

Multivariate Analysis

Discriminant Analysis and Classification

Shaobo Jin

Department of Mathematics

Intended Learning Outcome

Through this chapter, you should be able to

- ① apply different classifiers,
- ② compare different classifiers.

Motivation

Discrimination and **classification** are used to separate observations and allocate new subjects to previously defined classes. They are called **supervised learning** in machine learning.

- ① Discrimination: use labeled observations to build a classification rule (classifier).
- ② Classification: allocate new subjects to classes using classifiers.

Some examples are

- ① given the medical history, whether a patient has cancer or not.
- ② given your credit history, whether you can pay back the mortgage loan or not.
- ③ given your hotel booking history, which hotel brand you will choose, Scandic, Best Western, or Radisson.

Two-Class Problem

Suppose that there are in total two classes (pass or fail, healthy or sick, positive or negative). The classes are labeled as π_1 and π_2 . For each subject, we observe a $p \times 1$ vector of covariates \mathbf{X} .

- We need to partition the sample space into two disjoint region R_1 and R_2 , such that $R_1 \cup R_2 = \Omega$, where Ω is the sample space.
- If a new vector $\mathbf{x}_0 \in R_1$ ($\mathbf{x}_0 \in R_2$), then the subject is classified to class π_1 (π_2).

How we partition and classify the subjects can be either **probabilistic** or **deterministic**.

Misclassification

A good classifier should result in few misclassifications or low costs associated with misclassification. Let $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ be the density functions associated with the $p \times 1$ vector of covariates \mathbf{X} for the populations π_1 and π_2 . Then,

$$P(i | j) := P(\text{Classify as } \pi_i \mid \text{truth } \pi_j) = \int_{\mathbf{x} \in R_i} f_j(\mathbf{x}) d\mathbf{x}.$$

Hence, for $j \neq i$,

$$\begin{aligned} P(\text{Correctly classified as } \pi_i) &= P(\text{Classify as } \pi_i \mid \text{truth } \pi_i) P(\text{truth } \pi_i), \\ P(\text{Misclassified as } \pi_i) &= P(\text{Classify as } \pi_i \mid \text{truth } \pi_j) P(\text{truth } \pi_j), \end{aligned}$$

where $P(\text{truth } \pi_i)$ is some prior probability such that

$$P(\text{truth } \pi_1) + P(\text{truth } \pi_2) = 1.$$

Classification Table

A **classification table** or a **confusion matrix** cross-classifies the binary response y and the predicted binary response.

Predicted	Observed	
	1 (TRUE)	2 (FALSE)
1 (TRUE)	m_{11} (TP)	m_{12} (FP)
2 (FALSE)	m_{21} (FN)	m_{22} (TN)

From the classification table, we can compute

$$\text{sensitivity (true positive rate)} = \frac{m_{11}}{m_{11} + m_{21}},$$

$$\text{specificity (true negative rate)} = \frac{m_{22}}{m_{12} + m_{22}},$$

$$\text{false positive rate} = \frac{m_{12}}{m_{12} + m_{22}},$$

$$\text{false negative rate} = \frac{m_{21}}{m_{11} + m_{21}}.$$

F-Score of Binary Classification

Predicted	Observed	
	1 (TRUE)	2 (FALSE)
1 (TRUE)	m_{11}	m_{12}
2 (FALSE)	m_{21}	m_{22}

In machine learning,

$$\text{Precision} = \frac{m_{11}}{m_{11} + m_{12}},$$

$$\text{Recall} = \frac{m_{11}}{m_{11} + m_{21}}.$$

The **F-score** is

$$\text{F-score} = 2 \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2m_{11}}{2m_{11} + m_{12} + m_{21}} \in [0, 1].$$

Cost of Misclassification

Let $c(i | j)$ be the cost of classifying π_j to π_i . The **cost matrix** tabulates the costs of misclassification:

Predicted	Observed	
	1 (TRUE)	2 (FALSE)
1 (TRUE)	0	$c(1 2)$
2 (FALSE)	$c(2 1)$	0

The **expected cost of misclassification (ECM)** is

$$\text{ECM} = c(2 | 1) P(2 | 1) p_1 + c(1 | 2) P(1 | 2) p_2,$$

where $p_i = P(\text{truth } \pi_i)$. A reasonable classification rule should have a small ECM.

Loosely speaking, since it requires the prior probability p_i , it can be viewed as **Bayesian discrimination**.

Minimizing ECM

Result 11.1: Known Prior Probabilities

The regions R_1 and R_2 that minimize the ECM are defined by the values \mathbf{x} for which the following inequalities hold

$$\begin{aligned} R_1 &= \left\{ \mathbf{x}; \quad \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{c(1|2)p_2}{c(2|1)p_1} \right\}, \\ R_2 &= \left\{ \mathbf{x}; \quad \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{c(1|2)p_2}{c(2|1)p_1} \right\}. \end{aligned}$$

An Example Using ECM

Known Prior Probabilities

Suppose that $\mathbf{X} \sim N(\mathbf{0}, \Sigma)$ in π_1 , and $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$ in π_2 . It is known that $p_1 = 0.8$ and $p_2 = 0.2$. Let $c(2 | 1) = 5$ and $c(1 | 2) = 10$. Then,

$$\frac{c(1 | 2) p_2}{c(2 | 1) p_1} = \frac{10 \cdot 0.2}{5 \cdot 0.8} = 0.5.$$

If a new subject with $\mathbf{X} = \mathbf{x}_0$ satisfies

$$\frac{f_1(\mathbf{x}_0)}{f_2(\mathbf{x}_0)} = \frac{\exp\left\{-\frac{1}{2}\mathbf{x}_0^T \Sigma^{-1} \mathbf{x}_0\right\}}{\exp\left\{-\frac{1}{2}(\mathbf{x}_0 - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_0 - \boldsymbol{\mu})\right\}} \leq 0.5,$$

then it will be classified to π_1 . Otherwise, it is classified to π_2 .

Special Case I: Highest Likelihood

Suppose that

$$\frac{c(1|2)p_2}{c(2|1)p_1} = 1.$$

Then the classifier that minimizes ECM is

$$\begin{aligned} R_1 &= \{\mathbf{x}; f_1(\mathbf{x}) \geq f_2(\mathbf{x})\}, \\ R_2 &= \{\mathbf{x}; f_1(\mathbf{x}) < f_2(\mathbf{x})\}. \end{aligned}$$

The densities $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ can be viewed as conditional density (or [likelihood](#))

$$f_i(\mathbf{x}) = f(\mathbf{x} | \text{truth } \pi_i).$$

Special Case II: Total Probability of Misclassification

The total probability of misclassification is

$$\begin{aligned}\text{TPM} &= P(\text{Misclassifying a observation}) \\ &= P(\text{Classify as } \pi_2, \text{ truth } \pi_1) + P(\text{Classify as } \pi_1, \text{ truth } \pi_2) \\ &= P(2 | 1)p_1 + P(1 | 2)p_2,\end{aligned}$$

which is equivalent to ECM when $c(1 | 2) = c(2 | 1)$. Hence, the regions R_1 and R_2 that minimize the TPM are

$$\begin{aligned}R_1 &= \left\{ \mathbf{x}; \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{p_2}{p_1} \right\}, \\ R_2 &= \left\{ \mathbf{x}; \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{p_2}{p_1} \right\}.\end{aligned}$$

Special Case III: Highest Posterior Probability

The densities $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ can be viewed as conditional density (or [likelihood](#))

$$f_i(\mathbf{x}) = f(\mathbf{x} \mid \text{truth } \pi_i).$$

By the [Bayes rule](#),

$$\begin{aligned} P(\text{truth } \pi_i \mid \text{observe } \mathbf{x}_0) &= \frac{P(\text{observe } \mathbf{x}_0 \mid \text{truth } \pi_i) P(\text{truth } \pi_i)}{\sum_{j=1}^2 P(\text{observe } \mathbf{x}_0 \mid \text{truth } \pi_j) P(\text{truth } \pi_j)} \\ &= \frac{f_i(\mathbf{x}_0) p_i}{f_1(\mathbf{x}_0) p_1 + f_2(\mathbf{x}_0) p_2}. \end{aligned}$$

The [Naive Bayes](#) classifier classifies \mathbf{x}_0 to π_i if

$$P(\text{truth } \pi_i \mid \text{observe } \mathbf{x}_0) > P(\text{truth } \pi_j \mid \text{observe } \mathbf{x}_0).$$

This is equivalent to the classifier that minimizes TPM.

Naive Bayes: Gaussian Discriminant Analysis

So far, we assume $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ are known. If they are unknown, we need to estimate them from the data. We illustrate the idea with the naive Bayes classifier.

- Let $Z \sim \text{Bernoulli}(\phi)$ be a random variable that indicates the class, where $Z = 1$ or 0 :

$$\mathbf{X} \mid Z = 1 \sim N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}), \quad \mathbf{X} \mid Z = 0 \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}).$$

- The likelihood function is

$$L(\boldsymbol{\mu}_1, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}, \phi) = \prod_{j=1}^n P(\mathbf{x}_j, z_j) = \prod_{j=1}^n [f_1(\mathbf{x}_j) \phi]^{z_j} [f_0(\mathbf{x}_j) (1 - \phi)]^{1-z_j}.$$

- The MLEs are

$$\begin{aligned} \hat{\phi} &= n^{-1} \sum_{j=1}^n z_j, \quad \hat{\boldsymbol{\mu}}_1 = \frac{\sum_{j=1}^n z_j \mathbf{x}_j}{\sum_{j=1}^n z_j}, \quad \hat{\boldsymbol{\mu}}_0 = \frac{\sum_{j=1}^n (1 - z_j) \mathbf{x}_j}{\sum_{j=1}^n (1 - z_j)}, \\ \hat{\boldsymbol{\Sigma}} &= \frac{1}{n} \left[\sum_{j=1}^n z_j (\mathbf{x}_j - \hat{\boldsymbol{\mu}}_1)(\mathbf{x}_j - \hat{\boldsymbol{\mu}}_1)^T + \sum_{j=1}^n (1 - z_j) (\mathbf{x}_j - \hat{\boldsymbol{\mu}}_0)(\mathbf{x}_j - \hat{\boldsymbol{\mu}}_0)^T \right]. \end{aligned}$$

Gaussian Discriminant Analysis: Classification

We classify \mathbf{x}_0 to class 1, if $P(Z = 1 | \mathbf{x}_0) \geq P(Z = 0 | \mathbf{x}_0)$, where

$$P(Z = 1 | \mathbf{x}_0) = \frac{\frac{\hat{\phi}}{(2\pi)^{p/2} \sqrt{\det(\hat{\Sigma})}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_1)^T \hat{\Sigma}^{-1} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_1) \right\}}{P(\mathbf{x}_0)},$$

$$P(Z = 0 | \mathbf{x}_0) = \frac{\frac{1-\hat{\phi}}{(2\pi)^{p/2} \sqrt{\det(\hat{\Sigma})}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_0)^T \hat{\Sigma}^{-1} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_0) \right\}}{P(\mathbf{x}_0)}.$$

It is equivalent to

$$\frac{\hat{\phi} \exp \left\{ -\frac{1}{2} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_1)^T \hat{\Sigma}^{-1} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_1) \right\}}{(1 - \hat{\phi}) \exp \left\{ -\frac{1}{2} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_0)^T \hat{\Sigma}^{-1} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_0) \right\}} \geq 1.$$

Gaussian Discriminant Analysis: Decision Boundary

We can also say that we classify \mathbf{x}_0 to class 1, if

$$\begin{aligned}\log \left(\frac{\hat{\phi}}{1 - \hat{\phi}} \right) &\geq \frac{1}{2} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_1)^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_1) \\ &\quad - \frac{1}{2} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_0)^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_0) \\ &= (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1)^T \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{x}_0 + \frac{1}{2} \left(\hat{\boldsymbol{\mu}}_1^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_0^T \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_0 \right),\end{aligned}$$

which is a quadratic function in \mathbf{x}_0 . Hence, the **decision boundary** is linear in \mathbf{x}_0 .

We have assume $\boldsymbol{\Sigma}$ to be the same for different groups. You can use other assumptions, e.g., different $\boldsymbol{\Sigma}$, diagonal $\boldsymbol{\Sigma}$, etc.

Fisher's Linear Discriminant Analysis

We need to know or estimate $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ in order to use the classifiers based on ECM. Fisher's linear discriminant analysis (LDA) does not assume any distributional assumption. The idea is that a linear combination of \mathbf{x} ($y = \mathbf{a}^T \mathbf{x}$) should carry enough information for discrimination. If we observe

① $\mathbf{x}_{11}, \dots, \mathbf{x}_{1n_1}$ from π_1 ,

② $\mathbf{x}_{21}, \dots, \mathbf{x}_{2n_2}$ from π_2 ,

then we should find \mathbf{a} such that the linear combination of the first class should be as far from the linear combination of the second class as possible. We need to assume two classes have the same covariance matrix, i.e., $\Sigma_1 = \Sigma_2 = \Sigma$.

MANOVA-Like Idea

Let $\bar{\mathbf{x}}_i$ be the sample mean of class π_i , and $\bar{\mathbf{x}}$ be the sample mean ignoring classes. Then,

$$\begin{aligned} \sum_{i=1}^2 \sum_{j=1}^{n_j} (\mathbf{a}^T \mathbf{x}_{ij} - \mathbf{a}^T \bar{\mathbf{x}})^2 &= \sum_{i=1}^2 \sum_{j=1}^{n_j} \mathbf{a}^T (\mathbf{x}_{ij} - \bar{\mathbf{x}}) (\mathbf{x}_{ij} - \bar{\mathbf{x}})^T \mathbf{a} \\ &= \mathbf{a}^T \left[\sum_{i=1}^2 \sum_{j=1}^{n_j} (\mathbf{x}_{ij} - \bar{\mathbf{x}}) (\mathbf{x}_{ij} - \bar{\mathbf{x}})^T \right] \mathbf{a} \\ &= \mathbf{a}^T [\mathbf{W} + \mathbf{B}] \mathbf{a}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{W} &= \sum_{i=1}^2 \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i) (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)^T, \\ \mathbf{B} &= \sum_{i=1}^2 n_i (\bar{\mathbf{x}}_i - \bar{\mathbf{x}}) (\bar{\mathbf{x}}_i - \bar{\mathbf{x}})^T. \end{aligned}$$

Within Versus Between Variation

In LDA, we want to maximize

$$\frac{\mathbf{a}^T \mathbf{B} \mathbf{a}}{\mathbf{a}^T \mathbf{W} \mathbf{a}}.$$

The maximizer \mathbf{a} is the eigenvector of $\mathbf{W}^{-1} \mathbf{B}$ corresponding to the largest eigenvalue of $\mathbf{W}^{-1} \mathbf{B}$. In this case

$$\mathbf{a} = \mathbf{W}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2),$$

which is called [Fisher's linear discriminant function](#).

Fisher's LDA

We classify \mathbf{x}_0 to π_1 if $\mathbf{a}^T \mathbf{x}_0$ is closer to the mean of $\mathbf{a}^T \mathbf{X}_1$ (first class) than the mean of $\mathbf{a}^T \mathbf{X}_2$ (second class). That is,

$$\left| (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{W}^{-1} (\mathbf{x}_0 - \bar{\mathbf{x}}_1) \right| < \left| (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{W}^{-1} (\mathbf{x}_0 - \bar{\mathbf{x}}_2) \right|.$$

This is equivalent to classifying \mathbf{x}_0 to π_1 if

$$(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{W}^{-1} \mathbf{x}_0 \geq \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{W}^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2).$$

The decision boundary is linear in \mathbf{x}_0 . Hence, it is called Fisher's [linear discriminant analysis](#) ([LDA](#)).

Two Normal Populations

Now we assume that $\mathbf{X} \sim N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ in π_1 , and $\mathbf{X} \sim N(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ in π_2 . We will derive classifiers for two multivariate normal populations.

- ① Case I: $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}$,
- ② Case II: $\boldsymbol{\Sigma}_1 \neq \boldsymbol{\Sigma}_2$.

Case I: $\Sigma_1 = \Sigma_2 = \Sigma$

Suppose that

- ① $\Sigma_1 = \Sigma_2 = \Sigma$,
- ② μ_1 , μ_2 , and Σ are known.

The classifier that minimizes ECM is the one that classifies \mathbf{x}_0 to π_1 if

$$\frac{f_1(\mathbf{x}_0)}{f_2(\mathbf{x}_0)} = \frac{\exp\left\{-\frac{1}{2}(\mathbf{x}_0 - \mu_1)^T \Sigma^{-1}(\mathbf{x}_0 - \mu_1)\right\}}{\exp\left\{-\frac{1}{2}(\mathbf{x}_0 - \mu_2)^T \Sigma^{-1}(\mathbf{x}_0 - \mu_2)\right\}} \leq \frac{c(1|2)p_2}{c(2|1)p_1}.$$

It is equivalent to classifying \mathbf{x}_0 to π_1 if

$$(\mu_1 - \mu_2)^T \Sigma^{-1} \mathbf{x}_0 - \frac{1}{2}(\mu_1 - \mu_2)^T \Sigma^{-1}(\mu_1 + \mu_2) \leq \log \left[\frac{c(1|2)p_2}{c(2|1)p_1} \right].$$

Estimated μ and Σ

μ_1 , μ_2 , and Σ are typically unknown. We need to estimate them in order to classify \mathbf{x}_0 .

- If we observe $\mathbf{x}_{11}, \dots, \mathbf{x}_{1n_1}$ from π_1 and $\mathbf{x}_{21}, \dots, \mathbf{x}_{2n_2}$ from π_2 , then

$$\begin{aligned}\hat{\mu}_1 &= \bar{\mathbf{x}}_1, \quad \hat{\mu}_2 = \bar{\mathbf{x}}_2, \\ \mathbf{S}_{\text{pooled}} &= \frac{(n_1 - 1) \mathbf{S}_1 + (n_2 - 1) \mathbf{S}_2}{n_1 + n_2 - 2}.\end{aligned}$$

- We allocate \mathbf{x}_0 to π_1 if

$$\begin{aligned}& (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{\text{pooled}}^{-1} \mathbf{x}_0 - \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{\text{pooled}}^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) \\ & \geq \log \left[\frac{c(1|2) p_2}{c(2|1) p_1} \right].\end{aligned}$$

Connection to Fisher's LDA

Let $y = \mathbf{a}^T \mathbf{x}_0$ and $\hat{y} = \hat{\mathbf{a}}^T \mathbf{x}_0$, where $\mathbf{a} = \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ and $\hat{\mathbf{a}} = \mathbf{S}_{\text{pooled}}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$. Consider a special case where

$$\frac{c(1|2) p_2}{c(2|1) p_1} = 1.$$

We allocate \mathbf{x}_0 to π_1 if

$$\underbrace{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{\text{pooled}}^{-1}}_{\hat{\mathbf{a}}^T} \mathbf{x}_0 \geq \frac{1}{2} \underbrace{(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{\text{pooled}}^{-1}}_{\hat{\mathbf{a}}^T} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2).$$

Note that

$$\begin{aligned} \mathbf{W} &= \sum_{i=1}^2 \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i) (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)^T, \\ \mathbf{S}_{\text{pooled}} &= \frac{(n_1 - 1) \mathbf{S}_1 + (n_2 - 1) \mathbf{S}_2}{n_1 + n_2 - 2} = \frac{\mathbf{W}}{n_1 + n_2 - 2}. \end{aligned}$$

Hence, it is identical to [Fisher's LDA](#).

Connection to Naive Bayes

Consider another special case where

$$\frac{c(1 | 2)}{c(2 | 1)} = 1.$$

We allocate \mathbf{x}_0 to π_1 if

$$(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{\text{pooled}}^{-1} \mathbf{x}_0 - \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}_{\text{pooled}}^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) \geq \log \left(\frac{\hat{p}_2}{\hat{p}_1} \right).$$

This is equivalent to the [Gaussian discriminant analysis](#).

Case II: $\Sigma_1 \neq \Sigma_2$

Suppose that $\Sigma_1 \neq \Sigma_2$. The minimum ECM classifier allocates \mathbf{x}_0 to π_1 , if

$$\frac{[\det(\Sigma_1)]^{-p/2} \exp\left\{-\frac{1}{2}(\mathbf{x}_0 - \boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x}_0 - \boldsymbol{\mu}_1)\right\}}{[\det(\Sigma_2)]^{-p/2} \exp\left\{-\frac{1}{2}(\mathbf{x}_0 - \boldsymbol{\mu}_2)^T \Sigma^{-1}(\mathbf{x}_0 - \boldsymbol{\mu}_2)\right\}} \geq \frac{c(1|2)p_2}{c(2|1)p_1}.$$

This is equivalent to allocating \mathbf{x}_0 to π_1 , if

$$-\frac{1}{2}\mathbf{x}_0^T(\Sigma_1^{-1} - \Sigma_2^{-1})\mathbf{x}_0 + (\boldsymbol{\mu}_1^T \Sigma_1^{-1} - \boldsymbol{\mu}_2^T \Sigma_2^{-1})\mathbf{x}_0 - k \geq \log\left[\frac{c(1|2)p_2}{c(2|1)p_1}\right],$$

where

$$k = \frac{p}{2} \log\left[\frac{\det(\Sigma_1)}{\det(\Sigma_2)}\right] + \frac{1}{2}\boldsymbol{\mu}_1^T \Sigma_1^{-1} \boldsymbol{\mu}_1 - \frac{1}{2}\boldsymbol{\mu}_2^T \Sigma_2^{-1} \boldsymbol{\mu}_2.$$

The decision boundary is a quadratic function of \mathbf{x} . Hence, this is called **quadratic discriminant analysis (QDA)**.

Logistic Model For Two Populations

Let $Y_j \mid \mathbf{X}_j \sim \text{Bernoulli}(\phi_j)$ be a random variable that indicates the class of subject j , where $Y_j = 1$ for class π_1 and $Y_j = 0$ for class π_0 . In logistic regression, we model ϕ_j as

$$\log \left(\frac{\phi_j}{1 - \phi_j} \right) = \alpha + \boldsymbol{\beta}^T \mathbf{x}_j,$$

that is

$$\begin{aligned} \phi_j = P(Y_j = 1 \mid \mathbf{x}_j) &= \frac{\exp(\alpha + \boldsymbol{\beta}^T \mathbf{x}_j)}{1 + \exp(\alpha + \boldsymbol{\beta}^T \mathbf{x}_j)}, \\ 1 - \phi_j = P(Y_j = 0 \mid \mathbf{x}_j) &= \frac{1}{1 + \exp(\alpha + \boldsymbol{\beta}^T \mathbf{x}_j)}. \end{aligned}$$

Maximum Likelihood Estimator

We often assume that $\{Y_j\}$ are mutually independent given \mathbf{X} . Then, the log-likelihood function is

$$\begin{aligned}\ell(\alpha, \beta) &= \sum_{j=1}^n [I(y_j = 1) \log \phi_j + I(y_j = 0) \log (1 - \phi_j)] \\ &= \sum_{j=1}^n \{I(y_j = 1) (\alpha + \beta^T \mathbf{x}_j) - \log [1 + \exp (\alpha + \beta^T \mathbf{x}_j)]\},\end{aligned}$$

where $I(y_j = k) = 1$ if $y_j = k$, and 0 otherwise. The estimator $(\hat{\alpha}, \hat{\beta})$ maximizes $L(\alpha, \beta)$ using numerical methods.

We often allocate \mathbf{x}_0 to π_1 if

$$\frac{\exp(\hat{\alpha} + \hat{\beta}^T \mathbf{x}_0)}{1 + \exp(\hat{\alpha} + \hat{\beta}^T \mathbf{x}_0)} \geq 0.5.$$

It is equivalent to $\hat{\alpha} + \hat{\beta}^T \mathbf{x}_0 \geq 0$, a linear decision boundary.

Penalized Logistic Regression

It is also popular to consider **penalized logistic regression**. Instead of maximizing the log-likelihood $\ell(\alpha, \beta)$, we minimize

$$\min_{\alpha, \beta} - \underbrace{\sum_{j=1}^n \{I(y_j = 1)(\alpha + \beta^T \mathbf{x}_j) - \log[1 + \exp(\alpha + \beta^T \mathbf{x}_j)]\}}_{-1 \times \ell(\alpha, \beta)} + \text{penalty term}.$$

Two popular penalty terms are

$$\begin{aligned} \text{LASSO :} \quad & \lambda \sum_{i=1}^p |\beta_i|, \\ \text{Ridge :} \quad & \lambda \beta^T \beta, \end{aligned}$$

where $\lambda \geq 0$ is a tuning parameter, often selected by cross validation or information criterion. The decision boundary remains linear.

Limitation

A limitation of classic logistic regression is that we need to specify the form of the covariates \mathbf{x} in $\boldsymbol{\beta}^T \mathbf{x}$ ourselves.

To overcome the limitation, we need to know how $\boldsymbol{\beta}$ is estimated. The gradient of the log-likelihood function is

$$\frac{\partial L(\alpha, \boldsymbol{\beta})}{\partial \beta_k} = \sum_{j=1}^n (y_j - \phi_j) x_{jk}.$$

If the [gradient descend](#) algorithm is used to update $\boldsymbol{\beta}$, the initial guess $\boldsymbol{\beta}^{(0)} = \mathbf{0}$ of $\boldsymbol{\beta}$ is improved by

$$\begin{aligned}\boldsymbol{\beta}^{(1)} &= \boldsymbol{\beta}^{(0)} - \sum_{j=1}^n \left(y_j - \phi_j^{(0)} \right) \mathbf{x}_j, \\ \boldsymbol{\beta}^{(2)} &= \boldsymbol{\beta}^{(1)} - \sum_{j=1}^n \left(y_j - \phi_j^{(1)} \right) \mathbf{x}_j = \boldsymbol{\beta}^{(0)} - \sum_{j=1}^n \left(\phi_j^{(1)} - \phi_j^{(0)} \right) \mathbf{x}_j,\end{aligned}$$

until no updates can be made.

Euclidean Inner Product

- Hence, we can expect the final estimate of $\boldsymbol{\beta}$ to be a linear combination of \mathbf{x}_j , say

$$\hat{\boldsymbol{\beta}} = \sum_{j=1}^n d_j \mathbf{x}_j,$$

for some coefficients d_j .

- We allocate \mathbf{x}_0 to π_1 if $\hat{\alpha} + \hat{\boldsymbol{\beta}}^T \mathbf{x}_0 \geq 0$. It is equivalent to

$$\hat{\alpha} + \sum_{j=1}^n d_j \mathbf{x}_j^T \mathbf{x}_0 \geq 0,$$

which depends on the Euclidean inner product $\mathbf{x}_j^T \mathbf{x}_0 = \langle \mathbf{x}_j, \mathbf{x}_0 \rangle$.

New Features

If we are not satisfied with the results using \mathbf{x} , we can create a vector of new features $\boldsymbol{\delta}(\mathbf{x}_j)$ and fit a new logistic model

$$\log\left(\frac{\phi_j}{1-\phi_j}\right) = \psi + \boldsymbol{\gamma}^T \boldsymbol{\delta}(\mathbf{x}_j).$$

Then, we allocate \mathbf{x}_0 to π_1 if

$$\hat{\psi} + \sum_{j=1}^n d_j^* \langle \boldsymbol{\delta}(\mathbf{x}_j), \boldsymbol{\delta}(\mathbf{x}_0) \rangle \geq 0,$$

for some coefficients d_j^* .

Kernel Function and Kernel Matrix

A function $\kappa(\mathbf{x}, \mathbf{z})$ is a **kernel function** if

- ① it is symmetric, $\kappa(\mathbf{x}, \mathbf{z}) = \kappa(\mathbf{z}, \mathbf{x})$,
- ② the **kernel matrix** \mathbf{K} with (i, j) th entry $\kappa(\mathbf{x}_i, \mathbf{x}_j)$ is positive semi-definite for all $\mathbf{x}_1, \dots, \mathbf{x}_n$.

If $\kappa(\mathbf{x}, \mathbf{z})$ is a kernel function, then you must be able to find a function $\delta(\cdot)$ such that

$$\kappa(\mathbf{x}, \mathbf{z}) = \delta^T(\mathbf{x}) \delta(\mathbf{z}).$$

As a kernel function, we will have an **eigen-decomposition**

$$\kappa(\mathbf{x}, \mathbf{z}) = \sum_{m=1}^{\infty} \rho_m e_m(\mathbf{x}) e_m(\mathbf{z}),$$

for some eigenvalues ρ_k and eigenfunctions $e_m(\mathbf{x})$.

Kernel Trick

Kernel trick is a commonly used trick to create new features from your original observed features.

- If the new features are created using $\delta(\mathbf{x})$, then we allocate \mathbf{x}_0 to π_1 if

$$\hat{\psi} + \sum_{j=1}^n d_j^* \langle \delta(\mathbf{x}_j), \delta(\mathbf{x}_0) \rangle \geq 0.$$

- If κ is the corresponding kernel function, the classifier is equivalent to

$$\hat{\psi} + \sum_{j=1}^n d_j^* \kappa(\mathbf{x}_0, \mathbf{x}_j) \geq 0.$$

- This means that we only need to choose the kernel function κ , if we need to create new features from the raw covariates \mathbf{x} .

Benefits of Kernel Trick

The eigen-decomposition implies that

$$\begin{aligned}
 \sum_{j=1}^n d_j^* \kappa(\mathbf{x}_0, \mathbf{x}_j) &= \sum_{j=1}^n d_j^* \sum_{m=1}^{\infty} \rho_m e_m(\mathbf{x}_0) e_m(\mathbf{x}_j) \\
 &= \sum_{m=1}^{\infty} \underbrace{\left(\sum_{j=1}^n d_j^* \sqrt{\rho_m} e_m(\mathbf{x}_j) \right)}_{\text{our new coefficients } d_j^*} \underbrace{\sqrt{\rho_m} e_m(\mathbf{x}_0)}_{\text{our new feature } \phi_m(\mathbf{x})}
 \end{aligned}$$

That is, we are using the vector of new features $\boldsymbol{\delta}(\mathbf{x})$, possibly of infinite dimension, to model ϕ_j .

Reproducing Kernel Hilbert Space Regression

- The kernel trick means that we first choose a kernel function, compute the kernel matrix \mathbf{K} using the data matrix \mathbf{X} , and use \mathbf{K} as the data matrix to fit our logistic regression model.
- However, when the data matrix \mathbf{X} is of dimension $n \times p$, then the kernel matrix \mathbf{K} is $n \times n$ and we have n regression coefficients to estimate. We often cannot obtain a unique estimator.
- Hence, we often use **penalized maximum likelihood**: the coefficients can be estimated by minimizing

$$-1 \times \log\text{-likelihood} + \lambda \boldsymbol{\beta}^T \mathbf{K} \boldsymbol{\beta},$$

where λ is the **tuning parameter**. It is often chosen by cross validation.

Classify New Observation

- Suppose that $\kappa(\mathbf{x}, \mathbf{z}) = \boldsymbol{\delta}^T(\mathbf{x}) \boldsymbol{\delta}(\mathbf{z})$. Then, it is equivalent to transforming our original data \mathbf{x} to new data $\boldsymbol{\delta}(\mathbf{x})$.
- Our model is then

$$P(Y_j = 1 \mid \mathbf{x}_j) = \frac{\exp[\psi + \boldsymbol{\gamma}^T \boldsymbol{\delta}(\mathbf{x}_j)]}{1 + \exp[\psi + \boldsymbol{\gamma}^T \boldsymbol{\delta}(\mathbf{x}_j)]},$$

- We classify a new observation \mathbf{x}_0 to $y_0 = 1$ if

$$\begin{aligned} 0 &\leq \hat{\boldsymbol{\gamma}}^T \boldsymbol{\delta}(\mathbf{x}_0) + \hat{\psi} \\ &= \left[\sum_{j=1}^n d_j^* \boldsymbol{\delta}(\mathbf{x}_j) \right]^T \boldsymbol{\delta}(\mathbf{x}_0) + \hat{\psi} \\ &= \sum_{j=1}^n d_j^* [\boldsymbol{\delta}^T(\mathbf{x}_j) \boldsymbol{\delta}(\mathbf{x}_0)] + \hat{\psi} \\ &= \sum_{j=1}^n d_j^* \kappa(\mathbf{x}_j, \mathbf{x}_0) + \hat{\psi}. \end{aligned}$$

Motivation: Margin

Suppose that we code Y as $\{-1, 1\}$. We define a linear classifier as

$$\hat{y} = \begin{cases} 1, & \text{if } \mathbf{w}^T \mathbf{x} + b \geq 0, \\ -1, & \text{if } \mathbf{w}^T \mathbf{x} + b < 0. \end{cases}$$

The hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ is our decision boundary. The distance of a point \mathbf{x}_0 to the decision boundary is

$$\frac{|\mathbf{w}^T \mathbf{x}_0 + b|}{\|\mathbf{w}\|} = \frac{y_0 (\mathbf{w}^T \mathbf{x}_0 + b)}{\|\mathbf{w}\|}.$$

We can choose the decision boundary such that the distance ([margin](#)) is maximized.

Support Vector: Brief Intro

We want to maximize the margin such that we can correct classify all data points:

$$\max_{w,b} \left[\min_j \frac{y_j (\mathbf{w}^T \mathbf{x}_j + b)}{\|\mathbf{w}\|} \right], \quad \text{s.t.} \quad y_j (\mathbf{w}^T \mathbf{x}_j + b) \geq 0, \quad \forall j.$$

It is equivalent to

$$\min_{w,b} \mathbf{w}^T \mathbf{w}, \quad \text{s.t.} \quad y_j (\mathbf{w}^T \mathbf{x}_j + b) \geq 1, \quad \forall j.$$

Some data points must satisfy $y_j (\mathbf{w}^T \mathbf{x}_j + b) = 1$. These points are called **support vectors**.

Hinge Loss: Brief Intro

In practice, we often cannot correctly classify all points. Hence, we change our optimization problem and allow some small errors:

$$\min_{\mathbf{w}, b, \xi_j} \mathbf{w}^T \mathbf{w} + C \sum_{j=1}^n \xi_j, \quad \text{s.t.} \quad y_j (\mathbf{w}^T \mathbf{x}_j + b) \geq 1 - \xi_j, \quad \xi_j \geq 0, \quad \forall j,$$

where $\{\xi_j\}$ are known as the **slack variables**, and C is a tuning parameter.

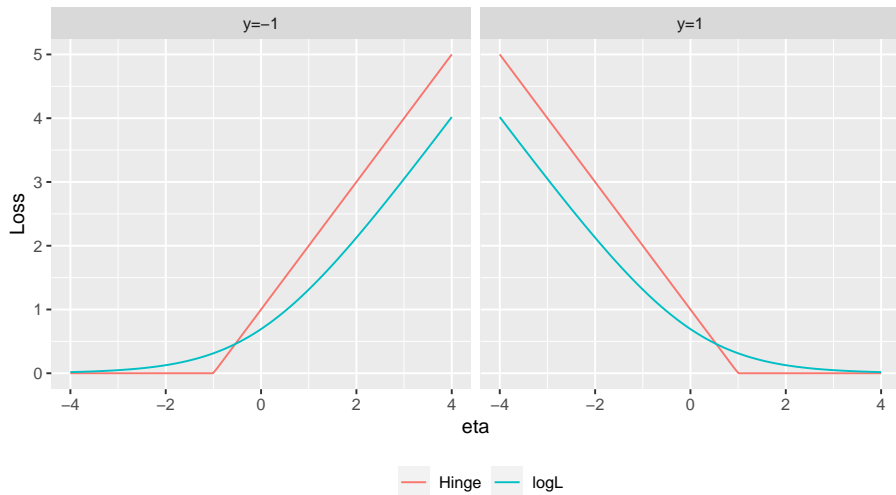
It is equivalent to

$$\min_{\mathbf{w}, b} \underbrace{\sum_{j=1}^n \max \{0, 1 - y_j (\mathbf{w}^T \mathbf{x}_j + b)\}}_{\text{Hinge loss}} + C^{-1} \mathbf{w}^T \mathbf{w}.$$

In contrast, the penalized logistic regression uses the log-likelihood loss

$$\min_{\alpha, \beta} - \sum_{j=1}^n \{ I(y_j = 1) (\alpha + \beta^T \mathbf{x}_j) - \log [1 + \exp (\alpha + \beta^T \mathbf{x}_j)] \} + \lambda \beta^T \beta.$$

Hinge Loss Vs Log-Likelihood Loss



Kernel Trick

The solution of

$$\min_{w,b} \sum_{j=1}^n \max \{0, 1 - y_j (\mathbf{w}^T \mathbf{x}_j + b)\} + C^{-1} \mathbf{w}^T \mathbf{w}$$

is of the form

$$\hat{\mathbf{w}} = \sum_{j=1}^n d_j y_j \mathbf{x}_j$$

for some coefficients d_j . We classify a new observation \mathbf{x}_0 to $y_0 = 1$ if

$$0 \leq \hat{\mathbf{w}}^T \mathbf{x}_0 + \hat{b} = \sum_{j=1}^n d_j y_j (\mathbf{x}_j^T \mathbf{x}_0) + \hat{b}.$$

Hence, the [support vector machine](#) often uses the kernel trick that replaces $\mathbf{x}_j^T \mathbf{x}_0$ by some kernel $\kappa(\mathbf{x}_j, \mathbf{x}_0)$.

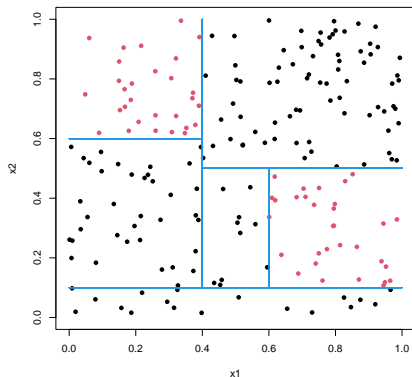
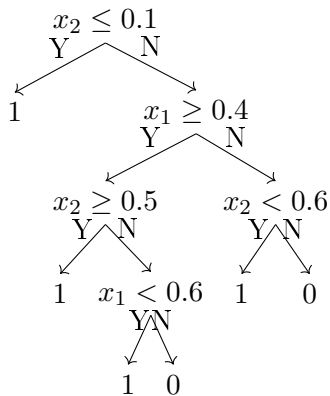
Support Vector Machine

- The kernel trick means that we first choose a kernel function, compute the kernel matrix \mathbf{K} using the data matrix \mathbf{X} , and use \mathbf{K} as the data matrix to train our support vector machine.
- Suppose that $\kappa(\mathbf{x}, \mathbf{z}) = \boldsymbol{\delta}^T(\mathbf{x}) \boldsymbol{\delta}(\mathbf{z})$. Then, it is equivalent to transforming our original data \mathbf{x} to new data $\boldsymbol{\delta}(\mathbf{x})$.
- We classify a new observation \mathbf{x}_0 to $y_0 = 1$ if

$$\begin{aligned}
 0 &\leq \hat{\mathbf{w}}^T \boldsymbol{\delta}(\mathbf{x}_0) + \hat{b} \\
 &= \left[\sum_{j=1}^n d_j^* y_j \boldsymbol{\delta}(\mathbf{x}_j) \right]^T \boldsymbol{\delta}(\mathbf{x}_0) + \hat{b} \\
 &= \sum_{j=1}^n d_j^* y_j [\boldsymbol{\delta}^T(\mathbf{x}_j) \boldsymbol{\delta}(\mathbf{x}_0)] + \hat{b} \\
 &= \sum_{j=1}^n d_j^* y_j \kappa(\mathbf{x}_j, \mathbf{x}_0) + \hat{b}.
 \end{aligned}$$

Classification Tree

A classification tree is a **hierarchical method** of partitioning the sample space. Every partition corresponds to a class.



Best Partition

The main idea of a tree method is that, at each step, we pick a binary split that best partitions the data. For classification trees, we often use the [Gini index](#)

$$G = \sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk}),$$

or the [cross entropy](#)

$$G = - \sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk},$$

to measure what is the best partition, where \hat{p}_{mk} is the proportion of observations in the m th region that are from the k th class.

- If \hat{p}_{mk} are close to 0 or 1, then G is small.
- We want to minimize G .

How To Grow A Tree

We still need to decide

- ① **tree depth**: when to stop split the tree,
- ② **leaf size**: number of observations in a leaf (node).

A tree can **overfit** the data. Hence, we can

- grow a large tree until each terminal node has few observations or the tree is tall enough.
- use cross validation to decide when we can stop growing the tree.

A Tree Versus A Forest

A **random forest** grows multiple trees and combine them together to make the final decision.

Algorithm 1: Random forest

- 1 Specify an integer B ;
 - 2 **for** *each integer i from 1 to B* **do**
 - 3 Sample with replacement a subset of the original data ;
 - 4 Sample without replacement a subset of covariates/features ;
 - 5 Grow a classification tree ;
 - 6 **end**
 - 7 Compute the proportion of trees that classify an object to a class.
-

Other Methods

There are many more classification methods such as

- Boosting methods,
- XGboost (state-of-the art),
- Neural networks.

Imbalanced Groups

Suppose that we have observed

Yes	No
100	9900

We can simply say we assign all classes to “No”, then the accuracy is 99%.

- ① **Downsampling**: We sample a subset from “No”.
- ② **Oversampling**: “duplicate” the group “Yes” until the proportions are approximately equal.
- ③ Weight the observations.

ECM for Several Populations

It is possible to generalize the above methods from two populations to several populations.

The ECM for two classes is

$$\text{ECM} = c(2 | 1) P(2 | 1) p_1 + c(1 | 2) P(1 | 2) p_2.$$

Suppose that we have more than two groups. The ECM of misclassifying π_i to a wrong class is

$$\text{ECM}(i) = \sum_{k \neq i} c(k | i) P(k | i) p_i.$$

The overall ECM for g classes is

$$\text{ECM} = \sum_{i=1}^g \left[\sum_{k \neq i} c(k | i) P(k | i) \right] p_i.$$

Minimizing ECM

Result 11.5: Min ECM for Several Classes

The classification regions that minimize the overall ECM are by allocating \mathbf{x}_0 to class π_k , $k = 1, \dots, g$, for which $\sum_{k \neq i} p_i f_i(\mathbf{x}_0) c(k | i)$ is smallest.

- If all $p_i c(k | i)$ are equal, then it becomes the highest likelihood classifier: we allocate \mathbf{x}_0 to π_k , if $f_k(\mathbf{x}_0) > f_i(\mathbf{x}_0)$, $\forall i \neq k$.
- If all misclassification costs are equal, then it becomes the minimum TPM classifier: we allocate \mathbf{x}_0 to π_k , if

$$p_k f_k(\mathbf{x}_0) > p_i f_i(\mathbf{x}_0), \quad \forall i \neq k.$$

- The minimum TPM classifier is identical to the [naive Bayes](#) classifier that maximizes the posterior,

$$P(\pi_k | \mathbf{x}_0) = \frac{f(\mathbf{x}_0 | \pi_k) P(\pi_k)}{P(\mathbf{x}_0)} = \frac{f_k(\mathbf{x}_0) p_k}{P(\mathbf{x}_0)} \propto f_k(\mathbf{x}_0) p_k.$$

Example: Classification with Normal Populations

Suppose that $f_i(\mathbf{x})$ is the density function of $N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$. The naive Bayes classifier allocates \mathbf{x} to π_k if

$$\log [p_k f_k(\mathbf{x})] = \max_i \log [p_i f_i(\mathbf{x})].$$

Here,

$$\begin{aligned} \log [p_k f_k(\mathbf{x})] &= \log p_k - \frac{p}{2} \log (2\pi) - \frac{1}{2} \log \det (\boldsymbol{\Sigma}) \\ &\quad - \frac{1}{2} (\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k - \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k). \end{aligned}$$

We call

$$d_i(\mathbf{x}) = \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \log p_i$$

the **linear discrimination score**. We allocate \mathbf{x} to π_k , if $d_k(\mathbf{x})$ is the largest of $d_1(\mathbf{x}), \dots, d_g(\mathbf{x})$. If we do not know $\boldsymbol{\mu}_i$ and $\boldsymbol{\Sigma}$, they are estimated by $\bar{\mathbf{x}}_i$ and $\mathbf{S}_{\text{pooled}}$, respectively.

Example: Classification with Normal Populations

Suppose that $f_i(\mathbf{x})$ is the density function of $N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$. The naive Bayes classifier allocates \mathbf{x} to π_k if

$$\log [p_k f_k(\mathbf{x})] = \max_i \log [p_i f_i(\mathbf{x})].$$

Here,

$$\begin{aligned} \log [p_k f_k(\mathbf{x})] &= \log p_k - \frac{p}{2} \log (2\pi) - \frac{1}{2} \log \det (\boldsymbol{\Sigma}_k) \\ &\quad - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k). \end{aligned}$$

We call

$$d_i^Q(\mathbf{x}) = -\frac{1}{2} \log \det (\boldsymbol{\Sigma}_i) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \log p_i$$

the **quadratic discrimination score**. We allocate \mathbf{x} to π_k , if $d_k^Q(\mathbf{x})$ is the largest of $d_1^Q(\mathbf{x}), \dots, d_g^Q(\mathbf{x})$. If we do not know $\boldsymbol{\mu}_i$ and $\boldsymbol{\Sigma}_i$, they are estimated by $\bar{\mathbf{x}}_i$ and \mathbf{S}_i , respectively.

Fisher's Discriminant Analysis

Fisher's discriminant analysis can be generalized to more than two groups. We still assume $\Sigma_1 = \cdots = \Sigma_g = \Sigma$.

Similar to LDA for two classes, we want to maximize

$$\frac{\mathbf{a}^T \mathbf{B} \mathbf{a}}{\mathbf{a}^T \mathbf{W} \mathbf{a}},$$

where

$$\begin{aligned}\mathbf{W} &= \sum_{\ell=1}^g \sum_{j=1}^{n_{\ell}} (\mathbf{x}_{\ell j} - \bar{\mathbf{x}}_{\ell}) (\mathbf{x}_{\ell j} - \bar{\mathbf{x}}_{\ell})^T, \\ \mathbf{B} &= \sum_{\ell=1}^g n_{\ell} (\bar{\mathbf{x}}_{\ell} - \bar{\mathbf{x}}) (\bar{\mathbf{x}}_{\ell} - \bar{\mathbf{x}})^T.\end{aligned}$$

The maximizer \mathbf{a} is the eigenvector of $\mathbf{W}^{-1} \mathbf{B}$ corresponding to the largest eigenvalue of $\mathbf{W}^{-1} \mathbf{B}$.

Fisher's Linear Discriminants

If we have more than two classes, using only one linear combination is often not enough. We often need to use $r > 1$ linear combinations. Let $(\lambda_i, \mathbf{e}_i)$ be the eigenvalue and eigenvector of $\mathbf{W}^{-1}\mathbf{B}$, where $\lambda_1 \geq \lambda_2 \geq \cdots \geq 0$.

- ① Fisher's first linear discriminant is $\mathbf{e}_1^T \mathbf{x}_0$,
- ② Fisher's second linear discriminant is $\mathbf{e}_2^T \mathbf{x}_0$,
- ③ Fisher's k th linear discriminant is $\mathbf{e}_k^T \mathbf{x}_0$.

We allocate \mathbf{x}_0 to π_k if the linear combinations $\{\mathbf{e}_j^T \mathbf{x}_0\}$ are closer to the mean of the linear combinations of the k th class. That is,

$$\sum_{j=1}^r (\mathbf{e}_j^T \mathbf{x}_0 - \mathbf{e}_j^T \bar{\mathbf{x}}_k)^2 \leq \sum_{j=1}^r (\mathbf{e}_j^T \mathbf{x}_0 - \mathbf{e}_j^T \bar{\mathbf{x}}_i)^2, \quad \forall i \neq k.$$

Linear Discriminants

If r equal to the rank of $\mathbf{W}^{-1}\mathbf{B}$, then it can be shown that

$$\begin{aligned}
 & \sum_{j=1}^r (\mathbf{e}_j^T \mathbf{x}_0 - \mathbf{e}_j^T \bar{\mathbf{x}}_i)^2 \\
 &= (\mathbf{x}_0 - \bar{\mathbf{x}}_i)^T \mathbf{W}^{-1} (\mathbf{x}_0 - \bar{\mathbf{x}}_i) \\
 &= -2 \left[\bar{\mathbf{x}}_i^T \mathbf{W}^{-1} \mathbf{x}_0 - \frac{1}{2} \bar{\mathbf{x}}_i^T \mathbf{W}^{-1} \bar{\mathbf{x}}_i + \log(p_i) \right] + \mathbf{x}_0^T \mathbf{W}^{-1} \mathbf{x}_0 + 2 \log(p_i).
 \end{aligned}$$

Let

$$d_i(\mathbf{x}) = \bar{\mathbf{x}}_i^T \mathbf{W}^{-1} \mathbf{x} - \frac{1}{2} \bar{\mathbf{x}}_i^T \mathbf{W}^{-1} \bar{\mathbf{x}}_i + \log(p_i).$$

It is the same as allocating \mathbf{x}_0 to π_k if $d_k(\mathbf{x})$ is the largest of $d_1(\mathbf{x}), \dots, d_g(\mathbf{x})$.