

Lecture 12

Discrimination and Classification

Readings: Zelterman, 2015, Chapters 10

DSA 8070 Multivariate Analysis

November 1 - November 5, 2021

Whitney Huang
Clemson University



Notes

Agenda

1 Overview

2 Binary Linear Classification



Notes

Discrimination and Classification

Main objectives behind discrimination and classification are:

- **Discrimination**: separate distinct "populations" of observations
- **Classification**: classify new objectives into pre-defined populations

Examples

- Given measurements on the concentrations of five elements in bullet lead, find combinations of those concentrations that best describe bullets made by Cascade, Federal, Winchester and Remington ⇒ **discrimination**
- Using information on prisoners eligible for parole (good behavior, history of drug use, job skills, etc) can we successfully allocate a prisoners eligible for parole into two groups: those who will commit another crime or those who will not commit another crime? ⇒ **classification**



Notes

Classification

- **Data:**
 $\{X_i, Y_i\}_{i=1}^n$,
where Y_i is the class information for the i_{th} observation $\Rightarrow Y$ is a qualitative variable
- **Classification** aims to classify a new observation (or several new observations) into one of those classes

Quantity of interest: $P(Y = k_{th} \text{ category} | X = x)$
- In this lecture we will focus on **binary linear classification**

Discrimination and Classification

CLEMSON UNIVERSITY

Overview

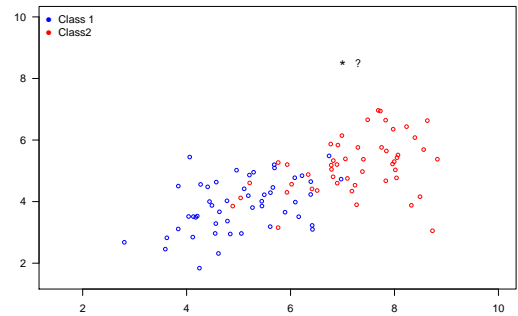
Binary Linear Classification

12.4

Notes

Toy Example

Wish to classify a new observation $x_i = (x_{1i}, x_{2i})$, denoted by (*), into one of the two groups (**class 1** or **class 2**)



Discrimination and Classification

CLEMSON UNIVERSITY

Overview

Binary Linear Classification

12.5

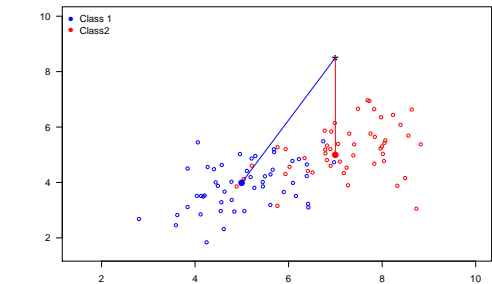
Notes

Toy Example Cont'd

We can compute the distances from this new observation $x = (x_1, x_2)$ to the groups, for example,

$$d_1 = \sqrt{(x_1 - \mu_{11})^2 + (x_2 - \mu_{12})^2},$$
$$d_2 = \sqrt{(x_1 - \mu_{21})^2 + (x_2 - \mu_{22})^2}.$$

We can assign x to the group with the smallest distance



Discrimination and Classification

CLEMSON UNIVERSITY

Overview

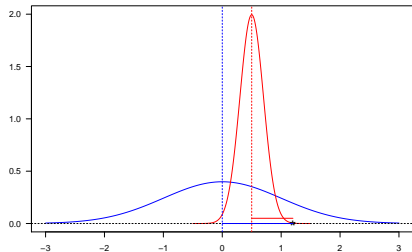
Binary Linear Classification

12.6

Notes

Variance Corrected Distance

In this one-dimensional example, $d_1 = |x - \mu_1| > |x - \mu_2|$. Does that mean x is "closer" to group 2 (red) than group 1 (blue)?



We should take the "spread" of each group into account. $\hat{d}_1 = |x - \mu_1|/\sigma_1 < \hat{d}_2 = |x - \mu_2|/\sigma_2$

Notes

General Covariance Adjusted Distance: Mahalanobis Distance

The Mahalanobis distance [Mahalanobis, 1936] is a measure of the distance between a point x and a multivariate distribution of X :

$$D_M(x) = \sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)},$$

where μ is the mean vector and Σ is the variance-covariance matrix of X

One can use the Mahalanobis distance, by computing the Mahalanobis distance between an observations x_i and the "center" of the k_{th} population μ_k , to carry out classification

Notes

Binary Classification with Multivariate Normal Populations

Assume $X_1 \sim \text{MVN}(\mu_1, \Sigma)$, $X_2 \sim \text{MVN}(\mu_2, \Sigma)$, that is, $\Sigma_1 = \Sigma_2 = \Sigma$

- Maximum Likelihood of group membership:

Group 1 if $\ell(x, \mu_1, \Sigma) > \ell(x, \mu_2, \Sigma)$

- Linear Discriminant Function:

Group 1 if $(\mu_1 - \mu_2)^T \Sigma^{-1} x - \frac{1}{2}(\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 + \mu_2) > 0$

- Minimize Mahalanobis distance:

Group 1 if $(x - \mu_1)^T \Sigma^{-1} (x - \mu_1) < (x - \mu_2)^T \Sigma^{-1} (x - \mu_2)$

All the methods above are equivalent

Notes

Priors and Misclassification Costs

In addition to the observed characteristics of units $\{x_i\}_{i=1}^n$, other features of classification rules are:

- **Prior probability:**

If one population is more prevalent than the other, chances are higher that a new unit came from the larger population. Stronger evidence would be needed to allocate the unit to the population with the smaller prior probability.

- **Costs of misclassification:**

It may be more costly to misclassify a seriously ill subject as healthy than to misclassify a healthy subject as being ill.

Discrimination and Classification
 CLEMSON UNIVERSITY
 Overview
 Binary Linear Classification

12.10

Notes

Classification Regions

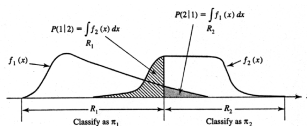
Ignore for now the prior probabilities of each population and the potentially different misclassification costs.

- The probability of misclassifying an object into π_2 when it belongs in π_1 is

$$P(2|1) = \mathbb{P}(X \in \mathcal{R}_2 | \pi_1)$$

- The probability of misclassifying an object into π_1 when it belongs in π_2 is

$$P(1|2) = \mathbb{P}(X \in \mathcal{R}_1 | \pi_2)$$



Source: Figure 11.3 from Applied Multivariate Statistical Analysis, 6th Ed (Johnson & Wichern). Visualization is for $p = 1$ variable.

Discrimination and Classification
 CLEMSON UNIVERSITY
 Overview
 Binary Linear Classification

12.11

Notes

Probability and Expected Cost of Misclassification

Let p_1 and p_2 denote the prior probabilities of π_1 , π_2 , and $c(1|2)$, $c(2|1)$ be the costs of misclassification:

- Then probabilities of the four possible outcomes are:

$$\begin{aligned} \mathbb{P}(\text{correctly classified as } \pi_1) &= \mathbb{P}(X \in \mathcal{R}_1 | \pi_1) \mathbb{P}(\pi_1) = P(1|1)p_1 \\ \mathbb{P}(\text{incorrectly classified as } \pi_1) &= \mathbb{P}(X \in \mathcal{R}_1 | \pi_2) \mathbb{P}(\pi_2) = P(1|2)p_2 \\ \mathbb{P}(\text{correctly classified as } \pi_2) &= \mathbb{P}(X \in \mathcal{R}_2 | \pi_2) \mathbb{P}(\pi_2) = P(2|2)p_2 \\ \mathbb{P}(\text{incorrectly classified as } \pi_2) &= \mathbb{P}(X \in \mathcal{R}_2 | \pi_1) \mathbb{P}(\pi_1) = P(2|1)p_1 \end{aligned}$$

- Classification rules are often evaluated in terms of the **expected cost of misclassification (ECM)**:

$$ECM = c(2|1)p(2|1)p_1 + c(1|2)P(1|2)P(1|2)p_2,$$

and we seek rules that **minimize the ECM**

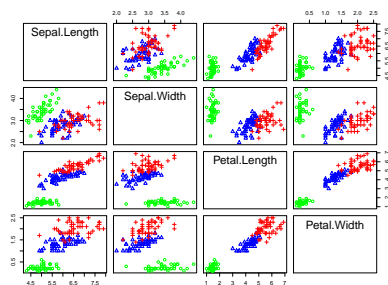
Discrimination and Classification
 CLEMSON UNIVERSITY
 Overview
 Binary Linear Classification

12.12

Notes

Example: Fisher's Iris Data

4 variables (sepal length and width and petal length and width), 3 species (setosa, versicolor, and virginica)

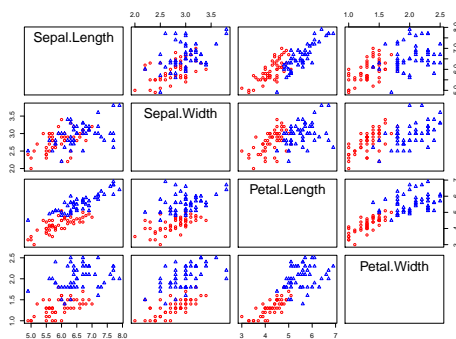


Task: Classify flowers into different species based on lengths and widths of sepal and petal

Notes

Fisher's Iris Data Cont'd

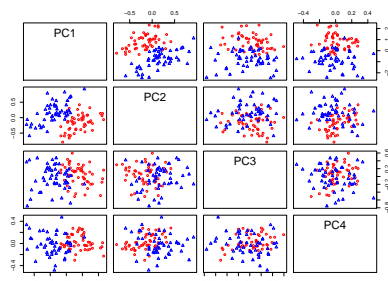
Let's focus on the latter two classes (versicolor, and virginica)



Notes

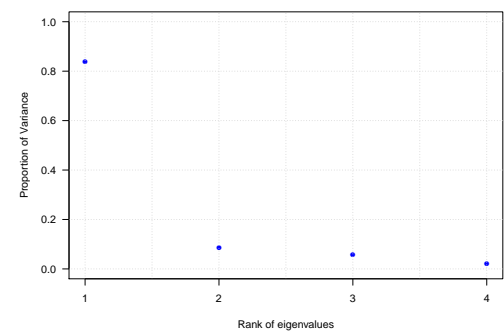
Fisher's iris Data Cont'd

To further simplify the matter, let's focus on the first two PCs of X



Notes

Screen Plot



Discrimination and Classification

CLEMSON UNIVERSITY

Overview

Binary Linear Classification

12.16

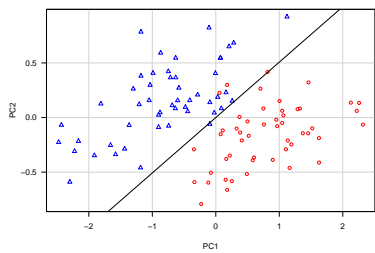
Notes

Linear Discriminant Analysis

Main idea: Use Bayes rule to compute

$$P(Y = k | \mathbf{X} = \mathbf{x}) = \frac{P(Y = k)P(\mathbf{X} = \mathbf{x} | Y = k)}{P(\mathbf{X} = \mathbf{x})} = \frac{\pi_k f_k(\mathbf{x})}{\sum_{k=1}^K \pi_k f_k(\mathbf{x})}$$

Assuming $f_k(\mathbf{x}) \sim \text{MVN}(\boldsymbol{\mu}_k, \Sigma)$, $k = 1, \dots, K$ and use $\hat{\pi}_k = \frac{n_k}{n} \Rightarrow$ it turns out the resulting classifier is linear in \mathbf{X}



Discrimination and Classification

CLEMSON UNIVERSITY

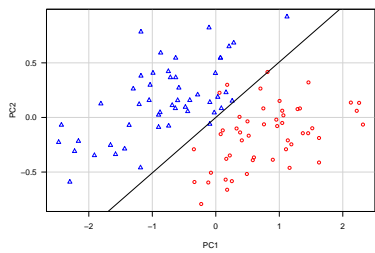
Overview

Binary Linear Classification

12.17

Notes

Classification Performance Evaluation



fit.LDA

	versicolor	virginica
versicolor	47	3
virginica	1	49

Discrimination and Classification

CLEMSON UNIVERSITY

Overview

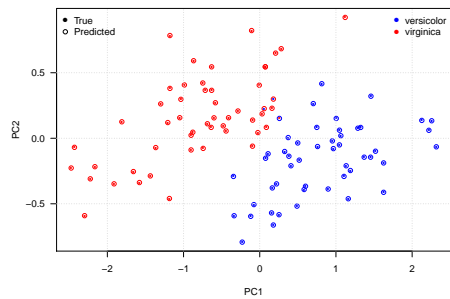
Binary Linear Classification

12.18

Notes

Logistic Regression Classifier

Main idea: Model the logit $\log\left(\frac{P(Y=1)}{1-P(Y=1)}\right)$ as a linear function in x



Discrimination and Classification

CLEMSON UNIVERSITY

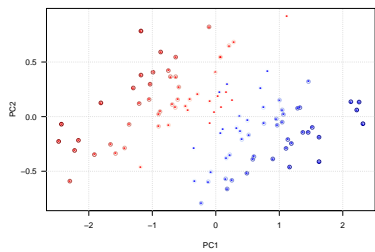
Overview

Binary Linear Classification

12.19

Notes

Logistic Regression Classifier Cont'd



logisticPred		
	versicolor	virginica
versicolor	48	2
virginica	1	49

Discrimination and Classification

CLEMSON UNIVERSITY

Overview

Binary Linear Classification

12.20

Notes

Linear Discriminant Analysis Versus Logistic Regression

For a binary classification problem, one can show that both linear discriminant analysis (LDA) and logistic regression are linear classifiers. The difference is in how the parameters are estimated:

- Logistic regression uses the conditional likelihood based on $P(Y|X = x)$
- LDA uses the full likelihood based on multivariate normal assumption on X
- Despite these differences, in practice the results are often very similar

Discrimination and Classification

CLEMSON UNIVERSITY

Overview

Binary Linear Classification

12.21

Notes

Quadratic Discriminant Analysis

In linear discriminant analysis, we **assume** $\{f_k(\mathbf{x})\}_{k=1}^K$ are normal densities and $\Sigma_1 = \Sigma_2$, therefore we obtain a linear classifier. What if $\Sigma_1 \neq \Sigma_2 \Rightarrow$ we get quadratic discriminant analysis

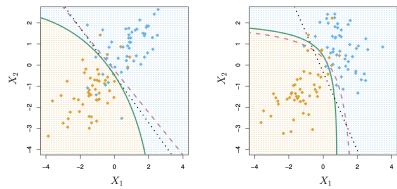


Figure: Figure courtesy of [An Introduction of Statistical Learning](#) by G. James et al. pp. 150

Discrimination and Classification

CLEMSON UNIVERSITY

Overview

Binary Linear Classification

12.22

Notes

Notes

Notes
