

朴素贝叶斯方法

1 模型及其说明

朴素贝叶斯方法

朴素贝叶斯方法就是频率学派在此公式上的一个巧妙的应用，通过最大化后验概率获取某一个确定的分类，在具体使用过程中效果良好

$$p(w|D) = \frac{p(D|w)p(w)}{p(D)}$$

贝叶斯公式

设输入空间  $X \in R$  为n维向量的集合，输出空间为类标记集合  $Y = \{c_1, c_2, \dots\}$ ，输入为特征向量  $x \in R$ ，输出为类标记(class label)  $y \in Y$ 。X是定义在输入空间x上的随机向量，Y是定义在输出空间D上的随机变量。P(X, Y)是X和Y的联合概率分布。训练数据集  $T = \{(X_1, Y_1), (X_2, Y_2), (X_3, Y_3), \dots, (X_n, Y_n)\}$ ，在给定的数据集下我们能够求出X和Y的联合概率分布 P(X, Y)

$$P(Y = c_k) (k = 1, 2, \dots, n)$$

为了使用贝叶斯公式，我们计算先验概率与似然函数

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

$$P(X, Y) = \sum_{i=0}^N P(X = x|Y = c_k) P(Y = c_k)$$

所求的联合密度函数就转换为了

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

在计算条件概率中P(X|Y)是一个非常大的指数级参数，在求取过程中非常复杂。于是，求解过程中给定一个条件独立的假设

1.1 模型推导

在已知  $X = (x^{(1)}, x^{(2)}, \dots, x^{(N)})$  的情况下，需要求得  $P(Y = C_k|X)$  的极大似然估计。根据朴素贝叶斯假设法， $P(X|Y) = \frac{P(X,Y)}{P(Y)}$ ，将(1)、(3)、(4)代入得到：

$$\begin{aligned} P(Y = C_k|X = x) &= \frac{P(X = x|Y = C_k)P(Y = C_k)}{P(X = x)} \\ &= \frac{P(X = x|Y = C_k)P(Y = C_k)}{\sum_{j=1}^N P(X = x|Y = C_j)P(Y = C_j)} \\ &= \frac{P(Y = C_k) \prod_{j=1}^N P(X^j = x^j|Y = C_k)}{\sum_{j=1}^N P(Y = C_j) \prod_{j=1}^N P(X^j = x^j|Y = C_j)} \end{aligned} \quad (5)$$

于是模型推导过程如下：

自此(5)所有的值都能通过已知数据获得，于是朴素贝叶斯分类器(6)可表示为：

$$y = f(x) = \arg \max_{\alpha} \frac{P(Y = \alpha_k) \prod_{j=1}^n P(X^j = x^j|Y = \alpha_k)}{\sum_{\alpha} P(Y = \alpha_k) \prod_{j=1}^n P(X^j = x^j|Y = \alpha_k)}$$

对于所有的分类而言，分类器的分母都是相同的，所以模型能够在次简化：

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

但是，在模型中还是过于复杂，对于所有的分类而言，分类器的分母都是相同的，所以模型能够在次简化：

选取0-1损失函数

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

则，经验风险函数可以表示为

$$R_{exp}(f) = E[L(Y, f(x))] = E_X \sum_{i=1}^K [L(c_k, f(x))] P(c_k|X)$$

在实际过程中我们需要最小化经验风险函数

$$\begin{aligned} \min R_{exp}(f) &= \arg \min_{\alpha_k} \sum_{i=1}^K L(c_k, y) P(c_k|X) \\ &= \arg \min_{\alpha_k} \sum_{i=1}^K P(y \neq c_k|X) \\ &= \arg \min_{\alpha_k} \sum_{i=1}^K 1 - P(y = c_k|X) \\ &= \arg \max_{\alpha_k} \sum_{i=1}^K P(y = c_k|X) \end{aligned}$$

从公式中能够看出最小化经验风险函数，与最大化决策分类函数是一样的，这也就是朴素贝叶斯方法的原理

最大化后验概率

部分验证中，准确率为100%，在全数据验证中，准确率为97.37%，且实验结果与sklearn的训练效果相同，这表明模型没有问题

Iris部分数验证  
Iris全数验证

1 GaussianNB 高斯朴素贝叶斯方法

2 Bernoulli Naive Bayes 伯努利朴素贝叶斯

3 Multinomial Naive Bayes 多项式朴素贝叶斯

3 实验验证

本实验在给定的数据集下实现了新闻分类，具体步骤如下：

- 读取文件，使用分词器'jieba'分词，得到全数据、训练数据、测试数据三部分  
\* 全数据用于特征向量的选取
- 生成断句文本，在创建特征向量时屏蔽掉一些不需要的文本  
\* 如果不讲这些词语屏蔽的话就会引入过多的噪声，这类词语在中文中常常含有'啊'、'的'、'嗯'等词语，这类词语出现频率高，而且对于分类的效果并不好，而且由于文本的处理，会引入空格、回车等字符，也会影响分类
- 将训练数据、测试数据向量化，使之能够使用'MultinomialNB'分类，分类器搭建完成

在测试中并不是特征值选取的越多越高，在最后选取了3500个特征值时能够达到100%的准确率

4 新闻分类器

2 具体应用

1 GaussianNB 高斯朴素贝叶斯方法

- $P(y)$  的估计：使用频率学派的极大似然估计——训练集中类别的相对频率： $P(Y = c_k) = \frac{\sum_{i=1}^N I(y=c_k)}{N}, k = 1, 2, \dots, K$
- $P(x_i|y)$  的估计：形成了不同的朴素贝叶斯模型：

GaussianNB 高斯朴素贝叶斯  
对于特征分量  $x_i$  是连续值情形，假设它服从高斯分布。  
概率密度函数：

$$P(x_i|y_k) = \frac{1}{\sqrt{2\pi\sigma_{y_k}^2}} \exp\left(-\frac{(x_i - \mu_{y_k})^2}{2\sigma_{y_k}^2}\right)$$

它的2个模型参数，可以使用频率学派的极大似然方法从训练集中进行估计：

- 均值： $\mu_{y_k} = \frac{\sum x_{y_k}}{N}$ ；
- 方差： $\sigma_{y_k}^2 = \frac{\sum (x_{y_k} - \mu_{y_k})^2}{N}$ ；

它们分别表示某个类别  $y_k$  中第  $i$  个特征  $x_i$  的均值与方差。因为朴素贝叶斯假设这些特征是条件独立的，所以可以把每一个特征分量的分布当作一维高斯分布。

通过已知数据训练得到均值与方法，在预测时只要使用均值与方差就能求出对应的条件概率然后就能通过模型求出其分类了

2 Bernoulli Naive Bayes 伯努利朴素贝叶斯

伯努利朴素贝叶斯：

- $P(x_i|y)$  的计算公式如下：

$$P(x_i|y) = P(i|y)x_i + (1 - P(i|y))(1 - x_i)$$

与高斯贝叶斯方法一样先计算条件概率，然后计算分类概率，求最大值就是其分类

该方法应用于文本分类(Text Classification)的2个经典朴素贝叶斯变体之一。它的特征数根据从多项式分布，使用单词计数向量(Word Count Vector)来表示输入数据；-分布参数由参数向量表示：  
 $\theta_y = (\theta_{y1}, \dots, \theta_{yn})$   
其中， $n$  表示特征个数，例如，在文本分类中， $n$  表示单词容量。 $\theta_{yi}$  表示特征  $i$  出现在属于类别  $y$  的样本  $x$  中的概率  $P(x_i|y)$ 。  
由于特征条件分布的独立性假设，参数  $\theta_y$  的每个分量都可以单独使用最大似然方法进行估计，常用情况下，使用平滑版本：

$$P(x_i|y) = \theta_{yi} = \frac{N_{yi} + \alpha}{N_y + \alpha n}$$

其中， $N_{yi} = \sum x_i$  是训练集  $T$  中，特征  $i$  出现在类别  $y$  中的个数， $N_y = \sum_{i=1}^n N_{yi}$  是类别  $y$  中所有特征的个数总和， $n$  表示特征  $x_i$  的取值个数。  
上式中， $\alpha \geq 0$ ，当  $\alpha = 0$  时，它等价于最大似然估计。 $\alpha > 0$  可以防止后估计的概率值出现0的情况，避免对后验计算产生影响。当  $\alpha = 1$  时，该平滑公式被称为拉普拉斯平滑(Laplace Smoothing)，当  $\alpha < 1$  时，被称为Lidstone平滑。

3 Multinomial Naive Bayes 多项式朴素贝叶斯

4.1 模型推导

假设有一个离散型变量，它可能取值1, 2, 3, ..., K，那么关于  $x$  的概率分布的集合可以使用一个参数为向量  $P$  的分布表示，其中  $P(X = k) = P_k$ ，则一种写法如下：  
$$p(x|P) = \prod_{i=1}^n p_i^{x_i} \quad (1)$$
  
这里  $p_i(x = k)$  是一个指示函数，当  $x = k$  时，它的取值为1，否则为0。那么  $N$  个独立同分布的样本  $X = \{x_1, x_2, \dots, x_N\}$  的联合概率分布函数：  
$$p(X|P) = \prod_{i=1}^N p_i^{x_i} \quad (2)$$
  
其实  $x_i$  的取值就是观测值  $x_i$  的变量的个数。构造分类器：  
$$f(x) = \arg \max_P p(X|P) = \prod_{i=1}^N p_i^{x_i}$$
  
且：
$$\sum_{i=1}^K p_i = 1$$
  
接下来，我们使用拉普拉斯平滑法，求得参数  $P$ 。  
构造似然函数：  
$$f(P, X) = \sum_{i=1}^N x_i \log p_i + \lambda \sum_{i=1}^K (p_i - 1)$$

在上述基础上我给出具体的推导过程：

对  $p_k$  与  $\lambda$  分别求偏导得：

$$\frac{\partial f(P, \lambda)}{\partial p_k} = \sum_{i=1}^N x_i - \sum_{i=1}^K \lambda = 0$$
$$\frac{\partial f(P, \lambda)}{\partial \lambda} = \sum_{i=1}^K p_i - 1 = 0$$

可以解得：

$$\lambda = \frac{\sum x_i}{p_k}$$
$$p_k = \frac{c_k}{\sum_{i=1}^K c_i}$$

用(3)代入(2)得到求最优参数。于是模型求取方法与高斯朴素贝叶斯方法为一样进行，但是在(2)中，我们能够得到，如果存在一个  $c_k = 0$ ，则求出的条件概率最终为0，显然这个概率是不符合实际情况的，于是我们加以变形：

$$p_k = \frac{c_k + \alpha}{\sum_{i=1}^K c_i + \alpha n}$$

上式中， $\alpha \geq 0$ ，当  $\alpha = 0$  时，它等价于最大似然估计。 $\alpha > 0$  可以防止后估计的概率值出现0的情况，避免对后验计算产生影响。当  $\alpha = 1$  时，该平滑公式被称为拉普拉斯平滑(Laplace Smoothing)，当  $\alpha < 1$  时，被称为Lidstone平滑。

4 对比与回顾

- 感知机
- KNN
- 朴素贝叶斯

	感知机	KNN	朴素贝叶斯
样本要求	线性可分	无	样本的各个属性间相互独立
模型类别	判别模型	判别模型	生成模型
有明显的训练模型	有	无	有
算法	梯度下降法	KD树搜索	概率计算，没有明确的计算过程
模型评价	效果一般般，不能通过训练次数确定模型的好坏，模型不唯一，模型有待改进。	训练效果较好，KD树的构造耗时较长，在数据变化较大的情况下不适用，不能及时的确定K值的大小，只能在已知数据集的情况下估算，但是这样的K值泛化能力不够强	训练效果好，在不同的情况下使用不同的概率计算方式，得到的效果往往较好。