

Bayesian ML.

1. Naive Bayes.

$$p(c_k | \vec{x}) = \frac{p(c_k) \cdot p(\vec{x} | c_k)}{p(\vec{x})} \propto p(c_k, \vec{x}) =$$

$$= p(x_1 | x_2, \dots, x_n, c_k) \cdot p(x_2 | x_3, \dots, x_n, c_k) \cdot \dots \cdot p(x_n | c_k) p(c_k) =$$

$$= \left(\prod_{i=1}^n p(x_i | c_k) \right) p(c_k)$$

Naiveness

\Downarrow

$$\hat{y}(\vec{x}) = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} \left(p(c_k) \prod_{i=1}^n p(x_i | c_k) \right)$$

$$1) p(x_i | c_k) = \frac{1}{\sqrt{2\pi} \sigma_k} \exp \left\{ -\frac{(x_i - \mu_k)^2}{2 \sigma_k^2} \right\} \quad \text{— Gaussian NB (numerical features)}$$

$$2) p(\vec{x} | c_k) = \frac{\left(\sum_{i=1}^n x_i \right)!}{\prod_{i=1}^n (x_{i0}!)} \prod_{i=1}^n p_{ki}^{x_i}, \quad p_{ki} - \text{probab. that event } i \text{ occurs in class } k.$$

x_i - number of occur. for event i

Multinomial NB (categorical features)

$$3) p(\vec{x} | c_k) = \prod_{i=1}^n p_{ki}^{x_i} \cdot (1 - p_{ki})^{1-x_i}, \quad p_{ki} - \text{—//—},$$

x_i - binary label if event i occurred.

Example ~~NB + Logistic~~ TF, IDF + NB + Logistic Regression (Wang, 2022) for sentiment analysis.

1) X - texts, y - ~~text~~ binary labels
 $\vec{x} = \text{tf-idf}(X, \text{gram_range}=(1,2)).$

$$2) p(c_0) =$$

$$\frac{p(c_1 | \vec{x})}{p(c_0 | \vec{x})} = \frac{p(c_1) \cdot \prod_{i=1}^n p_{1i}^{x_i}}{p(c_0) \cdot \prod_{i=1}^n p_{0i}^{x_i}}$$

$$2) \quad p(r_2 | \vec{x}) > p(r_0 | \vec{x})$$

$$\Downarrow$$

$$\log \left\{ \frac{p(r_2 | \vec{x})}{p(r_0 | \vec{x})} \right\} > 0$$

$$\Downarrow$$

$$\vec{w}^T \vec{x} + b > 0 \text{ where } w_i = \log \left(\frac{p_{2i}}{p_{0i}} \right), \quad b = \log \left(\frac{p(r_2)}{p(r_0)} \right).$$

Linear

$$\log \left(\frac{p(r_2 | \vec{x})}{1 - p(r_2 | \vec{x})} \right) > 0$$

$$\Downarrow$$

$$p(r_2 | \vec{x}) > 0.5$$

with MNB: $p(r_2 | \vec{x}) = p(r_2) \cdot \bar{\pi}$, $\bar{\pi} = \frac{(\sum_{i=2}^n x_i)!}{(\prod_{i=2}^n x_i!) \cdot p(\vec{x})}$

$$p(r_2 | \vec{x}) = p(r_2) \cdot \bar{\pi} \cdot \prod_{i=2}^n p_{2i}^{x_i}, \quad \bar{\pi} = \frac{(\sum_{i=2}^n x_i)!}{(\prod_{i=2}^n x_i!) \cdot p(\vec{x})} \quad \text{- hard to compute.}$$

idea: reexpress $p(r_2 | \vec{x})$ as the output of logistic regression:

$$p(r_2 | \vec{x}) = \frac{1}{1 + e^{-z}} \Leftrightarrow z = \ln \left(\frac{p(r_2 | \vec{x})}{1 - p(r_2 | \vec{x})} \right) = \ln \left(\frac{p(r_2 | \vec{x})}{p(r_0 | \vec{x})} \right) =$$

$$= \ln \left(\frac{p(r_2)}{p(r_0)} \right) + \sum_{i=2}^n \ln \left(\frac{p_{2i}}{p_{0i}} \right) \cdot x_i$$

$$p_{2i} \text{ is estimated as: } p_{2i} = \frac{\sum_{j: y^{(j)}=1} x_i^{(j)} + 1}{\sum_{j: y^{(j)}=1} 1 + 1} \quad \text{same with } p_{0i},$$

afterwards we apply logistic regression to $\tilde{x} = \left\{ \ln \left(\frac{p_{2i}}{p_{0i}} \right) \cdot x_i \right\}_{i=0}^n$.