

Machine Learning

Naive Bayes

1. Basics

Bayes' theorem:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

independence assumptions:

$$\begin{aligned} P(X|Y) &= P(X_1, X_2|Y) \\ &= P(X_1|X_2, Y)P(X_2|Y) \\ &= P(X_1|Y)P(X_2|Y) \end{aligned}$$

where $X = \langle X_1, X_2 \rangle$.

Normal distribution:

$$p(x, \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

2. Model

input:

$(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ where $x^{(i)} \in \mathbb{R}^n, y^{(i)} \in \{1, \dots, k\}$

hypothesis:

$$p(x|y) = \prod_{i=1}^n p(x_i|y)$$

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_y p(y|\hat{x}) = \operatorname{argmax}_y \frac{p(\hat{x}|y)p(y)}{p(\hat{x})} \\ &= \operatorname{argmax}_y p(\hat{x}|y)p(y) = \operatorname{argmax}_y p(y) \prod_{i=1}^n p(\hat{x}_i|y)\end{aligned}$$

parameter estimation(discrete):

$$\begin{aligned}p(\hat{x}_j|y) &= \frac{\sum_{i=1}^m \mathbf{1}\{x_j^{(i)} = \hat{x}_j \wedge y^{(i)} = y\}}{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = y\}} \\ p(y) &= \frac{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = y\}}{m}\end{aligned}$$

smooth(discrete):

$$\begin{aligned}p(\hat{x}_j|y) &= \frac{\sum_{i=1}^m \mathbf{1}\{x_j^{(i)} = \hat{x}_j \wedge y^{(i)} = y\} + l}{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = y\} + l \cdot \operatorname{distinct}\{x_j^{(1)}, \dots, x_j^{(m)}\}} \\ p(y) &= \frac{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = y\} + l}{m + lk}\end{aligned}$$

parameter estimation(continuous):

$$\begin{aligned}p(\hat{x}_j|y) &= \frac{1}{\sigma_{jy}\sqrt{2\pi}} e^{-\frac{(\hat{x}_j - \mu_{jy})^2}{2\sigma_{jy}^2}} \\ \mu_{jy} &= \frac{\sum_{i=1}^m x_j^{(i)} \cdot \mathbf{1}\{y^{(i)} = y\}}{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = y\}} \\ \sigma_{jy} &= \frac{\sum_{i=1}^m (x_j^{(i)} - \mu_{jy})^2 \cdot \mathbf{1}\{y^{(i)} = y\}}{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = y\}} \\ p(y) &= \frac{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = y\}}{m}\end{aligned}$$

unbiased estimation:

$$\sigma_{jy} = \frac{\sum_{i=1}^m (x_j^{(i)} - \mu_{jy})^2 \cdot \mathbf{1}\{y^{(i)} = y\}}{(\sum_{i=1}^m \mathbf{1}\{y^{(i)} = y\}) - 1}$$

output:

$$\hat{y} = \operatorname{argmax}_y p(y|\hat{x}) = \operatorname{argmax}_y p(y) \prod_{i=1}^n p(\hat{x}_i|y)$$