

机器学习入门

贝叶斯决策

河北师范大学软件学院 张朝晖

2017.04.16-04.26



监督式学习

- 分类 { A. 懒惰学习 -- KNN法(K最近邻法)分类模型
B. 概率学习 -- 贝叶斯分类器
C. 分而治之 -- 分类树
- 回归 { D. KNN法回归
E. 分而治之 -- 回归树
F. 最小二乘回归

集成学习

- 非监督式学习 { 聚类 -- G. K-Means Clustering
特征提取 -- H. PCA

模型评价

回顾两种典型的贝叶斯分类模型

模型1--最小错误率贝叶斯分类

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \quad \sum_{i=1}^c P(\omega_j | \mathbf{x}) = 1$$

观测空间为连续特征空间

$$P(\omega_j | \mathbf{x}) = \frac{p(\mathbf{x}, \omega_j)}{p(\mathbf{x})} = \frac{p(\mathbf{x} | \omega_j) P(\omega_j)}{p(\mathbf{x})} = \frac{p(\mathbf{x} | \omega_j) P(\omega_j)}{\sum_{i=1}^c p(\mathbf{x} | \omega_i) P(\omega_i)}$$

观测空间为离散特征空间

$$P(\omega_j | \mathbf{x}) = \frac{P(\mathbf{x}, \omega_j)}{P(\mathbf{x})} = \frac{P(\mathbf{x} | \omega_j) P(\omega_j)}{P(\mathbf{x})} = \frac{P(\mathbf{x} | \omega_j) P(\omega_j)}{\sum_{i=1}^c P(\mathbf{x} | \omega_i) P(\omega_i)}$$



两类情况下，基于最小错误率的判决规则，各等价形式：

1， 状态后验概率：

$$\begin{cases} P(\omega_1 | \mathbf{x}) > P(\omega_2 | \mathbf{x}), \text{则 } \mathbf{x} \in \omega_1 \text{类;} \\ P(\omega_1 | \mathbf{x}) < P(\omega_2 | \mathbf{x}), \text{则 } \mathbf{x} \in \omega_2 \text{类} \end{cases}$$

2， 似然值 \times 先验概率

$$\begin{cases} P(\omega_1)p(\mathbf{x} | \omega_1) > P(\omega_2)p(\mathbf{x} | \omega_2), \text{则 } \mathbf{x} \in \omega_1 \text{类;} \\ P(\omega_1)p(\mathbf{x} | \omega_1) < P(\omega_2)p(\mathbf{x} | \omega_2), \text{则 } \mathbf{x} \in \omega_2 \text{类} \end{cases}$$

3， 似然比

$$\begin{cases} \text{若 } l(\mathbf{x}) = \frac{p(\mathbf{x} | \omega_1)}{p(\mathbf{x} | \omega_2)} > \frac{P(\omega_2)}{P(\omega_1)}, \text{则 } \mathbf{x} \in \omega_1 \text{类;} \\ \text{若 } l(\mathbf{x}) = \frac{p(\mathbf{x} | \omega_1)}{p(\mathbf{x} | \omega_2)} < \frac{P(\omega_2)}{P(\omega_1)}, \text{则 } \mathbf{x} \in \omega_2 \text{类} \end{cases}$$

多类情况下，基于最小错误率的判决规则，各等价形式

(1)后验概率：

$$\text{若 } P(\omega_i | \mathbf{x}) = \max_{j=1,2,\dots,c} P(\omega_j | \mathbf{x})$$

则 $\mathbf{x} \in \omega_i$ 类

(2)似然值 \times 先验概率：

$$\text{若 } p(\mathbf{x} | \omega_i)P(\omega_i) = \max_{j=1,2,\dots,c} p(\mathbf{x} | \omega_j)P(\omega_j)$$

则 $\mathbf{x} \in \omega_i$ 类



对于观测 \mathbf{x}

(1) 计算后验概率:

$$P(\omega_j | \mathbf{x}) = \frac{P(\omega_j) p(\mathbf{x} | \omega_j)}{p(\mathbf{x})} = \frac{P(\omega_j) p(\mathbf{x} | \omega_j)}{\sum_{i=1}^c P(\omega_i) p(\mathbf{x} | \omega_i)}$$

$$j = 1, 2, \dots, c$$

(2) 计算 \mathbf{x} 的条件风险: $R(\alpha_i | \mathbf{x}) = \sum_{j=1}^c \lambda(\alpha_i, \omega_j) P(\omega_j | \mathbf{x})$

$$i = 1, 2, \dots, k$$

(3) 决策: 选择关于 \mathbf{x} 的条件风险最小的决策 α , 即

$$\alpha = \arg \min_{i=1,2,\dots,k} R(\alpha_i | \mathbf{x})$$

连续特征空间

若各类条件概率密度函数为正态分布

1. 引言

2. 贝叶斯决策模型

2.1 最小错误率贝叶斯决策

2.2 最小风险的贝叶斯决策

3. 正态分布的概率密度函数及性质

4. 概率/概率密度函数估计

5. 朴素贝叶斯(Naive Bayes)分类

问题的引入:

➤ 贝叶斯决策中, 涉及“连续随机变量或向量的**类条件概率密度函数**”。

➤ “正态分布”的概率密度函数, 特点:

物理上的合理性。如果在特征空间中的某一类样本, 较多地分布在这一类均值附近, 远离均值点的样本比较少, 此时用正态分布作为这一类的概率模型是合理的。

数学上, 比较简便。正态分布概率模型有许多好的性质, 有利于作数学分析。

➔ 简单、符合一些实际情况

主要内容

3.1. 单个随机变量的正态分布

3.2. 随机向量的正态分布

连续随机变量 X 概率密度函数定义

$$p_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$$

记 $p_X(x) \sim N(\mu, \sigma^2)$

标准正态分布: $p_X(x) \sim N(0, 1)$

其中 $\left\{ \begin{array}{l} \mu - \text{随机变量 } X \text{ 的期望, 或均值} \\ \mu \equiv E\{X\} = \int_{-\infty}^{\infty} xp_X(x)dx \\ \sigma - x \text{ 的标准差, 描述 } x \text{ 的分散程度} \\ \sigma^2 - x \text{ 的方差: } \sigma^2 \equiv E\{(X-\mu)^2\} = \int_{-\infty}^{\infty} (x-\mu)^2 p_X(x)dx \end{array} \right.$

性质: $p_X(x) \geq 0, -\infty < x < +\infty$

$$\int_{-\infty}^{\infty} p_X(x) dx = 1$$

$$\int_{\mu-\sigma}^{\mu+\sigma} p_X(x) dx = 0.683$$

$$\int_{\mu-2\sigma}^{\mu+2\sigma} p_X(x) dx = 0.9544$$

$$\int_{\mu-3\sigma}^{\mu+3\sigma} p_X(x) dx = 0.9974$$

其它: 熵 $H(X) = -\int p_X(x) \ln p_X(x) dx > 0$

*Mahalanobis*距离 $r = \frac{|x - \mu|}{\sigma}$

$$E\{f(X)\} = \int_{-\infty}^{\infty} f(x) p_X(x) dx$$

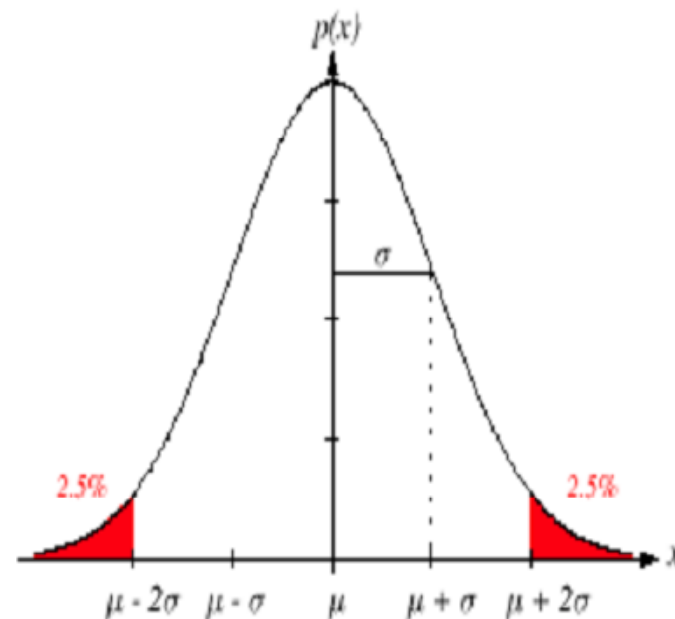


图: 单变量正态分布。

$$p_X(x=\mu) = \frac{1}{\sqrt{2\pi}\sigma}$$

主要内容

3.1 单个随机变量的正态分布

3.2 随机向量的正态分布



[1]多元正态分布的概率密度函数 -- $p_X(x)$

定义:
$$p_X(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

记为 $p_X(x) \sim N(\boldsymbol{\mu}, \Sigma)$

其中: $\left\{ \begin{array}{l} \mathbf{X} - d\text{维列向量}, \mathbf{X} = [X^{(1)}, X^{(2)}, \dots, X^{(d)}]^T \\ \quad X^{(i)} \text{为随机变量}, i = 1, 2, \dots, d; \\ \boldsymbol{\mu} - d\text{维均值向量}, \boldsymbol{\mu} = [\mu^{(1)}, \mu^{(2)}, \dots, \mu^{(d)}]^T \\ \quad \mu^{(i)} \text{为第} i \text{维随机变量} X^{(i)} \text{的期望}, i = 1, 2, \dots, d; \\ \Sigma - d \times d \text{维协方差矩阵 (covariance matrix)} \\ \Sigma^{-1} - \text{矩阵} \Sigma \text{的逆矩阵, 精度矩阵} \\ |\Sigma| - \text{矩阵} \Sigma \text{的行列式} \end{array} \right.$

[1]多元正态分布概率密度函数定义(续) -- μ

μ - d 维均值向量, $\mu = [\mu^{(1)}, \mu^{(2)}, \dots, \mu^{(d)}]^T$

$$\mu \equiv E\{X\} = \int_{\mathbb{R}^d} \mathbf{x} p_X(x) d\mathbf{x}$$

$$= [E\{X^{(1)}\}, E\{X^{(2)}\}, \dots, E\{X^{(d)}\}]^T = [\mu^{(1)}, \mu^{(2)}, \dots, \mu^{(d)}]^T$$

其中:

$$\left\{ \begin{array}{l} \mu_i = E\{X^{(i)}\} = \int_{\mathbb{R}^d} x^{(i)} p_X(x) d\mathbf{x} = \int_{-\infty}^{+\infty} x^{(i)} p_{X^{(i)}}(x^{(i)}) dx^{(i)} \\ \text{边缘密度函数 } p_{X^{(i)}}(x^{(i)}) = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} p_X(x) dx^{(1)} \dots dx^{(i-1)} dx^{(i+1)} \dots dx^{(d)} \end{array} \right.$$

[1]多元正态分布概率密度函数定义(续)-- Σ

Σ — $d \times d$ 维协方差矩阵(covariance matrix):

$$\Sigma \equiv E \{ (X - \mu)(X - \mu)^T \}$$

$$= E \left\{ \begin{bmatrix} X^{(1)} - \mu^{(1)} \\ X^{(2)} - \mu^{(2)} \\ \vdots \\ X^{(d)} - \mu^{(d)} \end{bmatrix} \begin{bmatrix} X^{(1)} - \mu^{(1)} & X^{(2)} - \mu^{(2)} & \cdots & X^{(d)} - \mu^{(d)} \end{bmatrix} \right\}$$

$$= E \begin{bmatrix} (X^{(1)} - \mu^{(1)})^2 & (X^{(1)} - \mu^{(1)})(X^{(2)} - \mu^{(2)}) & \cdots & (X^{(1)} - \mu^{(1)})(X^{(d)} - \mu^{(d)}) \\ (X^{(2)} - \mu^{(2)})(X^{(1)} - \mu^{(1)}) & (X^{(2)} - \mu^{(2)})^2 & \cdots & (X^{(2)} - \mu^{(2)})(X^{(d)} - \mu^{(d)}) \\ \vdots & \vdots & \vdots & \vdots \\ (X^{(d)} - \mu^{(d)})(X^{(1)} - \mu^{(1)}) & (X^{(d)} - \mu^{(d)})(X^{(2)} - \mu^{(2)}) & \cdots & (X^{(d)} - \mu^{(d)})^2 \end{bmatrix}$$

= (下页待续)

[1]多元正态分布概率密度函数定义(续)-- Σ

Σ — $d \times d$ 维协方差矩阵(covariance matrix):

$$\Sigma \equiv E[(X - \mu)(X - \mu)^T]$$

$$= \begin{bmatrix} E\left\{\left(X^{(1)} - \mu^{(1)}\right)^2\right\} & E\left\{\left(X^{(1)} - \mu^{(1)}\right)\left(X^{(2)} - \mu^{(2)}\right)\right\} & \cdots & E\left\{\left(X^{(1)} - \mu^{(1)}\right)\left(X^{(d)} - \mu^{(d)}\right)\right\} \\ E\left\{\left(X^{(2)} - \mu^{(2)}\right)\left(X^{(1)} - \mu^{(1)}\right)\right\} & E\left\{\left(X^{(2)} - \mu^{(2)}\right)^2\right\} & \cdots & E\left\{\left(X^{(2)} - \mu^{(2)}\right)\left(X^{(d)} - \mu^{(d)}\right)\right\} \\ \vdots & \vdots & \vdots & \vdots \\ E\left\{\left(X^{(d)} - \mu^{(d)}\right)\left(X^{(1)} - \mu^{(1)}\right)\right\} & E\left\{\left(X^{(d)} - \mu^{(d)}\right)\left(X^{(2)} - \mu^{(2)}\right)\right\} & \cdots & E\left\{\left(X^{(d)} - \mu^{(d)}\right)^2\right\} \end{bmatrix}$$

$$= \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2d} \\ \cdots & \cdots & \cdots & \cdots \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_{dd} \end{bmatrix}$$

Σ 对称非负定 $\begin{cases} \text{对角线元素 } \sigma_{ii} = \sigma^{(i)2}, & X_i \text{ 方差} \\ \text{非对角元素 } \sigma_{ij}, & X_i, X_j \text{ 协方差} \end{cases}$

这里: 只考虑 Σ 对称正定

[2]多元正态分布概率密度函数的几个典型性质

性质1 多元正态分布完全由 μ 、 Σ 决定

参数共计 $d + \frac{d(d+1)}{2}$ 个

$$\mu = [\mu^{(1)}, \mu^{(2)}, \dots, \mu^{(d)}]^T$$
$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2d} \\ \dots & \dots & \dots & \dots \\ \sigma_{d1} & \sigma_{d2} & \dots & \sigma_{dd} \end{bmatrix} \text{ 为对称矩阵}$$

$$p_X(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right]$$

$$p_X(\mathbf{x}) \sim N(\mu, \Sigma)$$



性质2 等概率密度点的轨迹为一超椭球面

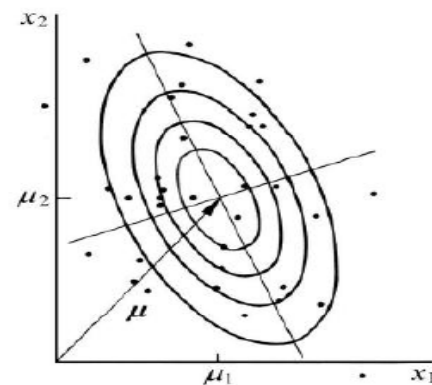
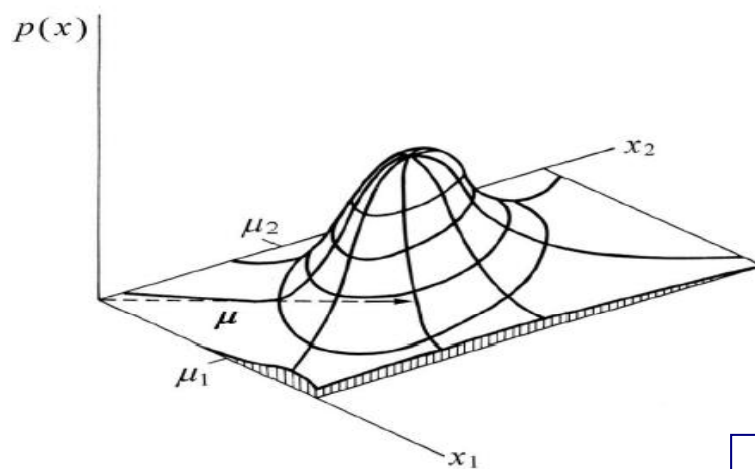


河北师范大学软件学院
Software College of Hebei Normal University

$$p_X(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

\mathbf{x} 到 $\boldsymbol{\mu}$ 的 **Mahalanobis距离**平方 $r^2 = (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$

等概率密度点: $r^2 = \text{常数}$



超椭球面中心 $\boldsymbol{\mu}$

二维正态分布，等概率密度点的轨迹为**椭圆周**。

超椭球面主轴 $\left\{ \begin{array}{l} \text{方向取决于} \Sigma \text{本征向量} \\ \text{长度与} \Sigma \text{本征值的平方根成正比。} \end{array} \right.$



性质3 不相关性等价于独立性。

$$\mathbf{x} \text{多元正态分布} \left\{ \begin{array}{l} \text{"}\mathbf{X} \text{任意两分量 } X_i, X_j \text{ 间互不相关"} \\ \Leftrightarrow \text{"各分量间相互独立"} \\ \text{协方差矩阵 } \Sigma \text{ 是对角矩阵} \Rightarrow \mathbf{X} \text{ 各分量相互独立} \end{array} \right.$$

定义：随机变量 $X^{(i)}, X^{(j)}$ 之间 $\boldsymbol{\mu} = [\mu^{(1)}, \mu^{(2)}, \dots, \mu^{(d)}]^T$

$$\left\{ \begin{array}{l} X^{(i)}, X^{(j)} \text{ 间不相关: } \mathbf{E} \{ X^{(i)} X^{(j)} \} = \mathbf{E} \{ X^{(i)} \} \mathbf{E} \{ X^{(j)} \} \\ X^{(i)}, X^{(j)} \text{ 相互独立: } p_{X^{(i)} X^{(j)}}(x^{(i)}, x^{(j)}) = p_{X^{(i)}}(x^{(i)}) p_{X^{(j)}}(x^{(j)}) \end{array} \right.$$

并且

$$\boxed{p_{X^{(i)} X^{(j)}}(x^{(i)}, x^{(j)}) = p_{X^{(i)}}(x^{(i)}) p_{X^{(j)}}(x^{(j)})} \Rightarrow \boxed{\mathbf{E} \{ X^{(i)} X^{(j)} \} = \mathbf{E} \{ X^{(i)} \} \mathbf{E} \{ X^{(j)} \}}$$

性质3 不相关性等价于独立性。

多元正态分布的任意随机变量间不相关性，等价于独立性。

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2d} \\ \cdots & \cdots & \cdots & \cdots \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_{dd} \end{bmatrix} = \begin{bmatrix} \sigma_{11} & & & \\ & \sigma_{22} & & \\ & & \cdots & \\ & & & \sigma_{dd} \end{bmatrix}$$

若 X_i, X_j 间互不相关, 则协方差矩阵 Σ 是对角的

性质3 不相关性等价于独立性。

$$(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^d \left(\frac{x^{(i)} - \mu^{(i)}}{\sigma^{(i)}} \right)^2$$

$$\begin{aligned} p_{\mathbf{X}}(\mathbf{x}) &= \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \\ &= \frac{1}{(2\pi)^{d/2} \prod_{i=1}^d \sigma^{(i)}} \exp \left[-\frac{1}{2} \sum_{i=1}^d \left(\frac{x^{(i)} - \mu^{(i)}}{\sigma^{(i)}} \right)^2 \right] = \prod_{i=1}^d \frac{1}{\sqrt{2\pi} \sigma^{(i)}} \exp \left[-\frac{1}{2} \left(\frac{x^{(i)} - \mu^{(i)}}{\sigma^{(i)}} \right)^2 \right] \\ &= \prod_{i=1}^d p_{X^{(i)}}(x^{(i)}) \end{aligned}$$

对于多元正态分布 \mathbf{X} ,

各分量 X_i, X_j 间互不相关 \Leftrightarrow 各分量相互独立。

协方差矩阵 $\boldsymbol{\Sigma}$ 是对角的 $\Rightarrow \mathbf{X}$ 各分量相互独立, 且正态分布

下一节 主题

如何估计类条件概率密度函数？

主要内容

1. 引言

2. 贝叶斯决策模型

2.1 最小错误率贝叶斯决策

2.2 最小风险的贝叶斯决策

3. 正态分布的概率密度函数及性质

4. 概率/概率密度函数估计

5. 朴素贝叶斯(Naive Bayes)分类



PART1. 问题的引入

➤ 如何基于贝叶斯决策解决实际问题？

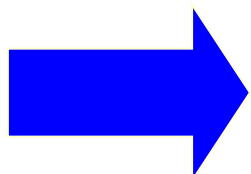
基本思路：

[1] 利用有限规模训练样本, **估计**

$$\begin{cases} \text{先验概率} & P(\omega_i) \\ \text{条件概率密度} & p(x|\omega_i) \end{cases} \quad i=1,2,\dots,c$$

[2] $\left\{ \begin{array}{l} \text{利用估计的 } \hat{P}(\omega_i), \hat{p}(x|\omega_i), \text{ 设计贝叶斯分类器;} \\ \text{对未知样本 } x \text{ 进行判决} \end{array} \right.$

上述过程又称**基于样本的两步贝叶斯决策**。



问题：如何估计有关概率/概率密度函数？

➤ 如何进行有关概率/概率密度函数的估计?

前提 { 训练样本数量足够多
训练样本分布能代表样本的真实分布(如 : *i.i.d*)

$$\Rightarrow \begin{cases} \hat{p}(x|\omega_i) \xrightarrow{N \rightarrow \infty} p(x|\omega_i) \\ \hat{P}(\omega_i) \xrightarrow{N \rightarrow \infty} P(\omega_i) \end{cases}$$

[1] 状态先验概率 $P(\omega_i)$, $i=1,2,\dots,c$ 的估计

比较容易 { **依靠经验**
如：异常细胞识别，医生可根据
以往细胞病理检查统计结果，
做出推断
利用训练数据中各类出现的频度估计

[2] 类条件概率密度 $p(x|\omega_i)$, $i=1,2,\dots,c$ 的估计

客观情况 { (1) 应包含随机变量的全部信息
(2) 可以是满足如下条件的任何函数
 $p(x) \geq 0, \quad \int p(x) dx = 1$

实际情况 { (1) 训练样本数目不够多
(2) 噪声污染、甚至部分特征丢失
(3) 特征维数对分类器计算复杂度的影响

⇒ 类条件概率密度估计非常困难

➤ 概率密度函数估计基本方法



[1] 参数估计 (*parametric estimation*)

类条件总体概率密度函数

- 形式已知 根据对问题一般认识，
假设随机变量 x 服从某种分布
- 参数未知 利用训练样本估计分布参数

如：若 $p(x | \omega_i) \sim N(\mu, \Sigma)$ ，则待估计参数为 $\theta = (\mu, \Sigma)$

两种类型的参数估计法

- 监督参数估计 (*supervised parametric estimation*)
各训练样本类别状态已知
- 非监督参数估计 (*nonsupervised parametric estimation*)
各训练样本类别状态未知



概率密度函数的估计 \Rightarrow 参数估计

[2] 非参数估计 (*nonparametric estimation*)

概率密度函数形式未知;

利用训练样本, 直接推断概率密度函数

PART2. 概率密度函数的参数估计法之一

----最大似然估计

(Maximum Likelihood Estimation, MLE)

1.最大似然估计的问题描述

假设条件

- [1] 参数 θ : 确定的未知量 (θ 不是随机变量)
- [2] c 类 $\left\{ \begin{array}{l} c \text{个确定类别的样本集 } \mathcal{X}_i, i = 1, \dots, c. \\ \text{各样本集的样本个数 } N_i, i = 1, \dots, c \end{array} \right.$
- [3] 独立同分布 (i.i.d : *independent identical distribution*)
样本集 \mathcal{X}_i 中, 各样本依 $p(x | \omega_i)$ 独立抽取, $i = 1, \dots, c$
- [4] 类条件概率密度 $p(x | \omega_i), i = 1, \dots, c$ $\left\{ \begin{array}{l} \text{函数形式确定} \\ \text{参数 } \theta_i \text{未知} \end{array} \right.$
 $p(x | \omega_i)$ 依赖于 θ_i , 记: $p(x | \omega_i; \theta_i)$, 或 $p(x; \theta_i)$
- [5] 各类样本只含本类的分布信息
各参数 θ_i 在函数上相互独立, θ_i 的估计不受 $\mathcal{X}_j (j \neq i)$ 影响。

⇒

目标 根据各训练样本集 $\mathcal{X}_1, \dots, \mathcal{X}_c$, 分别估计各参数(向量)
 $\theta_1, \dots, \theta_c$ 的“最可能”的取值, 记为 $\hat{\theta}_1, \dots, \hat{\theta}_c$

2.最大似然估计的基本思想

转化为c个独立问题，各问题均可表述为

已知：

某类训练样本集 $\mathcal{X} = \{x_1, \dots, x_N\}$;

各样本按已知整体分布形式 $p(x; \theta)$ 独立抽取 (i.i.d);

参数向量 $\theta = [\theta_1, \dots, \theta_s]^T$

目标：

采用最大似然估计法，确定该类参数 θ 的估计 $\hat{\theta}$.

2.最大似然估计的基本思想 (续1)



河北师范大学软件学院
Software College of Hebei Normal University

(1) 似然函数 (*likelihood function*)

$$l(\theta) = p(\mathcal{X}; \theta) = p(x_1, \dots, x_N; \theta) = \prod_{i=1}^N p(x_i; \theta)$$

参数 θ 下, 观测到样本集 $\mathcal{X} = \{x_1, \dots, x_N\}$ 的概率

或: 参数 θ 下, N 个独立随机样本 x_1, \dots, x_N 的联合概率。

样本集 \mathcal{X} 固定时, $l(\theta)$ 随 θ 的取值而变, 是 θ 的函数。

$\Rightarrow l(\theta)$ 体现了参数 θ 下取得样本集 \mathcal{X} 的可能性
是 θ 关于 \mathcal{X} 的似然函数

2.最大似然估计的基本思想 (续2)



河北师范大学软件学院
Software College of Hebei Normal University

(2)对数似然函数

$$H(\theta) = \ln l(\theta) = \ln p(\mathcal{X}; \theta) = \sum_{i=1}^N \ln p(x_i; \theta)$$

(3)最大似然估计量

$$\text{若 } \hat{\theta} = \arg \max_{\theta \in \Theta} [l(\theta)] = \arg \max_{\theta \in \Theta} [\ln l(\theta)]$$

则 $\hat{\theta}$ 是 θ 的最大似然估计量。

$$\text{记: } \hat{\theta} = d(x_1, \dots, x_N) = d(\mathcal{X})$$

$$\text{或: } \hat{\theta} = \hat{\theta}(x_1, \dots, x_N) = \hat{\theta}(\mathcal{X})$$

2.最大似然估计的基本思想 (续3)



河北师范大学软件学院
Software College of Hebei Normal University

最大似然估计的实质：

根据抽取的 N 个样本 x_1, \dots, x_N ，估计它们
"最可能"来自哪一个密度函数.

对数似然函数 $H(\theta)$ 关于 $l(\theta)$ 是单调
增加的，所以最大似然估计量 $\hat{\theta}$ 必然使
 $H(\theta)$ 、 $l(\theta)$ 同时最大.

3.最大似然估计的必要条件



对于密度函数

$$p(\mathbf{x}; \theta)$$

待估计的参数向量

$$\theta = [\theta_1, \dots, \theta_s]^T$$

似然函数

$$l(\theta) = p(\mathcal{X}; \theta) = p(\mathbf{x}_1, \dots, \mathbf{x}_N; \theta) = \prod_{i=1}^N p(\mathbf{x}_i; \theta)$$

对数似然函数

$$\begin{aligned} H(\theta) &= \ln l(\theta) = \ln p(\mathbf{x}_1, \dots, \mathbf{x}_N; \theta) \\ &= \sum_{k=1}^N \ln p(\mathbf{x}_k; \theta) = \sum_{k=1}^N \ln p(\mathbf{x}_k; \theta_1, \dots, \theta_s) \end{aligned}$$

最大似然估计量

$$\hat{\theta} = \arg \max_{\theta \in \Theta} [l(\theta)] = \arg \max_{\theta \in \Theta} [H(\theta)]$$

记 梯度算子 $\nabla_{\theta} = \left[\frac{\partial}{\partial \theta_1} \quad \cdots \quad \frac{\partial}{\partial \theta_s} \right]^T$

若 似然函数 $l(\theta)$, $H(\theta)$ 关于 θ **连续可微**, 则 $\nabla_{\theta} H(\theta) = 0$

$$\begin{aligned} \nabla_{\theta} H(\theta) &= \begin{bmatrix} \frac{\partial H(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial H(\theta)}{\partial \theta_s} \end{bmatrix} = \begin{bmatrix} \sum_{k=1}^N \frac{\partial}{\partial \theta_1} [\ln p(x_k; \theta)] \\ \vdots \\ \sum_{k=1}^N \frac{\partial}{\partial \theta_s} [\ln p(x_k; \theta)] \end{bmatrix} \\ &= \sum_{k=1}^N \begin{bmatrix} \frac{\partial}{\partial \theta_1} [\ln p(x_k; \theta)] \\ \vdots \\ \frac{\partial}{\partial \theta_s} [\ln p(x_k; \theta)] \end{bmatrix} = \sum_{k=1}^N \nabla_{\theta} [\ln p(x_k; \theta)] \end{aligned}$$

可得最大似然估计的必要条件：

最大似然估计量 $\hat{\theta}$ 必满足

$$\nabla_{\theta} H(\theta) = \sum_{k=1}^N \nabla_{\theta} [\ln p(x_k; \theta)] = 0$$

或 $\nabla_{\theta} l(\theta) = 0$

$$\nabla_{\theta} H(\theta) = 0, \quad \text{即方程组} \begin{cases} \sum_{k=1}^N \frac{\partial}{\partial \theta_1} [\ln p(x_k; \theta)] = 0 \\ \vdots \\ \sum_{k=1}^N \frac{\partial}{\partial \theta_s} [\ln p(x_k; \theta)] = 0 \end{cases}$$

若上述方程组的某个解 $\hat{\theta}$ 能使 $H(\theta)$ 最大，则 $\hat{\theta}$ 为 θ 的最大似然估计。

PART3. 单变量正态分布概率密度函数的最大似然估计

某类条件概率密度函数 $p(x | \omega_i; \boldsymbol{\theta}) = p(x; \boldsymbol{\theta})$

$$p(x; \boldsymbol{\theta}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right] = \frac{1}{\sqrt{2\pi\theta_2}} \exp\left[-\frac{(x - \theta_1)^2}{2\theta_2}\right]$$

待估计参数 $\boldsymbol{\theta} = [\theta_1, \theta_2]^T = [\mu, \sigma^2]^T$

该类样本集 $\mathcal{X} = \{x_1, \dots, x_N\}$ 各样本独立抽取 (*i.i.d*)

似然函数

$$l(\boldsymbol{\theta}) = p(\mathcal{X}; \boldsymbol{\theta}) = \prod_{i=1}^N p(x_i; \boldsymbol{\theta}) = \frac{1}{(2\pi\theta_2)^{\frac{N}{2}}} \exp\left[-\frac{\sum_{i=1}^N (x_i - \theta_1)^2}{2\theta_2}\right]$$

对数似然函数 $H(\boldsymbol{\theta}) = \ln l(\boldsymbol{\theta}) = -\frac{\sum_{i=1}^N (x_i - \theta_1)^2}{2\theta_2} - \frac{N}{2} \ln(2\pi\theta_2)$

由于 $H(\theta) = \ln l(\theta) = -\frac{\sum_{i=1}^N (x_i - \theta_1)^2}{2\theta_2} - \frac{N}{2} \ln(2\pi\theta_2)$

所以 $\nabla_{\theta} H(\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} H(\theta) \\ \frac{\partial}{\partial \theta_2} H(\theta) \end{bmatrix} = \begin{bmatrix} \frac{1}{\theta_2} \sum_{i=1}^N (x_i - \theta_1) \\ -\frac{N}{2\theta_2} + \frac{\sum_{i=1}^N (x_i - \theta_1)^2}{2\theta_2^2} \end{bmatrix}$

最大似然估计量 $\hat{\theta}$ 满足方程: $\nabla_{\theta} H(\theta) = 0$

即 $\begin{cases} \sum_{i=1}^N \frac{1}{\hat{\theta}_2} (x_i - \hat{\theta}_1) = 0 \\ \sum_{i=1}^N \left[-\frac{1}{2\hat{\theta}_2} + \frac{(x_i - \hat{\theta}_1)^2}{2\hat{\theta}_2^2} \right] = 0 \end{cases}$

解得

$$\begin{cases} \hat{\mu} = \hat{\theta}_1 = \frac{1}{N} \sum_{i=1}^N x_i \\ \hat{\sigma}^2 = \hat{\theta}_2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2 \end{cases}$$

样本集 $\mathcal{X} = \{x_1, \dots, x_N\}$ 各样本依 $p(x; \theta)$ 独立抽取

$$\begin{cases} \text{样本均值} & \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \\ \text{样本方差} & s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 \end{cases}$$



河北师范大学软件学院
Software College of Hebei Normal University

参数 μ 、 σ^2 的最大似然估计

$$\begin{cases} \hat{\mu} = \bar{x} \\ \hat{\sigma}^2 = \frac{1}{N} \sum_{k=1}^N (x_k - \hat{\mu})^2 \end{cases}$$

PART4. 多元正态分布概率密度函数的 最大似然估计

- 一般情况
- 各特征相互独立

对于多变量正态分布 μ, Σ 均未知



d 维特征空间

样本集 $\mathcal{X} = \{x_1, \dots, x_N\}$ $x_i = [x_i^{(1)}, \dots, x_i^{(d)}]^T, i = 1, \dots, N$

待估计的参数 θ $\left\{ \begin{array}{l} \text{均值向量} \quad \mu = [\mu^{(1)}, \mu^{(2)}, \dots, \mu^{(d)}]^T \\ \text{协方差矩阵} \quad \Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2d} \\ \dots & \dots & \dots & \dots \\ \sigma_{d1} & \sigma_{d2} & \dots & \sigma_{dd} \end{bmatrix} \end{array} \right.$

其中: $\sigma_{ij} = \sigma_{ji}$

概率密度函数

$$p(\mathbf{x}; \theta) = p(\mathbf{x}; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right]$$

样本集 $\mathcal{X} = \{x_1, \dots, x_N\}$ 各样本依 $p(x | \theta)$ 独立抽取

$$\left\{ \begin{array}{ll} \text{样本均值} & \bar{x} = \frac{1}{N} \sum_{k=1}^N x_k \\ \text{样本协方差矩阵} & C = \frac{1}{N-1} \sum_{k=1}^N (x_k - \bar{x})(x_k - \bar{x})^T \end{array} \right.$$

最大似然参数估计

$$\left\{ \begin{array}{ll} \text{期望} & \hat{\mu} = \frac{1}{N} \sum_{k=1}^N x_k = \bar{x} \\ \text{协方差矩阵} & \hat{\Sigma} = \frac{1}{N} \sum_{k=1}^N (x_k - \hat{\mu})(x_k - \hat{\mu})^T \end{array} \right.$$



PART4. 多元正态分布概率密度函数的 最大似然估计

- 一般情况
- 各特征相互独立

对于多变量正态分布

μ 、 Σ 均未知



$$\Sigma = \begin{bmatrix} \sigma_{11} & & & \\ & \sigma_{22} & & \\ & & \dots\dots & \\ & & & \sigma_{dd} \end{bmatrix} = \begin{bmatrix} \sigma^{(1)2} & & & \\ & \sigma^{(2)2} & & \\ & & \dots\dots & \\ & & & \sigma^{(d)2} \end{bmatrix}$$

概率密度函数

$$p_X(\mathbf{x}; \boldsymbol{\theta}) = p_X(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

$$= \prod_{i=1}^d \frac{1}{\sqrt{2\pi} \sigma^{(i)}} \exp \left[-\frac{1}{2} \left(\frac{x^{(i)} - \mu^{(i)}}{\sigma^{(i)}} \right)^2 \right] = \prod_{i=1}^d p_{X^{(i)}}(x^{(i)}; \mu^{(i)}, \sigma^{(i)2})$$

$$p_{X^{(i)}}(x^{(i)}; \mu^{(i)}, \sigma^{(i)2}) = \frac{1}{\sqrt{2\pi} \sigma^{(i)}} \exp \left[-\frac{1}{2} \left(\frac{x^{(i)} - \mu^{(i)}}{\sigma^{(i)}} \right)^2 \right]$$

随机变量 $\mathbf{X}^{(i)}$ 的边缘密度

$$p_{\mathbf{X}^{(i)}}(x^{(i)}; \mu^{(i)}, \sigma^{(i)2}) = \frac{1}{\sqrt{2\pi}\sigma^{(i)}} \exp\left[-\frac{1}{2}\left(\frac{x^{(i)} - \mu^{(i)}}{\sigma^{(i)}}\right)^2\right]$$

随机向量 \mathbf{X} 的概率密度函数: $p_{\mathbf{X}}(\mathbf{x}; \mu, \Sigma) = \prod_{i=1}^d p_{\mathbf{X}^{(i)}}(x^{(i)}; \mu^{(i)}, \sigma^{(i)2})$

边缘密度 $p_{\mathbf{X}^{(i)}}(x^{(i)}; \mu^{(i)}, \sigma^{(i)2})$ 最大似然参数估计

$$\begin{aligned} &\left\{ \begin{array}{l} \text{期望} \\ \text{方差} \end{array} \right. \quad \begin{aligned} \hat{\mu}^{(i)} &= \frac{1}{N} \sum_{k=1}^N x_k^{(i)} \\ \widehat{\sigma^{(i)2}} &= \frac{1}{N} \sum_{k=1}^N \left(x_k^{(i)} - \hat{\mu}^{(i)} \right)^2 \end{aligned} \\ \Rightarrow &\left\{ \begin{array}{l} \text{期望向量 } \mu \\ \text{协方差矩阵 } \hat{\Sigma} \end{array} \right. \quad \begin{aligned} \hat{\mu} &= \left[\hat{\mu}^{(1)}, \hat{\mu}^{(2)}, \dots, \hat{\mu}^{(d)} \right]^T \\ \hat{\Sigma} &= \text{diag} \left(\widehat{\sigma^{(1)2}}, \widehat{\sigma^{(2)2}}, \dots, \widehat{\sigma^{(d)2}} \right) \end{aligned} \end{aligned}$$

基于样本的两步贝叶斯决策

[1] 利用有限规模训练样本 $\{(x_j, y_j), j=1, \dots, N\}$,

设计贝叶斯分类器. 估计

$\left\{ \begin{array}{l} \text{先验概率} \quad P(\omega_i) \\ \text{条件概率密度 } p(x|\omega_i) \text{ 或 条件概率 } P(x|\omega_i) \end{array} \right.$
 $i=1, 2, \dots, c$

[2] $\left\{ \begin{array}{l} \text{利用估计的 } \hat{P}(\omega_i), \hat{p}(x|\omega_i) \text{ 或 } \hat{P}(x|\omega_i) \\ \text{对未知样本 } x \text{ 进行判决} \end{array} \right.$



若观测样本 \mathbf{x} 所在特征空间为连续的:

$$P(\omega_j | \mathbf{x}) = \frac{\hat{P}(\omega_j) \hat{p}(\mathbf{x} | \omega_j)}{p(\mathbf{x})} = \frac{\hat{P}(\omega_j) \hat{p}(\mathbf{x} | \omega_j)}{\sum_{i=1}^c \hat{P}(\omega_i) \hat{p}(\mathbf{x} | \omega_i)}$$
$$j = 1, 2, \dots, c$$

若观测样本 \mathbf{x} 所在特征空间为离散的:

$$P(\omega_j | \mathbf{x}) = \frac{P(\omega_j) P(\mathbf{x} | \omega_j)}{P(\mathbf{x})} = \frac{P(\omega_j) P(\mathbf{x} | \omega_j)}{\sum_{i=1}^c P(\omega_i) P(\mathbf{x} | \omega_i)}$$
$$j = 1, 2, \dots, c$$

1. 引言

2. 贝叶斯决策模型

2.1 最小错误率贝叶斯决策

连续情况；离散情况

2.2 最小风险的贝叶斯决策

连续情况；离散情况

3. 正态分布的概率密度函数及性质

4. 概率/概率密度函数估计

5. 朴素贝叶斯(Naive Bayes)分类

朴素贝叶斯法，基于**贝叶斯公式**、以及**各特征条件独立**的**假设**，实现贝叶斯决策。

连续特征空间： $p(x | \omega_j) = \prod_{k=1}^d p(x^{(k)} | \omega_j)$

离散特征空间： $P(x | \omega_j) = \prod_{k=1}^d P(x^{(k)} | \omega_j)$

PART1. 离散特征空间的 朴素贝叶斯(Naive Bayes)分类 ——算法描述

输入：

(1) 训练样本集 $\{(\mathbf{x}_j, y_j), j=1, \dots, N\}$,

$$\text{其中 } \mathbf{x}_j = [x_j^{(1)}, \dots, x_j^{(d)}]^T$$

$x_j^{(k)}$ 是第 j 个样本的第 k 个特征;

并且 $x_j^{(k)} \in \{a_{k1}, a_{k2}, \dots, a_{kS_k}\}, k \in \{1, 2, \dots, d\}$

$$y_j \in \{\omega_1, \omega_2, \dots, \omega_C\}$$

(2) 待决策的观测样本 $\mathbf{x} = [x^{(1)}, \dots, x^{(d)}]^T$

输出： 观测样本 \mathbf{x} 的类别

实现步骤：

(1)利用训练样本集估计先验概率及条件概率

估计方式1--最大似然估计法

$$X = [X^{(1)}, \dots, X^{(d)}]^T$$

$$\hat{P}(\omega = \omega_i) = \frac{\sum_{j=1}^N I(y_j = \omega_i)}{N}, \quad i \in \{1, 2, \dots, C\}$$

$$\hat{P}(X^{(k)} = a_{kl} | \omega_i) = \frac{\hat{P}(X^{(k)} = a_{kl}, \omega = \omega_i)}{\hat{P}(\omega = \omega_i)} = \frac{\sum_{j=1}^N I(x_j^{(k)} = a_{kl}, y_j = \omega_i) / N}{\sum_{j=1}^N I(y_j = \omega_i) / N}$$

$$= \frac{\sum_{j=1}^N I(x_j^{(k)} = a_{kl}, y_j = \omega_i)}{\sum_{j=1}^N I(y_j = \omega_i)} \quad \begin{cases} k \in \{1, 2, \dots, d\} \\ l \in \{1, 2, \dots, S_k\} \\ i \in \{1, 2, \dots, C\} \end{cases}$$

可能会导致某个估计结果为0，影响后续的后验概率计算，以及决策。

实现步骤：(1)利用训练样本集估计先验概率及条件边缘概率

估计方式2--贝叶斯估计法

($\lambda=1$, 称**LAPLACE平滑**; $0<\lambda<1$, **Lidstone平滑**)

$$\hat{P}(\omega=\omega_i) = \frac{\lambda + \sum_{j=1}^N I(y_j = \omega_i)}{C\lambda + N}, \quad i \in \{1, 2, \dots, C\}$$

MultinomialNB

$$\hat{P}(X^{(k)} = a_{kl} | \omega_i) = \frac{\lambda + \sum_{j=1}^N I(x_j^{(k)} = a_{kl}, y_j = \omega_i)}{S_k\lambda + \sum_{j=1}^N I(y_j = \omega_i)} \quad \begin{cases} k \in \{1, 2, \dots, d\} \\ l \in \{1, 2, \dots, S_k\} \\ i \in \{1, 2, \dots, C\} \end{cases}$$

(多项朴素贝叶斯, Multinomial Naive Bayes)

Multinomial Naive Bayes 基于各类多项分布假设(各特征只能取有限个离散值之一), 是面向文档分类中的两种经典朴素贝叶斯模型之一(另一种为 **Bernoulli Naive Bayes**)。

(2) 对于给定的观测样本 $\mathbf{x} = [x_1, \dots, x_d]^T$, 计算

$$\hat{P}(\omega = \omega_j) \hat{P}(X = \mathbf{x} | \omega = \omega_j) \quad j \in \{1, 2, \dots, C\}$$

$$\begin{aligned} & \hat{P}(\omega = \omega_j) \hat{P}(X = \mathbf{x} | \omega = \omega_j) \\ &= \hat{P}(\omega = \omega_j) \prod_{k=1}^d \hat{P}(X^{(k)} = x^{(k)} | \omega = \omega_j) \end{aligned}$$

(3) 确定观测样本 \mathbf{x} 的预测类别 y

$$y = \arg \max_{\omega_j} \left[\hat{P}(\omega = \omega_j) \hat{P}(X = \mathbf{x} | \omega = \omega_j) \right]$$

试由表的训练数据学习一个朴素贝叶斯分类器并确定 $x = (2, S)^T$ 类标记 y 。
表中 $X^{(1)}$, $X^{(2)}$ 为特征, 取值的集合分别为 $A_1 = \{1, 2, 3\}$, $A_2 = \{S, M, L\}$;
 Y 为类标记, $Y \in C = \{1, -1\}$

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$X^{(1)}$	1	1	1	1	1	2	2	2	2	2	3	3	3	3	3
$X^{(2)}$	S	M	M	S	S	S	M	M	L	L	L	M	M	L	L
Y	-1	-1	1	1	-1	-1	-1	1	1	1	1	1	1	1	-1

先验概率估计 ($\lambda=1$, 类别数=2): $\hat{P}(Y=c_i) = \frac{\lambda + \sum_{j=1}^{15} I(y_j=c_i)}{2\lambda+15} \quad \{c_1=1, c_2=-1\}$

条件边缘概率估计 $\hat{P}(X^{(k)}=a_{kl} | Y=c_i) = \frac{\lambda + \sum_{j=1}^{15} I(x_j^{(k)}=a_{kl}, y_j=c_i)}{s_k \lambda + \sum_{j=1}^{15} I(y_j=c_i)} \quad l \in \{1, \dots, s_k\}$

$X^{(1)} \in A_1 = \{1, 2, 3\}, s_1 = 3 \quad X^{(2)} \in A_2 = \{S, M, L\}, s_2 = 3$

解：首先进行先验概率、以及各类条件概率估计

方式1：最大似然估计

$$P(Y=1)=\frac{9}{15}, \quad P(Y=-1)=\frac{6}{15}$$

$$P(X^{(1)}=1|Y=1)=\frac{2}{9}, \quad P(X^{(1)}=2|Y=1)=\frac{3}{9}, \quad P(X^{(1)}=3|Y=1)=\frac{4}{9}$$

$$P(X^{(2)}=S|Y=1)=\frac{1}{9}, \quad P(X^{(2)}=M|Y=1)=\frac{4}{9}, \quad P(X^{(2)}=L|Y=1)=\frac{4}{9}$$

$$P(X^{(1)}=1|Y=-1)=\frac{3}{6}, \quad P(X^{(1)}=2|Y=-1)=\frac{2}{6}, \quad P(X^{(1)}=3|Y=-1)=\frac{1}{6}$$

$$P(X^{(2)}=S|Y=-1)=\frac{3}{6}, \quad P(X^{(2)}=M|Y=-1)=\frac{2}{6}, \quad P(X^{(2)}=L|Y=-1)=\frac{1}{6}$$

对于给定的 $x=(2,S)^T$ 计算：

基于上述概率
信息，对观测
样本 x 进行类
别决策：

$$P(Y=1)P(X^{(1)}=2|Y=1)P(X^{(2)}=S|Y=1)=\frac{9}{15} \cdot \frac{3}{9} \cdot \frac{1}{9} = \frac{1}{45}$$

$$P(Y=-1)P(X^{(1)}=2|Y=-1)P(X^{(2)}=S|Y=-1)=\frac{6}{15} \cdot \frac{2}{6} \cdot \frac{3}{6} = \frac{1}{15}$$

因为 $P(Y=-1)P(X^{(1)}=2|Y=-1)P(X^{(2)}=S|Y=-1)$ 最大，所以 $y=-1$.

解：首先进行先验概率、以及各类条件概率估计

$$P(Y=1)=\frac{10}{17}, \quad P(Y=-1)=\frac{7}{17}$$

方式2：LAPLACE平滑

$$P(X^{(1)}=1|Y=1)=\frac{3}{12}, \quad P(X^{(1)}=2|Y=1)=\frac{4}{12}, \quad P(X^{(1)}=3|Y=1)=\frac{5}{12}$$

$$P(X^{(2)}=S|Y=1)=\frac{2}{12}, \quad P(X^{(2)}=M|Y=1)=\frac{5}{12}, \quad P(X^{(2)}=L|Y=1)=\frac{5}{12}$$

$$P(X^{(1)}=1|Y=-1)=\frac{4}{9}, \quad P(X^{(1)}=2|Y=-1)=\frac{3}{9}, \quad P(X^{(1)}=3|Y=-1)=\frac{2}{9}$$

$$P(X^{(2)}=S|Y=-1)=\frac{4}{9}, \quad P(X^{(2)}=M|Y=-1)=\frac{3}{9}, \quad P(X^{(2)}=L|Y=-1)=\frac{2}{9}$$

基于上述概率信息，对观测样本 x 进行类别决策：

对于给定的 $x=(2,S)^T$ 计算：

$$P(Y=1)P(X^{(1)}=2|Y=1)P(X^{(2)}=S|Y=1)=\frac{10}{17} \cdot \frac{4}{12} \cdot \frac{2}{12} = \frac{5}{153} = 0.0327$$

$$P(Y=-1)P(X^{(1)}=2|Y=-1)P(X^{(2)}=S|Y=-1)=\frac{7}{17} \cdot \frac{3}{9} \cdot \frac{4}{9} = \frac{28}{459} = 0.0610$$

由于 $P(Y=-1)P(X^{(1)}=2|Y=-1)P(X^{(2)}=S|Y=-1)$ 最大，所以 $y=-1$ 。

PART2. 连续特征空间的 朴素贝叶斯(Naive Bayes)分类 ——算法描述

输入：

(1) 训练样本集 $\mathbf{D} = \{(\mathbf{x}_j, y_j), j=1, \dots, N\}$,

$$\text{其中 } \mathbf{x}_j = [x_j^{(1)}, \dots, x_j^{(d)}]^T$$

$x_j^{(k)}$ 是第 j 个样本的第 k 个特征;

$$k \in \{1, 2, \dots, d\}, \quad y_j \in \{\omega_1, \omega_2, \dots, \omega_C\}$$

(2) 待决策的观测样本 $\mathbf{x} = [x^{(1)}, \dots, x^{(d)}]^T$

输出： 观测样本 \mathbf{x} 的类别

实现步骤：

(1)利用训练样本集估计**先验概率**及**条件概率密度函数**

估计**先验概率**:
$$\hat{P}(\omega=\omega_i)=\frac{\sum_{j=1}^N I(y_j=\omega_i)}{N}, \quad i \in \{1, 2, \dots, C\}$$

估计**条件密度**:
$$\hat{p}(X=\mathbf{x}|\omega=\omega_i)=\prod_{k=1}^d \hat{p}(X^{(k)}=x^{(k)}|\omega=\omega_i)$$

分别利用训练样本集内第 ω_i 类所有样本的第 k 个特征取值
确定估计的**条件边缘密度**: $\hat{p}(X^{(k)}=x^{(k)}|\omega=\omega_i)$

$$k \in \{1, 2, \dots, d\}; \quad i \in \{1, 2, \dots, C\}$$

特别地，若为**Gaussian Naive Bayes**分类模型

GaussianNB

$$p(\mathbf{X} = \mathbf{x} \mid \omega = \omega_i) = \prod_{k=1}^d p(X^{(k)} = x^{(k)} \mid \omega = \omega_i)$$

$$\text{并且 } p(X^{(k)} = x^{(k)} \mid \omega = \omega_i) = \frac{1}{\sqrt{2\pi}\sigma_i^{(k)}} \exp\left[-\frac{1}{2}\left(\frac{x^{(k)} - \mu_i^{(k)}}{\sigma_i^{(k)}}\right)^2\right]$$

$$\text{所以 } \hat{p}(X^{(k)} = x^{(k)} \mid \omega = \omega_i) = \frac{1}{\sqrt{2\pi}\hat{\sigma}_i^{(k)}} \exp\left[-\frac{1}{2}\left(\frac{x^{(k)} - \hat{\mu}_i^{(k)}}{\hat{\sigma}_i^{(k)}}\right)^2\right]$$

期望

$$\hat{\mu}_i^{(k)} = \frac{1}{N_i} \sum_{\substack{(\mathbf{x}_j, y_j) \in D \\ \text{并且 } y_j = \omega_i}} \mathbf{x}_j^{(k)}$$

$$k \in \{1, 2, \dots, d\}$$

方差

$$\widehat{\sigma_i^{(k)}}^2 = \frac{1}{N_i} \sum_{\substack{(\mathbf{x}_j, y_j) \in D \\ \text{并且 } y_j = \omega_i}} \left(\mathbf{x}_j^{(k)} - \hat{\mu}_i^{(k)} \right)^2$$

$$i \in \{1, 2, \dots, C\}$$

(2) 对于给定的观测样本 $\mathbf{x} = [x_1, \dots, x_d]^T$, 计算

$$\hat{P}(\omega = \omega_j) \hat{p}(X = \mathbf{x} | \omega = \omega_j) \quad j \in \{1, 2, \dots, C\}$$

$$\hat{P}(\omega = \omega_j) \hat{p}(X = \mathbf{x} | \omega = \omega_j)$$

$$= \hat{P}(\omega = \omega_j) \prod_{k=1}^d \hat{p}(X^{(k)} = x^{(k)} | \omega = \omega_j)$$

(3) 确定观测样本 \mathbf{x} 的预测类别 y

$$y = \arg \max_{\omega_j} \left[\hat{P}(\omega = \omega_j) \hat{p}(X = \mathbf{x} | \omega = \omega_j) \right]$$



1. 面向两类别/多类别分类问题的两种贝叶斯分类模型决策规则是什么？
 - (1) 基于最小错误率的贝叶斯分类，决策规则；
 - (2) 基于最小风险的贝叶斯分类，决策规则。
2. 正确写出单变量/多元正态分布的概率密度函数表达式。
3. (1) 若 d 维随机向量正态分布，并且 N 个观测样本按照正态分布独立抽取得到，请写出 N 个观测样本组成的样本集的似然值；并写出概率密度函数的各参数最大似然估计结果。

(2) 若各特征正确写出征相互独立，请写出正态分布概率密度函数的最大似然估计结果

4. 掌握如下算法：

- (1) 离散特征空间，基于最小错误率的朴素贝叶斯分类模型；
- (2) 连续特征空间，各类别条件概率密度函数正态分布情况下，基于最小错误率的朴素贝叶斯分类模型。