

## Linear Discriminant Analysis

Winter 2023

Clayton Scott

## 1 Linear Classifiers

Linear Discriminant Analysis (LDA) is a classification method that produces a *linear classifier*. What is a linear classifier? When there are two classes, and the labels are -1 and 1, a linear classifier has the form

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b), \quad (1)$$

where  $\mathbf{w} \in \mathbb{R}^d$  is called the *weight vector* and  $b \in \mathbb{R}$  is called the *bias* or *offset*, and  $\text{sign}(t)$  is 1 if  $t \geq 0$  and -1 otherwise.<sup>1</sup> The set  $\{\mathbf{x} \mid \mathbf{w}^T \mathbf{x} + b = 0\}$  is the equation of a *hyperplane*.

For multiclass classification with labels  $1, 2, \dots, K$ , a linear classifier has the form

$$f(\mathbf{x}) = \arg \max_{k=1, \dots, K} \mathbf{w}_k^T \mathbf{x} + b_k, \quad (2)$$

where  $\mathbf{w}_1, \dots, \mathbf{w}_K \in \mathbb{R}^d$  and  $b_1, \dots, b_K \in \mathbb{R}$ . LDA is our first example of a linear classifier but we will see others. As an exercise one can easily show that the multiclass definition of linear classifier when  $K = 2$  agrees with the definition in (1).

## 2 The LDA Assumption

Unless otherwise noted, let the class labels be  $1, 2, \dots, K$ . LDA is a plug-in classifier based on the formula

$$f^*(\mathbf{x}) = \arg \max_k \pi_k g_k(\mathbf{x})$$

for a Bayes classifier. In LDA, it is assumed that  $\mathbf{X} \mid Y = k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma})$  for  $k \in \{1, \dots, K\}$ . In other words,

$$\begin{aligned} g_k(\mathbf{x}) &= \phi(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) \\ &:= (2\pi)^{-\frac{d}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right). \end{aligned}$$

Thus, the class conditional distributions are multivariate Gaussians with a common covariance matrix  $\boldsymbol{\Sigma}$ .

## 3 The LDA Classifier

Suppose we have training data  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ . LDA is the classifier obtained by estimating  $\pi_k, \boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}$  by

$$\begin{aligned} \hat{\pi}_k &= \frac{n_k}{n}, & n_k &= |\{i : y_i = k\}| \\ \hat{\boldsymbol{\mu}}_k &= \frac{1}{n_k} \sum_{i: y_i = k} \mathbf{x}_i \\ \hat{\boldsymbol{\Sigma}} &= \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_{y_i})(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_{y_i})^T. \end{aligned}$$

<sup>1</sup>A better name for a linear classifier would be an *affine classifier*, since  $\mathbf{w}^T \mathbf{x} + b$  is an affine function of  $\mathbf{x}$ , but the term “linear” is very entrenched.

and plugging these estimates into the Bayes classifier formula. In other words, the LDA classifier is

$$\hat{f}(\mathbf{x}) = \arg \max_{1 \leq k \leq K} \hat{\pi}_k \cdot \phi(\mathbf{x}; \hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\Sigma}})$$

The estimates  $\hat{\pi}_k$ ,  $\hat{\boldsymbol{\mu}}_k$ , and  $\hat{\boldsymbol{\Sigma}}$  can be shown to be maximum likelihood estimates. The estimate  $\hat{\boldsymbol{\Sigma}}$  is known as *the pooled sample covariance matrix*. This estimate can be written

$$\hat{\boldsymbol{\Sigma}} = \sum_{k=1}^K \frac{n_k}{n} \hat{\boldsymbol{\Sigma}}_k,$$

where  $\hat{\boldsymbol{\Sigma}}_k$  is the sample covariance matrix based on data from class  $k$ .

Since the natural logarithm is a strictly increasing function, the LDA classifier can also be expressed

$$\hat{f}(\mathbf{x}) = \arg \max_{1 \leq k \leq K} \log \hat{\pi}_k + \log \phi(\mathbf{x}; \hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\Sigma}}).$$

### 3.1 Binary Case

Let's show that  $\hat{f}$  is a linear classifier when  $K = 2$ . For this subsection only, let the labels be  $-1$  and  $1$ . Then the LDA classifier is

$$\begin{aligned} \hat{f}(\mathbf{x}) &= \text{sign}(\log \hat{\pi}_1 + \log \phi(\mathbf{x}; \hat{\boldsymbol{\mu}}_1, \hat{\boldsymbol{\Sigma}}) - \log \hat{\pi}_{-1} - \log \phi(\mathbf{x}; \hat{\boldsymbol{\mu}}_{-1}, \hat{\boldsymbol{\Sigma}})) \\ &= \text{sign} \left( \log \hat{\pi}_1 - \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\hat{\boldsymbol{\Sigma}}| - \frac{1}{2} (\mathbf{x} - \hat{\boldsymbol{\mu}}_1)^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_1) \right. \\ &\quad \left. - \log \hat{\pi}_{-1} + \frac{d}{2} \log 2\pi + \frac{1}{2} \log |\hat{\boldsymbol{\Sigma}}| + \frac{1}{2} (\mathbf{x} - \hat{\boldsymbol{\mu}}_{-1})^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_{-1}) \right) \\ &= \text{sign} \left( \log \frac{\hat{\pi}_1}{\hat{\pi}_{-1}} + \frac{1}{2} \left[ (\mathbf{x} - \hat{\boldsymbol{\mu}}_{-1})^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_{-1}) - (\mathbf{x} - \hat{\boldsymbol{\mu}}_1)^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_1) \right] \right) \\ &= \text{sign} \left( \log \frac{\hat{\pi}_1}{\hat{\pi}_{-1}} + \frac{1}{2} \left[ \mathbf{x}^T \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{x} - 2\mathbf{x}^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_{-1} + \hat{\boldsymbol{\mu}}_{-1}^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_{-1} - \mathbf{x}^T \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{x} + 2\mathbf{x}^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_1^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_1 \right] \right) \\ &= \text{sign} \left( \log \frac{\hat{\pi}_1}{\hat{\pi}_{-1}} + \mathbf{x}^T \hat{\boldsymbol{\Sigma}}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_{-1}) + \frac{1}{2} (\hat{\boldsymbol{\mu}}_{-1}^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_{-1} - \hat{\boldsymbol{\mu}}_1^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_1) \right) \\ &= \text{sign}(\mathbf{w}^T \mathbf{x} + b) \end{aligned}$$

where  $\mathbf{w} = \hat{\boldsymbol{\Sigma}}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_{-1})$  and  $b = \log(\hat{\pi}_1/\hat{\pi}_{-1}) + \frac{1}{2}(\hat{\boldsymbol{\mu}}_{-1}^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_{-1} - \hat{\boldsymbol{\mu}}_1^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_1)$ .

### 3.2 Multiclass Case

In the multiclass case, we can also show that LDA is linear. Because constant terms (that is, terms not depending on  $k$ ) do not affect which  $k$  achieves the maximum, we have

$$\begin{aligned}
\hat{f}(\mathbf{x}) &= \arg \max_k \log \hat{\pi}_k + \log \phi(\mathbf{x}; \hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\Sigma}}) \\
&= \arg \max_k \log \hat{\pi}_k - \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\hat{\boldsymbol{\Sigma}}| - \frac{1}{2} (\mathbf{x} - \hat{\boldsymbol{\mu}}_k)^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_k) \\
&= \arg \max_k \log \hat{\pi}_k - \frac{1}{2} (\mathbf{x} - \hat{\boldsymbol{\mu}}_k)^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_k) \\
&= \arg \max_k \log \hat{\pi}_k - \frac{1}{2} \left[ \mathbf{x}^T \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{x} - 2 \mathbf{x}^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_k + \hat{\boldsymbol{\mu}}_k^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_k \right] \\
&= \arg \max_k \log \hat{\pi}_k + \mathbf{x}^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_k - \frac{1}{2} \hat{\boldsymbol{\mu}}_k^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_k \\
&= \arg \max_k \mathbf{w}_k^T \mathbf{x} + b_k
\end{aligned}$$

where  $\mathbf{w}_k = \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_k$  and  $b_k = \log \hat{\pi}_k - \frac{1}{2} \hat{\boldsymbol{\mu}}_k^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_k$ .

## 4 Mahalanobis Distance

Let  $\mathbf{A}$  be a  $d \times d$  positive definite matrix, and define

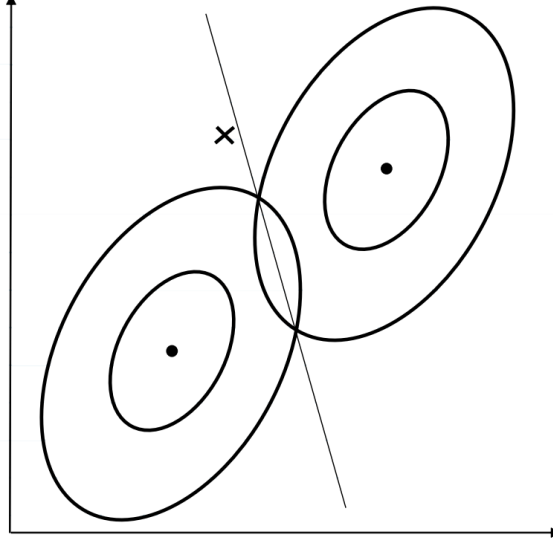
$$D_{\mathbf{A}}(\mathbf{x}, \mathbf{x}') := \sqrt{(\mathbf{x} - \mathbf{x}')^T \mathbf{A} (\mathbf{x} - \mathbf{x}')},$$

referred to as the *Mahalanobis distance* between  $\mathbf{x}$  and  $\mathbf{x}'$ . Then the LDA classifier can be expressed

$$\begin{aligned}
\hat{f}(\mathbf{x}) &= \arg \max_k \log \hat{\pi}_k - \frac{1}{2} D_{\hat{\boldsymbol{\Sigma}}^{-1}}^2(\mathbf{x}, \boldsymbol{\mu}_k) \\
&= \arg \min_k -\log \hat{\pi}_k + \frac{1}{2} D_{\hat{\boldsymbol{\Sigma}}^{-1}}^2(\mathbf{x}, \boldsymbol{\mu}_k).
\end{aligned}$$

Therefore, if we ignore the  $\log \hat{\pi}_k$  terms, LDA can be viewed as assigning a test point  $\mathbf{x}$  to the class whose mean  $\boldsymbol{\mu}_k$  is closest to  $\mathbf{x}$  according to the Mahalanobis distance associated to  $\hat{\boldsymbol{\Sigma}}^{-1}$ . See Figure 1.

Figure 1: This figure assumes  $\hat{\pi}_1 = \hat{\pi}_2$ . Points on one side of the line are closer to one class's mean in the Mahalanobis distance, while points on the other side are closer to the other mean. The ellipses are contours of the estimated class-conditional distributions.



## 5 Exercises

1. (★) Prove that if  $f$  is a linear classifier according to (2) with  $K = 2$ , then  $f$  is linear according to (1) with an appropriate reassignment of labels. Express  $\mathbf{w}$  and  $b$  in terms of  $\mathbf{w}_1, \mathbf{w}_2, b_1$ , and  $b_2$ . Verify your result in the special case of LDA by referencing the derivations above.
2. (★) Is LDA generative or discriminative? Parametric or nonparametric?
3. (★★) Show that if  $n < d$ , then  $\hat{\Sigma}$  is not invertible, meaning LDA cannot be applied.
4. (☆☆) Show that the Mahalanobis distance  $D_{\mathbf{A}}$  is the metric induced by a certain inner product. This inner product reduces to the dot product when  $\mathbf{A}$  is the identity matrix. Then show that with respect to this inner product (with  $\mathbf{A} = \Sigma^{-1}$ ), the LDA ( $K = 2$ ) decision boundary is orthogonal to the perpendicular bisector of the line segment joining the class sample means.
5. (★★) *Quadratic discriminant analysis* (QDA) is like LDA except each class is allowed to have its own covariance matrix  $\Sigma_k$ , which is estimated by the sample covariance matrix for that class, that is,

$$\hat{\Sigma}_k = \frac{1}{n_k} \sum_{i: y_i = k} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^T.$$

In the two class case with labels  $-1$  and  $1$ , determine  $\mathbf{A}$ ,  $\mathbf{b}$ , and  $c$  such that the QDA classifier can be expressed

$$\hat{f}(\mathbf{x}) = \text{sign} \left( \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c \right).$$

6. (☆☆☆) Consider binary classification where the class labels are  $-1$  and  $1$ . Let the training data be  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ . Let  $\boldsymbol{\mu}_{-1}$  and  $\boldsymbol{\mu}_1$  denote the sample means of the two classes. Define the “between class scatter matrix”

$$S_b = (\boldsymbol{\mu}_{-1} - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_{-1} - \boldsymbol{\mu}_1)^T,$$

and the “within class scatter matrix”

$$S_w = \sum_{k \in \{\pm 1\}} \sum_{i: y_i = k} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T.$$

Fisher’s Linear Discriminant (FLD) seeks a linear classifier with normal vector  $\mathbf{w}$  obtained by solving

$$\max_{\mathbf{w}} \frac{\mathbf{w}^T S_b \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}}.$$

The offset  $b$  of the linear classifier is determined by some other criterion after  $\mathbf{w}$  has been determined. FLD has nothing to say about  $b$ .

- (a) Determine the optimal  $\mathbf{w}$  according to the FLD criterion, and argue that it coincides with the LDA normal vector.
- (b) Justify FLD by interpreting the numerator and denominator of the objective function. Be as quantitative as possible. *Hint:* Think of FLD as a supervised version of PCA, but with just one principal component.