

# Linear Discriminant Analysis and Quadratic Discriminant Analysis



# Introduction

- Assume we have a dataset of instances  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  with sample size  $n$  and dimensionality  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ .
- The  $y_i$  's are the class labels.
- Linear Discriminant Analysis (LDA) and Quadratic discriminant Analysis (QDA) are two well-known supervised classification methods in statistical and probabilistic learning.

# Optimization for the Boundary of Classes

- First suppose the data is one dimensional,  $x \in \mathbb{R}$ .
- We assume that the two classes have normal (Gaussian) distribution which is the most common and default distribution in the real-world applications.
- An instance  $x \in \mathbb{R}$  belongs to one of these two classes:

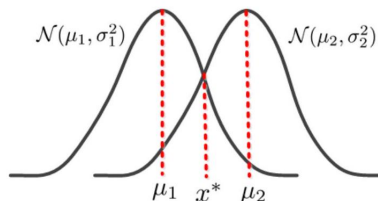
$$x \sim \begin{cases} \mathcal{N}(\mu_1, \sigma_1^2), & \text{if } x \in \mathcal{C}_1 \\ \mathcal{N}(\mu_2, \sigma_2^2), & \text{if } x \in \mathcal{C}_2 \end{cases}$$

where  $\mathcal{C}_1$  and  $\mathcal{C}_2$  denote the first and second class, respectively.

- For an instance  $x$ , we may have an error in estimation of the class it belongs to. At a point, which we denote by  $x^*$ , the probability of the two classes are equal; therefore, the point  $x^*$  is on the boundary of the two classes.

# Optimization for the Boundary of Classes

- Suppose  $\mu_1 < \mu_2$ .



- We can say  $\mu_1 < x^* < \mu_2$  as shown in Fig. 1
- Therefore, if  $x < x^*$  or  $x > x^*$  the instance  $x$  belongs to the first and second class, respectively.

# Optimization for the Boundary of Classes

- Estimating  $x < x^*$  or  $x > x^*$  to respectively belong to the second and first class is an error in estimation of the class. This probability of the error can be stated as:

$$\mathbb{P}(\text{error}) = \mathbb{P}(x > x^*, x \in \mathcal{C}_1) + \mathbb{P}(x < x^*, x \in \mathcal{C}_2).$$

- As we have  $\mathbb{P}(A, B) = \mathbb{P}(A | B)\mathbb{P}(B)$ , we can say:

$$\begin{aligned}\mathbb{P}(\text{error}) = & \mathbb{P}(x > x^* | x \in \mathcal{C}_1) \mathbb{P}(x \in \mathcal{C}_1) \\ & + \mathbb{P}(x < x^* | x \in \mathcal{C}_2) \mathbb{P}(x \in \mathcal{C}_2),\end{aligned}$$

which is what we want to minimize.

# Optimization for the Boundary of Classes

- According to the definition of CDF, we have:

$$\begin{aligned}\mathbb{P}(x < c, x \in \mathcal{C}_1) &= F_1(c), \\ \implies \mathbb{P}(x > x^*, x \in \mathcal{C}_1) &= 1 - F_1(x^*), \\ \mathbb{P}(x < x^*, x \in \mathcal{C}_2) &= F_2(x^*).\end{aligned}$$

- According to the definition of PDF, we have:

$$\begin{aligned}\mathbb{P}(x \in \mathcal{C}_1) &= f_1(x) = \pi_1, \\ \mathbb{P}(x \in \mathcal{C}_2) &= f_2(x) = \pi_2,\end{aligned}$$

where we denote the priors  $f_1(x)$  and  $f_2(x)$  by  $\pi_1$  and  $\pi_2$ , respectively.

- Hence, the problem becomes:

$$\underset{x^*}{\text{minimize}} (1 - F_1(x^*)) \pi_1 + F_2(x^*) \pi_2.$$

# Optimization for the Boundary of Classes

- We take derivative for the sake of minimization:

$$\begin{aligned}\frac{\partial \mathbb{P}(\text{error})}{\partial x^*} &= -f_1(x^*)\pi_1 + f_2(x^*)\pi_2 \stackrel{\text{set}}{=} 0, \\ \implies f_1(x^*)\pi_1 &= f_2(x^*)\pi_2.\end{aligned}$$

# Optimization for the Boundary of Classes

- Another way to obtain this expression is equating the posterior probabilities to have the equation of the boundary of classes:

$$\mathbb{P}(x \in \mathcal{C}_1 | X = x) \stackrel{\text{set}}{=} \mathbb{P}(x \in \mathcal{C}_2 | X = x).$$

- According to Bayes rule, the posterior is:

$$\begin{aligned}\mathbb{P}(x \in \mathcal{C}_1 | X = x) &= \frac{\mathbb{P}(X = x | x \in \mathcal{C}_1) \mathbb{P}(x \in \mathcal{C}_1)}{\mathbb{P}(X = x)} \\ &= \frac{f_1(x) \pi_1}{\sum_{k=1}^{|\mathcal{C}|} \mathbb{P}(X = x | x \in \mathcal{C}_k) \pi_k},\end{aligned}$$

- We can see that the boundary is also  $f_1(x^*) \pi_1 = f_2(x^*) \pi_2$ .



## $d$ -dimensional case

- Now let us think of data as multivariate data with dimensionality  $d$ . The PDF for multivariate Gaussian distribution,  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is:

$$f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left( -\frac{(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}{2} \right)$$

where  $\mathbf{x} \in \mathbb{R}^d$ ,  $\boldsymbol{\mu} \in \mathbb{R}^d$  is the mean,  $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$  is the covariance matrix, and  $|\cdot|$  is the determinant of matrix.

- The decision boundary becomes:

$$\begin{aligned} & \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_1|}} \exp \left( -\frac{(\mathbf{x} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)}{2} \right) \pi_1 \\ &= \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_2|}} \exp \left( -\frac{(\mathbf{x} - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)}{2} \right) \pi_2, \end{aligned}$$

# LDA for Binary Classification

- In Linear Discriminant Analysis (LDA), we assume that the two classes have equal covariance matrices:

$$\Sigma_1 = \Sigma_2 = \Sigma$$

- The decision boundary becomes:

$$\begin{aligned} & \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left( -\frac{(\mathbf{x} - \boldsymbol{\mu}_1)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)}{2} \right) \pi_1 \\ &= \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left( -\frac{(\mathbf{x} - \boldsymbol{\mu}_2)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)}{2} \right) \pi_2 \\ &\Rightarrow \exp \left( -\frac{(\mathbf{x} - \boldsymbol{\mu}_1)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)}{2} \right) \pi_1 = \exp \left( -\frac{(\mathbf{x} - \boldsymbol{\mu}_2)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)}{2} \right) \pi_2 \\ &\stackrel{(a)}{\Rightarrow} -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \ln(\pi_1) \\ &= -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \ln(\pi_2) \end{aligned}$$

# LDA for Binary Classification

■ Thus, we have:

$$\begin{aligned} & -\frac{1}{2}\mathbf{x}^\top \Sigma^{-1}\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}_1^\top \Sigma^{-1}\boldsymbol{\mu}_1 + \boldsymbol{\mu}_1^\top \Sigma^{-1}\mathbf{x} + \ln(\pi_1) \\ & = -\frac{1}{2}\mathbf{x}^\top \Sigma^{-1}\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}_2^\top \Sigma^{-1}\boldsymbol{\mu}_2 + \boldsymbol{\mu}_2^\top \Sigma^{-1}\mathbf{x} + \ln(\pi_2). \end{aligned}$$

■ Therefore, if we multiply the sides of equation by 2, we have:

$$\begin{aligned} & 2(\Sigma^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1))^\top \mathbf{x} \\ & + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + 2\ln\left(\frac{\pi_2}{\pi_1}\right) = 0 \end{aligned}$$

■ It is the equation of a line in the form of  $\mathbf{a}^\top \mathbf{x} + b = 0$ .

# LDA for Binary Classification

■ Let

$$\delta(\mathbf{x}) := 2 \left( \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) \right)^\top \mathbf{x} \\ + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + 2 \ln \left( \frac{\pi_2}{\pi_1} \right),$$

■ The class of an instance  $\mathbf{x}$  is estimated as:

$$\hat{\mathcal{C}}(\mathbf{x}) = \begin{cases} 1, & \text{if } \delta(\mathbf{x}) < 0 \\ 2, & \text{if } \delta(\mathbf{x}) > 0 \end{cases}$$

## QDA for Binary Classification

- we relax the assumption of equality of the covariance matrices:

$$\Sigma_1 \neq \Sigma_2$$

- By similar calculation as before, we can show that the decision boundary is:

$$\begin{aligned} & x^\top (\Sigma_1 - \Sigma_2)^{-1} x + 2 (\Sigma_2^{-1} \mu_2 - \Sigma_1^{-1} \mu_1)^\top x \\ & + \left( \mu_1^\top \Sigma_1^{-1} \mu_1 - \mu_2^\top \Sigma_2^{-1} \mu_2 \right) + \ln \left( \frac{|\Sigma_1|}{|\Sigma_2|} \right) \\ & + 2 \ln \left( \frac{\pi_2}{\pi_1} \right) = 0, \end{aligned}$$

- It is in the quadratic form  $x^\top \mathbf{A}x + \mathbf{b}^\top x + c = 0$ .

# LDA and QDA for Multi-class Classification

- Now we consider multiple classes, which can be more than two, indexed by  $k \in \{1, \dots, |\mathcal{C}|\}$ .
- the scaled posterior of the  $k$ -th class becomes:

$$\delta_k(\mathbf{x}) := -\frac{1}{2} \ln(|\Sigma_k|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \ln(\pi_k).$$

In QDA, the class of the instance  $\mathbf{x}$  is estimated as:

$$\hat{\mathcal{C}}(\mathbf{x}) = \arg \max_k \delta_k(\mathbf{x})$$

# LDA and QDA for Multi-class Classification

- In LDA, we assume that the covariance matrices of the  $k$  classes are equal:

$$\Sigma_1 = \dots = \Sigma_{|C|} = \Sigma$$

$$\begin{aligned}\delta_k(\mathbf{x}) &= -\frac{1}{2} \ln(|\Sigma|) \\ &\quad - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \ln(\pi_k) \\ &= -\frac{1}{2} \ln(|\Sigma|) - \frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^\top \Sigma^{-1} \boldsymbol{\mu}_k + \boldsymbol{\mu}_k^\top \Sigma^{-1} \mathbf{x} + \ln(\pi_k).\end{aligned}$$

## Estimation for LDA and QDA

- Usually, the prior of the  $k$ -th class is estimated according to the sample size of the  $k$ -th class:

$$\hat{\pi}_k = \frac{n_k}{n}$$

where  $n_k$  and  $n$  are the number of training instances in the  $k$ -th class and in total, respectively.

- The mean of the  $k$ -th class can be estimated using the Maximum Likelihood Estimation (MLE), or Method of Moments (MOM), for the mean of a Gaussian distribution:

$$\mathbb{R}^d \ni \hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^n x_i \mathbb{I}(\mathcal{C}(x_i) = k).$$



# Estimation for LDA and QDA

- In QDA, the covariance matrix of the  $k$ -th class is estimated using MLE:

$$\mathbb{R}^{d \times d} \ni \hat{\Sigma}_k = \frac{1}{n_k} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k) (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^\top \mathbb{I}(\mathcal{C}(\mathbf{x}_i) = k).$$

Or we can use the unbiased estimation of the covariance matrix:

$$\mathbb{R}^{d \times d} \ni \hat{\Sigma}_k = \frac{1}{n_k - 1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k) (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^\top \mathbb{I}(\mathcal{C}(\mathbf{x}_i) = k).$$

## Estimation for LDA and QDA

- In LDA, we assume that the covariance matrices of the classes are equal; therefore, we use the weighted average of the estimated covariance matrices as the common covariance matrix in LDA:

$$\mathbb{R}^{d \times d} \ni \hat{\Sigma} = \frac{\sum_{k=1}^{|C|} n_k \hat{\Sigma}_k}{\sum_{r=1}^{|C|} n_r} = \frac{\sum_{k=1}^{|C|} n_k \hat{\Sigma}_k}{n},$$

where the weights are the cardinality of the classes.

# LDA and QDA are Metric Learning!

- Assume that the covariance matrices are all equal (as we have in LDA) and they all are the identity matrix:

$$\Sigma_1 = \cdots = \Sigma_{|\mathcal{C}|} = \mathbf{I}$$

which means that all the classes are assumed to be spherically distributed in the  $d$  dimensional space.

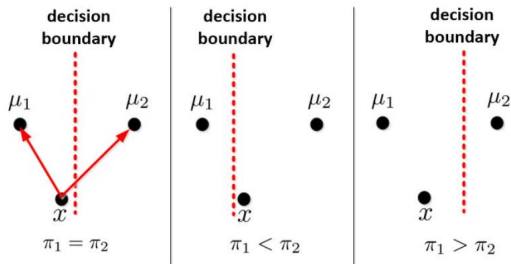


$$\delta_k(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top (\mathbf{x} - \boldsymbol{\mu}_k) + \ln(\pi_k),$$

$$\delta_k(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top (\mathbf{x} - \boldsymbol{\mu}_k) = -\frac{1}{2} d_k^2,$$

$$d_k = \|\mathbf{x} - \boldsymbol{\mu}_k\|_2 = \sqrt{(\mathbf{x} - \boldsymbol{\mu}_k)^\top (\mathbf{x} - \boldsymbol{\mu}_k)}.$$

Thus, the QDA or LDA reduce to simple Euclidean distance from the means of classes if the covariance matrices are all identity matrix and the priors are equal. Simple



- For general case when  $\Sigma \neq \mathbf{I}$ , we use Mahalanobis distance:

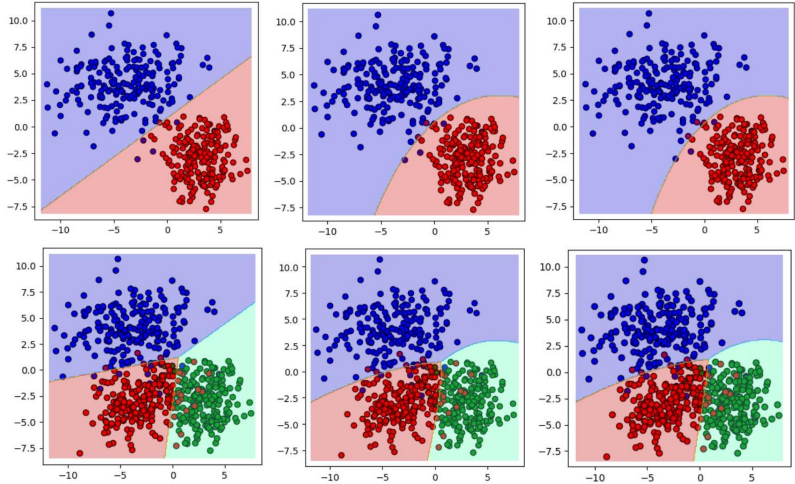
$$d_M^2(\mathbf{x}, \boldsymbol{\mu}) := \|\mathbf{x} - \boldsymbol{\mu}\|_M^2 = (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}),$$

.

## Relationship with Naive Bayes classifier

- Gaussian naive Bayes is equivalent to QDA where the covariance matrices are *di* agonal, i.e., the off-diagonal of the covariance matrices are ignored.
- Therefore, we can say that QDA is more powerful than Gaussian naive Bayes because Gaussian naive Bayes is a simplified version of QDA.
- Gaussian naive Bayes and QDA are equivalent for one dimensional data.

# Experiments



# Experiments

Experiments with equal class sample sizes: (1) LDA for two classes, (2) QDA for two classes, (3) Gaussian naive Bayes for two classes, (4) LDA for three classes, (5) QDA for three classes, (6) Gaussian naive Bayes for three classes.



## High-dimensional consideration

- For  $K$  classes case, the decision rule for LDA can be equivalently written as

$$\hat{Y} = \arg \max_k \left\{ \left( \mathbf{X} - \frac{\boldsymbol{\mu}_k}{2} \right)^T \boldsymbol{\beta}_k + \log \pi_k \right\},$$

where  $\boldsymbol{\beta}_k = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k$  for  $k = 1, \dots, K$ .

- To estimate  $\boldsymbol{\beta}_k$ , traditional approach estimates  $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}^{-1}$  separately.
- In high-dimensions, like LASSO, we can assume both  $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}^{-1}$  are sparse.
- This assumption can somehow decrease the estimation errors for the parameters and make  $\boldsymbol{\Sigma}^{-1}$  estimable.

# High-dimensional consideration

- A sparse estimation for  $\mu_k$  is easy to construct; sparse estimation for  $\Sigma^{-1}$  is more difficult and computationally expensive.
- What if we direct assume that  $\beta_k = \Sigma^{-1}\mu_k$  is sparse?
- Mai Q, Yang Y, Zou H. Multiclass sparse discriminant analysis[J]. Statistica Sinica, 2019, 29(1): 97-111.

## High-dimensional consideration

- Let  $\boldsymbol{\theta}_k^{\text{Bayes}} = \boldsymbol{\beta}_k - \boldsymbol{\beta}_1$  for  $k = 1, \dots, K$ . Then the Bayes decision rule can be written as

$$\hat{Y} = \arg \max_k \left\{ \left( \boldsymbol{\theta}_k^{\text{Bayes}} \right)^T \left( \mathbf{X} - \frac{\boldsymbol{\mu}_k}{2} \right) + \log \pi_k \right\}.$$

- We refer to the directions  $\boldsymbol{\theta}^{\text{Bayes}} = \left( \boldsymbol{\theta}_2^{\text{Bayes}}, \dots, \boldsymbol{\theta}_K^{\text{Bayes}} \right) \in \mathbb{R}^{p \times (K-1)}$  as the discriminant directions.

- On the population level, we have

$$\begin{aligned} & \left( \boldsymbol{\theta}_2^{\text{Bayes}}, \dots, \boldsymbol{\theta}_K^{\text{Bayes}} \right) \\ &= \arg \min_{\boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_K} \sum_{k=2}^K \left\{ \frac{1}{2} \boldsymbol{\theta}_k^T \boldsymbol{\Sigma} \boldsymbol{\theta}_k - (\boldsymbol{\mu}_k - \boldsymbol{\mu}_1)^T \boldsymbol{\theta}_k \right\}. \end{aligned}$$

# High-dimensional consideration

- In the classical low-dimension-large-sample-size setting, we can estimate  $(\theta_2^{\text{Bayes}}, \dots, \theta_K^{\text{Bayes}})$  via an empirical version

$$\begin{aligned} & (\hat{\theta}_2, \dots, \hat{\theta}_K) \\ &= \arg \min_{\theta_2, \dots, \theta_K} \sum_{k=2}^K \left\{ \frac{1}{2} \theta_k^T \hat{\Sigma} \theta_k - (\hat{\mu}_k - \hat{\mu}_1)^T \theta_k \right\}. \end{aligned}$$

- Still cannot handle high-dimensional data sets.

## High-dimensional consideration

- Write  $\boldsymbol{\theta}_{\cdot j} = (\theta_{2j}, \dots, \theta_{Kj})^T$  and define
$$\|\boldsymbol{\theta}_{\cdot j}\| = \left( \sum_{i=2}^K \theta_{ij}^2 \right)^{1/2}.$$
- For the high-dimensional case, consider the following penalized formulation for multiclass sparse discriminant analysis.

$$\begin{aligned} & \left( \hat{\boldsymbol{\theta}}_2, \dots, \hat{\boldsymbol{\theta}}_K \right) \\ &= \arg \min_{\boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_K} \sum_{k=2}^K \left\{ \frac{1}{2} \boldsymbol{\theta}_k^T \hat{\boldsymbol{\Sigma}} \boldsymbol{\theta}_k - (\hat{\boldsymbol{\mu}}_k - \hat{\boldsymbol{\mu}}_1)^T \boldsymbol{\theta}_k \right\} \\ & \quad + \lambda \sum_{j=1}^p \|\boldsymbol{\theta}_{\cdot j}\|, \end{aligned}$$

## Algorithm

- Given  $\{\boldsymbol{\theta}_{.j'}, j' \neq j\}$ , the solution of  $\boldsymbol{\theta}_{.j}$  to (10) is defined as

$$\arg \min_{\boldsymbol{\theta}_{.j}} \sum_{k=2}^K \frac{1}{2} \left( \theta_{kj} - \tilde{\theta}_{kj} \right)^2 + \frac{\lambda}{\hat{\sigma}_{jj}} \|\boldsymbol{\theta}_{.j}\|$$

where  $\tilde{\theta}_{k,j} = \frac{\hat{\delta}_j^k - \sum_{l \neq j} \hat{\sigma}_{lj} \theta_{kl}}{\hat{\sigma}_{jj}}.$

- Let  $\tilde{\boldsymbol{\theta}}_{.j} = \left( \tilde{\theta}_{2j}, \dots, \tilde{\theta}_{Kj} \right)^T$  and  $\|\tilde{\boldsymbol{\theta}}_{.j}\| = \left( \sum_{k=2}^K \tilde{\theta}_{kj}^2 \right)^{1/2}$ . The solution is given by

$$\hat{\boldsymbol{\theta}}_{.j} = \tilde{\boldsymbol{\theta}}_{.j} \left( 1 - \frac{\lambda}{\|\tilde{\boldsymbol{\theta}}_{.j}\|} \right)_+$$

# Algorithm

- write  $\hat{\delta}^k = \hat{\mu}_k - \hat{\mu}_1$ .
  1. Compute  $\hat{\Sigma}$  and  $\hat{\delta}^k, k = 1, 2, \dots, K$ ;
  2. Initialize  $\hat{\theta}_k^{(0)}$  and compute  $\tilde{\theta}_k^{(0)}$  accordingly;
  3. For  $m = 1, \dots$ , do the following loop until convergence: for  $j = 1, \dots, p$ , (a) compute

$$\hat{\theta}_{\cdot j}^{(m)} = \tilde{\theta}_{\cdot j}^{(m-1)} \left( 1 - \frac{\lambda}{\|\tilde{\theta}_{\cdot j}^{(m-1)}\|} \right)_+$$

(b) update

$$\tilde{\theta}_{kj} = \frac{\hat{\delta}_j^k - \sum_{l \neq j} \hat{\sigma}_{lj} \hat{\theta}_{kl}^{(m)}}{\hat{\sigma}_{jj}}.$$

4. Let  $\hat{\theta}_k$  be the solution at convergence. The output classifier is the usual linear discriminant classifier on  $(\mathbf{X}^T \hat{\theta}_2, \dots, \mathbf{X}^T \hat{\theta}_K)$ .

# Theorem

Under some technical conditions , there exists a generic constant  $M$  such that, if  $\lambda < \min \left\{ \frac{\theta_{\min}}{8\varphi}, M(1 - \kappa) \right\}$ , then with a probability greater than

$$1 - Cp d \exp \left( -Cn \frac{\epsilon^2}{Kd^2} \right) - CK \exp \left( -C \frac{n}{K^2} \right) - Cp(K-1) \exp \left( -Cn \frac{\epsilon^2}{K} \right)$$

we have that  $\hat{\mathcal{D}} = \mathcal{D}$ , and  $\left\| \hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^{\text{Bayes}} \right\|_{\infty} \leq 4\varphi\lambda$  for  $k = 2, \dots, K$ . 2. If we further assume conditions (C2)-(C3), we have that if  $\left\{ \frac{d^2 \log(pd)}{n} \right\}^{1/2} \ll \lambda \ll \theta_{\min}$ , then with probability tending to 1 , we have  $\hat{\mathcal{D}} = \mathcal{D}$ , and  $\left\| \hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^{\text{Bayes}} \right\|_{\infty} \leq 4\varphi\lambda$  for  $k = 2, \dots, K$ .



- What about high-dimensional QDA?
- It is more technically involved since we need to give a sparse estimation for  $\Sigma_k^{-1} - \Sigma_1^{-1}$ .
- Cai T T, Zhang L. A convex optimization approach to high-dimensional sparse quadratic discriminant analysis[J]. The Annals of Statistics, 2021, 49(3): 1537-1568.
- Jiang B, Wang X, Leng C. A direct approach for sparse quadratic discriminant analysis[J]. The Journal of Machine Learning Research, 2018, 19(1): 1098-1134.