

QUADRO: A Supervised Dimension Reduction Method via Rayleigh Quotient Optimization

Jianqing Fan, Tracy Ke, Han Liu, and Lucy Xia

Princeton University

November 22, 2013

Abstract

We propose a novel Rayleigh quotient based sparse quadratic dimension reduction method — named QUADRO (Quadratic Dimension Reduction via Rayleigh Optimization) — for analyzing high dimensional data. Unlike in the linear setting where Rayleigh quotient optimization coincides with classification, these two problems are very different under nonlinear settings. In this paper, we clarify this difference and show that Rayleigh quotient optimization may be of independent scientific interests. One major challenge of Rayleigh quotient optimization is that the variance of quadratic statistics involves all fourth cross-moments of predictors, which are infeasible to compute for high-dimensional applications and may accumulate too many stochastic errors. This issue is resolved by considering a family of elliptical models. Moreover, for heavy-tail distributions, robust estimates of mean vectors and covariance matrices are employed to guarantee uniform convergence in estimating nonpolynomially many parameters, even though the fourth moments are assumed. Methodologically, QUADRO is based on elliptical models which allow us to formulate the Rayleigh quotient maximization as a convex optimization problem. Computationally, we propose an efficient linearized augmented Lagrangian method to solve the constrained optimization problem. Theoretically, we provide explicit rates of convergence in terms of Rayleigh quotient under both Gaussian and general elliptical models. Thorough numerical results on both synthetic and real datasets are also provided to back up our theoretical results.

1 Introduction

Rapid developments of imaging technology, microarray data studies, and many other applications call for the analysis of high-dimensional binary-labeled data. We consider the problem of finding a “nice” projection $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that embeds all data into the real line. Such projection f has applications in many statistical problems for analyzing high-dimensional binary-labeled data, including:

- *Dimension Reduction:* f provides a data reduction tool for people to visualize the high-dimensional data in a one-dimensional space.
- *Classification:* f can be used to construct classification rules. With a carefully-chosen set $A \subset \mathbb{R}$, we can classify a new data point $\mathbf{x} \in \mathbb{R}^d$ by checking whether or not $f(\mathbf{x}) \in A$.

- *Feature Selection*: When $f(\mathbf{x})$ only depends on a small number of coordinates of \mathbf{x} , this projection selects just a few features from numerous observed ones.

A natural question is what kind of f is a “nice” projection? It depends on the goal of statistical analysis. For classification, a good f should yield to a small classification error. In feature selection, different criteria select distinct features and they may suit for different real problems. In this paper, we propose using the following criterion for finding f :

Under the mapping f , the data are as “separable” as possible between two classes, and as “coherent” as possible within each class.

It can be formulated as to maximize the *Rayleigh quotient* of f . Suppose all data are drawn independently from a joint distribution of (\mathbf{X}, Y) , where $\mathbf{X} \in \mathbb{R}^d$, and $Y \in \{0, 1\}$ is the label. The *Rayleigh quotient* of f is defined as

$$\text{Rq}(f) \equiv \frac{\text{var} \{ \mathbb{E}[f(\mathbf{X})|Y] \}}{\text{var} \{ f(\mathbf{X}) - \mathbb{E}[f(\mathbf{X})|Y] \}}. \quad (1)$$

Here, the numerator is the variance of \mathbf{X} explained by the class label, and the denominator is the remaining variance of \mathbf{X} . Simple calculation shows that $\text{Rq}(f) = \pi(1 - \pi)R(f)$, where $\pi \equiv \mathbb{P}(Y = 0)$ and

$$R(f) \equiv \frac{\{ \mathbb{E}[f(\mathbf{X})|Y = 0] - \mathbb{E}[f(\mathbf{X})|Y = 1] \}^2}{\pi \text{var}[f(\mathbf{X})|Y = 0] + (1 - \pi) \text{var}[f(\mathbf{X})|Y = 1]}. \quad (2)$$

Our goal is to develop a data-driven procedure to find \hat{f} such that $\text{Rq}(\hat{f})$ is large and \hat{f} is sparse in the sense that it depends on few coordinates of \mathbf{X} .

The Rayleigh quotient, as a new criterion for finding a projection f , serves well for different purposes of statistical analysis. First, for dimension reduction, it is a meaningful criterion which takes care of both variance explanation and label explanation. In contrast, many popular dimension reduction methods such as principal component analysis only consider variance explanation. Second, as we shall see in Section 6, a monotone transform of the Rayleigh quotient approximates the classification error. As a result, starting from an f with a large Rayleigh quotient, we can also construct a classification rule with a small classification error. In addition, Rayleigh quotient maximization is a convex optimization for quadratic discriminant functions. Third, with appropriate regularization, this criterion can select features that are different from those selected by many existing dimension reduction and classification methods. Thus, it is a new feature selection tool for statistical studies.

1.1 Rayleigh quotient and classification error

Many popular statistical methods for analyzing high-dimensional binary-labeled data are based on classification error minimization, which is closely related to the Rayleigh quotient maximization. We summarize their connections and differences as follows:

- (a) In an “ideal” setting where two classes follow multivariate normal distributions with a common covariance matrix and the class of linear functions f is considered, the two criteria are exactly the same, with one being a monotone transform of the other.

- (b) In other settings, the two criteria can be very different.
- (c) In a “relaxed” setting where two classes follow elliptical distributions (including multivariate normal as a special case) with possibly non-equal covariance matrices and the class of quadratic functions f (including linear functions as special cases) is considered, the two criteria are closely related in the sense that a monotone transform of the Rayleigh quotient is an approximation of the classification error.

From these observations, we conclude that the Rayleigh quotient maximization is indeed a new criterion for both dimension reduction and feature selection. In addition, if we use it for classification, the Rayleigh quotient also serves as a surrogate of the classification error. In the remaining of this section, we show (a) and (b). We will discuss (c) in Section 6.

For each f , we define a family of classifiers $h_c(\mathbf{x}) = I\{f(\mathbf{x}) < c\}$ indexed by c , where $I(\cdot)$ is the indicator function defined as $I(A) = 1$ if an event A happens and $I(A) = 0$ otherwise. For each given c , we define the classification error of h_c to be $\text{err}(h_c) \equiv \mathbb{P}(h_c(\mathbf{X}) \neq Y)$. The classification error of f is then defined by

$$\text{Err}(f) \equiv \min_{c \in \mathbb{R}} \{ \text{err}(h_c) \}.$$

Most existing classification procedures aim at finding a data-driven projection \hat{f} such that $\text{Err}(\hat{f})$ is small (the threshold c is usually easy to choose). Examples include linear discriminant analysis (LDA) and its variations in high dimensions (e.g., Guo et al. (2005); Fan and Fan (2008); Cai and Liu (2011); Shao et al. (2011); Fan et al. (2012); Han et al. (2013)), quadratic discriminant analysis (QDA), support vector machine (SVM), logistic regression, boosting, etc.

We now compare $\text{Rq}(f)$ and $\text{Err}(f)$. Denote by $\pi = \mathbb{P}(Y = 0)$, $\boldsymbol{\mu}_1 = \mathbb{E}(\mathbf{X}|Y = 0)$, $\boldsymbol{\Sigma}_1 = \text{cov}(\mathbf{X}|Y = 0)$, $\boldsymbol{\mu}_2 = \mathbb{E}(\mathbf{X}|Y = 1)$ and $\boldsymbol{\Sigma}_2 = \text{cov}(\mathbf{X}|Y = 1)$. We consider linear functions $\{f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x} + b : \mathbf{a} \in \mathbb{R}^d, b \in \mathbb{R}\}$, and write $\text{Rq}(\mathbf{a}) = \text{Rq}(\mathbf{a}^\top \mathbf{x})$, $\text{Err}(\mathbf{a}) = \text{Err}(\mathbf{a}^\top \mathbf{x})$ for short. By direct calculation, when the two classes have a common covariance matrix $\boldsymbol{\Sigma}$,

$$\text{Rq}(\mathbf{a}) = \pi(1 - \pi) \frac{[\mathbf{a}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)]^2}{\mathbf{a}^\top \boldsymbol{\Sigma} \mathbf{a}}.$$

Hence, the optimal $\mathbf{a}_R = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$. On the other hand, when data follow multivariate normal distributions, the optimal classifier is $h^*(\mathbf{x}) = I\{\mathbf{a}_E^\top \mathbf{x} < c\}$, where $\mathbf{a}_E = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ and $c = \frac{1}{2} \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \frac{1}{2} \boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \log(\frac{1-\pi}{\pi})$. It is observed that $\mathbf{a}_R = \mathbf{a}_E$ and the two criteria are the same. In fact, for all vectors \mathbf{a} such that $\mathbf{a}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) > 0$,

$$\text{Err}(\mathbf{a}) = 1 - \Phi \left(\frac{1}{2} \left[\frac{\text{Rq}(\mathbf{a})}{\pi(1 - \pi)} \right]^{1/2} \right),$$

where Φ is the distribution function of a standard normal random variable, and we fix $c = \mathbf{a}^\top (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)/2$. Therefore, the classification error is a monotone transform of the Rayleigh quotient.

When we move away from these ideal assumptions, the above two criteria can be very different. We illustrate this point using a bivariate distribution, i.e., $d = 2$, with different covariance matrices. Specifically, $\pi = 0.55$, $\boldsymbol{\mu}_1 = (0, 0)^\top$, $\boldsymbol{\mu}_2 = (1.28, 0.8)^\top$, $\boldsymbol{\Sigma}_1 = \text{diag}(1, 1)$ and $\boldsymbol{\Sigma}_2 = \text{diag}(3, 1/3)$. We still consider linear functions $f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x}$ but select only one out of

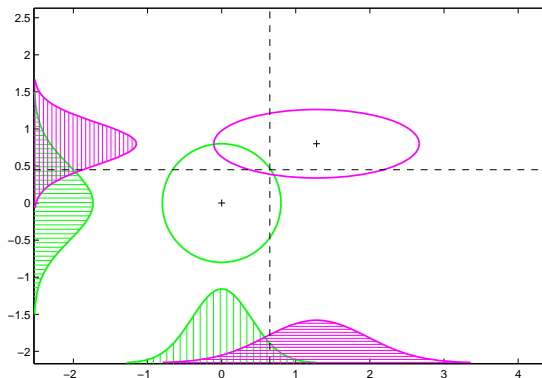


Figure 1: An example in \mathbb{R}^2 . The green and purple represent class 1 and class 2 respectively. The ellipses are contours of distributions. Probability densities after being projected to X_1 and X_2 are also displayed. The dot lines correspond to optimal thresholds for classification using each feature.

the two features, X_1 or X_2 . Then, the maximum Rayleigh quotients by using each of the two features alone are 0.853 and 0.923 respectively, whereas the minimum classification errors are 0.284 and 0.295 respectively. As a result, under the criterion of maximizing Rayleigh quotient, Feature 2 is selected; under the criterion of minimizing classification error, Feature 1 is selected. Figure 1 displays the distributions of data after being projected to each of the two features. It shows that since data from the second class has a much larger variability at Feature 1 than at Feature 2, the Rayleigh quotient maximization favors Feature 2, although Feature 1 yields a smaller classification error.

1.2 Objective of the paper

In this paper, we consider the Rayleigh quotient maximization problem in the following setting:

- We consider sparse quadratic functions, i.e., $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{\Omega} \mathbf{x} - 2\boldsymbol{\delta}^\top \mathbf{x}$, where $\mathbf{\Omega}$ is a sparse $d \times d$ symmetric matrix and $\boldsymbol{\delta}$ is a sparse d -dimensional vector.
- The two classes can have different covariance matrices.
- Data from these two classes follow *elliptical distributions*.
- The dimension is large (it is possible that $d \gg n$).

Rayleigh quotient maximization can trace back to Fisher's linear discriminant analysis. However, our setting has several new ingredients. First, we go beyond linear classifiers to enhance flexibility. It is well known that the linear classifiers are inefficient. For example, when two classes have the same mean, linear classifiers perform no better than random guess. Instead of exploring arbitrary nonlinear functions, we consider the class of quadratic functions, so that the Rayleigh quotient still has a nice parametric formulation, and at the same time it helps identify interaction effects between features. Second, we drop the requirement that the two classes share a common covariance matrix, which is a critical condition for Fisher's rule and many other high dimensional classification methods (e.g., Fan and Fan (2008); Fan et al. (2012); Cai and Liu (2011)). In fact, by using quadratic discriminant functions, we take advantage of the difference of covariance matrices between the two classes to enhance classification power.

Third, we generalize multivariate normal distributions to the elliptical family, which includes many heavy-tailed distributions, such as multivariate t-distributions, Laplace distributions, and Cauchy distributions. This family of distributions allows us to avoid estimating all $O(d^4)$ fourth cross-moments of d predictors in computing the variance of quadratic statistics and hence overcomes the computation and noise accumulation issues.

In our setting, Fisher’s rule, i.e., $\mathbf{a}_R = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$, is no longer a solution to the Rayleigh quotient maximization. We propose a new method, named Quadratic Dimension Reduction via Rayleigh Optimization (QUADRO). It is a *Rayleigh-quotient-oriented procedure*, and is a statistical tool for simultaneous dimension reduction and feature selection. QUADRO has several properties. First, it is a statistically efficient generalization of Fisher’s linear discriminant analysis to the quadratic setting. A naive generalization involves estimation of all fourth cross-moments of the two underlying distributions. In contrast, QUADRO only requires estimating a one-dimensional kurtosis parameter. Second, QUADRO adopts rank-based estimators and robust M-estimators of the covariance matrices and the means. Therefore, it is robust to possibly heavy-tail distributions. Third, QUADRO can be formulated as a convex programming and is computationally efficient.

Theoretically, we prove that under elliptical models, the Rayleigh quotient of the estimated quadratic function \hat{f} converges to population maximum Rayleigh quotient at rate $O_p(s\sqrt{\log(d)/n})$, where s is the number of important features (counting both single terms and interaction terms). In addition, we establish a connection between our method and quadratic discriminant analysis (QDA) under elliptical models.

The rest of this paper is organized as follows. Section 2 formulates Rayleigh quotient maximization as a convex optimization. Section 3 describes QUADRO. Section 4 discusses rank-based estimators and robust M-estimators used in QUADRO. Section 5 presents theoretical analysis. Section 6 discusses the application of QUADRO in elliptically distributed classification problems. Section 7 contains numerical studies and real data examples. Section 8 concludes the paper. All proofs are relegated to Section 9.

Throughout this paper, for $0 \leq q \leq \infty$, $|\mathbf{v}|_q$ denotes the L_q -norm of a vector \mathbf{v} , $|\mathbf{A}|_q$ denotes the elementwise L_q -norm of a matrix \mathbf{A} and $\|\mathbf{A}\|_q$ denotes the matrix L_q -norm of \mathbf{A} . When $q = 2$, we omit the subscript q . $\lambda_{\min}(\mathbf{A})$ and $\lambda_{\max}(\mathbf{A})$ denote the minimum and maximum eigenvalues of \mathbf{A} . $\det(\mathbf{A})$ denotes the determinant of \mathbf{A} . Let $I(\cdot)$ be the indicator function: for any event B , $I(B) = 1$ if B happens and $I(B) = 0$ otherwise. Let $\text{sign}(\cdot)$ be the sign function, where $\text{sign}(u) = 1$ when $u \geq 0$ and $\text{sign}(u) = -1$ when $u < 0$.

2 Rayleigh Quotient for Quadratic Functions

We first study the population form of Rayleigh quotient for an arbitrary quadratic function. We show that it has a simplified form under the elliptical family.

For a quadratic function

$$Q(\mathbf{X}) = \mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X} - 2\boldsymbol{\delta}^\top \mathbf{X},$$

using (2), its Rayleigh quotient is

$$R(\boldsymbol{\Omega}, \boldsymbol{\delta}) = \frac{\{\mathbb{E}[Q(\mathbf{X})|Y = 0] - \mathbb{E}[Q(\mathbf{X})|Y = 1]\}^2}{\pi \text{var}[Q(\mathbf{X})|Y = 0] + (1 - \pi) \text{var}[Q(\mathbf{X})|Y = 1]} \quad (3)$$

up to a constant multiplier. The Rayleigh quotient maximization can be expressed as

$$\max_{(\boldsymbol{\Omega}, \boldsymbol{\delta}): \boldsymbol{\Omega}=\boldsymbol{\Omega}^\top} R(\boldsymbol{\Omega}, \boldsymbol{\delta}).$$

2.1 General setting

Suppose $\mathbb{E}(\mathbf{Z}) = \boldsymbol{\mu}$ and $\text{cov}(\mathbf{Z}) = \boldsymbol{\Sigma}$, by direct calculation,

$$\begin{aligned} \mathbb{E}[Q(\mathbf{Z})] &= \text{tr}(\boldsymbol{\Omega}\boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\Omega}\boldsymbol{\mu} - 2\boldsymbol{\delta}^\top \boldsymbol{\mu}, \\ \text{var}[Q(\mathbf{Z})] &= \mathbb{E}[\text{tr}(\boldsymbol{\Omega}\mathbf{Z}\mathbf{Z}^\top \boldsymbol{\Omega}\mathbf{Z}\mathbf{Z}^\top)] - 4\mathbb{E}[\boldsymbol{\delta}^\top \mathbf{Z}\mathbf{Z}^\top \boldsymbol{\Omega}\mathbf{Z}] + 4\boldsymbol{\delta}^\top \boldsymbol{\Sigma}\boldsymbol{\delta} + 4(\boldsymbol{\delta}^\top \boldsymbol{\mu})^2 - \{\mathbb{E}[Q(\mathbf{Z})]\}^2. \end{aligned}$$

So $\mathbb{E}[Q(\mathbf{Z})]$ is a linear combination of the elements in $\{\Omega(i, j), 1 \leq i \leq j \leq d; \delta(i), 1 \leq i \leq d\}$, and $\text{var}[Q(\mathbf{Z})]$ is a quadratic form of these elements. The coefficients in $\mathbb{E}[Q(\mathbf{Z})]$ are functions of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ only. However, the coefficients in $\text{var}[Q(\mathbf{Z})]$ also depend on all the fourth cross-moments of \mathbf{Z} and there are $O(d^4)$ of them.

Let us define $M_1(\boldsymbol{\Omega}, \boldsymbol{\delta}) = \mathbb{E}[Q(\mathbf{X})|Y = 0]$, $L_1(\boldsymbol{\Omega}, \boldsymbol{\delta}) = \text{var}[Q(\mathbf{X})|Y = 0]$, and $M_2(\boldsymbol{\Omega}, \boldsymbol{\delta})$, $L_2(\boldsymbol{\Omega}, \boldsymbol{\delta})$ similarly. Also, let $\kappa = (1 - \pi)/\pi$. We have

$$R(\boldsymbol{\Omega}, \boldsymbol{\delta}) = \frac{[M_1(\boldsymbol{\Omega}, \boldsymbol{\delta}) - M_2(\boldsymbol{\Omega}, \boldsymbol{\delta})]^2}{L_1(\boldsymbol{\Omega}, \boldsymbol{\delta}) + \kappa L_2(\boldsymbol{\Omega}, \boldsymbol{\delta})}.$$

Therefore, both the numerator and denominator are quadratic combinations of the elements in $\boldsymbol{\Omega}$ and $\boldsymbol{\delta}$. We can stack the $d(d+1)/2$ elements in $\boldsymbol{\Omega}$ (assuming it is symmetric) and the d elements in $\boldsymbol{\delta}$ into a long vector \mathbf{v} . Then $R(\boldsymbol{\Omega}, \boldsymbol{\delta})$ can be written as

$$R(\mathbf{v}) = \frac{(\mathbf{a}^\top \mathbf{v})^2}{\mathbf{v}^\top \mathbf{A} \mathbf{v}},$$

where \mathbf{a} is a $d' \times 1$ vector, \mathbf{A} is a $d' \times d'$ positive semi-definite matrix, and $d' = d(d+1)/2 + d$. \mathbf{A} and \mathbf{a} are determined by the coefficients in the denominator and numerator of $R(\boldsymbol{\Omega}, \boldsymbol{\delta})$ respectively. Now, $\max_{(\boldsymbol{\Omega}, \boldsymbol{\delta})} R(\boldsymbol{\Omega}, \boldsymbol{\delta})$ is equivalent to $\max_{\mathbf{v}} R(\mathbf{v})$. It has explicit solutions. For example, when \mathbf{A} is positive definite, the function $R(\mathbf{v})$ is maximized at

$$\mathbf{v}^* = \mathbf{A}^{-1} \mathbf{a}.$$

We can then reshape \mathbf{v}^* to get the desired $(\boldsymbol{\Omega}^*, \boldsymbol{\delta}^*)$.

Practical implementation of the above idea is infeasible in high dimensions as it involves $O(d^4)$ cross moments of \mathbf{Z} . This not only poses computational challenges, but also accumulates noise in the estimation. Furthermore, good estimates of fourth moments usually require the existence of eighth moments, which is not realistic for many heavy tailed distributions. These problems can be avoided under the elliptical family, as we now illustrate in the next subsection.

2.2 Elliptical distributions

The elliptical family contains multivariate distributions whose densities have elliptical contours. It generalizes multivariate normal distributions and inherits many of their nice properties.

Given a $d \times 1$ vector $\boldsymbol{\mu}$ and a $d \times d$ positive definite matrix $\boldsymbol{\Sigma}$, a random vector \mathbf{Z} that follows an elliptical distribution admits

$$\mathbf{Z} = \boldsymbol{\mu} + \xi \boldsymbol{\Sigma}^{1/2} \mathbf{U}, \tag{4}$$

where \mathbf{U} is a random vector which follows the uniform distribution on unit sphere \mathcal{S}^{d-1} , and ξ is a nonnegative random variable independent of \mathbf{U} . Denote the elliptical distribution by $\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$, where g is the density of ξ . In this paper, we always assume that $\mathbb{E}\xi^4 < \infty$ and require that $\mathbb{E}(\xi^2) = d$ for the model identifiability. Then $\boldsymbol{\Sigma}$ is the covariance matrix of \mathbf{Z} .

Proposition 2.1. *Suppose \mathbf{Z} follows an elliptical distribution as in (4). Then*

$$\begin{aligned}\mathbb{E}[Q(\mathbf{Z})] &= \text{tr}(\boldsymbol{\Omega}\boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\Omega}\boldsymbol{\mu} - 2\boldsymbol{\mu}^\top \boldsymbol{\delta}, \\ \text{var}[Q(\mathbf{Z})] &= 2(1 + \gamma) \text{tr}(\boldsymbol{\Omega}\boldsymbol{\Sigma}\boldsymbol{\Omega}\boldsymbol{\Sigma}) + \gamma[\text{tr}(\boldsymbol{\Omega}\boldsymbol{\Sigma})]^2 + 4(\boldsymbol{\Omega}\boldsymbol{\mu} - \boldsymbol{\delta})^\top \boldsymbol{\Sigma}(\boldsymbol{\Omega}\boldsymbol{\mu} - \boldsymbol{\delta}),\end{aligned}$$

where $\gamma = \frac{E(\xi^4)}{d(d+2)} - 1$ is the kurtosis parameter.

The proof is given in Section 9. The variance of $Q(\mathbf{Z})$ does not involve any fourth cross-moments, but only the kurtosis parameter γ . For multivariate normal distributions, ξ^2 follows a χ^2 -distribution with d degrees of freedom, and $\gamma = 0$. For multivariate t distribution with degrees of freedom $\nu > 4$, we have $\gamma = 2/(\nu - 4)$.

2.3 Rayleigh optimization

We assume that the two classes both follow elliptical distributions: $\mathbf{X}|(Y = 0) \sim \mathcal{E}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, g_1)$ and $\mathbf{X}|(Y = 1) \sim \mathcal{E}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2, g_2)$. Without loss of generality, we assume the quantity γ is the same for both classes of conditional distributions. Let

$$\begin{aligned}M(\boldsymbol{\Omega}, \boldsymbol{\delta}) &= -\boldsymbol{\mu}_1^\top \boldsymbol{\Omega}\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2^\top \boldsymbol{\Omega}\boldsymbol{\mu}_2 + 2(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\delta} - \text{tr}(\boldsymbol{\Omega}(\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2)), \\ L_k(\boldsymbol{\Omega}, \boldsymbol{\delta}) &= 2(1 + \gamma) \text{tr}(\boldsymbol{\Omega}\boldsymbol{\Sigma}_k\boldsymbol{\Omega}\boldsymbol{\Sigma}_k) + \gamma[\text{tr}(\boldsymbol{\Omega}\boldsymbol{\Sigma}_k)]^2 + 4(\boldsymbol{\Omega}\boldsymbol{\mu}_k - \boldsymbol{\delta})^\top \boldsymbol{\Sigma}_k(\boldsymbol{\Omega}\boldsymbol{\mu}_k - \boldsymbol{\delta}),\end{aligned}\quad (5)$$

for $k = 1$ and 2 . Combining (3) with Proposition 2.1, we have

$$R(\boldsymbol{\Omega}, \boldsymbol{\delta}) = \frac{[M(\boldsymbol{\Omega}, \boldsymbol{\delta})]^2}{L_1(\boldsymbol{\Omega}, \boldsymbol{\delta}) + \kappa L_2(\boldsymbol{\Omega}, \boldsymbol{\delta})},\quad (6)$$

where $\kappa = (1 - \pi)/\pi$.

Note that if we multiply both $\boldsymbol{\Omega}$ and $\boldsymbol{\delta}$ by a common constant, $R(\boldsymbol{\Omega}, \boldsymbol{\delta})$ remains unchanged. Therefore, maximizing $R(\boldsymbol{\Omega}, \boldsymbol{\delta})$ is equivalent to solving the following constrained minimization problem

$$\min_{(\boldsymbol{\Omega}, \boldsymbol{\delta}): M(\boldsymbol{\Omega}, \boldsymbol{\delta})=1, \boldsymbol{\Omega}=\boldsymbol{\Omega}^\top} \{L_1(\boldsymbol{\Omega}, \boldsymbol{\delta}) + \kappa L_2(\boldsymbol{\Omega}, \boldsymbol{\delta})\}.\quad (7)$$

We call problem (7) the *Rayleigh optimization*. It is a convex problem whenever $\boldsymbol{\Sigma}_1$ and $\boldsymbol{\Sigma}_2$ are both positive semi-definite.

The formulation of the Rayleigh optimization only involves the means and covariance matrices, and the kurtosis parameter γ . Therefore, if we know γ (e.g., when we know which subfamily the distributions belong to) and have good estimates $(\hat{\boldsymbol{\mu}}_1, \hat{\boldsymbol{\mu}}_2, \hat{\boldsymbol{\Sigma}}_1, \hat{\boldsymbol{\Sigma}}_2)$, we can solve the empirical version of (7) to obtain $(\hat{\boldsymbol{\Omega}}, \hat{\boldsymbol{\delta}})$, which is the main idea of QUADRO. In addition, (7) is a convex problem, with a quadratic objective and equality constraints. Hence, it can be solved efficiently by many optimization algorithms.

3 Quadratic Dimension Reduction via Rayleigh Optimization

Now, we formally introduce the QUADRO procedure. We fix a model parameter $\gamma \geq 0$. Let \widehat{M} , \widehat{L}_1 and \widehat{L}_2 be the sample versions of M, L_1, L_2 in (5) by replacing $(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2)$ with their estimates. Details of these estimates will be given in Section 4. Let $\widehat{\pi} = n_1/(n_1 + n_2)$ and $\kappa = \widehat{\pi}/(1 - \widehat{\pi})$. Given tuning parameters $\lambda_1 > 0$ and $\lambda_2 > 0$, we solve

$$\min_{(\boldsymbol{\Omega}, \boldsymbol{\delta}): \widehat{M}(\boldsymbol{\Omega}, \boldsymbol{\delta})=1, \boldsymbol{\Omega}=\boldsymbol{\Omega}^\top} \{ \widehat{L}_1(\boldsymbol{\Omega}, \boldsymbol{\delta}) + \kappa \widehat{L}_2(\boldsymbol{\Omega}, \boldsymbol{\delta}) + \lambda_1 |\boldsymbol{\Omega}|_1 + \lambda_2 |\boldsymbol{\delta}|_1 \}. \quad (8)$$

We propose a linearized augmented Lagrangian method to solve (8). To simplify the notation, we write $\widehat{L} = \widehat{L}_1 + \kappa \widehat{L}_2$, and omit the hat symbol on M and L when there is no confusion. The optimization problem is then

$$\min_{(\boldsymbol{\Omega}, \boldsymbol{\delta}): M(\boldsymbol{\Omega}, \boldsymbol{\delta})=1, \boldsymbol{\Omega}=\boldsymbol{\Omega}^\top} \{ L(\boldsymbol{\Omega}, \boldsymbol{\delta}) + \lambda_1 |\boldsymbol{\Omega}|_1 + \lambda_2 |\boldsymbol{\delta}|_1 \}.$$

For an algorithm parameter $\rho > 0$, and a dual variable ν , we define the *augmented Lagrangian* as

$$F_\rho(\boldsymbol{\Omega}, \boldsymbol{\delta}, \nu) = L(\boldsymbol{\Omega}, \boldsymbol{\delta}) + \nu[M(\boldsymbol{\Omega}, \boldsymbol{\delta}) - 1] + (\rho/2)[M(\boldsymbol{\Omega}, \boldsymbol{\delta}) - 1]^2.$$

Using zero as initial values, we iteratively update

- $\boldsymbol{\delta}^{(k)} = \operatorname{argmin}_{\boldsymbol{\delta}} \{ F_\rho(\boldsymbol{\Omega}^{(k-1)}, \boldsymbol{\delta}, \nu^{(k-1)}) + \lambda_2 |\boldsymbol{\delta}|_1 \},$
- $\boldsymbol{\Omega}^{(k)} = \operatorname{argmin}_{\boldsymbol{\Omega}: \boldsymbol{\Omega}=\boldsymbol{\Omega}^\top} \{ F_\rho(\boldsymbol{\Omega}, \boldsymbol{\delta}^{(k)}, \nu^{(k-1)}) + \lambda_1 |\boldsymbol{\Omega}|_1 \},$
- $\nu^{(k)} = \nu^{(k-1)} + \rho[M(\boldsymbol{\Omega}^{(k)}, \boldsymbol{\delta}^{(k)}) - 1].$

Here, the first two steps are *primal updates* and the third step is a *dual update*.

First, we consider the update of $\boldsymbol{\delta}$. When $\boldsymbol{\Omega}$ and ν are fixed, we can write

$$F_\rho(\boldsymbol{\Omega}, \boldsymbol{\delta}, \nu) = \boldsymbol{\delta}^\top \mathbf{A} \boldsymbol{\delta} - 2\boldsymbol{\delta}^\top \mathbf{b} + c_\rho(\boldsymbol{\Omega}, \nu),$$

where

$$\begin{aligned} \mathbf{A} &= 4(\boldsymbol{\Sigma}_1 + \kappa \boldsymbol{\Sigma}_2) + 2\rho(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top, \\ \mathbf{b} &= 4(\boldsymbol{\Sigma}_1 \boldsymbol{\Omega} \boldsymbol{\mu}_1 + \kappa \boldsymbol{\Sigma}_2 \boldsymbol{\Omega} \boldsymbol{\mu}_2) + [\rho \operatorname{tr}(\boldsymbol{\Omega}(\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2)) + \rho \boldsymbol{\mu}_1^\top \boldsymbol{\Omega} \boldsymbol{\mu}_1 - \rho \boldsymbol{\mu}_2^\top \boldsymbol{\Omega} \boldsymbol{\mu}_2 + (\rho - \nu)](\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \end{aligned} \quad (9)$$

and $c_\rho(\boldsymbol{\Omega}, \nu)$ does not depend on $\boldsymbol{\delta}$. Note that \mathbf{A} is a positive semi-definite matrix. The update of $\boldsymbol{\delta}$ is indeed a Lasso problem.

Next, we consider the update of $\boldsymbol{\Omega}$. When $\boldsymbol{\delta}$ and ν are fixed, $F_\rho(\boldsymbol{\Omega}, \boldsymbol{\delta}, \nu)$ is a convex function of $\boldsymbol{\Omega}$. We propose an approximate update step: We first “linearize” F_ρ at $\boldsymbol{\Omega} = \boldsymbol{\Omega}^{(k-1)}$ to construct an upper envelope \bar{F}_ρ , and then minimize this upper envelope. In detail, at any $\boldsymbol{\Omega} = \boldsymbol{\Omega}_0$, we consider an upper bound of $F_\rho(\boldsymbol{\Omega}, \boldsymbol{\delta}, \nu)$:

$$\begin{aligned} \bar{F}_\rho(\boldsymbol{\Omega}, \boldsymbol{\delta}, \nu) &\equiv F_\rho(\boldsymbol{\Omega}_0, \boldsymbol{\delta}, \nu) + \sum_{1 \leq i \leq j \leq d} [\Omega(i, j) - \Omega_0(i, j)] \frac{\partial F_\rho(\boldsymbol{\Omega}_0, \boldsymbol{\delta}, \nu)}{\partial \Omega(i, j)} \\ &\quad + \frac{\tau}{2} \sum_{1 \leq i \leq j \leq d} [\Omega(i, j) - \Omega_0(i, j)]^2, \end{aligned}$$

where τ is a large enough constant (e.g., we can take $\tau = \sum_{1 \leq i \leq j \leq d} \frac{\partial^2 F_\rho(\boldsymbol{\Omega}_0, \boldsymbol{\delta}, \nu)}{\partial \Omega(i, j)^2}$). We then minimize $\bar{F}_\rho(\boldsymbol{\Omega}, \boldsymbol{\delta}, \nu) + \lambda_1 |\boldsymbol{\Omega}|_1$ to update $\boldsymbol{\Omega}$. This modified update step has an explicit solution

$$\Omega^*(i, j) = \mathcal{S}\left(\Omega_0(i, j) - \frac{1}{\tau} \frac{\partial F_\rho(\boldsymbol{\Omega}_0, \boldsymbol{\delta}, \nu)}{\partial \Omega(i, j)}, \frac{\lambda_1}{\tau}\right),$$

where $\mathcal{S}(x, a) \equiv (|x| - a)_+ \text{sign}(x)$ is the soft-thresholding function. We can write $\boldsymbol{\Omega}^*$ in a matrix form. Let

$$\begin{aligned} \mathbf{D} = & 4(1 + \gamma)(\boldsymbol{\Sigma}_1 \boldsymbol{\Omega} \boldsymbol{\Sigma}_1 + \kappa \boldsymbol{\Sigma}_2 \boldsymbol{\Omega} \boldsymbol{\Sigma}_2) + 2\gamma [\text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma}_1) \boldsymbol{\Sigma}_1 + \kappa \text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma}_2) \boldsymbol{\Sigma}_2] \\ & + 4 \text{sym}(\boldsymbol{\Sigma}_1(\boldsymbol{\Omega} \boldsymbol{\mu}_1 - \boldsymbol{\delta}) \boldsymbol{\mu}_1^\top + \kappa \boldsymbol{\Sigma}_2(\boldsymbol{\Omega} \boldsymbol{\mu}_2 - \boldsymbol{\delta}) \boldsymbol{\mu}_2^\top), \end{aligned} \quad (10)$$

where $\text{sym}(\mathbf{B}) = (\mathbf{B} + \mathbf{B}^\top)/2$ for any square matrix \mathbf{B} . By direct calculation,

$$\boldsymbol{\Omega}^* = \mathcal{S}\left(\boldsymbol{\Omega}_0 - \frac{1}{\tau} \mathbf{D}, \frac{\lambda_1}{\tau}\right).$$

We now describe our algorithm. Let us initialize $\boldsymbol{\Omega}^{(0)} = \mathbf{0}_{d \times d}$, $\boldsymbol{\delta}^{(0)} = \mathbf{0}$, and $\nu^{(0)} = 0$. At iteration k , the algorithm updates as follows:

- Compute $\mathbf{A} = \mathbf{A}(\boldsymbol{\Omega}^{(k-1)}, \boldsymbol{\delta}^{(k-1)}, \nu^{(k-1)})$ and $\mathbf{b} = \mathbf{b}(\boldsymbol{\Omega}^{(k-1)}, \boldsymbol{\delta}^{(k-1)}, \nu^{(k-1)})$ using (9). Update $\boldsymbol{\delta}^{(k)} = \arg\min_{\boldsymbol{\delta}} \{\boldsymbol{\delta}^\top \mathbf{A} \boldsymbol{\delta} - 2\boldsymbol{\delta}^\top \mathbf{b} + \lambda_2 |\boldsymbol{\delta}|_1\}$.
- Compute $\mathbf{D} = \mathbf{D}(\boldsymbol{\Omega}^{(k-1)}, \boldsymbol{\delta}^{(k)}, \nu^{(k-1)})$ using (10). Update $\boldsymbol{\Omega}^{(k)} = \mathcal{S}(\boldsymbol{\Omega}^{(k-1)} - \frac{1}{\tau} \mathbf{D}, \frac{\lambda_1}{\tau})$.
- Update $\nu^{(k)} = \nu^{(k-1)} + \rho[M(\boldsymbol{\Omega}^{(k)}, \boldsymbol{\delta}^{(k)}) - 1]$.

Stop until $\max\{\rho|\boldsymbol{\Omega}^{(k)} - \boldsymbol{\Omega}^{(k-1)}|, \rho|\boldsymbol{\delta}^{(k)} - \boldsymbol{\delta}^{(k-1)}|, |\nu^{(k)} - \nu^{(k-1)}|/\rho\} \leq \epsilon$ for some pre-specified precision ϵ .

This is a modified version of the augmented Lagrangian method, where in the step of updating $\boldsymbol{\Omega}$, we minimize an upper envelope which is obtained by locally linearizing the augmented Lagrangian.

4 Estimation of Mean and Covariance Matrix

QUADRO requires estimates of the mean vector and covariance matrix for each class as inputs. We will show in Section 5 that the performance of QUADRO is closely related to the max-norm estimation error on mean vectors and covariance matrices. Sample mean and sample covariance matrix work well for Gaussian data. However, when data are from elliptical distributions, they may have inferior performance as we estimate nonpolynomially many of means and variances. In Sections 4.1-4.2, we suggest a robust M-estimator to estimate the mean and a rank-based estimator to estimate the covariance matrix, which are more appropriate for nonGaussian data. Moreover, in Section 4.3 we discuss how to estimate the model parameter γ when it is unknown.

4.1 Estimation of the mean

Suppose $\mathbf{x}_1, \dots, \mathbf{x}_n$ are *iid* samples of a random vector $\mathbf{X} = (X_1, \dots, X_d)^\top$ from an elliptical distribution $\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$. Let us denote $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)^\top$ and $\mathbf{x}_i = (x_{i1}, \dots, x_{id})^\top$ for $i = 1, \dots, n$. We estimate each μ_j marginally using the data $\{x_{1j}, \dots, x_{nj}\}$.

One possible estimator is the sample median

$$\hat{\mu}_{Mj} = \text{median}(\{x_{1j}, \dots, x_{nj}\}).$$

It can be shown that even under heavy-tailed distributions, $P(|\hat{\mu}_{Mj} - \mu_j| > A\sqrt{\log(\delta^{-1})/n}) \leq \delta$ for small $\delta \in (0, 1)$, where A is a constant determined by the probability density at μ_j , for each fixed j . This combined with the union bound gives that $|\hat{\mu}_M - \mu|_\infty = O_p(\sqrt{\log(d)/n})$.

Catoni (2012) proposed another M -estimator for the mean of heavy-tailed distributions. It works for distributions where mean is not necessarily equal to median. We denote the diagonal elements of the covariance matrix Σ as $\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2$, and the off-diagonal elements as σ_{kj} for $k \neq j$. The estimator $\hat{\mu}_C = (\hat{\mu}_{C,1}, \dots, \hat{\mu}_{C,d})^\top$ is obtained as follows. For a strictly increasing function $h : \mathbb{R} \rightarrow \mathbb{R}$ such that $-\log(1 - y + y^2/2) \leq h(y) \leq \log(1 + y + y^2/2)$, and a value $\delta \in (0, 1)$ such that $n > 2 \log(1/\delta)$, we let

$$\alpha_\delta = \left\{ \frac{2 \log(\delta^{-1})}{n \left[v + \frac{2v \log(\delta^{-1})}{n - 2 \log(\delta^{-1})} \right]} \right\}^{1/2},$$

where v is an upper bound of $\max\{\sigma_1^2, \dots, \sigma_d^2\}$. For each j , we define $\hat{\mu}_{Cj}$ as the unique value that satisfies

$$\sum_{i=1}^n h(\alpha_\delta(x_{ij} - \hat{\mu}_{Cj})) = 0.$$

It was shown in Catoni (2012) that $P(|\hat{\mu}_{Cj} - \mu_j| > \sqrt{\frac{2v \log(\delta^{-1})}{n(1 - 2 \log(\delta^{-1})/n)}) \leq \delta$ when the variance of X_j exists. Therefore, by taking $\delta = 1/(n \vee d)^2$, $|\hat{\mu}_M - \mu|_\infty \leq C\sqrt{\log(d)/n}$ with probability at least $1 - (n \vee d)^{-1}$, which gives the desired convergence rate.

To implement this estimator, we take $h(y) = \log(1 + y + y^2/2)$ for $y \geq 0$, and $h(y) = -\log(1 - y + y^2/2)$ for $y < 0$ in practice. For the choice of v , any value larger than $\max\{\sigma_1^2, \dots, \sigma_d^2\}$ would work in theory. Catoni (2012) introduced a Lepski's adaptation method to choose v . For simplicity, we take $v = 3 \max\{\tilde{\sigma}_1^2, \dots, \tilde{\sigma}_d^2\}$, where $\tilde{\sigma}_j^2$ is the sample covariance of X_j .

The two estimators, the median and the M -estimator, both have a convergence rate of $O_p(\sqrt{\log(d)/n})$ in terms of the max-norm error. In our numerical experiments, the M -estimator has a better numerical performance, and we stick to this estimator.

4.2 Estimation of the covariance matrix

To estimate the covariance matrix Σ , we estimate the marginal covariances $\{\sigma_j^2, 1 \leq j \leq d\}$ and the correlation matrix \mathbf{C} separately. Again, we need robust estimates even though the data have fourth moments, as we simultaneously estimate nonpolynomial number of covariance parameters.

First, we consider estimating σ_j^2 . Note that $\sigma_j^2 = \mathbb{E}(X_j^2) - \mathbb{E}^2(X_j)$. We estimate $\mathbb{E}(X_j^2)$ and $\mathbb{E}(X_j)$ separately. To estimate $\mathbb{E}(X_j^2)$, we use the M -estimator described above on the squared data $\{x_{1j}^2, \dots, x_{nj}^2\}$ and denote the estimator by $\hat{\eta}_{Cj}$. This works as in our setting, $\mathbb{E}(X_j^4)$ is finite for each j ; in addition, the M -estimator applies to asymmetric distributions. We then define

$$\hat{\sigma}_{Cj}^2 = \max\{\hat{\eta}_{Cj} - \hat{\mu}_{Cj}^2, \delta_0\},$$

where $\hat{\mu}_{C_j}$ is the M -estimator of $\mathbb{E}(X_j)$ and $\delta_0 > 0$ is a small constant ($\delta_0 < \min\{\sigma_1^2, \dots, \sigma_d^2\}$). It is easy to see that when the fourth moments of X_j are uniformly upper bounded by a constant and $n \geq 4 \log(d^2)$, $\max\{|\hat{\sigma}_{C_j} - \sigma_j|, 1 \leq j \leq d\} = O_p(\sqrt{\log(d)/n})$.

Next, we consider estimating the correlation matrix \mathbf{C} . For this, we use the Kendall's tau correlation matrix proposed by Han and Liu (2012). The Kendall's tau correlation coefficients (Kendall, 1938) are defined as

$$\tau_{jk} = \mathbb{P}((X_j - \tilde{X}_j)(X_k - \tilde{X}_k) > 0) - \mathbb{P}((X_j - \tilde{X}_j)(X_k - \tilde{X}_k) < 0),$$

where $\tilde{\mathbf{X}}$ is an independent copy of \mathbf{X} . They have the following relationship to the true coefficients: $C_{jk} = \sin(\frac{\pi}{2}\tau_{jk})$ for the elliptical family. Based on this equality, we first estimate the Kendall's tau correlation coefficients using rank-based estimators:

$$\hat{\tau}_{jk} = \begin{cases} \frac{2}{n(n-1)} \sum_{1 \leq i < i' \leq n} \text{sign}((x_{ij} - x_{i'j})(x_{ik} - x_{i'k})), & j \neq k, \\ 1 & j = k; \end{cases}$$

and then estimate the correlation matrix by $\hat{\mathbf{C}} = (\hat{C}_{jk})$ with

$$\hat{C}_{jk} = \sin\left(\frac{\pi}{2}\hat{\tau}_{jk}\right).$$

It is shown in Han and Liu (2012) that $|\hat{\mathbf{C}} - \mathbf{C}|_\infty = O_p(\sqrt{\log(d)/n})$.

Finally, we combine $\{\hat{\sigma}_j^2, 1 \leq j \leq d\}$ and $\hat{\mathbf{C}}$ to get $\tilde{\Sigma}$. Let

$$\tilde{\Sigma}_{jk} = \hat{\sigma}_j \hat{\sigma}_k \hat{C}_{jk}, \quad 1 \leq j, k \leq d.$$

It follows immediately that $|\tilde{\Sigma} - \Sigma|_\infty = O_p(\sqrt{\log(d)/n})$. However, this estimator is not necessarily positive semi-definite. To implement QUADRO, we need that $\hat{\Sigma}$ to be positive semi-definite so that the optimization in (8) is a convex problem. We obtain $\hat{\Sigma}$ by projecting $\tilde{\Sigma}$ onto the cone of positive semi-definite matrices through the convex optimization:

$$\hat{\Sigma} = \underset{\mathbf{A}: \mathbf{A} \text{ is positive semidefinite}}{\text{argmin}} \{|\mathbf{A} - \tilde{\Sigma}|_\infty\}. \quad (11)$$

Note that $|\hat{\Sigma} - \tilde{\Sigma}|_\infty \leq |\Sigma - \tilde{\Sigma}|_\infty$ by definition. Therefore, $|\hat{\Sigma} - \Sigma|_\infty \leq |\hat{\Sigma} - \tilde{\Sigma}|_\infty + |\tilde{\Sigma} - \Sigma|_\infty \leq 2|\tilde{\Sigma} - \Sigma|_\infty = O_p(\sqrt{\log(d)/n})$. To compute $\hat{\Sigma}$, we note that the optimization problem in (11) can be formulated as the dual of a graphical lasso problem corresponding to the smallest possible tuning parameter that still guarantees a feasible solution (Liu et al., 2012). Zhao et al. (2013) provides more algorithmic details.

4.3 Estimation of kurtosis parameter

When the kurtosis parameter γ is unknown, we could estimate it from data. Recall that $\gamma = \frac{1}{d(d+2)}\mathbb{E}(\xi^4) - 1$. Using the decomposition (4) and properties of \mathbf{U} , we have

$$\mathbb{E}(\xi^4) = \mathbb{E}\{[(\mathbf{X} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{X} - \boldsymbol{\mu})]^2\}.$$

Motivated by this equality, we propose the estimator

$$\hat{\gamma} = \max \left\{ \frac{1}{d(d+2)} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \tilde{\boldsymbol{\mu}})^\top \tilde{\boldsymbol{\Omega}} (\mathbf{x}_i - \tilde{\boldsymbol{\mu}}) - 1, \quad 0 \right\},$$

where $\tilde{\boldsymbol{\mu}}$ and $\tilde{\boldsymbol{\Omega}}$ are estimators of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}^{-1}$, respectively. Maruyama and Seo (2003) considered a similar estimator in low dimensional settings, where they used the sample mean and sample covariance matrix. In high dimensions, we need robust estimate to guarantee uniform convergence. In particular, we take $\tilde{\boldsymbol{\mu}} = \hat{\boldsymbol{\mu}}_C$ and $\tilde{\boldsymbol{\Omega}} = \hat{\boldsymbol{\Omega}}_{\text{clime}}$ where $\hat{\boldsymbol{\Omega}}_{\text{clime}}$ is the CLIME estimator proposed in Cai et al. (2011). We can also take the covariance estimator in Section 4.2, but then need to establish its sampling property as a precision matrix estimator. We decide to use the CLIME estimator since such a property has already been established by Cai et al. (2011). Denote by $\boldsymbol{\Sigma}^{-1} = (\Omega_{jk})_{d \times d}$. From simple algebra,

$$|\hat{\gamma} - \gamma| \leq \max_{1 \leq j, k \leq d} |\tilde{\mu}_j \tilde{\Omega}_{jk} \tilde{\mu}_k - \mu_j \Omega_{jk} \mu_k| \leq C \max \{ |\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}|_{\infty}, |\tilde{\boldsymbol{\Omega}} - \boldsymbol{\Sigma}^{-1}|_{\infty} \}.$$

In Section 4.1, we have seen that $\|\hat{\boldsymbol{\mu}}_C - \boldsymbol{\mu}\|_{\infty} = O_p(\sqrt{\log(d)/n})$. Moreover, Cai et al. (2011) showed that $|\tilde{\boldsymbol{\Omega}} - \boldsymbol{\Sigma}^{-1}|_{\infty} = \|\boldsymbol{\Sigma}^{-1}\|_1 \cdot O_p(\sqrt{\log(d)/n})$ under mild conditions, where $\|\cdot\|_1$ is the matrix L_1 -norm. Therefore, provided that $\|\boldsymbol{\Sigma}^{-1}\|_1 \leq C$, we immediately have $|\hat{\gamma} - \gamma| = O_p(\sqrt{\log(d)/n})$.

5 Theoretical Properties

In this section, we establish an oracle inequality for the Rayleigh quotient of the QUADRO estimates $(\hat{\boldsymbol{\Omega}}, \hat{\boldsymbol{\delta}})$. We assume that π and γ are known. For notational simplicity, we set $\lambda_1 = \lambda_2 = \lambda$. The results can be easily generalized to the case $\lambda_1 \neq \lambda_2$. Moreover, we drop the symmetry constraint $\boldsymbol{\Omega} = \boldsymbol{\Omega}^{\top}$ in all optimization problems involved. This simplifies the expression of the regularity conditions. The analysis with the symmetry constraint is a trivial extension of current analysis.

Recall the definition of M , L_1 and L_2 in (5) and $\kappa = (1 - \pi)/\pi$ and $L = L_1 + \kappa L_2$, the Rayleigh quotient of $(\boldsymbol{\Omega}, \boldsymbol{\delta})$ is equal to (up to a multiplicative constant)

$$R(\boldsymbol{\Omega}, \boldsymbol{\delta}) = \frac{[M(\boldsymbol{\Omega}, \boldsymbol{\delta})]^2}{L(\boldsymbol{\Omega}, \boldsymbol{\delta})}.$$

The QUADRO estimates are

$$(\hat{\boldsymbol{\Omega}}, \hat{\boldsymbol{\delta}}) = \underset{(\boldsymbol{\Omega}, \boldsymbol{\delta}): \widehat{M}(\boldsymbol{\Omega}, \boldsymbol{\delta})=1}{\operatorname{argmin}} \{ \widehat{L}(\boldsymbol{\Omega}, \boldsymbol{\delta}) + \lambda |\boldsymbol{\Omega}|_1 + \lambda |\boldsymbol{\delta}|_1 \}.$$

We shall compare the Rayleigh quotient of $(\hat{\boldsymbol{\Omega}}, \hat{\boldsymbol{\delta}})$ with the Rayleigh quotients of a class of “oracle solutions”. This class includes the one that maximizes the true Rayleigh quotient, which we denote by $(\boldsymbol{\Omega}_0^*, \boldsymbol{\delta}_0^*)$. Here we adopt a class of solutions as the “oracle” instead of only $(\boldsymbol{\Omega}_0^*, \boldsymbol{\delta}_0^*)$, because we want the results not tied to the sparsity assumption on $(\boldsymbol{\Omega}_0^*, \boldsymbol{\delta}_0^*)$ but a weaker assumption: at least one solution in this class is sparse.

Our theoretical development is technically nontrivial. Conventional oracle inequalities are derived in a setting of minimizing a data-dependent loss without constraint, and the risk function is the expectation of the loss. Here we minimize a data-dependent loss with a data-dependent equality constraint, and the risk function — the Rayleigh quotient — is not equal to the expectation of the loss. A similar setting was considered in Fan et al. (2012), where they introduced a data-dependent intermediate solution to deal with such equality constraint. But the rate they

obtained depends on this intermediate solution, which is very hard to quantify. In contrast, the rate in our results purely depends on the oracle solution. To get rid of the intermediate solution in the rate, we need to carefully quantify its difference from both the QUADRO solution and the oracle solution. The technique is new, and potentially useful for other problems.

5.1 Oracle solutions, the restricted eigenvalue condition

For any $\lambda_0 \geq 0$, we define the *oracle solution associated with λ_0* to be

$$(\boldsymbol{\Omega}_{\lambda_0}^*, \boldsymbol{\delta}_{\lambda_0}^*) = \underset{(\boldsymbol{\Omega}, \boldsymbol{\delta}): M(\boldsymbol{\Omega}, \boldsymbol{\delta})=1}{\operatorname{argmin}} \{L(\boldsymbol{\Omega}, \boldsymbol{\delta}) + \lambda_0 |\boldsymbol{\Omega}|_1 + \lambda_0 |\boldsymbol{\delta}|_1\}. \quad (12)$$

We shall compare the Rayleigh quotient of $(\widehat{\boldsymbol{\Omega}}, \widehat{\boldsymbol{\delta}})$ to that of $(\boldsymbol{\Omega}_{\lambda_0}^*, \boldsymbol{\delta}_{\lambda_0}^*)$, for an arbitrary λ_0 . In particular, when $\lambda_0 = 0$, the associated oracle solution (may not be unique) becomes

$$(\boldsymbol{\Omega}_0^*, \boldsymbol{\delta}_0^*) = \underset{(\boldsymbol{\Omega}, \boldsymbol{\delta}): M(\boldsymbol{\Omega}, \boldsymbol{\delta})=1}{\operatorname{argmin}} \{L(\boldsymbol{\Omega}, \boldsymbol{\delta})\}.$$

It maximizes the true Rayleigh quotient.

Next, we introduce a restricted eigenvalue (RE) condition jointly on $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2, \boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$. For any matrices \mathbf{A} and \mathbf{B} , let $\operatorname{vec}(\mathbf{A})$ be the vectorization of \mathbf{A} by stacking all the elements of \mathbf{A} column by column, and $\mathbf{A} \otimes \mathbf{B}$ be the Kronecker product of \mathbf{A} and \mathbf{B} . We define the matrices

$$\mathbf{Q}_k = \begin{bmatrix} (2(1 + \gamma)\boldsymbol{\Sigma}_k + 4\boldsymbol{\mu}_k\boldsymbol{\mu}_k^\top) \otimes \boldsymbol{\Sigma}_k + \gamma \operatorname{vec}(\boldsymbol{\Sigma}_k) \operatorname{vec}(\boldsymbol{\Sigma}_k)^\top & -4\boldsymbol{\mu}_k \otimes \boldsymbol{\Sigma}_k \\ -4\boldsymbol{\mu}_k^\top \otimes \boldsymbol{\Sigma}_k & 4\boldsymbol{\Sigma}_k \end{bmatrix}, \quad k = 1, 2.$$

We note that there are $(d^2 + d)$ coefficients to decide when maximizing $R(\boldsymbol{\Omega}, \boldsymbol{\delta})$: d^2 elements of $\boldsymbol{\Omega}$ and d elements of $\boldsymbol{\delta}$. We can stack all these coefficients into a long vector $\mathbf{x} = \mathbf{x}(\boldsymbol{\Omega}, \boldsymbol{\delta})$ in \mathbb{R}^{d^2+d} defined as:

$$\mathbf{x}(\boldsymbol{\Omega}, \boldsymbol{\delta}) \equiv \left[\operatorname{vec}(\boldsymbol{\Omega})^\top, \boldsymbol{\delta}^\top \right]^\top. \quad (13)$$

It can be shown that $L_k(\boldsymbol{\Omega}, \boldsymbol{\delta}) = \mathbf{x}^\top \mathbf{Q}_k \mathbf{x}$, for $k = 1, 2$ (see Lemma 9.2). Therefore, $L(\boldsymbol{\Omega}, \boldsymbol{\delta}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x}$, where $\mathbf{Q} = \mathbf{Q}_1 + \kappa \mathbf{Q}_2$. Our RE condition is then imposed on the $(d^2 + d) \times (d^2 + d)$ matrix \mathbf{Q} , and hence implicitly on $(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2)$.

We now formally introduce the RE condition. For a set $S \subset \{1, 2, \dots, d^2 + d\}$ and a nonnegative value \bar{c} , we define the *restricted eigenvalue*:

$$\Theta(S; \bar{c}) = \min_{\mathbf{v}: |\mathbf{v}_{S^c}|_1 \leq \bar{c} |\mathbf{v}_S|_1} \frac{\mathbf{v}^\top \mathbf{Q} \mathbf{v}}{|\mathbf{v}_S|^2}.$$

Generally speaking, $\Theta(S; \bar{c})$ depends on $(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2)$ in a complicated way. For $\bar{c} = 0$, the following proposition builds a connection between $\Theta(S; 0)$ and $(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2)$. For each $S \subset \{1, 2, \dots, d^2 + d\}$, there exist sets $U \subset \{1, \dots, d\} \times \{1, \dots, d\}$ and $V \subset \{1, \dots, d\}$ such that the support of $\mathbf{x}(\boldsymbol{\Omega}, \boldsymbol{\delta})$ is S if and only if the support of $\boldsymbol{\Omega}$ is U and the support of $\boldsymbol{\delta}$ is V . Let

$$U' = \cup_{(i,j) \in U} \{i, j\},$$

and it is easy to see that $U \subset U' \times U'$. The following result is proved in Section 9.

Proposition 5.1. For any set $S \subset \{1, \dots, d^2 + d\}$, suppose U' and V are defined as above. Let $\tilde{\Sigma}_k$ be the submatrix of Σ_k by restricting rows and columns to $U' \cup V$, $\tilde{\mu}_k$ be the subvector of μ_k by constraining elements to $U' \cup V$, for $k = 1, 2$. If there exist constants $v_1, v_2 > 0$ such that $\lambda_{\min}(\tilde{\Sigma}_k - v_1 \tilde{\mu}_k \tilde{\mu}_k^\top) \geq \frac{1}{2} \lambda_{\min}(\tilde{\Sigma}_k) \geq \frac{v_2}{2}$ for $k = 1, 2$, then

$$\Theta(S, 0) \geq (1 + \gamma)(1 + \kappa)v_2 \min \left\{ v_2, \frac{4v_1}{2 + v_1(1 + \gamma)} \right\} > 0.$$

5.2 Oracle inequality on Rayleigh quotient

We shall assume that $\max\{|\Sigma_k|_\infty, |\mu_k|_\infty, k = 1, 2\} \leq 1$, and $|\hat{\Sigma}_k - \Sigma_k|_\infty \leq |\Sigma_k|_\infty$, $|\hat{\mu}_k - \mu_k|_\infty \leq |\mu_k|_\infty$ for $k = 1, 2$, without loss of generality. For any $\lambda_0 \geq 0$, let $(\Omega_{\lambda_0}^*, \delta_{\lambda_0}^*)$ be the associated oracle solution and S be the support of $\mathbf{x}_{\lambda_0}^* = [\text{vec}(\Omega_{\lambda_0}^*)^\top, (\delta_{\lambda_0}^*)^\top]^\top$. Let $\Delta_n = \max\{|\hat{\Sigma}_k - \Sigma_k|_\infty, |\hat{\mu}_k - \mu_k|_\infty, k = 1, 2\}$. We have the following result for any given estimators, which proof can be found in Section 9.

Theorem 5.1. Given $\lambda_0 \geq 0$, let S be the support of $\mathbf{x}_{\lambda_0}^*$, $s_0 = |S|$ and $k_0 = \max\{s_0, R(\Omega_{\lambda_0}^*, \delta_{\lambda_0}^*)\}$. Suppose that $\Theta(S, 0) \geq c_0$, $\Theta(S, 3) \geq a_0$, and $R(\Omega_{\lambda_0}^*, \delta_{\lambda_0}^*) \geq u_0$, for some positive constants a_0 , c_0 and u_0 . We assume $4s_0\Delta_n^2 \leq a_0c_0$ and $\max\{s_0\Delta_n, s_0^{1/2}k_0^{1/2}\lambda_0\} < 1$ without loss of generality. Then, there exist positive constants $C = C(a_0, c_0, u_0)$ and $A = A(a_0, c_0, u_0)$ such that for any $\eta > 1$, by taking $\lambda = C\eta \max\{s_0^{1/2}\Delta_n, k_0^{1/2}\lambda_0\} [R(\Omega_{\lambda_0}^*, \delta_{\lambda_0}^*)]^{-1/2}$,

$$\frac{R(\hat{\Omega}, \hat{\delta})}{R(\Omega_{\lambda_0}^*, \delta_{\lambda_0}^*)} \geq 1 - A\eta^2 \max\{s_0\Delta_n, s_0^{1/2}k_0^{1/2}\lambda_0\}.$$

In Theorem 5.1, the rate of convergence has two parts. The term $s_0\Delta_n$ reflects how the stochastic errors of estimating $(\Sigma_1, \Sigma_2, \mu_1, \mu_2)$ affect the Rayleigh quotient. The term $s_0^{1/2}k_0^{1/2}\lambda_0$ is an extra term that depends on the oracle solution we aim to compare with. In particular, if we compare $R(\hat{\Omega}, \hat{\delta})$ with $R_{\max} \equiv R(\Omega_0^*, \delta_0^*)$, the population maximum Rayleigh quotient with $\lambda_0 = 0$, this extra term disappears. If we further use the estimators in Section 4, $\Delta_n = O_p(\sqrt{\log(d)/n})$. We summarize the result as follows.

Corollary 5.1. Suppose that Theorem 5.1 holds with $\lambda_0 = 0$, then for some positive constants A and C , when $\lambda > Cs_0^{1/2}R_{\max}^{-1/2}\Delta_n$, we have

$$R(\hat{\Omega}, \hat{\delta}) \geq (1 - As_0\Delta_n)R_{\max}.$$

Furthermore, if the mean vectors and covariance matrices are estimated by using the robust methods in Section 4, then when $\lambda > Cs_0^{1/2}R_{\max}^{-1/2}\sqrt{\log(d)/n}$,

$$R(\hat{\Omega}, \hat{\delta}) \geq (1 - As_0\sqrt{\log(d)/n})R_{\max},$$

with probability at least $1 - (n \vee d)^{-1}$.

From Corollary 5.1, when (Ω_0^*, δ_0^*) is truly sparse, $R(\hat{\Omega}, \hat{\delta})$ is close to the population maximum Rayleigh quotient R_{\max} . However, we note that Theorem 5.1 considers more general situations including cases where (Ω_0^*, δ_0^*) is not sparse. As long as there exists an ‘‘approximately optimal’’ and sparse solution, i.e. for a small λ_0 the associated oracle solution $(\Omega_{\lambda_0}^*, \delta_{\lambda_0}^*)$ is sparse, Theorem 5.1 guarantees that $R(\hat{\Omega}, \hat{\delta})$ is close to $R(\Omega_{\lambda_0}^*, \delta_{\lambda_0}^*)$, and hence close to R_{\max} .

6 Application to Classification

One application of QUADRO is high-dimensional classification for elliptically-distributed data. Suppose $(\widehat{\boldsymbol{\Omega}}, \widehat{\boldsymbol{\delta}})$ are the QUADRO estimates. It yields a classification rule

$$\widehat{h}(\mathbf{x}) = I\{\mathbf{x}^\top \widehat{\boldsymbol{\Omega}} \mathbf{x} - 2\widehat{\boldsymbol{\delta}}^\top \mathbf{x} < c\}.$$

In this section, we first show that the Rayleigh quotient is a proxy of the classification error and then derive an analytic choice of c . Comparing with many other high-dimensional classification methods, QUADRO produces quadratic boundaries and can handle both nonGaussian distributions and non-equal covariance matrices.

6.1 Approximation of classification errors

Given $(\boldsymbol{\Omega}, \boldsymbol{\delta})$ and a threshold c , a general quadratic rule $h(\mathbf{x}) = h(\mathbf{x}; \boldsymbol{\Omega}, \boldsymbol{\delta}, c)$ is defined as

$$h(\mathbf{x}; \boldsymbol{\Omega}, \boldsymbol{\delta}, c) = I\{\mathbf{x}^\top \boldsymbol{\Omega} \mathbf{x} - 2\mathbf{x}^\top \boldsymbol{\delta} < c\}. \quad (14)$$

We reparametrize c as

$$c = tM_1(\boldsymbol{\Omega}, \boldsymbol{\delta}) + (1-t)M_2(\boldsymbol{\Omega}, \boldsymbol{\delta}). \quad (15)$$

Here $M_k(\boldsymbol{\Omega}, \boldsymbol{\delta}) = \boldsymbol{\mu}_k^\top \boldsymbol{\Omega} \boldsymbol{\mu}_k - 2\boldsymbol{\mu}_k^\top \boldsymbol{\delta} + \text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma}_k)$ is the mean of $Q(\mathbf{X})$ in class k , for $k = 1, 2$. After the reparametrization, t is *scale-free*. As we will see below, in most cases, given $\boldsymbol{\Omega}$ and $\boldsymbol{\delta}$, the optimal t that minimizes the classification error takes values on $(0, 1)$.

From now on, we write $h(\mathbf{x}; \boldsymbol{\Omega}, \boldsymbol{\delta}, c) = h(\mathbf{x}; \boldsymbol{\Omega}, \boldsymbol{\delta}, t)$. Let $\text{Err}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t)$ be the classification error of $h(\cdot; \boldsymbol{\Omega}, \boldsymbol{\delta}, t)$. Due to technical difficulties, we only give results for Gaussian distributions. Suppose $\mathbf{X}|(Y=0) \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $\mathbf{X}|(Y=1) \sim \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$. For $k = 1, 2$, we write

$$\boldsymbol{\Sigma}_k^{1/2} \boldsymbol{\Omega} \boldsymbol{\Sigma}_k^{1/2} = \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^\top,$$

where \mathbf{S}_k is a diagonal matrix containing the nonzero eigenvalues, and the columns of \mathbf{K}_k are corresponding eigenvectors. Let $\boldsymbol{\beta}_k = \mathbf{K}_k^\top \boldsymbol{\Sigma}_k (\boldsymbol{\Omega} \boldsymbol{\mu}_k - \boldsymbol{\delta})$. When $\max\{|\mathbf{S}_k|_\infty, |\boldsymbol{\beta}_k|_\infty, k = 1, 2\}$ is bounded, the following proposition shows that an approximation of $\text{Err}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t)$ is

$$\overline{\text{Err}}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) \equiv \pi \bar{\Phi} \left(\frac{(1-t)M(\boldsymbol{\Omega}, \boldsymbol{\delta})}{\sqrt{L_1(\boldsymbol{\Omega}, \boldsymbol{\delta})}} \right) + (1-\pi) \bar{\Phi} \left(\frac{tM(\boldsymbol{\Omega}, \boldsymbol{\delta})}{\sqrt{L_2(\boldsymbol{\Omega}, \boldsymbol{\delta})}} \right),$$

where M , L_1 and L_2 are defined in (5), Φ is the distribution function of a standard normal variable and $\bar{\Phi} = 1 - \Phi$.

Proposition 6.1. *Suppose that $\max\{|\mathbf{S}_k|_\infty, |\boldsymbol{\beta}_k|_\infty, k = 1, 2\} \leq C_0$ for some constant $C_0 > 0$ and let q be the rank of $\boldsymbol{\Omega}$. Then as d goes to infinity,*

$$|\text{Err}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) - \overline{\text{Err}}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t)| = \frac{O(q) + o(d)}{[\min\{L_1(\boldsymbol{\Omega}, \boldsymbol{\delta}), L_2(\boldsymbol{\Omega}, \boldsymbol{\delta})\}]^{3/2}}.$$

In particular, if we consider all such $(\boldsymbol{\Omega}, \boldsymbol{\delta})$ that the variance of $Q(\mathbf{X}; \boldsymbol{\Omega}, \boldsymbol{\delta})$ under both classes are lower bounded by $c_0 d^\theta$ for some constants $\theta > 2/3$ and $c_0 > 0$, then we have $|\text{Err} - \overline{\text{Err}}| = o(1)$.

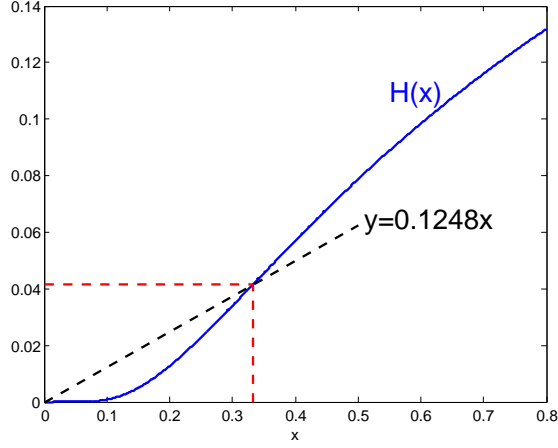


Figure 2: Function $H(x) = \bar{\Phi}(1/\sqrt{x})$.

We now take a closer look at $\bar{\text{Err}}$. Let $H(x) = \bar{\Phi}(1/\sqrt{x})$, which is monotone increasing on $(0, \infty)$. Writing for short $M = M_1 - M_2$, $M_k = M_k(\boldsymbol{\Omega}, \boldsymbol{\delta})$ and $L_k = L_k(\boldsymbol{\Omega}, \boldsymbol{\delta})$ for $k = 1, 2$, we have

$$\bar{\text{Err}}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) = \pi H\left(\frac{L_1}{(1-t)^2 M^2}\right) + (1-\pi)H\left(\frac{L_2}{t^2 M^2}\right).$$

Figure 2 shows that $H(\cdot)$ is nearly linear on an important range. This suggests the following approximation

$$\bar{\text{Err}}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) \approx H\left(\pi \frac{L_1}{(1-t)^2 M^2} + (1-\pi) \frac{L_2}{t^2 M^2}\right) = H\left(\frac{\pi R^{(t)}}{(1-t)^2}\right), \quad (16)$$

where $R^{(t)} = R^{(t)}(\boldsymbol{\Omega}, \boldsymbol{\delta})$ is the $R(\boldsymbol{\Omega}, \boldsymbol{\delta})$ in (6) corresponding to the κ value

$$\kappa(t) \equiv \frac{1-\pi}{\pi} \frac{(1-t)^2}{t^2}.$$

The approximation in (16) is quantified in the following proposition.

Proposition 6.2. *Given $(\boldsymbol{\Omega}, \boldsymbol{\delta}, t)$, we write $R_k = R_k(\boldsymbol{\Omega}, \boldsymbol{\delta}) = [M(\boldsymbol{\Omega}, \boldsymbol{\delta})]^2 / L_k(\boldsymbol{\Omega}, \boldsymbol{\delta})$, for $k = 1, 2$, and define*

$$V_1 = V_1(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) = \min\left\{(1-t)^2 R_1, \frac{1}{(1-t)^2 R_1}\right\},$$

$$V_2 = V_2(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) = \min\left\{t^2 R_2, \frac{1}{t^2 R_2}\right\},$$

$$V = V(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) = \max\{V_1/V_2, V_2/V_1\}.$$

Then there exists a constant $C > 0$ such that

$$\left|\bar{\text{Err}}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) - H\left(\frac{\pi R^{(t)}(\boldsymbol{\Omega}, \boldsymbol{\delta})}{(1-t)^2}\right)\right| \leq C [\max\{V_1, V_2\}]^{1/2} \cdot |V - 1|^2.$$

In particular, when $t = 1/2$,

$$\left|\bar{\text{Err}}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) - H\left(\frac{\pi R^{(t)}(\boldsymbol{\Omega}, \boldsymbol{\delta})}{(1-t)^2}\right)\right| \leq C R_0^{1/2} \cdot \left(\frac{\Delta R}{R_0}\right)^2,$$

where $R_0 = \max\{\min\{R_1, 1/R_1\}, \min\{R_2, 1/R_2\}\}$ and $\Delta R = |R_1 - R_2|$.

Note that L_1 and L_2 are the variances of $Q(\mathbf{X}) = \mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X} - 2\mathbf{X}^\top \boldsymbol{\delta}$ for two classes, respectively. In cases where $|L_1 - L_2| \ll \min\{L_1, L_2\}$, $\Delta R \ll R_0$. Also, R_0 is always bounded by 1; and it tends to 0 in many situations, for example, when $R_1, R_2 \rightarrow \infty$, or $R_1, R_2 \rightarrow 0$, or $R_1 \rightarrow 0, R_2 \rightarrow \infty$. Proposition 6.2 then implies that the approximation in (16) when $t = 1/2$ is good.

Combining Propositions 6.1 and 6.2, the classification error of a general quadratic rule $h(\cdot; \boldsymbol{\Omega}, \boldsymbol{\delta}, t)$ is approximately a monotone decreasing transform of the Rayleigh quotient $R^{(t)}(\boldsymbol{\Omega}, \boldsymbol{\delta})$, corresponding to $\kappa = \kappa(t)$. In particular, when $t = 1/2$ (i.e., $c = (M_1 + M_2)/2$), $R^{(1/2)}(\boldsymbol{\Omega}, \boldsymbol{\delta})$ is exactly the one used in QUADRO. Consequently, if we fix the threshold to be $c = (M_1 + M_2)/2$, then the Rayleigh quotient (upon with a monotone transform) is a good proxy for classification error. This explains why Rayleigh-quotient based procedures can be used for classification.

We remark that even in the region that $H(\cdot)$ is far from being linear such that the upper bound in Proposition 6.2 is not $o(1)$, we can still find a monotone transform of the Rayleigh quotient as an *upper bound* of the classification error. To see this, note that for $x \in [1/3, \infty)$, $H(x)$ is a concave function. Therefore, the approximation in (16) becomes an inequality, i.e., $\overline{\text{Err}}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) \leq H\left(\frac{\pi R^{(t)}}{(1-t)^2}\right)$. For $x \in (0, 1/3)$, $H(x) \leq 0.1248x$. It follows that $\overline{\text{Err}}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) \leq 0.1248 \cdot \frac{\pi R^{(t)}}{(1-t)^2}$.

6.2 QUADRO as a classification method

Results in Section 6.1 suggest an analytic method to choose the threshold c , or equivalently t , with given $(\boldsymbol{\Omega}, \boldsymbol{\delta})$. Let

$$\hat{t} \in \min_t \left\{ \pi \bar{\Phi} \left(\frac{(1-t)\widehat{M}(\boldsymbol{\Omega}, \boldsymbol{\delta})}{\sqrt{\widehat{L}_1(\boldsymbol{\Omega}, \boldsymbol{\delta})}} \right) + (1-\pi) \bar{\Phi} \left(\frac{t\widehat{M}(\boldsymbol{\Omega}, \boldsymbol{\delta})}{\sqrt{\widehat{L}_2(\boldsymbol{\Omega}, \boldsymbol{\delta})}} \right) \right\} \quad (17)$$

and set

$$\hat{c} = (1-\hat{t})\widehat{M}_1(\boldsymbol{\Omega}, \boldsymbol{\delta}) + \hat{t}\widehat{M}_2(\boldsymbol{\Omega}, \boldsymbol{\delta}). \quad (18)$$

Here (17) is a one-dimensional optimization problem, and can be solved easily. The resulting QUADRO classification rule is

$$\hat{h}^{\text{Quad}}(\mathbf{x}) = I\{\mathbf{x}^\top \widehat{\boldsymbol{\Omega}} \mathbf{x} - 2\mathbf{x}^\top \widehat{\boldsymbol{\delta}} - \hat{c} < 0\}.$$

As a by-product, the method to decide c , described in (17) and (18), can be used in other classification procedures on Gaussian data, such as logistic regression, quadratic discriminant analysis (QDA) and kernel support vector machine, once $(\widehat{\boldsymbol{\Omega}}, \widehat{\boldsymbol{\delta}})$ are given. It provides a fast and purely data-driven way to decide the threshold value in quadratic classification rules. In our numerical experiments, it performs well.

7 Numerical Studies

In this section, we investigate the performance of QUADRO in several simulation examples and a real data example. The simulation studies contain both Gaussian models and general elliptical models. We compare QUADRO with several *classification-oriented procedures*. Performances are evaluated in terms of both the Rayleigh quotient and classification error.

Table 1: Classification error and Rayleigh quotient for Gaussian examples. Means and standard deviations (in the parenthesis) of 100 replications are reported.

	Classification error			
	QUADRO	SLR	L-SLR	ROAD
Model 1	0.179 (0.016)	0.235 (0.028)	0.191 (0.017)	0.246 (0.074)
Model 2	0.144 (0.016)	0.224 (0.018)	0.470 (0.008)	0.491 (0.010)
Model 3	0.109 (0.013)	0.164 (0.018)	0.176 (0.016)	0.235 (0.068)
	Rayleigh quotient			
	QUADRO	SLR	L-SLR	ROAD
Model 1	3.016 (0.405)	1.874 (0.499)	2.897 (0.407)	2.193 (1.055)
Model 2	3.081 (0.389)	1.508 (0.300)	0 (0)	0 (0)
Model 3	5.377 (0.628)	2.681 (0.385)	3.027 (0.437)	2.184 (1.173)

7.1 Simulations under Gaussian models

Let $n_1 = n_2 = 50$ and $d = 40$. For each given $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_1$ and $\boldsymbol{\Sigma}_2$, we generate 100 training datasets independently, each with n_1 data from $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and n_2 data from $\mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$. In QUADRO, we input the sample means and sample covariance matrices. We set $\lambda_2 = r\lambda_1$ and work with λ_1 and r from now on. The two tuning parameters $\lambda_1 \geq 0$ and $r > 0$ are selected in the following way. For various pairs of (λ_1, r) , we apply QUADRO for each pair and evaluate the criteria (Rayleigh quotient or classification error) via 4,000 newly generated testing data; we then choose the (λ_1, r) that optimize the criteria.

We compare QUADRO with three *classification-oriented procedures*:

- Sparse Logistic Regression (SLR): We apply the sparse logistic regression to the augmented feature space $\{X_i, 1 \leq i \leq d; X_i X_j, 1 \leq i \leq j \leq d\}$. The resulting estimator then gives a quadratic projection with $(\boldsymbol{\Omega}, \boldsymbol{\delta}, c)$ decided from the fitted regression coefficients. We implement the sparse logistic regression using the R package glmnet.
- Linear Sparse Logistic Regression (L-SLR): We apply the sparse logistic regression directly to the original feature space $\{X_i, 1 \leq i \leq d\}$.
- ROAD: ROAD (Fan et al., 2012) is a linear classification method for Gaussian distributed data. It is equivalent to a modified version of QUADRO by enforcing $\boldsymbol{\Omega}$ to be the zero matrix. Hence, ROAD only produces linear rules.

To make a fair comparison, the tuning parameters in SLR and L-SLR are selected in the same way as in QUADRO based on 4,000 testing data. ROAD is self-tuned by its package.

We consider three models:

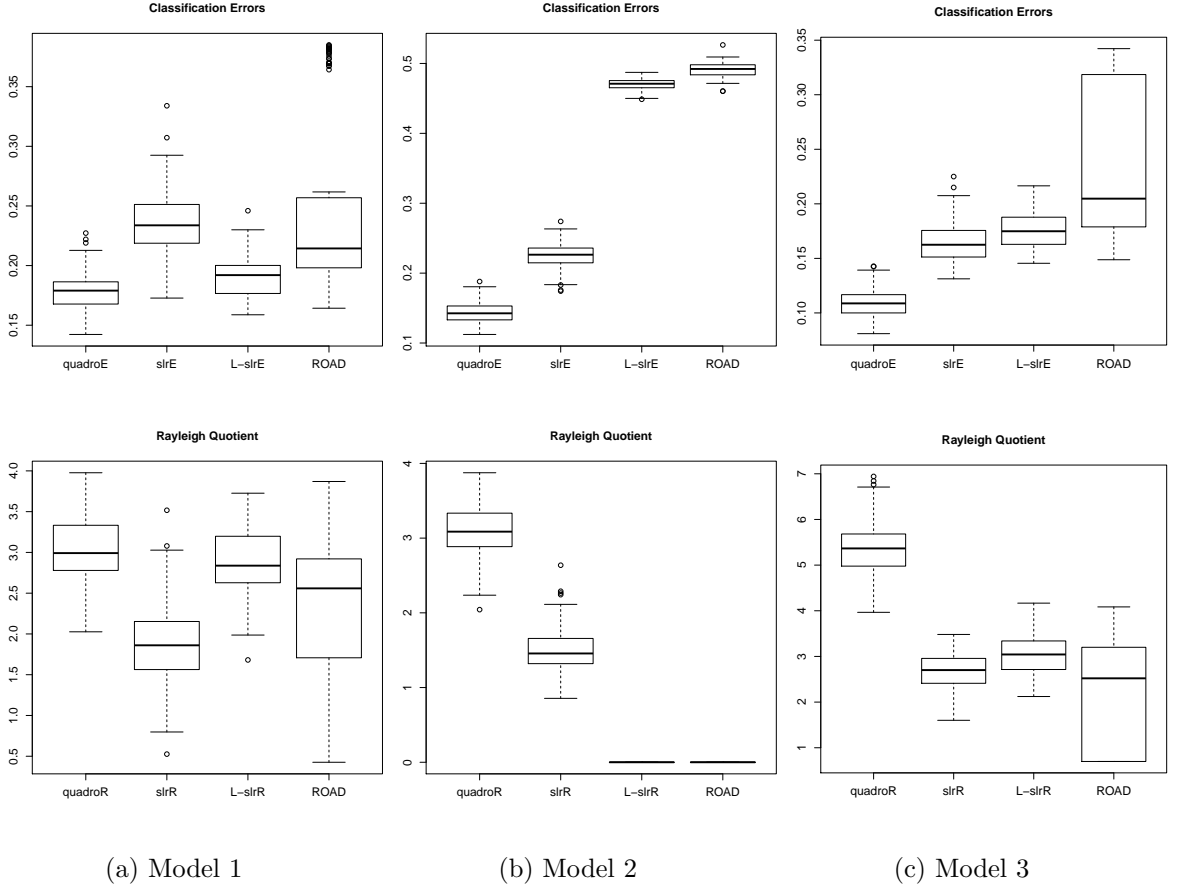


Figure 3: Distributions of minimum classification errors (top panel) and maximum Rayleigh quotients (bottom panel) based on 100 replications for three different normal models. The tuning parameters for quadroE, slrE and L-slrE, are chosen to minimize the classification errors. They are also chosen to maximize Rayleigh quotients for quadroR, slrR, and L-slrR.

- *Model 1*: Σ_1 is the identity matrix. Σ_2 is a diagonal matrix in which the first 10 elements are equal to 1.3 and the rest are equal to 1. $\mu_1 = \mathbf{0}$, and $\mu_2 = (0.7, \dots, 0.7, 0, \dots, 0)^\top$ with the first 10 elements of μ_2 being nonzero.
- *Model 2*: Σ_1 is a block-diagonal matrix. Its upper left 20×20 block is an equal correlation matrix with $\rho = 0.4$; and its lower right 20×20 block is an identity matrix. $\Sigma_2 = (\Sigma_1^{-1} + \mathbf{I})^{-1}$. We also set $\mu_1 = \mu_2 = \mathbf{0}$. In this model, neither Σ_1^{-1} nor Σ_2^{-1} is sparse, but $\Sigma_1^{-1} - \Sigma_2^{-1}$ is.
- *Model 3*: Σ_1 , Σ_2 and μ_1 are the same as in *Model 2* and μ_2 is taken from *Model 1*.

Table 1 reports the classification errors and Rayleigh quotients of all methods. Figures 3 contains the corresponding boxplots. In all three models, QUADRO outperforms other methods in terms of both classification error and Rayleigh quotient. In Model 2, since $\mu_1 = \mu_2$, no linear methods can do better in classification than the random guess. Therefore, both ROAD and L-SLR have poor performances in terms of classification error. In Models 1 and 3, $\mu_1 \neq \mu_2$ and $\Sigma_1 \neq \Sigma_2$. So, in the Bayes classification rule, both “linear” parts and “quadratic” parts play important roles. Comparing SLR and L-SLR, we see the former considers a broader class, while

Table 2: Classification error and Rayleigh quotient for elliptical examples. Means and standard deviations (in the parenthesis) of 100 replications are reported.

	Classification Error			
	QUADRO	QUADRO-0	SLR	L-SLR
Model 4	0.136 (0.015)	0.144 (0.015)	0.167 (0.019)	0.157 (0.017)
Model 5	0.161 (0.012)	0.173 (0.014)	0.184 (0.014)	0.184 (0.016)
Model 6	0.130 (0.017)	0.129 (0.016)	0.152 (0.018)	0.211 (0.018)
	Rayleigh quotient			
	QUADRO	QUADRO-0	SLR	L-SLR
Model 4	3.179 (0.413)	2.975 (0.410)	1.984 (0.511)	2.846 (0.450)
Model 5	2.415 (0.224)	2.191 (0.274)	1.625 (0.329)	2.166 (0.276)
Model 6	2.374 (0.194)	2.160 (0.213)	1.363 (0.246)	1.669 (0.223)

the latter is more robust, but neither of them perform uniformly well. However, QUADRO performs well in all cases.

7.2 Simulations under elliptical models

Let $n_1 = n_2 = 50$ and $d = 40$. For each given $\boldsymbol{\mu}_1$, $\boldsymbol{\mu}_2$, $\boldsymbol{\Sigma}_1$ and $\boldsymbol{\Sigma}_2$, data are generated from multivariate t distribution with degrees of freedom 5. In QUADRO, we input the robust M-estimators for means and the rank-based estimators for covariance matrices as described in Section 4. We compare the performance of QUADRO with that of SLR and L-SLR. ROAD is not designed for elliptical models, so we do not report its results here. We also implement QUADRO with inputs of sample means and sample covariance matrices. We name this method QUADRO-0 to differentiate it from QUADRO.

We consider three models:

- *Model 4*: The same parameters as those in *Model 1*.
- *Model 5*: $\boldsymbol{\Sigma}_1$, $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$ are the same as in *Model 1*. $\boldsymbol{\Sigma}_2$ is the covariance matrix of a fractional white noise process, where the difference parameter $l = 0.2$. In other words, $\boldsymbol{\Sigma}_2$ has a polynomial off-diagonal decay: $|\Sigma_2(i, j)| = O(|i - j|^{1-2l})$.
- *Model 6*: $\boldsymbol{\Sigma}_1$, $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$ are the same as in *Model 1*. $\boldsymbol{\Sigma}_2$ is a matrix such that $\Sigma_2(i, j) = 0.6^{|i-j|}$, i.e., $\boldsymbol{\Sigma}_2$ has an exponential off-diagonal decay.

Table 2 presents the average classification error and Rayleigh quotient over 100 replications. Figure 4 contains the boxplots. We see that in all models, QUADRO outperforms SLR and L-SLR. In addition, QUADRO is better than QUADRO-0, which illustrates the advantage of using the robust M-estimators for means and the rank-based estimators for covariance matrices.

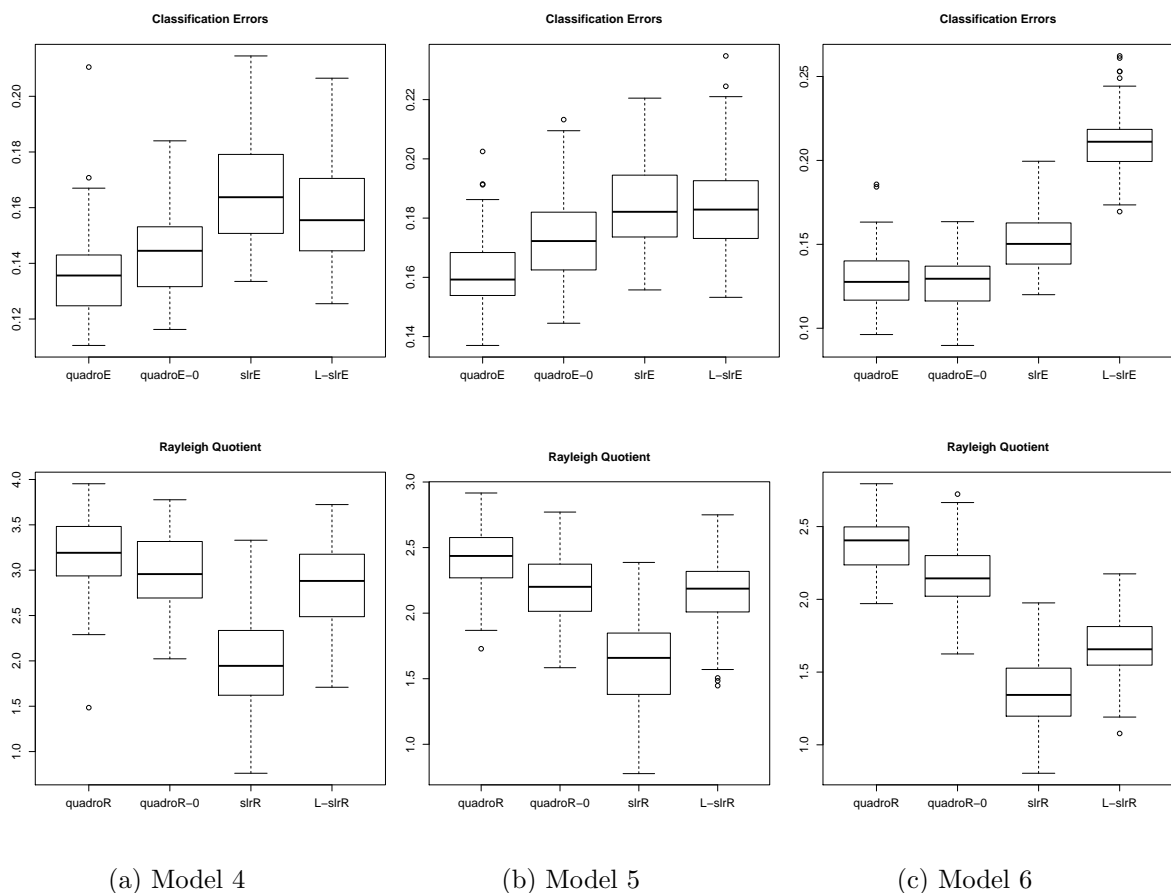


Figure 4: Distributions of minimum classification errors (top panel) and maximum Rayleigh quotients (bottom panel) based on 100 replications across different elliptical distribution models. The tuning parameters for quadroE, slrE and L-slrE, are chosen to minimize the classification errors. They are also chosen to maximize Rayleigh quotients for quadroR, slrR, and L-slrR.

7.3 Real data analysis

We apply QUADRO to a large-scale genomic dataset, GPL96, and compare the performance of QUADRO with SLR and L-SLR. The GPL96 data set contains 20,263 probes and 8,124 samples from 309 tissues. Among the tissues, breast tumor has 1,142 samples, which is the largest set. We merge the probes from the same gene by averaging them, and finally get 12,679 genes and 8,124 samples. We divide all samples into two groups: breast tumor or non-breast tumor.

First, we look at the classification errors. We replicate our experiment 100 times. Each time, with QUADRO, SLR and L-SLR we proceed with the following steps.

- Randomly choose a training set of 400 samples, 200 from breast tumor and 200 from non-breast tumor.
- For each training set, we use half of the samples to compute $(\hat{\Omega}, \hat{\delta})$ and the other half to select the tuning parameters by minimizing the classification error.
- Use the rest 942 samples from breast tumor and another randomly chosen 942 samples from non-breast tumor as testing set, and calculate the testing error.

Table 3: Classification errors on GPL96 dataset, across methods QUADRO, SLR and L-SLR. Means and standard deviations (in the parenthesis) of 100 replications are reported.

Classification Errors		
QUADRO	SLR	L-SLR
0.014	0.025	0.025
(0.007)	(0.007)	(0.009)

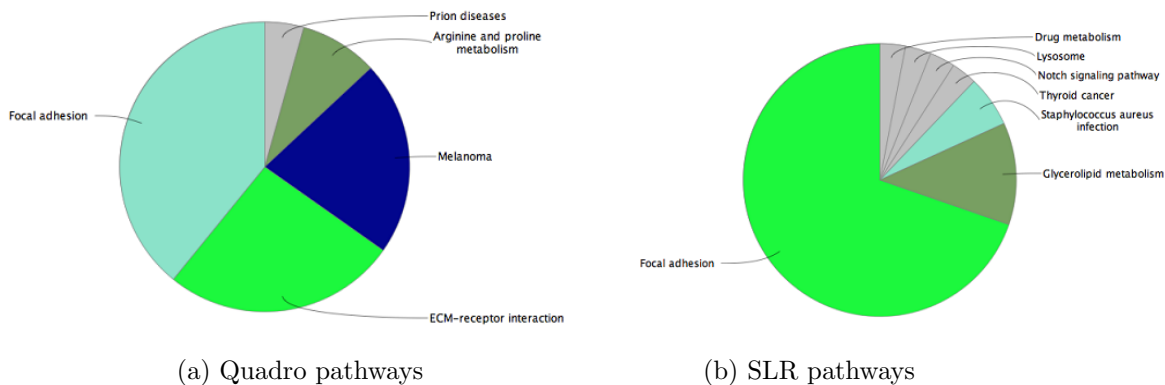


Figure 5: Overall KEGG enrichment chart, using (a) QUADRO; (b) SLR.

The results are summarized in Table 3. We see that QUADRO outperforms SLR and L-SLR.

Next, we look at gene selection. We apply two-fold cross-validation to both QUADRO and SLR. In the results, Quadro selects 139 genes and SLR selects 128 genes. According to KEGG database, genes selected by QUADRO belong to 5 of the pathways that contain more than two genes; correspondingly, genes selected by SLR belong to 7 pathways. Using the tool ClueGo (Bindea et al., 2009), we display the overall KEGG enrichment chart in Figure 5. We see from Figure 5 that both QUADRO and SLR have *focal adhesion* as its most important functional group. Nevertheless, QUADRO finds *ECM-receptor interaction* as another important functional group. *ECM-receptor interaction* is a class consisting of a mixture of structural and functional macromolecules, and it plays an important role in maintaining cell and tissue structures and functions. Massive studies (Luparello, 2013; Wei and Li, 2007) have found evidence that this class is closely related to breast cancer.

Besides the pathway analysis, we also perform the Gene Ontology (GO) enrichment analysis on genes selected by QUADRO. The analysis is done by “DAVID Bioinformatics Resources” and the results are shown in Table 4. We present the biological processes with p-values smaller than 10^{-3} . According to the table, we see that many biological processes are significantly enriched, and they are related to previously selected pathways. For instance, the biological process *cell adhesion* is known to be highly related to *cell communication pathways*, including *focal adhesion* and *ECM-receptor interaction*.

Table 4: Enrichment analysis results according to Gene Ontology for genes selected by Quadro. The four columns represent GO ID, GO attribute, number of selected genes having the attribute and their corresponding p-values. We rank them according to p-values in increasing order.

GO ID	GO attribute	No. of Genes	p-value
0048856	anatomical structure development	58	3.7E-12
0032502	developmental process	62	2.9E-10
0048731	system development	52	3.1E-10
0007275	multicellular organismal development	55	1.8E-8
0001501	skeletal system development	15	1.3E-6
0032501	multicellular organismal process	66	1.4E-6
0048513	organ development	37	1.4E-6
0009653	anatomical structure morphogenesis	28	8.7E-6
0048869	cellular developmental process	34	1.9E-5
0030154	cell differentiation	33	2.1E-5
0007155	cell adhesion	18	2.4E-4
0022610	biological adhesion	18	2.2E-4
0042127	regulation of cell proliferation	19	2.9E-4
0009888	tissue development	17	3.7E-4
0007398	ectoderm development	9	4.8E-4
0048518	positive regulation of biological process	34	5.6E-4
0009605	response to external stimulus	20	6.3E-4
0043062	extracellular structure organization	8	7.4E-4
0007399	nervous system development	22	8.4E-4

8 Conclusions and Extensions

QUADRO is a robust sparse high-dimensional classifier, which allows us to use differences in covariance matrices to enhance discriminability. It is based on Rayleigh quotient optimization. The variance of quadratic statistics involves all fourth cross moments and this can create both computational and statistical problems. These problems are avoided by limiting our applications to the elliptical class of distributions. Robust M-estimator and rank-based estimation of correlations allow us to obtain the uniform convergence for nonpolynomially many of parameters even when the underlying distributions have the finite fourth moments. This allows us to establish oracle inequalities under relatively weaker conditions.

The Rayleigh optimization framework developed in this paper can also be extended to the multi-class case. Suppose the data are drawn independently from a joint distribution of (\mathbf{X}, Y) , where $\mathbf{X} \in \mathbb{R}^d$ and Y takes values in $\{0, 1, \dots, K-1\}$. The definition (1) for the Rayleigh quotient of a projection $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is still well-defined. Let $\pi_k = \mathbb{P}(Y = k)$, for $k = 0, 1, \dots, K-1$. In this K -class situation,

$$\text{Rq}(f) = \frac{\sum_{0 \leq k < l \leq K-1} \pi_k \pi_l \{ \mathbb{E}[f(\mathbf{X})|Y = k] - \mathbb{E}[f(\mathbf{X})|Y = l] \}^2}{\sum_{0 \leq k \leq K-1} \pi_k \text{var}[f(\mathbf{X})|Y = k]}. \quad (19)$$

Let $M_k(f) = \mathbb{E}[f(\mathbf{X})|Y = k]$ and $L_k(f) = \text{var}[f(\mathbf{X})|Y = k]$. Similar to the two-class case,

maximizing $\text{Rq}(f)$ is equivalent to solving the following optimization problem

$$\min_f \sum_{k=0}^{K-1} \pi_k L_k(f), \quad \text{s.t.} \quad \sum_{0 \leq k < l \leq K-1} \pi_k \pi_l |M_k(f) - M_l(f)|^2 = 1.$$

However, this is not a convex problem. We consider an approximate Rayleigh-quotient-maximization problem as follows:

$$\min_f \sum_{k=0}^{K-1} \pi_k L_k(f) \quad \text{s.t.} \quad \sqrt{\pi_k \pi_l} |M_k(f) - M_l(f)| \geq 1, \quad 0 \leq k < l \leq K-1.$$

To solve this problem, we first pick an order of $M_1(f), \dots, M_K(f)$ to remove the absolute values in the constraints. Then it becomes a convex problem. Therefore, the whole optimization can be carried out by simultaneously solving $K!$ convex problems. When K is small, the computational cost is reasonable. In practice, we can apply more efficient algorithms to speed up the computation.

9 Proofs

9.1 Proof of Proposition 2.1

We first present a lemma which is proved in Section A.

Lemma 9.1. *If \mathbf{U} follows a uniform distribution on \mathcal{S}^{d-1} , for any $d \times d$ diagonal matrix \mathbf{S} and any vector $\boldsymbol{\beta} \in \mathbb{R}^d$, we have*

- $\mathbb{E}(\mathbf{U}^\top \mathbf{S} \mathbf{U}) = \frac{\text{tr}(\mathbf{S})}{d}$, $\mathbb{E}[(\mathbf{U}^\top \mathbf{S} \mathbf{U})^2] = \frac{2 \text{tr}(\mathbf{S}^2) + [\text{tr}(\mathbf{S})]^2}{d^2 + 2d}$;
- $\mathbb{E}(\mathbf{U}^\top \boldsymbol{\beta}) = 0$, $\mathbb{E}[(\mathbf{U}^\top \boldsymbol{\beta})^2] = \frac{\|\boldsymbol{\beta}\|^2}{d}$;
- $\mathbb{E}(\mathbf{U}^\top \mathbf{S} \mathbf{U} \mathbf{U}^\top \boldsymbol{\beta}) = 0$.

Now, we show the claim of Proposition 2.1. Let $\mathbf{Y} = \boldsymbol{\Sigma}^{-1/2}(\mathbf{Z} - \boldsymbol{\mu})$, then $\mathbf{Y} = \xi \mathbf{U}$ where \mathbf{U} follows a uniform distribution on \mathcal{S}^{d-1} and is independent of ξ . The quadratic form $Q(\mathbf{Z})$ can be rewritten as

$$\begin{aligned} Q(\mathbf{Z}) &= \mathbf{Z}^\top \boldsymbol{\Omega} \mathbf{Z} - 2\boldsymbol{\delta}^\top \mathbf{Z} \\ &= \mathbf{Y}^\top \boldsymbol{\Sigma}^{1/2} \boldsymbol{\Omega} \boldsymbol{\Sigma}^{1/2} \mathbf{Y} + 2\mathbf{Y}^\top \boldsymbol{\Sigma}^{1/2} (\boldsymbol{\Omega} \boldsymbol{\mu} - \boldsymbol{\delta}) + \boldsymbol{\mu}^\top \boldsymbol{\Omega} \boldsymbol{\mu} - 2\boldsymbol{\mu}^\top \boldsymbol{\delta} \\ &= \bar{Q}(\mathbf{Y}) + c, \end{aligned}$$

where $c = \boldsymbol{\mu}^\top \boldsymbol{\Omega} \boldsymbol{\mu} - 2\boldsymbol{\mu}^\top \boldsymbol{\delta}$. Therefore, $\mathbb{E}[Q(\mathbf{Z})] = \mathbb{E}[\bar{Q}(\mathbf{Y})] + c$ and $\text{var}[Q(\mathbf{Z})] = \text{var}[\bar{Q}(\mathbf{Y})]$.

Furthermore, we let

$$\boldsymbol{\Sigma}^{1/2} \boldsymbol{\Omega} \boldsymbol{\Sigma}^{1/2} = \mathbf{K} \mathbf{S} \mathbf{K}^\top$$

be the eigenvalue decomposition of $\boldsymbol{\Sigma}^{1/2} \boldsymbol{\Omega} \boldsymbol{\Sigma}^{1/2}$, where \mathbf{K} is an orthogonal matrix and \mathbf{S} is a diagonal matrix. We also define

$$\boldsymbol{\beta} = \mathbf{K}^\top \boldsymbol{\Sigma}^{1/2} (\boldsymbol{\Omega} \boldsymbol{\mu} - \boldsymbol{\delta}).$$

Notice that $\mathbf{Y}^\top \boldsymbol{\Sigma}^{1/2} \boldsymbol{\Omega} \boldsymbol{\Sigma}^{1/2} \mathbf{Y} = \xi^2 \mathbf{U}^\top \mathbf{K} \mathbf{S} \mathbf{K}^\top \mathbf{U} = \xi^2 \mathbf{U}_1^\top \mathbf{S} \mathbf{U}_1$, where $\mathbf{U}_1 = \mathbf{K}^\top \mathbf{U}$. Since \mathbf{K} is an orthogonal matrix, \mathbf{U}_1 follows the same distribution as \mathbf{U} and is also independent of ξ . Moreover, we can write $\mathbf{Y}^\top \boldsymbol{\Sigma}^{1/2} (\boldsymbol{\Omega} \boldsymbol{\mu} - \boldsymbol{\delta}) = \xi \mathbf{U}_1^\top \boldsymbol{\beta}$. To save notation, we still use \mathbf{U} to represent \mathbf{U}_1 . It follows that

$$\bar{Q}(\mathbf{Y}) = \xi^2 \mathbf{U}^\top \mathbf{S} \mathbf{U} + 2\xi \mathbf{U}^\top \boldsymbol{\beta}.$$

Let's calculate $\mathbb{E}[\bar{Q}(\mathbf{Y})]$ first.

$$\begin{aligned} \mathbb{E}[\bar{Q}(\mathbf{Y})] &= \mathbb{E}(\xi^2) \mathbb{E}(\mathbf{U}^\top \mathbf{S} \mathbf{U}) + 2\mathbb{E}(\xi) \mathbb{E}[\mathbf{U}^\top \boldsymbol{\beta}] \\ &= \frac{\mathbb{E}(\xi^2)}{d} \text{tr}(\mathbf{S}) = \text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma}). \end{aligned}$$

The first equality is due to the fact that ξ and \mathbf{U} are independent; the second equality is from Lemma 9.1; and the last inequality is because $\mathbb{E}(\xi^2) = d$ and $\text{tr}(\mathbf{S}) = \text{tr}(\boldsymbol{\Sigma}^{1/2} \boldsymbol{\Omega} \boldsymbol{\Sigma}^{1/2}) = \text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma})$. It follows that

$$\mathbb{E}[Q(\mathbf{Z})] = \mathbb{E}[\bar{Q}(\mathbf{Y})] + c = \text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\Omega} \boldsymbol{\mu} - 2\boldsymbol{\mu}^\top \boldsymbol{\delta}.$$

Next, we calculate $\text{var}[\bar{Q}(\mathbf{Y})]$. It follows that

$$\begin{aligned} \text{var}[\bar{Q}(\mathbf{Y})] &= \text{var}(\xi^2 \mathbf{U}^\top \mathbf{S} \mathbf{U} + 2\xi \mathbf{U}^\top \boldsymbol{\beta}) \\ &= \text{var}(\xi^2 \mathbf{U}^\top \mathbf{S} \mathbf{U}) + 4 \text{var}(\xi \mathbf{U}^\top \boldsymbol{\beta}) + 4 \text{cov}(\xi^2 \mathbf{U}^\top \mathbf{S} \mathbf{U}, \xi \mathbf{U}^\top \boldsymbol{\beta}). \end{aligned}$$

Let's look at them term by term. First,

$$\begin{aligned} \text{var}(\xi^2 \mathbf{U}^\top \mathbf{S} \mathbf{U}) &= \mathbb{E}[\xi^4 (\mathbf{U}^\top \mathbf{S} \mathbf{U})^2] - \mathbb{E}^2(\xi^2 \mathbf{U}^\top \mathbf{S} \mathbf{U}) \\ &= \mathbb{E}(\xi^4) \mathbb{E}[(\mathbf{U}^\top \mathbf{S} \mathbf{U})^2] - \mathbb{E}^2(\xi^2) \mathbb{E}^2(\mathbf{U}^\top \mathbf{S} \mathbf{U}) \\ &= \mathbb{E}(\xi^4) \frac{2 \text{tr}(\mathbf{S}^2) + \text{tr}^2(\mathbf{S})}{2d + d^2} - \mathbb{E}^2(\xi^2) \frac{\text{tr}^2(\mathbf{S})}{d^2} \\ &= 2(\gamma + 1) \text{tr}(\mathbf{S}^2) + \gamma \text{tr}^2(\mathbf{S}) \\ &= 2(\gamma + 1) \text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma} \boldsymbol{\Omega} \boldsymbol{\Sigma}) + \gamma [\text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma})]^2. \end{aligned}$$

The third equality comes from Lemma 9.1; the last equality follows from the fact that $\text{tr}(\mathbf{S}^2) = \text{tr}(\boldsymbol{\Sigma}^{1/2} \boldsymbol{\Omega} \boldsymbol{\Sigma} \boldsymbol{\Omega} \boldsymbol{\Sigma}^{1/2}) = \text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma} \boldsymbol{\Omega} \boldsymbol{\Sigma})$. Second,

$$\begin{aligned} \text{var}(\xi \mathbf{U}^\top \boldsymbol{\beta}) &= \mathbb{E}(\xi^2 (\mathbf{U}^\top \boldsymbol{\beta})^2) - \mathbb{E}^2(\xi \mathbf{U}^\top \boldsymbol{\beta}) \\ &= \mathbb{E}(\xi^2) \mathbb{E}[(\mathbf{U}^\top \boldsymbol{\beta})^2] - \mathbb{E}^2(\xi) \mathbb{E}^2(\mathbf{U}^\top \boldsymbol{\beta}) \\ &= \mathbb{E}(\xi^2) \frac{\|\boldsymbol{\beta}\|^2}{d} = (\boldsymbol{\Omega} \boldsymbol{\mu} - \boldsymbol{\delta})^\top \boldsymbol{\Sigma} (\boldsymbol{\Omega} \boldsymbol{\mu} - \boldsymbol{\delta}). \end{aligned}$$

In the last equality, we have used $\mathbb{E}(\xi^2) = d$, $\boldsymbol{\beta} = \mathbf{K} \boldsymbol{\Sigma}^{1/2} (\boldsymbol{\Omega} \boldsymbol{\mu} - \boldsymbol{\delta})$ and $\mathbf{K}^\top \mathbf{K} = \mathbf{I}_d$. Last,

$$\begin{aligned} \text{cov}(\xi^2 \mathbf{U}^\top \mathbf{S} \mathbf{U}, \xi \mathbf{U}^\top \boldsymbol{\beta}) &= \mathbb{E}(\xi^3) \mathbb{E}(\mathbf{U}^\top \mathbf{S} \mathbf{U} \mathbf{U}^\top \boldsymbol{\beta}) - \mathbb{E}(\xi^2) \mathbb{E}(\mathbf{U}^\top \mathbf{S} \mathbf{U}) \mathbb{E}(\xi \mathbf{U}^\top \boldsymbol{\beta}) \\ &= \mathbb{E}(\xi^3) \mathbb{E}(\mathbf{U}^\top \mathbf{S} \mathbf{U} \mathbf{U}^\top \boldsymbol{\beta}) = 0 \end{aligned}$$

Combining the above gives

$$\text{var}[Q(\mathbf{Z})] = \text{var}[\bar{Q}(\mathbf{Y})] = 2(\gamma + 1) \text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma} \boldsymbol{\Omega} \boldsymbol{\Sigma}) + \gamma [\text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma})]^2 + 4(\boldsymbol{\Omega} \boldsymbol{\mu} - \boldsymbol{\delta})^\top \boldsymbol{\Sigma} (\boldsymbol{\Omega} \boldsymbol{\mu} - \boldsymbol{\delta}).$$

□

9.2 Proof of Proposition 5.1

For any $d \times 1$ vector \mathbf{v} and $d \times d$ matrix \mathbf{A} , we denote by $\text{Supp}(\mathbf{v})$ the support of \mathbf{v} , which is contained in $\{1, \dots, d\}$, and by $\text{Supp}(\mathbf{A})$ the support of \mathbf{A} , which is contained in $\{1, \dots, d\} \times \{1, \dots, d\}$. Let $\theta = \sqrt{1 + v_1(1 + \gamma)/2} > 1$ and

$$c = (1 + \kappa) \min \{ (1 + \gamma)v_2^2, 4(1 - 1/\theta^2)v_2 \}.$$

The claim then becomes $\Theta(S, 0) \geq c$, or in other words,

$$\mathbf{x}^\top \mathbf{Q} \mathbf{x} \geq c |\mathbf{x}|^2, \quad \text{when } \text{Supp}(\mathbf{x}) \subset S.$$

First, using (13) and Lemma 9.2, we find that for each \mathbf{x} , there exists unique $(\mathbf{\Omega}, \boldsymbol{\delta})$ such that $\mathbf{x} = \mathbf{x}(\mathbf{\Omega}, \boldsymbol{\delta})$ and $\mathbf{x}^\top \mathbf{Q} \mathbf{x} = L(\mathbf{\Omega}, \boldsymbol{\delta})$. Second, by definition of U' and V , $\text{Supp}(\mathbf{x}) \subset S$ implies that $\text{Supp}(\mathbf{\Omega}) \subset U' \times U'$ and $\text{Supp}(\boldsymbol{\delta}) \subset V$. Therefore, it suffices to show

$$L(\mathbf{\Omega}, \boldsymbol{\delta}) \geq c(|\mathbf{\Omega}|^2 + |\boldsymbol{\delta}|^2), \quad \text{when } \text{Supp}(\mathbf{\Omega}) \subset U' \times U' \text{ and } \text{Supp}(\boldsymbol{\delta}) \subset V. \quad (20)$$

Now, we show (20). From (5) and that $\gamma \geq 0$,

$$L_k(\mathbf{\Omega}, \boldsymbol{\delta}) \geq 2(1 + \gamma) \text{tr}(\mathbf{\Omega} \boldsymbol{\Sigma}_k \mathbf{\Omega} \boldsymbol{\Sigma}_k) + 4(\mathbf{\Omega} \boldsymbol{\mu}_k - \boldsymbol{\delta})^\top \boldsymbol{\Sigma}_k (\mathbf{\Omega} \boldsymbol{\mu}_k - \boldsymbol{\delta}), \quad k = 1, 2.$$

Let $\tilde{\mathbf{\Omega}}$ be the submatrix of $\mathbf{\Omega}$ by restricting rows and columns to the set $U' \cup V$, and $\tilde{\boldsymbol{\delta}}$ be the subvector of $\boldsymbol{\delta}$ by restricting the elements to the set $U' \cup V$. It is easy to see that when $\text{Supp}(\mathbf{\Omega}) \subset U' \times U'$ and $\text{Supp}(\boldsymbol{\delta}) \subset V$,

$$\begin{aligned} \text{tr}(\mathbf{\Omega} \boldsymbol{\Sigma}_k \mathbf{\Omega} \boldsymbol{\Sigma}_k) &= \text{tr}(\tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\Sigma}}_k \tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\Sigma}}_k), \\ (\mathbf{\Omega} \boldsymbol{\mu}_k - \boldsymbol{\delta})^\top \boldsymbol{\Sigma}_k (\mathbf{\Omega} \boldsymbol{\mu}_k - \boldsymbol{\delta}) &= (\tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\mu}}_k - \tilde{\boldsymbol{\delta}})^\top \tilde{\boldsymbol{\Sigma}}_k (\tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\mu}}_k - \tilde{\boldsymbol{\delta}}), \end{aligned}$$

where we recall that $\tilde{\boldsymbol{\Sigma}}_k$ is the submatrix of $\boldsymbol{\Sigma}_k$ by restricting rows and columns to the set $U' \cup V$, $\tilde{\boldsymbol{\mu}}_k$ to the set $U' \cup V$. It follows that

$$\begin{aligned} &L_k(\mathbf{\Omega}, \boldsymbol{\delta}) \\ &\geq 2(1 + \gamma) \text{tr}(\tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\Sigma}}_k \tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\Sigma}}_k) + 4(\tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\mu}}_k - \tilde{\boldsymbol{\delta}})^\top \tilde{\boldsymbol{\Sigma}}_k (\tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\mu}}_k - \tilde{\boldsymbol{\delta}}) \\ &= 2(1 + \gamma) \text{tr}(\tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\Sigma}}_k \tilde{\mathbf{\Omega}} (\tilde{\boldsymbol{\Sigma}}_k - v_1 \tilde{\boldsymbol{\mu}}_k \tilde{\boldsymbol{\mu}}_k^\top)) + 4(1 - 1/\theta^2) \tilde{\boldsymbol{\delta}}^\top \tilde{\boldsymbol{\Sigma}}_k \tilde{\boldsymbol{\delta}} \\ &\quad + 4(\theta \tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\mu}}_k - \theta^{-1} \tilde{\boldsymbol{\delta}})^\top \tilde{\boldsymbol{\Sigma}}_k (\theta \tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\mu}}_k - \theta^{-1} \tilde{\boldsymbol{\delta}}) \\ &\geq 2(1 + \gamma) \text{tr}(\tilde{\mathbf{\Omega}} \tilde{\boldsymbol{\Sigma}}_k \tilde{\mathbf{\Omega}} (\tilde{\boldsymbol{\Sigma}}_k - v_1 \tilde{\boldsymbol{\mu}}_k \tilde{\boldsymbol{\mu}}_k^\top)) + 4(1 - 1/\theta^2) \lambda_{\min}(\tilde{\boldsymbol{\Sigma}}_k) |\tilde{\boldsymbol{\delta}}|^2. \end{aligned} \quad (21)$$

Denote by I_1 the first term in (21). We aim to derive a lower bound for I_1 . It is well known that $\text{tr}(\mathbf{A}^\top \mathbf{B} \mathbf{C} \mathbf{D}^\top) = \text{vec}(\mathbf{A})^\top (\mathbf{D} \otimes \mathbf{B}) \text{vec}(\mathbf{C})$, where $\text{vec}(\mathbf{A})$ be the vectorization of \mathbf{A} by stacking all the columns, $\mathbf{D} \otimes \mathbf{B}$ is the Kronecker product of \mathbf{D} and \mathbf{B} . Using this formula and that $\boldsymbol{\Sigma}_k$ is symmetric, we find that

$$\begin{aligned} I_1 &= 2(1 + \gamma) \text{vec}(\tilde{\mathbf{\Omega}})^\top [(\tilde{\boldsymbol{\Sigma}}_k - v_1 \tilde{\boldsymbol{\mu}}_k \tilde{\boldsymbol{\mu}}_k^\top) \otimes \tilde{\boldsymbol{\Sigma}}_k] \text{vec}(\tilde{\mathbf{\Omega}}) \\ &\geq 2(1 + \gamma) |\tilde{\mathbf{\Omega}}|^2 \lambda_{\min}((\tilde{\boldsymbol{\Sigma}}_k - v_1 \tilde{\boldsymbol{\mu}}_k \tilde{\boldsymbol{\mu}}_k^\top) \otimes \tilde{\boldsymbol{\Sigma}}_k) \\ &\geq (1 + \gamma) \lambda_{\min}^2(\tilde{\boldsymbol{\Sigma}}_k) |\tilde{\mathbf{\Omega}}|^2. \end{aligned}$$

The last inequality is from the property that $\lambda_{\min}(\mathbf{A} \otimes \mathbf{B}) = \lambda_{\min}(\mathbf{A})\lambda_{\min}(\mathbf{B})$ when \mathbf{A} and \mathbf{B} are positive semi-definite, and also the assumption that $\lambda_{\min}(\tilde{\Sigma}_k - v_1 \tilde{\mu}_k \tilde{\mu}_k^\top) \geq \frac{1}{2} \lambda_{\min}(\tilde{\Sigma}_k)$. Plugging I_1 into (21), we have

$$\begin{aligned}
& L_1(\Omega, \delta) + \kappa L_2(\Omega, \delta) \\
& \geq (1 + \gamma)[\lambda_{\min}^2(\tilde{\Sigma}_1) + \kappa \lambda_{\min}^2(\tilde{\Sigma}_2)]|\tilde{\Omega}|^2 + 4(1 - 1/\theta^2)[\lambda_{\min}(\tilde{\Sigma}_1) + \kappa \lambda_{\min}(\tilde{\Sigma}_2)]|\tilde{\delta}|^2 \\
& \geq (1 + \gamma)(1 + \kappa)v_2^2|\tilde{\Omega}|^2 + 4(1 - 1/\theta^2)(1 + \kappa)v_2|\tilde{\delta}|^2 \\
& \geq c(|\tilde{\Omega}|^2 + |\tilde{\delta}|^2) = c(|\Omega|^2 + |\delta|^2).
\end{aligned}$$

This proves (20). \square

9.3 Proof of Theorem 5.1

We prove the claim by first rewriting the optimization problem (8) into a vector form. For any (Ω, δ) , write $\mathbf{x} = [\text{vec}(\Omega)^\top, \delta^\top]^\top$. Let \mathbf{Q} be as defined in Section 5, and

$$\mathbf{q} = \left[\text{vec}(\Sigma_2 + \mu_2 \mu_2^\top - \Sigma_1 - \mu_1 \mu_1^\top)^\top, 2(\mu_1 - \mu_2)^\top \right]^\top.$$

We introduce the following lemma which is proved in Section A.

Lemma 9.2. $M(\Omega, \delta) = \mathbf{q}^\top \mathbf{x}$ and $L(\Omega, \delta) = \mathbf{x}^\top \mathbf{Q} \mathbf{x}$.

Let $\mathbf{x}_{\lambda_0}^* = [\text{vec}(\Omega_{\lambda_0}^*)^\top, (\delta_{\lambda_0}^*)^\top]^\top$ and $\hat{\mathbf{x}} = [\text{vec}(\hat{\Omega})^\top, \hat{\delta}^\top]^\top$. Using Lemma 9.2,

$$\begin{aligned}
\mathbf{x}_{\lambda_0}^* &= \min_{\mathbf{x} \in \mathbb{R}^d: \mathbf{q}^\top \mathbf{x} = 1} \{ \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \lambda_0 |\mathbf{x}|_1 \}, \\
\hat{\mathbf{x}} &= \underset{\mathbf{x} \in \mathbb{R}^d: \hat{\mathbf{q}}^\top \mathbf{x} = 1}{\text{argmin}} \{ \mathbf{x}^\top \hat{\mathbf{Q}} \mathbf{x} + \lambda |\mathbf{x}|_1 \},
\end{aligned}$$

where $\hat{\mathbf{Q}}$ and $\hat{\mathbf{q}}$ are counter parts of \mathbf{Q} and \mathbf{q} respectively, by replacing μ_1, μ_2, Σ_1 and Σ_2 with their estimates. Moreover, the Rayleigh quotient

$$R(\Omega, \delta) = R(\mathbf{x}) \equiv \frac{(\mathbf{q}^\top \mathbf{x})^2}{\mathbf{x}^\top \mathbf{Q} \mathbf{x}}.$$

In addition, we have the following lemma, which is proved in Section A.

Lemma 9.3. $\max\{|\hat{\mathbf{Q}} - \mathbf{Q}|_\infty, |\hat{\mathbf{q}} - \mathbf{q}|_\infty\} \leq C_0 \max\{|\hat{\Sigma}_k - \Sigma_k|_\infty, |\hat{\mu}_k - \mu_k|_\infty, k = 1, 2\}$ for some constant $C_0 > 0$.

Combining the above results, the claim follows immediately from the following theorem:

Theorem 9.1. For any $\lambda_0 \geq 0$, let S be the support of $\mathbf{x}_{\lambda_0}^*$. Suppose $\Theta(S, 0) \geq c_0$, $\Theta(S, 3) \geq a_0$ and $R(\mathbf{x}_{\lambda_0}^*) \geq u_0$, for positive constants a_0, c_0 and u_0 . Let $\Delta_n = \max\{|\hat{\mathbf{Q}} - \mathbf{Q}|_\infty, |\hat{\mathbf{q}} - \mathbf{q}|_\infty\}$, $s_0 = |S|$ and $k_0 = \max\{s_0, R(\mathbf{x}_{\lambda_0}^*)\}$. Suppose $4s_0 \Delta_n^2 < c_0 u_0$ and $\max\{s_0 \Delta_n, s_0^{1/2} k_0^{1/2} \lambda_0\} < 1$. Then, there exist positive constants $C = C(a_0, c_0, u_0)$ and $A = A(a_0, c_0, u_0)$, such that for any $\eta > 1$, by taking $\lambda = C\eta \max\{s_0^{1/2} \Delta_n, k_0^{1/2} \lambda_0\} [R(\mathbf{x}_{\lambda_0}^*)]^{-1/2}$,

$$\frac{R(\hat{\mathbf{x}})}{R(\mathbf{x}_{\lambda_0}^*)} \geq 1 - A\eta^2 \max\{s_0 \Delta_n, s_0^{1/2} k_0^{1/2} \lambda_0\}.$$

The main part of the proof is to show Theorem 9.1. Write for short $\mathbf{x}^* = \mathbf{x}_{\lambda_0}^*$, $R^* = R(\mathbf{x}^*)$, $V^* = (R^*)^{-1} = (\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^*$, $\bar{V}^* = (V^*)^{1/2}$. Let $\alpha_n = \Delta_n |\mathbf{x}^*|_0^{1/2}$, $\beta_n = \Delta_n |\mathbf{x}^*|_0$ and $T_n(\mathbf{x}^*) = \max\{s_0 \Delta_n, s_0^{1/2} k_0^{1/2} \lambda_0\}$. We define the quantity

$$\Gamma(\mathbf{x}) = \frac{|\mathbf{Q}\mathbf{x} - (\mathbf{x}^\top \mathbf{Q}\mathbf{x})\mathbf{q}|_\infty}{(\mathbf{x}^\top \mathbf{Q}\mathbf{x})^{1/2}}, \quad \text{for any } \mathbf{x}.$$

Step 1: We introduce \mathbf{x}_1^* , a multiple of \mathbf{x}^* , and use it to bound $|\hat{\mathbf{x}}|_1$.

Let \mathbf{Q}_{SS} be the submatrix of \mathbf{Q} formed by rows and columns corresponding to S . Since $\lambda_{\min}(\mathbf{Q}_{SS}) = \Theta(S, 0) \geq c_0$, we have $(\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^* \geq c_0 |\mathbf{x}^*|^2$. Using this fact and by Cauchy-Schwartz inequality,

$$|\mathbf{x}^*|_1 \leq \sqrt{|\mathbf{x}^*|_0} |\mathbf{x}^*| \leq c_0^{-1/2} \sqrt{|\mathbf{x}^*|_0} \bar{V}^* \quad (22)$$

It follows that

$$|\hat{\mathbf{q}}^\top \mathbf{x}^* - \mathbf{q}^\top \mathbf{x}^*| \leq |\hat{\mathbf{q}} - \mathbf{q}|_\infty |\mathbf{x}^*|_1 \leq c_0^{-1/2} \Delta_n \sqrt{|\mathbf{x}^*|_0} \bar{V}^* = c_0^{-1/2} \alpha_n \bar{V}^*. \quad (23)$$

Let $t_n = \hat{\mathbf{q}}^\top \mathbf{x}^*$. Then (23) says that $|t_n - 1| \leq c_0^{-1/2} \alpha_n \bar{V}^*$. Noting that $\bar{V}^* = (R^*)^{1/2} \leq u_0^{-1/2}$, we have $|t_n - 1| \leq (c_0 u_0)^{-1/2} s_0^{1/2} \Delta_n < 1/2$ by assumption. In particular, $t_n > 0$. Let

$$\mathbf{x}_1^* = t_n^{-1} \mathbf{x}^*.$$

Then $\hat{\mathbf{q}}^\top \mathbf{x}_1^* = 1$. From the definition of $\hat{\mathbf{x}}$,

$$\hat{\mathbf{x}}^\top \hat{\mathbf{Q}} \hat{\mathbf{x}} + \lambda |\hat{\mathbf{x}}|_1 \leq (\mathbf{x}_1^*)^\top \hat{\mathbf{Q}} \mathbf{x}_1^* + \lambda |\mathbf{x}_1^*|_1. \quad (24)$$

By direct calculation,

$$\begin{aligned} \hat{\mathbf{x}}^\top \hat{\mathbf{Q}} \hat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \hat{\mathbf{Q}} \mathbf{x}_1^* &= (\hat{\mathbf{x}} - \mathbf{x}_1^*)^\top \hat{\mathbf{Q}} (\hat{\mathbf{x}} - \mathbf{x}_1^*) + 2(\hat{\mathbf{x}} - \mathbf{x}_1^*)^\top \hat{\mathbf{Q}} \mathbf{x}_1^* \\ &= (\hat{\mathbf{x}} - \mathbf{x}_1^*)^\top \hat{\mathbf{Q}} (\hat{\mathbf{x}} - \mathbf{x}_1^*) + 2(\hat{\mathbf{x}} - \mathbf{x}_1^*)^\top (\hat{\mathbf{Q}} \mathbf{x}_1^* - V^* \hat{\mathbf{q}}) \\ &\geq 2(\hat{\mathbf{x}} - \mathbf{x}_1^*)^\top (\hat{\mathbf{Q}} \mathbf{x}_1^* - V^* \hat{\mathbf{q}}) \end{aligned} \quad (25)$$

where the second equality is because $\hat{\mathbf{q}}^\top \hat{\mathbf{x}} = \hat{\mathbf{q}}^\top \mathbf{x}_1^* = 1$. We aim to bound $|\hat{\mathbf{Q}} \mathbf{x}_1^* - V^* \hat{\mathbf{q}}|_\infty$. The following lemma is proved in Section A.

Lemma 9.4. *When $\Theta(S, 0) \geq c_0$, there exists a positive constant $C_1 = C_1(c_0)$ such that $\Gamma(\mathbf{x}_{\lambda_0}^*) \leq C_1 \lambda_0 [\max\{s_0, R(\mathbf{x}_{\lambda_0}^*)\}]^{1/2}$ for any $\lambda_0 \geq 0$.*

Since $\mathbf{x}_1^* = t_n^{-1} \mathbf{x}^*$ and $t_n^{-1} < 2$,

$$\begin{aligned} |\hat{\mathbf{Q}} \mathbf{x}_1^* - V^* \hat{\mathbf{q}}|_\infty &\leq t_n^{-1} |\hat{\mathbf{Q}} \mathbf{x}^* - V^* \hat{\mathbf{q}}|_\infty + V^* |t_n^{-1} - 1| |\hat{\mathbf{q}}|_\infty \\ &\leq 2(|\mathbf{Q} \mathbf{x}^* - V^* \hat{\mathbf{q}}|_\infty + |\hat{\mathbf{Q}} - \mathbf{Q}|_\infty |\mathbf{x}^*|_1 + V^* |\hat{\mathbf{q}} - \mathbf{q}|_\infty + V^* |t_n - 1| |\hat{\mathbf{q}}|_\infty) \\ &\leq 2[\Gamma(\mathbf{x}^*) \bar{V}^* + c_0^{-1/2} \alpha_n \bar{V}^* + u_0^{-1/2} \Delta_n \bar{V}^* + |\hat{\mathbf{q}}|_\infty c_0^{-1/2} u_0^{-1} \alpha_n \bar{V}^*] \\ &\leq C_2 (\lambda_0 k_0^{1/2} + s_0^{1/2} \Delta_n) \bar{V}^*. \end{aligned}$$

Here the third inequality follows from (22)-(23) and $V^* = \bar{V}^* (R^*)^{-1/2} \leq u_0^{-1/2} \bar{V}^*$. The last inequality is obtained as follows: from Lemma 9.3, we know that $|\hat{\mathbf{q}}|_\infty \leq |\mathbf{q}|_\infty + |\hat{\mathbf{q}} - \mathbf{q}|_\infty \leq$

$2C_0$ (see also the assumptions in the beginning of Section 5.2); we also use Lemma 9.4 and $\alpha_n \bar{V}^* \leq u_0^{-1/2} s_0^{1/2} \Delta_n$. By letting $C = 8C_2$, the choice of $\lambda = C\eta \max\{s_0^{1/2} \Delta_n, k_0^{1/2} \lambda_0\} \bar{V}^*$ for $\eta > 1$ ensures that

$$|\widehat{\mathbf{Q}}\mathbf{x}_1^* - \widehat{\mathbf{q}}|_\infty \leq \lambda/4.$$

Plugging this result into (25) gives

$$\widehat{\mathbf{x}}^\top \widehat{\mathbf{Q}}\widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \widehat{\mathbf{Q}}\mathbf{x}_1^* \geq -\frac{\lambda}{2} |\widehat{\mathbf{x}} - \mathbf{x}_1^*|_1. \quad (26)$$

Combining (24) and (26) gives

$$\lambda |\widehat{\mathbf{x}}|_1 - \frac{\lambda}{2} |\widehat{\mathbf{x}} - \mathbf{x}_1^*|_1 \leq \lambda |\mathbf{x}_1^*|_1. \quad (27)$$

First, since $|\widehat{\mathbf{x}}|_1 = |\widehat{\mathbf{x}}_S|_1 + |\widehat{\mathbf{x}}_{S^c}|_1 \geq |\mathbf{x}_{1S}^*|_1 - |\widehat{\mathbf{x}}_S - \mathbf{x}_{1S}^*|_1 + |\widehat{\mathbf{x}}_{S^c}|_1$ and $|\widehat{\mathbf{x}} - \mathbf{x}_1^*|_1 = |\widehat{\mathbf{x}}_S - \mathbf{x}_{1S}^*|_1 + |\widehat{\mathbf{x}}_{S^c}|_1$, we immediately see from (27) that

$$|(\widehat{\mathbf{x}} - \mathbf{x}_1^*)_{S^c}|_1 \leq 3|(\widehat{\mathbf{x}} - \mathbf{x}_1^*)_S|_1. \quad (28)$$

Second, note that $|\widehat{\mathbf{x}} - \mathbf{x}_1^*|_1 \leq |\widehat{\mathbf{x}}|_1 + |\mathbf{x}_1^*|_1$. Plugging it into (27) gives

$$|\widehat{\mathbf{x}}|_1 \leq 3|\mathbf{x}_1^*|_1 = 3t_n^{-1} |\mathbf{x}^*|_1 \leq 6c_0^{-1/2} \sqrt{|\mathbf{x}^*|_0} \bar{V}^*. \quad (29)$$

Step 2: We use (28)-(29) to derive an upper bound for $(\widehat{\mathbf{x}})^\top \widehat{\mathbf{Q}}\widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \widehat{\mathbf{Q}}\mathbf{x}_1^*$.

Note that

$$\begin{aligned} & \widehat{\mathbf{x}}^\top \widehat{\mathbf{Q}}\widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \widehat{\mathbf{Q}}\mathbf{x}_1^* \\ & \geq \widehat{\mathbf{x}}^\top \mathbf{Q}\widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \mathbf{Q}\mathbf{x}_1^* - (|\widehat{\mathbf{x}}^\top \widehat{\mathbf{Q}}\widehat{\mathbf{x}} - \widehat{\mathbf{x}}^\top \mathbf{Q}\widehat{\mathbf{x}}| + |(\mathbf{x}_1^*)^\top \widehat{\mathbf{Q}}\mathbf{x}_1^* - (\mathbf{x}_1^*)^\top \mathbf{Q}\mathbf{x}_1^*|) \\ & \geq \widehat{\mathbf{x}}^\top \mathbf{Q}\widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \mathbf{Q}\mathbf{x}_1^* - (|\widehat{\mathbf{Q}} - \mathbf{Q}|_\infty |\widehat{\mathbf{x}}|_1^2 + |\widehat{\mathbf{Q}} - \mathbf{Q}|_\infty |\mathbf{x}_1^*|_1^2) \\ & \geq \widehat{\mathbf{x}}^\top \mathbf{Q}\widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \mathbf{Q}\mathbf{x}_1^* - 10t_n^{-2} |\widehat{\mathbf{Q}} - \mathbf{Q}|_\infty |\mathbf{x}^*|_1^2 \\ & \geq \widehat{\mathbf{x}}^\top \mathbf{Q}\widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \mathbf{Q}\mathbf{x}_1^* - C_3 \beta_n V^*, \end{aligned} \quad (30)$$

where the last two inequalities are direct results of (29). Combining (24) and (30),

$$\widehat{\mathbf{x}}^\top \widehat{\mathbf{Q}}\widehat{\mathbf{x}} + \lambda |\widehat{\mathbf{x}}|_1 \leq (\mathbf{x}_1^*)^\top \mathbf{Q}\mathbf{x}_1^* + \lambda |\mathbf{x}_1^*|_1 + C_3 \beta_n V^*. \quad (31)$$

Similarly to (25), we have

$$\widehat{\mathbf{x}}^\top \widehat{\mathbf{Q}}\widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \mathbf{Q}\mathbf{x}_1^* = (\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*) + 2(\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top (\mathbf{Q}\mathbf{x}_1^* - V^* \widehat{\mathbf{q}}), \quad (32)$$

where

$$\begin{aligned} |\mathbf{Q}\mathbf{x}_1^* - V^* \widehat{\mathbf{q}}|_\infty & \leq t_n^{-1} (|\mathbf{Q}\mathbf{x}^* - V^* \mathbf{q}|_\infty + V^* |\widehat{\mathbf{q}} - \mathbf{q}|_\infty) + V^* |t_n^{-1} - 1| |\widehat{\mathbf{q}}|_\infty \\ & \leq 2[\Gamma(\mathbf{x}^*) \bar{V}^* + u_0^{-1/2} \Delta_n \bar{V}^* + |\widehat{\mathbf{q}}|_\infty c_0^{-1/2} u_0^{-1} \alpha_n \bar{V}^*] \\ & \leq \lambda/4. \end{aligned}$$

It follows that

$$\widehat{\mathbf{x}}^\top \widehat{\mathbf{Q}}\widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \mathbf{Q}\mathbf{x}_1^* \geq (\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*) - \frac{\lambda}{2} |\widehat{\mathbf{x}} - \mathbf{x}_1^*|_1.$$

Plugging it into (31), we obtain

$$(\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*) + \lambda |\widehat{\mathbf{x}}|_1 - \frac{\lambda}{2} |\widehat{\mathbf{x}} - \mathbf{x}_1^*|_1 \leq \lambda |\mathbf{x}_1^*|_1 + C_3 \beta_n V^*. \quad (33)$$

We can rewrite the second and third terms on the left hand side of (33) as

$$\lambda |\widehat{\mathbf{x}}_S|_1 - \frac{\lambda}{2} |\widehat{\mathbf{x}}_S - \mathbf{x}_{1S}^*|_1 + \frac{\lambda}{2} |\widehat{\mathbf{x}}_{S^c}|_1.$$

Plugging it into (33), and by the triangular inequality $|\mathbf{x}_{1S}^*|_1 - |\widehat{\mathbf{x}}_S|_1 \leq |\widehat{\mathbf{x}}_S - \mathbf{x}_{1S}^*|_1$, we find that

$$(\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*) + \frac{\lambda}{2} |\widehat{\mathbf{x}}_{S^c}|_1 \leq \frac{3\lambda}{2} |\widehat{\mathbf{x}}_S - \mathbf{x}_{1S}^*|_1 + C_3 \beta_n V^*.$$

We drop the term $\frac{\lambda}{2} |\widehat{\mathbf{x}}_{S^c}|_1$ on the left hand side, and apply the Cauchy-Schwartz inequality to the term $|\widehat{\mathbf{x}}_S - \mathbf{x}_{1S}^*|_1$. It gives

$$(\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*) \leq \frac{3\lambda}{2} \sqrt{|\mathbf{x}_1^*|_0} |\widehat{\mathbf{x}}_{1S} - \mathbf{x}_{1S}^*| + C_3 \beta_n V^*. \quad (34)$$

Since (28) holds, by definition of $\Theta(S, 3)$,

$$(\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*) \geq a_0 |\widehat{\mathbf{x}}_S - \mathbf{x}_{1S}^*|^2.$$

We write temporarily $Y = (\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*)$ and $b = C_3 \beta_n V^*$. Combining these to (34),

$$Y \leq \frac{3\lambda}{2\sqrt{a_0}} \sqrt{|\mathbf{x}_1^*|_0} Y + b.$$

Note that when $u^2 \leq au + b$, we have $(u - \frac{a}{2})^2 \leq b + \frac{a^2}{4}$, and hence $u^2 \leq 2[\frac{a^2}{4} + (u - \frac{a}{2})^2] \leq a^2 + 2b$. As a result, the above inequality implies

$$(\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*) \leq \frac{9\lambda^2}{4a_0} |\mathbf{x}_1^*|_0 + 2C_3 \beta_n V^*, \quad (35)$$

where we have used $|\mathbf{x}_{1S}^*|_0 = |\mathbf{x}_1^*|_0$. Furthermore, (32) yields that

$$\begin{aligned} \widehat{\mathbf{x}}^\top \mathbf{Q} \widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \mathbf{Q} \mathbf{x}_1^* &\leq (\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*) + \frac{\lambda}{2} |\widehat{\mathbf{x}} - \mathbf{x}_1^*|_1 \\ &\leq (\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*) + 2\lambda |\mathbf{x}_1^*|_1 \\ &\leq (\widehat{\mathbf{x}} - \mathbf{x}_1^*)^\top \mathbf{Q}(\widehat{\mathbf{x}} - \mathbf{x}_1^*) + 4c_0^{-1/2} \bar{V}^* \lambda \sqrt{|\mathbf{x}_1^*|_0}. \end{aligned} \quad (36)$$

where the second inequality is because $|\widehat{\mathbf{x}} - \mathbf{x}_1^*|_1 \leq |\widehat{\mathbf{x}}|_1 + |\mathbf{x}_1^*|_1 \leq 4|\mathbf{x}_1^*|_1$, and the last inequality is from (29). Recall that $\lambda = C\eta \max\{k_0^{1/2} \lambda_0, s_0^{1/2} \Delta_n\} \bar{V}^*$. As a result,

$$\lambda \sqrt{|\mathbf{x}_1^*|_0} = C\eta \max\{k_0^{1/2} s_0^{1/2} \lambda_0, s_0 \Delta_n\} \bar{V}^* = C\eta T_n(\mathbf{x}^*) \bar{V}^*. \quad (37)$$

Combining (35), (36) and (37) gives

$$\begin{aligned} &\widehat{\mathbf{x}}^\top \mathbf{Q} \widehat{\mathbf{x}} - (\mathbf{x}_1^*)^\top \mathbf{Q} \mathbf{x}_1^* \\ &\leq \frac{9C^2}{4a_0} \eta^2 [T_n(\mathbf{x}^*)]^2 V^* + 4C c_0^{-1/2} \eta T_n(\mathbf{x}^*) V^* + 2C_3 \beta_n V^* \end{aligned}$$

$$\leq C_4\eta^2 T_n(\mathbf{x}^*)V^*. \quad (38)$$

Step 3: We use (38) to give a lower bound of $R(\widehat{\mathbf{x}})$.

Note that $R(\widehat{\mathbf{x}}) = (\mathbf{q}^\top \widehat{\mathbf{x}})^2 / (\widehat{\mathbf{x}}^\top \mathbf{Q} \widehat{\mathbf{x}})$. First, we look at the denominator $\widehat{\mathbf{x}}^\top \mathbf{Q} \widehat{\mathbf{x}}$. From (23) and that $t_n > 1/2$,

$$|t_n^{-2} - 1| = t_n^{-1}(1 + t_n^{-1})|t_n - 1| \leq 6c_0^{-1/2}\alpha_n \bar{V}^*.$$

Combining it to (38) and noting that $(\mathbf{x}_1^*)^\top \mathbf{Q} \mathbf{x}_1^* = t_n^{-2}(\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^* = t_n^{-2}V^*$, we have

$$\begin{aligned} \widehat{\mathbf{x}}^\top \mathbf{Q} \widehat{\mathbf{x}} &\leq [t_n^{-2} + C_4\eta^2 T_n(\mathbf{x}^*)](\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^* \\ &\leq [1 + 6c_0^{-1/2}\alpha_n \bar{V}^* + C_4\eta^2 T_n(\mathbf{x}^*)](\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^* \\ &\leq [1 + C_5\eta^2 T_n(\mathbf{x}^*)](\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^*. \end{aligned} \quad (39)$$

Second, we look at the numerator $\mathbf{q}^\top \widehat{\mathbf{x}}$. Since $\widehat{\mathbf{q}}^\top \widehat{\mathbf{x}} = 1$, by (29),

$$|\mathbf{q}^\top \widehat{\mathbf{x}} - 1| \leq |\widehat{\mathbf{q}} - \mathbf{q}|_\infty |\widehat{\mathbf{x}}|_1 \leq 6c_0^{-1/2}\alpha_n \bar{V}^* \leq C_6 T_n(\mathbf{x}^*). \quad (40)$$

Combining (39) and (40) gives

$$\begin{aligned} R(\widehat{\mathbf{x}}) &= \frac{(\mathbf{q}^\top \widehat{\mathbf{x}})^2}{\widehat{\mathbf{x}}^\top \mathbf{Q} \widehat{\mathbf{x}}} \geq \frac{[1 - C_6 T_n(\mathbf{x}^*)]^2}{1 + C_5\eta^2 T_n(\mathbf{x}^*)} \frac{1}{(\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^*} \\ &\geq [1 - A\eta^2 T_n(\mathbf{x}^*)] \frac{(\mathbf{q}^\top \mathbf{x}^*)^2}{(\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^*} \\ &= [1 - A\eta^2 T_n(\mathbf{x}^*)] R(\mathbf{x}^*), \end{aligned} \quad (41)$$

where $A = A(a_0, c_0, u_0)$ is a positive constant. \square

9.4 Proof of Proposition 6.1

Denote by $\mathbb{P}(i|j)$ is the probability that a new sample from class j is misclassified to class i , for $i, j \in \{1, 2\}$ and $i \neq j$. The classification error of h is

$$\text{err}(h) = \pi \mathbb{P}(2|1) + (1 - \pi) \mathbb{P}(1|2).$$

Write $M_k = M_k(\boldsymbol{\Omega}, \boldsymbol{\delta})$ and $L_k = L_k(\boldsymbol{\Omega}, \boldsymbol{\delta})$ for short. It suffices to show that

$$\begin{aligned} \mathbb{P}(2|1) &= \bar{\Phi} \left(\frac{(1-t)M}{\sqrt{L_1}} \right) + \frac{O(q) + o(d)}{L_1^{3/2}}, \\ \mathbb{P}(1|2) &= \bar{\Phi} \left(\frac{tM}{\sqrt{L_2}} \right) + \frac{O(q) + o(d)}{L_2^{3/2}}. \end{aligned}$$

We only consider $\mathbb{P}(2|1)$. The analysis of $\mathbb{P}(1|2)$ is similar. Suppose $\mathbf{X}|\text{class } 1 \stackrel{(d)}{=} \mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$. Define

$$\mathbf{Y} = \boldsymbol{\Sigma}_1^{-1/2}(\mathbf{Z} - \boldsymbol{\mu}_1),$$

so that $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and $\mathbf{Z} = \boldsymbol{\Sigma}_1^{1/2} \mathbf{Y} + \boldsymbol{\mu}_1$. Note that

$$Q(\mathbf{Z}) = (\boldsymbol{\Sigma}_1^{1/2} \mathbf{Y} + \boldsymbol{\mu}_1)^\top \boldsymbol{\Omega} (\boldsymbol{\Sigma}_1^{1/2} \mathbf{Y} + \boldsymbol{\mu}_1) - 2(\boldsymbol{\Sigma}_1^{1/2} \mathbf{Y} + \boldsymbol{\mu}_1)^\top \boldsymbol{\delta}$$

$$= \mathbf{Y}^\top \Sigma_1^{1/2} \Omega \Sigma_1^{1/2} \mathbf{Y} + 2\mathbf{Y}^\top \Sigma_1^{1/2} (\Omega \boldsymbol{\mu}_1 - \boldsymbol{\delta}) + \boldsymbol{\mu}_1^\top \Omega \boldsymbol{\mu}_1 - 2\boldsymbol{\mu}_1^\top \boldsymbol{\delta}. \quad (42)$$

Recall that $\Sigma_1^{1/2} \Omega \Sigma_1^{1/2} = \mathbf{K}_1 \mathbf{S}_1 \mathbf{K}_1^\top$ is the eigen-decomposition by excluding the 0 eigenvalues. Since Σ_1 has full rank and the rank of Ω is q , the rank of $\Sigma_1^{1/2} \Omega \Sigma_1^{1/2}$ is q . Therefore, \mathbf{S}_1 is a $q \times q$ diagonal matrix and \mathbf{K}_1 is a $d \times q$ matrix satisfying $\mathbf{K}_1^\top \mathbf{K}_1 = \mathbf{I}_q$. Let $\tilde{\mathbf{K}}_1$ be any $d \times (d-q)$ matrix such that $\mathbf{K} = [\mathbf{K}_1, \tilde{\mathbf{K}}_1]$ is a $d \times d$ orthogonal matrix. Since $\mathbf{I}_d = \mathbf{K} \mathbf{K}^\top = \mathbf{K}_1 \mathbf{K}_1^\top + \tilde{\mathbf{K}}_1 \tilde{\mathbf{K}}_1^\top$, we have

$$\mathbf{Y}^\top \Sigma_1^{1/2} (\Omega \boldsymbol{\mu}_1 - \boldsymbol{\delta}) = \mathbf{Y}^\top \mathbf{K}_1 \mathbf{K}_1^\top \Sigma_1^{1/2} (\Omega \boldsymbol{\mu}_1 - \boldsymbol{\delta}) + \mathbf{Y}^\top \tilde{\mathbf{K}}_1 \tilde{\mathbf{K}}_1^\top \Sigma_1^{1/2} (\Omega \boldsymbol{\mu}_1 - \boldsymbol{\delta}).$$

We recall that $\boldsymbol{\beta}_1 = \mathbf{K}_1^\top \Sigma_1^{1/2} (\Omega \boldsymbol{\mu}_1 - \boldsymbol{\delta})$. Let $\tilde{\boldsymbol{\beta}}_1 = \tilde{\mathbf{K}}_1^\top \Sigma_1^{1/2} (\Omega \boldsymbol{\mu}_1 - \boldsymbol{\delta})$, $\mathbf{W} = \mathbf{K}_1^\top \mathbf{Y}$, $\tilde{\mathbf{W}} = \tilde{\mathbf{K}}_1^\top \mathbf{Y}$ and $c_1 = \boldsymbol{\mu}_1^\top \Omega \boldsymbol{\mu}_1 - 2\boldsymbol{\mu}_1^\top \boldsymbol{\delta}$. It follows from (42) that

$$\begin{aligned} Q(\mathbf{Z}) &= \mathbf{Y}^\top \mathbf{K}_1 \mathbf{S}_1 \mathbf{K}_1^\top \mathbf{Y} + 2\mathbf{Y}^\top \mathbf{K}_1 \boldsymbol{\beta}_1 + 2\mathbf{Y}^\top \tilde{\mathbf{K}}_1 \tilde{\boldsymbol{\beta}}_1 + c_1 \\ &= \mathbf{W}^\top \mathbf{S}_1 \mathbf{W} + 2\mathbf{W}^\top \boldsymbol{\beta}_1 + 2\tilde{\mathbf{W}}^\top \tilde{\boldsymbol{\beta}}_1 + c_1 \\ &\equiv \bar{Q}_1(\mathbf{W}) + \bar{F}_1(\tilde{\mathbf{W}}) + c_1, \end{aligned}$$

where $\bar{Q}_1(\mathbf{w}) = \mathbf{w}^\top \mathbf{S}_1 \mathbf{w} + 2\mathbf{w}^\top \boldsymbol{\beta}_1$ and $\bar{F}_1(\tilde{\mathbf{w}}) = 2\tilde{\mathbf{w}}^\top \tilde{\boldsymbol{\beta}}_1$. Therefore,

$$\mathbb{P}(2|1) = \mathbb{P}(Q(\mathbf{Z}) > c) = \mathbb{P}(\bar{Q}_1(\mathbf{W}) + \bar{F}_1(\tilde{\mathbf{W}}) > c - c_1).$$

We write for convenience $\mathbf{W} = (W_1, \dots, W_q)^\top$, $\tilde{\mathbf{W}} = (W_{q+1}, \dots, W_d)^\top$, $\boldsymbol{\beta}_1 = (\beta_{11}, \dots, \beta_{1q})^\top$ and $\tilde{\boldsymbol{\beta}}_1 = (\beta_{1(q+1)}, \dots, \beta_{1d})^\top$, and notice that $W_i \stackrel{iid}{\sim} N(0, 1)$ for $1 \leq i \leq d$. Moreover,

$$\bar{Q}_1(\mathbf{W}) + \bar{F}_1(\tilde{\mathbf{W}}) = \sum_{i=1}^q (s_i W_i^2 + 2W_i \beta_{1i}) + \sum_{i=q+1}^d 2W_i \beta_{1i} \equiv \sum_{i=1}^d \xi_i, \quad (43)$$

where $\xi_i = s_i W_i^2 I\{1 \leq i \leq q\} + 2W_i \beta_{1i}$, for $1 \leq i \leq d$. The right hand side of (43) is a sum of independent variables, so we can apply the Edgeworth expansion to its distribution function, as described in detail below.

Note that $\mathbb{E}(W_i^2) = 1$, $\mathbb{E}(W_i^4) = 3$, $\mathbb{E}(W_i^6) = 15$ and $\mathbb{E}(W_i^{2j+1}) = 0$ for nonnegative integers j . By direct calculation,

$$\begin{aligned} \eta_1 &\equiv \sum_{i=1}^d \mathbb{E}(\xi_i) = \sum_{i=1}^q s_i = \text{tr}(\mathbf{S}_1) = \text{tr}(\Omega \Sigma_1), \\ \eta_2 &\equiv \sum_{i=1}^d \text{var}(\xi_i) = \sum_{i=1}^q (2s_i^2 + 4\beta_{1i}^2) + \sum_{i=q+1}^d 4\beta_{1i}^2 \\ &= 2 \text{tr}(\mathbf{S}_1^2) + 4|\boldsymbol{\beta}_1|^2 + 4|\tilde{\boldsymbol{\beta}}_1|^2 \\ &= 2 \text{tr}(\Omega \Sigma_1 \Omega \Sigma_1) + 4(\Omega \boldsymbol{\mu}_1 - \boldsymbol{\delta})^\top \Sigma_1 (\Omega \boldsymbol{\mu}_1 - \boldsymbol{\delta}), \\ \eta_3 &\equiv \sum_{i=1}^d \mathbb{E}[\xi_i - \mathbb{E}(\xi_i)]^3 = \sum_{i=1}^q (8s_i^3 + 24\beta_{1i}^2 s_i) \\ &= 8 \text{tr}(\mathbf{S}_1^3) + 24\boldsymbol{\beta}_1^\top \mathbf{S}_1 \boldsymbol{\beta}_1 \\ &= 8 \text{tr}[(\Omega \Sigma_1)^3] + 24(\Omega \boldsymbol{\mu}_1 - \boldsymbol{\delta})^\top \Sigma_1 \Omega \Sigma_1 (\Omega \boldsymbol{\mu}_1 - \boldsymbol{\delta}). \end{aligned} \quad (44)$$

Notice that $\mathbb{E}(|\xi_i - \mathbb{E}(\xi_i)|^3) < \infty$, as $\max\{|s_i|, |\beta_{1i}|, 1 \leq i \leq d\} \leq C_0$ by assumption. Therefore, using results in Chapter XVI of Feller (1966), we know

$$\begin{aligned} \mathbb{P}(2|1) &= \mathbb{P}\left(\sum_{i=1}^d \xi_i > c - c_1\right) \\ &= \mathbb{P}\left(\frac{\sum_{i=1}^d \xi_i - \mathbb{E}(\sum_{i=1}^d \xi_i)}{\sqrt{\sum_{i=1}^d \text{var}(\xi_i)}} > \frac{c - c_1 - \mathbb{E}(\sum_{i=1}^d \xi_i)}{\sqrt{\sum_{i=1}^d \text{var}(\xi_i)}}\right) \\ &= \bar{\Phi}\left(\frac{c - c_1 - \eta_1}{\sqrt{\eta_2}}\right) + \frac{\eta_3(1 - (\frac{c_1 - c + \eta_1}{\eta_2})^2)}{6\eta_2^{3/2}}\phi\left(\frac{c_1 - c + \eta_1}{\sqrt{\eta_2}}\right) + o\left(\frac{d}{\eta_2^{3/2}}\right), \end{aligned} \quad (45)$$

where ϕ is the probability density function of the standard normal distribution. It is observed that $\eta_2 = L_1(\boldsymbol{\Omega}, \boldsymbol{\delta})$ and $c_1 + \eta_1 = M_1(\boldsymbol{\Omega}, \boldsymbol{\delta})$. Also, $c = tM_1(\boldsymbol{\Omega}, \boldsymbol{\delta}) + (1-t)M_2(\boldsymbol{\Omega}, \boldsymbol{\delta})$. As a result,

$$\frac{c - c_1 - \eta_1}{\sqrt{\eta_2}} = \frac{[tM_1 + (1-t)M_2] - M_1}{\sqrt{L_1}} = \frac{(1-t)(M_2 - M_1)}{\sqrt{L_1}} = (1-t)\frac{M}{\sqrt{L_1}}.$$

Plugging it into (45), the first term is $\bar{\Phi}((1-t)\frac{M}{\sqrt{L_1}})$. Moreover, since the function $(1-u^2)\phi(u)$ is uniformly bounded, the second term is $O(\frac{\eta_3}{\eta_2^{3/2}})$. Here $\eta_2 = L_1$, and $\eta_3 = O(q)$ as s_i 's and β_{1i} 's are abounded in magnitude. Combining the above gives

$$\mathbb{P}(2|1) = \bar{\Phi}\left(\frac{(1-t)M}{\sqrt{L_1}}\right) + \frac{O(q) + o(d)}{L_1^{3/2}}.$$

The proof is now complete. \square

9.5 Proof of Proposition 6.2

Given $(\boldsymbol{\Omega}, \boldsymbol{\delta}, t)$, recall that $R_k = R_k(\boldsymbol{\Omega}, \boldsymbol{\delta})$, for $k = 1, 2$. Let $x_1 = [(1-t)^2 R_1]^{-1}$, $x_2 = [t^2 R_2]^{-1}$, and $x = \pi x_1 + (1-\pi)x_2$. By direct calculation,

$$\overline{\text{Err}}(\boldsymbol{\Omega}, \boldsymbol{\delta}, t) = \pi H(x_1) + (1-\pi)H(x_2), \quad H\left(\frac{\pi}{(1-t)^2 R_{\kappa(t)}}\right) = H(x). \quad (46)$$

Since H is twice continuously differentiable, from the Taylor expansion,

$$\begin{aligned} H(x_1) &= H(x) + H'(x)(x_1 - x) + \frac{1}{2}H''(z_1)(x_1 - x)^2, \\ H(x_2) &= H(x) + H'(x)(x_2 - x) + \frac{1}{2}H''(z_2)(x_2 - x)^2, \end{aligned}$$

where z_1 is a number between x_1 and x , and z_2 is a number between x_2 and x . Noticing that $\pi(x_1 - x) + (1-\pi)(x_2 - x) = 0$, we further obtain

$$\begin{aligned} \pi H(x_1) + (1-\pi)H(x_2) &= H(x) + \frac{\pi}{2}H''(z_1)(x_1 - x)^2 + \frac{1-\pi}{2}H''(z_2)(x_2 - x)^2 \\ &= H(x) + \frac{\pi(1-\pi)}{2}[(1-\pi)H''(z_1) + \pi H''(z_2)](x_1 - x_2)^2. \end{aligned}$$

Here, the second equality is because $x_1 - x = (1 - \pi)(x_1 - x_2)$ and $x_2 - x = \pi(x_2 - x_1)$. Let $A = \sup_{z \in [\min\{x_1, x_2\}, \max\{x_1, x_2\}]} |H''(z)|$. It follows that

$$|\pi H(x_1) + (1 - \pi)H(x_2) - H(x)| \leq \frac{\pi(1 - \pi)}{2} \cdot A(x_1 - x_2)^2. \quad (47)$$

Now, we bound $A(x_1 - x_2)^2$. Write for short $x_{\min} = \min\{x_1, x_2\}$ and $x_{\max} = \max\{x_1, x_2\}$. It is easy to see that

$$A(x_1 - x_2)^2 \leq \sup_{z \in [x_{\min}, x_{\max}]} |z^2 H''(z)| \cdot \left(\frac{x_1 - x_2}{x_{\min}}\right)^2. \quad (48)$$

By direct calculation, the function $z^2 H''(z)$ has the expression

$$z^2 H''(z) = \frac{1}{4\sqrt{2\pi}} \left(\frac{1}{z} - 3\right) e^{-\frac{1}{2z}} \frac{1}{z^{1/2}}.$$

On one hand, as $z \rightarrow \infty$, $|z^2 H''(z)| \leq Cz^{-1/2}$; on the other hand, as $z \rightarrow 0$, $|z^2 H''(z)| \leq Ce^{-1/(2z)} z^{-3/2} \leq Cz^{1/2}$. Therefore, we have $|z^2 H''(z)| \leq C \min\{\sqrt{z}, 1/\sqrt{z}\}$. Plugging it into (48), we have

$$A(x_1 - x_2)^2 \leq C \left[\max\{x_1 \wedge (1/x_1), x_2 \wedge (1/x_2)\} \right]^{1/2} \cdot \left(\frac{x_1 - x_2}{x_{\min}}\right)^2.$$

Note that $x_1 \wedge (1/x_1) = V_1$, $x_2 \wedge (1/x_2) = V_2$ and $|x_1 - x_2|/x_{\min} = |V - 1|$. It follows that

$$A(x_1 - x_2)^2 \leq C \left[\max\{V_1, V_2\} \right]^{1/2} \cdot |V - 1|^2. \quad (49)$$

We combine (46), (47) and (49), and note that $\pi(1 - \pi) \leq 1/4$. The first claim then follows.

When $t = 1/2$, we find $V - 1 = \Delta R/R_0$. The second claim follows immediately. \square

A Supplementary Proofs

A.1 Proof of Lemma 9.1

Consider $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Then $\mathbf{Y} \stackrel{(d)}{=} R\mathbf{U}$, where $R^2 \sim \chi_d^2$ and it is independent of \mathbf{U} . Since $\mathbb{E}(R^2) = d$ and $\text{var}(R^2) = 2d$, it follows that

$$\begin{aligned} \mathbb{E}(\mathbf{Y}^\top \mathbf{S}\mathbf{Y}) &= \mathbb{E}(R^2) \mathbb{E}(\mathbf{U}^\top \mathbf{S}\mathbf{U}) = d \mathbb{E}(\mathbf{U}^\top \mathbf{S}\mathbf{U}), \\ \mathbb{E}[(\mathbf{Y}^\top \mathbf{S}\mathbf{Y})^2] &= \mathbb{E}(R^4) \mathbb{E}[(\mathbf{U}^\top \mathbf{S}\mathbf{U})^2] = (d^2 + 2d) \mathbb{E}[(\mathbf{U}^\top \mathbf{S}\mathbf{U})^2], \\ \mathbb{E}(\mathbf{Y}^\top \boldsymbol{\beta}) &= \mathbb{E}(R) \mathbb{E}(\mathbf{U}^\top \boldsymbol{\beta}), \\ \mathbb{E}[(\mathbf{Y}^\top \boldsymbol{\beta})^2] &= \mathbb{E}(R^2) \mathbb{E}[(\mathbf{U}^\top \boldsymbol{\beta})^2] = d \mathbb{E}[(\mathbf{U}^\top \boldsymbol{\beta})^2], \\ \mathbb{E}(\mathbf{Y}^\top \mathbf{S}\mathbf{Y}\mathbf{Y}^\top \boldsymbol{\beta}) &= \mathbb{E}(R^3) \mathbb{E}(\mathbf{U}^\top \mathbf{S}\mathbf{U}\mathbf{U}^\top \boldsymbol{\beta}). \end{aligned}$$

First, note that $\mathbf{Y}^\top \boldsymbol{\beta} \sim N(0, \|\boldsymbol{\beta}\|^2)$. So $\mathbb{E}(\mathbf{Y}^\top \boldsymbol{\beta}) = 0$ and $\mathbb{E}[(\mathbf{Y}^\top \boldsymbol{\beta})^2] = \|\boldsymbol{\beta}\|^2$. We immediately have

$$\mathbb{E}(\mathbf{U}^\top \boldsymbol{\beta}) = 0, \quad \mathbb{E}[(\mathbf{U}^\top \boldsymbol{\beta})^2] = \|\boldsymbol{\beta}\|^2/d.$$

Second, write $\mathbf{S} = \text{diag}(s_1, \dots, s_d)$ and $\mathbf{Y}^\top \mathbf{S} \mathbf{Y} = \sum_{i=1}^d s_i Y_i^2$, where $Y_i \stackrel{iid}{\sim} N(0, 1)$. Therefore,

$$\begin{aligned}\mathbb{E}(\mathbf{Y}^\top \mathbf{S} \mathbf{Y}) &= \sum_{i=1}^d s_i \mathbb{E}(Y_i^2) = \sum_{i=1}^d s_i = \text{tr}(\mathbf{S}), \\ \text{var}(\mathbf{Y}^\top \mathbf{S} \mathbf{Y}) &= \sum_{i=1}^d s_i^2 \text{var}(Y_i^2) = \sum_{i=1}^d 2s_i^2 = 2 \text{tr}(\mathbf{S}^2).\end{aligned}$$

We immediately have

$$\mathbb{E}(\mathbf{U}^\top \mathbf{S} \mathbf{U}) = \text{tr}(\mathbf{S})/d, \quad \mathbb{E}[(\mathbf{U}^\top \mathbf{S} \mathbf{U})^2] = \frac{[\text{tr}(\mathbf{S})]^2 + 2 \text{tr}(\mathbf{S}^2)}{d^2 + 2d}.$$

Last, note that

$$\mathbf{Y}^\top \mathbf{S} \mathbf{Y} \mathbf{Y}^\top \boldsymbol{\beta} = \left(\sum_{i=1}^d s_i Y_i^2 \right) \left(\sum_{j=1}^d \beta_j Y_j \right) = \sum_{i=1}^d s_i \beta_i Y_i^3 + \sum_{i \neq j} s_i \beta_j Y_i^2 Y_j.$$

So $\mathbb{E}(\mathbf{Y}^\top \mathbf{S} \mathbf{Y} \mathbf{Y}^\top \boldsymbol{\beta}) = 0$. Since $\mathbb{E}(R^3) \neq 0$, we immediately have $\mathbb{E}(\mathbf{U}^\top \mathbf{S} \mathbf{U} \mathbf{U}^\top \boldsymbol{\beta}) = 0$. \square

A.2 Proof of Lemma 9.2

Recall that

$$\begin{aligned}M(\boldsymbol{\Omega}, \boldsymbol{\delta}) &= -\boldsymbol{\mu}_1^\top \boldsymbol{\Omega} \boldsymbol{\mu}_1 + \boldsymbol{\mu}_2^\top \boldsymbol{\Omega} \boldsymbol{\mu}_2 + 2(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\delta} - \text{tr}(\boldsymbol{\Omega}(\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2)) \\ &= \text{tr}[\boldsymbol{\Omega}(\boldsymbol{\Sigma}_2 + \boldsymbol{\mu}_2 \boldsymbol{\mu}_2^\top - \boldsymbol{\Sigma}_1 - \boldsymbol{\mu}_1 \boldsymbol{\mu}_1^\top)] + 2(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\delta}.\end{aligned}$$

It is well known that for any matrices \mathbf{A} and \mathbf{B} , $\text{tr}(\mathbf{A}^\top \mathbf{B}) = \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})$. So we have

$$\begin{aligned}M(\boldsymbol{\Omega}, \boldsymbol{\delta}) &= \text{vec}(\boldsymbol{\Omega})^\top \text{vec}(\boldsymbol{\Sigma}_2 + \boldsymbol{\mu}_2 \boldsymbol{\mu}_2^\top - \boldsymbol{\Sigma}_1 - \boldsymbol{\mu}_1 \boldsymbol{\mu}_1^\top) + 2\boldsymbol{\delta}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ &= \mathbf{x}^\top \mathbf{q}.\end{aligned}$$

Moreover, for $k = 1, 2$,

$$\begin{aligned}L_k(\boldsymbol{\Omega}, \boldsymbol{\delta}) &= 2(1 + \gamma) \text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma}_k \boldsymbol{\Omega} \boldsymbol{\Sigma}_k) + \gamma [\text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma}_k)]^2 + 4(\boldsymbol{\Omega} \boldsymbol{\mu}_k - \boldsymbol{\delta})^\top \boldsymbol{\Sigma}_k (\boldsymbol{\Omega} \boldsymbol{\mu}_k - \boldsymbol{\delta}) \\ &= 2(1 + \gamma) \text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma}_k \boldsymbol{\Omega} \boldsymbol{\Sigma}_k) + \gamma [\text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma}_k)]^2 + 4\boldsymbol{\mu}_k^\top \boldsymbol{\Omega} \boldsymbol{\Sigma}_k \boldsymbol{\Omega} \boldsymbol{\mu}_k - 8\boldsymbol{\delta}^\top \boldsymbol{\Sigma}_k \boldsymbol{\Omega} \boldsymbol{\mu}_k + 4\boldsymbol{\delta}^\top \boldsymbol{\Sigma}_k \boldsymbol{\delta} \\ &= 2 \text{tr}[\boldsymbol{\Omega} \boldsymbol{\Sigma}_k \boldsymbol{\Omega} ((1 + \gamma) \boldsymbol{\Sigma}_k + 2\boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top)] - 8 \text{tr}(\boldsymbol{\delta}^\top \boldsymbol{\Sigma}_k \boldsymbol{\Omega} \boldsymbol{\mu}_k) + 4\boldsymbol{\delta}^\top \boldsymbol{\Sigma}_k \boldsymbol{\delta} + \gamma [\text{tr}(\boldsymbol{\Omega} \boldsymbol{\Sigma}_k)]^2.\end{aligned}$$

From linear algebra, $\text{tr}(\mathbf{A}^\top \mathbf{B} \mathbf{C} \mathbf{D}^\top) = \text{vec}(\mathbf{A})^\top (\mathbf{D} \otimes \mathbf{B}) \text{vec}(\mathbf{C})$. It follows that

$$\begin{aligned}L_k(\boldsymbol{\Omega}, \boldsymbol{\delta}) &= 2 \text{vec}(\boldsymbol{\Omega})^\top [((1 + \gamma) \boldsymbol{\Sigma}_k + 2\boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top) \otimes \boldsymbol{\Sigma}_k] \text{vec}(\boldsymbol{\Omega}) - 8\boldsymbol{\delta}^\top (\boldsymbol{\mu}_k^\top \otimes \boldsymbol{\Sigma}_k) \text{vec}(\boldsymbol{\Omega}) \\ &\quad + 4\boldsymbol{\delta}^\top \boldsymbol{\Sigma}_k \boldsymbol{\delta} + \gamma \text{vec}(\boldsymbol{\Omega})^\top \text{vec}(\boldsymbol{\Sigma}_k) \text{vec}(\boldsymbol{\Sigma}_k)^\top \text{vec}(\boldsymbol{\Omega}) \\ &= 2 \begin{bmatrix} \text{vec}(\boldsymbol{\Omega})^\top & \boldsymbol{\delta}^\top \end{bmatrix} \begin{bmatrix} ((1 + \gamma) \boldsymbol{\Sigma}_k + 2\boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top) \otimes \boldsymbol{\Sigma}_k & -2(\boldsymbol{\mu}_k \otimes \boldsymbol{\Sigma}_k) \\ -2(\boldsymbol{\mu}_k^\top \otimes \boldsymbol{\Sigma}_k) & 2\boldsymbol{\Sigma}_k \end{bmatrix} \begin{bmatrix} \text{vec}(\boldsymbol{\Omega}) \\ \boldsymbol{\delta} \end{bmatrix} \\ &\quad + \gamma \text{vec}(\boldsymbol{\Omega})^\top \text{vec}(\boldsymbol{\Sigma}_k) \text{vec}(\boldsymbol{\Sigma}_k)^\top \text{vec}(\boldsymbol{\Omega}) \\ &= \mathbf{x}^\top \mathbf{Q}_k \mathbf{x}.\end{aligned}$$

Note that $\mathbf{Q} = \mathbf{Q}_1 + \kappa \mathbf{Q}_2$ and $L = L_1 + \kappa L_2$. This immediately implies $L(\boldsymbol{\Omega}, \boldsymbol{\delta}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x}$. \square

A.3 Proof of Lemma 9.3

Let $\Delta_n = \max\{|\widehat{\Sigma}_k - \Sigma_k|_\infty, |\widehat{\boldsymbol{\mu}}_k - \boldsymbol{\mu}_k|_\infty, k = 1, 2\}$ to save notation. We recall that $|\widehat{\boldsymbol{\mu}}_k - \boldsymbol{\mu}_k|_\infty \leq |\boldsymbol{\mu}_k|_\infty \leq 1$ and $|\widehat{\Sigma}_k - \Sigma_k|_\infty \leq |\Sigma_k|_\infty$, for $k = 1, 2$, as assumed in the beginning of Section 5.2.

Consider $|\widehat{\mathbf{q}} - \mathbf{q}|_\infty$ first. Note that $|\text{vec}(\mathbf{A}) - \text{vec}(\mathbf{B})|_\infty = |\mathbf{A} - \mathbf{B}|_\infty$ for any matrices \mathbf{A} and \mathbf{B} . It follows that

$$\begin{aligned} |\widehat{\mathbf{q}} - \mathbf{q}|_\infty &\leq \sum_{k=1}^2 (|\widehat{\Sigma}_k - \Sigma_k|_\infty + |\widehat{\boldsymbol{\mu}}_k \widehat{\boldsymbol{\mu}}_k^\top - \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top|_\infty + 2|\widehat{\boldsymbol{\mu}}_k - \boldsymbol{\mu}_k|_\infty) \\ &\leq C\Delta_n + \sum_{k=1}^2 |\widehat{\boldsymbol{\mu}}_k \widehat{\boldsymbol{\mu}}_k^\top - \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top|_\infty. \end{aligned}$$

Write $\boldsymbol{\mu}_k = \boldsymbol{\mu}$ for short. We have $|\widehat{\boldsymbol{\mu}}(i)\widehat{\boldsymbol{\mu}}(j) - \boldsymbol{\mu}(i)\boldsymbol{\mu}(j)| \leq |\widehat{\boldsymbol{\mu}}(j)||\widehat{\boldsymbol{\mu}}(i) - \boldsymbol{\mu}(i)| + |\boldsymbol{\mu}(i)||\widehat{\boldsymbol{\mu}}(j) - \boldsymbol{\mu}(j)| \leq (|\widehat{\boldsymbol{\mu}}|_\infty + |\boldsymbol{\mu}|_\infty)|\widehat{\boldsymbol{\mu}} - \boldsymbol{\mu}|_\infty$. Since $|\widehat{\boldsymbol{\mu}}|_\infty \leq |\boldsymbol{\mu}|_\infty + |\widehat{\boldsymbol{\mu}} - \boldsymbol{\mu}|_\infty \leq 2|\boldsymbol{\mu}|_\infty \leq 2$, it follows that $|\widehat{\boldsymbol{\mu}}(i)\widehat{\boldsymbol{\mu}}(j) - \boldsymbol{\mu}(i)\boldsymbol{\mu}(j)| \leq 3|\widehat{\boldsymbol{\mu}} - \boldsymbol{\mu}|_\infty \leq 3\Delta_n$. As a result,

$$|\widehat{\boldsymbol{\mu}}_k \widehat{\boldsymbol{\mu}}_k^\top - \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top|_\infty \leq 3\Delta_n, \quad k = 1, 2. \quad (50)$$

We immediately have $|\widehat{\mathbf{q}} - \mathbf{q}|_\infty \leq C\Delta_n$.

Next, consider $|\widehat{\mathbf{Q}} - \mathbf{Q}|_\infty$. It is easy to see that

$$\begin{aligned} |\widehat{\mathbf{Q}}_k - \mathbf{Q}_k|_\infty &\leq 2(1 + \gamma)|\widehat{\Sigma}_k \otimes \widehat{\Sigma}_k - \Sigma_k \otimes \Sigma_k|_\infty + 4|(\widehat{\boldsymbol{\mu}}_k \widehat{\boldsymbol{\mu}}_k^\top) \otimes \widehat{\Sigma}_k - (\boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top) \otimes \Sigma_k|_\infty \\ &\quad + \gamma|\text{vec}(\widehat{\Sigma}_k) \text{vec}(\widehat{\Sigma}_k)^\top - \text{vec}(\Sigma_k) \text{vec}(\Sigma_k)^\top|_\infty \\ &\quad + 4|\widehat{\boldsymbol{\mu}}_k \otimes \widehat{\Sigma}_k - \boldsymbol{\mu}_k \otimes \Sigma_k|_\infty + 4|\widehat{\Sigma}_k - \Sigma_k|_\infty. \end{aligned}$$

Here $|\widehat{\Sigma}_k - \Sigma_k|_\infty \leq \Delta_n$, and using a similar argument as in (50), we can show that $|\text{vec}(\widehat{\Sigma}_k) \text{vec}(\widehat{\Sigma}_k)^\top - \text{vec}(\Sigma_k) \text{vec}(\Sigma_k)^\top|_\infty \leq C|\text{vec}(\widehat{\Sigma}_k) - \text{vec}(\Sigma_k)|_\infty \leq C\Delta_n$. To bound the other terms, it suffices to show that

$$|\mathbf{A} \otimes \mathbf{B} - \mathbf{A}' \otimes \mathbf{B}'|_\infty \leq C(|\mathbf{A} - \mathbf{A}'|_\infty + |\mathbf{B} - \mathbf{B}'|_\infty), \quad (51)$$

when $|\mathbf{A} - \mathbf{A}'|_\infty \leq |\mathbf{A}|_\infty \leq C$ and $|\mathbf{B} - \mathbf{B}'|_\infty \leq |\mathbf{B}|_\infty \leq C$. To see this, note that $|A(i, j)B(k, l) - A'(i, j)B'(k, l)| \leq |B'(k, l)||A(i, j) - A'(i, j)| + |A(i, j)||B(k, l) - B'(k, l)| \leq (|\mathbf{A}|_\infty + |\mathbf{B}'|_\infty)(|\mathbf{A} - \mathbf{A}'|_\infty + |\mathbf{B} - \mathbf{B}'|_\infty) \leq C(|\mathbf{A} - \mathbf{A}'|_\infty + |\mathbf{B} - \mathbf{B}'|_\infty)$. This proves (51). It follows immediately that $|\widehat{\mathbf{Q}} - \mathbf{Q}|_\infty \leq C\Delta_n$. \square

A.4 Proof of Lemma 9.4

Write $\mathbf{x}^* = \mathbf{x}_{\lambda_0}^*$ for short. From KKT conditions, there exists a dual variable θ such that

$$\mathbf{q}^\top \mathbf{x}^* = 1, \quad 2\mathbf{Q}\mathbf{x}^* + \lambda_0 \widetilde{\text{sign}}(\mathbf{x}^*) = \theta \mathbf{q}.$$

Here $\widetilde{\text{sign}}(\mathbf{x}^*)$ is the vector whose j -th coordinate is $-1, 1$ or any value between $[-1, 1]$ when $x_j^* < 0, x_j^* > 0$ and $x_j^* = 0$, respectively. We multiply both sides of the second equation by $(\mathbf{x}^*)^\top$, and note that $(\mathbf{x}^*)^\top \widetilde{\text{sign}}(\mathbf{x}^*) = |\mathbf{x}^*|_1$. It follows that $\theta = 2(\mathbf{x}^*)^\top \mathbf{Q}\mathbf{x}^* + \lambda_0 |\mathbf{x}^*|_1$ and

$$\begin{aligned} \mathbf{Q}\mathbf{x}^* &= \left[(\mathbf{x}^*)^\top \mathbf{Q}\mathbf{x}^* + \frac{\lambda_0}{2} |\mathbf{x}^*|_1 \right] \mathbf{q} - \frac{\lambda_0}{2} \widetilde{\text{sign}}(\mathbf{x}^*) \\ &= V^* \mathbf{q} + \frac{\lambda_0}{2} [|\mathbf{x}^*|_1 \mathbf{q} - \widetilde{\text{sign}}(\mathbf{x}^*)], \end{aligned} \quad (52)$$

where $V^* = (\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^* = [R(\mathbf{x}^*)]^{-1}$. Let $\bar{V}^* = (V^*)^{1/2}$. Then $\Gamma(\mathbf{x}) = \frac{1}{\bar{V}^*} |\mathbf{Q} \mathbf{x}^* - V^* \mathbf{q}|_\infty$, and (52) implies

$$\Gamma(\mathbf{x}^*) \leq \frac{\lambda_0}{2\bar{V}^*} (|\mathbf{x}^*|_1 |\mathbf{q}|_\infty + 1) = \frac{|\mathbf{q}|_\infty}{2\bar{V}^*} \lambda_0 |\mathbf{x}^*|_1 + \frac{\lambda_0}{2} [R(\mathbf{x}^*)]^{1/2}. \quad (53)$$

From the Cauchy-Schwartz inequality,

$$|\mathbf{x}^*|_1 \leq |\mathbf{x}^*|_0^{1/2} |\mathbf{x}^*| \leq |\mathbf{x}^*|_0^{1/2} \frac{1}{\sqrt{c_0}} \sqrt{(\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^*} = \frac{\bar{V}^*}{\sqrt{c_0}} |\mathbf{x}^*|_0^{1/2},$$

where we have used the fact that $(\mathbf{x}^*)^\top \mathbf{Q} \mathbf{x}^* \geq \lambda_{\min}(\mathbf{Q}_{SS}) |\mathbf{x}^*|^2 \geq c_0 |\mathbf{x}^*|^2$. Plugging it into (53) gives

$$\begin{aligned} \Gamma(\mathbf{x}^*) &\leq \frac{|\mathbf{q}|_\infty}{2\sqrt{c_0}} \lambda_0 |\mathbf{x}^*|_0^{1/2} + \frac{\lambda_0}{2} [R(\mathbf{x}^*)]^{1/2} \\ &\leq C_1 \lambda_0 [\max\{|\mathbf{x}^*|_0, R(\mathbf{x}^*)\}]^{1/2}, \end{aligned}$$

where C_1 is a constant that only depends on c_0 (noting that $|\mathbf{q}|_\infty \leq C \max\{|\boldsymbol{\mu}_k|_\infty, |\boldsymbol{\Sigma}_k|_\infty, k = 1, 2\} \leq C$).

References

- Bindea, G., Mlecnik, B., Hackl, H., Charoentong, P., Tosolini, M., Kirilovsky, A., Fridman, W.-H., Pagès, F., Trajanoski, Z., and Galon, J. (2009), “ClueGO: a Cytoscape plug-in to decipher functionally grouped gene ontology and pathway annotation networks,” *Bioinformatics*, 25, 1091–1093.
- Cai, T. and Liu, W. (2011), “A direct estimation approach to sparse linear discriminant analysis,” *J. Amer. Statist. Assoc.*, 106, 1566–1577.
- Cai, T., Liu, W., and Luo, X. (2011), “A Constrained L1 Minimization Approach to Sparse Precision Matrix Estimation,” *Journal of the American Statistical Association*, 106, 594–607.
- Catoni, O. (2012), “Challenging the empirical mean and empirical variance: a deviation study,” in *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, Institut Henri Poincaré, vol. 48, pp. 1148–1185.
- Fan, J. and Fan, Y. (2008), “High-dimensional classification using features annealed independence rules,” *Ann. Statist.*, 36, 2605–2637.
- Fan, J., Feng, Y., and Tong, X. (2012), “A ROAD to classification in high dimensional space,” *J. Roy. Statist. Soc. B*, 74.
- Feller, W. (1966), *An Introduction to Probability Theory and Its Applications*, vol. II, John Wiley and Sons, Inc. New York, London.
- Guo, Y., Hastie, T., and Tibshirani, R. (2005), “Regularized discriminant analysis and its application in microarrays,” *Biostatistics*, 1, 1–18.
- Han, F. and Liu, H. (2012), “Transelliptical component analysis,” *Advances in Neural Information Processing Systems*, 25, 368–376.

- Han, F., Zhao, T., and Liu, H. (2013), “CODA: High Dimensional Copula Discriminant Analysis,” *Journal of Machine Learning Research*, 14, 629–671.
- Kendall, M. (1938), “A New Measure Of Rank Correlation,” *Biometrika*, 30, 81–93.
- Liu, H., Han, F., Yuan, M., Lafferty, J., and Wasserman, L. (2012), “High-dimensional Semiparametric Gaussian Copula Graphical Models,” *The Annals of Statistics*, 40, 2293–2326.
- Luparello, C. (2013), “Aspects of Collagen Changes in Breast Cancer,” *J Carcinogene Mutagene S*, 13.
- Maruyama, Y. and Seo, T. (2003), “Estimation of moment parameter in elliptical distributions,” *Journal of Japanese Statistical Society*, 33, 215–229.
- Shao, J., Wang, Y., Deng, X., and Wang, S. (2011), “Sparse linear discriminant analysis by thresholding for high dimensional data,” *Ann. Statist.*, 39, 1241–1265.
- Wei, Z. and Li, H. (2007), “A Markov random field model for network-based analysis of genomic data,” *Bioinformatics*, 23, 1537–1544.
- Zhao, T., Roeder, K., and Liu, H. (2013), “Positive Semidefinite Rank-based Correlation Matrix Estimation with Application to Semiparametric Graph Estimation,” *manuscript*.