

MATH 829: Introduction to Data Mining and
Analysis
Linear Discriminant Analysis

Dominique Guillot

Departments of Mathematical Sciences
University of Delaware

March 9, 2016

Linear discriminant analysis (LDA)

- Categorical data Y . Predictors X_1, \dots, X_p .

Linear discriminant analysis (LDA)

- Categorical data Y . Predictors X_1, \dots, X_p .
- We saw how *logistic regression* can be used to predict Y by modelling the log-odds

$$\log \frac{P(Y = 1|X = x)}{P(Y = 0|X = x)} = x^T \beta.$$

Linear discriminant analysis (LDA)

- Categorical data Y . Predictors X_1, \dots, X_p .
- We saw how *logistic regression* can be used to predict Y by modelling the log-odds

$$\log \frac{P(Y = 1|X = x)}{P(Y = 0|X = x)} = x^T \beta.$$

- More now examine other models for $P(Y = i|X = x)$.

Linear discriminant analysis (LDA)

- Categorical data Y . Predictors X_1, \dots, X_p .
- We saw how *logistic regression* can be used to predict Y by modelling the log-odds

$$\log \frac{P(Y = 1|X = x)}{P(Y = 0|X = x)} = x^T \beta.$$

- More now examine other models for $P(Y = i|X = x)$.

Recall: Bayes' theorem (Rev. Thomas Bayes, 1701–1761). Given two events A, B :

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$



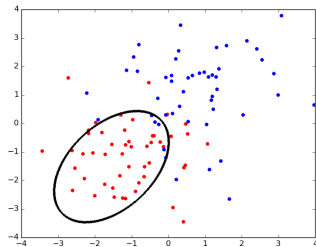
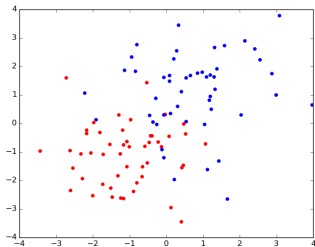
Source: Wikipedia (Public Domain).

Using Bayes' theorem

- $P(Y = i|X = x)$ harder to model.
- $P(X = x|Y = i)$ easier to model.

Using Bayes' theorem

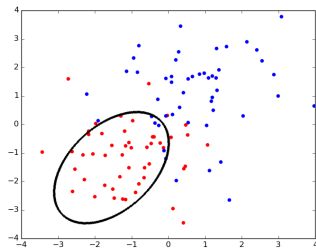
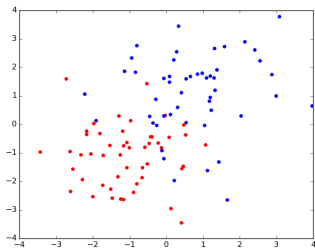
- $P(Y = i|X = x)$ harder to model.
- $P(X = x|Y = i)$ easier to model.



$$P(X = x|Y = \text{red}).$$

Using Bayes' theorem

- $P(Y = i|X = x)$ harder to model.
- $P(X = x|Y = i)$ easier to model.



$$P(X = x|Y = \text{red}).$$

Going back to our prediction using Bayes' theorem:

$$P(Y = i|X = x) = \frac{P(X = x|Y = i)P(Y = i)}{P(X = x)}$$

Using Bayes' theorem

More precisely, suppose

- $Y \in \{1, \dots, k\}$.
- $P(Y = i) = \pi_i \quad (i = 1, \dots, k)$.
- $P(X = x|Y = i) \sim f_i(x) \quad (i = 1, \dots, k)$.

Using Bayes' theorem

More precisely, suppose

- $Y \in \{1, \dots, k\}$.
- $P(Y = i) = \pi_i \quad (i = 1, \dots, k)$.
- $P(X = x|Y = i) \sim f_i(x) \quad (i = 1, \dots, k)$.

Then

$$\begin{aligned} P(Y = i|X = x) &= \frac{P(X = x|Y = i)P(Y = i)}{P(X = x)} \\ &= \frac{P(X = x|Y = i)P(Y = i)}{\sum_{j=1}^k P(X = x|Y = j)P(Y = j)} \\ &= \frac{f_i(x)\pi_i}{\sum_{j=1}^k f_j(x)\pi_j}. \end{aligned}$$

Using Bayes' theorem

More precisely, suppose

- $Y \in \{1, \dots, k\}$.
- $P(Y = i) = \pi_i \quad (i = 1, \dots, k)$.
- $P(X = x|Y = i) \sim f_i(x) \quad (i = 1, \dots, k)$.

Then

$$\begin{aligned} P(Y = i|X = x) &= \frac{P(X = x|Y = i)P(Y = i)}{P(X = x)} \\ &= \frac{P(X = x|Y = i)P(Y = i)}{\sum_{j=1}^k P(X = x|Y = j)P(Y = j)} \\ &= \frac{f_i(x)\pi_i}{\sum_{j=1}^k f_j(x)\pi_j}. \end{aligned}$$

- We can easily estimate π_i using the proportion of observations in category i .

Using Bayes' theorem

More precisely, suppose

- $Y \in \{1, \dots, k\}$.
- $P(Y = i) = \pi_i \quad (i = 1, \dots, k)$.
- $P(X = x|Y = i) \sim f_i(x) \quad (i = 1, \dots, k)$.

Then

$$\begin{aligned} P(Y = i|X = x) &= \frac{P(X = x|Y = i)P(Y = i)}{P(X = x)} \\ &= \frac{P(X = x|Y = i)P(Y = i)}{\sum_{j=1}^k P(X = x|Y = j)P(Y = j)} \\ &= \frac{f_i(x)\pi_i}{\sum_{j=1}^k f_j(x)\pi_j}. \end{aligned}$$

- We can easily estimate π_i using the proportion of observations in category i .
- We need a model for $f_i(x)$.

Using a Gaussian model: LDA and QDA

A natural model for the f_j s is the multivariate Gaussian distribution:

$$f_j(x) = \frac{1}{\sqrt{(2\pi)^p \det \Sigma_j}} e^{-\frac{1}{2}(x-\mu_j)^T \Sigma_j^{-1} (x-\mu_j)}.$$

Using a Gaussian model: LDA and QDA

A natural model for the f_j s is the multivariate Gaussian distribution:

$$f_j(x) = \frac{1}{\sqrt{(2\pi)^p \det \Sigma_j}} e^{-\frac{1}{2}(x-\mu_j)^T \Sigma_j^{-1} (x-\mu_j)}.$$

Linear discriminant analysis (LDA): We assume $\Sigma_j = \Sigma$ for all $j = 1, \dots, k$.

Quadratic discriminant analysis (QDA): general case, i.e., Σ_j can be distinct.

Using a Gaussian model: LDA and QDA

A natural model for the f_j s is the multivariate Gaussian distribution:

$$f_j(x) = \frac{1}{\sqrt{(2\pi)^p \det \Sigma_j}} e^{-\frac{1}{2}(x-\mu_j)^T \Sigma_j^{-1} (x-\mu_j)}.$$

Linear discriminant analysis (LDA): We assume $\Sigma_j = \Sigma$ for all $j = 1, \dots, k$.

Quadratic discriminant analysis (QDA): general case, i.e., Σ_j can be distinct.

Note: When p is large, using QDA instead of LDA can dramatically increase the number of parameters to estimate.

Using a Gaussian model: LDA and QDA

A natural model for the f_j s is the multivariate Gaussian distribution:

$$f_j(x) = \frac{1}{\sqrt{(2\pi)^p \det \Sigma_j}} e^{-\frac{1}{2}(x-\mu_j)^T \Sigma_j^{-1} (x-\mu_j)}.$$

Linear discriminant analysis (LDA): We assume $\Sigma_j = \Sigma$ for all $j = 1, \dots, k$.

Quadratic discriminant analysis (QDA): general case, i.e., Σ_j can be distinct.

Note: When p is large, using QDA instead of LDA can dramatically increase the number of parameters to estimate.

In order to use LDA or QDA, we need:

- An estimate of the class probabilities π_j .
- An estimate of the mean vectors μ_j .
- An estimate of the covariance matrices Σ_j (or Σ for LDA).

Estimating the parameters

LDA: Suppose we have N observations, and N_j of these observations belong to the j category ($j = 1, \dots, k$). We use

- $\hat{\pi}_j = N_j/N$.
- $\hat{\mu}_j = \frac{1}{N_j} \sum_{y_i=j} x_i$ (average of x over each category).
- $\hat{\Sigma} = \frac{1}{N-k} \sum_{j=1}^k \sum_{y_i=j} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T$. (Pooled variance.)

Estimating the parameters

LDA: Suppose we have N observations, and N_j of these observations belong to the j category ($j = 1, \dots, k$). We use

- $\hat{\pi}_j = N_j/N$.
- $\hat{\mu}_j = \frac{1}{N_j} \sum_{y_i=j} x_i$ (average of x over each category).
- $\hat{\Sigma} = \frac{1}{N-k} \sum_{j=1}^k \sum_{y_i=j} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T$. (Pooled variance.)

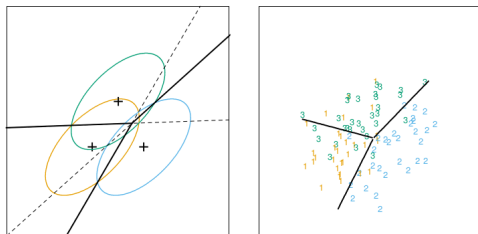


FIGURE 4.5. The left panel shows three Gaussian distributions, with the same covariance and different means. Included are the contours of constant density enclosing 95% of the probability in each case. The Bayes decision boundaries between each pair of classes are shown (broken straight lines), and the Bayes decision boundaries separating all three classes are the thicker solid lines (a subset of the former). On the right we see a sample of 30 drawn from each Gaussian distribution, and the fitted LDA decision boundaries.

LDA: linearity of the decision boundary

In the previous figure, we saw that the decision boundary is linear. Indeed, examining the *log-odds*:

$$\begin{aligned}\log \frac{P(Y = l|X = x)}{P(Y = m|X = x)} &= \log \frac{f_l(x)}{f_m(x)} + \log \frac{\pi_l}{\pi_m} \\ &= \log \frac{\pi_l}{\pi_m} - \frac{1}{2}(\mu_l + \mu_m)^T \Sigma^{-1}(\mu_l - \mu_m) + x^T \Sigma^{-1}(\mu_l - \mu_m) \\ &= \beta_0 + x^T \beta.\end{aligned}$$

LDA: linearity of the decision boundary

In the previous figure, we saw that the decision boundary is linear. Indeed, examining the *log-odds*:

$$\begin{aligned}\log \frac{P(Y = l|X = x)}{P(Y = m|X = x)} &= \log \frac{f_l(x)}{f_m(x)} + \log \frac{\pi_l}{\pi_m} \\ &= \log \frac{\pi_l}{\pi_m} - \frac{1}{2}(\mu_l + \mu_m)^T \Sigma^{-1}(\mu_l - \mu_m) + x^T \Sigma^{-1}(\mu_l - \mu_m) \\ &= \beta_0 + x^T \beta.\end{aligned}$$

Note that the previous expression is **linear** in x .

LDA: linearity of the decision boundary

In the previous figure, we saw that the decision boundary is linear. Indeed, examining the *log-odds*:

$$\begin{aligned}\log \frac{P(Y = l|X = x)}{P(Y = m|X = x)} &= \log \frac{f_l(x)}{f_m(x)} + \log \frac{\pi_l}{\pi_m} \\ &= \log \frac{\pi_l}{\pi_m} - \frac{1}{2}(\mu_l + \mu_m)^T \Sigma^{-1}(\mu_l - \mu_m) + x^T \Sigma^{-1}(\mu_l - \mu_m) \\ &= \beta_0 + x^T \beta.\end{aligned}$$

Note that the previous expression is **linear** in x . Recall that for logistic regression, we model

$$\log \frac{P(Y = l|X = x)}{P(Y = m|X = x)} = \beta_0 + x^T \beta.$$

LDA: linearity of the decision boundary

In the previous figure, we saw that the decision boundary is linear. Indeed, examining the *log-odds*:

$$\begin{aligned}\log \frac{P(Y = l|X = x)}{P(Y = m|X = x)} &= \log \frac{f_l(x)}{f_m(x)} + \log \frac{\pi_l}{\pi_m} \\ &= \log \frac{\pi_l}{\pi_m} - \frac{1}{2}(\mu_l + \mu_m)^T \Sigma^{-1}(\mu_l - \mu_m) + x^T \Sigma^{-1}(\mu_l - \mu_m) \\ &= \beta_0 + x^T \beta.\end{aligned}$$

Note that the previous expression is **linear** in x . Recall that for logistic regression, we model

$$\log \frac{P(Y = l|X = x)}{P(Y = m|X = x)} = \beta_0 + x^T \beta.$$

How is this different from LDA?

LDA: linearity of the decision boundary

In the previous figure, we saw that the decision boundary is linear. Indeed, examining the *log-odds*:

$$\begin{aligned}\log \frac{P(Y = l|X = x)}{P(Y = m|X = x)} &= \log \frac{f_l(x)}{f_m(x)} + \log \frac{\pi_l}{\pi_m} \\ &= \log \frac{\pi_l}{\pi_m} - \frac{1}{2}(\mu_l + \mu_m)^T \Sigma^{-1}(\mu_l - \mu_m) + x^T \Sigma^{-1}(\mu_l - \mu_m) \\ &= \beta_0 + x^T \beta.\end{aligned}$$

Note that the previous expression is **linear** in x . Recall that for logistic regression, we model

$$\log \frac{P(Y = l|X = x)}{P(Y = m|X = x)} = \beta_0 + x^T \beta.$$

How is this different from LDA?

- In LDA, the parameters are more constrained and are not estimated the same way.
- Can lead to smaller variance if the Gaussian model is correct.
- In practice, logistic regression is considered *safer* and *more robust*.
- LDA and logistic regression often return similar results.

QDA: quadratic decision boundary

Let us now examine the log-odds for QDA: in that case no simplification occurs as before

$$\begin{aligned} & \log \frac{P(Y = l|X = x)}{P(Y = m|X = x)} \\ &= \log \frac{\pi_l}{\pi_m} + \frac{1}{2} \log \frac{\det \Sigma_m}{\det \Sigma_l} \\ & \quad - \frac{1}{2} (x - \mu_l)^T \Sigma_l^{-1} (x - \mu_l) - \frac{1}{2} (x - \mu_m)^T \Sigma_l^{-1} (x - \mu_m). \end{aligned}$$

QDA: quadratic decision boundary

Let us now examine the log-odds for QDA: in that case no simplification occurs as before

$$\begin{aligned} & \log \frac{P(Y = l|X = x)}{P(Y = m|X = x)} \\ &= \log \frac{\pi_l}{\pi_m} + \frac{1}{2} \log \frac{\det \Sigma_m}{\det \Sigma_l} \\ & - \frac{1}{2} (x - \mu_l)^T \Sigma_l^{-1} (x - \mu_l) - \frac{1}{2} (x - \mu_m)^T \Sigma_l^{-1} (x - \mu_m). \end{aligned}$$

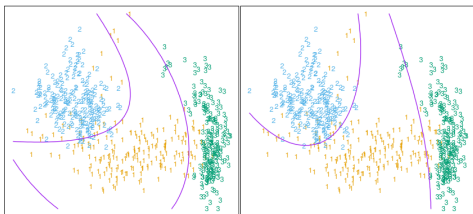


FIGURE 4.6. Two methods for fitting quadratic boundaries. The left plot shows the quadratic decision boundaries for the data in Figure 4.1 (obtained using LDA in the five-dimensional space $X_1, X_2, X_1X_2, X_1^2, X_2^2$). The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.

- Despite their simplicity, LDA and QDA often perform very well.
- Both techniques are widely used.

- Despite their simplicity, LDA and QDA often perform very well.
- Both techniques are widely used.

Problems when $n < p$:

- Estimating covariance matrices when n is small compared to p is challenging.
- The *sample covariance* (MLE for Gaussian)
$$S = \frac{1}{n-1} \sum_{j=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$
 has rank at most $\min(n, p)$ so is singular when $n < p$.
- This is a problem since Σ needs to be inverted in LDA and QDA.

- Despite their simplicity, LDA and QDA often perform very well.
- Both techniques are widely used.

Problems when $n < p$:

- Estimating covariance matrices when n is small compared to p is challenging.
- The *sample covariance* (MLE for Gaussian)
$$S = \frac{1}{n-1} \sum_{j=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$
 has rank at most $\min(n, p)$ so is singular when $n < p$.
- This is a problem since Σ needs to be inverted in LDA and QDA.

Many strategies exist to obtain better estimates of Σ (or Σ_j).

Among them:

- Regularization methods. E.g. $\hat{\Sigma}(\lambda) = \hat{\Sigma} + \lambda I$.
- Graphical modelling (discussed later during the course).

LDA:

```
from sklearn lda import LDA
```

QDA:

```
from sklearn qda import QDA
```