Big Data Analysis (MA60306)

Bibhas Adhikari

Spring 2022-23, IIT Kharagpur

Lecture 15 March 3, 2023

Bibhas Adhikari (Spring 2022-23, IIT Kharag

Big Data Analysis

Lecture 15 March 3, 2023 1/8

3

< □ > < □ > < □ > < □ > < □ > < □ >

Observation The set of all vectors \mathbf{x} such that $L(\mathbf{x}) = 0$ defines a hyperplane which divides the classes!! The term $\mathbf{b}^T \mathbf{x}$ in $L(\mathbf{x})$ is called the Fisher's linear discriminant function (LDF).

▲ □ ▶ ▲ □ ▶ ▲ □ ▶

Observation The set of all vectors \mathbf{x} such that $L(\mathbf{x}) = 0$ defines a hyperplane which divides the classes!! The term $\mathbf{b}^T \mathbf{x}$ in $L(\mathbf{x})$ is called the Fisher's linear discriminant function (LDF).

Let us denote the partition of the feature space \mathbb{R}^d into two regions R_1 and R_2 . If **x** falls into R_1 then it is classified as belonging to Π_1 , whereas if **x** falls into region R_2 , it is classified into Π_2 .

Observation The set of all vectors \mathbf{x} such that $L(\mathbf{x}) = 0$ defines a hyperplane which divides the classes!! The term $\mathbf{b}^T \mathbf{x}$ in $L(\mathbf{x})$ is called the Fisher's linear discriminant function (LDF).

Let us denote the partition of the feature space \mathbb{R}^d into two regions R_1 and R_2 . If **x** falls into R_1 then it is classified as belonging to Π_1 , whereas if **x** falls into region R_2 , it is classified into Π_2 .

Total misclassification probability Misclassification happens if \mathbf{x} is assigned to Π_2 , but actually it belongs to Π_1 or if \mathbf{x} is assigned to Π_1 , but actually it belongs to Π_2 .

< □ > < 同 > < 回 > < 回 > < 回 >

Squared Mahalanobis distance The squared Mahalanobis distance between Π_1 and Π_2 is defined as

$$riangle^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

3

(日) (四) (日) (日) (日)

Squared Mahalanobis distance The squared Mahalanobis distance between Π_1 and Π_2 is defined as

$$\triangle^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

Question What does Mahalanobis distance measure?

э

A B A A B A

Image: A matrix

Squared Mahalanobis distance The squared Mahalanobis distance between Π_1 and Π_2 is defined as

$$\triangle^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

Question What does Mahalanobis distance measure? Recall Let X a random matrix. For the random matrix $Y = AXB^T + C$, where A, B, C are compatible matrices, $\mathbb{E}(Y) = A\mathbb{E}(X)B^T$ and the covariance matrix of vec(Y) is $\Sigma_{YY} = (A \otimes B)\Sigma_{XX}(A \otimes B)^T$

Squared Mahalanobis distance The squared Mahalanobis distance between Π_1 and Π_2 is defined as

$$\triangle^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

Question What does Mahalanobis distance measure? Recall Let X a random matrix. For the random matrix $Y = AXB^T + C$, where A, B, C are compatible matrices, $\mathbb{E}(Y) = A\mathbb{E}(X)B^T$ and the covariance matrix of vec(Y) is $\Sigma_{YY} = (A \otimes B)\Sigma_{XX}(A \otimes B)^T$ Set $U = \mathbf{b}^T \mathbf{x}$. Then

$$\mathbb{E}(U|x \in \Pi_i) = \mathbf{b}^T \mu_i = (\mu_1 - \mu_2)^T \Sigma^{-1} \mu_i$$

$$Var(U|x \in \Pi_i) = \mathbf{b}^T \Sigma \mathbf{b} = (\mu_1 - \mu_2)^T \Sigma^{-1} \Sigma \Sigma^{-1} (\mu_1 - \mu_2) = \triangle^2.$$

Squared Mahalanobis distance The squared Mahalanobis distance between Π_1 and Π_2 is defined as

$$\triangle^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

Question What does Mahalanobis distance measure? Recall Let X a random matrix. For the random matrix $Y = AXB^T + C$, where A, B, C are compatible matrices, $\mathbb{E}(Y) = A\mathbb{E}(X)B^T$ and the covariance matrix of vec(Y) is $\Sigma_{YY} = (A \otimes B)\Sigma_{XX}(A \otimes B)^T$ Set $U = \mathbf{b}^T \mathbf{x}$. Then

$$\mathbb{E}(U|x \in \Pi_i) = \mathbf{b}^T \mu_i = (\mu_1 - \mu_2)^T \Sigma^{-1} \mu_i$$

$$Var(U|x \in \Pi_i) = \mathbf{b}^T \Sigma \mathbf{b} = (\mu_1 - \mu_2)^T \Sigma^{-1} \Sigma \Sigma^{-1} (\mu_1 - \mu_2) = \triangle^2.$$

Then the total misclassification probability is:

$$P(\mathbf{x} \in R_2 | \mathbf{x} \in \Pi_1) \pi_1 + P(\mathbf{x} \in R_1 | \mathbf{x} \in \Pi_2) \pi_2$$

LDA Now

$$P(\mathbf{x} \in R_2 | \mathbf{x} \in \Pi_1) = P(L(\mathbf{x}) < 0 | \mathbf{x} \in \Pi_1).$$

Bibhas Adhikari (Spring 2022-23, IIT Kharag

◆□▶ ◆圖▶ ◆臣▶ ◆臣▶ ─ 臣

Now

$$P(\mathbf{x} \in R_2 | \mathbf{x} \in \Pi_1) = P(L(\mathbf{x}) < 0 | \mathbf{x} \in \Pi_1).$$

Note that $L(\mathbf{x}) = b_0 + U$. Also,

$$Z = \frac{U - \mathbb{E}(U | \mathbf{x} \in \Pi_i)}{\sqrt{var(U | \mathbf{x} \in \Pi_i)}} \sim \mathcal{N}(0, 1)$$

æ

<ロト <問ト < 目と < 目と

Now

$$P(\mathbf{x} \in R_2 | \mathbf{x} \in \Pi_1) = P(L(\mathbf{x}) < 0 | \mathbf{x} \in \Pi_1).$$

Note that $L(\mathbf{x}) = b_0 + U$. Also,

$$Z = \frac{U - \mathbb{E}(U | \mathbf{x} \in \Pi_i)}{\sqrt{var(U | \mathbf{x} \in \Pi_i)}} \sim \mathcal{N}(0, 1)$$

Then from using the expressions for $\mathbb{E}(U|x \in \Pi_i)$, $Var(U|x \in \Pi_i)$, and b_0 as derived above,

$$P(L(\mathbf{x}) < 0 | \mathbf{x} \in \Pi_1) = P(U < b_0 | \mathbf{x} \in \Pi_1)$$

$$= P(Z < \frac{b_0 - (\mu_1 - \mu_2)^T \Sigma^{-1} \mu_i}{\Delta})$$

$$= P(Z < -\frac{\Delta}{2} - \frac{1}{\Delta} \log_e \frac{\pi_2}{\pi_1})$$

$$= \Phi\left(-\frac{\Delta}{2} - \frac{1}{\Delta} \log_e \frac{\pi_2}{\pi_1}\right)$$

Similarly we can obtain

$$P(\mathbf{x} \in R_1 | \mathbf{x} \in \Pi_2) = P(L(\mathbf{x}) > 0 | \mathbf{x} \in \Pi_2)$$

= $P(Z < -\frac{\Delta}{2} - \frac{1}{\Delta} \log_e \frac{\pi_2}{\pi_1})$
= $\Phi\left(-\frac{\Delta}{2} + \frac{1}{\Delta} \log_e \frac{\pi_2}{\pi_1}\right)$

3

<ロト <問ト < 目と < 目と

Similarly we can obtain

$$P(\mathbf{x} \in R_1 | \mathbf{x} \in \Pi_2) = P(L(\mathbf{x}) > 0 | \mathbf{x} \in \Pi_2)$$

= $P(Z < -\frac{\Delta}{2} - \frac{1}{\Delta} \log_e \frac{\pi_2}{\pi_1})$
= $\Phi\left(-\frac{\Delta}{2} + \frac{1}{\Delta} \log_e \frac{\pi_2}{\pi_1}\right)$

If $\pi_1 = \pi_2 = 1/2$ then

 $P(\mathbf{X} \in R_2 | \mathbf{X} \in \Pi_1) = P(\mathbf{X} \in R_1 | \mathbf{X} \in \Pi_2) = \Phi(-\triangle/2)$

Bibhas Adhikari (Spring 2022-23, IIT Kharag

イロト 不得 トイヨト イヨト 二日

Similarly we can obtain

$$P(\mathbf{x} \in R_1 | \mathbf{x} \in \Pi_2) = P(L(\mathbf{x}) > 0 | \mathbf{x} \in \Pi_2)$$

= $P(Z < -\frac{\Delta}{2} - \frac{1}{\Delta} \log_e \frac{\pi_2}{\pi_1})$
= $\Phi\left(-\frac{\Delta}{2} + \frac{1}{\Delta} \log_e \frac{\pi_2}{\pi_1}\right)$

If $\pi_1 = \pi_2 = 1/2$ then

$$P(\mathbf{X} \in R_2 | \mathbf{X} \in \Pi_1) = P(\mathbf{X} \in R_1 | \mathbf{X} \in \Pi_2) = \Phi(-\triangle/2)$$

Observation Since miscalculation probability depends on \triangle , we can write the probability of miscalculation as $P(\triangle)$. Plotting the graph for $\pi_1 = \pi_2 = 1/2$, what is your conclusion?

イロト 不得 トイラト イラト 一日

Question How do we implement the method in real data?

(日) (四) (日) (日) (日)

э

Question How do we implement the method in real data? Observations

 $\rightarrow\,$ Note that $\mu_1,\mu_1,\boldsymbol{\Sigma}$ are not known

3

(日) (四) (日) (日) (日)

Question How do we implement the method in real data? Observations

- $\rightarrow~$ Note that $\mu_1, \mu_1, \boldsymbol{\Sigma}$ are not known
- \rightarrow In general there are 2d + d(d + 2) distinct parameters in μ_1, μ_2, Σ that can possibly be estimated from learning the data

A B b A B b

Question How do we implement the method in real data? Observations

- $\rightarrow~$ Note that $\mu_1, \mu_1, \boldsymbol{\Sigma}$ are not known
- \rightarrow In general there are 2d + d(d + 2) distinct parameters in μ_1, μ_2, Σ that can possibly be estimated from learning the data
- → Suppose we have a random sample $\mathbf{X}_{1j}, 1 \leq j \leq n_1$, and $\mathbf{X}_{2l}, 1 \leq l \leq n_2$ with values \mathbf{x}_{1j} and \mathbf{x}_{2l} from Π_1 and Π_2 respectively

くぼう くほう くほう しほ

Question How do we implement the method in real data? Observations

- $\rightarrow~$ Note that $\mu_1, \mu_1, \boldsymbol{\Sigma}$ are not known
- \rightarrow In general there are 2d + d(d + 2) distinct parameters in μ_1, μ_2, Σ that can possibly be estimated from learning the data
- → Suppose we have a random sample $X_{1j}, 1 \le j \le n_1$, and $X_{2l}, 1 \le l \le n_2$ with values x_{1j} and x_{2l} from Π_1 and Π_2 respectively

Sampling methods from a population

 \rightarrow Mixture sampling - a sample of $n=n_1+n_2$ is selected so that n_1 and n_2 are randomly selected

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Question How do we implement the method in real data? Observations

- $\rightarrow~$ Note that $\mu_1, \mu_1, \boldsymbol{\Sigma}$ are not known
- \rightarrow In general there are 2d + d(d + 2) distinct parameters in μ_1, μ_2, Σ that can possibly be estimated from learning the data
- → Suppose we have a random sample $X_{1j}, 1 \le j \le n_1$, and $X_{2l}, 1 \le l \le n_2$ with values x_{1j} and x_{2l} from Π_1 and Π_2 respectively

Sampling methods from a population

- \rightarrow Mixture sampling a sample of $n=n_1+n_2$ is selected so that n_1 and n_2 are randomly selected
- \rightarrow Separate sampling a sample of n_i is randomly selected from $\Pi_i, i=1,2$ and $n=n_1+n_2$

イロト 不得 トイヨト イヨト 二日



Estimation of parameters The ML estimates of μ_i , i = 1, 2 and Σ are

3

(日) (四) (日) (日) (日)

Estimation of parameters The ML estimates of μ_i , i = 1, 2 and Σ are

$$\widehat{\mu}_{i} = \overline{\mathbf{x}}_{i} = \frac{1}{n_{i}} \sum_{j=1}^{n_{i}} \mathbf{x}_{ij}, i = 1, 2 \text{ and}$$

$$\widehat{\Sigma} = \frac{1}{n} S, \ S = S_{1} + S_{2}, \ S_{i} = \sum_{j=1}^{n_{i}} (\mathbf{x}_{ij} - \overline{\mathbf{x}}_{i}) (\mathbf{x}_{ij} - \overline{\mathbf{x}}_{i})^{T}$$

3

(日) (四) (日) (日) (日)

Estimation of parameters The ML estimates of μ_i , i = 1, 2 and Σ are

$$\widehat{\mu}_{i} = \overline{\mathbf{x}}_{i} = \frac{1}{n_{i}} \sum_{j=1}^{n_{i}} \mathbf{x}_{ij}, i = 1, 2 \text{ and}$$

$$\widehat{\Sigma} = \frac{1}{n} S, \ S = S_{1} + S_{2}, \ S_{i} = \sum_{j=1}^{n_{i}} (\mathbf{x}_{ij} - \overline{\mathbf{x}}_{i}) (\mathbf{x}_{ij} - \overline{\mathbf{x}}_{i})^{T}$$

Note that for unbiased estimator of Σ , we can divide S by its degree of freedom $n-2 = n_1 + n_2 - 2$ rather than n to make $\widehat{\Sigma}$

- 3

イロト イポト イヨト イヨト

The probabilities π_1, π_2 can be chosen based on past experiences or can be estimated as

$$\widehat{\pi}_i = \frac{n_i}{n}, i = 1, 2$$

3

イロト イポト イヨト イヨト

The probabilities π_1, π_2 can be chosen based on past experiences or can be estimated as

$$\widehat{\pi}_i = \frac{n_i}{n}, i = 1, 2$$

Then $\widehat{L}(\mathbf{x}) = \widehat{b}_0 + \widehat{\mathbf{b}}^T \mathbf{x}$, where

$$\widehat{\mathbf{b}} = \widehat{\mathbf{\Sigma}}^{-1}(\overline{\mathbf{x}}_1 - \overline{\mathbf{x}}_2) \widehat{b}_0 = -\frac{1}{2} \left[\overline{\mathbf{x}}_1^T \widehat{\mathbf{\Sigma}}^{-1} \overline{\mathbf{x}}_1 - \overline{\mathbf{x}}_2^T \widehat{\mathbf{\Sigma}}^{-1} \overline{\mathbf{x}}_2 \right] + \log_e \frac{n_1}{n} - \log_e \frac{n_2}{n}$$

are ML estimates of \mathbf{b} and b_0 respectively.

★ ∃ ► < ∃ ►</p>

Image: Image:

The probabilities π_1, π_2 can be chosen based on past experiences or can be estimated as

$$\widehat{\pi}_i = \frac{n_i}{n}, i = 1, 2$$

Then $\widehat{L}(\mathbf{x}) = \widehat{b}_0 + \widehat{\mathbf{b}}^T \mathbf{x}$, where

$$\widehat{\mathbf{b}} = \widehat{\mathbf{\Sigma}}^{-1}(\overline{\mathbf{x}}_1 - \overline{\mathbf{x}}_2) \widehat{b}_0 = -\frac{1}{2} \left[\overline{\mathbf{x}}_1^T \widehat{\mathbf{\Sigma}}^{-1} \overline{\mathbf{x}}_1 - \overline{\mathbf{x}}_2^T \widehat{\mathbf{\Sigma}}^{-1} \overline{\mathbf{x}}_2 \right] + \log_e \frac{n_1}{n} - \log_e \frac{n_2}{n}$$

are ML estimates of **b** and b_0 respectively.

Classification rule The classification rule assigns \mathbf{x} to Π_1 if $\widehat{L}(\mathbf{x}) > 0$, and assigns \mathbf{x} to Π_2 otherwise.