

Classification

Daniel Mao

Contents

1	The First Chapter	1
1.1	Classification	1
1.2	Linear Discriminant Analysis (LDA)	2
1.3	Fisher's Linear Discriminant Analysis (FDA)	3

Chapter 1

The First Chapter

1.1 Classification

Definition (Classification). We define **classification** to be the problem of predicting a discrete random variable \mathbf{y} from another random variable \mathbf{x} .

Definition (Error Rate). Let h be a classifier. We define the **error rate** of h , denoted by $L(h)$, to be a probability given by

$$L(h) := \Pr(h(\mathbf{x}) \neq \mathbf{y}).$$

Definition (Empirical Error Rate). Let h be a classifier. We define the **empirical error rate** of h , denoted by $\hat{L}_n(h)$, to be an average given by

$$\hat{L}_n(h) := \frac{1}{n} \sum_{i=1}^n I(h(\mathbf{x}_i) \neq \mathbf{y}_i).$$

Definition (Bayes Classifier). Let $\mathcal{X} = \mathbb{R}^d$ for some $d \in \mathbb{N}$. Let $\mathcal{Y} = \{1..K\}$. We define the **Bayes classifier**, denoted by h^* , to be a function from \mathcal{X} to \mathcal{Y} given by

$$h^*(\mathbf{x}) := \operatorname{argmax}_{k \in \mathcal{Y}} \{\mathbb{P}(Y = k \mid \mathbf{X} = \mathbf{x})\}.$$

Definition (Decision Boundary). Let h be a classifier. We define the **decision boundary** of h , denoted by $D(h)$, to be a set given by

$$D(h) := \{\mathbf{x} : \Pr(\mathbf{y} = 1 \mid \mathbf{x} = \mathbf{x}_0) = \Pr(\mathbf{y} = 0 \mid \mathbf{x} = \mathbf{x}_0)\}.$$

Theorem 1. The Bayes rule is optimal. i.e., if h^* is the Bayes rule and h is any other classification rule, then $L(h^*) \leq L(h)$.

1.2 Linear Discriminant Analysis (LDA)

Theorem 2. Define for each $k \in \mathcal{Y}$ the class conditional f_k by $f_k(\mathbf{x}) := \mathbb{P}(\mathbf{X} = \mathbf{x} \mid Y = k)$. Define for each $k \in \mathcal{Y}$ the prior π_k by $\pi_k := \mathbb{P}(Y = k)$. Assume that the random variables $\mathbf{X} \mid Y = k$ follows a Gaussian distribution. i.e.,

$$\forall k \in \mathcal{Y}, \quad f_k(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right).$$

Then the Bayes classifier h^* is

$$h^*(\mathbf{x}) = \operatorname{argmax}_{k \in \mathcal{Y}} \delta_k(\mathbf{x})$$

where

$$\delta_k(\mathbf{x}) := -\frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \log(\pi_k).$$

Proof.

$$\begin{aligned} & \mathbb{P}(Y = i \mid \mathbf{X} = \mathbf{x}) \sim \mathbb{P}(Y = j \mid \mathbf{X} = \mathbf{x}) \\ \iff & \frac{f_i(\mathbf{x})\pi_i}{\sum_{k=1}^K f_k(\mathbf{x})\pi_k} \sim \frac{f_j(\mathbf{x})\pi_j}{\sum_{k=1}^K f_k(\mathbf{x})\pi_k} \\ \iff & f_i(\mathbf{x})\pi_i \sim f_j(\mathbf{x})\pi_j \\ \iff & \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right) \pi_i \sim \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^\top \Sigma_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) \right) \pi_j \\ \iff & \frac{1}{|\Sigma_i|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right) \pi_i \sim \frac{1}{|\Sigma_j|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^\top \Sigma_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) \right) \pi_j \\ \iff & -\frac{1}{2} \log(|\Sigma_i|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \log(\pi_i) \sim -\frac{1}{2} \log(|\Sigma_j|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^\top \Sigma_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) + \log(\pi_j). \end{aligned}$$

So

$$\begin{aligned} h^*(\mathbf{x}) &= \operatorname{argmax}_{k \in \mathcal{Y}} \mathbb{P}(Y = k \mid \mathbf{X} = \mathbf{x}) \\ &= \operatorname{argmax}_{k \in \mathcal{Y}} \left(-\frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \log(\pi_k) \right) \\ &= \operatorname{argmax}_{k \in \mathcal{Y}} \delta_k(\mathbf{x}). \end{aligned}$$

■

Remark (Euclidian Distance). Suppose $\forall k \in \mathcal{Y}, \Sigma_k = I$. Then

$$\begin{aligned} \delta_k(\mathbf{x}) &= -\frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \log(\pi_k) \\ &= -\frac{1}{2} \log(|I|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top I (\mathbf{x} - \boldsymbol{\mu}_k) + \log(\pi_k) \\ &= -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top (\mathbf{x} - \boldsymbol{\mu}_k) + \log(\pi_k). \end{aligned}$$

So a new point can be classified by its distance from the center of a class, adjusted by some prior. In a two-class problem with equal prior, the discriminating function would be the perpendicular bisector of the two classes' means.

Remark (Mahalanobis Distance). Say the singular value decomposition of Σ_k is $\Sigma_k = U_k S_k V_k^\top$. Since covariance matrices are symmetric, $U = V$. So $\Sigma_k = U_k S_k U_k^\top$. Now

$$\begin{aligned} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) &= (\mathbf{x} - \boldsymbol{\mu}_k)^\top U_k S_k^{-1} U_k^\top (\mathbf{x} - \boldsymbol{\mu}_k) \\ &= (U_k^\top \mathbf{x} - U_k^\top \boldsymbol{\mu}_k)^\top S_k^{-1} (U_k^\top \mathbf{x} - U_k^\top \boldsymbol{\mu}_k) \\ &= (U_k^\top \mathbf{x} - U_k^\top \boldsymbol{\mu}_k)^\top S_k^{-1/2} S_k^{-1/2} (U_k^\top \mathbf{x} - U_k^\top \boldsymbol{\mu}_k) \\ &= (S_k^{-1/2} U_k^\top \mathbf{x} - S_k^{-1/2} U_k^\top \boldsymbol{\mu}_k)^\top I (S_k^{-1/2} U_k^\top \mathbf{x} - S_k^{-1/2} U_k^\top \boldsymbol{\mu}_k) \\ &= (A\mathbf{x} - A\boldsymbol{\mu}_k)^\top I (A\mathbf{x} - A\boldsymbol{\mu}_k), \text{ where } A := S_k^{-1/2} U_k^\top. \end{aligned}$$

That is,

$$(\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) = (A\mathbf{x} - A\boldsymbol{\mu}_k)^\top I (A\mathbf{x} - A\boldsymbol{\mu}_k)$$

where $A = S_k^{-1/2} U_k^\top$.

1.3 Fisher's Linear Discriminant Analysis (FDA)

FDA is a method of dimensionality reduction. Projects the data on the direction of maximum separation.

Goal:

- Minimize the variance of each class.
- Maximize the distance between projected means.

$$\max (\mathbf{w}^\top \boldsymbol{\mu}_0 - \mathbf{w}^\top \boldsymbol{\mu}_1)^2.$$

$$\begin{aligned} (\mathbf{w}^\top \boldsymbol{\mu}_0 - \mathbf{w}^\top \boldsymbol{\mu}_1)^2 &= (\mathbf{w}^\top \boldsymbol{\mu}_0 - \mathbf{w}^\top \boldsymbol{\mu}_1)^\top (\mathbf{w}^\top \boldsymbol{\mu}_0 - \mathbf{w}^\top \boldsymbol{\mu}_1) \\ &= (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w} \mathbf{w}^\top (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) \\ &= \mathbf{w}^\top (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w} \end{aligned}$$

The variance of projected points in each class are

$$\mathbf{w}^\top \Sigma_0 \mathbf{w}, \text{ and } \mathbf{w}^\top \Sigma_1 \mathbf{w}.$$

To minimize them,

$$\min \quad \mathbf{w}^\top \Sigma_0 \mathbf{w} + \mathbf{w}^\top \Sigma_1 \mathbf{w}.$$

$$\iff \min \quad \mathbf{w}^\top (\Sigma_0 + \Sigma_1) \mathbf{w}.$$

So we can do

$$\max \quad \mathbf{w}^\top S_B \mathbf{w}$$

subject to $\mathbf{w}^\top S_w \mathbf{w} = 1$.

$$L(\mathbf{w}, \lambda) = \mathbf{w}^\top S_B \mathbf{w} - \lambda(\mathbf{w}^\top S_w \mathbf{w} - 1).$$

$$\frac{\partial}{\partial \mathbf{w}} L = 2S_B \mathbf{w} - 2\lambda S_w \mathbf{w}$$

Setting this to zero we get

$$S_B \mathbf{w} = \lambda S_w \mathbf{w}.$$

$$\iff S_w^{-1} S_B \mathbf{w} = \lambda \mathbf{w}.$$

Note that $\text{rank}(S_B) = 1$. So $\text{rank}(S_w^{-1} S_B) = 1$. So $S_w^{-1} S_B$ has only one eigenvector.