

基础数理统计

第二十二章 分类

1 22.1 引言

2 22.2 错误率与贝叶斯分类器

3 22.3 高斯分类器与线性分类器

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

22.1 引言

1 22.1 引言

2 22.2 错误率与贝叶斯分类器

3 22.3 高斯分类器与线性分类器

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

[22.1 引言](#)[22.2 错误率与贝叶斯分类器](#)[22.3 高斯分类器与线性分类器](#)

定义 1 (分类)

从一个随机变量 X 来预测另一个离散的随机变量 Y 的问题被称作是分类，或有指导的学习，或判别，或者称为模式识别。

具体来说，考虑 i.i.d 数据 $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$ ，其中 $\mathbf{X}_i = (X_{i1}, X_{i2}, \dots, X_{id})^\top \in \mathcal{X} \subset \mathbb{R}^d$ 为一个 d 维向量且 Y_i 在某个有限集 \mathcal{Y} 中取值。一个分类规则就是一个函数 $h: \mathcal{X} \rightarrow \mathcal{Y}$ 。当观测到一个新的 \mathbf{X} ，预测 Y 为 $h(\mathbf{X})$ 。

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

我们首先从理论上研究分类问题：

分类问题

已知两个总体分布 $F(\mathbf{x})$ 和 $G(\mathbf{x})$, 对一个新的观察值 $\mathbf{X} = \mathbf{x}$ 如何分类？

- 例如 $N(0, 1)$ 与 $Exp(1)$;
- 例如 $N(0, 1)$ 与 $U[-1, 1]$;
- ...

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

思考 1

对于两个正态总体 $N(0, 1)$ 和 $N(1, 1)$, 直观上应该如何分类?

思考 2

对于两个正态总体 $N(0, 1)$ 和 $N(1, 2)$, 应该如何分类?

22.2 错误率与贝叶斯分类器

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

1 22.1 引言

2 22.2 错误率与贝叶斯分类器

3 22.3 高斯分类器与线性分类器

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

每一个分类准则或者分类方法是把整个样本空间 $\mathbf{x} \in \mathbb{R}^p$ 分成两部分 \mathcal{A} 和 \mathcal{A}^c .

- 当 $\mathbf{x} \in \mathcal{A}$ 时候, 认为新的观察值来自于 $F(\mathbf{x})$.
- 当 $\mathbf{x} \in \mathcal{A}^c$ 时候, 认为新的观察值来自于 $G(\mathbf{x})$

好的分类器即找出好的区域 \mathcal{A} .

[22.1 引言](#)[22.2 错误率与贝叶斯分类器](#)[22.3 高斯分类器与线性分类器](#)

给定一个分类方法，即给定区域 \mathcal{A} ,

$$\begin{aligned} & P(\mathbf{X} \text{ 被错分}) \\ &= P(\mathbf{X} \in \mathcal{A}^c | \mathbf{X} \sim F) P(\mathbf{X} \sim F) + P(\mathbf{X} \in \mathcal{A} | \mathbf{X} \sim G) P(\mathbf{X} \sim G) \\ &= \pi_1 \int_{\mathcal{A}^c} dF(\mathbf{x}) + \pi_2 \int_{\mathcal{A}} dG(\mathbf{x}) \end{aligned}$$

其中 π_1, π_2 为先验概率。

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

更一般的可以延伸为带有错分成本的损失函数：

$$L(\mathcal{A}) = P(X \in \mathcal{A}^c | X \sim F) P(X \sim F) * \text{loss}_1 \\ + P(X \in \mathcal{A} | X \sim G) P(X \sim G) * \text{loss}_2.$$

这里 loss_1 是把来自 F 的样本误判为 G 的损失， loss_2 是把来自 G 的样本误判为 F 的损失。形式上，

$$L(\mathcal{A}) = \pi_1 \int_{\mathcal{A}^c} dF(x) + \pi_2 \int_{\mathcal{A}} dG(x),$$

其中 π_1, π_2 为先验概率乘以对应损失。

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

最优的分类方法 (Bayes 分类器) 为:

$$\begin{aligned}\mathcal{A}^* &= \arg \min_{\mathcal{A}} \pi_1 \int_{\mathcal{A}^c} dF(x) + \pi_2 \int_{\mathcal{A}} dG(x) \\ &= \arg \min_{\mathcal{A}} \{ \pi_1 P(\mathcal{A}^c) + \pi_2 Q(\mathcal{A}) \},\end{aligned}$$

其中 P, Q 为分别为分布函数 F, G 对应的概率函数, 即

$$P(\mathcal{A}) = \int_{\mathcal{A}} dF(x), Q(\mathcal{A}) = \int_{\mathcal{A}} dG(x).$$

(1) $F(x)$ 和 $G(x)$ 有密度函数 $f(x)$ 和 $g(x)$,

$$\mathcal{A}^* = \{x : \pi_1 f(x) > \pi_2 g(x)\}.$$

(2) $F(x)$ 和 $G(x)$ 是离散分布:

$$\mathcal{A}^* = \{x_i : \pi_1 P(X = x_i) > \pi_2 Q(X = x_i)\}.$$

证明.

以 (1) 为例, 记 $h(\mathcal{A}) = \pi_1 P(\mathcal{A}^c) + \pi_2 Q(\mathcal{A})$, 则易证明

$$\begin{aligned} h(\mathcal{A}) - h(\mathcal{A}^*) &= \int_{\mathcal{A}^c \cap \mathcal{A}^*} [\pi_1 f(x) - \pi_2 g(x)] dx \\ &\quad - \int_{\mathcal{A}^{*c} \cap \mathcal{A}} [\pi_1 f(x) - \pi_2 g(x)] dx. \end{aligned}$$

□

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

在 $\pi_1 = \pi_2 = 1/2$ 的条件下, 计算

例 1

假定 F, G 分别是正态分布 $N(\mu_1, \sigma_1^2)$ 和 $N(\mu_2, \sigma_2^2)$ 的分布函数, 计算理论分类器及错分率。

$$\mathcal{A}^* = \left\{ x : \frac{1}{\sigma_1} \exp \left\{ -\frac{(x - \mu_1)^2}{2\sigma_1^2} \right\} > \frac{1}{\sigma_2} \exp \left\{ -\frac{(x - \mu_2)^2}{2\sigma_2^2} \right\} \right\}.$$

例 2

假定 F, G 分别是抛硬币和扔色子的分布函数, 即

$$P(X = 1) = P(X = 2) = 1/2;$$

$$Q(X = 1) = Q(X = 2) = \cdots = Q(X = 6) = 1/6,$$

计算理论分类器及错分率。

$$\mathcal{A}^* = \{1, 2\}.$$

错分率是 $1/6$.

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

例 3

假定 F, G 分别是正态分布 $N(0, 1)$ 和指数分布 $\text{Exp}(1)$ 的分布函数, 计算理论分类器及错分率。

$$\mathcal{A}^* = \{x : x < 0\}.$$

错分率是 $1/4$ 。

例 4

假定 F, G 分别是正态分布 $N(0, 1)$ 和均匀分布 $U[-1, 1]$ 的分布函数, 计算理论分类器及错分率。

$$\mathcal{A}^* = \{x : |x| > 1\}$$

错分率是 $1/2(\Phi(1) - \Phi(-1))$ 。

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

定义 2 (误差率)

一个分类器 h 的真实误差率为

$$L(h) = P(h(\mathbf{X}) \neq Y),$$

而经验误差率或训练误差率为

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

考虑特殊情况, $\mathcal{Y} = \{0, 1\}$, 令

$$r(x) = E(Y|X=x) = \frac{\pi f_1(x)}{\pi f_1(x) + (1-\pi)f_0(x)},$$

其中

$$f_0(x) = f(x|Y=0), f_1(x) = f(x|Y=1), \pi = P(Y=1).$$

定义 3 (贝叶斯分类规则)

贝叶斯分类规则 h^* 为

$$h^*(x) = \begin{cases} 1 & r(x) > 1/2 \\ 0 & \text{其他} \end{cases}$$

集合 $\mathcal{D}(h) = \{x: P(Y=1|X=x) = P(Y=0|X=x)\}$ 称为决策边界。

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

贝叶斯分类规则可以写成一些等价形式

$$\begin{aligned}h^*(x) &= \begin{cases} 1 & P(Y=1|X=x) > P(Y=0|X=x) \\ 0 & \text{其他} \end{cases} \\&= \begin{cases} 1 & \pi f_1(x) > (1-\pi)f_0(x) \\ 0 & \text{其他} \end{cases}\end{aligned}$$

定理 1

贝叶斯规则是最优的，即若 h 是任何其他分类准则，则 $L(h^*) \leq L(h)$.

实际中贝叶斯规则是基于数据估计得到的，

- 经验风险极小化：选择一组分类器集 \mathcal{H} 并且找到 $\hat{h} \in \mathcal{H}$ 使得能够极小化 $L(h)$ 的某个估计；
- 回归：找到回归函数 r 的一个估计 \hat{r} 并定义

$$\hat{h}(x) = \begin{cases} 1 & \hat{r}(x) > \frac{1}{2} \\ 0 & \text{其他} \end{cases}$$

- 密度估计：估计 f_0, f_1 并且令 $\hat{\pi} = n^{-1} \sum_{i=1}^n Y_i$ 。定义

$$\hat{r}(x) = \hat{P}(Y = 1 | X = x) = \frac{\hat{\pi} \hat{f}_1(x)}{\hat{\pi} \hat{f}_1(x) + (1 - \hat{\pi}) \hat{f}_0(x)},$$

$$\hat{h}(x) = \begin{cases} 1 & \hat{r}(x) > \frac{1}{2} \\ 0 & \text{其他} \end{cases}$$

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

定理 2

假设 $Y \in \mathcal{Y} = \{1, 2, \dots, K\}$, 最优规则为

$$h(x) = \operatorname{argmax}_k P(Y = k | X = x) = \operatorname{argmax}_k (\pi_k f_k(x)),$$

其中

$$P(Y = k | X = x) = \frac{\pi_k f_k(x)}{\sum_r \pi_r f_r(x)},$$

$$\pi_r = P(Y = r), \quad f_r(x) = f(x | Y = r).$$

22.3 高斯分类器与线性分类器

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

1 22.1 引言

2 22.2 错误率与贝叶斯分类器

3 22.3 高斯分类器与线性分类器

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

假设 $\mathbf{x}|y = k$ 服从多元正态分布 $k = 0, 1$:

$$f_k(\mathbf{x}) = (2\pi)^{-\frac{p}{2}} |\Sigma_k|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}.$$

其中 $\boldsymbol{\mu}_k$ 为 p 维均值, Σ_k 为 $p \times p$ 维的正定对称矩阵。

定理 3

若 $\mathbf{x}|y=0 \sim N_p(\boldsymbol{\mu}_0, \Sigma_0)$, $\mathbf{x}|y=1 \sim N_p(\boldsymbol{\mu}_1, \Sigma_1)$, 则贝叶斯规则为

$$h^*(\mathbf{x}) = \begin{cases} 1 & r_1^2 < r_0^2 + 2 \log \left(\frac{\pi_1}{\pi_0} \right) + \log \left(\frac{|\Sigma_0|}{|\Sigma_1|} \right) \\ 0 & \text{其他} \end{cases}$$

其中 $r_i^2(\mathbf{x}) = (\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)$, $i=0,1$ 为 Mahalanobis 距离。等价于

$$h^*(\mathbf{x}) = \operatorname{argmax}_k \delta_k(\mathbf{x}),$$

其中

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

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

[22.1 引言](#)[22.2 错误率与贝叶斯分类器](#)[22.3 高斯分类器与线性分类器](#)

整理得到 QDA(Quadratic Discriminant Analysis):

$$\begin{aligned} & \mathbf{x}^\top (\Sigma_0^{-1} - \Sigma_1^{-1}) \mathbf{x} - 2\mathbf{x}^\top (\Sigma_0^{-1} \boldsymbol{\mu}_0 - \Sigma_1^{-1} \boldsymbol{\mu}_1) \\ & + \boldsymbol{\mu}_0^\top \Sigma_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1^\top \Sigma_1^{-1} \boldsymbol{\mu}_1 + 2 \log \frac{\pi_1}{\pi_0} + \log \frac{|\Sigma_0|}{|\Sigma_1|} \leq 0. \end{aligned}$$

分类器具有二次型的形式:

$$\mathbf{x}^\top \Omega \mathbf{x} + \mathbf{x}^\top \boldsymbol{\beta} + \alpha \leq 0.$$

[22.1 引言](#)[22.2 错误率与贝叶斯分类器](#)[22.3 高斯分类器与线性分类器](#)

$$\begin{aligned}n_0 &= \sum_{i=1}^n (1 - Y_i), \quad n_1 = \sum_{i=1}^n Y_i, \\ \hat{\pi}_0 &= \frac{1}{n} \sum_{i=1}^n (1 - Y_i), \quad \hat{\pi}_1 = \frac{1}{n} \sum_{i=1}^n Y_i, \\ \hat{\boldsymbol{\mu}}_0 &= \frac{1}{n_0} \sum_{i: Y_i=0} \mathbf{x}_i, \quad S_0 = \frac{1}{n_0} \sum_{i: Y_i=0} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_0)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_0)^\top, \\ \hat{\boldsymbol{\mu}}_1 &= \frac{1}{n_1} \sum_{i: Y_i=1} \mathbf{x}_i, \quad S_1 = \frac{1}{n_1} \sum_{i: Y_i=1} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_1)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_1)^\top\end{aligned}$$

[22.1 引言](#)[22.2 错误率与贝叶斯分类器](#)[22.3 高斯分类器与线性分类器](#)

特别的，如果协方差结构相同 $\Sigma_0 = \Sigma_1 = \Sigma$ ，我们得到 LDA(Linear Discriminant Analysis):

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

即

$$\left(\mathbf{x} - \frac{\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1}{2}\right)^\top \Sigma^{-1}(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) + \log \frac{\pi_0}{\pi_1} \geq 0,$$

[22.1 引言](#)[22.2 错误率与贝叶斯分类器](#)[22.3 高斯分类器与线性分类器](#)

$$S = \frac{n_0 S_0 + n_1 S_1}{n_0 + n_1},$$

分类规则等价为

$$h^*(\mathbf{x}) = \begin{cases} 1 & \delta_1(\mathbf{x}) > \delta_0(\mathbf{x}) \\ 0 & \text{其他} \end{cases}$$

其中

$$\delta_k(\mathbf{x}) = \mathbf{x}^\top S^{-1} \hat{\boldsymbol{\mu}}_k - \frac{1}{2} \hat{\boldsymbol{\mu}}_j^\top S^{-1} \hat{\boldsymbol{\mu}}_j + \log \hat{\pi}_j.$$

作为判别函数。

Fisher 的线性判别分析

- (1). 将数据投影到一条直线上, 即用 $U = \omega^\top \mathbf{x}$ 代替 \mathbf{x} 来进行分类;
- (2). 不同类尽可能分开, Σ 为 \mathbf{x} 的协方差阵, 定义分离为

$$\begin{aligned} J(\omega) &= \frac{(\mathbb{E}(U|Y=0) - \mathbb{E}(U|Y=1))^2}{\omega^\top \Sigma \omega} \\ &= \frac{\omega^\top (\mu_0 - \mu_1)(\mu_0 - \mu_1)^\top \omega}{\omega^\top \Sigma \omega} \end{aligned}$$

- (3) 实际中 $J(\omega)$ 用 $\hat{J}(\omega)$ 进行估计

$$\hat{J}(\omega) = \frac{\omega^\top S_B \omega}{\omega^\top S_W \omega},$$

其中

$$S_B = (\bar{\mathbf{x}}_0 - \bar{\mathbf{x}}_1)(\bar{\mathbf{x}}_0 - \bar{\mathbf{x}}_1)^\top, \quad S_W = \frac{(n_0 - 1)S_0 + (n_1 - 1)S_1}{(n_0 - 1) + (n_1 - 1)}.$$

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

定理 4

向量 $\omega = S_W^{-1}(\bar{\mathbf{x}}_0 - \bar{\mathbf{x}}_1)$ 为 $\hat{J}(\omega)$ 的极小值点, 称

$$U = \omega^\top \mathbf{x} = (\bar{\mathbf{x}}_0 - \bar{\mathbf{x}}_1)^\top S_W^{-1} \mathbf{x}$$

为 *Fisher* 线性判别函数。定义

$$m = \frac{1}{2}(\bar{\mathbf{x}}_0 - \bar{\mathbf{x}}_1)^\top S_W^{-1}(\bar{\mathbf{x}}_0 + \bar{\mathbf{x}}_1).$$

Fisher 分类规则为

$$h(\mathbf{x}) = \begin{cases} 0 & \omega^\top \mathbf{x} \geq m \\ 1 & \omega^\top \mathbf{x} < m. \end{cases}$$

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器

- 逻辑回归;
- 支持向量机;
- 随机森林;
- 神经网络;
- 决策树;
- 最小近邻法;
- ...

22.1 引言

22.2 错误率与贝叶斯分类器

22.3 高斯分类器与线性分类器