

Naive Bayes Classifier

Andersen Ang

First created: 2014. Last update : 2017

- Consider a supervised learning problem : given data, the aim is try to construct a parametric probabilistic models, assume the model is static.
- When all th attributes of the data (feature vector for example) are independent of each other.
- For a data vector $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T$ with class label C .
- Using Bayes rule, the probability of class label C given the features \mathbf{x} is :

$$\mathbb{P}(C | x_1, x_2, \dots, x_n) = \frac{\mathbb{P}(x_1, \dots, x_n | C) \mathbb{P}(C)}{\mathbb{P}(x_1, \dots, x_n)}$$

- Because now all the parameters are assumed to be independent of each other

$$\mathbb{P}(x_1, \dots, x_n | C) = \mathbb{P}(x_1 | C) \mathbb{P}(x_2 | C) \dots \mathbb{P}(x_n | C)$$

- Using short hand notation

$$\mathbb{P}(C | x_1, x_2, \dots, x_n) = \frac{\prod_i \mathbb{P}(x_i | C) \mathbb{P}(C)}{\mathbb{P}(x_1, \dots, x_n)}$$

- $\mathbb{P}(C)$: Prior probability distribution, which can be estimated from data.
- $\mathbb{P}(x_i | C)$: likelihood function, class conditional probability distribution, which can be estimated from data.
- $\mathbb{P}(x_1, \dots, x_n)$: __ distribution, which is a constant

In this case, the class label C is choosen such that the $\mathbb{P}(C | x_1, x_2, \dots, x_n)$ is maximized.

$$\text{maximize } \mathbb{P}(C | x_1, x_2, \dots, x_n) \iff \text{maximize } \frac{\prod_i \mathbb{P}(x_i | C) \mathbb{P}(C)}{\mathbb{P}(x_1, \dots, x_n)}$$

Since $\mathbb{P}(x_1, \dots, x_n)$ and $\mathbb{P}(C)$ is a constant.

$$\text{maximize } \mathbb{P}(C | x_1, x_2, \dots, x_n) \quad \iff \quad \text{maximize } \mathbb{P}(C) \prod_i \mathbb{P}(x_i | C)$$

Therefore, in Naive Bayes Classifier, we need to estimate the Prior probability and conditional probability.

Example. Normal distribution

$$\mathcal{N}(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \left(\frac{x - \mu}{\sigma}\right)^2\right)$$

Based on Bayes assumption, all the parameters are independent, so the pdf can be decomposed as

$$\mathbb{P}(x_1, x_2, \dots, x_n | \mu, \sigma) = \prod_{i=1}^N \mathbb{P}(x_i | \mu, \sigma)$$

$$\mathbb{P}(x_1, x_2, \dots, x_n | \mu, \sigma) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma}\right)^2\right)$$

Log likelihood

$$\log \mathbb{P}(x_1, x_2, \dots, x_n | \mu, \sigma) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma}\right)^2\right)$$

$$L = \log \mathbb{P}(x_1, x_2, \dots, x_n | \mu, \sigma) = N(-\sqrt{2\pi} - \log \sigma) + \sum_{i=1}^N \frac{1}{2} \left(\frac{x_i - \mu}{\sigma}\right)^2$$

Take derivative

$$\frac{\partial L}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu) = 0 \quad \Rightarrow \quad \mu = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\frac{\partial L}{\partial \sigma} = -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^N (x_i - \mu)^2 = 0 \quad \Rightarrow \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$