

Linear discriminant analysis

Linear discriminant analysis (LDA) is a common “plug-in” method for classification which operates by estimating $\pi_k f_{X|Y}(\mathbf{x}|k)$ for each class $k = 0, \dots, K - 1$ and then simply plugging these into the formula for the Bayes classifier in order to make a decision. In LDA we make the (strong) assumption that class conditional pdfs are given by the multivariate normal distribution, but with differing means. Mathematically, this corresponds to the assumption that

$$f_{X|Y}(\mathbf{x}|k) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_k)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_k)}$$

for $k = 0, \dots, K - 1$. Note that under this assumption, each class has a distinct mean $\boldsymbol{\mu}_k$, but all classes share the same covariance matrix Σ .

In LDA, we assume that $\boldsymbol{\mu}_0, \dots, \boldsymbol{\mu}_{K-1}$ and Σ , as well as the prior probabilities π_0, \dots, π_{K-1} are all unknown, but can be estimated from the data. In particular, we can use the estimates

$$\begin{aligned} \hat{\pi}_k &= \frac{|\{i : y_i = k\}|}{n} \\ \hat{\boldsymbol{\mu}}_k &= \frac{1}{|\{i : y_i = k\}|} \sum_{i:y_k=k} \mathbf{x}_i \\ \hat{\Sigma} &= \frac{1}{n} \sum_{k=0}^{K-1} \sum_{i:y_i=k} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^T. \end{aligned}$$

The LDA classifier is then defined by

$$\hat{h}(\mathbf{x}) = \arg \max_k \hat{\pi}_k \cdot \frac{1}{(2\pi)^{d/2} |\hat{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\hat{\boldsymbol{\mu}}_k)^T \hat{\Sigma}^{-1}(\mathbf{x}-\hat{\boldsymbol{\mu}}_k)}.$$

Since the log is a monotonic transformation (meaning that if $x > y$ then $\log(x) > \log(y)$), we can equivalently state the classifier as

$$\begin{aligned}\widehat{h}(\mathbf{x}) &= \arg \max_k \log(\widehat{\pi}_k) + \log\left(\frac{1}{(2\pi)^{d/2}|\widehat{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\widehat{\boldsymbol{\mu}}_k)^T\widehat{\Sigma}^{-1}(\mathbf{x}-\widehat{\boldsymbol{\mu}}_k)}\right) \\ &= \arg \max_k \log(\widehat{\pi}_k) - \frac{1}{2}(\mathbf{x}-\widehat{\boldsymbol{\mu}}_k)^T\widehat{\Sigma}^{-1}(\mathbf{x}-\widehat{\boldsymbol{\mu}}_k) \\ &= \arg \min_k \frac{1}{2}(\mathbf{x}-\widehat{\boldsymbol{\mu}}_k)^T\widehat{\Sigma}^{-1}(\mathbf{x}-\widehat{\boldsymbol{\mu}}_k) - \log(\widehat{\pi}_k)\end{aligned}$$

where the second equality above follows from the fact that

$$\log\left(\frac{1}{(2\pi)^{d/2}|\widehat{\Sigma}|^{1/2}}\right)$$

is constant across all k and so does not affect which k maximizes the expression.

It is enlightening to consider what happens in the special case of $K = 2$ (i.e., binary classification). In this case, LDA results in a classifier such that $\widehat{h}(\mathbf{x}) = 1$ when

$$(\mathbf{x}-\widehat{\boldsymbol{\mu}}_0)^T\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\widehat{\boldsymbol{\mu}}_0) - 2\log\widehat{\pi}_0 \geq (\mathbf{x}-\widehat{\boldsymbol{\mu}}_1)^T\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\widehat{\boldsymbol{\mu}}_1) - 2\log\widehat{\pi}_1.$$

We can rewrite this as

$$(\mathbf{x}-\widehat{\boldsymbol{\mu}}_0)^T\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\widehat{\boldsymbol{\mu}}_0) - (\mathbf{x}-\widehat{\boldsymbol{\mu}}_1)^T\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\widehat{\boldsymbol{\mu}}_1) + 2\log\frac{\widehat{\pi}_1}{\widehat{\pi}_0} \geq 0.$$

Using the fact that $\boldsymbol{\Sigma}$ is symmetric, which implies that we have

$(\boldsymbol{\Sigma}^{-1})^T = \boldsymbol{\Sigma}^{-1}$, we can simplify this rule to

$$\begin{aligned}
0 &\leq (\mathbf{x} - \hat{\boldsymbol{\mu}}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_0) - (\mathbf{x} - \hat{\boldsymbol{\mu}}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_1) + 2 \log \frac{\hat{\pi}_1}{\hat{\pi}_0} \\
&= \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2 \hat{\boldsymbol{\mu}}_0^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \hat{\boldsymbol{\mu}}_0^T \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_0 \\
&\quad - \left(\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2 \hat{\boldsymbol{\mu}}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \hat{\boldsymbol{\mu}}_1^T \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_1 \right) + 2 \log \frac{\hat{\pi}_1}{\hat{\pi}_0} \\
&= 2(\hat{\boldsymbol{\mu}}_1^T - \hat{\boldsymbol{\mu}}_0^T) \boldsymbol{\Sigma}^{-1} \mathbf{x} + \hat{\boldsymbol{\mu}}_0^T \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_0 - \hat{\boldsymbol{\mu}}_1^T \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_1 + 2 \log \frac{\hat{\pi}_1}{\hat{\pi}_0} \\
&= (\boldsymbol{\Sigma}^{-1} (\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_0))^T \mathbf{x} + \frac{1}{2} \hat{\boldsymbol{\mu}}_0^T \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_0 - \frac{1}{2} \hat{\boldsymbol{\mu}}_1^T \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_1 + \log \frac{\hat{\pi}_1}{\hat{\pi}_0}.
\end{aligned}$$

Thus, if

$$\mathbf{w} = \boldsymbol{\Sigma}^{-1} (\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_0)$$

and

$$b = \frac{1}{2} \hat{\boldsymbol{\mu}}_0^T \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_0 - \frac{1}{2} \hat{\boldsymbol{\mu}}_1^T \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_1 + \log \frac{\hat{\pi}_1}{\hat{\pi}_0},$$

we can re-write this as

$$\mathbf{w}^T \mathbf{x} + b \geq 0.$$

This is the expression of a simple linear classifier, and thus LDA will always result in a linear classifier.