

# Lecture 10: Bayes Classifier, LDA and QDA

Statistical Learning and Data Mining

Xingye Qiao

Department of Mathematical Sciences  
Binghamton University

E-mail: [qiao@math.binghamton.edu](mailto:qiao@math.binghamton.edu)

Read: ELSII Ch. 4.3, ISLR Chs. 4.4–4.6, and SLS Ch. 8.4

# Outline

- 1 Bayes Rule motivated LDA & QDA
- 2 Fisher's Linear Discriminant (optimization and geometry views)
- 3 Related Methods
- 4 Regularized & Sparse LDA

The next section would be .....

- 1 Bayes Rule motivated LDA & QDA
- 2 Fisher's Linear Discriminant (optimization and geometry views)
- 3 Related Methods
- 4 Regularized & Sparse LDA

## Setup

Assume  $K = 2$  for simplicity

$$\mathbf{X} \mid (Y = 1) \sim f_1(\mathbf{x}), \quad \mathbf{X} \mid (Y = 2) \sim f_2(\mathbf{x}).$$

Denote a random observation from this population (a mixture of two sub-populations) by  $(\mathbf{X}, Y)$ . We assume

$$P(Y = 1) = \pi_1, \quad P(Y = 2) = \pi_2, \quad \pi_1 + \pi_2 = 1$$

If the observed value of  $\mathbf{X}$  is  $\mathbf{x}$ , then Bayes theorem yields the posterior probability that  $\mathbf{X}$  was from Class 1 is

$$\eta(\mathbf{x}) := P(Y = 1 \mid \mathbf{X} = \mathbf{x}) = \frac{f_1(\mathbf{x})\pi_1}{f_1(\mathbf{x})\pi_1 + f_2(\mathbf{x})\pi_2}.$$

## 0-1 Loss and Bayes (decision) rule

- Use a decision theory framework, if we care about the 0-1 loss:  $\mathbb{1}\{Y \neq \delta(\mathbf{X})\}$ , we would hope to choose a decision function to minimize the risk associated with the 0-1 loss.

$$E[\mathbb{1}\{Y \neq \delta(\mathbf{X})\}] = P[Y \neq \delta(\mathbf{X})]$$

- Only need to choose the best decision  $\delta(\mathbf{x})$  for each given  $\mathbf{x}$ , which minimizes  $P[Y \neq \delta(\mathbf{x}) \mid \mathbf{X} = \mathbf{x}]$
- Rewrite  $P[Y \neq \delta(\mathbf{x}) \mid \mathbf{X} = \mathbf{x}]$  as

$$\mathbb{1}\{\delta(\mathbf{x}) = 1\}P(Y \neq 1 \mid \mathbf{X} = \mathbf{x}) + \mathbb{1}\{\delta(\mathbf{x}) = 2\}P(Y \neq 2 \mid \mathbf{X} = \mathbf{x})$$

- The blue parts are either 0 or 1. Hence we only need to choose  $\delta(\mathbf{x})$  to be  $j$  with a greater  $P(Y = j \mid \mathbf{X} = \mathbf{x})$

## Bayes Rule for Gaussian Data

Bayes Rule classifier assigns  $\mathbf{x}$  to the class label (1 or 2) which gives the higher posterior probability:

$$\phi_{\text{Bayes}}(\mathbf{x}) = \underset{k=1,2}{\operatorname{argmax}} P(Y = k \mid \mathbf{X} = \mathbf{x}).$$

Now assume Gaussian data as example:

$$\mathbf{X} \mid (Y = 1) \sim N_p(\mu_1, \Sigma), \quad \mathbf{X} \mid (Y = 2) \sim N_p(\mu_2, \Sigma), \quad \mu_1 \neq \mu_2.$$

Recall that

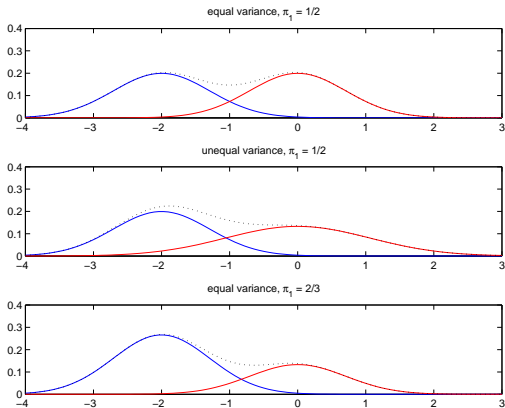
$$\eta(\mathbf{x}) := P(Y = 1 \mid \mathbf{X} = \mathbf{x}) = \frac{f_1(\mathbf{x})\pi_1}{f_1(\mathbf{x})\pi_1 + f_2(\mathbf{x})\pi_2}.$$

## Bayes Rule Classifier: Gaussian 1-D

As a special case, assume  $p = 1$  and

$$X|(Y = 1) \sim N_1(\mu_1, \sigma_1^2), \quad X|(Y = 2) \sim N_1(\mu_2, \sigma_2^2), \quad \mu_1 \neq \mu_2.$$

$$\text{Then } P(Y = i | X = x) \propto \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(x-\mu_i)^2}{2\sigma_i^2}\right) \times \pi_i$$



## Bayes Rule Classifier: Gaussian 1-D

Bayes rule classifies  $x$  into class 1 (the blue class)

- 1 case I ( $\sigma_1 = \sigma_2, \pi_1 = \pi_2$ ): if

$$|x - \mu_1| < |x - \mu_2|.$$

(Linear; bisector cutoff.)

- 2 case II ( $\sigma_1 < \sigma_2, \pi_1 = \pi_2$ ): if

$$\left(\frac{x - \mu_1}{\sigma_1}\right)^2 < \left(\frac{x - \mu_2}{\sigma_2}\right)^2 + \log(\sigma_2^2/\sigma_1).$$

(Quadratic; biased cutoff.)

- 3 case III ( $\sigma_1 = \sigma_2, \pi_1 > \pi_2$ ): if

$$\left(\frac{x - \mu_1}{\sigma_1}\right)^2 < \left(\frac{x - \mu_2}{\sigma_1}\right)^2 + 2 \log(\pi_1/\pi_2).$$

$$\left(\frac{-2x\mu_1 + \mu_1^2}{\sigma_1^2}\right) < \left(\frac{-2x\mu_2 + \mu_2^2}{\sigma_1^2}\right) + 2 \log(\pi_1/\pi_2).$$

(Linear; biased cutoff.)



## Multiclass Bayes Rule

In general, for  $k = 1, \dots, K > 2$ , if observed value of  $\mathbf{X} = \mathbf{x}$ , then

$$\eta_k(\mathbf{x}) := P(Y = k \mid \mathbf{X} = \mathbf{x}) = \frac{f_k(\mathbf{x})\pi_k}{f(\mathbf{x})},$$

where  $f_k(\mathbf{x})$  is density function of  $\mathbf{X} \mid (Y = k)$  and  $f(\mathbf{x})$  is the marginal density of  $\mathbf{X}$ . Bayes rule assigns  $\mathbf{x}$  to class label with highest  $\eta_k(\mathbf{x})$ :

$$\phi_{\text{Bayes}}(\mathbf{x}) = \operatorname{argmax}_{k=1,\dots,K} \frac{f_k(\mathbf{x})\pi_k}{f(\mathbf{x})} = \operatorname{argmax}_{k=1,\dots,K} f_k(\mathbf{x})\pi_k.$$

# (True) Quadratic Discriminant Analysis

## Gaussian Example:

$$\mathbf{X} \mid (Y = k) \sim N_p(\boldsymbol{\mu}_k, \Sigma_k), \quad k = 1, \dots, K,$$

$$P(Y = k) = \pi_k, \quad \sum_{k=1}^K \pi_k = 1,$$

$\mathbf{X} \sim f(\mathbf{x})$ , a mixture density for multivariate normals

Then

$$\phi_{\text{Bayes}}(\mathbf{x}) = \underset{k}{\operatorname{argmax}} \eta_k(\mathbf{x}) = \underset{k}{\operatorname{argmax}} \delta_k(\mathbf{x})$$

where

$$\delta_k(\mathbf{x}) = \left[ \log(\pi_k) - \frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)' \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right].$$

## (True) Quadratic Discriminant Analysis - Binary Case

For binary classification ( $k = 1, 2$ ), it is equivalent to calculating the sign of

$$\begin{aligned}\delta_1(\mathbf{x}) - \delta_2(\mathbf{x}) &= -\frac{1}{2} [(\mathbf{x} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) - (\mathbf{x} - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)] \\ &\quad - \frac{1}{2} \log(|\boldsymbol{\Sigma}_1|/|\boldsymbol{\Sigma}_2|) + \log(\pi_1/\pi_2) \\ &= -\frac{1}{2} \mathbf{x}' \boldsymbol{\nabla} \mathbf{x} + (\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2)' \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1' \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2' \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2 \\ &\quad - \frac{1}{2} \log(|\boldsymbol{\Sigma}_1|/|\boldsymbol{\Sigma}_2|) + \log(\pi_1/\pi_2)\end{aligned}$$

where  $\boldsymbol{\nabla} = \boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1}$

## Quadratic Discriminant Analysis

- In practice, we do not know the parameters  $(\mu_k, \Sigma_k, \pi_k)$ .
- Given  $n$  observations  $(\mathbf{x}_{ij}, i), (i = 1, \dots, K), (j = 1, \dots, n_k), n = \sum_{k=1}^K n_k$ , QDAs are obtained by substituting
  - 1  $\mu_k$  with  $\hat{\mu}_k = \bar{\mathbf{x}}_k$ ,
  - 2  $\Sigma_k$  with  $\hat{\Sigma}_k = \mathbf{S}_k$ ,
  - 3  $\hat{\pi}_k = n_k/n$ .

## Quadratic Discriminant Analysis

Quadratic Discriminant Analysis is essentially an estimate of the Bayes rule classifier for Gaussian data, which is

$$\phi(\mathbf{x}) = \operatorname{argmax}_{k=1,\dots,K} \hat{\delta}_k(\mathbf{x}), \text{ where}$$

$$\hat{\delta}_k(\mathbf{x}) = \left[ \log(n_k/n) - \frac{1}{2} \log(|\mathbf{S}_k|) - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_k)' \mathbf{S}_k^{-1} (\mathbf{x} - \bar{\mathbf{x}}_k) \right]$$

For binary classification, QDA is  $\phi(\mathbf{x}) = 1$  if

$$\begin{aligned} & -\frac{1}{2} \mathbf{x}' \hat{\nabla} \mathbf{x} + (\mathbf{S}_1^{-1} \bar{\mathbf{x}}_1 - \mathbf{S}_2^{-1} \bar{\mathbf{x}}_2)' \mathbf{x} - \frac{1}{2} \bar{\mathbf{x}}_1' \mathbf{S}_1^{-1} \bar{\mathbf{x}}_1 + \frac{1}{2} \bar{\mathbf{x}}_2' \mathbf{S}_2^{-1} \bar{\mathbf{x}}_2 \\ & - \frac{1}{2} \log(|\mathbf{S}_1|/|\mathbf{S}_2|) + \log(\hat{\pi}_1/\hat{\pi}_2) > 0 \text{ where } \hat{\nabla} = \mathbf{S}_1^{-1} - \mathbf{S}_2^{-1} \end{aligned}$$

or  $= -1$  otherwise (assuming  $\pm 1$  coding)

## Mahalanobis distance

In a special case where

$\Sigma_k \equiv \Sigma, \pi_k = 1/K$ , for all  $k$ :

The Bayes rule classifier boils down to comparing the quantity  $d_M^2(\mathbf{x}, \boldsymbol{\mu}_k) = (\mathbf{x} - \boldsymbol{\mu}_k)' \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)$ , which is called (squared) Mahalanobis distance between  $\mathbf{x}$  and  $\boldsymbol{\mu}_k$ .

- Mahalanobis distance  $d_M(\mathbf{x}, \boldsymbol{\mu})$  measures how much  $\mathbf{x}$  is away from the center of the distribution  $N_p(\boldsymbol{\mu}, \Sigma)$ .
- The set of points with the same Mahalanobis distance away from  $\boldsymbol{\mu}$  is an ellipsoid.
- Replacing  $\Sigma$  with its estimator  $\mathbf{S}$ , and  $\mathbf{x}$  with the sample mean  $\bar{\mathbf{x}}$ , the squared Mahalanobis distance is proportional to Hotelling's  $T^2$  statistic.

## Now consider an even more special case

Recall that the True QDA, which is the Bayes rule classifier for the Gaussian data, is

$$\phi_{\text{QDA}}(\mathbf{x}) = \underset{k}{\operatorname{argmax}} \eta_k(\mathbf{x}) = \underset{k}{\operatorname{argmax}} \delta_k(\mathbf{x})$$

where

$$\begin{aligned} \delta_k(\mathbf{x}) &= \left[ \log(\pi_k) - \frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)' \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right] \\ &= \left[ \log(\pi_k) - \frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} \mathbf{x}' \Sigma_k^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k' \Sigma_k^{-1} \boldsymbol{\mu}_k + \boldsymbol{\mu}_k' \Sigma_k^{-1} \mathbf{x} \right] \end{aligned}$$

Now we assume Equal covariance:  $\Sigma_k \equiv \Sigma$ , then

$-\frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} \mathbf{x}' \Sigma_k^{-1} \mathbf{x}$  will be same for all classes.

By removing  $-\frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} \mathbf{x}' \Sigma_k^{-1} \mathbf{x}$ ,  $\delta_k(\mathbf{x})$  is simplified to

$$\begin{aligned}\tilde{\delta}_k(\mathbf{x}) &= \log(\pi_k) - \frac{1}{2} \boldsymbol{\mu}_k' \Sigma^{-1} \boldsymbol{\mu}_k + \boldsymbol{\mu}_k' \Sigma^{-1} \mathbf{x} \\ &= b_{0k} + \mathbf{b}_k' \mathbf{x}, \text{ where} \\ b_{0k} &= \log(\pi_k) - \frac{1}{2} \boldsymbol{\mu}_k' \Sigma^{-1} \boldsymbol{\mu}_k, \quad \mathbf{b}_k = \Sigma^{-1} \boldsymbol{\mu}_k\end{aligned}$$

In this case, for binary classification ( $K = 2$ ):  $\phi(\mathbf{x}) = 1$  when

$$(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)' \Sigma^{-1} \left( \mathbf{x} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right) < \log(\pi_1/\pi_2),$$

that is

$$\mathbf{v}' \mathbf{x} - \mathbf{v}' \bar{\boldsymbol{\mu}} < \log(\pi_1/\pi_2), \quad \text{where } \mathbf{v} = \Sigma^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)$$

or  $= -1$  otherwise.



# Linear Discriminant Analysis

In practice, since we do not know the parameters  $(\mu_k, \Sigma, \pi_k)$ , substitute

- 1  $\mu_k$  with  $\hat{\mu}_k = \bar{\mathbf{x}}_k$ ,
- 2  $\Sigma$  with  $\hat{\Sigma} = \mathbf{S}_P$  (pooled sample covariance matrix), and
- 3  $\hat{\pi}_k = n_k/n$ .

Then we have the (sample) Linear Discriminant Analysis - estimated Bayes rule for Gaussian data with equal covariance assumption.

## Linear Discriminant Analysis

Binary case,

$$\phi(\mathbf{x}) = 1 \text{ if } \mathbf{v}'(\mathbf{x} - \frac{\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2}{2}) < \log(\frac{n_1}{n_2}), \quad \mathbf{v} = \mathbf{S}_P^{-1}(\bar{\mathbf{x}}_2 - \bar{\mathbf{x}}_1).$$

In general,  $\phi(\mathbf{x}) = \operatorname{argmax}_k (b_{0k} + \mathbf{b}'_k \mathbf{x})$ , where  $\mathbf{b}_k = \mathbf{S}_P^{-1} \bar{\mathbf{x}}_k$ .

## Quadratic Discriminant Analysis

Binary case,  $\phi(\mathbf{x}) = 1$  if

$$-\frac{1}{2} \mathbf{x}' \hat{\nabla} \mathbf{x} + (\mathbf{S}_1^{-1} \bar{\mathbf{x}}_1 - \mathbf{S}_2^{-1} \bar{\mathbf{x}}_2)' \mathbf{x} - \frac{1}{2} \bar{\mathbf{x}}_1' \mathbf{S}_1^{-1} \bar{\mathbf{x}}_1 + \frac{1}{2} \bar{\mathbf{x}}_2' \mathbf{S}_2^{-1} \bar{\mathbf{x}}_2 \\ - \frac{1}{2} \log(|\mathbf{S}_1|/|\mathbf{S}_2|) + \log(\hat{\pi}_1/\hat{\pi}_2) > 0 \text{ where } \hat{\nabla} = \mathbf{S}_1^{-1} - \mathbf{S}_2^{-1}$$

Generally  $\phi(\mathbf{x}) = \operatorname{argmax}_k \left[ \log(n_k/n) - \frac{1}{2} \log(|\mathbf{S}_k|) - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_k)' \mathbf{S}_k^{-1} (\mathbf{x} - \bar{\mathbf{x}}_k) \right]$

Bayes rule  $\xrightarrow{\text{estimate, Gaussian data}}$  QDA  $\xrightarrow{\text{equal covariance}}$  LDA

## The next section would be .....

- 1 Bayes Rule motivated LDA & QDA
- 2 Fisher's Linear Discriminant (optimization and geometry views)**
- 3 Related Methods
- 4 Regularized & Sparse LDA

## Fisher's LDA

- LDA is often referred to as R.A. Fisher's Linear Discriminant Analysis.
- His original work did not involve any distributional assumption, and developed LDA through a geometric understanding of PCA.
- The LDA direction  $\mathbf{v}_0 \in \mathbb{R}^p$  is a direction vector orthogonal to the separating hyperplane, and is found by maximizing the between-group variance while minimizing the within-group variance of the projected scores.

$$\mathbf{v}_0 = \operatorname{argmax}_{\mathbf{u} \in \mathbb{R}^p} \frac{(\mathbf{u}'\bar{\mathbf{x}}_1 - \mathbf{u}'\bar{\mathbf{x}}_2)^2}{\mathbf{u}'\mathbf{S}_P\mathbf{u}} = \frac{\mathbf{u}'(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'\mathbf{u}}{\mathbf{u}'\mathbf{S}_P\mathbf{u}}$$

$$\mathbf{v}_0 = \operatorname{argmax}_{\mathbf{u} \in \mathbb{R}^p} \frac{(\mathbf{u}'\bar{\mathbf{x}}_1 - \mathbf{u}'\bar{\mathbf{x}}_2)^2}{\mathbf{u}'\mathbf{S}_p\mathbf{u}} = \frac{\mathbf{u}'(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'\mathbf{u}}{\mathbf{u}'\mathbf{S}_p\mathbf{u}}$$

Because the objective is invariant with respect to the scaling of the vector  $\mathbf{u}$ , we can just assume that  $\mathbf{u}'\mathbf{S}_p\mathbf{u} = 1$ . Then we can transform the problem into the following constrained optimization problem,

$$\begin{aligned} \min \quad & \mathbf{u}'(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'\mathbf{u} \\ \text{s.t.} \quad & \mathbf{u}'\mathbf{S}_p\mathbf{u} = 1 \end{aligned}$$

This is a generalized eigenvalue problem.

From Theorem 2.5 in *Applied Multivariate Statistical Analysis* by Härdle and Simar, we have

$$\mathbf{v}_0 = \text{first eigenvector}\{S_P^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'\}$$

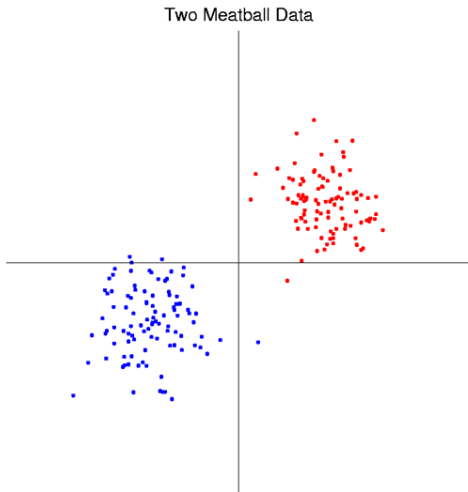
Observe the following identity

$$[S_P^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'] S_P^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) = \lambda S_P^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$$

where  $\lambda = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' S_P^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$  which happens to be the greatest eigenvalue of  $[S_P^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)']$ . Hence the solution of  $\mathbf{v}_0$  is actually

$$\mathbf{v}_0 = S_P^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$$

## Fisher's LDA–Geometric understanding

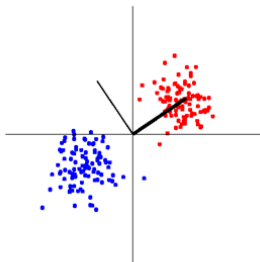


Two data clouds, each with  $\mathbf{S}_i = \mathbb{I}_2 = \mathbf{S}_P$ .

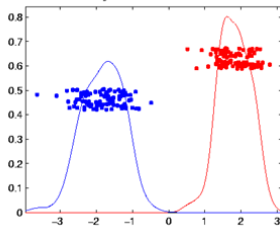


$$\mathbf{v}_0 \propto \mathbf{S}_P^{-1}(\bar{\mathbf{x}}_2 - \bar{\mathbf{x}}_1) = (\bar{\mathbf{x}}_2 - \bar{\mathbf{x}}_1). \text{ (direction of mean difference)}$$

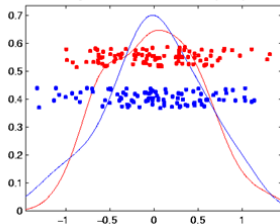
2 Meatballs, with Mean Diff. directions



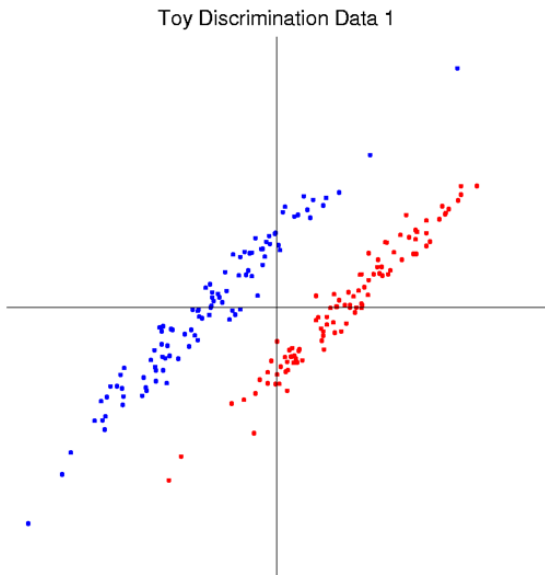
Projection onto MD



Projection onto MD perp.

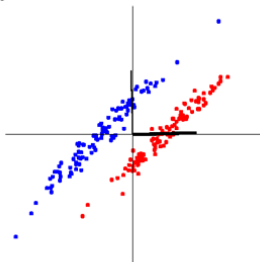


Slanted clouds. Assumed to have equal covariance.

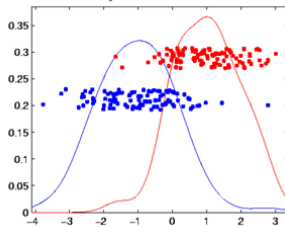


Mean difference direction not efficient, as  $\mathbf{S}_P \neq c\mathbb{I}_2$ .

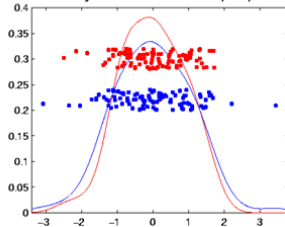
Toy Data 1, with Mean Diff. directions



Projection onto MD



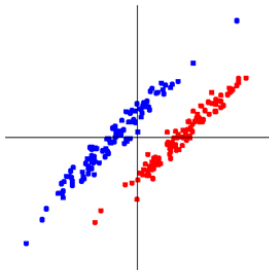
Projection onto MD perp.



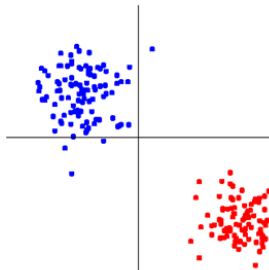
Individually transform subpopulations so that both are '*spherical*' about their means.

$$\tilde{\mathbf{x}}_{ij} = \mathbf{S}_P^{-1/2} \mathbf{x}_{ij}.$$

Toy Discrimination Data 1



Class by class sphered



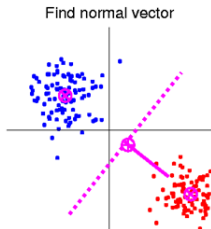
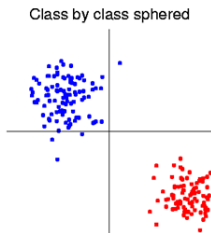
In **transformed space**, best separating hyperplane is the perpendicular bisector of line between means. Transformed mean diff. direction:

$$\tilde{\mathbf{v}} \propto \tilde{\mathbf{x}}_2 - \tilde{\mathbf{x}}_1 = \mathbf{S}_P^{-1/2}(\bar{\mathbf{x}}_2 - \bar{\mathbf{x}}_1).$$

Transformed center:

$$\tilde{\mathbf{b}}_0 = (\tilde{\mathbf{x}}_1 + \tilde{\mathbf{x}}_2)/2 = \mathbf{S}_P^{-1/2}(\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2)/2.$$

Transformed input  $\tilde{\mathbf{x}} = \mathbf{S}_P^{-1/2} \mathbf{x}$  is classified to 1 if  $\tilde{\mathbf{v}}'(\tilde{\mathbf{x}} - \tilde{\mathbf{b}}_0) < 0$



Original input  $\mathbf{x} = \mathbf{S}_P^{1/2} \tilde{\mathbf{x}}$  is classified to 1 if

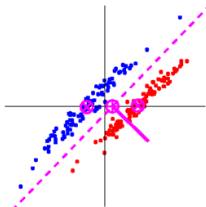
$$\tilde{\mathbf{v}}'(\tilde{\mathbf{x}} - \tilde{\mathbf{b}}_0) < 0$$

$$\Leftrightarrow \{\mathbf{S}_P^{-1/2}(\bar{\mathbf{x}}_2 - \bar{\mathbf{x}}_1)\}'(\mathbf{S}_P^{-1/2}\mathbf{x} - \mathbf{S}_P^{-1/2}(\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2)/2) < 0$$

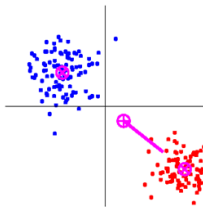
$$\Leftrightarrow (\bar{\mathbf{x}}_2 - \bar{\mathbf{x}}_1)'\mathbf{S}_P^{-1}(\mathbf{x} - (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2)/2) < 0$$

THIS IS THE LDA RULE!!!

Fisher Linear Discrimination

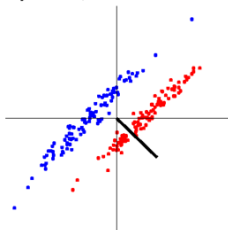


Find normal vector

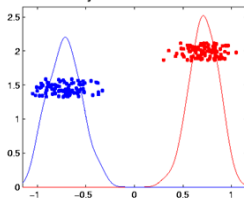


Leads to Fisher's LDA ( $\mathbf{v}_0 \propto \mathbf{S}_P^{-1}(\bar{\mathbf{x}}_2 - \bar{\mathbf{x}}_1)$ ) by actively using covariance structure.

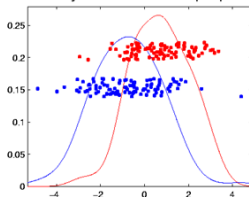
Toy Data 1, with FLD. directions



Projection onto FLD



Projection onto FLD perp.



## LDA vs QDA – Examples

In the next four sets of examples,

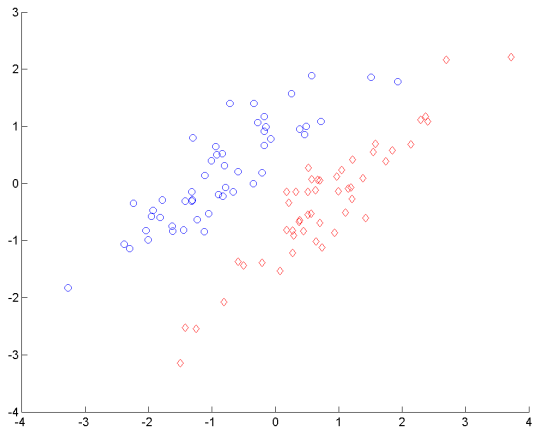
- Blue and red points represent observations from two classes.
- Blue line is the *separating hyperplane* given by computing the sample LDA, and is a line perpendicular to *LDA direction*  $\mathbf{b} = \mathbf{S}_p^{-1}(\bar{\mathbf{x}}_2 - \bar{\mathbf{x}}_1)$ , and is the set

$$\{\mathbf{x} \in \mathbb{R}^2 : \mathbf{b}'(\mathbf{x} - \frac{\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2}{2}) = \log(\frac{n_1}{n_2})\}.$$

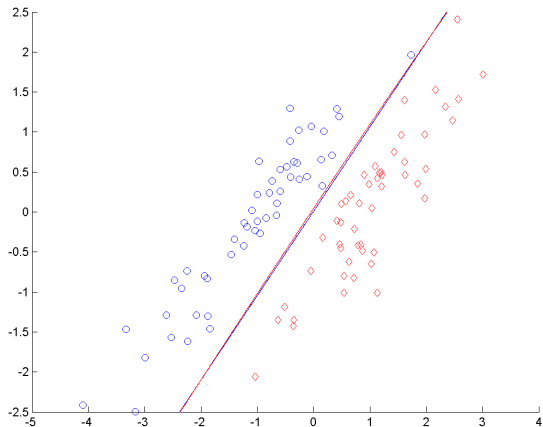
- Red curve represents classification boundary of sample QDA  $\{\mathbf{x} \in \mathbb{R}^2 : 0 = -\frac{1}{2}\mathbf{x}'\nabla\mathbf{x} + (\Sigma_1^{-1}\boldsymbol{\mu}_1 - \Sigma_2^{-1}\boldsymbol{\mu}_2)'\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}_1'\Sigma_1^{-1}\boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2'\Sigma_2^{-1}\boldsymbol{\mu}_2 - \frac{1}{2}\log(|\Sigma_1|/|\Sigma_2|) + \log(\pi_1/\pi_2)\}$ .



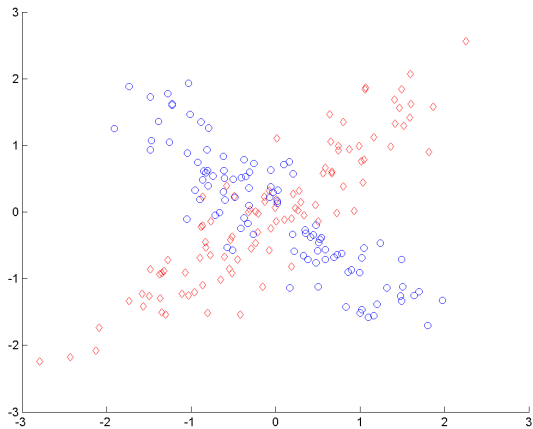
## LDA vs QDA – Ex.1



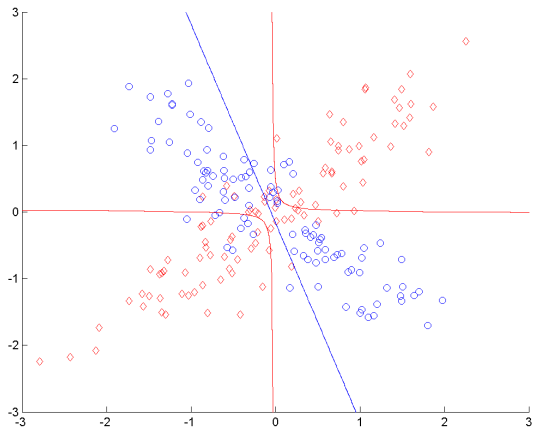
## LDA vs QDA – Ex.1



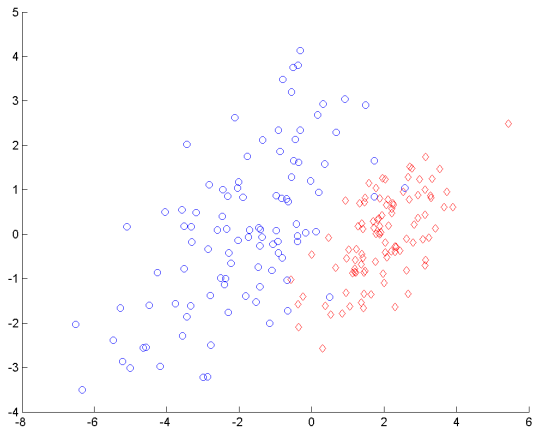
## LDA vs QDA – Ex.2



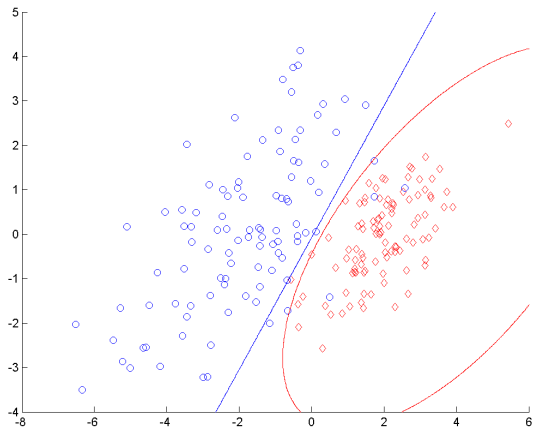
## LDA vs QDA – Ex.2



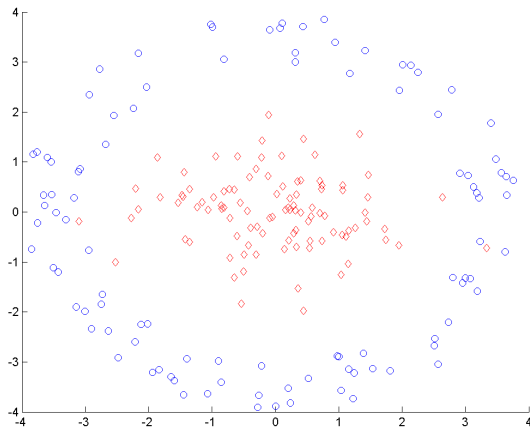
## LDA vs QDA – Ex.3



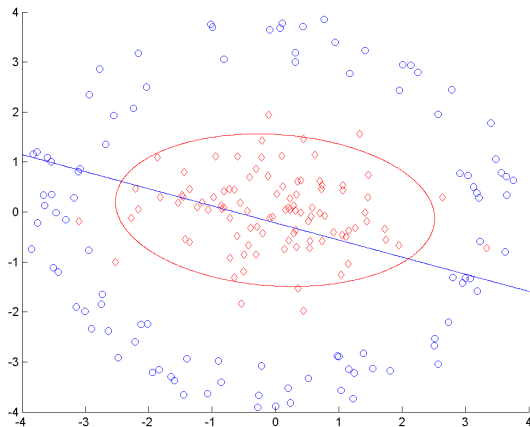
## LDA vs QDA – Ex.3



## LDA vs QDA – Ex.4



## LDA vs QDA – Ex.4





## The next section would be .....

- 1 Bayes Rule motivated LDA & QDA
- 2 Fisher's Linear Discriminant (optimization and geometry views)
- 3 Related Methods**
- 4 Regularized & Sparse LDA

## Nearest Centroid (Mean Difference) rule

A simplification of LDA where  $\mathbf{S}_P$  is replaced by  $\mathbb{I}$

- the binary nearest centroid classifier is

$$\phi(\mathbf{x}) = 1 \text{ if } \mathbf{b}'(\mathbf{x} - \frac{\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2}{2}) < 0, \quad \mathbf{b} = (\bar{\mathbf{x}}_2 - \bar{\mathbf{x}}_1),$$

sometimes called mean difference classifier;

## Naive Bayes classifier

A simplification of LDA where  $\mathbf{S}_P$  is replaced by  $\mathbf{D}_p = \text{Diag}(\mathbf{S}_P)$  the diagonal matrix consisting of diagonal elements of  $\mathbf{S}_P$

- the binary **Naive Bayes** classifier is

$$\phi(\mathbf{x}) = 1 \text{ if } \mathbf{b}'(\mathbf{x} - \frac{\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2}{2}) < 0, \quad \mathbf{b} = \mathbf{D}_p^{-1}(\bar{\mathbf{x}}_2 - \bar{\mathbf{x}}_1),$$

This is an (estimated) Bayes rule classifier which assumes that the (common) covariance matrix  $\Sigma$  is diagonal (which is quite a naive assumption).

LDA supress off-diagnol entries of covariance to 0 → naive Bayes  
force covariance to identity → nearest centroid (mean difference)  
rule

## Supervised dimension reduction

Recall that Fisher tried to maximize

$$\frac{(\mathbf{u}'\bar{\mathbf{x}}_1 - \mathbf{u}'\bar{\mathbf{x}}_2)^2}{\mathbf{u}'\mathbf{S}_P\mathbf{u}} = \frac{\mathbf{u}'(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'\mathbf{u}}{\mathbf{u}'\mathbf{S}_P\mathbf{u}}$$

In general for  $K \geq 3$ , defined  $\mathbf{B}$  as the covariance matrix of the class centroids and  $\mathbf{W}$  the within-class covariance.  $\mathbf{T} = \mathbf{W} + \mathbf{B}$  is the total covariance matrix of  $\mathbf{X}$ . Then we may iteratively solve

$$\max_{\mathbf{a}} \frac{\mathbf{a}^T \mathbf{B} \mathbf{a}}{\mathbf{a}^T \mathbf{W} \mathbf{a}}$$

Find the optimal  $\mathbf{a}_1$ , then the next direction  $\mathbf{a}_2$  orthogonal in  $\mathbf{W}$  to  $\mathbf{a}_1$  such that  $\mathbf{a}_2^T \mathbf{B} \mathbf{a}_2 / \mathbf{a}_2^T \mathbf{W} \mathbf{a}_2$  is maximized. . .

Similar to PCA, this is a way of dimension reduction. It finds directions that best separate the classes.

- These  $\mathbf{a}_k$ 's form a subspace  $\text{span}\{\mathbf{a}_k\}_{k=1}^K$ .
- It is equivalent to apply standard LDA to, instead of the original data, the projected data onto subspace  $\text{span}\{\mathbf{a}_k\}_{k=1}^K$ .
- Moreover, it can be shown that<sup>1</sup> this subspace is the same as

$$\mathbf{W}^{-1} \text{span}\{\hat{\boldsymbol{\mu}}_k - \hat{\boldsymbol{\mu}}_\ell\}_{1 \leq k \neq \ell \leq K}.$$

---

<sup>1</sup>Niu, Y.S., Hao, N. and Dong, B. (2018) "A New Reduced-Rank Linear Discriminant Analysis Method and Its Applications." *Statistica Sinica*, 28, 189-202.

## Reduced-rank LDA

- Therefore, the original LDA can be viewed as first projecting the data to the subspace  $\text{span}\{\mathbf{a}_k\}_{k=1}^K$  (as a form of dimension reduction), then applying standard LDA to the projected data.
- We may reduce the dimension more aggressively. The idea of the **reduced-rank LDA** is to find and project the data to a smaller subspace ( $\text{span}\{\mathbf{a}_k\}_{k=1}^L$ ) (with  $L < K$ ) before applying standard LDA.
- Treat  $L$  as a measure of model complexity and think about the trade-off between model fit and model complexity.

## Strengths & Weaknesses of LDA.

- Logistic regression is less sensitive to non-Gaussian data, since there is no Gaussian assumption. Beats LDA when non-Gaussian or the covariance is not equal
- LDA is subject to outliers.
- LDA exploits the full likelihood while logistic regression uses conditional likelihood.
- Logistic regression less efficient than LDA (needs a large sample to work well.) LDA can be quite flexible.
- Both have big problem when  $p \gg n$ .



## The next section would be .....

- 1 Bayes Rule motivated LDA & QDA
- 2 Fisher's Linear Discriminant (optimization and geometry views)
- 3 Related Methods
- 4 Regularized & Sparse LDA

## Regularized Discriminant Analysis

- Friedman (1989) proposed a compromise between LDA and QDA: replace  $\mathbf{S}_k$  by

$$\mathbf{S}_k(\alpha) = \alpha \mathbf{S}_k + (1 - \alpha) \mathbf{S}_P$$

- Similar modifications allow  $\mathbf{S}_P$  to be shrunk toward the scalar covariance,,: replace  $\mathbf{S}_P$  by

$$\mathbf{S}_P(\alpha) = \alpha \mathbf{S}_P + (1 - \alpha) \hat{\sigma}^2 \mathbb{I}$$

- Regularization with sparsity effect (next page)

## Sparse LDA

FSDA (Wu et al., 2008):

$$\operatorname{argmin}_{\beta} \beta^T \hat{\Sigma} \beta + \lambda \|\beta\|_1, \text{ subject to } (\hat{\mu}_2 - \hat{\mu}_1)^T \beta = 1$$

DSDA (Mai et al. 2012):

$$\beta \propto \operatorname{argmin}_{\beta} \sum_{i=1}^n (\tilde{y}_i - \beta^T x_i)^2 + \lambda \|\beta\|_1$$

PDA (Witten and Tibshirani 2011):

$$\operatorname{argmax}_{\beta} \beta^T \mathbf{B} \beta - \lambda \sum_{j=1}^p \hat{\sigma}_j |\beta_j|, \text{ subject to } \beta^t \mathbf{W} \beta \leq 1$$

SOS (Clemmensen et al., 2011), ROAD (Fan et al. 2012), and many many others.