

朴素贝叶斯原理推导 naive bayes

朴素贝叶斯是基于贝叶斯定理和特征条件独立性假设的分类算法。

朴素贝叶斯的概率计算公式：

$$P(C|X) = \frac{P(X|C) P(C)}{\sum P(X|C) P(C)}$$

似然函数 类先验概率
类后验概率 全概率

朴素贝叶斯的学习步骤如下：

首先计算类先验概率分布：

$$P(Y=c_k) = \frac{1}{N} \sum_{i=1}^N I(\tilde{y}_i = c_k), \quad k=1, 2, \dots, K$$

c_k 表示第 k 个类别， \tilde{y}_i 表示第 i 个样本的类标记。类先验概率分布可通过极大似然估计得到

然后计算类条件概率分布：

$$P(X=x | Y=c_k) = P(X^{(1)}=x^{(1)}, \dots, X^{(n)}=x^{(n)} | Y=c_k), \quad k=1, 2, \dots, K$$

直接对 $P(X=x | Y=c_k)$ 进行估计不太行，因为参数太多。但是朴素贝叶斯的一个最重要的假设就是条件独立性假设，即：

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

有了条件独立性假设之后，便可基于极大似然估计计算类条件概率：

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

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

基于上式便可学习一个朴素贝叶斯分类模型。给定新的数据样本时，计算最大后验概率即可：

$$\hat{y} = \arg \max_{c_k} \frac{\prod_{j=1}^n P(X^{(j)}=x^{(j)} | Y=c_k) P(Y=c_k)}{\sum_k \prod_{j=1}^n P(X^{(j)}=x^{(j)} | Y=c_k) P(Y=c_k)}$$

其中分母对所有 c_k 都一样，所以进一步化简为：

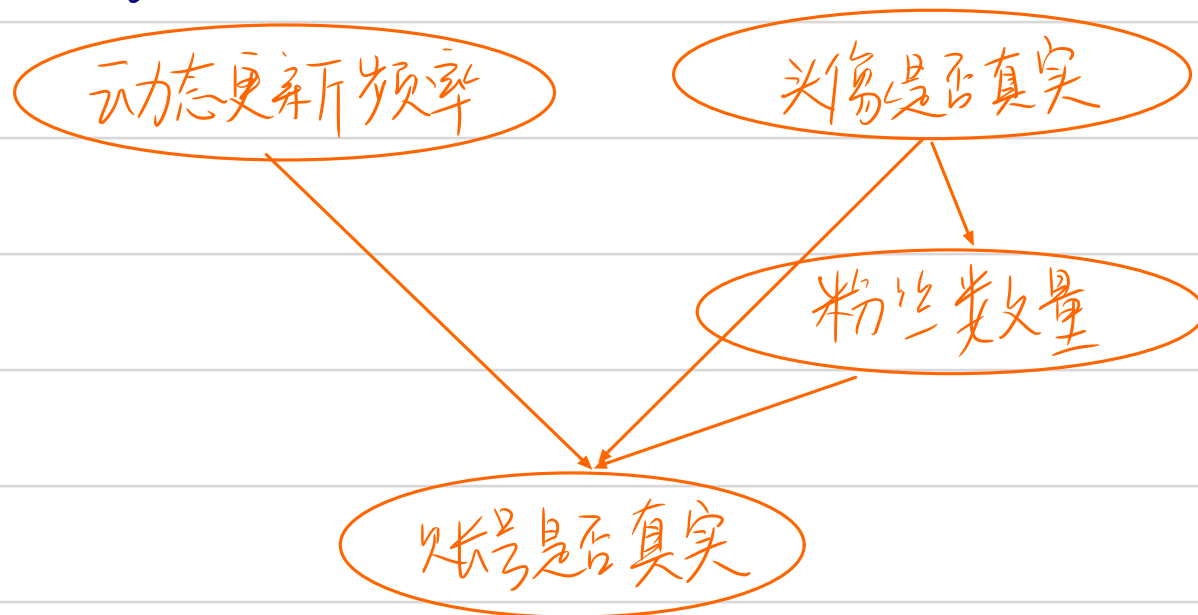
$$\hat{y} = \arg \max_{c_k} \prod_{j=1}^n P(X^{(j)}=x^{(j)} | Y=c_k) P(Y=c_k)$$

以上就是朴素贝叶斯的简单推导过程。

贝叶斯网络原理推导 Bayesian network

将朴素贝叶斯的条件独立性假设去掉，认为特征之间存在相关性的贝叶斯模型就是贝叶斯网络模型

一个例子作为引子，假设我们通过头像真实性、粉丝数量和动态更新频率来判断一个微博账号是否为真实账号。



上图是一个有向无环图 (directed graph, DAG)，贝叶斯网络中每个结点还有一个与之对应的概率表。假设下面是账号是否真实和头像是否真实之间的概率表。

A=0	A=1
0.13	0.87

账号是否真实

	H=0	H=1
A=0	0.88	0.12
A=1	0.25	0.75

账号真实性对于头像真实性的条件概率

已知某微博账号使用了虚假头像，那么其账号为虚假账号的概率可以推断为：

$$\begin{aligned} P(A=0|H=0) &= \frac{P(H=0|A=0)P(A=0)}{P(H=0)} \\ &= \frac{P(H=0|A=0)P(A=0)}{P(H=0|A=0)P(A=0) + P(H=0|A=1)P(A=1)} = \frac{0.88 \times 0.13}{0.88 \times 0.13 + 0.25 \times 0.87} \approx 0.35 \end{aligned}$$

上面例子展示了贝叶斯网络的用法，一个贝叶斯网络通常由 DAG 和结点对应的概率表组成。其中 DAG 由结点 (node) 和有向边 (edge) 组成，结点表示特征属性或随机变量，有向边表示各变量之间的依赖关系。贝叶斯网络的一个重要特征是：当一个结点的父结点概率分布确定之后，该结点条件独立于所有非直接父结点。该性质方便我们计算变量之间联合概率分布。

一般来说,多变量非独立随机变量的联合概率分布计算公式如下:

$$P(X_1, X_2, \dots, X_n) = P(X_1) P(X_2 | X_1) P(X_3 | X_1, X_2) \dots P(X_n | X_1, X_2, \dots, X_{n-1})$$

有了结点条件独立性质后,上式可化简为:

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Parents}(X_i))$$

当由DAG表示的结点关系和概率表确定后,相关的概率分布、条件概率分布就能确定,然后基于贝叶斯公式,我们就可以使用贝叶斯网络进行推断了。

使用pgmpy的贝叶斯网络实现

D0	D1
0.6	0.4

难度(D)

天赋(L)

L0	L1
0.7	0.3

成绩(G)

SAT(S)

推荐信(L)

	L0		L1	
	G0	G1	G0	G1
G0	0.3	0.65	0.9	0.5
G1	0.4	0.25	0.08	0.3
G2	0.3	0.7	0.02	0.2

	L0	L1
S0	0.95	0.2
S1	0.05	0.8

	G0	G1	G2
L0	0.1	0.4	0.99
L1	0.9	0.6	0.01