

Statistics 101C - Week 2 - Thursday

Shirong Xu

University of California, Los Angeles

shirong@stat.ucla.edu

October 9, 2024

Logistic Regression and Logistic Loss function

- Logistic Loss function:

$$L(f(\mathbf{x}_i), y_i) = \log \left(1 + \exp(-f(\mathbf{x}_i)y_i) \right)$$

- The expected logistic loss (Logistic Risk):

$$\begin{aligned} R_{\log}(f) &= \mathbb{E}_{\mathbf{X}, Y} (L(f(\mathbf{X}), Y)) \\ &= \mathbb{E}_{\mathbf{X}} \left[\mathbb{P}(Y = 1 | \mathbf{X}) \log \left(1 + \exp(-f(\mathbf{X})) \right) + \mathbb{P}(Y = -1 | \mathbf{X}) \log \left(1 + \exp(f(\mathbf{X})) \right) \right] \end{aligned}$$

- The optimal function minimizing $R_{\log}(f)$ is defined as

$$f_{\log}^*(\mathbf{X}) = \log \left(\frac{\mathbb{P}(Y = 1 | \mathbf{X})}{1 - \mathbb{P}(Y = 1 | \mathbf{X})} \right)$$

Logistic Loss for classification

- We can construct a classifier as $\text{sign}(f_{\log}^*(\mathbf{x})) \in \{-1, 1\}$. This classifier is identical to the Bayes classifier.

$$\text{sign}(f_{\log}^*(\mathbf{x})) = \text{sign}(\eta(\mathbf{x}) - 1/2)$$

Logistic Loss for classification

- Suppose we have a dataset $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, and we employ the logistic loss for classification. Then the training error in terms of logistic loss can be written as

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp(-f(\mathbf{x}_i)y_i) \right)$$

where $y_i \in \{-1, 1\}$.

Logistic Loss for classification

- Suppose we have a dataset $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, and we employ the logistic loss for classification. Then the training error in terms of logistic loss can be written as

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp(-f(\mathbf{x}_i)y_i) \right)$$

where $y_i \in \{-1, 1\}$.

- Here we do not make any assumption of f .

Logistic Loss for classification

- Suppose we have a dataset $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, and we employ the logistic loss for classification. Then the training error in terms of logistic loss can be written as

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp(-f(\mathbf{x}_i)y_i) \right)$$

where $y_i \in \{-1, 1\}$.

- Here we do not make any assumption of f .
- **Question:** What if we assume that $f(\mathbf{x}_i) = \beta_0 + \beta^T \mathbf{x}_i$

If $f(\mathbf{x}_i) = \beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i$

- $R_n(f)$ can be further written as

$$\begin{aligned} R_n(f) &= \frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp(-f(\mathbf{x}_i)y_i) \right) \\ &= \frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp \left(-(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i)y_i \right) \right) \end{aligned}$$

- We suppose that $y_i = 1 \Leftrightarrow \tilde{y}_i = 1$ and $y_i = -1 \Leftrightarrow \tilde{y}_i = 0$

$$\begin{aligned} R_n(f) &= \frac{1}{n} \sum_{i=1}^n \left[\tilde{y}_i \log \left(1 + \exp \left(-(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i) \right) \right) \right. \\ &\quad \left. + (1 - \tilde{y}_i) \log \left(1 + \exp \left(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i \right) \right) \right] \end{aligned}$$

$$\text{If } f(\mathbf{x}_i) = \beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i$$

- We suppose that $y_i = 1 \Leftrightarrow \tilde{y}_i = 1$ and $y_i = -1 \Leftrightarrow \tilde{y}_i = 0$

$$\begin{aligned} R_n(f) &= \frac{1}{n} \sum_{i=1}^n \left[\tilde{y}_i \log \left(1 + \exp \left(- (\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i) \right) \right) \right. \\ &\quad \left. + (1 - \tilde{y}_i) \log \left(1 + \exp \left(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i \right) \right) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \left[\tilde{y}_i \log \left(\exp \left(- (\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i) \right) \right) \right. \\ &\quad \left. + \log \left(1 + \exp \left(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i \right) \right) \right] \end{aligned}$$

Logistic Regression

- **Logistic regression** estimates the conditional probability probability:

$$\mathbb{P}(\tilde{Y} = 1 | \mathbf{X})$$

- In logistic regression, it is assumed that

$$\mathbb{P}(\tilde{Y} = 1 | \mathbf{X} = \mathbf{x}) = \frac{\exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})},$$
$$\mathbb{P}(\tilde{Y} = 0 | \mathbf{X} = \mathbf{x}) = \frac{1}{1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})},$$

where

- $\mathbf{x} = (x_1, \dots, x_p)^T$ is a p -dimensional predictor
- β_0 and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$ are unknown parameters
- $\boldsymbol{\beta}^T \mathbf{x} = \sum_{i=1}^p \beta_i x_i$

Estimation in Logistic Regression

- Suppose a dataset in logistic regression is $\{(\mathbf{x}_i, \tilde{y}_i)\}_{i=1}^n$, here $\tilde{y}_i \in \{0, 1\}$.
- Likelihood function $L(\beta_0, \boldsymbol{\beta})$:

$$L(\beta_0, \boldsymbol{\beta}) = \prod_{i=1}^n \left(\mathbb{P}(\tilde{Y} = 1 | \mathbf{X} = \mathbf{x}) \right)^{\tilde{y}_i} \left(\mathbb{P}(\tilde{Y} = 0 | \mathbf{X} = \mathbf{x}) \right)^{1-\tilde{y}_i}$$

Negative Log-likelihood

Negative Log-likelihood:

$$\begin{aligned} & -\log L(\beta_0, \beta) \\ &= -\sum_{i=1}^n \left[\tilde{y}_i \log \left(\frac{\exp(\beta_0 + \beta^T \mathbf{x}_i)}{1 + \exp(\beta_0 + \beta^T \mathbf{x}_i)} \right) + (1 - \tilde{y}_i) \log \left(\frac{1}{1 + \exp(\beta_0 + \beta^T \mathbf{x}_i)} \right) \right] \\ &= \sum_{i=1}^n \left[\tilde{y}_i \log \left(\exp(-\beta_0 - \beta^T \mathbf{x}_i) \right) + \log \left(1 + \exp(\beta_0 + \beta^T \mathbf{x}_i) \right) \right] \end{aligned}$$

If we take the expectation with respect to y_i , we have

$$\begin{aligned} & \mathbb{P}(\tilde{y}_i = 1 | \mathbf{x}_i) \log(1 + \exp(-\beta_0 - \beta^T \mathbf{x}_i)) \\ & + \mathbb{P}(\tilde{y}_i = 0 | \mathbf{x}_i) \log(1 + \exp(\beta_0 + \beta^T \mathbf{x}_i)) \end{aligned}$$

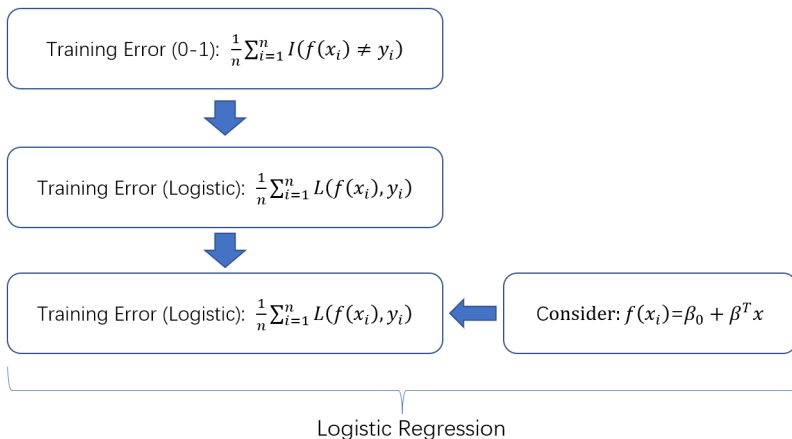
Conclusion

- Logistic Regression is an special example of using **Logistic loss** for classification.
- Logistic Regression assumes that the $\mathbb{P}(\tilde{Y} = 1|\mathbf{X} = \mathbf{x})$ as

$$\eta(\mathbf{x}) = \mathbb{P}(\tilde{Y} = 1|\mathbf{X} = \mathbf{x}) = \frac{\exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}.$$

- The optimal function minimize the logistic risk is $f^*(\mathbf{x}) = \log(\frac{\eta(\mathbf{x})}{1-\eta(\mathbf{x})}) = \beta_0 + \boldsymbol{\beta}^T \mathbf{x}.$

Relationship: Logistic Loss and Logistic Regression



- Discriminant Analysis
 - Linear Discriminant Analyses
 - Quadratic Discriminant Analysis

Classification

- A typical dataset in classification $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$.
 - \mathbf{x}_i : the covariate vector of i -th instance
 - $y_i \in \{0, 1\}$: binary label of i -th instance
- Bayes classifier f^* :

$$f^*(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) > 1/2 \\ 0 & \text{if } \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) < 1/2 \end{cases}$$

- Minimal risk $R(f^*)$:

$$R(f^*) = \mathbb{E}[f^*(\mathbf{X}) \neq Y] = \mathbb{E}[\min(\eta(\mathbf{X}), 1 - \eta(\mathbf{X}))],$$

where $\eta(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$.

How can we construct classifier?

- **Discriminative models**

- Discriminative modeling studies the $P(Y|\mathbf{X})$
- Examples: Logistic regression (LR)

- **Generative models**

- Generative models studies the joint probability distribution $\mathbb{P}(\mathbf{X}, Y)$
- Examples: linear discriminant analysis and quadratic discriminant analysis

Discriminant Analysis

- 1 Introduction
- 2 Linear and Quadratic Discriminant Analyses
- 3 LDA and QDA in practice

Basics of Generative models

- LDA and QDA are **generative models**, we need to consider the structure of $\mathbb{P}(\mathbf{X}, Y)$

$$\mathbb{P}(\mathbf{X}, Y) = \mathbb{P}(\mathbf{X}|Y)\mathbb{P}(Y)$$

$$\mathbb{P}(\mathbf{X}, Y) = \mathbb{P}(Y|\mathbf{X})\mathbb{P}(\mathbf{X})$$

An alternative look

Let $k \in \{0, 1\}$. We can develop an alternative formulation of $\mathbb{P}(Y = k | \mathbf{X} = \mathbf{x})$ from the definition of conditional probability.

$$\begin{aligned}\mathbb{P}(Y = k | \mathbf{X} = \mathbf{x}) &= \frac{\mathbb{P}(\mathbf{X} = \mathbf{x}, Y = k)}{\mathbb{P}(\mathbf{X} = \mathbf{x})} = \frac{\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k) \cdot \mathbb{P}(Y = k)}{\mathbb{P}(\mathbf{X} = \mathbf{x})} \\ &= \frac{\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k) \cdot \mathbb{P}(Y = k)}{\sum_{k=0}^1 \mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k) \cdot \mathbb{P}(Y = k)}\end{aligned}$$

- $\mathbb{P}(\mathbf{X} = \mathbf{x})$ the marginal distribution
- $\mathbb{P}(Y = k | \mathbf{X} = \mathbf{x})$: given $\mathbf{X} = \mathbf{x}$ the probability that outcome $Y = k$.

Banknote Dataset

conterfeit	Length	Left	Right	Bottom	Top	Diagonal
0	214.70000	129.70000	129.30000	8.60000	9.60000	141.60000
0	215.40000	130.00000	129.90000	8.50000	9.70000	141.40000
0	214.90000	129.40000	129.50000	8.20000	9.90000	141.50000
0	214.50000	129.50000	129.30000	7.40000	10.70000	141.50000
0	214.70000	129.60000	129.50000	8.30000	10.00000	142.00000
0	215.60000	129.90000	129.90000	9.00000	9.50000	141.70000
0	215.00000	130.40000	130.30000	9.10000	10.20000	141.10000
0	214.40000	129.70000	129.50000	8.00000	10.30000	141.20000
0	215.10000	130.00000	129.80000	9.10000	10.20000	141.50000
0	214.70000	130.00000	129.40000	7.80000	10.00000	141.20000
1	214.40000	130.10000	130.30000	9.70000	11.70000	139.80000
1	214.90000	130.50000	130.20000	11.00000	11.50000	139.50000
1	214.90000	130.30000	130.10000	8.70000	11.70000	140.20000
1	215.00000	130.40000	130.60000	9.90000	10.90000	140.30000
1	214.70000	130.20000	130.30000	11.80000	10.90000	139.70000
1	215.00000	130.20000	130.20000	10.60000	10.70000	139.90000
1	215.30000	130.30000	130.10000	9.30000	12.10000	140.20000

Discriminant Analysis (DA)

Discriminant Analysis models $\mathbb{P}(Y|\mathbf{X})$ as follows:

- Step 1: Make assumptions on data structure

Discriminant Analysis (DA)

Discriminant Analysis models $\mathbb{P}(Y|\mathbf{X})$ as follows:

- Step 1: Make assumptions on data structure
 - Let $\pi_k = \mathbb{P}(Y = k)$ be the prior probability of category $k = 0, 1$

Discriminant Analysis (DA)

Discriminant Analysis models $\mathbb{P}(Y|\mathbf{X})$ as follows:

- Step 1: Make assumptions on data structure
 - Let $\pi_k = \mathbb{P}(Y = k)$ be the prior probability of category $k = 0, 1$
 - Suppose that $\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k)$ is a multivariate normal distribution with mean vector $\boldsymbol{\mu}_k$ and covariance matrix $\boldsymbol{\Sigma}_k$

Discriminant Analysis (DA)

Discriminant Analysis models $\mathbb{P}(Y|\mathbf{X})$ as follows:

- Step 1: Make assumptions on data structure
 - Let $\pi_k = \mathbb{P}(Y = k)$ be the prior probability of category $k = 0, 1$
 - Suppose that $\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k)$ is a multivariate normal distribution with mean vector $\boldsymbol{\mu}_k$ and covariance matrix $\boldsymbol{\Sigma}_k$

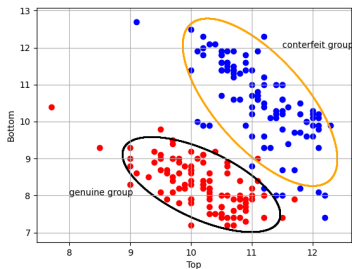


Figure: Black ellipsoid: covariance structure of genuine group. Green ellipsoid: covariance structure of the counterfeit group

Discriminant Analysis (DA)

- $\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k)$ is a multivariate normal distribution with mean $\boldsymbol{\mu}_k$ and covariance matrix $\boldsymbol{\Sigma}_k$.

$$\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right),$$

where

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix}, \boldsymbol{\mu}_k = \begin{pmatrix} \mu_{1,k} \\ \mu_{2,k} \\ \vdots \\ \mu_{p,k} \end{pmatrix}, \boldsymbol{\Sigma}_k = \begin{pmatrix} \sigma_{1,1,k}^2 & \sigma_{1,2,k}^2 & \cdots & \sigma_{1,p,k}^2 \\ \sigma_{2,1,k}^2 & \sigma_{2,2,k}^2 & \cdots & \sigma_{2,p,k}^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p,1,k}^2 & \sigma_{p,2,k}^2 & \cdots & \sigma_{p,p,k}^2 \end{pmatrix}.$$

Discriminant Analysis (DA)

- Step 2: We use the Bayes' theorem to compute $\mathbb{P}(Y = k|\mathbf{X} = \mathbf{x})$, $k = 0, 1$.

$$\mathbb{P}(Y = k|\mathbf{X} = \mathbf{x}) = \frac{\pi_k \mathbb{P}(\mathbf{X} = \mathbf{x}|Y = k)}{\pi_1 \mathbb{P}(\mathbf{X} = \mathbf{x}|Y = 1) + \pi_0 \mathbb{P}(\mathbf{X} = \mathbf{x}|Y = 0)}$$

Discriminant Analysis (DA)

- Step 2: We use the Bayes' theorem to compute $\mathbb{P}(Y = k|\mathbf{X} = \mathbf{x})$, $k = 0, 1$.

$$\mathbb{P}(Y = k|\mathbf{X} = \mathbf{x}) = \frac{\pi_k \mathbb{P}(\mathbf{X} = \mathbf{x}|Y = k)}{\pi_1 \mathbb{P}(\mathbf{X} = \mathbf{x}|Y = 1) + \pi_0 \mathbb{P}(\mathbf{X} = \mathbf{x}|Y = 0)}$$

Question: What is the difference between Linear and Quadratic discriminant analyses?

Discriminant Analysis (DA)

- Step 2: We use the Bayes' theorem to compute $\mathbb{P}(Y = k|\mathbf{X} = \mathbf{x})$, $k = 0, 1$.

$$\mathbb{P}(Y = k|\mathbf{X} = \mathbf{x}) = \frac{\pi_k \mathbb{P}(\mathbf{X} = \mathbf{x}|Y = k)}{\pi_1 \mathbb{P}(\mathbf{X} = \mathbf{x}|Y = 1) + \pi_0 \mathbb{P}(\mathbf{X} = \mathbf{x}|Y = 0)}$$

Question: What is the difference between Linear and Quadratic discriminant analyses?

- **Linear** Discriminant Analysis (LDA) assumes that the classes have a common covariance matrix. In other words, that is $\Sigma = \Sigma_0 = \Sigma_1$
- **Quadratic** Discriminant Analysis (QDA) does not assume this. So, we have a covariance matrix Σ_0 for class 0 and Σ_1 for class 1.

Linear Discriminant Analysis (LDA)

Three Assumptions in LDA

Linear Discriminant Analysis (LDA)

Three Assumptions in LDA

- 1 Multivariate normal distribution for each group, that $\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k)$ is multivariate normal

Linear Discriminant Analysis (LDA)

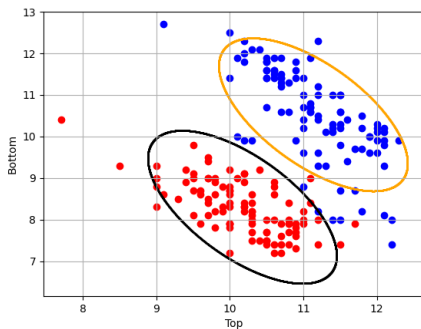
Three Assumptions in LDA

- 1 Multivariate normal distribution for each group, that $\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k)$ is multivariate normal
- 2 They have different mean vectors

Linear Discriminant Analysis (LDA)

Three Assumptions in LDA

- 1 Multivariate normal distribution for each group, that $\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k)$ is multivariate normal
- 2 They have different mean vectors
- 3 **Same covariance matrices**



Use LDA for classification

We make predictions using LDA as follows:

$$f_{LDA}(\mathbf{x}) = \begin{cases} 1, & \text{if } \frac{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=1)}{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=1) + \pi_0 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=0)} > 0.5 \\ 0, & \text{if } \frac{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=1)}{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=1) + \pi_0 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=0)} \leq 0.5 \end{cases}$$

Use LDA for classification

We make predictions using LDA as follows:

$$f_{LDA}(\mathbf{x}) = \begin{cases} 1, & \text{if } \frac{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=1)}{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=1) + \pi_0 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=0)} > 0.5 \\ 0, & \text{if } \frac{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=1)}{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=1) + \pi_0 \mathbb{P}(\mathbf{X}=\mathbf{x} | Y=0)} \leq 0.5 \end{cases}$$

Conclusions we can make

- 1 Similar to the Bayes classifier, we classify to the most probable class using the posterior probability

Use LDA for classification

We make predictions using LDA as follows:

$$f_{LDA}(\mathbf{x}) = \begin{cases} 1, & \text{if } \frac{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x}|Y=1)}{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x}|Y=1) + \pi_0 \mathbb{P}(\mathbf{X}=\mathbf{x}|Y=0)} > 0.5 \\ 0, & \text{if } \frac{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x}|Y=1)}{\pi_1 \mathbb{P}(\mathbf{X}=\mathbf{x}|Y=1) + \pi_0 \mathbb{P}(\mathbf{X}=\mathbf{x}|Y=0)} \leq 0.5 \end{cases}$$

Conclusions we can make

- 1 Similar to the Bayes classifier, we classify to the most probable class using the posterior probability
- 2 The decision boundary can be easily derived as

$$\begin{aligned} \frac{\pi_1 \mathbb{P}(\mathbf{X} = \mathbf{x} | Y = 1)}{\pi_1 \mathbb{P}(\mathbf{X} = \mathbf{x} | Y = 1) + \pi_0 \mathbb{P}(\mathbf{X} = \mathbf{x} | Y = 0)} &= 1/2 \\ \Leftrightarrow \log \frac{\pi_1}{\pi_0} + \log \frac{\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = 1)}{\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = 0)} &= 0. \end{aligned}$$

Decision boundary in LDA

A closer look at the decision boundary.

$$\log \frac{\pi_1}{\pi_0} + \log \frac{\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = 1)}{\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = 0)} = 0$$

\Updownarrow

$$\log \frac{\pi_1}{\pi_0} + \mathbf{x}^T \Sigma^{-1} (\mu_1 - \mu_0) - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 = 0$$

\Updownarrow

$$\log \frac{\pi_1}{\pi_0} + \mathbf{x}^T \Sigma^{-1} (\mu_1 - \mu_0) - \frac{1}{2} (\mu_1 + \mu_0)^T \Sigma^{-1} (\mu_1 - \mu_0) = 0.$$

Decision boundary in LDA

A closer look at the decision boundary.

$$\log \frac{\pi_1}{\pi_0} + \log \frac{\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = 1)}{\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = 0)} = 0$$

\Updownarrow

$$\log \frac{\pi_1}{\pi_0} + \mathbf{x}^T \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - \frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_0^T \Sigma^{-1} \boldsymbol{\mu}_0 = 0$$

\Updownarrow

$$\log \frac{\pi_1}{\pi_0} + \mathbf{x}^T \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - \frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0)^T \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) = 0.$$

The decision boundary can be written as (a linear equation)

$$\mathbf{x}^T C_1(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) + C_2(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) = 0,$$

where $C_1(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$ and

$C_2(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) = \log \frac{\pi_1}{\pi_0} - \frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0)^T \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0).$

Parameter estimation in LDA

Thanks to the formulation of LDA, we can easily estimate its parameters.

- The prior probability π_0 and π_1 .

$$\hat{\pi}_0 = \frac{n_0}{n_0 + n_1} \text{ and } \hat{\pi}_1 = \frac{n_1}{n_0 + n_1},$$

where n_k is the number of observations in the training data set that belong to class.

Parameter estimation in LDA

Thanks to the formulation of LDA, we can easily estimate its parameters.

- The prior probability π_0 and π_1 .

$$\hat{\pi}_0 = \frac{n_0}{n_0 + n_1} \text{ and } \hat{\pi}_1 = \frac{n_1}{n_0 + n_1},$$

where n_k is the number of observations in the training data set that belong to class.

- The means are estimated as

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} \mathbf{x}_i, k = 0, 1$$

Parameter estimation in LDA

Thanks to the formulation of LDA, we can easily estimate its parameters.

- The prior probability π_0 and π_1 .

$$\hat{\pi}_0 = \frac{n_0}{n_0 + n_1} \text{ and } \hat{\pi}_1 = \frac{n_1}{n_0 + n_1},$$

where n_k is the number of observations in the training data set that belong to class.

- The means are estimated as

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} \mathbf{x}_i, k = 0, 1$$

- The covariance matrices are estimated as

$$\hat{\Sigma} = \frac{1}{n-2} \sum_{k=0}^1 \sum_{i:y_i=k} (\mathbf{x}_i - \hat{\mu}_k)(\mathbf{x}_i - \hat{\mu}_k)^T$$

Quadratic Discriminant Analysis (QDA)

Three Assumptions in QDA

- 1 Multivariate normal distribution for each group, that $\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = k)$ is multivariate normal
- 2 They have different mean vectors
- 3 **Different covariance matrices**

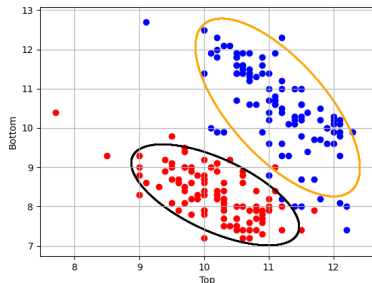


Figure: Different covariance structures

Decision boundary in QDA

We follow a similar analysis of QDA as with LDA. After some algebra, we arrive to the following (interesting) equation:

$$\log \frac{\pi_1}{\pi_0} - \frac{1}{2} \mathbf{x}^T (\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_0^{-1}) \mathbf{x} + \mathbf{x}^T (\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0) + \dots = 0$$

Conclusion

- The decision boundary in QDA is a quadratic function

LDA vs QDA

The difference between LDA and QDA can be summarized as

- LDA is simpler than QDA. (LDA is a special case of QDA)
- QDA needs to estimate more parameters. One covariance matrix for each class.
- LDA is much less flexible than QDA, but this also means that it has low variance
- If the assumptions of LDA do not hold, then it can lead to poor estimates and so, a high bias.

Exercise: Prediction of counterfeit banknotes

counterfeit	Length	Left	Right	Bottom	Top	Diagonal
0	214.70000	129.70000	129.30000	8.60000	9.60000	141.60000
0	215.40000	130.00000	129.90000	8.50000	9.70000	141.40000
0	214.90000	129.40000	129.50000	8.20000	9.90000	141.50000
0	214.50000	129.50000	129.30000	7.40000	10.70000	141.50000
0	214.70000	129.60000	129.50000	8.30000	10.00000	142.00000
0	215.60000	129.90000	129.90000	9.00000	9.50000	141.70000
0	215.00000	130.40000	130.30000	9.10000	10.20000	141.10000
0	214.40000	129.70000	129.50000	8.00000	10.30000	141.20000
0	215.10000	130.00000	129.80000	9.10000	10.20000	141.50000
0	214.70000	130.00000	129.40000	7.80000	10.00000	141.20000
1	214.40000	130.10000	130.30000	9.70000	11.70000	139.80000
1	214.90000	130.50000	130.20000	11.00000	11.50000	139.50000
1	214.90000	130.30000	130.10000	8.70000	11.70000	140.20000
1	215.00000	130.40000	130.60000	9.90000	10.90000	140.30000
1	214.70000	130.20000	130.30000	11.80000	10.90000	139.70000
1	215.00000	130.20000	130.20000	10.60000	10.70000	139.90000
1	215.30000	130.30000	130.10000	9.30000	12.10000	140.20000

- Length: length of banknote (mm)
- Left: length of left edge (mm)
- Right: length of right edge (mm)
- Top: distance from the image to top edge
- Bottom: distance from image to bottom
- Diagonal: length of diagonal (mm)
- counterfeit: 1 means counterfeit and 0 means genuine

Exercise: Prediction of counterfeit banknotes using R

- Step 1: Loading the dataset and split the dataset into training set and testing set:

```
library(mclust)
# Load the data set.
data(banknote)
banknote$Status<-factor(banknote$Status,levels=c("genuine", "counterfeit"))
# Split into training and test data.
set.seed(123) # Set seed to reproduce results.
i <- 1:dim(banknote)[1]
# Generate a random sample.
i.train <- sample(i, 130, replace = F) # 130 samples are used for training
bn.train <- banknote[i.train,] # training dataset
bn.test <- banknote[-i.train,] # testing dataset
```

- Step 2: Implement LDA and make prediction by LDA

```
library(MASS)
lda.mod <- lda(Status~Length + Right + Left + Top, data = bn.train) # Fit a LDA model
pred.lda.test <- predict(lda.mod,bn.test[,1])

table('Reference' = bn.test[,1], "Predicted" =
      pred.lda.test$class)
```

Exercise: Prediction of counterfeit banknotes using R

- Result:

```
> table('Reference' = bn.test[,1], "Predicted" =  
+       pred.lda.test$class)  
      Predicted  
Reference  genuine counterfeit  
genuine      30           3  
counterfeit   4          33
```

- Conclusion: The prediction accuracy of LDA is $(30+33)/70=0.9$.

Exercise: Prediction of counterfeit banknotes using R

- Implementation of QDA

```
qda.mod <- qda(Status~Length + Right + Left + Top, data = bn.train)
pred.qda.test <- predict(qda.mod, bn.test[, -1])
table('Reference' = bn.test[, 1], |
      "Predicted" = pred.qda.test$class)
```

- Result:

```
> table('Reference' = bn.test[, 1],
+       "Predicted" = pred.qda.test$class)
      Predicted
Reference      genuine counterfeit
genuine         29             4
counterfeit      3            34
```

- Conclusion: The prediction accuracy of LDA is $(29+34)/70=0.9$. No improvement is observed.

Some questions

- Can you finish an implementation of LDA and QDA in R or Python?
- Can you summarize the difference between Logistic regression, LDA, and QDA?
- What is the difference between generative model and discriminative model?
- Is K-nearest neighbor classifier a generative model or a discriminative model?