

# MACHINE LEARNING

## LINEAR AND QUADRATIC DISCRIMINANT ANALYSIS

Sebastian Engelke

MASTER IN BUSINESS ANALYTICS



**UNIVERSITÉ  
DE GENÈVE**

## Linear decision boundaries: hyperplanes

- ▶ Suppose we have only two classes, denoted by  $\mathcal{G} = \{0, 1\}$ .
- ▶ Linear classification methods produce **linear decision boundaries**, called **hyperplanes**.
- ▶ In a  $p$ -dimensional space, a hyperplane is a flat **affine subspace of dimension  $p - 1$** , and it can be written, for parameters  $b_0, b_1, \dots, b_p \in \mathbb{R}$ , as the set

$$x \in \mathbb{R}^p : \quad b_0 + b_1 x_1 + \dots + b_p x_p = 0.$$

- ▶ Suppose a point  $x \in \mathbb{R}^p$  is not on the hyperplane, then

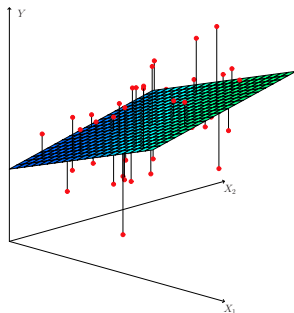
$$b_0 + b_1 x_1 + \dots + b_p x_p > 0,$$

tells us that it lies on the **one side** of the hyperplane, or with “ $<$ ” on the other.

### Example:

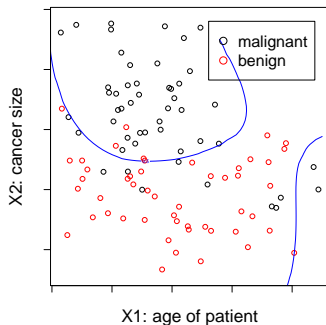
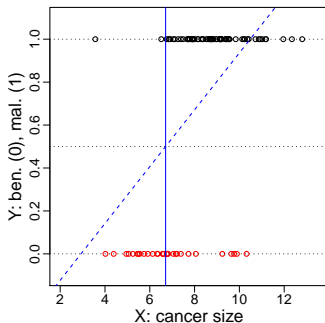
- ▶ For  $p = 3$ , the hyperplane is

$$(x_1, x_2, x_3) \in \mathbb{R}^3 : \quad b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 = 0.$$



## Linear methods for classification

- ▶ Example: direct linear regression on  $\mathcal{G} = \{0, 1\}$ ; usually not a good idea!
- ▶ Many other methods will produce linear decision boundaries.
- ▶ Decision boundaries of linear methods are linear only in the **feature space**.
- ▶ **Augmenting the feature** space by basis functions might lead to non-linear decision boundaries in the original space of  $(X_1, \dots, X_p)$ .
- ▶ Right: the model is linear in  $X_1, X_2, X_1^2, X_2^2, X_1^3, X_2^3, X_1X_2, X_1X_2^2, X_1^2X_2, X_1^2X_2^2$ .
- ▶ For more than two classes, linear classification methods produce **piecewise linear decision boundaries**.



# Linear discriminant analysis

- **Setup:** Classify  $Y \in \mathcal{G} = \{1, \dots, q\}$  given predictors  $X = (X_1, \dots, X_p)$
- **Idea:** model the **class-conditional densities**  $g_j$  of  $X \mid Y = j$ , and the marginals  $\pi_j = \Pr(Y = j)$ ,  $j = 1, \dots, q$ . Deduce  $\Pr(Y \mid X)$  from **Bayes' formula**:

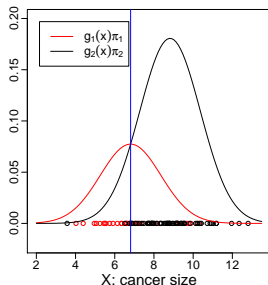
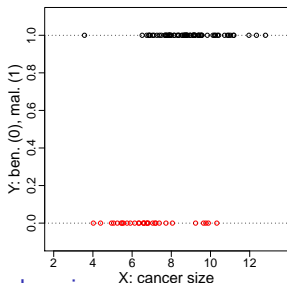
$$\Pr(Y = j \mid X = x) = \frac{\Pr(X = x \mid Y = j)\Pr(Y = j)}{\Pr(X = x)} \propto g_j(x)\pi_j.$$

- Linear discriminant analysis (LDA):

- (1) model  $g_j(x)$  by a **multivariate normal distribution**,

$$g_j(x) = \frac{1}{(2\pi)^{p/2}(\det \Sigma_j)^{1/2}} \exp \left\{ -\frac{1}{2}(x - \mu_j)^T \Sigma_j^{-1}(x - \mu_j) \right\};$$

- (2) assume the  $g_j$  have the **same covariance matrix**  $\Sigma_j = \Sigma$  for all  $j$ .



► **Model parameters:**

- the class/prior probabilities  $\pi_1, \dots, \pi_q \in (0, 1)$ ;
- the mean vectors  $\mu_1, \dots, \mu_q \in \mathbb{R}^P$ ;
- the covariance matrix  $\Sigma \in \mathbb{R}^{P \times P}$ .

- **Estimation:** From a training sample  $(y_1, x_1), \dots, (y_n, x_n)$  we estimate the parameters of LDA by the usual **maximum likelihood estimators**

$$\begin{aligned}\hat{\pi}_j &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{y_i = j\} \\ \hat{\mu}_j &= \frac{\sum_{i=1}^n x_i \mathbf{1}\{y_i = j\}}{\sum_{i=1}^n \mathbf{1}\{y_i = j\}} \\ \hat{\Sigma} &= \frac{\sum_{j=1}^q \sum_{i=1}^n (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T \mathbf{1}\{y_i = j\}}{n - q}\end{aligned}$$

- This is very simple and **computationally fast**!
- Note: Different to linear regression, we don't use the numerical value of the  $y_i$ 's but only the class!

## LDA: prediction

- **Goal:** for a new input value  $x_0 \in \mathbb{R}^p$  predict its unobserved class  $y_0$ .
- **Method:** compute the **posterior class probabilities**  $\widehat{\Pr}(Y = j \mid X = x_0)$  and use the **Bayes classifier** resulting from these estimates (plug-in estimate), i.e.,

$$\hat{y}_0 = \arg \max_{j=1,\dots,q} \widehat{\Pr}(Y = j \mid \mathbf{X} = x_0) = \arg \max_{j=1,\dots,q} \log \widehat{\Pr}(Y = j \mid X = x_0).$$

- Note that

$$\begin{aligned} \log \widehat{\Pr}(Y = j \mid X = x_0) &= \log \widehat{g}_j(x_0) + \log \widehat{\pi}_j + \text{constant} \\ &= -\frac{1}{2}(x_0 - \widehat{\mu}_j)^T \widehat{\Sigma}^{-1}(x_0 - \widehat{\mu}_j) + \log \widehat{\pi}_j + \text{constant} \\ &= x_0^T \widehat{\Sigma}^{-1} \widehat{\mu}_j - \frac{1}{2} \widehat{\mu}_j^T \widehat{\Sigma}^{-1} \widehat{\mu}_j + \log \widehat{\pi}_j + \text{constant}, \end{aligned} \tag{1}$$

where the constant includes terms that do not depend on  $j$ .

- **Geometric interpretation:** eq. (1) implies that LDA classifies  $x_0$  in terms of the **Mahalanobis distance**,  $d_{\widehat{\Sigma}^{-1}}(x, y) = x^T \widehat{\Sigma}^{-1} y$ , of  $x_0$  from the centers  $\widehat{\mu}_1, \dots, \widehat{\mu}_q$  (with a correction for the prior class probabilities).

## LDA: decision boundaries

- ▶ The **decision boundary** that separates two classes  $j, m \in \{1, \dots, q\}$  are all  $x \in \mathbb{R}^p$  with

$$\Pr(Y = j \mid X = x) = \Pr(Y = m \mid X = x),$$

which is equivalent to

$$\log \Pr(Y = j \mid X = x) - \log \Pr(Y = m \mid X = x) = 0.$$

Replacing the probabilities by their expressions in terms of the Gaussian densities and the prior class probabilities we get the decision boundary

$$\log \frac{\pi_j}{\pi_m} - \frac{1}{2}(\mu_j + \mu_m)^T \Sigma^{-1}(\mu_j - \mu_m) + x^T \Sigma^{-1}(\mu_j - \mu_m) = 0,$$

a linear equation in  $x$ .

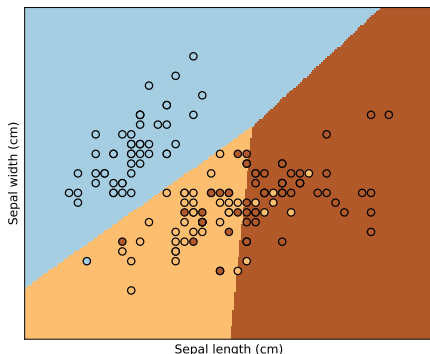
- ▶ Hence the decision boundaries separating any two classes  $j$  and  $m$  are hyperplanes in  $\mathbb{R}^p$  of dimension  $p - 1$ .
- ▶ Thus the decision boundary that separates the  $q$  classes is **piecewise linear**.

## Iris data set

For 150 measured plants;  $p = 4$  predictors as response the three Iris species

$\mathcal{G} = \{\text{setosa}, \text{versicolor}, \text{virginica}\}$  the plant belongs. We only use two predictors here.

```
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
iris = sklearn.datasets.load_iris()
X, y = iris.data[:, :2], iris.target # we only keep the first two features: 'Sepal length', 'Sepal width'
lda = LinearDiscriminantAnalysis()
lda.fit(X, y)
x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5
h = .02 # step size in the mesh
xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
Z = lda.predict(np.c_[xx.ravel(), yy.ravel()]).reshape(xx.shape)
```





## Extensions of LDA

- There are many possibilities to extend LDA. The Bayes' construction

$$\Pr(Y = j \mid X = x) = \frac{\Pr(X = x \mid Y = j)\Pr(Y = j)}{\Pr(X = x)} = \frac{g_j(x)\pi_j}{\sum_{j=1}^q g_j(x)\pi_j}.$$

can be used with other densities  $g_j$ .

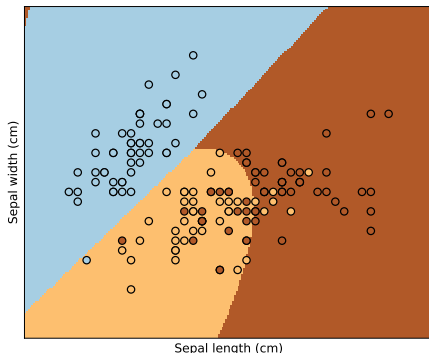
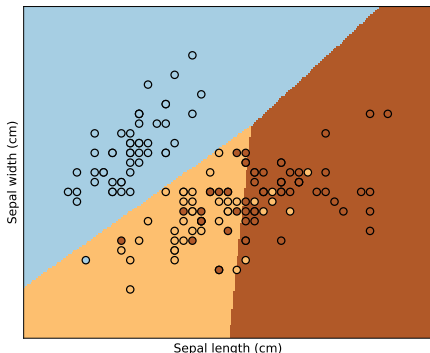
- **Quadratic discriminant analysis** (QDA): model  $g_j$  by multivariate normals with different covariance matrices  $\Sigma_j$ .

## Iris data set: QDA

For 150 measured plants;  $p = 4$  predictors as response the three Iris species

$\mathcal{G} = \{\text{setosa}, \text{versicolor}, \text{virginica}\}$  the plant belongs. We only use two predictors here.

```
from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
iris = sklearn.datasets.load_iris()
X, y = iris.data[:, :2], iris.target # we only keep the first two features: 'Sepal length', 'Sepal width'
qda = QuadraticDiscriminantAnalysis()
qda.fit(X, y)
x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5
h = .02 # step size in the mesh
xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
Z = qda.predict(np.c_[xx.ravel(), yy.ravel()]).reshape(xx.shape)
```



## Extensions of LDA

- ▶ There are many possibilities to extend LDA. The Bayes' construction

$$\Pr(Y = j \mid X = x) = \frac{\Pr(X = x \mid Y = j)\Pr(Y = j)}{\Pr(X = x)} = \frac{g_j(x)\pi_j}{\sum_{j=1}^q g_j(x)\pi_j}.$$

can be used with other densities  $g_j$ .

- ▶ **Quadratic discriminant analysis** (QDA): model  $g_j$  by multivariate normals with different covariance matrices  $\Sigma_j$ .
- ▶ **Mixture** models: use mixtures of Gaussians for the  $g_j$ .
- ▶ **Naive Bayes**: conditional independence assumption,  $g_j(x) = \prod_{l=1}^p g_{j,l}(x_l)$ .
- ▶ When some features are **categorical**, for example  $X_1$ , we can use

$$g_j(x) = \Pr(X_1 = x_1 \mid Y = j)g_j(x_2, \dots, x_p \mid x_1),$$

where  $g_j(x_2, \dots, x_p \mid x_1)$  is a density for  $X_2, \dots, X_p \mid X_1 = x_1, Y = j$ .

- ▶ **Non-parametric** approaches: estimate  $g_j$  non-parametrically.

## Comments on LDA

- ▶ LDA is easy and **fast to compute**.
- ▶ LDA provides estimates of the **posterior class probabilities**

$$\widehat{\Pr}(Y = j \mid X = x_0) = \frac{\hat{g}_j(x_0)\hat{\pi}_j}{\sum_{j=1}^q \hat{g}_j(x_0)\pi_j}.$$

- ▶ LDA usually performs well and can compete with more complex methods.
- ▶ Why? Clearly not because the normality assumption is realistic, but because a linear decision boundary is often the best the data can support and LDA is a **stable estimation method**.
- ▶ We don't need to have good models for the  $g_j$  to obtain a good model for the posterior: complex  $g_j$  may not be visible in the decision (see figure).

