输入空间：$X \in R^n$

输出空间 $Y = \{c_1, c_2 \ldots, c_K\}$

输入： $x \in X$    feature vector

输出： $y \in Y$    class label

先验概率   $P(Y = c_k) \quad k = 1, 2 \ldots, K$

条件概率   $P(X = x \mid Y = c_k) = P(X^{(1)} = x^{(1)}, X^{(2)} = x^{(2)} \ldots, X^{(n)} = x^{(n)} \mid Y = c_k)$

条件独立性假设：   $P(X = x \mid Y = c_k)$

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

后验概率： $P(Y = c_k \mid X = x) = \dfrac{P(X = x \mid Y = c_k) P(Y = c_k)}{\sum\limits_{k=1}^{K} P(X = x \mid Y = c_k)}$

(应用 Bayes 公式)

$$= \dfrac{\prod\limits_{i=1}^{n} P(X^{(i)} = x^{(i)} \mid Y = c_k) P(Y = c_k)}{\sum\limits_{k=1}^{K} \prod\limits_{i=1}^{n} P(X^{(i)} = x^{(i)} \mid Y = c_k)}$$

朴素贝叶斯分类： $y = f(x) = \underset{c_k}{\text{argmax}} \dfrac{\prod\limits_{i=1}^{n} P(X^{(i)} = x^{(i)} \mid Y = c_k) P(Y = c_k)}{\sum\limits_{k=1}^{K} \prod\limits_{i=1}^{n} P(X^{(i)} = x^{(i)} \mid Y = c_k)}$

$$= \underset{c_k}{\text{argmax}} \prod_{i=1}^{n} P(X^{(i)} = x^{(i)} \mid Y = c_k) P(Y = c_k)$$

(分母被所有 $c_k$ 共享).

后验概率最大化 ⟺ 期望风险最小化

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

0~1 损失函数

$$R_{exp}(f) = E[L(Y, f(x))]$$

期望风险函数
$$= E_x \sum_{k=1}^{K} [L(c_k, f(x))] P(c_k | x)$$

$$f(x) = \underset{y \in Y}{\arg\min} \sum_{k=1}^{K} L(c_k, y) P(c_k | X = x)$$

$$= \underset{y \in Y}{\arg\min} \sum_{k=1}^{K} P(y \neq c_k | X = x)$$

$$= \underset{y \in Y}{\arg\min} 1 - P(y = c_k | X = x)$$

$$= \underset{y \in Y}{\arg\max} P(y = c_k | X = x)$$

极大似然估计:

$$P(Y = c_k) = \frac{\sum_{i=1}^{N} I\{y_i = c_k\}}{N}$$

$$P(X^{(j)} = a_{jr} | Y = c_k) = \frac{\sum_{i=1}^{N} I\{x_i^{(j)} = a_{jr}, y_i = c_k\}}{\sum_{i=1}^{N} I\{y_i = c_k\}}$$

( $a_{jr}$ 表示 x 第 j 个特征的 j 旅取值)

朴素贝叶斯算法： 本质上 Bernoulli Event Model

输入： $T = \{ (x_1, y_1), (x_2, y_2) \cdots (x_N, y_N) \}$ N个样本

$x_i = (x_i^{(1)}, x_i^{(2)} \cdots, x_i^{(n)})$ n个特征

$x_i^{(j)} \in \{ a_{j1}, a_{j2} \cdots, a_{jp} \}$ 每个特征 P 种取值

$y_i \in \{ c_1, c_2 \cdots c_K \}$ K个分类

(1) 计算先验概率

$$P(Y = c_k) = \frac{\sum\limits_{i=1}^{N} \mathcal{L} \{ y_i = c_k \}}{N}$$

$$P(X^{(j)} = a_{jp} \mid Y = c_k) = \frac{\sum\limits_{i=1}^{N} \mathcal{L} \{ x^{(j)} = a_{jp}, y_i = c_k \}}{\sum\limits_{i=1}^{N} \mathcal{L} \{ y_i = c_k \}}$$

$(j = 1, 2 \cdots, n$
$\quad p = 1, 2 \cdots, P$
$\quad k = 1, 2 \cdots, N )$

(2) 对于给定 $x$, 计算

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

(3) 确定 $x$ 的类

$$y = \arg\max_{a_k} P(Y = c_k) \prod_{i=1}^{n} P(X^{(i)} = x^{(i)} \mid Y = c_k)$$

贝叶斯估计

$$P(X^{(j)} = a_{jp} | Y = c_k) = \frac{\sum\limits_{i=1}^{N} \mathbb{1}\{x_i^{(j)} = a_{jp}, y_i = c_k\} + \lambda}{\sum\limits_{i=1}^{N} \mathbb{1}\{y_i = c_k\} + P \cdot \lambda}$$

$P$: 特征可能取值数

$\lambda = 1$ 时: Laplace Smoothing

① $P(X^{(j)} = a_{jp} | Y = c_k) > 0$

② $\sum\limits_{p=1}^{P} P(X^{(j)} = a_{jp} | Y = c_k) = 1$

同理, 先验概率变为 $P(Y = c_k) = \dfrac{\sum\limits_{i=1}^{N} \mathbb{1}\{y_i = c_k\} + \lambda}{N + K\lambda}$

实际中通常不用. 因为 数据集中 $P(Y = c_k) \neq 0$
$k = 1, 2 \cdots, K$

补充: 将连续值转化为离散值(如设立 12 个左右区间)即

可以应用 NB 算法

补充: 多项式概率分布  Multinomial Event Model

$$P(X^{(j)} = a_{jp} | y = c_k) = \frac{\sum\limits_{i=1}^{N} \sum\limits_{m=1}^{n} \mathbb{1}\{X_m^{(j)} = a_{jp}, y = c_k\}}{\sum\limits_{i=1}^{N} \mathbb{1}\{y = c_k\} \cdot n}$$

$n$: $x$ 的特征维度

(以文本为例. 假设每个单词与其出现的位置无关)

应用 Laplace 变换: 同上