

## NaiveBayes

19 July 2018 15:47

Given a training dataset which are sampled from a hidden  $P(x, y)$  distribution

$$D = \{ (x^1, y_1), \dots, (x^N, y_N) \}$$

Goal: (during the training process) to estimate  $\hat{P}(x|y)$ ,  $\hat{P}(y)$

(1) How to estimate  $\hat{P}(y)$   $P(x, y) = P(y)P(x|y)$

$D \Rightarrow \{N_1, N_2, \dots, N_k\}$   $\sum_j N_j = N$   
 $\nearrow$  # of examples in each class

$$\hat{P}(y) = \left[ \frac{N_1}{N}, \frac{N_2}{N}, \dots, \frac{N_k}{N} \right]$$

$$\hat{P}(y=j) = \frac{N_j}{N} \quad \text{eg: students} \begin{array}{c} 2 \quad 3 \quad 4 \\ \hline 0.2 \quad 0.7 \quad 0.1 \end{array}$$

(2) How to estimate  $\hat{P}(x|y) = \hat{P}(x_1, \dots, x_d|y)$ ?

eg: student:  $x_1 = \text{age}$

$$\hat{p}(x_1 | y=2) \quad \hat{p}(x_1 | y=3) \quad \hat{p}(x_1 | y=4)$$

$$p(x_1 | y) \sim N(\mu_y, \sigma_y)$$

$$\hat{p}(x_1 | y=2) \sim N(\mu_2 = 18.7, \sigma_2 = 0.5)$$

$$\hat{p}(x_1 | y=3) \sim N(\mu_3 = 19.6, \sigma_3 = 0.4)$$

$$\hat{p}(x_1 | y=4) \sim N(\mu_4 = 20.5, \sigma_4 = 0.45)$$

Testing!  $x: x_1 = 19$

$$p(y | x) \propto \frac{\hat{p}(x | y) \hat{p}(y)}{\hat{p}(x)}$$

$$= \frac{1}{\sqrt{2\pi}\sigma_y} e^{-\frac{(x_1 - \mu_y)^2}{2\sigma_y^2}} \hat{p}(y)$$

$$p(y=2 | x: x_1=19) \propto \frac{1}{0.5} e^{-\frac{(19-18.7)^2}{2 \times 0.5^2}} \times 0.2$$

$$p(y=3 | x_1=19) \propto \frac{1}{0.4} e^{-\frac{(19-19.6)^2}{2 \times 0.4^2}} \times 0.7$$

$$p(y=4 | x_1=19) \propto \frac{1}{0.45} e^{-\frac{(19-20.5)^2}{2 \times 0.45^2}} \times 0.1$$

$$0.45 \quad \hookrightarrow \quad 2 \neq 0.45$$

$P(y|x)$  ← posterior probability of

$$d > 1 \quad y$$

$$P(x_1, x_2, \dots, x_d | y)$$

↑ age      ↑ late or not      ↑ # of questions answered      ↑ row in class

$\mathcal{E}_v$

Conditional independence assumption  
Example from slides