

朴素贝叶斯分类

朴素贝叶斯分类算法的核心是贝叶斯概率公式

贝叶斯概率公式

有监督的分类

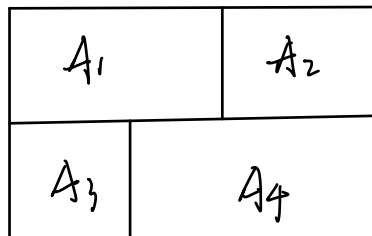
1. 完备事件组

$A_1 \cup A_2 \cup \dots \cup A_n = \Omega$, 且 $A_i \cap A_j = \emptyset$, $1 \leq i \neq j \leq n$

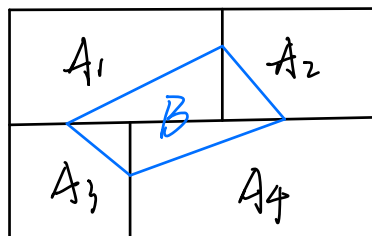
则称 A_1, A_2, \dots, A_n 为一个完备事件组

也就是说 A_1, A_2, \dots, A_n 的并集是完整的事件空间

且两两之间相互独立



假设有一事件 B . 只有在 A_i 事件发生后, B 事件才发生
并且 A_i 事件发生后, B 事件发生的概率可能不相等



先验概率：由简单分析可获得的概率

条件概率：事件 A_i 发生的前提下， B 发生的概率

联合概率：事件 A_i 和 B 同时发生的概率

全概率公式：

$$P(B) = \sum_{i=1}^n P(A_i, B) = \sum_{i=1}^n P(A_i) \cdot P(B|A_i) \quad (1)$$

2. 贝叶斯概率公式

另一个场景： B 事件发生后， A_i 事件发生的概率

$$P(A_i|B) = \frac{P(A_i, B)}{P(B)} = \frac{P(B|A_i) \cdot P(A_i)}{\sum_{i=1}^n P(A_i) \cdot P(B|A_i)} \quad (2)$$

上式即为贝叶斯概率公式。

作用：已知事件 B 发生了，去探求是某个原因 A_i 导致这一结果发生的概率 $P(A_i|B)$

3. 朴素贝叶斯分类算法

3.1 单属性的朴素贝叶斯分类算法

★ 3.2 多属性的朴素贝叶斯分类算法

假设有一个容量为 n 、包含 m 个特征的数据集

$$D = \{ (x_1^{(1)}, x_2^{(1)}, \dots, x_m^{(1)}, y_1), (x_1^{(2)}, x_2^{(2)}, \dots, x_m^{(2)}, y_2), \dots, (x_1^{(n)}, x_2^{(n)}, \dots, x_m^{(n)}, y_n) \}$$

其中 y 是分类标签, $y_i \in \{c_1, c_2, \dots, c_k\}$

通过对 D 进行统计分析, 得到每个类别的

先验概率: $P(Y=c_k), k=1, 2, \dots, k$

以及在每个类别下不同特征的条件概率:

$$P(X_j=x_j | Y=c_k) \quad (\text{在 } c_k \text{ 类中第 } j \text{ 维特征为 } x_j \text{ 的概率})$$

朴素贝叶斯算法有一个前提, 即各个特征间独立

$$\begin{aligned} \text{表示为 } P(X=\vec{x} | Y=c_k) &= P(X_1=x_1, X_2=x_2, \dots, X_m=x_m | c_k) \\ &= \prod_{j=1}^m P(X_j=x_j | c_k) \end{aligned} \quad (3)$$

利用贝叶斯概率公式, 可得

$$\begin{aligned} P(Y=c_k | X=\vec{x}) &= \frac{P(X=\vec{x}, Y=c_k)}{P(X=\vec{x})} = \frac{P(X=\vec{x} | Y=c_k) \cdot P(Y=c_k)}{\sum_{k=1}^K P(Y=c_k) \cdot P(X=\vec{x} | Y=c_k)} \\ &= \frac{\prod_{j=1}^m P(X_j=x_j | Y=c_k) \cdot P(Y=c_k)}{\sum_{k=1}^K P(Y=c_k) \cdot \prod_{j=1}^m P(X_j=x_j | Y=c_k)} \end{aligned} \quad (4)$$

比较各个概率对应的大小, 选择概率最大的一个类别作为分类结果

注意: 在计算过程中无需计算分母, 因为分母不变

全概率公式

$$P(B) = \sum_{i=1}^n P(A_i, B) = \sum_{i=1}^n P(A_i) \cdot P(B|A_i)$$

贝叶斯概率公式

$$\begin{aligned} P(A_i|B) &= \frac{P(A_i, B)}{P(B)} = \frac{P(A_i) \cdot P(B|A_i)}{P(B)} \\ &= \frac{P(A_i) \cdot P(B|A_i)}{\sum_{i=1}^n P(A_i) \cdot P(B|A_i)} \end{aligned}$$

朴素: 各特征之间相互独立