

Quadratic Discriminant Revisited

Wenbo Cao, Robert Haralick

Dept. of Computer Science, The Graduate Center, City University of New York, New York, NY 10011 USA

Email: wenbo.cao@gmail.com, haralick@aim.com

Abstract—In this study, we revisit quadratic discriminant analysis (QDA). For this purpose, we present a majorize-minimize (MM) optimization algorithm to estimate parameters for generative classifiers, of which conditional distributions are from the exponential family. Furthermore, we propose a block-coordinate descent algorithm to sequentially update parameters of QDA in each iteration of the MM algorithm; for each update, we apply a trust region method, of which each iteration has a simple closed form solution. Numerical experiments show that: when compared with conjugate gradient method, the new proposed method is faster in 9 of 10 benchmark data sets; when compared with other widely used quadratic classifiers in the literature, QDA trained with the proposed method is either the best or not statistically significantly different from the best ones in 8 of 10 benchmark data sets.

I. INTRODUCTION

Quadratic discriminant analysis (QDA) is a classical generative probabilistic method for classification problems. Essentially speaking, conditional distributions of QDA are assumed to be multivariate Gaussian; posterior distributions are derived via Bayes theorem, and used to classify observations into different classes [1], [2]. Traditionally, parameters of QDA are estimated by maximizing the joint likelihood of observations and their associate class labels. Though computational efficient, this generative approach does not aim at reducing classification error, and is not robust to model mis-specification. For recent studies on QDA, we refer to [3] and references thereof.

Generative and discriminative learning are two commonly used approaches for estimating parameters of given classifiers. It is well understood that generative learning is more statistically efficient when models are well specified; discriminative learning can achieve better classification accuracy when models are mis-specified [4], [5]. In the hope of taking advantages of both generative and discriminative learning, several authors proposed hybrid generative/discriminative learning for classification methods [6]–[8].

Motivated by recent works on hybrid learning, we propose to estimate parameters of QDA by maximizing a convex combination of the joint log-likelihood and conditional log-likelihood of given observations and their class labels. Our main contributions of this study are: 1) we present a general iterative algorithm for parameter estimation when conditional distributions of classifiers are from the well known exponential family; and 2) we revisit QDA, and present a block coordinate descent algorithm for solving a specific optimization problem derived for QDA.

The remainder of the paper is organized as follows. In section II, we give a brief introduction of generative classifiers and commonly used approaches for parameter estimation. In section III, we present an iterative optimization algorithm for

parameter estimation when conditional distribution are from the exponential family. Multivariate Gaussian is a specific distribution of the exponential family; we revisit QDA and present a specific block-coordinate descent algorithm for estimating parameters of QDA in section IV. We present experimental results in section V, and conclude the study with the summary of our work, and possible future directions in section VI.

Throughout this paper, we use the following conventions unless otherwise specified: (1) vectors are column vectors and are represented as bold small letters, e.g. \mathbf{x} ; (2) matrices are represented as bold capital letters, e.g. \mathbf{X} ; (3) \mathbf{I} is an identity matrix; (4) \square^T is the transpose of a vector or matrix \square ; (5) $\det(\mathbf{X})$ is the determinant of \mathbf{X} ; (6) \mathbf{X}^{-1} is the inverse of \mathbf{X} ; (7) $\mathbf{X} \succeq \mathbf{Y}$ means that matrix $\mathbf{X} - \mathbf{Y}$ is positive semidefinite; and (8) we define the inner product of matrices $\mathbf{X} \in \mathbb{R}^{m \times n}$ and $\mathbf{Y} \in \mathbb{R}^{m \times n}$ as

$$\langle \mathbf{X}, \mathbf{Y} \rangle = \sum_{i=1}^m \sum_{j=1}^n X_{ij} Y_{ij}.$$

II. GENERATIVE MODELS FOR CLASSIFICATION

We consider a generative probabilistic model for classification problems. For this purpose, we assume, for the k -th class, a parametric density, $p(\mathbf{x}|k; \Theta_k)$, which is described by a set of parameters Θ_k ; here $k = 1, 2, \dots, K$. We intend to keep the densities abstract at this moment, and will specify them in section III. Define prior probabilities as follows

$$p(k) = \frac{e^{v_k}}{\sum_{m=1}^K e^{v_m}}. \quad (1)$$

Posterior probabilities can be derived via Bayes theorem:

$$p(k|\mathbf{x}; \Theta, \Upsilon) = \frac{p(\mathbf{x}|k; \Theta_k) e^{v_k}}{\sum_{k=1}^K p(\mathbf{x}|k; \Theta_k) e^{v_k}}, \quad (2)$$

where $\Theta = (\Theta_1, \Theta_2, \dots, \Theta_K)$ and $\Upsilon = (v_1, v_2, \dots, v_K)$ are unknown parameters. To minimize classification loss for 0-1 cost, we would assign an observation to a class which has the highest posterior probability [1]. Formally speaking, the decision rule to classify \mathbf{x} is as follows

$$\hat{y} = \arg \max_k p(k|\mathbf{x}; \Theta, \Upsilon). \quad (3)$$

Consider a given training data set

$$\mathcal{D} = \{(\mathbf{x}_n, y_n) | n = 1, 2, \dots, N, \mathbf{x}_n \in \mathbb{R}^d, y_n \in \{1, 2, \dots, K\}\},$$

where \mathbf{x}_n is the n -th observation, and y_n is its class label. Define index set \mathcal{C}_k that contains the indices of all observations of class k in \mathcal{D} ; that is,

$$\mathcal{C}_k = \{n | y_n = k, n = 1, 2, \dots, N\}, \quad (4)$$

where $k = 1, 2, \dots, K$. Let N_k be the number of observations of class k in the data set \mathcal{D} .

Our main focus is how to estimate the parameters of a generative model, Θ and Υ , using observations in the data set \mathcal{D} . Let \mathcal{L}_G and \mathcal{L}_D be the scaled negative joint log-likelihood and the scaled negative conditional log-likelihood, which are defined as follows

$$\begin{aligned}\mathcal{L}_G(\Theta, \Upsilon) &= -\frac{1}{N} \sum_{k=1}^K \sum_{n \in \mathcal{C}_k} \log[p(\mathbf{x}_n | k; \Theta_k) \frac{e^{v_k}}{\sum_{m=1}^K e^{v_m}}] \\ \mathcal{L}_D(\Theta, \Upsilon) &= -\frac{1}{N} \sum_{n=1}^N \log p(y_n | \mathbf{x}_n; \Theta, \Upsilon).\end{aligned}$$

One approach for estimating Θ and Υ is to maximize \mathcal{L}_G . Alternatively, Θ and Υ can be estimated by maximizing \mathcal{L}_D .

Motivated by the recent works on hybrid generative/discriminative learning [6]–[8], we minimize the following objective function for estimating Θ and Υ in this study.

$$\mathcal{L}(\Theta, \Upsilon) = \beta \mathcal{L}_D(\Theta, \Upsilon) + (1 - \beta) \mathcal{L}_G(\Theta, \Upsilon), \quad (5)$$

where $0 \leq \beta \leq 1$. As β varies from 0 to 1, we can interpolate between the two approaches mentioned above.

III. LEARNING FOR EXPONENTIAL FAMILY

In this study, we assume conditional distributions are from the exponential family. That is,

$$p(\mathbf{x} | k; \Theta_k) = \exp\{h_k(\mathbf{x}) + \langle \Theta_k, T_k(\mathbf{x}) \rangle - A_k(\Theta_k)\}, \quad (6)$$

where $T_k(\mathbf{x})$ is the potential function or sufficient statistics; $h_k(\mathbf{x})$ is a function that does not depend on Θ_k ; $A_k(\Theta_k)$ is the log partition function. For details about the exponential family, we refer to [9]. Canonical parameter $\Theta_k \in \mathcal{P}_k$ needs to be estimated; \mathcal{P}_k is the domain of Θ_k , which we assume to be convex in this study. We intend to keep $h_k(\mathbf{x})$'s and $T_k(\mathbf{x})$'s, $A_k(\Theta_k)$'s and \mathcal{P}_k 's abstract at this moment, and will specify them when needed.

Following the definition in equation (5), we can write the objective function to be minimized as follows

$$\begin{aligned}\mathcal{L}(\Theta, \Upsilon) &= -\frac{1}{N} \sum_{k=1}^K \sum_{n \in \mathcal{C}_k} \log[p(\mathbf{x}_n | k; \Theta_k) \frac{e^{v_k}}{\sum_{m=1}^K e^{v_m}}] \\ &\quad + \frac{\beta}{N} \sum_{n=1}^N \log \sum_{k=1}^K [e^{v_k} p(\mathbf{x}_n | k; \Theta_k) \frac{e^{v_k}}{\sum_{m=1}^K e^{v_m}}].\end{aligned} \quad (7)$$

The optimization problem for estimating parameters Θ and Υ can be written as

$$\begin{aligned}\min_{\Theta, \Upsilon} \quad & \mathcal{L}(\Theta, \Upsilon) \\ \text{s.t.} \quad & \Theta_k \in \mathcal{P}_k \quad k = 1, \dots, K\end{aligned} \quad (8)$$

We note that the optimization problem is not convex, Thus we will be satisfied to find local optimums. In the following, we shall present a majorize-minimize (MM) algorithm in which we iteratively minimize a convex upper bound of the original objective function. For an introduction of the MM optimization method, we refer to [10] and references therein.

We note that, due to log-partition functions $A_k(\Theta_k)$'s, the log-sum terms of \mathcal{L} are generally not convex. Several lower bounds of log-partition functions have been proposed in the literature, for example, [11], [12]. As global bounds, they are complex and loose. Note that one important property of log-partition functions $A_k(\Theta_k)$, ($k = 1, 2, \dots, K$), is that they are convex functions of Θ_k . Therefore we have

$$A_k(\Theta_k) \geq A_k(\tilde{\Theta}_k) + \left\langle \frac{\partial A_k}{\partial \Theta_k} \Big|_{\tilde{\Theta}_k}, \Theta_k - \tilde{\Theta}_k \right\rangle. \quad (9)$$

Define an auxiliary function for $p(\mathbf{x} | k; \Theta_k)$ at $\tilde{\Theta}_k$ as follows,

$$q(\mathbf{x} | k; \Theta_k, \tilde{\Theta}_k) = e^{h_k(\mathbf{x}) + \langle \Theta_k, T_k(\mathbf{x}) \rangle - A_k(\tilde{\Theta}_k) - \left\langle \frac{\partial A_k}{\partial \Theta_k} \Big|_{\tilde{\Theta}_k}, \Theta_k - \tilde{\Theta}_k \right\rangle}. \quad (10)$$

Consequently, we have $p(\mathbf{x} | k; \Theta_k) \leq q(\mathbf{x} | k; \Theta_k, \tilde{\Theta}_k)$. When $\Theta_k = \tilde{\Theta}_k$, $p(\mathbf{x} | k; \Theta_k) = q(\mathbf{x} | k; \Theta_k, \tilde{\Theta}_k)$. Define

$$\begin{aligned}\bar{\mathcal{L}}(\Theta, \Upsilon; \tilde{\Theta}) &= \frac{1}{N} \sum_{k=1}^K \sum_{n \in \mathcal{C}_k} \log[p(\mathbf{x}_n | k; \Theta_k) \frac{e^{v_k}}{\sum_{m=1}^K e^{v_m}}] \\ &\quad - \frac{\beta}{N} \sum_{n=1}^N \log \sum_{k=1}^K [q(\mathbf{x}_n | k; \Theta_k, \tilde{\Theta}_k) \frac{e^{v_k}}{\sum_{m=1}^K e^{v_m}}]\end{aligned} \quad (11)$$

Then we have $\bar{\mathcal{L}}(\Theta, \Upsilon; \tilde{\Theta})$ majorizes $\mathcal{L}(\Theta, \Upsilon)$ at $\tilde{\Theta}$: $\mathcal{L}(\Theta, \Upsilon) \geq \bar{\mathcal{L}}(\Theta, \Upsilon; \tilde{\Theta})$, and $\mathcal{L}(\tilde{\Theta}, \Upsilon) = \bar{\mathcal{L}}(\tilde{\Theta}, \Upsilon; \tilde{\Theta})$. Moreover, $\bar{\mathcal{L}}$ is a convex function.

Thus we have the following convex optimization problem

$$\begin{aligned}\min_{\Theta, \Upsilon} \quad & \bar{\mathcal{L}}(\Theta, \Upsilon; \tilde{\Theta}) \\ \text{s.t.} \quad & \Theta_k \in \mathcal{P}_k, \quad k = 1, \dots, K\end{aligned} \quad (12)$$

Let Θ^* and Υ^* be optimal solutions of optimization problem (12). We note that, for any $\tilde{\Theta}$ and $\tilde{\Upsilon}$,

$$\mathcal{L}(\Theta^*, \Upsilon^*) \leq \bar{\mathcal{L}}(\Theta^*, \Upsilon^*; \tilde{\Theta}) \leq \bar{\mathcal{L}}(\tilde{\Theta}, \tilde{\Upsilon}; \tilde{\Theta}) = \mathcal{L}(\tilde{\Theta}, \tilde{\Upsilon}).$$

That is $\mathcal{L}(\Theta^*, \Upsilon^*)$ is less than or equal to $\mathcal{L}(\tilde{\Theta}, \tilde{\Upsilon})$. In principal, we can iteratively find a new point (Θ, Υ) to reduce the function value of \mathcal{L} until a local optimum is reached. We present in algorithm 1 the pseudo-code of the MM algorithm in which the hybrid objective function for the exponential family is minimized. We note that, instead of finding optimums of the optimization problem (12) in each iteration, we can find a feasible point that reduces the objective function.

- 1: $\tau \leftarrow 0$ { τ : iteration count}
- 2: Initialize feasible $\Theta^{(\tau)}$, and $\Upsilon^{(\tau)}$ for problem (8)
- 3: **while** not converge **do**
- 4: $\tilde{\Theta} \leftarrow \Theta^{(\tau)}$
- 5: Obtain $(\Theta^{(\tau+1)}, \Upsilon^{(\tau+1)})$ by solving convex optimization problem (12);
- 6: $\tau \leftarrow \tau + 1$
- 7: **end while**

Fig. 1. A MM optimization algorithm for estimating parameters of the exponential family

IV. QDA REVISITED

In this section, we shall investigate how to apply algorithm 1 for estimation parameters of QDA. Define

$$\Theta_k = (\Omega_k, \theta_k) = (\Sigma_k^{-1}, \Sigma_k^{-1} \mu_k), \quad (13)$$

where Σ_k and μ_k are, respectively, the covariance matrix and the mean for the k -th conditional distribution. In the following discussion, we shall use Ω_k , θ_k and Σ_k , μ_k interchangeably. In QDA, conditional distributions are specified as multivariate Gaussian; thus for $k = 1, 2, \dots, K$, we have the densities of conditional distribution expressed in the standard form of the exponential family

$$p(\mathbf{x} | k; \Theta_k) = e^{h_k(\mathbf{x}) + \langle \Theta_k, T_k(\mathbf{x}) \rangle - A_k(\Theta_k)},$$

where

$$h_k(\mathbf{x}) = -\frac{d}{2} \log(2\pi), \quad T_k(\mathbf{x}) = \left(-\frac{1}{2} \mathbf{x} \mathbf{x}^T, \quad \mathbf{x}\right), \quad (14)$$

$$A_k(\Theta_k) = \frac{1}{2} \theta_k^T \Omega_k^{-1} \theta_k - \frac{1}{2} \log \det \Omega_k. \quad (15)$$

Note that

$$\frac{\partial}{\partial \Theta_k} A_k(\Theta_k) = \left(-\frac{1}{2} \Omega_k^{-1} \theta_k \theta_k^T \Omega_k^{-1} - \frac{1}{2} \Omega_k^{-1}, \quad \mu_k \right).$$

We thus have the auxiliary function for $p(\mathbf{x}|k; \Theta_k)$ at $\tilde{\Theta}_k = (\tilde{\Omega}_k, \tilde{\theta}_k)$ as following,

$$q(\mathbf{x}|k; \Theta_k, \tilde{\Theta}_k) = e^{\frac{1}{2} \langle \tilde{\mathbf{M}}_k - \mathbf{x} \mathbf{x}^T, \Omega_k \rangle + \langle \mathbf{x} - \tilde{\mu}_k, \theta_k \rangle + c_k}, \quad (16)$$

where

$$\begin{aligned} \tilde{\mathbf{M}}_k &= \tilde{\Omega}_k^{-1} + \tilde{\Omega}_k^{-1} \tilde{\theta}_k \tilde{\theta}_k^T \tilde{\Omega}_k^{-1} \\ c_k &= -\frac{d}{2} \log(2\pi) - A_k(\tilde{\Theta}_k) + \left\langle \tilde{\mu}_k, \tilde{\theta}_k \right\rangle - \frac{1}{2} \left\langle \tilde{\mathbf{M}}_k, \tilde{\Omega}_k \right\rangle. \end{aligned}$$

The convex upper bound of $\mathcal{L}(\Theta, \Upsilon)$ is as follows

$$\begin{aligned} \bar{\mathcal{L}}(\Theta, \Upsilon; \tilde{\Theta}) &= \frac{\beta}{N} \sum_{n=1}^N \log \left[\sum_{k=1}^K e^{v_k} q(\mathbf{x}_n|k; \Theta_k, \tilde{\Theta}_k) \right] + (1-\beta) \log \sum_{k=1}^K e^{v_k} \\ &\quad - \frac{1}{N} \sum_{k=1}^K \sum_{n \in \mathcal{C}_k} [\langle \mathbf{x}_n, \theta_k \rangle - \frac{1}{2} \langle \mathbf{x}_n \mathbf{x}_n^T, \Omega_k \rangle] \\ &\quad - \sum_{k=1}^K \frac{N_k}{N} \left[v_k - \frac{1}{2} \theta_k^T \Omega_k^{-1} \theta_k + \frac{1}{2} \log \det \Omega_k \right]. \end{aligned}$$

We require covariance matrices being positive definite; i.e. $\Sigma_k \succ \mathbf{0}$. Equivalently, we can have $\Omega_k \succ \mathbf{0}$. In practice, we can have $\Omega_k \succeq \alpha \mathbf{I}$, where $\alpha > 0$.

We can apply algorithm 1 to estimate Θ and Υ . In each iteration of algorithm 1, we need solve the following convex optimization problem

$$\begin{aligned} \min_{\Theta, \Upsilon} \quad & \bar{\mathcal{L}}(\Theta, \Upsilon; \tilde{\Theta}) \\ \text{s.t.} \quad & \Omega_k \succeq \alpha \mathbf{I}, \quad k = 1, 2, \dots, K \end{aligned} \quad (18)$$

where $\alpha > 0$ is given. To the best of our knowledge, there is no existing optimization method to solve this problem. In the following, we will present a block-coordinate descent algorithm for problem (18) that sequentially updates Υ , θ_1 , $\theta_2, \dots, \theta_K$ and $\Omega_1, \Omega_2, \dots, \Omega_K$. We will discuss the details of optimization algorithms for Υ in subsection IV-A, for θ_k in subsection IV-B, and for Ω_k in subsection IV-C; we will discuss selecting initial points, early stopping and β in subsection IV-D.

A. Updating v_k 's

We assume Θ_k 's are fixed: for $k = 1, 2, \dots, K$, $\theta_k = \tilde{\theta}_k$, and $\Omega_k = \tilde{\Omega}_k$. Since $q(\mathbf{x}_n|k; \tilde{\Theta}_k, \tilde{\Theta}_k) = p(\mathbf{x}_n|k; \tilde{\Theta}_k)$, we have the optimization problem for $\Upsilon = (v_1, v_2, \dots, v_K)$ as follows

$$\min_{\Upsilon} - \sum_{k=1}^K \frac{N_k}{N} v_k + \frac{\beta}{N} \sum_{n=1}^N \log \sum_{k=1}^K p(\mathbf{x}_n|k; \tilde{\Theta}_k) e^{v_k} + (1-\beta) \log \sum_{k=1}^K e^{v_k}.$$

B. Updating θ_k 's

In this subsection, we consider v_k 's and Ω_k 's are fixed; that is, for $k = 1, \dots, K$, $v_k = v_k^*$, and $\Omega_k = \tilde{\Omega}_k$. Furthermore, we consider all θ_k 's but θ_m are fixed; that is $\forall k < m, \theta_k = \theta_k^*$; and $\forall k > m, \theta_k = \tilde{\theta}_k$. We define $f(\theta_m)$ in this subsection as follows,

$$f(\theta_m) = \bar{\mathcal{L}}(\theta_1^*, \dots, \theta_{m-1}^*, \theta_m, \tilde{\theta}_{m+1}, \dots, \tilde{\theta}_K, \tilde{\Omega}_1, \dots, \tilde{\Omega}_K, \Upsilon^*). \quad (19)$$

In the following, we will present a trust region method to minimize $f(\theta_m)$. Due to the space limit, we refer to [13], [14] for details about trust region method.

Define

$$\phi_m = \theta_m - \theta_m^{(\tau)}, \quad (20)$$

and

$$q^{(\tau)}(m|\mathbf{x}_n) = \frac{e^{v_m^*} q(\mathbf{x}_n|m; \Theta_m^{(\tau)}, \tilde{\Theta}_m)}{\sum_{k=1}^K e^{v_k^*} q(\mathbf{x}_n|k; \Theta_k^{(\tau)}, \tilde{\Theta}_k)}$$

where $\theta_m^{(\tau)}$ is an estimation of θ_m at τ -th iteration, and

$$\Theta_k^{(\tau)} = \begin{cases} (\tilde{\Omega}_k, \theta_k^*), & k = 1, 2, \dots, m-1 \\ (\tilde{\Omega}_k, \theta_k^{(\tau)}), & k = m \\ (\tilde{\Omega}_k, \tilde{\theta}_k), & k = m+1, m+2, \dots, K. \end{cases}$$

We need solve the following trust region subproblem in the τ -th iteration,

$$\begin{aligned} \min_{\phi_m} \quad & \frac{\lambda}{2} \phi_m^T \phi_m + \mathbf{g}_m^T \phi_m \\ \text{s.t.} \quad & \phi_m^T \phi_m \leq \epsilon^{(\tau)}, \end{aligned} \quad (21)$$

where

$$\mathbf{g}_m = \frac{1}{N} \left\{ \sum_{n=1}^N \beta q^{(\tau)}(m|\mathbf{x}_n) (\mathbf{x}_n - \tilde{\mu}_m) - \sum_{n \in \mathcal{C}_m} (\mathbf{x}_n - \mu_m^{(\tau)}) \right\},$$

and λ is the 1-norm of the Hessian matrix $\frac{\beta}{N} \sum_{n=1}^N q^{(\tau)}(m|\mathbf{x}_n) [1 - q^{(\tau)}(m|\mathbf{x}_n)] (\mathbf{x}_n - \tilde{\mu}_m) (\mathbf{x}_n - \tilde{\mu}_m)^T + \frac{N_m}{N} \tilde{\Omega}_m^{-1}$. It can be shown that the optimal solution of problem (21) is

$$\phi_m^* = - \frac{\mathbf{g}_m}{\max(\sqrt{\frac{\mathbf{g}_m^T \mathbf{g}_m}{\epsilon^{(\tau)}}}, \lambda)}. \quad (22)$$

To summarize, we present the pseudo code of the trust region algorithm for finding θ_m^* in algorithm 2.

- 1: Specify $\eta_0, \eta_1, \eta_2, \eta_3, \eta_4$ {see [14]}
- 2: $\tau \leftarrow 0$ { τ : iteration count}
- 3: $\theta_m^{(\tau)} \leftarrow \theta_m$
- 4: **while** not converge **do**
- 5: $\phi_m^* = -\mathbf{g}_m / \max(\sqrt{\frac{\mathbf{g}_m^T \mathbf{g}_m}{\epsilon^{(\tau)}}}, \lambda)$
- 6: $r = \frac{f(\theta_m + \phi_m^*) - f(\theta_m)}{\mathbf{g}_m^T \phi_m^* + 0.5\lambda(\phi_m^*)^T \phi_m^*}$
- 7: if $r > \eta_0$, $\theta_m^{(\tau+1)} \leftarrow \theta_m^{(\tau)} + \phi_m^*$; otherwise, $\theta_m^{(\tau+1)} \leftarrow \theta_m^{(\tau)}$
- 8: if $r \geq \eta_2$, update $\epsilon^{(\tau+1)} \in [\epsilon^{(\tau)}, \eta_1 \epsilon^{(\tau)}]$; otherwise, update $\epsilon^{(\tau+1)} \in [\eta_3 \|\phi_m^*\|_2, \eta_4 \epsilon^{(\tau)}]$;
- 9: $\tau \leftarrow \tau + 1$
- 10: **end while**

Fig. 2. A trust-region algorithm for finding θ_m^*

C. Updating Ω_k 's

In this subsection, we consider v_k 's and θ_k 's are fixed; that is, for $k = 1, \dots, K$, $v_k = v_k^*$, and $\theta_k = \theta_k^*$. Furthermore, consider all Ω_k 's but Ω_m are fixed; that is $\forall k < m, \Omega_k = \Omega_k^*$; $\forall k > m, \Omega_k = \tilde{\Omega}_k$. We define $h(\Omega_m)$ in this subsection as follows,

$$h(\Omega_m) = \bar{L}(\theta_1^*, \dots, \theta_K^*, \Omega_1^*, \dots, \Omega_{m-1}^*, \Omega_m, \tilde{\Omega}_{m+1}, \dots, \tilde{\Omega}_K, \mathbf{Y}^*). \quad (23)$$

So the optimization problem for Ω_m is

$$\begin{aligned} \min_{\Omega_m} \quad & h(\Omega_m) \\ \text{s.t.} \quad & \Omega_m \succeq \alpha \mathbf{I}. \end{aligned}$$

To our best knowledge, there is no existing method to solve this optimization problem. We will present an iterative trust region method to solve this optimization problem in the remaining of this subsection.

Let $\Omega_m^{(\tau)}$ be an estimation of Ω_m at τ -th iteration. We define

$$\Phi_m = (\Omega_m^{(\tau)})^{-\frac{1}{2}} (\Omega_m - \Omega_m^{(\tau)}) (\Omega_m^{(\tau)})^{-\frac{1}{2}}. \quad (24)$$

Then by $\Omega_m \succeq \alpha \mathbf{I}$, we have

$$\Phi_m \succeq \alpha (\Omega_m^{(\tau)})^{-1} - \mathbf{I}.$$

Define

$$q^{(\tau)}(m|\mathbf{x}_n) = \frac{e^{v_m^* q(\mathbf{x}_n|m; \Theta_m^{(\tau)}, \tilde{\Theta}_m)}}}{\sum_{k=1}^K e^{v_k^* q(\mathbf{x}_n|k; \Theta_k^{(\tau)}, \tilde{\Theta}_k)}}$$

where

$$\Theta_k^{(\tau)} = \begin{cases} (\Omega_k^*, \theta_k^*), & k = 1, 2, \dots, m \\ (\Omega_k^{(\tau)}, \theta_k^*), & k = m \\ (\tilde{\Omega}_k, \theta_k^*), & k = m+1, m+2, \dots, K, \end{cases}$$

We need solve the following trust region subproblem in the τ -th iteration

$$\begin{aligned} \min_{\Phi_m} \quad & \frac{\lambda}{2} \langle \Phi_m, \Phi_m \rangle + \langle \mathbf{G}_m, \Phi_m \rangle \\ \text{s.t.} \quad & \Phi_m \succeq \alpha (\Omega_m^{(\tau)})^{-1} - \mathbf{I} \\ & \langle \Phi_m, \Phi_m \rangle \leq \epsilon^{(\tau)}, \end{aligned}$$

where

$$\lambda = \frac{N_m}{N} \{ (\theta_m^*)^T (\Omega_m^{(\tau)})^{-1} \theta_m^* + \frac{1}{2} \}, \quad (25)$$

and

$$\mathbf{G}_m = \frac{1}{2N} \sqrt{\Omega_m^{(\tau)}} \{ \beta \sum_{n=1}^N q^{(\tau)}(m|\mathbf{x}_n) [\tilde{\mathbf{M}}_m - \mathbf{x}_n \mathbf{x}_n^T] + \sum_{n \in \mathcal{C}_m} \mathbf{S}_n^{(\tau)} \} \sqrt{\Omega_m^{(\tau)}}, \quad (26)$$

where $\mathbf{S}_n^{(\tau)} = (\Omega_m^{(\tau)})^{-1} + (\Omega_m^{(\tau)})^{-1} \theta_m^* (\theta_m^*)^T (\Omega_m^{(\tau)})^{-1} - \mathbf{x}_n \mathbf{x}_n^T$. In appendix A, we show that we can approximately solve this problem by ¹

$$\Phi_m^* = \sqrt{\min\left(\frac{\epsilon^{(\tau)}}{\rho}, 1\right) \{ \mathbf{C}_-^{(\tau)} + \alpha (\Omega_m^{(\tau)})^{-1} - \mathbf{I} \}}, \quad (27)$$

where

$$\begin{aligned} \mathbf{C}^{(\tau)} &= \frac{1}{\lambda} \mathbf{G}_m + \alpha (\Omega_m^{(\tau)})^{-1} - \mathbf{I} \\ \rho &= \| \mathbf{C}_-^{(\tau)} + \alpha (\Omega_m^{(\tau)})^{-1} - \mathbf{I} \|_F^2, \end{aligned}$$

¹To be precise, Φ_m^* is a solution of the following feasibility problem: $\frac{\lambda}{2} \langle \Phi_m, \Phi_m \rangle + \langle \mathbf{G}_m, \Phi_m \rangle \leq 0$, $\Phi_m \succeq \alpha (\Omega_m^{(\tau)})^{-1} - \mathbf{I}$, and $\langle \Phi_m, \Phi_m \rangle \leq \epsilon$.

and $\mathbf{C}_-^{(\tau)}$ is the negative semidefinite parts of $\mathbf{C}^{(\tau)}$; let $\mathbf{C}^{(\tau)} = \sum_{i=1}^d \lambda_i \mathbf{e}_i \mathbf{e}_i^T$ be the eigen-decomposition of $\mathbf{C}^{(\tau)}$, then $\mathbf{C}_-^{(\tau)} = \sum_{i=1}^d \max(-\lambda_i, 0) \mathbf{e}_i \mathbf{e}_i^T$, (see [15]).

In summary, we present in Algorithm 3 the pseudo-code of the trust region algorithm for finding Ω_m^* .

- 1: Specify $\eta_0, \eta_1, \eta_2, \eta_3, \eta_4$ {see [14]}
- 2: $\tau \leftarrow 0$ {iteration count}
- 3: $\Omega_m^{(\tau)} \leftarrow \tilde{\Omega}_m$
- 4: **while** not converge **do**
- 5: calculate Φ_m^* as described in equation (27)
- 6: $\tilde{\Omega}_m = \Omega_m^{(\tau)} + (\Omega_m^{(\tau)})^{1/2} \Phi_m^* (\Omega_m^{(\tau)})^{1/2}$
- 7: $r = \frac{h(\Omega_m) - h(\tilde{\Omega}_m)}{0.5\lambda \langle \Phi_m^*, \Phi_m^* \rangle + \langle \mathbf{G}_m, \Phi_m^* \rangle}$
- 8: if $r > \eta_0$, $\Omega_m^{(\tau+1)} \leftarrow \tilde{\Omega}_m$; otherwise, $\Omega_m^{(\tau+1)} \leftarrow \Omega_m^{(\tau)}$
- 9: if $r \geq \eta_2$, update $\epsilon^{(\tau+1)} \in [\epsilon^{(\tau)}, \eta_1 \epsilon^{(\tau)}]$; otherwise, update $\epsilon^{(\tau+1)} \in [\eta_3 \|\Phi_m^*\|_2, \eta_4 \epsilon^{(\tau)}]$;
- 10: $\tau \leftarrow \tau + 1$
- 11: **end while**

Fig. 3. A trust-region algorithm for finding Ω_m^*

D. Initial Points, Early Stopping and β

We use the following as our initial points,

$$\Omega_k^{(0)} = \text{diag}(\hat{\Sigma}_k)^+ + \alpha \mathbf{I} \quad \theta_k^{(0)} = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \Omega_k^{(0)} \mathbf{x}_n \quad v_k^{(0)} = \log \frac{N_k}{N}.$$

where $\text{diag}(\hat{\Sigma}_k)$ is the diagonal sample covariance matrices for class k ; \square^+ is the pseudo-inverse of a matrix \square .

Early stopping has been shown to be an effective regularization method for boosting algorithms [16], [17]. Motivated by their works, we use early stopping in algorithm 1. We use cross-validation to determine β and the early stopping step in the iterative MM algorithm.

V. NUMERICAL EXPERIMENTS

We downloaded 10 datasets from UCI machine learning repository [18]. The statistics of the selected data sets are listed in table I. In our numerical experiments, we followed the following procedure: for all datasets, we randomly selected 40% of the data as training sets, and the remaining data are used as test sets; we standardized training data using the sample mean and standard deviation of training data; all methods are trained with training sets, and evaluated on the corresponding test sets unless otherwise specified. We repeated this procedure 20 times for each dataset.

We examined the influence of β on classification performance by varying β from 0 (discriminative) to 1 (generative). For Datasets diabetes, heart, ionosphere and sonar, we showed box plots of test set error rates in figure 4. For diabetes and heart, discriminative learning ($\beta = 0$) gives the best performance; for ionosphere and sonar, combinations of discriminative and generative learning gave best performance.

Bouchard and Triggs [6] used conjugate gradient (CG) method to optimize the same objective function as ours. We compared the new proposed method against CG for training QDA. For this purpose, we disabled early stopping in our method, and performed numerical experiments on an Intel

	d	N	K
Diabetes	8	768	2
Heart	13	270	2
Parkinsons	22	195	2
Ionosphere	33	351	2
Sonar	60	208	2
Thyroid	5	215	3
Vertebral Column	6	310	3
Wine	13	178	3
Splice	60	3175	3
Breast Tissue	9	106	6

TABLE I. DESCRIPTIONS OF 10 BENCHMARK DATASETS (K : NUMBER OF CLASSES, d : DIMENSION, N : NUMBER OF OBSERVATIONS)

Xeon (3.30GHz) computer with 16GB RAM. CG method reported local optimums within 3000 iterations for all datasets; our methods reported local optimums within 3000 iterations for all datasets except ionosphere. Computational speed and objective function values for both methods are reported in table V. In 9 of 10 datasets, our method is faster than CG method; in 7 of 10 datasets, our method terminates with equal or better local optimums.

To evaluate the classification performance of the new proposed method, we compared it with some widely used quadratic classifiers in the literature: naive Bayes, QDA², regularized discriminant analysis (RDA) [19]³, Bayesian quadratic discriminant analysis (BDA7) [3]⁴, L2-regularized logistic regression with quadratic terms (LR), SVM with polynomial kernel ($d = 2$) (QSVM). For RDA, we used leave-one-out approach to select optimal parameters on a $(0, 0.1, \dots, 1) \times (0, 0.1, \dots, 1)$ grid. For BDA7, we used the same settings as described in [3], and used leave-one-out approach to select optimal parameter from 42 parameter choices. For our method, we let $\alpha = 10^{-6}$ and $\beta = (0, 0.1, \dots, 1.0)$, and used 10 fold cross validation to select optimal β and early stopping iteration. For CG QDA [6], we used the optimal β selected from our methods, and trained QDA with $\frac{1}{2}$ -, 1- and 2-times of iterations determined from our methods; we recorded the smallest test set error rates. For LR and QSVM, we searched regularization coefficients in $[0, 10]$ by 10-fold cross validation.

Mean error rates for 10 benchmark datasets are reported in table V. For Naive Bayes and QDA, we reported “—” when the matlab function `classify` failed. The smallest mean error rates are emphasized with bold numbers. We also reported other mean error rates in bold, which are not statistically significantly different from the best ones; for this purpose, we used Wilcoxon signed-rank test with significant level 0.05 to compare error rates. As shown in table V, in 8 of the 10 datasets our method is either the best one, or not statistically significantly different from the best ones. The experiments shows that our method is competitive to other widely used quadratic classifiers in the literature.

VI. CONCLUSION

Quadratic discriminant analysis is a standard tool for classification problems. Classical QDA sacrifices classification accuracy for computation efficiency, which is not desired for classification problems. Motivated by the recent studies in

²We use matlab function `classify` for naive Bayes and QDA

³We use code from <http://www-stat.stanford.edu/~jhf>

⁴We modified code from <http://www.ee.washington.edu/research/guptalab>. $\Phi_m^* = \sqrt{\min(\frac{\epsilon(\tau)}{\rho}, 1)} \{C_-^{(\tau)} + \alpha(\Omega_m^{(\tau)})^{-1} - I\}$,

	CG QDA		Our Method	
	CPU Time (sec.)	Obj. Value	CPU Time (sec.)	Obj. Value
Diabetes	0.19 ± 0.08	5.70 ± 0.05	0.06 ± 0.01	5.70 ± 0.05
Heart	0.18 ± 0.01	8.32 ± 0.10	0.09 ± 0.02	8.32 ± 0.10
Parkinsons	10.46 ± 1.12	-2.47 ± 0.60	2.91 ± 1.35	-9.09 ± 1.70
Ionosphere	10.99 ± 2.18	4.88 ± 1.03	10.13 ± 2.26	5.52 ± 1.01
Sonar	23.24 ± 2.04	-17.29 ± 4.80	3.19 ± 4.24	-92.01 ± 14.37
Thyroid	1.13 ± 0.33	1.04 ± 0.22	0.40 ± 0.23	1.07 ± 0.23
Vertebral	4.01 ± 1.49	1.54 ± 0.13	0.39 ± 0.11	-0.41 ± 0.08
Wine	1.73 ± 0.96	4.67 ± 0.24	0.89 ± 0.22	4.87 ± 0.25
Splice	67.59 ± 31.53	39.38 ± 0.34	84.06 ± 80.61	39.35 ± 0.70
Breast Tissue	10.89 ± 0.67	-3.86 ± 0.62	6.44 ± 0.92	-9.53 ± 2.54

TABLE II. CPU TIME (SECONDS) AND FINAL OBJECTIVE FUNCTION VALUES OF TRAINING SETS ($\beta = 0.5$): CG V.S. OUR METHOD

hybrid generative/discriminative learning, we argue that, in order to obtain better classification accuracy, parameters of QDA can be estimated by maximizing a convex combination joint log-likelihood and conditional log-likelihood of given observations and their labels. For this purpose, we presented a MM optimization algorithm to estimate parameters for generative classifiers, of which conditional distributions are from the well known exponential family. Furthermore, we proposed a block-coordinate descent algorithm to sequentially update parameters of QDA in each iteration of the MM algorithm. For each update, we used a trust region method, of which each iteration has a simple closed form solution. Our numerical experiments show that our method is competitive with other well-known quadratic classification methods in the literature.

Our method can be easily adapted for linear discriminant analysis (LDA). Moreover in spirit, our MM algorithm and block-coordinate descent algorithm can be applied for any generative classifiers of which conditional distributions are from the exponential family.

Sparse parameters are desired when handling small sample problems, e.g. [20], [21]. Though early stopping provides regularization in our algorithm, we think explicitly adding sparse constraint via L_1 or L_0 regularization in the block-coordinate descent algorithm might provide additional robustness for QDA.

Recently, there is a growing interest in semi-supervised learning; that is to use unlabelled data in training classification methods [22]. Part of our future work is to include unlabelled data and missing data for estimating parameters for LDA and QDA.

APPENDIX

Let us first consider the following relaxed optimization problem,

$$\begin{aligned} \min_{\Phi_m} \quad & \frac{\lambda}{2} \langle \Phi_m, \Phi_m \rangle + \langle G_m, \Phi_m \rangle \\ \text{s.t.} \quad & \Phi_m \succeq \alpha(\Omega_m^{(\tau)})^{-1} - I \end{aligned}$$

Following [15], [23], [24], we can show that optimal solution for the relaxed problem is

$$\hat{\Phi}_m = C_-^{(\tau)} + \alpha(\Omega_m^{(\tau)})^{-1} - I$$

where $C^{(\tau)} = \frac{1}{\lambda} G_m + \alpha(\Omega_m^{(\tau)})^{-1} - I$. Define

$$\rho = \|C_-^{(\tau)} + \alpha(\Omega_m^{(\tau)})^{-1} - I\|_F^2,$$

If $\rho \leq \epsilon^{(\tau)}$, then we can let $\Phi_m^* = \hat{\Phi}_m$; if $\rho > \epsilon^{(\tau)}$, then we can let $\Phi_m^* = \sqrt{\frac{\epsilon^{(\tau)}}{\rho}} \hat{\Phi}_m$. Therefore, we can let

$$\Phi_m^* = \sqrt{\min(\frac{\epsilon^{(\tau)}}{\rho}, 1)} \{C_-^{(\tau)} + \alpha(\Omega_m^{(\tau)})^{-1} - I\},$$

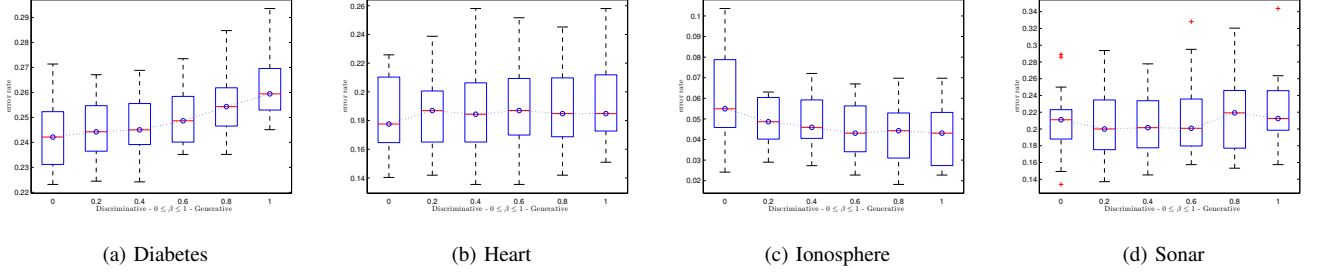


Fig. 4. Influence of coefficient of β on test set error rates: $\beta = 0$ is discriminative, and $\beta = 1$ is generative.

	Naive Bayes	QDA	RDA	BDA7	LR	QSVM	CG QDA	Our Method
Diabetes	25.61 \pm 1.44	26.68 \pm 1.88	23.86 \pm 1.37	25.75 \pm 1.21	23.56 \pm 1.38	23.34 \pm 1.34	24.39 \pm 1.11	24.09 \pm 1.06
Heart	17.44 \pm 2.68	22.68 \pm 3.39	16.56 \pm 1.93	17.87 \pm 2.31	18.52 \pm 2.74	21.64 \pm 3.38	20.38 \pm 3.66	18.49 \pm 2.53
Parkinsons	30.28 \pm 3.04	—	17.67 \pm 3.29	12.96 \pm 3.18	16.41 \pm 5.28	14.24 \pm 1.90	13.69 \pm 2.26	14.68 \pm 4.07
Ionosphere	—	—	6.75 \pm 1.60	7.45 \pm 2.07	9.02 \pm 1.92	14.40 \pm 2.73	5.64 \pm 1.89	4.90 \pm 1.29
Sonar	32.25 \pm 5.01	—	25.10 \pm 4.47	22.65 \pm 4.90	22.11 \pm 3.11	22.58 \pm 2.98	20.75 \pm 3.59	19.85 \pm 4.01
Thyroid	3.92 \pm 1.00	5.65 \pm 2.27	6.01 \pm 2.16	4.87 \pm 1.54	8.58 \pm 2.93	7.86 \pm 3.92	4.35 \pm 0.97	4.00 \pm 1.27
Vertebral Column	19.17 \pm 2.80	18.35 \pm 2.77	19.15 \pm 3.72	16.81 \pm 2.04	16.90 \pm 2.48	17.71 \pm 3.33	17.22 \pm 2.42	16.60 \pm 2.18
Wine	3.76 \pm 1.88	—	2.67 \pm 1.37	3.75 \pm 2.56	3.88 \pm 1.73	4.28 \pm 2.18	3.51 \pm 2.44	2.85 \pm 1.62
Splice	—	—	11.06 \pm 1.00	7.21 \pm 1.86	10.88 \pm 0.81	11.94 \pm 0.75	8.93 \pm 2.55	8.16 \pm 3.22
Breast Tissue	38.03 \pm 5.55	—	41.91 \pm 5.95	39.75 \pm 6.45	45.51 \pm 5.86	48.04 \pm 3.68	37.00 \pm 5.41	37.49 \pm 6.49

TABLE III. MEAN AND STANDARD DEVIATION OF TEST SET ERROR RATES FOR 10 DATASETS (SHOWN IN PERCENTAGE)

Since $\alpha(\Omega_m^{(\tau)})^{-1} \preceq \mathbf{I}$, and $\min(\frac{\epsilon(\tau)}{\rho}, 1) \leq 1$, we have $\sqrt{\min(\frac{\epsilon(\tau)}{\rho}, 1)[\alpha\Omega_m^{-1} - \mathbf{I}]} \succeq \alpha\Omega_m^{-1} - \mathbf{I}$. Therefore, $\Phi_m^* \succeq \alpha\Omega_m^{-1} - \mathbf{I}$.

Note that when $\rho > \epsilon(\tau)$, Φ_m^* is not an optimal solution. Since $\rho > \epsilon(\tau)$, we have

$$\frac{\lambda}{2} \langle \Phi_m^*, \Phi_m^* \rangle + \langle \mathbf{G}_m, \Phi_m^* \rangle \leq \sqrt{\frac{\epsilon(\tau)}{\rho}} \left(\frac{\lambda}{2} \langle \hat{\Phi}_m, \hat{\Phi}_m \rangle + \langle \mathbf{G}_m, \hat{\Phi}_m \rangle \right) \leq 0.$$

That is Φ_m^* does reduce the objective function value.

REFERENCES

- [1] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern classification*, 2nd ed. Wiley, Nov. 2000.
- [2] T. Hastie, R. Tibshirani, and J. H. Friedman, *The Elements of Statistical Learning*, corrected ed. Springer, Jul. 2003.
- [3] S. Srivastava, M. R. Gupta, and B. A. Frigiyik, "Bayesian quadratic discriminant analysis," *J. Mach. Learn. Res.*, vol. 8, pp. 1277–1305, Dec. 2007.
- [4] A. Y. Ng and M. I. Jordan, "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes," in *NIPS*, 2001, pp. 841–848.
- [5] P. Liang and M. I. Jordan, "An asymptotic analysis of generative, discriminative, and pseudolikelihood estimators," in *Proceedings of the 25th international conference on Machine learning*, ser. ICML '08, 2008, pp. 584–591.
- [6] G. Bouchard and B. Triggs, "The trade-off between generative and discriminative classifiers," in *Proceedings in Computational Statistics, 16th Symposium of IASC*. Physica-Verlag, 2004, pp. 721–728.
- [7] J. A. Lasserre, C. M. Bishop, and T. P. Minka, "Principled hybrids of generative and discriminative models," ser. CVPR '06, 2006, pp. 87–94.
- [8] A. McCallum, C. Pal, G. Druck, and X. Wang, "Multi-conditional learning: generative/discriminative training for clustering and classification," ser. AAAI'06, 2006, pp. 433–439.
- [9] M. J. Wainwright and M. I. Jordan, "Graphical models, exponential families, and variational inference," *Found. Trends Mach. Learn.*, vol. 1, no. 1-2, pp. 1–305, Jan. 2008.
- [10] D. R. Hunter and K. Lange, "A tutorial on mm algorithms," *Amer. Statist.*, vol. 58, no. 1, pp. 30–37, 2004.
- [11] M. Wainwright and M. Jordan, "Log-determinant relaxation for approximate inference in discrete markov random fields," *Trans. Sig. Proc.*, vol. 54, no. 6, pp. 2099–2109, Jun. 2006.
- [12] T. Jebara and A. Pentland, "On reversing jensen's inequality," in *In Advances in Neural Information Processing Systems 13*. MIT Press, 2000, p. 2000.
- [13] A. R. Conn, N. I. M. Gould, and P. L. Toint, *Trust-region methods*. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 2000.
- [14] Y. Yuan, "A review of trust region algorithms for optimization," in *ICIAM 99: Proceedings of the Fourth International Congress on Industrial & Applied Mathematics, Edinburgh*, 2000, pp. 271–282.
- [15] S. Boyd and L. Xiao, "Least-squares covariance matrix adjustment," *SIAM J. Matrix Anal. Appl.*, vol. 27, no. 2, pp. 532–546, June 2005.
- [16] T. Zhang and B. Yu, "Boosting with early stopping: Convergence and consistency," *Annals of Statistics*, vol. 33, no. 4, pp. 1538–1579, 2005.
- [17] S. Rosset, J. Zhu, and T. Hastie, "Boosting as a regularized path to a maximum margin classifier," *J. Mach. Learn. Res.*, vol. 5, pp. 941–973, Dec. 2004.
- [18] A. Frank and A. Asuncion, "UCI machine learning repository," 2010. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [19] J. H. Friedman, "Regularized discriminant analysis," *Journal of the American Statistical Association*, vol. 84, no. 405, pp. 165–175, 1989.
- [20] O. Banerjee, L. El Ghaoui, and A. d'Aspremont, "Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data," *J. Mach. Learn. Res.*, vol. 9, pp. 485–516, June 2008.
- [21] J. Friedman, T. Hastie, and R. Tibshirani, "Sparse inverse covariance estimation with the graphical lasso," *Biostatistics*, vol. 9, no. 3, pp. 432–441, Jul. 2008.
- [22] O. Chapelle, B. Schölkopf, and A. Zien, Eds., *Semi-Supervised Learning*. Cambridge, MA: MIT Press, 2006.
- [23] J. Malick, "A dual approach to semidefinite least-squares problems," *SIAM J. Matrix Anal. Appl.*, vol. 26, no. 1, pp. 272–284, Jan 2005.
- [24] H. Qi and D. Sun, "A quadratically convergent newton method for computing the nearest correlation matrix," *SIAM J. Matrix Anal. Appl.*, vol. 28, no. 2, pp. 360–385, Jun. 2006.