

第4章, 朴素贝叶斯法

提前学习部分

贝叶斯公式:
$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} = P(A) \cdot \frac{P(B|A)}{P(B)}$$

$$P(\text{规律}|\text{现象}) = \frac{P(\text{现象}|\text{规律}) P(\text{规律})}{P(\text{现象})}$$

思考模式

① 先验分布 $\pi(\theta)$ + 样本信息 $X \Rightarrow$ 后验分布 $\pi(\theta|X)$

意味: 新观察到的样本信息将修正人们对事物的认知.

↓
从 $\pi(\theta)$ 修正为 后验分布 $\pi(\theta|X)$

② 边缘概率 (即先验概率)

理解为: 在联合概率中, 把最终结果中那些不需要的事件通过合并成它们的全概率, 从而消去它们。

(离散型随机变量 \rightarrow 求和)
连续 \rightarrow 积分

③ 条件概率 (即后验概率) \leftarrow 指我们观测到的样本的分布, 也就是似然函数.

④ 联合概率 表示 2 个事件共同发生的概率

⑤ 后验概率 = 先验^{概率} ~~概率~~ ^率 \times 调整因子

$$\begin{array}{ccc} \downarrow & \downarrow & \downarrow \\ P(A|B) & P(A) & \frac{P(B|A)}{P(B)} \end{array}$$

先预估一个“先验概率”, 然后加入实验结果, 看是增强还是削弱了“先验概率”

宏观上来看

↓ 决策

$$Y = f(X)$$

条件概率
max

$$P(Y|X)$$

生成

$$P(Y|X) = \frac{P(X, Y)}{P(X)}$$

判别

$$Y = f(X), P(Y|X)$$

$$Y = f(X) \Rightarrow P(Y|X) \Rightarrow P(Y|X) = \frac{P(X, Y)}{P(X)}$$

↓
不考虑 X, Y
随机性

感知机

↓
考虑 Y 随机性

↓
考虑 X 与 Y 随机性

朴素贝叶斯公式

推导: 用极大似然法估计贝叶斯参数
~~贝叶斯估计 (4.2.3)~~ 为什么 $P(Y=c_k) = \frac{m_k}{N}$

① $Y \in \{c_1, c_2, \dots, c_k\}$

$\begin{matrix} | & | & \dots & | \\ \theta_1 & \theta_2 & & \theta_k \end{matrix}$

θ (向量)

意为 Y 取 c_i 时概率为

$\sum_{i=1}^k \theta_i = 1$

θ_i

② $P(Y=y|\theta) = \frac{I(y=c_1)}{\theta_1} \times \frac{I(y=c_2)}{\theta_2} \times \dots \times \frac{I(y=c_k)}{\theta_k}$

意为给定 θ , 问你 y 是不是等于 c_i , $I(y=c_i)$ 是一个事件函数

③ 在极大似然估计中, 拿到 y_1, y_2, \dots, y_n 样本

根据样本写出样本的概率分布 y_i 是独立的

$P(y_1, y_2, \dots, y_n|\theta) = P(y_1|\theta) \times P(y_2|\theta) \times \dots \times P(y_n|\theta)$

$= \theta_1^{m_1} \times \theta_2^{m_2} \times \dots \times \theta_k^{m_k}$

出现结果 c_i 的次数是 m_i , $\therefore m_1 + \dots + m_k = N$

④ 最大化这个联合概率分布

$\max P(y_1, y_2, \dots, y_n|\theta)$

$= \max \ln P(\sim) = m_1 \ln \theta_1 + m_2 \ln \theta_2 + \dots + m_k \ln \theta_k$

要求使 \uparrow 最大的 θ 的向量 还有一个约束 s.t. $\theta_1 + \dots + \theta_k = 1$

⑤ 要引入拉格朗日乘子 (带约束的优化问题)

$\Rightarrow \max_{\theta_1, \dots, \theta_k} m_1 \ln \theta_1 + m_2 \ln \theta_2 + \dots + m_k \ln \theta_k + \lambda(\theta_1 + \dots + \theta_k - 1)$

分别对 $\theta_1, \dots, \theta_k$ 求导

$\frac{m_1}{\theta_1} + \lambda = 0 \Rightarrow \theta_1 = -\frac{m_1}{\lambda}$

$\frac{m_k}{\theta_k} + \lambda = 0 \Rightarrow \theta_k = -\frac{m_k}{\lambda}$

根据 $\theta_1 + \dots + \theta_k = 1$

$\therefore \text{即 } -\frac{m_1 + \dots + m_k}{\lambda} = 1 \Rightarrow \lambda = -(m_1 + \dots + m_k)$

$\frac{m_1}{N}$

贝叶斯估计 (4.2.3)

例子: 抛硬币 $y \in \{正, 反\}$

推导: 为什么用贝叶斯估计

朴素贝叶斯参数时

$$P_{\lambda}(Y = c_k) = \frac{\sum_{i=1}^N I(y_i = c_k) + \lambda}{N + K\lambda}$$

① 设多项分布时

$$y \in \{c_1, \dots, c_K\}$$

$$\theta_1, \dots, \theta_K$$

先验分布 $\rightarrow p(\theta) = \frac{\Gamma(d_1 + d_k) \theta_1^{d_1-1} \dots \theta_K^{d_K-1}}{\Gamma(d_1) \dots \Gamma(d_K)}$

Dirichlet (狄立克雷分布)

类比

二项分布的先验信息 $p(\theta) = \frac{\Gamma(d+\beta) \theta^{d-1} (1-\theta)^{\beta-1}}{\Gamma(d)\Gamma(\beta)}$

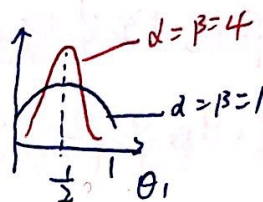
(Beta分布)

②

直观解释

根据右图, 我们^{倾向}认为先验

开始时 $d_1 = d_2 = \dots = d_K = d$



③ ~~拿到~~ 拿到样本后, 算后验

即 $P(\theta | y_1, \dots, y_N)$ 根据贝叶斯公式

$$= \frac{P(\theta, y_1, \dots, y_N)}{P(y_1, \dots, y_N)}$$

—— 联合分布

—— 边缘概率, 与 θ 无关

$$\propto P(\theta) \cdot P(y_1, \dots, y_N | \theta)$$

—— 前面似然函数已求出

$$\propto \theta_1^{d-1} \theta_2^{d-1} \dots \theta_K^{d-1} \cdot \theta_1^{m_1} \theta_1^{m_2} \dots \theta_K^{m_K}$$

$$\propto \theta_1^{m_1+d-1} \theta_2^{m_2+d-1} \dots \theta_K^{m_K+d-1}$$

(狄立克雷分布)

即 $P(\theta | y_1, y_2, \dots, y_N) = \frac{\Gamma(m_1 + d + m_2 + d + \dots + m_K + d)}{\Gamma(m_1 + d) \cdot \Gamma(m_2 + d) \dots \Gamma(m_K + d)} \times \theta_1^{m_1} \theta_2^{m_2} \dots \theta_K^{m_K}$

(LDA)

(4) 最大化 $\theta_1^{m_1+d-1} \theta_2^{m_2+d-1} \dots \theta_k^{m_k+d-1}$

$$\Rightarrow \theta_1 = \frac{m_1+d-1}{m_1+d-1 + m_2+d-1 + \dots + m_k+d-1}$$

$$= \frac{m_1+d-1}{N+Kd-K}$$

$$\lambda = d-1$$

$$= \frac{m_1+d-1}{N+K(d-1)}$$

推导: 为什么后验概率最大化即期望风险最小化

① 在分类问题中的损失函数

$$L(Y, f(x)) = \begin{cases} 1, & Y \neq f(x) \\ 0, & Y = f(x) \end{cases}$$

\uparrow \uparrow
真实 预测

$$f(x) = \arg \max_{c_k} P(Y = c_k | x)$$

把决策函数与条件

概率分布联合起来

$$\min E L(Y, f(x))$$

期望风险就是

损失函数的期望

$$= \sum_Y \sum_x [L(Y, f(x)) P(x, Y)]$$

$$= \sum_Y \sum_x L(Y, f(x)) P(Y|x) \cdot P(x)$$

$$= \sum_x \left[\sum_Y L(Y, f(x)) P(Y|x) \right] P(x)$$

最小化上面式子即最小化

$$\min \sum_Y L(Y, f(x)) P(Y|x)$$

$$\Rightarrow \min \sum_{c_k} L(Y = c_k, f(x)) P(Y = c_k | x)$$

$$\Rightarrow \min \sum_{c_k} I(f(x) \neq c_k) P(Y = c_k | x)$$

事情函数

$Y \neq f(x)$ 为1

$$\Rightarrow \min \sum_{c_k} [1 - I(f(x) = c_k)] P(Y = c_k | x)$$

条件概率之和为1

$$\Rightarrow \min \sum_{c_k} P(Y = c_k | x) - \sum_{c_k} I(f(x) = c_k) \cdot P(Y = c_k | x)$$

$$\Rightarrow \min \left[1 - \sum_{c_k} I(f(x) = c_k) P(Y = c_k | x) \right]$$

$$\text{等价于 } \max \sum_{c_k} I(f(x) = c_k) \cdot P(Y = c_k | x)$$

$$\max \sum_{c_k} I(f(x) = c_k) \cdot \underbrace{P(Y = c_k | x)}$$

★ 等价于 从

中找到一个 c_k ,

使 $P(Y = c_k | x)$ 概率最大且

~~使 $I(f(x) = c_k)$ 事情概率为 1~~

因为 $f(x)$ 一次只能取一个值.

使 $I(f(x) = c_k)$ 事情概率为 1

即等价于

$$f(x) = \arg \max_{c_k} P(Y = c_k | x)$$

∴ 得证

总结 朴素贝叶斯法

生成模型

$$P(Y=c_k | X=x) = \frac{P(Y=c_k) \cdot P(X=x | Y=c_k)}{P(X=x)}$$

模型假设 - 条件独立性 (∴ 朴素)

$$P(X=x | Y=c_k) = \prod_{i=1}^n P(x^{(i)} = x^{(i)} | Y=c_k)$$

预测准则: 后验概率最大

$$y = \arg \max_{c_k} P(Y=c_k | X=x)$$