 | - | -3.625 | -4.237 |
| Full-Dimension Upper Crit. Val. | 8.582 | 3.550 | 3.871 | 4.100 |
| Full-Dimensionla Test Value | 2.872 | 1.982 | 2.619 | 5.046 |
| Full-Dimension p-Value | 0.534 | 0.448 | 0.356 | 0.035 |
| Full-Dimension Test Decision | FTR $H_0$ | FTR $H_0$ | FTR $H_0$ | Reject $H_0$ |

| Reduced-Dimension test | $T_{RC}$ | $T_{RI}$ | $T_{RS10}$ | $T_{RA}$ |
| --- | --- | --- | --- | --- |
| Reduced-Dimension Lower Crit. Val. | -1.711 | - | -6.898 | -9.669 |
| Reduced-Dimension Upper Crit. Val. | 1.643 | 2.551 | 7.051 | 9.523 |
| Reduced-Dimension Test Value | 1.909 | 6.143 | 6.475 | 5.523 |
| Reduced-Dimension p-Value | 0.012 | 0.001 | 0.125 | 0.786 |
| Reduced-Dimension Test Decision | Reject $H_0$ | Reject $H_0$ | FTR $H_0$ | FTR $H_0$ |

After performing *LDR* on the original data, we reduced the data dimension from $p = 2000$ to $q = 22$, which is the group-minimal sample size. The tests $T_{RC}$ and $T_{RI}$ yielded empirical test scores greater than the permutation tests' larger critical values at the $\alpha = 0.05$ level. The reduction of the p-values (*PVs*) for these two *HPCHDS* tests was markedly substantial. More specifically, We determined that the *PV* for $T_C$ calculated with the original data was $PV(T_C) = 0.534$ and the *PV* for the reduced-data test was $PV(T_{RC}) = 0.012$. Also, for the full-dimensional-data, the *PV* for the original-data test $T_I$ was $PV(T_I) = 0.448$, while $PV(T_{RI}) = 0.001$. Again, for the Alon dataset, the differences $[PV(T_{RC}) - PV(T_C)]$ and $[PV(T_{RI}) - PV(T_I)]$ were considerable. Also, $PV(T_{RS10})$ was substantially smaller than $PV(T_{S10})$ but did not change the full-data-dimension hypothesis-test decision.

Finally, for the post-*LDR* test, $T_{RA}$, $PV(T_{RA})$ actually increased considerably over $PV(T_A)$. As a result, the test decision for the no-*LDR* test $T_A$ was to reject $H_0$,

while the test decision for the test $T_{RA}$ was to fail to reject $H_0$ at the $\alpha = 0.05$ level – thus $POW(T_A) > POW(T_{RA})$.

## 3.8   Discussion

In summary, we have derived and applied a new $LDR$ method for two covariance matrices to four tests for testing $HPCHDS$, yielding the new $HPCHDS$ tests $T_{RA}$, $T_{RC}$, $T_{RI}$, and $T_{RS10}$. Using a Monte Carlo simulation, we have compared and contrasted the powers of the four $HPCHDS$ tests for two population covariance matrices calculated before and after $LDR$ was performed on the original data for all possible reduced dimensions. Here, using Monte Carlo simulations and a real dataset, we have shown that the $HPCHDS$ tests proposed by Srivastava and Yanagihara (2010) and Chaipitak and Chongcharoen (2013), when combined with our proposed $LDR$ method for two covariance matrices, can be considerably more powerful than the power of these tests used without first applying $LDR$ to the original data. Also, we have discovered that the post-$LDR$ test $T_{RA}$ can be considerably less powerful than the no-$LDR$ test $T_A$ calculated from the full-dimensional data.

Finally, after applying $LDR$, we have applied the four $HPCHDS$ tests considered here to a real dataset by using permutation test procedures. The application of our $LDR$ method for two covariance matrices appears to have yielded a considerable power increase over the original-data power for the $HPCHDS$ tests $T_{RC}$, $T_{RI}$, and $T_{RS10}$.

CHAPTER FOUR

Linear Dimension Reduction for Power Improvement of tests for Homogeneity of
Three of More Population Covariance Matrices for a High-Dimensional Scenario

ABSTRACT

We derive and apply a linear dimension reduction ($LDR$) technique for multiple high-dimensional covariance matrices when testing the homogeneity of three or more multivariate-normal population covariance matrices under a high-dimensional setting. Using Monte Carlo simulations, we examine the change in power for testing homogeneity of population covariance matrices under a high-dimensional setting ($HPCHDS$). That is, we examine the difference in the powers for four tests calculated from the full-dimensional data versus the powers of the same tests after applying $LDR$ to the original data. We also perform permutation tests on real data and determine that the use of our $LDR$ method prior to the actual test can considerably improve power. We conclude that a test proposed by Ahmad (2017), when calculated with data from our $LDR$ method, yields a remarkable power improvement for the parameter and sample-size configurations considered here.

### 4.1   Introduction

In many scientific fields of study, such as biomedical imaging, magnetic resonance imaging, tomography, and financial portfolio analysis, one collects high-dimensional data that is greater than the group sample size or sample sizes. In general, increasing the sample size decreases estimator variability, which improves statistical inference. Alternatively, for a fixed sample size, increasing the data dimension increases estimator variability, thus making statistical inference more uncertain. Also, if the data dimension is greater than the group sample size, the corresponding sample covariance matrix is singular and, therefore, non-invertible. Thus, one cannot perform many clas-

47

sical multivariate statistical analyses, and new methods of analysis are often necessary to analyze data in high dimensions. Here, we are concerned with testing hypotheses of homogeneity for $(k > 2)$ population covariance matrices in a high-dimensional setting, that is when the sample size is less than the data dimension.

Several new approaches have been proposed for testing the equality of high-dimensional covariance structures. These methods include hard-thresholding such as the method defined by Chen et al. (2010), partitioning the covariance matrices into testing the diagonals and the off-diagonals, using random projections to reduce the data dimension, and using banded estimators and transformations, as in Peng et al. (2016), to attempt to increase the power of their tests.

In this paper, we derive and explore the efficacy of an $LDR$ matrix for improving the powers of four tests for testing the homogeneity of $(k > 2)$ population covariance matrices under a high-dimensional setting ($HPCHDS$). We perform Monte Carlo power simulations to compare and contrast the powers for three current and one proposed $(k > 2)$-population $HPCHDS$ test after the application of $LDR$ to the original data. In particular, in our power simulations we use five covariance matrix structures that have been used in the $HPCHDS$ literature. We restrict the differences in the covariance matrices to the differences in the hyper-volume as measured by the determinant while ignoring differences in the eigenvector orientation. In addition, we also contrast the power of the four $HPCHDS$ tests calculated with post-$LDR$ data and no-$LDR$ data on a real high-dimensional dataset using permutation tests. We find that when our $LDR$ technique is applied prior to calculating the tests Srivastava and Yanagihara (2010), Ahmad (2017), and Schott (2007), one can gain a considerable increase in the powers.

The first $HPCHDS$ test, generally attributed to Schott (2007), was based on a high-dimensional squared Frobenius norm ($HDSFN$) for two covariance matrices

given in Ledoit and Wolf (2004). This squared norm is

$$HDSFN := \frac{1}{p}\operatorname{tr}\left(\mathbf{\Sigma}_i^2\right) + \frac{1}{p}\operatorname{tr}\left(\mathbf{\Sigma}_j^2\right) - \frac{2}{p}\operatorname{tr}\left(\mathbf{\Sigma}_i\mathbf{\Sigma}_j\right), \quad i \neq j. \tag{4.1}$$

The inclusion of the divisor $p$ in (4.1) yields several desirable properties that one can find in Ledoit and Wolf (2004). Three *HPCHDS* tests motivated by the *HDSFN* have been proposed by Srivastava and Yanagihara (2010), Srivastava et al. (2014), and Ahmad (2017). In addition, Srivastava and Yanagihara (2010) and Chaipitak and Chongcharoen (2013) have proposed tests that use ratios of the summands of the *HDSFN*.

The remainder of the paper is structured as follows. In Section 4.2 we define notation used throughout the paper, and describe consistent estimators used in the considered *HPCHDS* tests. In Section 4.3 we describe four *HPCHDS* tests, and in Section 4.4 we propose a new *LDR* method to reduce the original data dimension before calculating $(k > 2)$-population *HPCHDS* tests. We then describe our power-simulation design in Section 4.5. Next, we present our simulation results contrasting the power of the *HPCHDS* tests using no-*LDR* and using post-*LDR* data in Section 4.6. We then contrast the efficacy of four *HPCHDS* tests calculated with post-*LDR* data to the four tests calculated with the original data on a real high-dimensional dataset in Section 4.7. Finally, we briefly discuss our power-contrast results from the application of *LDR* in Section 4.8.

## *4.2   Notation*

We use the notation $\mathbb{R}_{m \times n}$ to represent the vector space of all $m \times n$ matrices over the real field $\mathbb{R}$. The symbol $\mathbb{R}_{n \times n}^S$ represents all $n \times n$ symmetric matrices of real numbers. The symbol $\mathbb{R}_n^>$ represents the cone of all symmetric positive-definite matrices in $\mathbb{R}_{n \times n}$, and the notation $\mathbb{R}_n^\geq$ represents all symmetric nonnegative-definite matrices in $\mathbb{R}_{n \times n}$. In addition, $\mathcal{C}(\mathbf{A})$ represents the column space of $\mathbf{A} \in \mathbb{R}_{m \times n}$ and we use $SVD(\mathbf{A})$ to represent the singular value decomposition of $\mathbf{A}$.

Also, let $\mathbf{X}_i := \begin{bmatrix} \mathbf{x}_{i1} \vdots \mathbf{x}_{i2} \vdots \cdots \vdots \mathbf{x}_{in_i} \end{bmatrix} \in \mathbb{R}^{p \times n_i}$, represent a data matrix randomly sampled from the $i^{th}$ population so that $\mathbf{X}_i \sim MN_{pn}(\mathbf{M}_i, \mathbf{I}_n \otimes \boldsymbol{\Sigma}_i)$ with $\mathbf{M}_i \in \mathbb{R}^{p \times n}$ and $\boldsymbol{\Sigma}_i \in \mathbb{R}_p^>$. Therefore, $\mathbf{x}_{ij} \sim N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ are independent $p$-dimensional random vectors for $i = 1, 2, \ldots, k$, and $j = 1, 2, \ldots, n_i$. In addition, we use the horizontally concatenated matrix $\mathbf{X} := \begin{bmatrix} \mathbf{X}_1 \vdots \mathbf{X}_2 \vdots \ldots \vdots \mathbf{X}_k \end{bmatrix}$ to represent the complete data matrix.

Consider the following estimators for the multivariate-normal parameters $\boldsymbol{\mu}_i$ and $\boldsymbol{\Sigma}_i$, $i = 1, 2, \ldots, k$, namely, the $i^{th}$ sample mean

$$\overline{\mathbf{x}}_i := \frac{1}{n_i} \mathbf{X}_i \mathbf{j}$$

and the sample covariance matrix

$$\mathbf{S}_i := \frac{1}{n_i - 1} \mathbf{X}_i \left( \mathbf{I}_{n_i} - \frac{1}{n_i} \mathbf{J}_{n_i} \right) \mathbf{X}_i^T,$$

where $\mathbf{J}_{n_i} \in \mathbb{R}^{n_i \times n_i}$ is a matrix of ones and $\mathbf{j} \in \mathbb{R}^{n_i \times 1}$ is a vector of ones. Finally, let

$$\mathbf{V}_i := (n_i - 1) \mathbf{S}_i$$

be the Gram matrix of (4.2) and

$$\mathbf{S} := \frac{\sum\limits_{i=1}^{k} \mathbf{V}_i}{\sum\limits_{i=1}^{k} (n_i - 1)}, \quad i = 1, \ldots, k,$$

be the pooled sample covariance matrix.

### 4.3  Four tests for Homogeneity of $(k > 2)$ Covariance Matrices Under a High-Dimensional Setting

We now describe four $(k > 2)$-population *HPCHDS* tests. Three of these tests have been proposed by Schott (2007), Srivastava and Yanagihara (2010), and Ahmad (2017). Furthermore, we extend the test of Chaipitak and Chongcharoen (2013) from a two-population *HPCHDS* test to a $(k > 2)$-population *HPCHDS* test. We describe these $(k > 2)$-population *HPCHDS* tests in the following four subsections.

Schott (2007) has proposed an *HPCHDS* test based on the *HDSFN* in (4.1).

The test statistic from Schott (2007) is

$$Q_{Sc} := \sum_{i>j}^{k} \frac{\hat{a}_{2i} + \hat{a}_{2j} - \frac{2}{p}\,\mathrm{tr}\,(\mathbf{S}_i\mathbf{S}_j)}{\sqrt{\widehat{\mathrm{Var}(q_{Sc})}}},$$

where $\hat{a}_{2m}$ is defined in (C.3), $m = 1, 2, \ldots, k$, $q_{Sc} := \hat{a}_{2i} + \hat{a}_{2j} - \frac{2}{p}\,\mathrm{tr}\,(\mathbf{S}_i\mathbf{S}_j)$ estimates the sum of squared elements of the *HDSFN*, and

$$\widehat{\mathrm{Var}(q_{Sc})} = 4\hat{a}_2^2 \left\{ \sum_{i<j}^{k} \left( \frac{p}{(n_i - 1)} + \frac{p}{(n_j - 1)} \right) + (k-1)(k-2)\sum_{i=1}^{k} \left( \frac{p}{(n_i - 1)} \right)^2 \right\}^2,$$

where $\hat{a}_2$ is defined in (C.4). We refer to the *HPCHDS* test performed using the test statistic $Q_{Sc}$ by $T_{Sc}$. Schott (2007) has shown that $\widehat{\mathrm{Var}(q_{Sc})} \xrightarrow{P} \mathrm{Var}\,(q_{Sc})$, assuming $\mathrm{H}_0$ is true. Schott (2007) has also shown that $Q_{Sc} \overset{\cdot}{\sim} N(0,1)$ under $\mathrm{H}_0$ as $(p, n_i) \to \infty$, $i = 1, 2, ..., k$.

*4.3.2  Srivastava and Yanagihara* (2010)

Extrapolating from their two-population *HPCHDS* test, Srivastava and Yanagihara (2010) have extended their test to $(k > 2)$-population covariance matrices. Their test statistic is

$$Q_{S10} := \sum_{i=1}^{k} \frac{\left( \hat{\gamma}_i - \bar{\hat{\gamma}} \right)^2}{\hat{\xi}_i^2},$$

where

$$\hat{\gamma}_i := \frac{\hat{a}_{2i}}{\hat{a}_{1i}^2}, \tag{4.2}$$

$\hat{a}_{1i}$, is defined in (C.1) , $\hat{a}_{2i}$ is defined in (C.3), $i = 1, 2, \ldots, k$,

$$\bar{\hat{\gamma}} := \frac{\sum_{i=1}^{k} \hat{\gamma}_i / \hat{\xi}_i^2}{\sum_{i=1}^{k} 1 / \hat{\xi}_i^2}$$

is a pooled estimator for the right-hand side of (4.2),

$$\hat{\xi}_i^2 := \frac{4}{(n_i - 1)^2} \left\{ \frac{\hat{a}_2^2}{\hat{a}_1^4} + \frac{2(n_i - 1)}{p} \left( \frac{\hat{a}_2^3}{\hat{a}_1^6} - \frac{2\hat{a}_2\hat{a}_3}{\hat{a}_1^5} + \frac{\hat{a}_4}{\hat{a}_1^4} \right) \right\},$$

where $\hat{a}_1$ is defined in (C.2), $\hat{a}_2$ is defined in (C.4), $\hat{a}_3$ is defined as (C.5), and $\hat{a}_4$ is defined as (C.6). We refer to the *HPCHDS* test conducted with the test statistic $Q_{S10}$ by $T_{S10}$. Assuming $\mathrm{H}_0$ holds, they have also shown that $\hat{\gamma}_i \xrightarrow{P} \frac{\mathrm{tr}(\boldsymbol{\Sigma}_i^2)}{\mathrm{tr}(\boldsymbol{\Sigma}_i)^2}$, $\hat{\xi}_i^2 \xrightarrow{P} \mathrm{Var}\left(\hat{\gamma}_i - \overline{\hat{\gamma}}\right)$ and $Q_{S10} \overset{\cdot}{\sim} \chi^2_{k-1}$ as $(p, n_i) \to \infty, i = 1, 2, ..., k$. The power-contrast simulations in Srivastava and Yanagihara (2010) have compared the powers of the *HPCHDS* tests $T_{Sc}$ and $T_{S7}$ to $T_{S10}$ for only one population covariance-matrix structure.

### *4.3.3 Chaipitak-Barnard (2018)*

Next, we utilize the pooled estimator $\hat{a}_2$, given in (C.4), to extend the *HPCHDS* test introduced by Chaipitak and Chongcharoen (2013) from the two-covariance-matrix case to the $(k > 2)$-covariance-matrix case. Our proposed $(k > 2)$-covariance-matrix *HPCHDS* test statistic is

$$Q_{CB} := \sum_{i=1}^{k} \frac{\hat{a}_{2i}/\hat{a}_2 - 1}{\sqrt{\widehat{\mathrm{Var}(q_{CB})}}},$$

where $\hat{a}_{2i}$ is defined in (C.3), $i = 1, 2, \ldots, k$, $\hat{a}_2$ is given in (C.4), $q_{CB} := \sum_{i=1}^{k} \hat{a}_{2i}/\hat{a}_2 - 1$, and

$$\widehat{\mathrm{Var}(q_{CB})} = 4 \left\{ \frac{2\hat{a}_4^*}{p\hat{a}_2^2} \sum_{i=1}^{k} \frac{1}{n_i - 1} + \sum_{i=1}^{k} \frac{1}{(n_i - 1)^2} \right\},$$

where $\hat{a}_4^*$ is defined in (C.7). We refer to the *HPCHDS* test conducted with the test statistic $Q_{CB}$ by $T_{CB}$.

### *4.3.4 Ahmad (2017)*

Finally, Ahmad (2017) has extended his *HPCHDS* test from $k = 2$ to $k > 2$ covariance matrices. His test statistic is

$$Q_A := \frac{(k-1) \sum_{i=1}^{k} E_i - 2 \sum_{i \neq j}^{k} E_{ij}}{\sqrt{\widehat{\mathrm{Var}(q_A)}}},$$

where $E_{ij} := \operatorname{tr}(\mathbf{S}_i \mathbf{S}_j), 1 \le i < j \le k$. The estimators $E_i$, $q_A$, and $\widehat{\operatorname{Var}(q_A)}$ are defined as

$$E_i := \frac{(n_i - 1)}{n_i (n_i - 2)(n_i - 3)} \left\{ (n_i - 1)(n_i - 2) \operatorname{tr}(\mathbf{S}_i^2) + \left[ \operatorname{tr}(\mathbf{S}_i) \right]^2 \right.$$
$$\left. - \frac{n_i}{(n_i - 1)} \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)^T (\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)(\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)^T (\mathbf{x}_{ij} - \overline{\mathbf{x}}_i) \right\},$$

$$q_A := (k - 1) \sum_{i=1}^{k} E_i - 2 \sum_{i \ne j}^{k} E_{ij},$$

and

$$\widehat{\operatorname{Var}(q_A)} = 4 \left[ \operatorname{tr}(\mathbf{S}^2) \right]^2 \left\{ (k-1)^2 \sum_{i=1}^{k} \frac{1}{n_i^2} + \sum_{\substack{i=1 \\ i<j}}^{k} \sum_{j=1}^{k} \frac{2}{n_i n_j} \right\}.$$

We refer to the *HPCHDS* test performed using the test statistic $Q_A$ by $T_A$.

The test $q_A$ estimates the sums of squared elements of $[\mathbf{\Sigma}_j - \mathbf{\Sigma}_i]$. Ahmad (2017) has also proven that $\widehat{\operatorname{Var}(q_A)} \xrightarrow{P} \operatorname{Var}(q_A)$, assuming $\mathrm{H}_0$ is true. Additionally, Ahmad (2017) has shown that for the $(k > 2)$-*HPCHDS* case, we have $Q_A \overset{\cdot}{\sim} N(0,1)$ as $(p, n_i) \to \infty, i = 1, 2, 3$, provided $\mathrm{H}_0$ holds. However, Ahmad (2017) has not contrasted the power of $T_A$ with any other competing *HPCHDS* tests.

### 4.4 Linear Dimension Reduction for Power Improvement for Tests for the Homogeneity of $(k > 2)$ Covariance Matrices in a High-Dimensional Setting

Below, we propose an *LDR* method and prove a separability theorem concerning homogeneity for $(k > 2)$ population covariance matrices. Then, using the singular value decomposition (*SVD*), we derive an *LDR* matrix for the sample data that allows us to retain much of the distinguishing information in the $(k - 1)$ differences $[\mathbf{S}_k - \mathbf{S}_1], [\mathbf{S}_{k-1} - \mathbf{S}_1], \ldots, [\mathbf{S}_2 - \mathbf{S}_1]$. Thus, we propose employing an *LDR* matrix derived specifically for $k$ sample covariance matrices to increase the power of $(k > 2)$-population *HPCHDS* tests. Our new *LDR* matrix for $(k > 2)$ population covariance matrices is based on a property of a linear sufficient matrix for $(k > 2)$ population covariance matrices proposed by Peters et al. (1978).

For the proof of the theorem, we need the symmetrized Bregman log-determinant divergence among $k$ positive-definite matrices, which is

$$D(\mathbf{\Sigma}_1, \mathbf{\Sigma}_2, \ldots, \mathbf{\Sigma}_k) := \sum_{j \neq i} D(\mathbf{\Sigma}_i, \mathbf{\Sigma}_j), \tag{4.3}$$

where $1 \leq i, j \leq k$ , and

$$D(\mathbf{\Sigma}_i, \mathbf{\Sigma}_j) := \left[\text{tr}\left(\mathbf{\Sigma}_i \mathbf{\Sigma}_j^{-1}\right)\right] - \log\left[\det\left(\mathbf{\Sigma}_i \mathbf{\Sigma}_j^{-1}\right)\right] - p, \tag{4.4}$$

with $i \neq j$.

We now prove a theorem demonstrating that under certain conditions, our proposed parameter-based $LDR$ matrix for $(k > 2)$ full-rank population covariance matrices preserves the full-dimensional separability measure (4.3) for some reduced dimension $r$, where $1 \leq r < p$.

**Theorem.** *Suppose we have $k$ $(k > 2)$ multivariate normal populations with covariance matrices $\mathbf{\Sigma}_i \in \mathbb{R}_p^>$, $i = 1, 2, \ldots, k$, and let*

$$\mathbf{H} := \left[\mathbf{\Sigma}_2 - \mathbf{\Sigma}_1 \vdots \mathbf{\Sigma}_3 - \mathbf{\Sigma}_1 \vdots \ldots \vdots \mathbf{\Sigma}_k - \mathbf{\Sigma}_1\right]. \tag{4.5}$$

*Next, let $SVD(\mathbf{H}) = \mathbf{F\Lambda G} \in \mathbb{R}_{p \times (k-1)p}$, where $\mathbf{F} \in \mathbb{R}_{p \times r}$ and $\text{rank}(\mathbf{F}) = \text{rank}(\mathbf{H})$. Also, let the symmetrized Bregman log-determinant divergence $D(\mathbf{\Sigma}_1, \mathbf{\Sigma}_2, \ldots, \mathbf{\Sigma}_k)$ be defined in (4.3) and (4.4). Then,*

$$D(\mathbf{\Sigma}_1, \mathbf{\Sigma}_2, \ldots, \mathbf{\Sigma}_k) = D(\mathbf{F}^+ \mathbf{\Sigma}_1 \mathbf{F}^{+T}, \mathbf{F}^+ \mathbf{\Sigma}_2 \mathbf{F}^{+T}, \ldots, \mathbf{F}^+ \mathbf{\Sigma}_k \mathbf{F}^{+T}).$$

*Proof.* The proof follows from Lemma C.2.2 in Appendix C.

Next, let

$$\widehat{\mathbf{H}} := \left[\mathbf{S}_2 - \mathbf{S}_1 \vdots \mathbf{S}_3 - \mathbf{S}_1 \vdots \ldots \vdots \mathbf{S}_k - \mathbf{S}_1\right]$$

be an estimator of $\mathbf{H}$, where $\mathbf{H}$ is given in (4.5). Because $\text{rank}(\widehat{\mathbf{H}}) \neq \text{rank}(\mathbf{H})$ and $\text{rank}(\widehat{\mathbf{H}})$ is unknown, one cannot directly apply the above theorem to obtain

an operational *LDR* matrix $\mathbf{F}^+ \in \mathbb{R}_{p \times q}$ that preserves the full-feature covariance matrix separability measure $D(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2, \ldots, \boldsymbol{\Sigma}_k)$. Moreover, we often wish to obtain a very low dimensional representation of dimension $q$ for the original data, where $1 \leq q \ll \text{rank}(\widehat{\mathbf{H}}) \leq \text{rank}(\mathbf{H})$. That is, we desire to construct an *LDR* matrix that yields low-dimensional data from the original data that preserves almost all of the original $p$-dimensional information in the data concerning the differences $[\mathbf{S}_k - \mathbf{S}_1]$, $[\mathbf{S}_{k-1} - \mathbf{S}_1], \ldots, [\mathbf{S}_2 - \mathbf{S}_1]$.

First, let $SVD(\widehat{\mathbf{H}}) = \mathbf{U}\boldsymbol{\Lambda}\mathbf{V}^T$, let $\mathbf{U}^{(q)}$ denote the concatenated $q$ eigenvectors of $\widehat{\mathbf{H}}$ corresponding to the $q$ largest singular values, and let $\widehat{\mathbf{F}}_{p \times q} := \mathbf{U}^{(q)}$. Then $\left[\widehat{\mathbf{F}}_{p \times q}\right]^+ \in \mathbb{R}_{q \times p}$ is an *LDR* matrix for reducing the feature dimension from $p$ to $q$, where $1 \leq q < \text{rank}(\widehat{\mathbf{H}})$, while preserving most of the separability information in the estimated covariance matrices $\mathbf{S}_i \in \mathbb{R}_p^{\geq}$, $i = 1, 2, \ldots, k$.

Mapping the high-dimensional data matrix $\mathbf{X}$ onto $\mathcal{C}\left([\widehat{\mathbf{F}}_{p \times q}^{(q)}]^+\right)$ could enhance our ability to detect differences in the covariance matrices $\mathbf{S}_i \in \mathbb{R}_p^{\geq}$ and, thus, differences in $\boldsymbol{\Sigma}_i \in \mathbb{R}_p^{>}, i = 1, 2, \ldots, k$, because of the decreased number of parameters that must be estimated. However, the *LDR* matrices $\left[\widehat{\mathbf{F}}_{p \times q}\right]^+ \in \mathbb{R}_{q \times p}$ and $\left[\widehat{\mathbf{F}}_{p \times q}\right]^T \in \mathbb{R}_{q \times p}$ span the same subspace. Therefore, we use the computationally simpler matrix $\left[\widehat{\mathbf{F}}_{p \times q}\right]^T \in \mathbb{R}_{q \times p}$ as our *LDR* matrix to reduce the data from $p$ to $q$ dimensions.

### 4.5   Monte Carlo Power Simulation Design

#### 4.5.1   Simulation Covariance Structures

The covariance matrix structures in our Monte Carlo simulations were selected from the *HPCHDS* literature. We have compared test powers across five covariance matrix structures. The group sample sizes used in these simulations were $n_1 = n_2 = n_3$.

First, we use the constant-times-identity covariance-matrix structure. For our simulation, the parameters for the null and alternative hypotheses are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{\Sigma}_3 = \mathbf{I}_p$$

and

$$H_A : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = \mathbf{I}_p \text{ and } \mathbf{\Sigma}_2 = \sigma^2 \mathbf{I}_p,$$

respectively.

Second, we use the compound-symmetric covariance matrix class for our simulation. The parameters for the null and alternative hypotheses are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{\Sigma}_3 = \sigma_1^2 \mathbf{I}_p + \sigma_2^2 \mathbf{J}_p$$

and

$$H_A : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = \sigma_1^2 \mathbf{I}_p + \sigma_2^2 \mathbf{J}_p \text{ and } \mathbf{\Sigma}_2 = \sigma_{1A}^2 \mathbf{I}_p + \sigma_{2A}^2 \mathbf{J}_p,$$

respectively, where $\mathbf{J}_p \in \mathbb{R}^{p \times p}$ is a matrix of ones.

Third, we use the autoregressive covariance-matrix structure for our simulation. The parameters for the null and alternative hypotheses are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{\Sigma}_3 = \mathbf{U}_0$$

and

$$H_A : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = \mathbf{U}_0 \text{ and } \mathbf{\Sigma}_2 = \mathbf{U}_1,$$

where $\mathbf{U}_0 = \sigma_{ij} = 0.1^{|i-j|}$, $\mathbf{U}_1 = \sigma_{ij} = 0.3^{|i-j|}$, $1 \leq i, j \leq k$.

Fourth, we use the heterogeneous autoregressive covariance-matrix structure. For our simulation, we use the heterogeneous autoregressive covariance-matrix structures similar to those in Srivastava et al. (2014). These are created as follows. First, let $\sigma_l := 1 + (-1)^{l+1} Q_l / 2$ where, $Q_l \sim Unif(0,1)$ and $l = 1, 2, \ldots, p$. The simulation parameters are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{\Sigma}_3 = \sigma_i \sigma_j 0.1^{|i-j|^{\frac{1}{10}}}$$

and

$$H_A : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = \sigma_i \sigma_j 0.1^{|i-j|^{\frac{1}{10}}} \text{ and } \mathbf{\Sigma}_2 = \sigma_i \sigma_j 0.3^{|i-j|^{\frac{1}{10}}},$$

respectively, where $1 \leq i, j \leq k$.

Last, we examine an unstructured covariance-matrix structure, which we model as

$$\mathbf{U}_2 = \sigma_{ij} = \sigma_{ji} := \begin{cases} (-1)^{i+j} \left( \frac{0.10}{j} \right), & i < j \\ 1, & i = j \end{cases}$$

and

$$\mathbf{U}_3 = \sigma_{ij} = \sigma_{ji} := \begin{cases} (-1)^{i+j} \left( \frac{0.05}{j} \right), & i < j \\ 1, & i = j. \end{cases}$$

Then, the simulation covariance matrices are

$$H_0 : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{\Sigma}_3 = \mathbf{U}_2$$

and

$$H_A : \mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = \mathbf{U}_2 \text{ and } \mathbf{\Sigma}_2 = \mathbf{U}_3.$$

### 4.5.2 Monte Carlo Power Simulation Description

We now describe the simulation design used to compare the powers for the three-population no-*LDR HPCHDS* tests $T_{Sc}$, $T_{S10}$, $T_{CB}$, and $T_A$ and the post-*LDR HPCHDS* tests $T_{ScR}$, $T_{S10R}$, $T_{CBR}$, and $T_{AR}$. Our goal is to discover evidence that our proposed *LDR* method improves the power of these four tests. Using R, we generated 10,000 independent matrix-normal sample-data matrices from $MN_p(\mathbf{0}, \mathbf{I}_{n_i} \otimes \mathbf{\Sigma}_i)$ populations, where $\mathbf{\Sigma}_i \in \mathbb{R}_p^>$, $i = 1, 2, 3$, are the covariance matrices under $H_A$. We then determined $SVD(\widehat{\mathbf{M}}^{(j)})$ to extract the matrix $\left[ \widehat{\mathbf{F}}_{p \times q}^{(j)} \right]^T \in \mathbb{R}_{q \times p}$. Next, for $N := \sum_{i=1}^{3} n_i$, we reduced the dimensions of the sample data by mapping the full-dimensional data matrix $\mathbf{X}^{(j)} \in \mathbb{R}_{p \times N}$ onto $C([\widehat{\mathbf{F}}_{p \times q}^{(j)}]^T)$ so that the $j^{th}$ reduced-data

set is $\mathbf{Y}^{(j)} = \left[\widehat{\mathbf{F}}_{p\times q}^{(j)}\right]^T \mathbf{X}^{(j)}$, where $q$ is the targeted reduced-data dimension. We then calculated $T_{(*R),j}$ from each simulated reduced data matrix $\mathbf{Y}^{(j)}$ and calculated the simulated critical values ($SCVs$) for $T_{AR}$, $T_{CBR}$, and $T_{ScR}$ by

$$SCV_{1-\alpha/2} = \inf\left\{x \in \mathbb{R} : 1 - \alpha/2 \leq \widehat{F}_{T_{(*R)}}(x)\right\}$$

and

$$SCV_{\alpha/2} = \sup\left\{x \in \mathbb{R} : \alpha/2 \geq \widehat{F}_{T_{(*R)}}(x)\right\},$$

where $\widehat{F}_{T_{(*R)}}(x)$ is the empirical distribution function of $T_{(*R)}$ with $\alpha = 0.05$. For $T_{S10R}$ the SCV was determined with

$$SCV_{1-\alpha} = \inf\left\{x \in \mathbb{R} : 1 - \alpha \leq \widehat{F}_{T_{(S10R)}}(x)\right\}.$$

We used a similar approach to determine SCVs for the no-$LDR$ tests $T_{Sc}$, $T_{S10}$, $T_{CB}$, and $T_A$ except using the unreduced datasets $\mathbf{X}^{(j)}$.

Then, for each considered reduced dimension $q$, we generated 10,000 independent multivariate normal samples from $N_p(\mathbf{0}, \boldsymbol{\Sigma}_i)$ distributions for $i = 1, 2, \ldots, k$, where $\boldsymbol{\Sigma}_i \in \mathbb{R}_p^>$ is the $i^{th}$ population covariance matrix under the $\mathrm{H_A}$. Next, we reduced the original-data dimension by mapping the $j^{th}$ complete dataset $\mathbf{X}^{(j)}$ onto $\mathcal{C}\left([\widehat{\mathbf{F}}_{p\times q}^{(j)}]^T\right)$. That is, we applied the linear transformation $\mathbf{Y}^{(j)} = \left[\widehat{\mathbf{F}}_{p\times q}^{(j)}\right]^T \mathbf{X}^{(\mathbf{j})}$ and then calculated the tests $T_{(*R),j}$ from the reduced data-matrix $\mathbf{Y}^{(j)}$ for each $j$, where $1 \leq j \leq 10,000$. We then estimated the power by

$$POW(T_{(*R)}) := \frac{\sum\limits_{j=1}^{10,000} I[T_{(*R),j} \in RR(T_{(*R)})]}{10,000},$$

where $RR(T_{(*R)})$ is the rejection region for the test $T_{(*R)}$ and $I[\cdot]$ is the indicator function. We also calculated $POW(T_{Sc})$, $POW(T_{S10})$, $POW(T_{CB})$, and $POW(T_A)$ in a similar manner, but using the unreduced datasets $\mathbf{X}^{(j)}$.

We determined $POW(T_{ScR})$, $POW(T_{S10R})$, $POW(T_{CBR})$, and $POW(T_{AR})$ and $POW(T_{Sc})$, $POW(T_{S10})$, $POW(T_{CB})$, and $POW(T_A)$ for the common sample sizes

of $n_i \in 5, 10, 15, 20, i = 1, 2, 3$, and complete-data dimensions of $p \in 20, 40, 80, 160$. These power values are presented in Table 4.1. Also, power curve simulation results, plotted versus $q$, are displayed for $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ using post-$LDR$ data. In addition, power-difference curves, plotted versus $q$, are shown in Figures 4.1–4.5 for $n_i = 10$, $i = 1, 2, 3$, and $p = 160$ for each $q$, where $q \in \{1, 2, \ldots, 159\}$. We performed the power simulations in parallel using R and the covTestR package.

### 4.6  Monte Carlo Simulation Power-Contrast Results

In this section, we present our simulated power-contrast results. Let the notation $POW(T_{(*)})$ represent the power of the test $T_{(*)}$, and let $DPOW\left(T_{(*R)}, T_{(*)}\right) := \left[POW(T_{(*R)}) - POW(T_{(*)})\right]$. We provide a table of power values and we provide five figures displaying the curves for $POW(T_*)$ and power-difference curves for $DPOW\left(T_{(*R)}, T_{(*)}\right)$. The power and power difference are are reported for each $HPCHDS$ test at each $q$ such that $q \in \{1, 2, \ldots, 159\}$. We fit all power curves using generalized linear models with b-splines.

### 4.6.1  Power-Simulation Summary Table

In Table 4.1 we present the post-$LDR$ values for $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ calculated with post-$LDR$ data. We also report $POW(T_A)$, $POW(T_{CB})$, $POW(T_{Sc})$, and $POW(T_{S10})$ for the same $HPCHDS$ tests calculated with the unreduced data. The powers were reported for the scenarios with $p \in \{80, 160\}$, $n_i \in \{5, 10\}$, $i = 1, 2, 3$, and for $q \in \{5, 10, 15, 20\}$ with $n_i = 5$, $i = 1, 2, 3$ and $q \in \{10, 20, 30, 40\}$ with $n_i = 10$, $i = 1, 2, 3$.

The overwhelming result in the table is that our proposed $LDR$ method, applied to the original data, can yield a surprisingly large increase in the post-$LDR$ powers $POW(T_{AR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ compared to $POW(T_A)$, $POW(T_{Sc})$, and $POW(T_{S10})$ even though one may not map to the optimal $q$ for a particular

sample size or sizes. As the table shows, the application of $LDR$ to the original data can cause a decrease in $POW(T_{ScR})$ as compared to $POW(T_{Sc})$. Hence, the application of $LDR$ to the original high-dimensional data does not guarantee increased power.

Table 4.1: A Table Contrasting $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ with $POW(T_A)$, $POW(T_{CB})$, $POW(T_{Sc})$, and $POW(T_{S10})$ for Testing $HPCHDS$ for three autoregressive population covariance matrices with parameters $\Sigma_1 = \Sigma_3 = 0.1^{|i-j|}$ and $\Sigma_2 = 0.3^{|i-j|}$.

| $p$ | $n_1 = n_2 = n_3$ | $q$ | $POW(T_{AR})$ | $POW(T_{CBR})$ | $POW(T_{ScR})$ | $POW(T_{S10R})$ |
|---|---|---|---|---|---|---|
| 80 | 5 | 5 | 0.89 | 0.08 | 0.01 | 0.84 |
| | | 10 | 0.99 | 0.34 | 0.97 | 0.90 |
| | | 15 | 0.63 | 0.02 | 1.00 | 0.09 |
| | | 20 | 0.58 | 0.03 | 1.00 | 0.15 |
| | 10 | 10 | 1.00 | 0.00 | 0.00 | 0.92 |
| | | 20 | 1.00 | 0.15 | 0.01 | 0.58 |
| | | 30 | 0.97 | 0.01 | 1.00 | 0.01 |
| | | 40 | 0.89 | 0.01 | 1.00 | 0.01 |
| 160 | 5 | 5 | 0.76 | 0.36 | 0.00 | 0.91 |
| | | 10 | 1.00 | 0.45 | 1.00 | 0.08 |
| | | 15 | 0.74 | 0.01 | 1.00 | 0.16 |
| | | 20 | 0.70 | 0.01 | 1.00 | 0.24 |
| | 10 | 10 | 1.00 | 0.00 | 0.00 | 0.99 |
| | | 20 | 1.00 | 0.63 | 0.12 | 0.99 |
| | | 30 | 1.00 | 0.05 | 1.00 | 0.11 |
| | | 40 | 0.99 | 0.00 | 1.00 | 0.02 |

| $p$ | $n_1 = n_2 = n_3$ | $p$ | $POW(T_A)$ | $POW(T_{CB})$ | $POW(T_{Sc})$ | $POW(T_{S10})$ |
|---|---|---|---|---|---|---|
| 80 | 5 | 80 | 0.08 | 0.07 | 0.03 | 0.05 |
| | 10 | 80 | 0.13 | 0.09 | 0.02 | 0.04 |
| 160 | 5 | 160 | 0.13 | 0.09 | 0.02 | 0.04 |
| | 10 | 160 | 0.13 | 0.12 | 0.01 | 0.06 |

### 4.6.2 Power Curves and Power-Difference Curves for the $T_{AR}$, $T_{CBR}$, $T_{ScR}$, and $T_{S10R}$ Tests for Three Constant-Times-Identity Covariance Structures

In Figure 4.1, we present plots for the post-$LDR$ power curves for $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ and the power-difference curves described in Subsection 4.6.1. The original-dimension powers were $POW(T_A) = 0.78$,

$POW(T_{CB}) = 0.71$, $POW(T_{Sc}) = 0.01$, and $POW(T_{S10}) = 0.01$. The $HPCHDS$ hypothesis tests in this section were performed for three constant-times-identity covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = \mathbf{I}_p$ and $\mathbf{\Sigma}_2 = (1.5)\mathbf{I}_p$.



Figure 4.1: Reduced-dimension curves and power-difference curves for $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ with three constant-times-identity covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = \mathbf{I}_p$ and $\mathbf{\Sigma}_2 = (1.5)\mathbf{I}_p$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 10, i = 1, 2, 3$.

Comparing the power and power-difference plots, we had similar trends in the two sets of curves. We also observed that the powers $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ had curves with maximums between the individual group sample size $n_i, i = 1, 2, 3$, denoted by the solid vertical line and the sum of the sample sizes denoted by the dotted vertical line. The curves $POW(T_{ScR})$ and

$POW(T_{S10R})$ displayed maximum powers near 1.0. Also, $T_{CBR}$ displayed a maximum power of $POW(T_{CBR}) \approx 0.70$ at $q = 155$. Furthermore, $POW(T_{AR}) = 1.0$ for $q$, where $q \in \{8, 9, \ldots, 105\}$. For the power-difference curves, the maximum $DPOW(T_{AR}, T_A)$ value was 0.24 for $q$, where $q \in \{6, 7, \ldots, 100\}$. Also, the maximum $DPOW(T_{S10R}, T_{S10})$ value occurred at $q = 12$ with a maximum power increase near 0.99, and we observed that $DPOW(T_{ScR}, T_{Sc})$ attained an increase of 0.99 for $q$, where $q \in \{25, 26, \ldots, 112\}$. However, $DPOW(T_{CBR}, T_{CB})$ was negative for every considered value of $q$.

### 4.6.3   Power Curves and Power-Difference Curves for the $T_{AR}$, $T_{CBR}$, $T_{ScR}$, and $T_{S10R}$ Tests for Three Compound-Symmetric Covariance Structures

In Figure 4.2, we present plots for $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ tests and $DPOW(T_{S10R}, T_{S10})$, $DPOW(T_{AR}, T_A)$, $DPOW(T_{ScR}, T_{Sc})$, and $DPOW(T_{CBR}, T_{CB})$. The full-dimension powers at $p = 160$ were $POW(T_A) = 0.26$, $POW(T_{CB}) = 0.19$, $POW(T_{Sc}) = 0.1$, and $POW(T_{S10}) = 0.03$. The $HPCHDS$ hypothesis tests in this section were performed for three compound-symmetric covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = 0.99\mathbf{I}_p + 0.01\mathbf{J}_p$ and $\mathbf{\Sigma}_2 = 0.95\mathbf{I}_p + 0.05\mathbf{J}_p$.

We see that $POW(T_{ScR})$ and $POW(T_{AR})$ produced similar power curves that increased as the reduced-data dimension decreased. Also, $POW(T_{ScR}) = 1.0$ for $q \in \{25, 26, \ldots, 100\}$ but decreased rapidly for $q < 25$. Additionally, $POW(T_{S10R}) = 1.0$ for $q \in \{10, 11, \ldots, 18\}$. The test $T_{S10R}$ yielded a maximum power increase at $q = 18$ with $DPOW(T_{S10R}, T_{S10}) = 0.96$ while $DPOW(T_{ScR}, T_{Sc})$ 0.99 was the maximum power increase at $q$, where $q \in \{31, 32, \ldots, 100\}$. Also, $DPOW(T_{AR}, T_A) = 0.75$ was the maximum which was at $q = 50$. Also, $DPOW(T_{CBR}, T_{CB})$ peaked at $q = 18$ with a value of 0.35. Thus, all tests considered here yielded an increase in power for some $q$ because of the application of $LDR$ for multiple covariance matrices to the original data.

Figure 4.2: Reduced-dimension curves and power-difference curves for POW($T_{AR}$), POW($T_{CBR}$), POW($T_{ScR}$), and POW($T_{S10R}$) with three compound symmetric covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = 0.99\mathbf{I}_p + 0.01\mathbf{J}_p$ and $\mathbf{\Sigma}_2 = 0.95\mathbf{I}_p + 0.05\mathbf{J}_p$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, $p = 160$, and the common sample size was $n_i = 10, i = 1, 2, 3$.

### 4.6.4 Power Curves and Power-Difference Curves for the $T_{AR}$, $T_{CBR}$, $T_{ScR}$, and $T_{S10R}$ Tests for Three Autoregressive Covariance Structures

In Figure 4.3, we present plots for the reduced-dimension power curves for $T_{AR}$, $T_{CBR}$, $T_{ScR}$, and $T_{S10R}$ and for the difference in the reduced-dimension power curves from the original-dimension powers. The reduced-dimension powers and the power-differences were plotted versus each positive integer $q$ for $q \in \{1, 2, ..., 159\}$. The original-dimension powers for $p = 160$ were $POW(T_A) = 0.13$, $POW(T_{CB}) = 0.12$, $POW(T_{Sc}) = 0.01$, and $POW(T_{S10}) = 0.06$. The covariance matrices for which

the power comparisons in this section were performed were for three autoregressive covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = 0.1^{|i-j|}$ and $\mathbf{\Sigma}_2 = 0.3^{|i-j|}$. In the power-curve plot, $POW(T_{ScR}) = 1.0$ was attained at $q = 100$,



Figure 4.3: Reduced-dimension curves and power-difference curves for $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ with three autoregressive covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = 0.1^{|i-j|}$ and $\mathbf{\Sigma}_2 = 0.3^{|i-j|}$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 10, i = 1, 2, 3$.

and $POW(T_{AR}) = 1.0$ occurred at $q = 45$. Also, $POW(T_{ScR})$ decreased rapidly as $q$ decreased for $q < 25$ and $POW(T_{ScR}) = 0.0$ at $q = 20$. In addition, $POW(T_{S10R})$ had a sharp peak at $POW(T_{S10R}) = 1.0$ for $q = 15$.

In the power-difference plot, $DPOW(T_{AR}, T_A) = 0.87$ was the maximum increase, which occurred at $q = 8$, and $DPOW(T_{ScR}, T_{Sc}) \approx 1.0$ was attained for $q$

where $q \in \{30, 31, \ldots, 120\}$. Also, $DPOW(T_{CBR}, T_{CB}) \approx 0.50$ was the maximum which occurred at $q = 18$. Thus, $DPOW(T_{AR}, T_A) > 0$, $DPOW(T_{CBR}, T_{CB}) > 0$, and $DPOW(T_{S10R}, T_{S10}) > 0$, which implied a power increase when $LDR$ was first applied to the original data.

### 4.6.5 Power Curves and Power-Difference Curves for the $T_{AR}$, $T_{CBR}$, $T_{ScR}$, and $T_{S10R}$ Tests for Three Heterogeneous Autoregressive Covariance Structures

In Figure 4.4, we display curves for $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ plotted versus $q$ where $q \in \{1, 2, \ldots, 159\}$, along with the corresponding power-difference curves. The full-dimension powers at $p = 160$ were $POW(T_A) = 0.61$, $POW(T_{CB}) = 0.49$, $POW(T_{Sc}) = 0.01$, and $POW(T_{S10}) = 0.67$. The hypothesis tests for $HPCHDS$ were performed for three heterogeneous autoregressive covariance matrices at $p = 160$. The power plots reveal that $POW(T_{AR})$ and $POW(T_{S10R})$ yielded power maximums at or near the common sample size denoted by the solid vertical line.

Here, $POW(T_{AR}) = 1.0$ occurred at $q = 20$, and $POW(T_{ScR}) = 0.62$ was the maximum which was attained at $q = 35$. Also, $POW(T_{CBR})$ generally decreased as $q$ was reduced. In addition, $POW(T_{S10R})$ peaked at 0.95 at $q = 10$. In the power-difference plot, $DPOW(T_{AR}, T_A)$ peaked at $q = 12$ with a value near 0.39. The curve $DPOW(T_{S10R})$ yielded a maximum power increase of 0.28 at $q = 10$, while $DPOW(T_{ScR})$ peaked at $q = 35$ with a value of 0.62. Last, $DPOW(T_{CBR})$ was negative for all $q$.

Figure 4.4: Reduced-dimension curves and power-difference curves for POW($T_{AR}$), $POW(T_{ScR})$, POW($T_{CBR}$)and pow($T_{S10R}$) with three heterogeneous autoregressive covariance matrices. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 10, i = 1, 2, 3$.

### 4.6.6 Power Curves and Power-Difference Curves for the $T_{AR}$, $T_{CBR}$, $T_{ScR}$, and $T_{S10R}$ Tests for Three Unstructured Covariance Structures

In Figure 4.5, we present plots for $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ and plots for the difference in the post-LDR power curves from the original-dimension power curves plotted versus $q$. For $p = 160$, $POW(T_A) = 0.67$, $POW(T_{CB}) = 0.21$, $POW(T_{Sc}) = 0.01$, and $POW(T_{S10}) = 0.77$. The $HPCHDS$ tests were performed for three unstructured population covariance matrices.

66

Figure 4.5: Reduced-dimension curves and power-difference curves for POW($T_{AR}$), POW($T_{CBR}$), POW($T_{ScR}$), and POW($T_{S10R}$) with three unstructured covariance matrices. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 10, i = 1, 2, 3$.

In the power curve plot, we see that $POW(T_{AR})$ and $POW(T_{CBR})$ yielded essentially parallel power curves with little or no increase in power as $q$ was decreased. The maximums for these two tests were $POW(T_{AR}) = 0.75$ and $POW(T_{CBR}) = 0.24$. Also, $POW(T_{Sc}) = 0.29$ was the maximum at $q = 1$. Additionally, $POW(T_{S10R})$ decreased monotonically in power as the $q$ was decreased and attained a maximum power of 0.76 at $q = 159$.

In the power-difference plot, $DPOW(T_{AR}, T_A)$ yielded a maximum power increase at $q = 5$ with a value of 0.08. Also, $DPOW(T_{AR}, T_A)$ had almost no change

for $q$, $q \in \{1, 2, \ldots, 159\}$, and $DPOW(T_{ScR}, T_{Sc}) = 0.25$ was the maximum power increase which occurred at $q = 1$. Lastly, The curve $DPOW(T_{S10R,T_{S10}})$ was negative for all considered $q$.

*4.6.7 Selection of the Reduced Dimension $q$*

We have shown that our proposed sample $LDR$ method for three covariance matrices can significantly improve $POW(T_{AR})$, $POW(T_{CBR})$, and $POW(T_{ScR})$. However, the practitioner cannot feasibly check every possible power value for all $q$ to determine the optimal reduced dimension. In Figures 4.1 – 4.5, we have displayed the common sample size $n_i$ with a solid vertical line and the total sample size $\sum_{i=1}^{k} n_i$ with a dashed line. For the tests that showed improvement in power, the $HPCHDS$ tests considered here usually attained maximum power for $q$ between the individual group sample size and the sum of the common sample size. Even though some tests plateaued in power increase before the common sample size, we recommend reducing to a dimension $q$, where $\min_i n_i \leq q \leq \sum_{i=1}^{k} n_i, i = 1, 2, \ldots, k$, when the class sample sizes are approximately equal. For very small sample sizes, we found evidence that one should choose $q$, where $\sum_{i=1}^{k} n_i$. Using these guidelines for determining $q$ will not necessarily yield the largest possible power increase. However, one should generally obtain a power increase, provided one does not use a test when the use of $LDR$ is deleterious or with a spiked power curve. More work should be performed on determining a reduced dimension $q$ for the case where the sample sizes are markedly different.

We emphasize that the optimal dimension $q$ that yields maximum power depends on the test employed, the population covariance structure being tested, and the group sample sizes. For these reasons, we do not recommend using only a rank estimation method, such as those found in Luo and Li (2016), Cook and Forzani (2009), and Rohde and Tsybakov (2011).

## 4.7 A Contrast of the Tests $T_A$, $T_{CB}$, $T_{Sc}$, and $T_{S10}$ and $T_{AR}$, $T_{CBR}$, $T_{ScR}$, and $T_{S10R}$ on a Real Dataset

In this section we contrast $POW(T_A)$ with $POW(T_{AR})$, $POW(T_{CB})$ with $POW(T_{CBR})$, $POW(T_{Sc})$ with $POW(T_{ScR})$, and $POW(T_{S10})$ with $POW(T_{S10R})$ on a real-world data from Khan et al. (2001) as curated by Ramey (2016). The Khan dataset contains 63 sample observations with $p = 2,308$ features and $k = 4$ groups. We performed permutation tests to determine appropriate empirical critical values and then compared the results of four *HPCHDS* tests, calculated with and without the application of our *LDR* matrix. Using our *LDR* method on the original data, we reduced the data dimension from $p = 2000$ to $q = 22$, which is the group-minimal sample size.

To accomplish this *HPCHDS* test characteristic contrast, we created 1,000 permutations of the Khan dataset by randomly assigning the 63 observation vectors into one of the four classes as if the data vectors were exchangeable. We then reduced the original feature dimension of the datasets to the common sample size and obtained the observed tests for each of these 1,000 permuted datasets. Finally, we compared the empirical test scores from the reduced and original datasets with the corresponding empirical critical values calculated from the respective permuted-test distributions.

Table 4.2 presents the *HPCHDS* test comparison results of the permutation tests. in Table 4.2 we see that the reduction of the p-values (*PVs*) of the three *HPCHDS* tests $T_{SC}$, $T_{S10}$, and $T_A$ was substantial. The *PV* for the test $T_{SC}$ calculated with the original data was $PV(T_{SC}) = 0.456$, while the *PV* for the reduced-data test $T_{SCR}$ was $PV(T_{RC}) = 0.041$. Also, for the full-dimensional-data, the *PV* for the original-data test $T_{S10}$ was $PV(T_{S10}) = 0.432$, while the *PV* for the reduced-data test $T_{S10R}$ was $PV(T_{S10R}) = 0.013$. In addition, for the full-dimensional-data, the *PV* for the original-data test $T_A$ was $PV(T_A) = 0.325$, while for the reduced-data test $PV(T_{AR}) = 0.032$.

Thus, our real-data *LDR* application demonstrated that for certain *HPCHDS* tests and with an appropriate reduced dimension $q$, one can attain a considerable power increase. However, a power decrease can occur if an inappropriate *HPCHDS* test is combined with our *LDR* method. As an example, the post-*LDR* test, $T_{CB}$, $PV(T_{CBR})$ actually increased considerably over $PV(T_{CB})$. As a result, the test decision for the no-*LDR* test $T_{CB}$ was to reject $H_0$, while the test decision for $T_{CBR}$ was to fail to reject $H_0$ at the $\alpha = 0.05$ level.

Table 4.2: A table contrasting the characteristics of the full-dimensional tests $T_A$, $T_{CB}$, $T_{Sc}$, and $T_{S10}$ with the reduced-dimension tests $T_{AR}$, $T_{CBR}$, $T_{ScR}$, and $T_{S10R}$ for *HPCHDS* when applied to the Alon dataset.

| Original-Dimension test | $T_{CB}$ | $T_{Sc}$ | $T_{S10}$ | $T_A$ |
|---|---|---|---|---|
| Full-Dim. Lower Crit. Val. | -2.188 | -3.501 | - | -2.001 |
| Full-Dim. Upper Crit. Val. | 3.240 | 3.675 | 3.611 | 1.418 |
| Full-Dim. Empirical test | 2.661 | 1.359 | 1.184 | 1.093 |
| Full-Dim. p-Value | 0.561 | 0.456 | 0.432 | 0.325 |
| Full-Dim. Test Decision | FTR $H_0$ | FTR $H_0$ | FTR $H_0$ | FTR $H_0$ |
| | | | | |
| Reduced-Dimension test | $T_{CBR}$ | $T_{ScR}$ | $T_{S10R}$ | $T_{AR}$ |
| Reduced-Dim. Lower Crit. Val. | -10.179 | -2.894 | - | -3.255 |
| Reduced-Dim. Upper Crit. Val. | 9.828 | 3.385 | 2.369 | 3.265 |
| Reduced-Dim. Empirical test | 2.583 | 3.789 | 3.261 | 3.735 |
| Reduced-Dim. p-Value | 0.785 | 0.041 | 0.013 | 0.032 |
| Reduced-Dim. Test Decision | FTR $H_0$ | Reject $H_0$ | Reject $H_0$ | Reject $H_0$ |

## 4.8 Discussion

In summary, we have derived and applied a new *LDR* method for preserving the information for detecting the differences among $k > 2$ covariance matrices. We have used our proposed *LDR* method to reduce the original-data dimension and have then considered the change in powers for $POW(T_{AR})$, $POW(T_{ScR})$, $POW(T_{S10R})$, and $POW(T_{CBR})$. Using Monte Carlo simulations, we have contrasted the powers of these $(k > 2)$-population *HPCHDS* tests calculated with post-*LDR* data and have determined that $T_{AR}$ and $T_{ScR}$ yielded superior power in most cases examined here for

various parameter and sample size considerations. Finally, we have demonstrated the efficacy of the $HPCHDS$ tests $T_{AR}, T_{ScR}, and T_{S10R}$, calculated with post-$LDR$ data, on a real dataset using permutation tests. We determined that $POW(T_{AR}) \gg POW(T_A)$ and $POW(T_{ScR}) \gg POW(T_{Sc})$ because of the reduction in the number of parameters required to be estimated.

APPENDICES

# APPENDIX A

## Chapter Two Appendix

### A.1   Definitions of Statistics

The following estimators are incorporated into several of the *HPCHDS* test statistics considered in Chapter 2. First,

$$\hat{a}_{1i} := \frac{1}{p(n_i - 1)} \operatorname{tr} \mathbf{V}_i \xrightarrow{P} \frac{\operatorname{tr} \mathbf{\Sigma}_i}{p} \tag{A.1}$$

and

$$\hat{a}_1 := \frac{1}{pN} \operatorname{tr} \mathbf{V} \xrightarrow{P} \frac{\operatorname{tr} \mathbf{\Sigma}}{p}$$

denote estimators of the average of the eigenvalues of $\mathbf{\Sigma}_i$, $i = 1, 2$, and $\mathbf{\Sigma} \in \mathbb{R}_p^>$, respectively, where $\mathbf{\Sigma}$ is the pooled population covariance matrix. Next,

$$\hat{a}_{2i} := \frac{1}{p(n_i - 2)(n_i + 1)} \left\{ \operatorname{tr} \mathbf{V}_i^2 - \frac{1}{n_i - 1} (\operatorname{tr} \mathbf{V}_i)^2 \right\} \xrightarrow{P} \frac{\operatorname{tr} \mathbf{\Sigma}_i^2}{p} \tag{A.2}$$

and

$$\hat{a}_2 := \frac{1}{p(n - 1)(n + 2)} \left\{ \operatorname{tr} \mathbf{V}^2 - \frac{1}{n} (\operatorname{tr} \mathbf{V})^2 \right\} \xrightarrow{P} \frac{\operatorname{tr} \mathbf{\Sigma}^2}{p}. \tag{A.3}$$

Also, for $i = 1, 2$ and $n := N - 2$, we have that

$$\hat{a}_3 := \frac{1}{n(n^2 + 3n + 4)} \left\{ \frac{1}{p} \operatorname{tr} \mathbf{V}^3 - 3n(n + 1) p\hat{a}_2 \hat{a}_1 - np^2 \hat{a}_1^3 \right\} \tag{A.4}$$

and

$$\hat{a}_4 := \frac{1}{n(n^3 + 6n^2 + 21n + 18)} \left( \frac{1}{p} \operatorname{tr} \mathbf{V}^4 - 2pn(2n^2 + 6n + 9) \hat{a}_1 - \tag{A.5} \right.$$

$$\left. 2p^2 n(3n + 2) \hat{a}_1^2 \hat{a}_2 - pn(2n^2 + 5n + 7) \hat{a}_2^2 - np^3 \hat{a}_1^4 \right)$$

denote the consistent estimators for $\operatorname{tr}\boldsymbol{\Sigma}^3/p$ and $\operatorname{tr}\boldsymbol{\Sigma}^4/p$, respectively. Last, from Chaipitak and Chongcharoen (2013), we have

$$\hat{a}_4^* := \frac{(n+1)(n+2)(n+4)(n+6)(n-1)(n-2)(n-3)}{n^5(n^2+n+2)p} \times \tag{A.6}$$
$$\left( \operatorname{tr}\mathbf{S}^4 - \frac{4}{n}\operatorname{tr}\mathbf{S}^2\operatorname{tr}\mathbf{S} - \frac{2n^2+3n-6}{n(n^2+n+2)}\left(\operatorname{tr}\mathbf{S}^2\right)^2 + \right.$$
$$\left. \frac{2(5n+6)}{n(n^2+n+2)}\operatorname{tr}\mathbf{S}^2\left(\operatorname{tr}\mathbf{S}\right)^2 - \frac{5n+6}{n^2(n^2+n+2)}\left(\operatorname{tr}\mathbf{S}\right)^4 \right).$$

### A.2 Simulated Significance Levels for Suggested Asymptotic Critical Values and Power Contrast Tables

For the significance-level simulations, we generated 10,000 independent multivariate normal datasets from $N_p(\mathbf{0}, \boldsymbol{\Sigma}_i)$, where $\boldsymbol{\Sigma}_i$, $i = 1, 2$, are the covariance matrices of interest under $H_0$. The test statistics $T_A$, $T_{Sc}$, $T_{S10}$, $T_C$, and $T_{S14}$ were each calculated from the datasets for each iteration assuming $H_0$ was true. We calculated the simulated significance level using

$$SSL := \frac{\sum_{j=1}^{10,000} I[T_{(*),j} \in RR(T_{(*)})]}{10,000},$$

where $RR(T_{(*)})$ is the rejection region for the test statistic $T_{(*)}$, $T_{(*),j}$ is the test-statistic value corresponding to the $j^{th}$ simulated dataset, and $I[*]$ is the indicator function.

Tables A.1–A.5 display the simulated significance levels of the test statistics $T_A$, $T_C$, $T_{Sc}$, $T_{S10}$, and $T_{S14}$ using the suggested asymptotic critical values associated with the significance level $\alpha = 0.05$. The tables show that the simulated significance levels were often considerably different from the assumed asymptotic significance level of $\alpha = 0.05$. Because of the relatively poor accuracy of the asymptotic significance levels for the considered sample sizes, we use simulated critical values for each *HPCHDS* test statistic for all power-comparison simulations.

Table A.1: Simulated actual significance levels of the tests $T_A$, $T_C$, $T_{Sc}$, $T_{S10}$, and $T_{S14}$ for *HPCHDS* for two constant-times-identity covariance-matrix structures with parameters $\boldsymbol{\Sigma}_1 = \mathbf{I}_p$ and $\boldsymbol{\Sigma}_2 = (1.5)\mathbf{I}_p$ using the suggested asymptotic standard normal critical values corresponding to $\alpha = .05$.

| $p$ | $n_1 = n_2$ | $SSL(T_A)$ | $SSL(T_C)$ | $SSL(T_{Sc})$ | $SSL(T_{S10})$ | $SSL(T_{S14})$ |
|---|---|---|---|---|---|---|
| 20 | 5 | 0.00 | 0.00 | 0.03 | 0.01 | 0.09 |
| | 10 | 0.00 | 0.01 | 0.03 | 0.02 | 0.03 |
| | 15 | 0.00 | 0.02 | 0.03 | 0.01 | 0.02 |
| 40 | 5 | 0.00 | 0.00 | 0.03 | 0.01 | 0.01 |
| | 10 | 0.00 | 0.02 | 0.03 | 0.02 | 0.01 |
| | 15 | 0.00 | 0.03 | 0.03 | 0.02 | 0.01 |
| | 20 | 0.00 | 0.05 | 0.03 | 0.02 | 0.01 |
| 80 | 5 | 0.00 | 0.00 | 0.03 | 0.01 | 0.01 |
| | 10 | 0.00 | 0.03 | 0.03 | 0.03 | 0.02 |
| | 15 | 0.00 | 0.04 | 0.04 | 0.02 | 0.02 |
| | 20 | 0.00 | 0.05 | 0.04 | 0.03 | 0.01 |
| 160 | 5 | 0.00 | 0.00 | 0.03 | 0.01 | 0.01 |
| | 10 | 0.00 | 0.04 | 0.05 | 0.06 | 0.03 |
| | 15 | 0.00 | 0.05 | 0.05 | 0.04 | 0.01 |
| | 20 | 0.00 | 0.05 | 0.05 | 0.04 | 0.01 |

Table A.2: Simulated actual significance levels of the tests $T_A$, $T_C$, $T_{Sc}$, $T_{S10}$, and $T_{S14}$ for *HPCHDS* for two compound-symmetric covariance-matrix structures with parameters $\boldsymbol{\Sigma}_1 = (0.99)\mathbf{I}_p + (0.01)\mathbf{J}_p$ and $\boldsymbol{\Sigma}_2 = (0.95)\mathbf{I}_p + (0.05)\mathbf{J}_p$ using the suggested asymptotic standard normal critical values corresponding to $\alpha = .05$.

| $p$ | $n_1 = n_2$ | $SSL(T_A)$ | $SSL(T_C)$ | $SSL(T_{Sc})$ | $SSL(T_{S10})$ | $SSL(T_{S14})$ |
|---|---|---|---|---|---|---|
| 20 | 5 | 0.00 | 0.01 | 0.03 | 0.11 | 0.09 |
| | 10 | 0.00 | 0.01 | 0.03 | 0.03 | 0.03 |
| | 15 | 0.00 | 0.01 | 0.03 | 0.08 | 0.03 |
| 40 | 5 | 0.00 | 0.00 | 0.03 | 0.08 | 0.10 |
| | 10 | 0.00 | 0.01 | 0.03 | 0.02 | 0.10 |
| | 15 | 0.00 | 0.03 | 0.03 | 0.02 | 0.06 |
| | 20 | 0.00 | 0.04 | 0.03 | 0.02 | 0.03 |
| 80 | 5 | 0.00 | 0.00 | 0.03 | 0.12 | 0.11 |
| | 10 | 0.00 | 0.01 | 0.03 | 0.04 | 0.09 |
| | 15 | 0.00 | 0.03 | 0.03 | 0.01 | 0.03 |
| | 20 | 0.00 | 0.04 | 0.05 | 0.01 | 0.03 |
| 160 | 5 | 0.00 | 0.01 | 0.03 | 0.11 | 0.10 |
| | 10 | 0.00 | 0.03 | 0.03 | 0.06 | 0.06 |
| | 15 | 0.00 | 0.04 | 0.06 | 0.04 | 0.03 |
| | 20 | 0.00 | 0.04 | 0.05 | 0.04 | 0.03 |

Table A.3: Simulated actual significance levels of the tests $T_A$, $T_C$, $T_{Sc}$, $T_{S10}$, and $T_{S14}$ for *HPCHDS* for two autoregressive covariance-matrix structures with parameters $\mathbf{\Sigma}_1 = 0.1^{|i-j|}$ and $\mathbf{\Sigma}_2 = 0.3^{|i-j|}$ using the suggested asymptotic standard normal critical values corresponding to $\alpha = .05$.

| $p$ | $n_1 = n_2$ | $SSL(T_A)$ | $SSL(T_C)$ | $SSL(T_{Sc})$ | $SSL(T_{S10})$ | $SSL(T_{S14})$ |
|---|---|---|---|---|---|---|
| 20 | 5 | 0.00 | 0.01 | 0.03 | 0.11 | 0.01 |
| | 10 | 0.00 | 0.01 | 0.03 | 0.07 | 0.01 |
| | 15 | 0.00 | 0.01 | 0.03 | 0.05 | 0.02 |
| 40 | 5 | 0.00 | 0.02 | 0.03 | 0.10 | 0.01 |
| | 10 | 0.00 | 0.02 | 0.03 | 0.02 | 0.01 |
| | 15 | 0.00 | 0.01 | 0.03 | 0.02 | 0.01 |
| | 20 | 0.00 | 0.04 | 0.03 | 0.03 | 0.01 |
| 80 | 5 | 0.00 | 0.02 | 0.03 | 0.16 | 0.06 |
| | 10 | 0.00 | 0.02 | 0.03 | 0.04 | 0.01 |
| | 15 | 0.00 | 0.02 | 0.03 | 0.03 | 0.00 |
| | 20 | 0.00 | 0.02 | 0.03 | 0.02 | 0.00 |
| 160 | 5 | 0.00 | 0.02 | 0.03 | 0.13 | 0.01 |
| | 10 | 0.00 | 0.02 | 0.05 | 0.05 | 0.01 |
| | 15 | 0.00 | 0.02 | 0.05 | 0.03 | 0.00 |
| | 20 | 0.00 | 0.02 | 0.03 | 0.04 | 0.00 |

Table A.4: Simulated actual significance levels of the tests $T_A$, $T_C$, $T_{Sc}$, $T_{S10}$, and $T_{S14}$ for *HPCHDS* for two heterogeneous autoregressive covariance-matrix structures using the suggested asymptotic standard normal critical values corresponding to $\alpha = .05$.

| $p$ | $n_1 = n_2$ | $SL(T_A)$ | $SL(T_C)$ | $SL(T_{Sc})$ | $SL(T_{S10})$ | $SL(T_{S14})$ |
|---|---|---|---|---|---|---|
| 20 | 5 | 0.00 | 0.01 | 0.12 | 0.03 | 0.08 |
| | 10 | 0.00 | 0.03 | 0.11 | 0.03 | 0.07 |
| | 15 | 0.00 | 0.03 | 0.11 | 0.03 | 0.03 |
| 40 | 5 | 0.00 | 0.02 | 0.11 | 0.02 | 0.09 |
| | 10 | 0.00 | 0.03 | 0.10 | 0.03 | 0.04 |
| | 15 | 0.00 | 0.03 | 0.09 | 0.03 | 0.03 |
| | 20 | 0.00 | 0.03 | 0.09 | 0.03 | 0.01 |
| 80 | 5 | 0.00 | 0.01 | 0.10 | 0.02 | 0.11 |
| | 10 | 0.00 | 0.02 | 0.09 | 0.03 | 0.10 |
| | 15 | 0.00 | 0.02 | 0.05 | 0.02 | 0.09 |
| | 20 | 0.00 | 0.03 | 0.05 | 0.03 | 0.08 |
| 160 | 5 | 0.00 | 0.01 | 0.10 | 0.02 | 0.10 |
| | 10 | 0.00 | 0.03 | 0.08 | 0.03 | 0.08 |
| | 15 | 0.00 | 0.03 | 0.05 | 0.05 | 0.06 |
| | 20 | 0.00 | 0.04 | 0.05 | 0.05 | 0.05 |

Table A.5: Simulated actual significance levels of the tests $T_A$, $T_C$, $T_{Sc}$, $T_{S10}$, and $T_{S14}$ for *HPCHDS* for two unstructured covariance-matrix structures using the suggested asymptotic standard normal critical values corresponding to $\alpha = .05$.

| $p$ | $n_1 = n_2$ | $SSL(T_A)$ | $SSL(T_C)$ | $SSL(T_{Sc})$ | $SSL(T_{S10})$ | $SSL(T_{S14})$ |
|---|---|---|---|---|---|---|
| 20 | 5 | 0.00 | 0.10 | 0.03 | 0.04 | 0.06 |
| | 10 | 0.00 | 0.07 | 0.03 | 0.03 | 0.03 |
| | 15 | 0.00 | 0.07 | 0.03 | 0.03 | 0.04 |
| 40 | 5 | 0.00 | 0.07 | 0.03 | 0.02 | 0.07 |
| | 10 | 0.00 | 0.02 | 0.03 | 0.02 | 0.03 |
| | 15 | 0.00 | 0.02 | 0.03 | 0.03 | 0.01 |
| | 20 | 0.00 | 0.02 | 0.03 | 0.01 | 0.01 |
| 80 | 5 | 0.00 | 0.11 | 0.03 | 0.03 | 0.07 |
| | 10 | 0.00 | 0.05 | 0.03 | 0.04 | 0.10 |
| | 15 | 0.00 | 0.06 | 0.03 | 0.02 | 0.08 |
| | 20 | 0.00 | 0.05 | 0.04 | 0.03 | 0.07 |
| 160 | 5 | 0.00 | 0.10 | 0.03 | 0.04 | 0.08 |
| | 10 | 0.00 | 0.06 | 0.04 | 0.06 | 0.06 |
| | 15 | 0.00 | 0.05 | 0.06 | 0.05 | 0.03 |
| | 20 | 0.00 | 0.04 | 0.05 | 0.04 | 0.03 |

Table A.6: A table contrasting $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing *HPCHDS* on two constant-times-identity population covariance structures with parameters $\boldsymbol{\Sigma}_1 = \mathbf{I}_p$ and $\boldsymbol{\Sigma}_2 = (1.5)\mathbf{I}_p$.

| $p$ | $n_1 = n_2$ | $POW(T_A)$ | $POW(T_C)$ | $POW(T_{Sc})$ | $POW(T_{S10})$ | $POW(T_{S14})$ |
|---|---|---|---|---|---|---|
| 20 | 5 | 0.27 | 0.08 | 0.03 | 0.05 | 0.37 |
| | 10 | 0.71 | 0.41 | 0.08 | 0.01 | 0.32 |
| | 15 | 0.90 | 0.80 | 0.12 | 0.02 | 0.40 |
| 40 | 5 | 0.28 | 0.05 | 0.02 | 0.05 | 0.65 |
| | 10 | 0.76 | 0.54 | 0.04 | 0.00 | 0.74 |
| | 15 | 0.97 | 0.78 | 0.08 | 0.00 | 0.70 |
| | 20 | 1.00 | 0.97 | 0.12 | 0.00 | 0.63 |
| 80 | 5 | 0.30 | 0.07 | 0.01 | 0.02 | 0.93 |
| | 10 | 0.81 | 0.59 | 0.04 | 0.00 | 0.98 |
| | 15 | 0.90 | 0.88 | 0.06 | 0.00 | 0.98 |
| | 20 | 1.00 | 0.99 | 0.10 | 0.00 | 0.99 |
| 160 | 5 | 0.33 | 0.25 | 0.00 | 0.02 | 1.00 |
| | 10 | 0.86 | 0.74 | 0.01 | 0.00 | 1.00 |
| | 15 | 0.99 | 0.98 | 0.05 | 0.00 | 1.00 |
| | 20 | 1.00 | 1.00 | 0.07 | 0.00 | 1.00 |

Table A.7: A table contrasting $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing $HPCHDS$ on two compound-symmetric population covariance structures with parameters $\mathbf{\Sigma}_1 = (.99)\mathbf{I}_p + (0.01)\mathbf{J}_p$ and $\mathbf{\Sigma}_2 = (0.95)\mathbf{I}_p + (0.05)\mathbf{J}_p$.

| $p$ | $n_1 = n_2$ | $POW(T_A)$ | $POW(T_C)$ | $POW(T_{Sc})$ | $POW(T_{S10})$ | $POW(T_{S14})$ |
|---|---|---|---|---|---|---|
| 20 | 5 | 0.07 | 0.04 | 0.06 | 0.05 | 0.07 |
| | 10 | 0.09 | 0.06 | 0.05 | 0.05 | 0.08 |
| | 15 | 0.10 | 0.07 | 0.05 | 0.06 | 0.09 |
| 40 | 5 | 0.07 | 0.06 | 0.07 | 0.05 | 0.08 |
| | 10 | 0.09 | 0.08 | 0.08 | 0.04 | 0.08 |
| | 15 | 0.10 | 0.08 | 0.07 | 0.04 | 0.12 |
| | 20 | 0.13 | 0.09 | 0.07 | 0.04 | 0.14 |
| 80 | 5 | 0.11 | 0.06 | 0.05 | 0.05 | 0.10 |
| | 10 | 0.15 | 0.10 | 0.07 | 0.05 | 0.13 |
| | 15 | 0.17 | 0.12 | 0.10 | 0.04 | 0.19 |
| | 20 | 0.25 | 0.14 | 0.13 | 0.02 | 0.29 |
| 160 | 5 | 0.14 | 0.08 | 0.07 | 0.04 | 0.11 |
| | 10 | 0.30 | 0.15 | 0.12 | 0.02 | 0.31 |
| | 15 | 0.41 | 0.19 | 0.17 | 0.02 | 0.44 |
| | 20 | 0.52 | 0.30 | 0.22 | 0.02 | 0.55 |

Table A.8: A table contrasting $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing $HPCHDS$ on two autoregressive population covariance structures with parameters $\mathbf{\Sigma}_1 = 0.1^{|i-j|}$ and $\mathbf{\Sigma}_2 = 0.3^{|i-j|}$.

| $p$ | $n_1 = n_2$ | $POW(T_A)$ | $POW(T_C)$ | $POW(T_{Sc})$ | $POW(T_{S10})$ | $POW(T_{S14})$ |
|---|---|---|---|---|---|---|
| 20 | 5 | 0.09 | 0.07 | 0.06 | 0.07 | 0.08 |
| | 10 | 0.13 | 0.09 | 0.07 | 0.08 | 0.13 |
| | 15 | 0.17 | 0.12 | 0.06 | 0.08 | 0.20 |
| 40 | 5 | 0.08 | 0.07 | 0.04 | 0.04 | 0.08 |
| | 10 | 0.12 | 0.08 | 0.07 | 0.04 | 0.14 |
| | 15 | 0.17 | 0.09 | 0.07 | 0.06 | 0.25 |
| | 20 | 0.25 | 0.11 | 0.07 | 0.06 | 0.33 |
| 80 | 5 | 0.10 | 0.08 | 0.06 | 0.06 | 0.08 |
| | 10 | 0.16 | 0.09 | 0.06 | 0.05 | 0.16 |
| | 15 | 0.20 | 0.12 | 0.07 | 0.05 | 0.24 |
| | 20 | 0.26 | 0.15 | 0.08 | 0.05 | 0.34 |
| 160 | 5 | 0.10 | 0.07 | 0.05 | 0.05 | 0.07 |
| | 10 | 0.14 | 0.09 | 0.06 | 0.05 | 0.11 |
| | 15 | 0.20 | 0.12 | 0.06 | 0.05 | 0.20 |
| | 20 | 0.28 | 0.18 | 0.07 | 0.06 | 0.31 |

Table A.9: A table contrasting $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing $HPCHDS$ on two unstructured population covariance structures.

| $p$ | $n_1 = n_2$ | $POW(T_A)$ | $POW(T_C)$ | $POW(T_{Sc})$ | $POW(T_{S10})$ | $POW(T_{S14})$ |
|-----|-----|-----|-----|-----|-----|-----|
| 20 | 5 | 0.10 | 0.07 | 0.06 | 0.11 | 0.14 |
| | 10 | 0.17 | 0.08 | 0.09 | 0.23 | 0.26 |
| | 15 | 0.26 | 0.09 | 0.10 | 0.38 | 0.35 |
| 40 | 5 | 0.13 | 0.08 | 0.07 | 0.13 | 0.16 |
| | 10 | 0.21 | 0.08 | 0.11 | 0.27 | 0.30 |
| | 15 | 0.25 | 0.12 | 0.13 | 0.38 | 0.35 |
| | 20 | 0.32 | 0.15 | 0.17 | 0.44 | 0.46 |
| 80 | 5 | 0.15 | 0.09 | 0.08 | 0.21 | 0.21 |
| | 10 | 0.24 | 0.10 | 0.13 | 0.31 | 0.33 |
| | 15 | 0.30 | 0.13 | 0.15 | 0.42 | 0.45 |
| | 20 | 0.39 | 0.18 | 0.19 | 0.50 | 0.59 |
| 160 | 5 | 0.19 | 0.10 | 0.10 | 0.25 | 0.25 |
| | 10 | 0.28 | 0.13 | 0.13 | 0.35 | 0.39 |
| | 15 | 0.40 | 0.16 | 0.17 | 0.44 | 0.56 |
| | 20 | 0.46 | 0.20 | 0.20 | 0.54 | 0.63 |

Figure A.1: Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing *HPCHDS* with constant-times-identity covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{I}_p$ and $\mathbf{\Sigma}_2 = (1.5)\mathbf{I}_p$. The data dimension was increased from $p = 11$ to $p = 159$ and $n_i = 10$, i = 1, 2.

Figure A.2: Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing *HPCHDS* with compound-symmetric covariance matrices with parameters $\mathbf{\Sigma}_1 = 0.99\mathbf{I}_p + 0.01\mathbf{J}_p$ and $\mathbf{\Sigma}_2 = 0.95\mathbf{I}_p + 0.05\mathbf{J}_p$. The data dimension was increased from $p = 11$ to $p = 159$, and the common sample size was $n_i = 10$.

81

Figure A.3: Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing *HPCHDS* with autoregressive covariance matrices with parameters $\boldsymbol{\Sigma}_1 = 0.1^{|i-j|}$ and $\boldsymbol{\Sigma}_2 = 0.3^{|i-j|}$. The data dimension was increased from $p = 11$ to $p = 159$, and the common sample size was $n_i = 10$.

Figure A.4: Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing $HPCHDS$ with heterogeneous autoregressive covariance matrices. The data dimension was increased from $p = 11$ to $p = 159$, and the common sample size was $n_i = 10$.

Figure A.5: Curves for $POW(T_A)$, $POW(T_C)$, $POW(T_{Sc})$, $POW(T_{S10})$, and $POW(T_{S14})$ for testing $HPCHDS$ with unstructured covariance matrices. The data dimension was increased from $p = 11$ to $p = 159$, and the common sample size was $n_i = 10$.

# APPENDIX B

## Chapter Three Appendix

### *B.1   Definitions of Statistics*

The following four estimators of the summands for the *HDSFN* in (3.1) are used in two *HPCHDS* test statistics in Chapter 3. First,

$$\hat{a}_{1i} := \frac{1}{p(n_i - 1)} \operatorname{tr} \mathbf{V}_i \xrightarrow{P} \frac{\operatorname{tr} \boldsymbol{\Sigma}_i}{p}$$

and

$$\hat{a}_1 := \frac{1}{pn} \operatorname{tr} \mathbf{V} \xrightarrow{P} \frac{\operatorname{tr} \boldsymbol{\Sigma}}{p}$$

are consistent estimators of the average of the eigenvalues of the individual population covariance matrices and the pooled covariance matrix, respectively. Next,

$$\hat{a}_{2i} := \frac{1}{p(n_i - 2)(n_i + 1)} \left\{ \operatorname{tr} \mathbf{V}_i^2 - \frac{1}{n_i - 1} (\operatorname{tr} \mathbf{V}_i)^2 \right\} \xrightarrow{P} \frac{\operatorname{tr} \boldsymbol{\Sigma}_i^2}{p} \qquad (B.1)$$

and

$$\hat{a}_2 := \frac{1}{p(n-1)(n+2)} \left\{ \operatorname{tr} \mathbf{V}^2 - \frac{1}{n} (\operatorname{tr} \mathbf{V})^2 \right\} \xrightarrow{P} \frac{\operatorname{tr} \boldsymbol{\Sigma}^2}{p}, \qquad (B.2)$$

where $n = n_1 + n_2 - 2$. Srivastava (2005) has shown that the estimators (B.1) and (B.2) are consistent for estimating the average of the squared elements for the individual population covariance matrices and the pooled covariance matrix, respectively.

For the *HPCHDS* test statistic given in Chaipitak and Chongcharoen (2013), we have

$$\hat{a}_4^* := \frac{(n+1)(n+2)(n+4)(n+6)(n-1)(n-2)(n-3)}{n^5(n^2+n+2)p} \qquad (B.3)$$

$$\left( \operatorname{tr} \mathbf{S}^4 - \frac{4}{n} \operatorname{tr} \mathbf{S}^2 \operatorname{tr} \mathbf{S} - \frac{2n^2 + 3n - 6}{n(n^2+n+2)} \left( \operatorname{tr} \mathbf{S}^2 \right)^2 + \right.$$

$$\left. \frac{2(5n+6)}{n(n^2+n+2)} \operatorname{tr} \mathbf{S}^2 (\operatorname{tr} \mathbf{S})^2 - \frac{5n+6}{n^2(n^2+n+2)} (\operatorname{tr} \mathbf{S})^4 \right).$$

**Lemma B.2.1.** *Let*

$$\mathbf{H} := [\boldsymbol{\Sigma}_2 - \boldsymbol{\Sigma}_1],$$

*where* $\boldsymbol{\Sigma}_i \in \mathbb{R}_p^>$ *for* $i = 1, 2$. *Additionally, let* $SVD(\mathbf{H}) = \mathbf{F}\boldsymbol{\Lambda}\mathbf{G} \in \mathbb{R}_{p \times p}$. *Further, let* $\mathbf{F} \in \mathbb{R}_{p \times q}$ *with* $\mathrm{rank}(\mathbf{F}) = q < p$, *and* $\mathbf{C} = \mathbf{R}\left[\mathbf{I} - \mathbf{F}\mathbf{F}^+\right]$, *where* $\mathbf{R} \in \mathbb{R}_{(p-q) \times p}$ *such that* $\mathrm{rank}\,(\mathbf{C}) = p - q$. *Then, for* $i = 1, 2$,

(*a*) $\mathbf{F}\mathbf{F}^+\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}_i\mathbf{F}\mathbf{F}^+$,

(*b*) $\mathbf{C}\boldsymbol{\Sigma}_2\mathbf{C}^T = \mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^T$,

(*c*) $\left(\mathbf{F}^+\boldsymbol{\Sigma}_i\mathbf{F}^{+T}\right)^{-1} = \mathbf{F}^T\boldsymbol{\Sigma}_i^{-1}\mathbf{F}$,

(*d*) $\mathbf{F}^+\boldsymbol{\Sigma}_i\mathbf{C}^T = \mathbf{0}$.

*Proof.* The proof of (a) – (c) follows from the fact that $\left(\mathbf{I} - \mathbf{F}\mathbf{F}^+\right)\left(\boldsymbol{\Sigma}_i - \boldsymbol{\Sigma}_1\right) = 0$. The proof of (d) is trivial.

**Lemma B.2.2.** *Let* $\mathbf{F} \in \mathbb{R}_{p \times q}$ *such that* $\mathrm{rank}\,(\mathbf{F}) = q$ *and let* $\mathbf{C} = \mathbf{R}\left[\mathbf{I} - \mathbf{F}\mathbf{F}^+\right]$, *where* $\mathbf{R} \in \mathbb{R}_{(p-q) \times p}$ *and* $\mathrm{rank}\,(\mathbf{C}) = p - q$, *and let* $\boldsymbol{\Sigma}_i \in \mathbb{R}_p^>$ *for* $i = 1, 2$, *such that properties* (*b*) *and* (*d*) *of Lemma B.2.1 hold. and let* $\mathbf{W} := \left[\mathbf{F}^{+T} \vdots \mathbf{C}^T\right]^T$ *so that* $\mathrm{rank}\,(\mathbf{W}) = p$. *Then,*

$$\frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} = \frac{\left|\mathbf{F}^+\boldsymbol{\Sigma}_2\mathbf{F}^{+T}\right|}{\left|\mathbf{F}^+\boldsymbol{\Sigma}_1\mathbf{F}^{+T}\right|}.$$

*Proof.* Because $\mathbf{W}$ is full rank, we have that

$$\frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} = \frac{|\mathbf{W}\boldsymbol{\Sigma}_2\mathbf{W}^T|}{|\mathbf{W}\boldsymbol{\Sigma}_1\mathbf{W}^T|}$$

$$= \frac{\begin{vmatrix} \mathbf{F}^+\boldsymbol{\Sigma}_2\mathbf{F}^{+T} & \mathbf{F}^+\boldsymbol{\Sigma}_2\mathbf{C}^T \\ \mathbf{C}\boldsymbol{\Sigma}_2\mathbf{F}^{+T} & \mathbf{C}\boldsymbol{\Sigma}_2\mathbf{C}^T \end{vmatrix}}{\begin{vmatrix} \mathbf{F}^+\boldsymbol{\Sigma}_1\mathbf{F}^{+T} & \mathbf{F}^+\boldsymbol{\Sigma}_1\mathbf{C}^T \\ \mathbf{C}\boldsymbol{\Sigma}_1\mathbf{F}^{+T} & \mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^T \end{vmatrix}}$$

$$= \frac{\begin{vmatrix} \mathbf{F}^+\boldsymbol{\Sigma}_2\mathbf{F}^{+T} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^T \end{vmatrix}}{\begin{vmatrix} \mathbf{F}^+\boldsymbol{\Sigma}_1\mathbf{F}^{+T} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^T \end{vmatrix}}$$

$$= \frac{|\mathbf{F}^+\boldsymbol{\Sigma}_2\mathbf{F}^{+T}|\,|\mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^T|}{|\mathbf{F}^+\boldsymbol{\Sigma}_1\mathbf{F}^{+T}|\,|\mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^T|}$$

$$= \frac{|\mathbf{F}^+\boldsymbol{\Sigma}_2\mathbf{F}^{+T}|}{|\mathbf{F}^+\boldsymbol{\Sigma}_1\mathbf{F}^{+T}|}.$$

**Lemma B.2.3.** *Let* $\mathbf{F} \in \mathbb{R}_{p \times q}$ *such that* $\operatorname{rank}(\mathbf{F}) = q$ *and let* $\mathbf{C} = \mathbf{R}\left[\mathbf{I} - \mathbf{F}\mathbf{F}^+\right]$, *where* $\mathbf{R} \in \mathbb{R}_{(p-q) \times p}$ *and* $\operatorname{rank}(\mathbf{C}) = p - q$, *and let* $\boldsymbol{\Sigma}_{\boldsymbol{i}} \in \mathbb{R}_p^>$ *for* $i = 1, 2$, *such that properties* (a) *and* (c) *of Lemma B.2.1 hold. Then,*

$$\operatorname{tr}\left(\boldsymbol{\Sigma}_2^{-1}\boldsymbol{\Sigma}_1\right) = \operatorname{tr}\left[\left(\mathbf{F}^+\boldsymbol{\Sigma}_2\mathbf{F}^{+T}\right)^{-1}\left(\mathbf{F}^+\boldsymbol{\Sigma}_1\mathbf{F}^{+T}\right)\right].$$

*Proof.* Using properties (a) and (c) of Lemma B.2.1, we have that

$$\operatorname{tr}\left[\left(\mathbf{F}^{+}\mathbf{\Sigma}_{2}\mathbf{F}^{+T}\right)^{-1}\left(\mathbf{F}^{+}\mathbf{\Sigma}_{1}\mathbf{F}^{+T}\right)\right] = \operatorname{tr}\left[\left(\mathbf{F}^{T}\mathbf{\Sigma}_{2}^{-1}\mathbf{F}\right)\left(\mathbf{F}^{+}\mathbf{\Sigma}_{1}\mathbf{F}^{+T}\right)\right]$$

$$= \operatorname{tr}\left[\left(\mathbf{F}^{T}\mathbf{\Sigma}_{2}^{-1}\mathbf{F}\right)\left(\mathbf{F}^{+}\mathbf{\Sigma}_{1}\mathbf{F}^{+T}\right)\right]$$

$$= \operatorname{tr}\left(\mathbf{F}^{T}\mathbf{\Sigma}_{2}^{-1}\mathbf{\Sigma}_{1}\mathbf{F}\mathbf{F}^{+}\mathbf{F}^{+T}\right)$$

$$= tr\left(\mathbf{F}^{T}\mathbf{\Sigma}_{2}^{-1}\mathbf{\Sigma}_{1}\mathbf{F}^{+T}\mathbf{F}^{T}\mathbf{F}^{+T}\right)$$

$$= tr\left(\mathbf{F}^{T}\mathbf{\Sigma}_{2}^{-1}\mathbf{\Sigma}_{1}\mathbf{F}^{+T}\mathbf{F}^{+}\mathbf{F}\right)$$

$$= tr\left(\mathbf{F}^{T}\mathbf{\Sigma}_{2}^{-1}\mathbf{\Sigma}_{1}\mathbf{F}^{+T}\right)$$

$$= tr\left(\mathbf{\Sigma}_{2}^{-1}\mathbf{\Sigma}_{1}\mathbf{F}^{T}\mathbf{F}^{+T}\right)$$

$$= tr\left(\mathbf{\Sigma}_{2}^{-1}\mathbf{\Sigma}_{1}\mathbf{F}^{+}\mathbf{F}\right)$$

$$= tr\left(\mathbf{\Sigma}_{2}^{-1}\mathbf{\Sigma}_{1}\right).$$

Table B.1: A table contrasting POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) and POW($T_A$), POW($T_C$), POW($T_I$), and POW($T_{S10}$) when testing for *HPCHDS* for two constant-times-identity population covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{I}_p$ and $\mathbf{\Sigma}_2 = (1.5)\mathbf{I}_p$.

| $p$ | $n_1 = n_2$ | $q$ | POW($T_{RA}$) | POW($T_{RC}$) | POW($T_{RI}$) | POW($T_{RS10}$) |
|-----|-------------|-----|---------------|---------------|---------------|-----------------|
| 80  | 5           | 5   | 0.00          | 0.99          | 0.02          | 0.00            |
|     |             | 10  | 0.01          | 0.71          | 0.00          | 0.00            |
|     |             | 15  | 0.01          | 0.47          | 0.29          | 0.00            |
|     |             | 20  | 0.02          | 0.29          | 0.24          | 0.01            |
|     | 10          | 10  | 0.00          | 1.00          | 0.01          | 1.00            |
|     |             | 20  | 0.00          | 0.57          | 0.00          | 0.04            |
|     |             | 30  | 0.00          | 0.15          | 0.08          | 0.01            |
|     |             | 40  | 0.00          | 0.03          | 0.03          | 0.02            |
| 160 | 5           | 5   | 0.00          | 0.99          | 0.03          | 1.00            |
|     |             | 10  | 0.00          | 0.94          | 0.00          | 0.00            |
|     |             | 15  | 0.01          | 0.86          | 0.38          | 0.00            |
|     |             | 20  | 0.01          | 0.75          | 0.34          | 0.00            |
|     | 10          | 10  | 0.00          | 1.00          | 0.02          | 1.00            |
|     |             | 20  | 0.00          | 0.99          | 0.35          | 0.03            |
|     |             | 30  | 0.00          | 0.90          | 0.16          | 0.01            |
|     |             | 40  | 0.00          | 0.64          | 0.17          | 0.00            |

| $p$ | $n_1 = n_2$ | $p$ | POW($T_A$) | POW($T_C$) | POW($T_I$) | POW($T_{S10}$) |
|-----|-------------|-----|------------|------------|------------|----------------|
| 80  | 5           | 80  | 0.30       | 0.01       | 0.05       | 0.01           |
|     | 10          | 80  | 0.81       | 0.03       | 0.08       | 0.00           |
| 160 | 5           | 160 | 0.33       | 0.01       | 0.06       | 0.01           |
|     | 10          | 160 | 0.86       | 0.13       | 0.11       | 0.01           |

Table B.2: A table contrasting POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) and POW($T_A$), POW($T_C$), POW($T_I$), and POW($T_{S10}$) when testing for *HPCHDS* for two compound-symmetric population covariance matrices with parameters $\boldsymbol{\Sigma}_1 = (0.99)\mathbf{I}_p + (0.01)\mathbf{J}_p$ and $\boldsymbol{\Sigma}_2 = (0.95)\mathbf{I}_p + (0.05)\mathbf{J}_p$.

| $p$ | $n_1 = n_2$ | $q$ | POW($T_{RA}$) | POW($T_{RC}$) | POW($T_{RI}$) | POW($T_{RS10}$) |
|-----|-------------|-----|---------------|---------------|---------------|-----------------|
| 80 | 5 | 5 | 0.00 | 1.00 | 0.02 | 0.00 |
| | | 10 | 0.00 | 0.91 | 0.00 | 0.00 |
| | | 15 | 0.01 | 0.81 | 0.32 | 0.00 |
| | | 20 | 0.01 | 0.66 | 0.28 | 0.00 |
| | 10 | 10 | 0.00 | 1.00 | 0.02 | 1.00 |
| | | 20 | 0.00 | 0.96 | 0.01 | 0.37 |
| | | 30 | 0.00 | 0.80 | 0.12 | 0.16 |
| | | 40 | 0.01 | 0.54 | 0.08 | 0.08 |
| 160 | 5 | 5 | 0.00 | 1.00 | 0.02 | 0.99 |
| | | 10 | 0.00 | 0.98 | 0.00 | 0.00 |
| | | 15 | 0.00 | 0.97 | 0.44 | 0.00 |
| | | 20 | 0.00 | 0.94 | 0.38 | 0.00 |
| | 10 | 10 | 0.00 | 1.00 | 0.02 | 1.00 |
| | | 20 | 0.00 | 1.00 | 0.25 | 0.19 |
| | | 30 | 0.00 | 0.99 | 0.21 | 0.09 |
| | | 40 | 0.00 | 0.96 | 0.23 | 0.04 |

| $p$ | $n_1 = n_2$ | $p$ | POW($T_A$) | POW($T_C$) | POW($T_I$) | POW($T_{S10}$) |
|-----|-------------|-----|------------|------------|------------|----------------|
| 80 | 5 | 80 | 0.11 | 0.04 | 0.06 | 0.03 |
| | 10 | 80 | 0.15 | 0.03 | 0.06 | 0.04 |
| 160 | 5 | 160 | 0.15 | 0.02 | 0.06 | 0.04 |
| | 10 | 160 | 0.30 | 0.02 | 0.05 | 0.02 |

Table B.3: A table contrasting POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) and POW($T_A$), POW($T_C$), POW($T_I$), and POW($T_{S10}$) when testing for *HPCHDS* for two heterogeneous autoregressive population covariance matrices.

| $p$ | $n_1 = n_2$ | $q$ | POW($T_{RA}$) | POW($T_{RC}$) | POW($T_{RI}$) | POW($T_{RS10}$) |
|---|---|---|---|---|---|---|
| 80 | 5 | 5 | 0.01 | 0.95 | 0.01 | 0.00 |
| | | 10 | 0.03 | 0.46 | 0.01 | 0.01 |
| | | 15 | 0.05 | 0.36 | 0.25 | 0.01 |
| | | 20 | 0.06 | 0.25 | 0.19 | 0.02 |
| | 10 | 10 | 0.02 | 0.79 | 0.01 | 0.64 |
| | | 20 | 0.18 | 0.20 | 0.01 | 0.12 |
| | | 30 | 0.25 | 0.12 | 0.06 | 0.16 |
| | | 40 | 0.34 | 0.06 | 0.04 | 0.20 |
| 160 | 5 | 5 | 0.00 | 0.99 | 0.02 | 0.99 |
| | | 10 | 0.03 | 0.64 | 0.01 | 0.00 |
| | | 15 | 0.04 | 0.56 | 0.35 | 0.00 |
| | | 20 | 0.04 | 0.49 | 0.30 | 0.00 |
| | 10 | 10 | 0.00 | 0.94 | 0.01 | 0.83 |
| | | 20 | 0.14 | 0.38 | 0.24 | 0.12 |
| | | 30 | 0.18 | 0.25 | 0.13 | 0.11 |
| | | 40 | 0.23 | 0.22 | 0.12 | 0.21 |

| $p$ | $n_1 = n_2$ | $p$ | POW($T_A$) | POW($T_C$) | POW($T_I$) | POW($T_{S10}$) |
|---|---|---|---|---|---|---|
| 80 | 5 | 80 | 0.24 | 0.03 | 0.05 | 0.11 |
| | 10 | 80 | 0.52 | 0.01 | 0.06 | 0.33 |
| 160 | 5 | 160 | 0.34 | 0.03 | 0.07 | 0.12 |
| | 10 | 160 | 0.67 | 0.01 | 0.07 | 0.52 |

Table B.4: A table contrasting POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) and POW($T_A$), POW($T_C$), POW($T_I$), and POW($T_{S10}$) when testing for *HPCHDS* for two unstructured population covariance matrices.

| $p$ | $n_1 = n_2$ | $q$ | POW($T_{RA}$) | POW($T_{RC}$) | POW($T_{RI}$) | POW($T_{RS10}$) |
|-----|-------------|-----|---------------|---------------|---------------|-----------------|
| 80  | 5           | 5   | 0.01          | 0.27          | 0.00          | 0.07            |
|     |             | 10  | 0.08          | 0.10          | 0.00          | 0.09            |
|     |             | 15  | 0.09          | 0.09          | 0.18          | 0.09            |
|     |             | 20  | 0.09          | 0.08          | 0.14          | 0.10            |
|     | 10          | 10  | 0.16          | 0.07          | 0.01          | 0.07            |
|     |             | 20  | 0.19          | 0.02          | 0.01          | 0.12            |
|     |             | 30  | 0.20          | 0.02          | 0.08          | 0.13            |
|     |             | 40  | 0.21          | 0.02          | 0.05          | 0.14            |
| 160 | 5           | 5   | 0.05          | 0.30          | 0.00          | 0.08            |
|     |             | 10  | 0.10          | 0.10          | 0.00          | 0.10            |
|     |             | 15  | 0.10          | 0.10          | 0.23          | 0.10            |
|     |             | 20  | 0.10          | 0.09          | 0.24          | 0.10            |
|     | 10          | 10  | 0.17          | 0.07          | 0.01          | 0.06            |
|     |             | 20  | 0.21          | 0.02          | 0.23          | 0.13            |
|     |             | 30  | 0.22          | 0.02          | 0.06          | 0.13            |
|     |             | 40  | 0.22          | 0.02          | 0.05          | 0.13            |

| $p$ | $n_1 = n_2$ | $p$ | POW($T_A$) | POW($T_C$) | POW($T_I$) | POW($T_{S10}$) |
|-----|-------------|-----|------------|------------|------------|----------------|
| 80  | 5           | 80  | 0.15       | 0.04       | 0.04       | 0.12           |
|     | 10          | 80  | 0.24       | 0.02       | 0.06       | 0.15           |
| 160 | 5           | 160 | 0.19       | 0.04       | 0.04       | 0.17           |
|     | 10          | 160 | 0.28       | 0.01       | 0.06       | 0.17           |

Figure B.1: Reduced-dimension power curves and power-difference curves for POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) for two constant-times-identity covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{I}_p$ and $\mathbf{\Sigma}_2 = (1.5)\mathbf{I}_p$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 5, i = 1, 2, 3$.

Figure B.2: Reduced-dimension power curves and power-difference curves for $\mathrm{POW}(T_{RA})$, $\mathrm{POW}(T_{RC})$, $\mathrm{POW}(T_{RI})$, and $\mathrm{POW}(T_{RS10})$ for two compound-symmetric covariance matrices with parameters $\boldsymbol{\Sigma}_1 = (0.99)\mathbf{I}_p + (0.01)\mathbf{J}_p$ and $\boldsymbol{\Sigma}_2 = (0.95)\mathbf{I}_p + (0.05)\mathbf{J}_p$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 5, i = 1, 2, 3$.
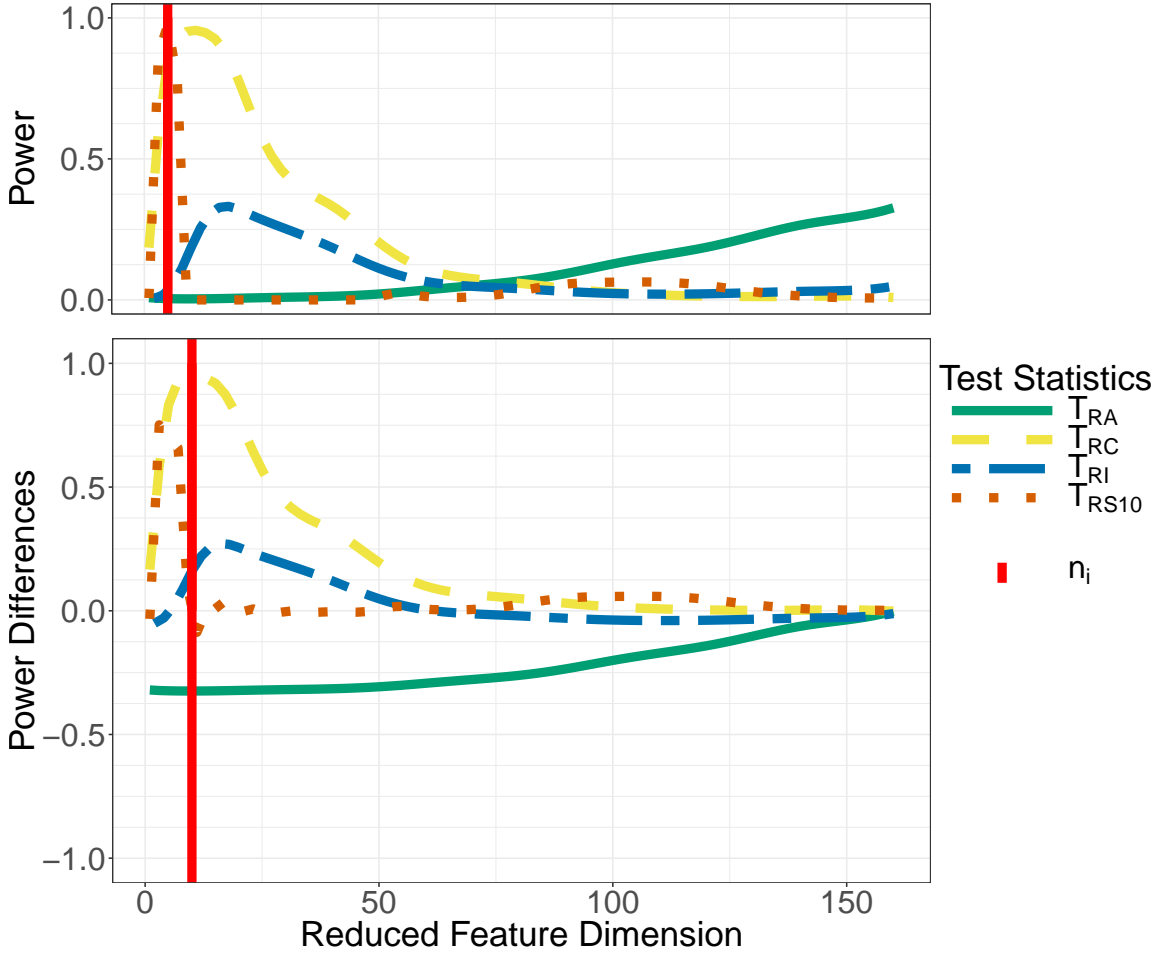
Figure B.3: Reduced-dimension power curves and power-difference curves for $\text{POW}(T_{RA})$, $\text{POW}(T_{RC})$, $\text{POW}(T_{RI})$, and $\text{POW}(T_{RS10})$ for two autoregressive covariance matrices with parameters $\Sigma_1 = 0.1^{|i-j|}$ and $\Sigma_2 = 0.3^{|i-j|}$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 5, i = 1, 2, 3$.

Figure B.4: Reduced-dimension power curves and power-difference curves for POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) for two heterogeneous autoregressive covariance matrices. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 5, i = 1, 2, 3$.

Figure B.5: Reduced-dimension power curves and power-difference curves for POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) for two unstructured covariance matrices. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 5, i = 1, 2, 3$.
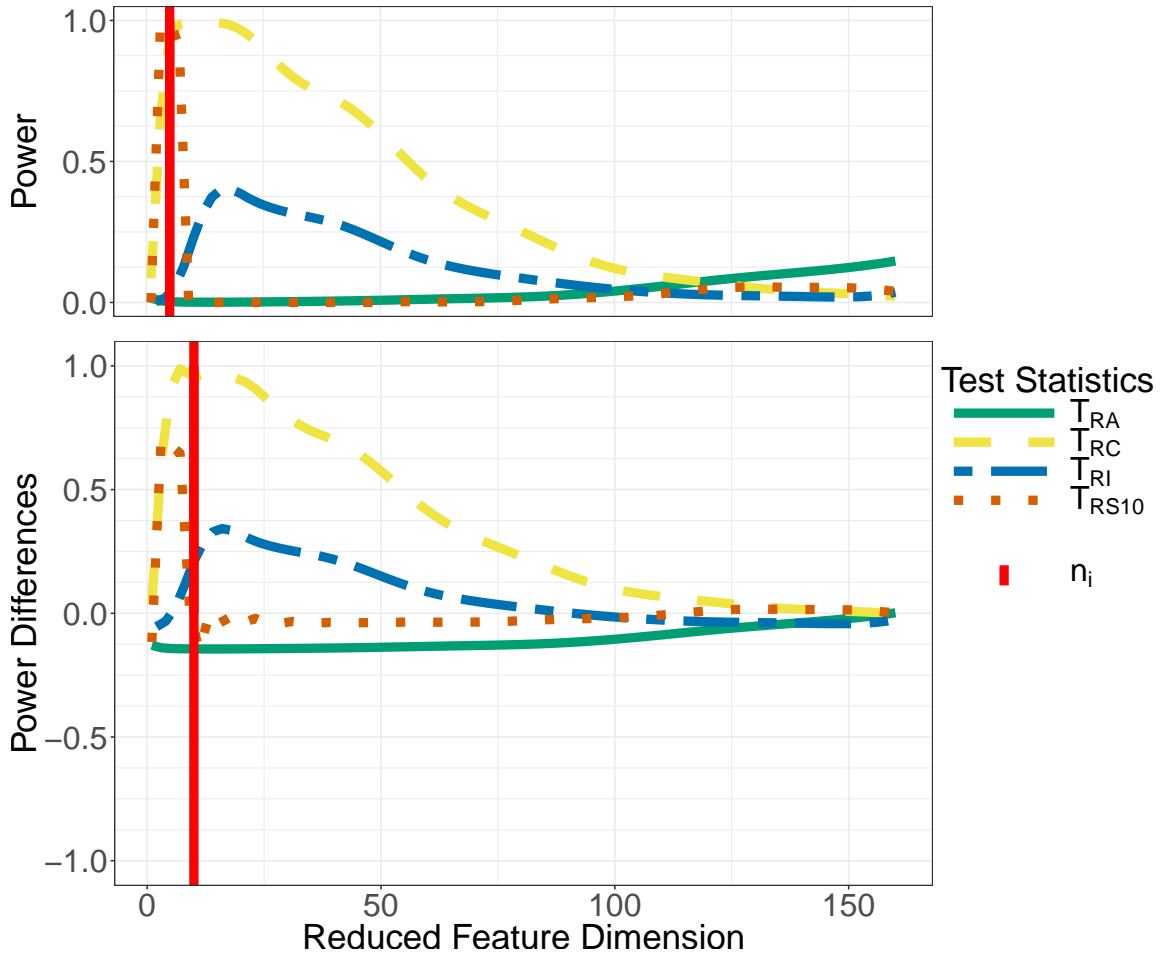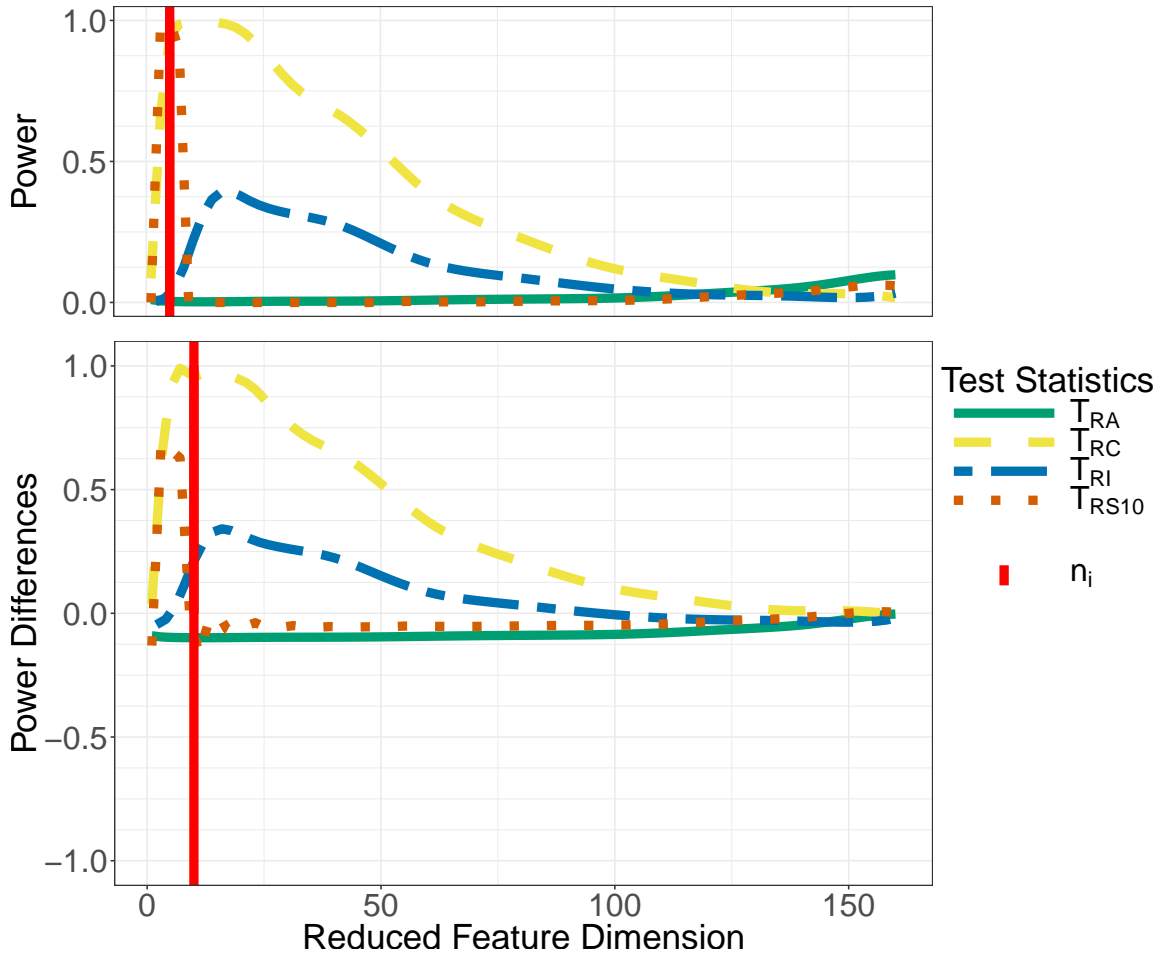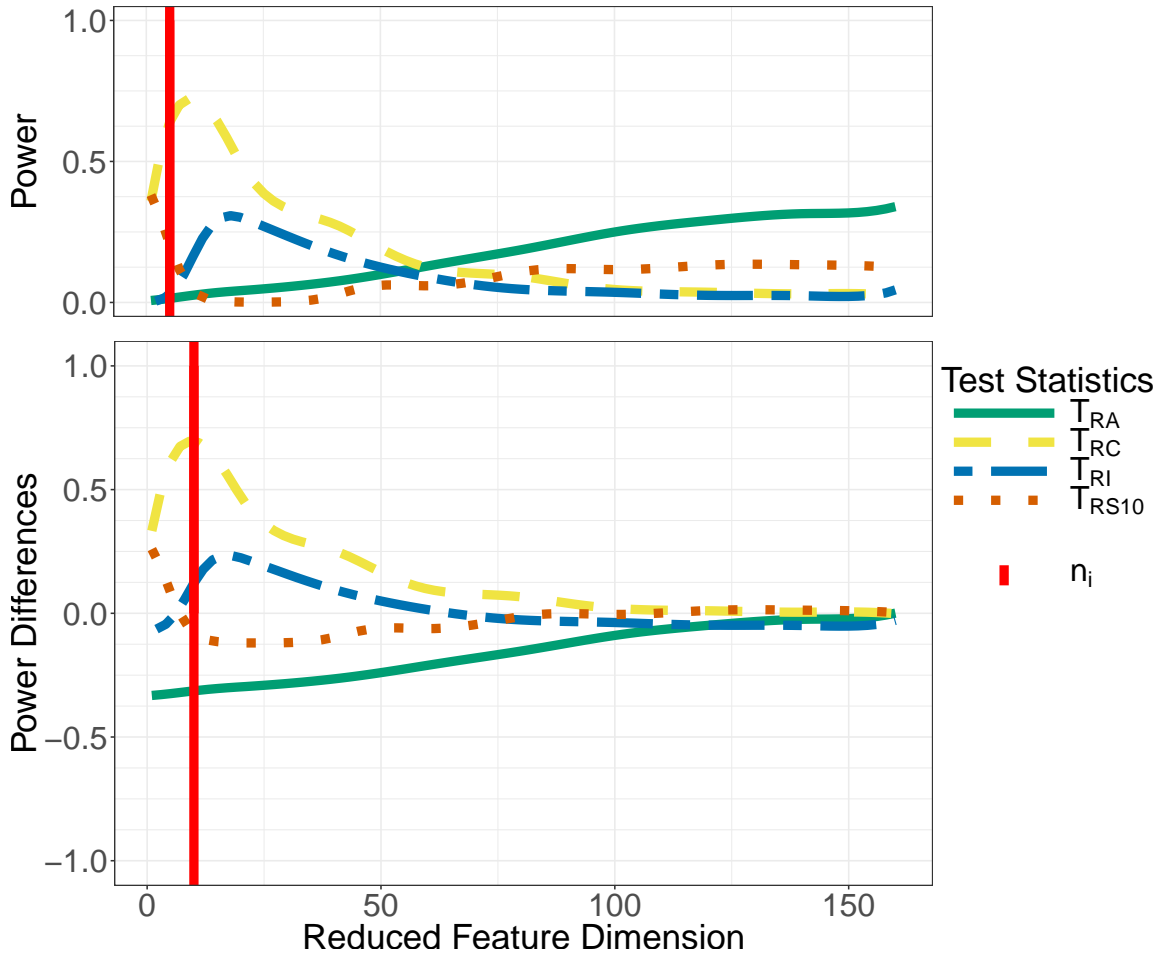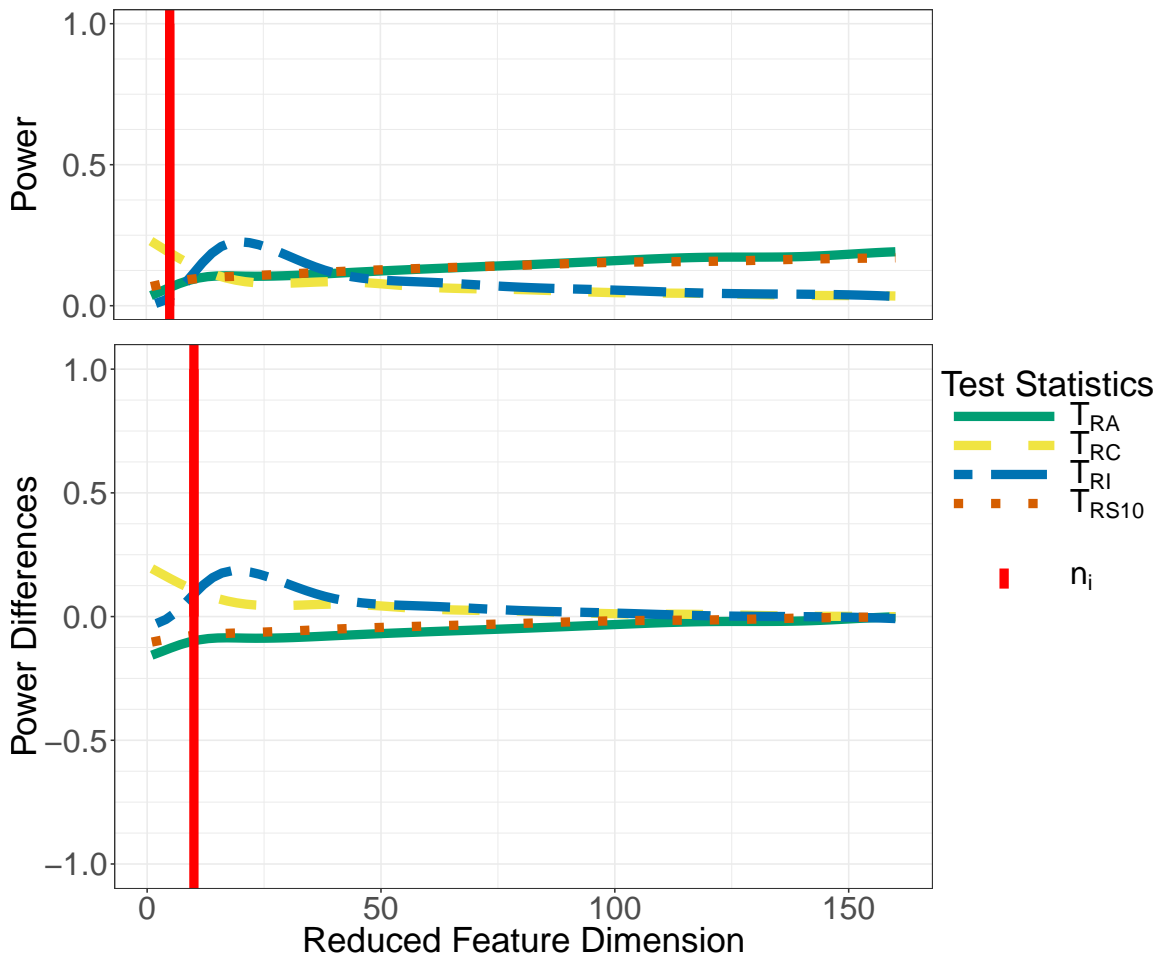
APPENDIX C

Chapter Four Appendix

## C.1 Definitions of Statistics

The following estimators of six parameters of certain summands of the *HDSFN* in (4.1) are used in several of the *HPCHDS* test statistics considered here. First, for $i = 1, 2, \ldots, k$,

$$\hat{a}_{1i} := \frac{1}{p(n_i - 1)} \operatorname{tr} \mathbf{V}_i \xrightarrow{P} \frac{\operatorname{tr} \mathbf{\Sigma}_i}{p} \tag{C.1}$$

and

$$\hat{a}_1 := \frac{1}{pn} \operatorname{tr} \mathbf{V} \xrightarrow{P} \frac{\operatorname{tr} \mathbf{\Sigma}}{p} \tag{C.2}$$

denote estimators of the average of the eigenvalues of the $i^{th}$ sample covariance matrices and the pooled covariance matrix, respectively. The consistency of these estimators has been proven by Srivastava (2005). Next,

$$\hat{a}_{2i} := \frac{1}{p(n_i - 2)(n_i + 1)} \left\{ \operatorname{tr} \mathbf{V}_i^2 - \frac{1}{n_i - 1} (\operatorname{tr} \mathbf{V}_i)^2 \right\} \xrightarrow{P} \frac{\operatorname{tr} \mathbf{\Sigma}_i^2}{p}, \tag{C.3}$$

and

$$\hat{a}_2 := \frac{1}{p(n - 1)(n + 2)} \left\{ \operatorname{tr} \mathbf{V}^2 - \frac{1}{n} (\operatorname{tr} \mathbf{V})^2 \right\} \xrightarrow{P} \frac{\operatorname{tr} \mathbf{\Sigma}^2}{p} \tag{C.4}$$

with

$$n = \sum_{i=1}^{k} (n_i - 1),$$

where $n_i$ is the sample size for population $i, i = 1, 2, \ldots, k$. The consistency of these estimators has been shown by Srivastava and Yanagihara (2010). Next, for $i = 1, 2 \ldots, k$,

$$\hat{a}_3 := \frac{1}{n(n^2 + 3n + 4)} \left\{ \frac{1}{p} \operatorname{tr} \mathbf{V}^3 - 3n(n + 1) p \hat{a}_2 \hat{a}_1 - np^2 \hat{a}_1^3 \right\} \tag{C.5}$$

and

$$\hat{a}_4 := \frac{1}{n\left(n^3 + 6n^2 + 21n + 18\right)} \left(\frac{1}{p} \operatorname{tr} \mathbf{V}^4 - 2pn\left(2n^2 + 6n + 9\right)\hat{a}_1 - \right. \tag{C.6}$$

$$\left. 2p^2 n\left(3n + 2\right)\hat{a}_1^2 \hat{a}_2 - pn\left(2n^2 + 5n + 7\right)\hat{a}_2^2 - np^3 \hat{a}_1^4\right)$$

denote the consistent estimators for $\operatorname{tr}\mathbf{\Sigma}^3/p$ and $\operatorname{tr}\mathbf{\Sigma}^4/p$, respectively. Also, a statistic based on a form from Chaipitak and Chongcharoen (2013) for $i = 1, 2\ldots, k$, is

$$\hat{a}_4^* := \frac{(n+1)(n+2)(n+4)(n+6)(n-1)(n-2)(n-3)}{n^5(n^2+n+2)p} \tag{C.7}$$

$$\left(\operatorname{tr}\mathbf{S}^4 - \frac{4}{n}\operatorname{tr}\mathbf{S}^2 \operatorname{tr}\mathbf{S} - \frac{2n^2 + 3n - 6}{n\left(n^2+n+2\right)}\left(\operatorname{tr}\mathbf{S}^2\right)^2 + \right.$$

$$\left. \frac{2\left(5n+6\right)}{n\left(n^2+n+2\right)}\operatorname{tr}\mathbf{S}^2\left(\operatorname{tr}\mathbf{S}\right)^2 - \frac{5n+6}{n^2\left(n^2+n+2\right)}\left(\operatorname{tr}\mathbf{S}\right)^4\right).$$

### C.2    Lemmas for the Theorem

**Lemma C.2.1.** *Let*

$$\mathbf{H} := [\mathbf{\Sigma}_2 - \mathbf{\Sigma}_1 \vdots \mathbf{\Sigma}_3 - \mathbf{\Sigma}_1 \vdots \ldots \vdots \mathbf{\Sigma}_k - \mathbf{\Sigma}_1],$$

*where $\mathbf{\Sigma}_i \in \mathbb{R}_p^S$ for $i = 1, 2, \ldots, k$. Additionally, let $SVD(\mathbf{H}) = \mathbf{F}\mathbf{\Lambda}\mathbf{G} \in \mathbb{R}_{p\times(k-1)p}$, and let $\mathbf{F} \in \mathbb{R}_{p\times q}$ with $\operatorname{rank}(\mathbf{F}) = q < p$, and $\mathbf{C} = \mathbf{R}\left[\mathbf{I} - \mathbf{F}\mathbf{F}^+\right]$, where $\mathbf{R} \in \mathbb{R}_{(p-q)\times p}$ such that $\operatorname{rank}(\mathbf{C}) = p - q$. Then, for $1 \leq i, j \leq k, i \neq j$, we have*

*(a)* $\mathbf{C}\mathbf{\Sigma}_i\mathbf{C}^T = \mathbf{C}\mathbf{\Sigma}_j\mathbf{C}^T$

*(b)* $\mathbf{F}^+\mathbf{\Sigma}_i\mathbf{C}^T = \mathbf{0}$.

*Proof.* The proof of follows from the facts that from the fact that $\left(\mathbf{I} - \mathbf{F}\mathbf{F}^+\right)\left(\mathbf{\Sigma}_i - \mathbf{\Sigma}_1\right) = 0$ and that $\mathbf{\Sigma}_i \perp \left[\mathbf{I} - \mathbf{F}\mathbf{F}^+\right]$.

**Lemma C.2.2.** *Let $\mathbf{F} \in \mathbb{R}_{p\times q}$ where $\operatorname{rank}(\mathbf{F}) = q$, and let $\mathbf{C} = \mathbf{R}\left[\mathbf{I} - \mathbf{F}\mathbf{F}^+\right]$, where $\mathbf{R} \in \mathbb{R}_{(p-q)\times p}$ and $\operatorname{rank}(\mathbf{C}) = p - q$, and let $\mathbf{\Sigma_i} \in \mathbb{R}_p^>$ for $i = 1, 2, \ldots, k$, such that*

*Lemma C.2.1 holds. Also, let* $\mathbf{W} := \left[\mathbf{F}^{+T} \vdots \mathbf{C}^T\right]^T$ *so that* $\operatorname{rank}(\mathbf{W}) = p$. *Then, for* $1 \leq i, j \leq k,, i \neq j$, *we have* $\det\left(\boldsymbol{\Sigma_i}\boldsymbol{\Sigma_j}^{-1}\right) = \det\left[\left(\mathbf{F}^+\boldsymbol{\Sigma}_i\mathbf{F}^{+T}\right)\left(\mathbf{F}^+\boldsymbol{\Sigma}_j^{-1}\mathbf{F}^{+T}\right)^{-1}\right]$.

*Proof.* Using Lemma C.2.1 and the fact that $\operatorname{rank}(\mathbf{W}) = p$, we have

$$
\det\left(\boldsymbol{\Sigma}_i\boldsymbol{\Sigma}_j^{-1}\right) = \det\left([\mathbf{W}\boldsymbol{\Sigma}_i\mathbf{W}^T][\mathbf{W}\boldsymbol{\Sigma}_j\mathbf{W}]^{-1}\right)
$$

$$
= \det\left(\begin{bmatrix} \mathbf{F}^+\boldsymbol{\Sigma}_i\mathbf{F}^{+T} & \mathbf{F}^+\boldsymbol{\Sigma}_i\mathbf{C} \\ \mathbf{C}\boldsymbol{\Sigma}_i\mathbf{F}^+ & \mathbf{C}\boldsymbol{\Sigma}_i\mathbf{C}^T \end{bmatrix}\right)\left(\begin{bmatrix} \mathbf{F}^+\boldsymbol{\Sigma}_j\mathbf{F}^{+T} & \mathbf{F}^+\boldsymbol{\Sigma}_j\mathbf{C} \\ \mathbf{C}\boldsymbol{\Sigma}_j\mathbf{F}^+ & \mathbf{C}\boldsymbol{\Sigma}_j\mathbf{C}^T \end{bmatrix}^{-1}\right)
$$

$$
= \det\left(\begin{bmatrix} \mathbf{F}^+\boldsymbol{\Sigma}_i\mathbf{F}^{+T} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^T \end{bmatrix}\right)\left(\begin{bmatrix} \mathbf{F}^+\boldsymbol{\Sigma}_j\mathbf{F}^{+T} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^T \end{bmatrix}^{-1}\right)
$$

$$
= \det\left(\begin{bmatrix} \left(\mathbf{F}^+\boldsymbol{\Sigma}_i\mathbf{F}^{+T}\right)\left(\mathbf{F}^+\boldsymbol{\Sigma}_j\mathbf{F}^{+T}\right)^{-1} & \mathbf{0} \\ \mathbf{0} & \left(\mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^T\right)\left(\mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^T\right)^{-1} \end{bmatrix}\right)
$$

$$
= \det\left(\begin{bmatrix} \left(\mathbf{F}^+\boldsymbol{\Sigma}_i\mathbf{F}^{+T}\right)\left(\mathbf{F}^+\boldsymbol{\Sigma}_j^{-1}\mathbf{F}^{+T}\right)^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}\right)
$$

$$
= \det\left[\left(\mathbf{F}^+\boldsymbol{\Sigma}_i\mathbf{F}^{+T}\right)\left(\mathbf{F}^+\boldsymbol{\Sigma}_j^{-1}\mathbf{F}^{+T}\right)^{-1}\right].
$$

Table C.1: A table contrasting $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ with $POW(T_A)$, $POW(T_{CB})$, $POW(T_{Sc})$, and $POW(T_{S10})$ for testing *HPCHDS* for three constant-times-identity population covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = \mathbf{I}_p$ and $\mathbf{\Sigma}_2 = (1.5)\mathbf{I}_p$.

| $p$ | $n_1 = n_2 = n_3$ | $q$ | $POW(T_{AR})$ | $POW(T_{CBR})$ | $POW(T_{ScR})$ | $POW(T_{S10R})$ |
|---|---|---|---|---|---|---|
| 80 | 5 | 5 | 0.91 | 0.08 | 0.04 | 0.93 |
| | | 10 | 0.95 | 0.06 | 0.99 | 0.95 |
| | | 15 | 0.84 | 0.01 | 1.00 | 0.02 |
| | | 20 | 0.80 | 0.01 | 1.00 | 0.03 |
| | 10 | 10 | 1.00 | 0.00 | 0.20 | 0.99 |
| | | 20 | 1.00 | 0.01 | 0.86 | 0.86 |
| | | 30 | 1.00 | 0.06 | 1.00 | 0.00 |
| | | 40 | 0.99 | 0.20 | 1.00 | 0.00 |
| 160 | 5 | 5 | 0.88 | 0.12 | 0.07 | 0.96 |
| | | 10 | 1.00 | 0.10 | 1.00 | 0.77 |
| | | 15 | 0.86 | 0.00 | 1.00 | 0.03 |
| | | 20 | 0.85 | 0.00 | 1.00 | 0.04 |
| | 10 | 10 | 1.00 | 0.03 | 0.40 | 0.99 |
| | | 20 | 1.00 | 0.01 | 1.00 | 1.00 |
| | | 30 | 1.00 | 0.00 | 1.00 | 0.00 |
| | | 40 | 1.00 | 0.00 | 1.00 | 0.00 |

| $p$ | $n_1 = n_2 = n_3$ | $p$ | $POW(T_A)$ | $POW(T_{CB})$ | $POW(T_{Sc})$ | $POW(T_{S10})$ |
|---|---|---|---|---|---|---|
| 80 | 5 | 80 | 0.24 | 0.42 | 0.00 | 0.01 |
| | 10 | 80 | 0.72 | 0.65 | 0.01 | 0.00 |
| 160 | 5 | 160 | 0.27 | 0.02 | 0.00 | 0.00 |
| | 10 | 160 | 0.78 | 0.71 | 0.01 | 0.01 |

Table C.2: A table contrasting $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ statistics with $POW(T_A)$, $POW(T_{CB})$, $POW(T_{Sc})$, and $POW(T_{S10})$ for testing $HPCHDS$ for three compound-symmetric population covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = 0.99\mathbf{I}_p + 0.01\mathbf{J}_p$ and $\mathbf{\Sigma}_2 = 0.95\mathbf{I}_p + 0.05\mathbf{J}_p$.

| $p$ | $n_1 = n_2 = n_3$ | $q$ | $POW(T_{AR})$ | $POW(T_{CBR})$ | $POW(T_{ScR})$ | $POW(T_{S10R})$ |
|---|---|---|---|---|---|---|
| 80 | 5 | 5 | 0.89 | 0.08 | 0.01 | 0.85 |
| | | 10 | 0.99 | 0.35 | 0.97 | 0.87 |
| | | 15 | 0.66 | 0.02 | 0.99 | 0.09 |
| | | 20 | 0.56 | 0.03 | 1.00 | 0.14 |
| | 10 | 10 | 1.00 | 0.00 | 0.00 | 0.93 |
| | | 20 | 1.00 | 0.11 | 0.01 | 0.61 |
| | | 30 | 0.97 | 0.01 | 0.99 | 0.00 |
| | | 40 | 0.86 | 0.02 | 0.99 | 0.01 |
| 160 | 5 | 5 | 0.77 | 0.38 | 0.00 | 0.92 |
| | | 10 | 0.99 | 0.46 | 0.99 | 0.18 |
| | | 15 | 0.73 | 0.02 | 1.00 | 0.13 |
| | | 20 | 0.71 | 0.02 | 1.00 | 0.24 |
| | 10 | 10 | 1.00 | 0.00 | 0.00 | 0.99 |
| | | 20 | 1.00 | 0.62 | 0.11 | 0.99 |
| | | 30 | 1.00 | 0.00 | 1.00 | 0.02 |
| | | 40 | 0.99 | 0.00 | 1.00 | 0.04 |

| $p$ | $n_1 = n_2 = n_3$ | $p$ | $POW(T_A)$ | $POW(T_{CB})$ | $POW(T_{Sc})$ | $POW(T_{S10})$ |
|---|---|---|---|---|---|---|
| 80 | 5 | 80 | 0.11 | 0.08 | 0.05 | 0.05 |
| | 10 | 80 | 0.12 | 0.12 | 0.02 | 0.04 |
| 160 | 5 | 160 | 0.12 | 0.10 | 0.04 | 0.05 |
| | 10 | 160 | 0.26 | 0.19 | 0.01 | 0.03 |

Table C.3: A table contrasting $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ statistics with $POW(T_A)$, $POW(T_{CB})$, $POW(T_{Sc})$, and $POW(T_{S10})$ for testing $HPCHDS$ for three heterogeneous autoregressive population covariance matrices with parameters $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_3 = 0.1^{|i-j|}$ and $\boldsymbol{\Sigma}_2 = 0.3^{|i-j|}$.

| $p$ | $n_1 = n_2 = n_3$ | $q$ | $POW(T_{AR})$ | $POW(T_{CBR})$ | $POW(T_{ScR})$ | $POW(T_{S10R})$ |
|---|---|---|---|---|---|---|
| 80 | 5 | 5 | 0.84 | 0.05 | 0.01 | 0.64 |
|  |  | 10 | 0.87 | 0.33 | 0.34 | 0.46 |
|  |  | 15 | 0.53 | 0.18 | 0.74 | 0.07 |
|  |  | 20 | 0.49 | 0.21 | 0.65 | 0.08 |
|  | 10 | 10 | 1.00 | 0.01 | 0.00 | 0.74 |
|  |  | 20 | 0.97 | 0.24 | 0.02 | 0.42 |
|  |  | 30 | 0.72 | 0.30 | 0.26 | 0.30 |
|  |  | 40 | 0.64 | 0.32 | 0.12 | 0.40 |
| 160 | 5 | 5 | 0.89 | 0.14 | 0.00 | 0.83 |
|  |  | 10 | 0.98 | 0.45 | 0.62 | 0.76 |
|  |  | 15 | 0.68 | 0.15 | 0.93 | 0.09 |
|  |  | 20 | 0.66 | 0.17 | 0.91 | 0.14 |
|  | 10 | 10 | 1.00 | 0.01 | 0.00 | 0.95 |
|  |  | 20 | 1.00 | 0.24 | 0.23 | 0.62 |
|  |  | 30 | 0.96 | 0.28 | 0.61 | 0.18 |
|  |  | 40 | 0.92 | 0.25 | 0.54 | 0.26 |

| $p$ | $n_1 = n_2 = n_3$ | $p$ | $POW(T_A)$ | $POW(T_{CB})$ | $POW(T_{Sc})$ | $POW(T_{S10})$ |
|---|---|---|---|---|---|---|
| 80 | 5 | 80 | 0.24 | 0.15 | 0.02 | 0.14 |
|  | 10 | 80 | 0.48 | 0.33 | 0.01 | 0.51 |
| 160 | 5 | 160 | 0.32 | 0.25 | 0.01 | 0.24 |
|  | 10 | 160 | 0.61 | 0.49 | 0.01 | 0.67 |

Table C.4: A table contrasting $POW(T_{AR})$, $POW(T_{CBR})$, $POW(T_{ScR})$, and $POW(T_{S10R})$ statistics with $POW(T_A)$, $POW(T_{CB})$, $POW(T_{Sc})$, and $POW(T_{S10})$ for testing $HPCHDS$ for three unstructured population covariance matrices.

| $p$ | $n_1 = n_2 = n_3$ | $q$ | $POW(T_{AR})$ | $POW(T_{CBR})$ | $POW(T_{ScR})$ | $POW(T_{S10R})$ |
|---|---|---|---|---|---|---|
| 80 | 5 | 5 | 0.47 | 0.15 | 0.01 | 0.05 |
| | | 10 | 0.42 | 0.15 | 0.02 | 0.09 |
| | | 15 | 0.40 | 0.14 | 0.02 | 0.04 |
| | | 20 | 0.39 | 0.14 | 0.02 | 0.07 |
| | 10 | 10 | 0.69 | 0.20 | 0.00 | 0.28 |
| | | 20 | 0.65 | 0.19 | 0.00 | 0.39 |
| | | 30 | 0.63 | 0.18 | 0.00 | 0.39 |
| | | 40 | 0.62 | 0.17 | 0.00 | 0.58 |
| 160 | 5 | 5 | 0.53 | 0.16 | 0.00 | 0.04 |
| | | 10 | 0.49 | 0.18 | 0.01 | 0.06 |
| | | 15 | 0.47 | 0.15 | 0.02 | 0.01 |
| | | 20 | 0.45 | 0.15 | 0.02 | 0.02 |
| | 10 | 10 | 0.75 | 0.21 | 0.00 | 0.11 |
| | | 20 | 0.72 | 0.24 | 0.00 | 0.12 |
| | | 30 | 0.70 | 0.22 | 0.00 | 0.11 |
| | | 40 | 0.69 | 0.22 | 0.00 | 0.19 |

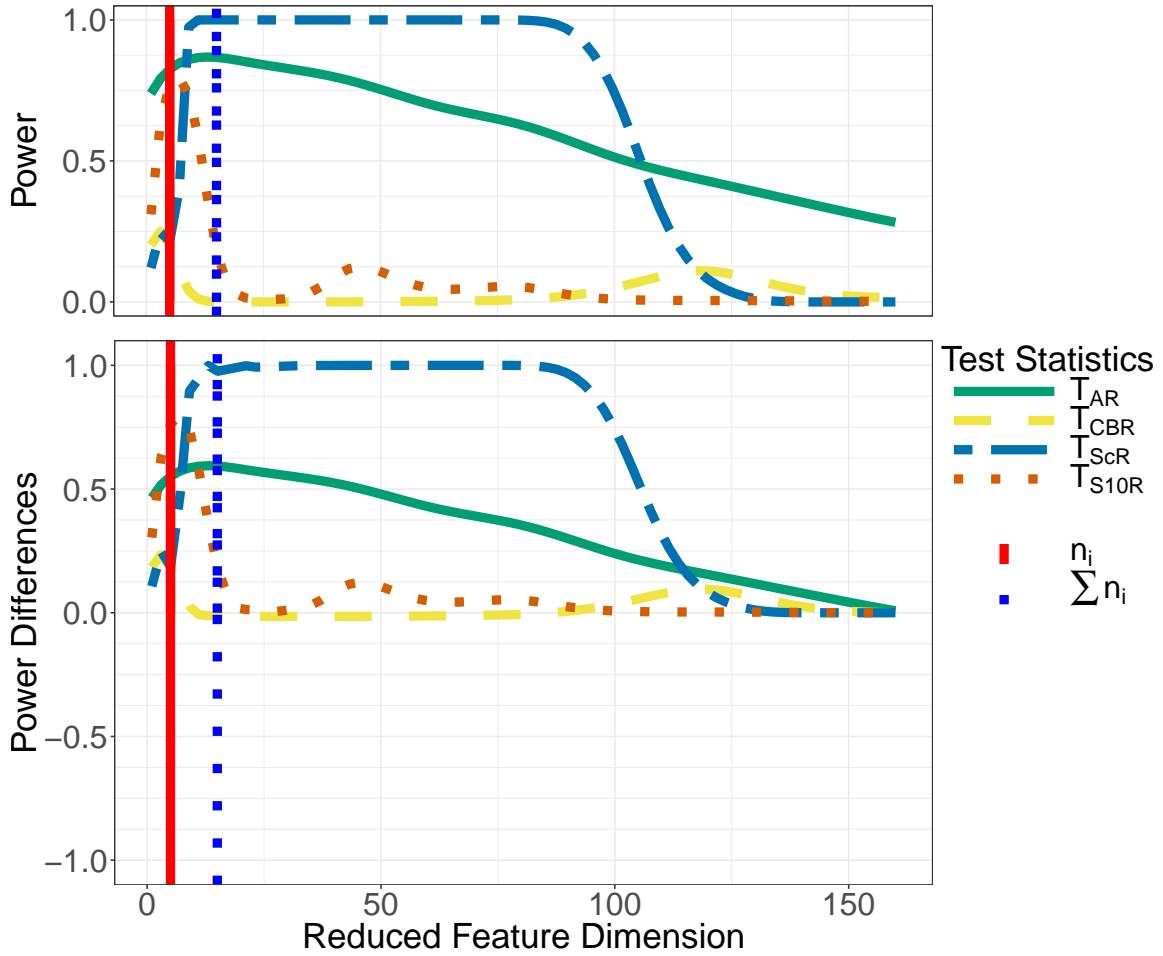| $p$ | $n_1 = n_2 = n_3$ | $p$ | $POW(T_A)$ | $POW(T_{CB})$ | $POW(T_{Sc})$ | $POW(T_{S10})$ |
|---|---|---|---|---|---|---|
| 80 | 5 | 80 | 0.37 | 0.12 | 0.00 | 0.34 |
| | 10 | 80 | 0.61 | 0.16 | 0.00 | 0.75 |
| 160 | 5 | 160 | 0.42 | 0.14 | 0.00 | 0.43 |
| | 10 | 160 | 0.67 | 0.21 | 0.01 | 0.77 |

104

Figure C.1: Reduced-dimension power curves and power-difference curves for POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) for three constant-times-identity covariance matrices with parameters $\mathbf{\Sigma}_1 = \mathbf{\Sigma}_3 = \mathbf{I}_p$ and $\mathbf{\Sigma}_2 = (1.5)\mathbf{I}_p$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 5, i = 1, 2, 3$.
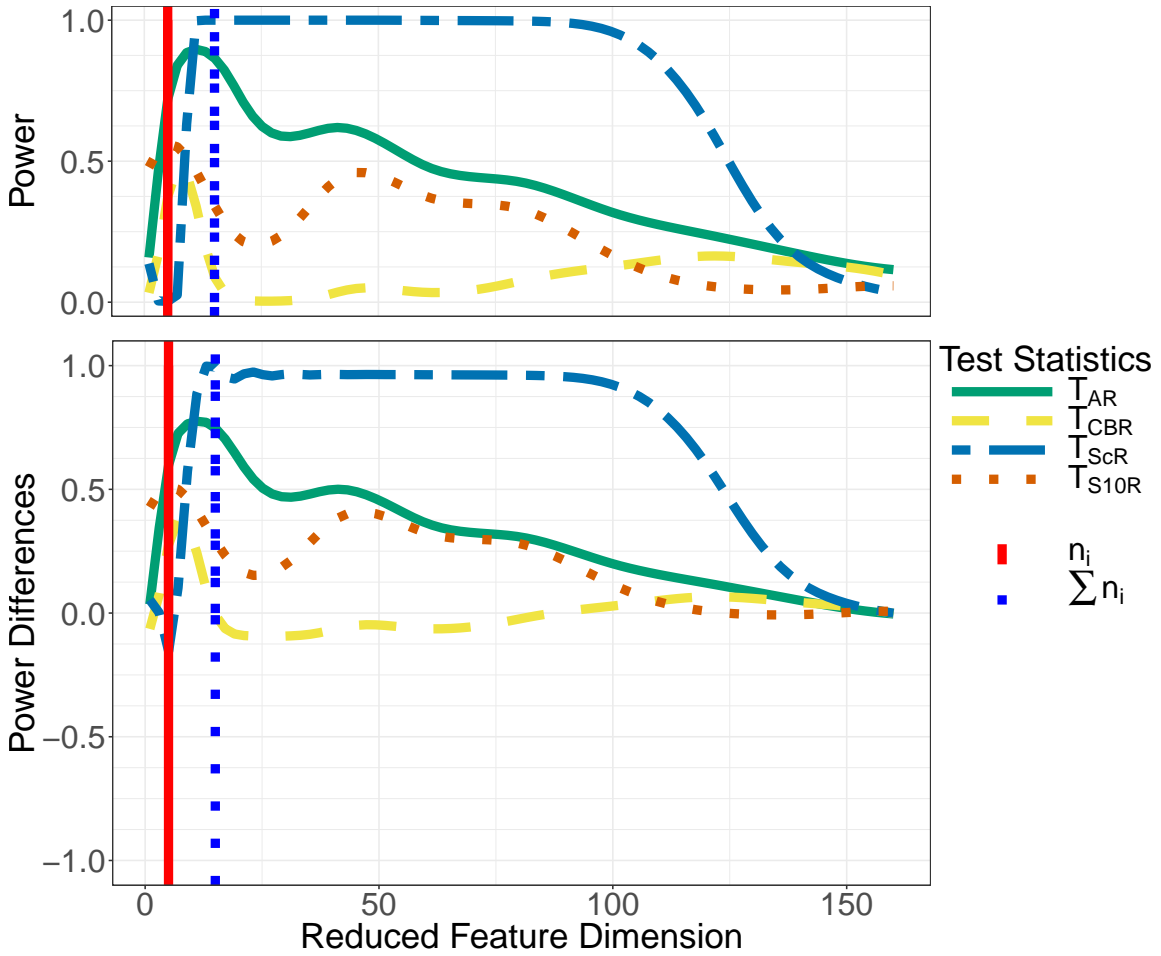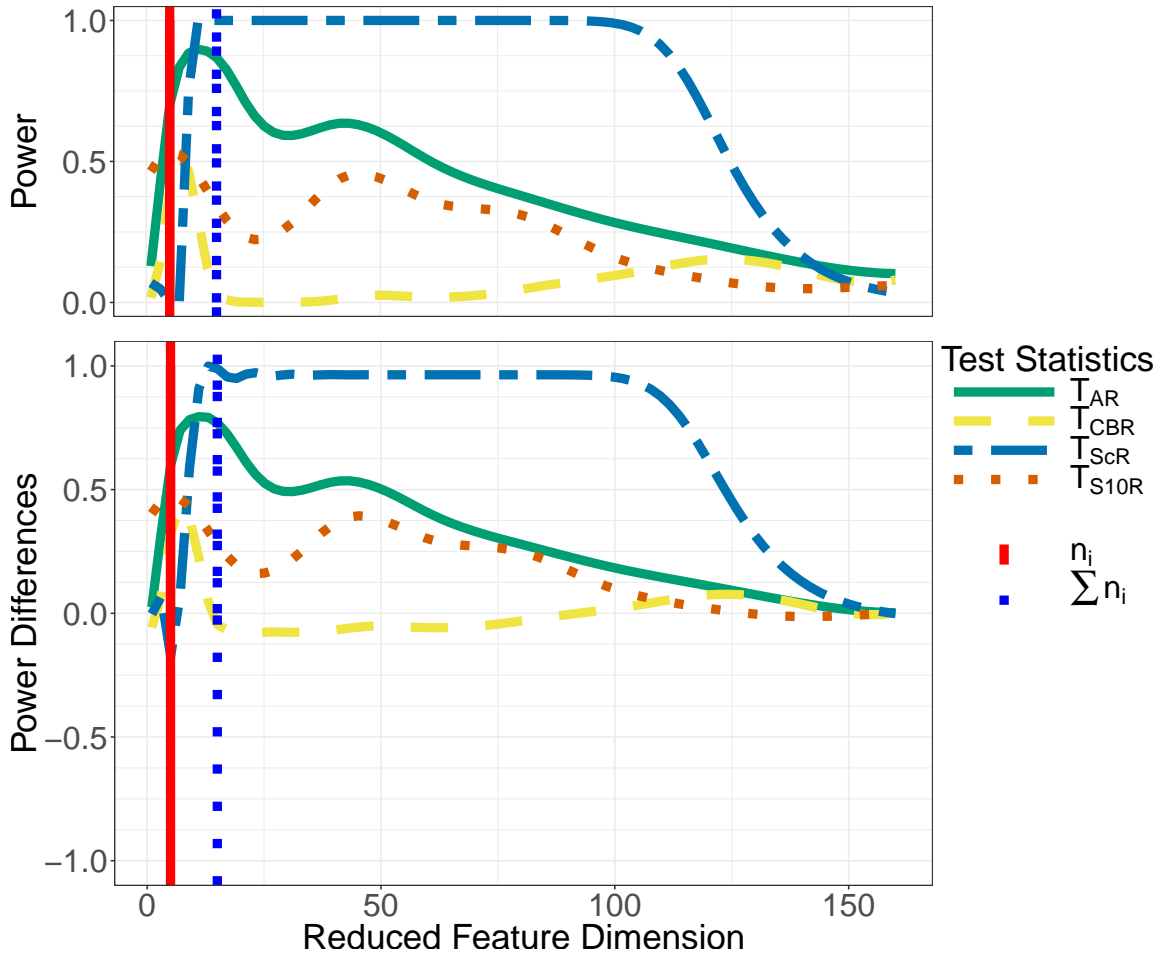
Figure C.2: Reduced-dimension power curves and power-difference curves for $\text{POW}(T_{RA})$, $\text{POW}(T_{RC})$, $\text{POW}(T_{RI})$, and $\text{POW}(T_{RS10})$ for three compound-symmetric covariance matrices with parameters $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_3 = (0.99)\mathbf{I}_p + (0.01)\mathbf{J}_p$ and $\boldsymbol{\Sigma}_2 = (0.95)\mathbf{I}_p + (0.05)\mathbf{J}_p$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 5, i = 1, 2, 3$.

Figure C.3: Reduced-dimension power curves and power-difference curves for POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) for three autoregressive covariance matrices with parameters $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_3 = 0.1^{|i-j|}$ and $\boldsymbol{\Sigma}_2 = 0.3^{|i-j|}$. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 5, i = 1, 2, 3$.
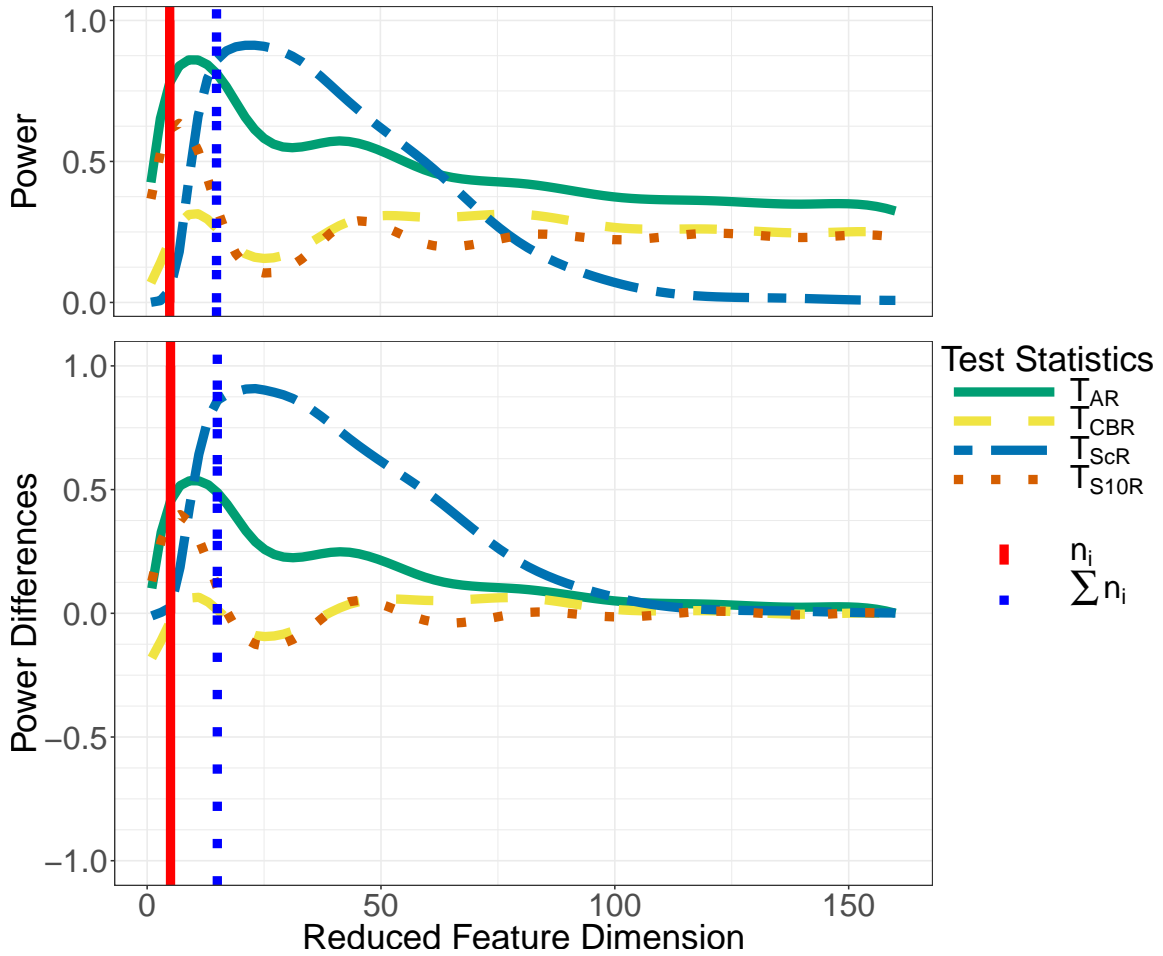
Figure C.4: Reduced-dimension power curves and power-difference curves for POW($T_{RA}$), POW($T_{RC}$), POW($T_{RI}$), and POW($T_{RS10}$) for three heterogeneous autoregressive covariance matrices. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 5, i = 1, 2, 3$.

Figure C.5: Reduced-dimension power curves and power-difference curves for $\text{POW}(T_{RA})$, $\text{POW}(T_{RC})$, $\text{POW}(T_{RI})$, and $\text{POW}(T_{RS10})$ for three unstructured covariance matrices. The reduced-data dimensions were $q \in \{1, 2, ..., 159\}$, the original data dimension was $p = 160$, and the common sample size was $n_i = 5, i = 1, 2, 3$.
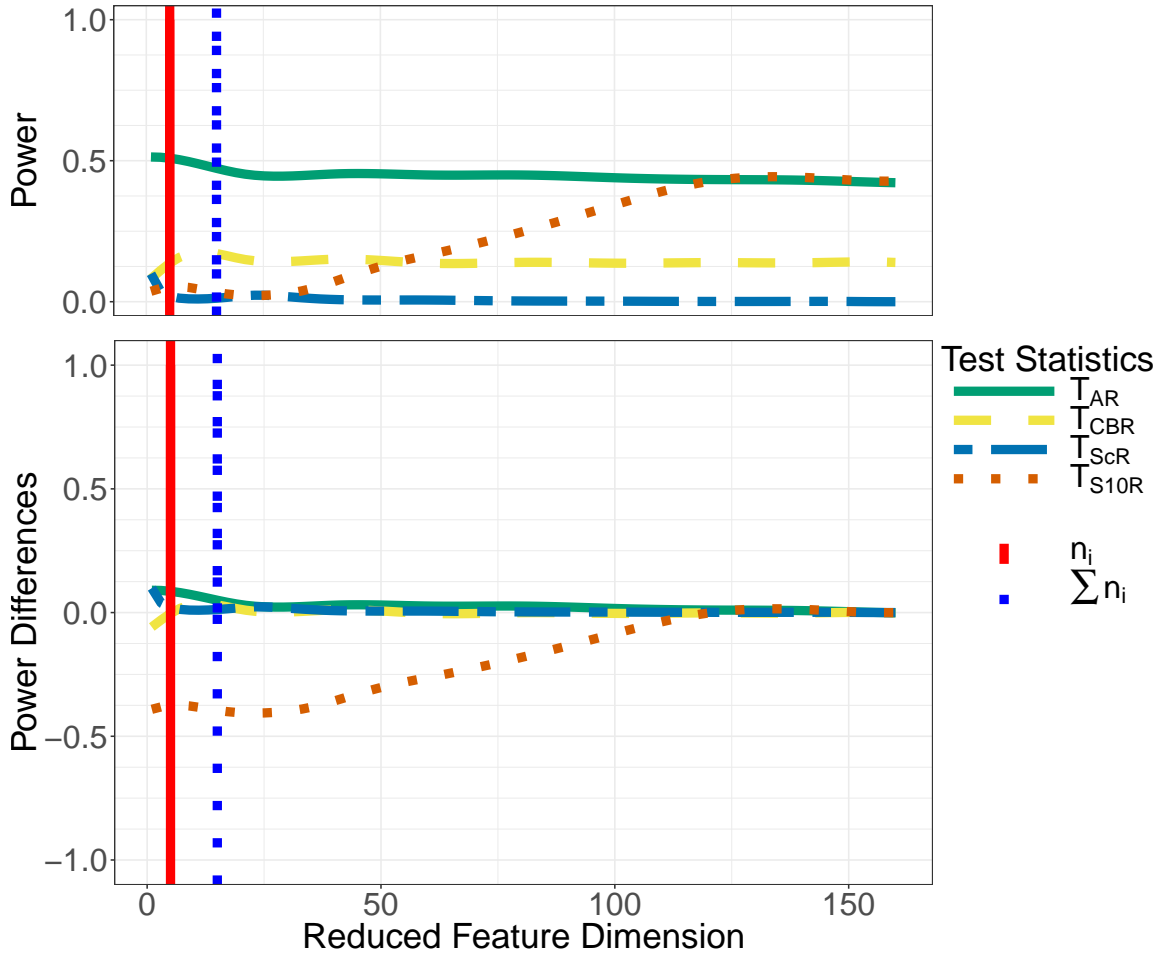
REFERENCES

Ahmad, M. R. (2017). Location-invariant tests of homogeneity of large-dimensional covariance matrices. *Journal of Statistical Theory and Practice*, pages 1–15.

Alon, U., Barkai, N., Notterman, D. A., Gish, K., Ybarra, S., Mack, D., and Levine, A. J. (1999). Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proceedings of the National Academy of Sciences of the United States of America*, 96(12):6745–50.

Chaipitak, S. and Chongcharoen, S. (2013). A test for testing the equality of two covariance matrices for high-dimensional data. *Journal of Applied Sciences*, 13(2):270–277.

Chen, S. X., Zhang, L.-X., and Zhong, P.-S. (2010). Tests for high-dimensional covariance matrices. *Journal of the American Statistical Association*, 105(490):810–819.

Cook, R. D. and Forzani, L. (2009). Likelihood-based sufficient dimension reduction. *Journal of the American Statistical Association*, 104(485):197–208.

Eckart, C. and Young, G. (1936). The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218.

Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J. P., Coller, H., Loh, M. L., Downing, J. R., Caligiuri, M. A., Bloomfield, C. D., and Lander, E. S. (1999). Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science*, 286(5439):531–7.

Ishii, A., Yata, K., and Aoshima, M. (2016). Asymptotic properties of the first principal component and equality tests of covariance matrices in high-dimension, low-sample-size context. *Journal of Statistical Planning and Inference*, 170:186–199.

Khan, J., Wei, J. S., Ringnér, M., Saal, L. H., Ladanyi, M., Westermann, F., Berthold, F., Schwab, M., Antonescu, C. R., Peterson, C., and Meltzer, P. S. (2001). Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks. *Nature Medicine*, 7(6):673–679.

Ledoit, O. and Wolf, M. (2002). Some hypothesis tests for the covariance matrix when the dimension Is large compared to the sample size. *The Annals of Statistics*, 30(4):1081–1102.

Ledoit, O. and Wolf, M. (2004). A well-conditioned estimator for large-dimensional covariance matrices. *Journal of Multivariate Analysis*, 88(2):365–411.

Luo, W. and Li, B. (2016). Combining eigenvalues and variation of eigenvectors for order determination. *Biometrika*, 103(4):875–887.

Peng, L., Chen, S. X., and Zhou, W. (2016). More powerful tests for sparse high-dimensional covariances matrices. *Journal of Multivariate Analysis*, 149:124–143.

Peters, B. C., Redner, R., and Decell, H. P. (1978). Characterizations of linear sufficient statistics. *Sankhyā: The Indian Journal of Statistics, Series A*, 40(3):303–309.

Ramey, J. A. (2016). datamicroarray: collection of data sets for classification.

Rohde, A. and Tsybakov, A. B. (2011). Estimation of high-dimensional low-rank matrices. *Ann. Statist.*, 39(2):887–930.

Schott, J. R. (2007). A test for the equality of covariance matrices when the dimension is large relative to the sample sizes. *Computational Statistics & Data Analysis*, 51(12):6535–6542.

Srivastava, M. S. (2005). Some tests concerning the covariance matrix in high dimensional data. *Journal of the Japan Statistical Society*, 35(2):251–272.

Srivastava, M. S. (2007). Testing the equality of two covariance matrices and independence of two sub-vectors with fewer observations than the dimension. In *International Conference on Advances in Interdisciplinary Stistics and Combinatorics*, University of North Carolina at Greensboro, NC, USA.

Srivastava, M. S. and Yanagihara, H. (2010). Testing the equality of several covariance matrices with fewer observations than the dimension. *Journal of Multivariate Analysis*, 101(6):1319–1329.

Srivastava, M. S., Yanagihara, H., and Kubokawa, T. (2014). Tests for covariance matrices in high dimension with less sample size. *Journal of Multivariate Analysis*, 130:289–309.