

Категор. текстоб

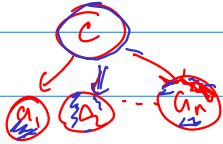
Болл [0...0 10...0]

Naive Bayes

$$p(C|\bar{x}) = \frac{p(C)p(\bar{x}|C)}{p(\bar{x})} \propto p(C) \cdot \prod_{i=1}^n p(a_i|C)$$

$$p(C) = \frac{|C|}{|D|}$$

$$p(w|C) = \frac{\sum_{w \in C} \#w}{|C|}$$



$$\bar{x} = (a_1, a_2, \dots, a_n)$$

Кр. Пот. нап. б. кр. w_i w_j

$$p(C_{noisy}|\bar{s}) \propto p(C_{noisy}) \cdot p(K_{noisy}|c_n) \cdot p(a_n|c_n) \cdot p(w|c_n)$$

$$p(Noisy|\bar{s}) \propto p(Noisy) \cdot p(\dots) \cdot p(\dots)$$

Multinomial NB



Multivariate NB

$$p(s|C) = \prod_{w \in V} p(w|C)$$

$$p(C_{noisy}|\bar{s}) \propto p(C_{noisy}) \cdot p(K_{noisy}|c_n) \cdot \dots \cdot [1 - p(K_{noisy}|c_n)] \cdot \dots$$

Мам. оценок

плоскостная $\Rightarrow p(a_{ij})$

Кр. Пот. нап. с. кр. \rightarrow катег. на 30 мин. едво
 $p(K|d) \cdot p(w|K)$
 $0.001 \times 0.001 = 0.001$

$\bar{s} \rightarrow p(Noisy|\bar{s}) = 75\%$
 $p(Noisy|\bar{s}) = 25\%$
 [0.999] 0.99 0.95 0.93
 [0.001] 0.01 0.05 0.07

- Болл ✓
- Naive assumption ✓
- это independent data set
- $\{ \pm \}$ \rightarrow $\{ \text{yes} \}$ - $\{ \text{no} \}$

$$D = \begin{bmatrix} 0.1 \\ 0.01 \\ 0.8 \\ \dots \\ \vdots \\ \vdots \end{bmatrix} \quad \begin{bmatrix} w \\ z \\ w \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

1) Мам. оценок \Rightarrow кросс-валидация

$$p(C|s) \propto p(C) \cdot \prod_{w \in s} p(w|C)$$

$$\pi_c = p(C), \varphi_{wt} = p(w|C_t)$$

- M-max: $\pi_t := \frac{\mathbb{E} \# \{c \in C_t\}}{|D|}$

$$\varphi_{wt} := \frac{\mathbb{E} \# \{w \in d \in C_t\}}{\mathbb{E} \sum_{d \in C_t} |d|}$$

- E-max: $p(s_n \in C_t) \propto p(C_t) \prod_{w \in s} \varphi_{wt} = \pi_t \cdot \prod_{w \in s} \varphi_{wt}$

$$\forall d \in D \quad z_{dt} = \frac{p(C_t|d)}{\sum_{t'} (\pi_{t'} \prod_{w \in d} \varphi_{wt'})}$$

$$\mathbb{E} \# \{d \in C_t\} = \sum_{d \in D} z_{dt}$$

$$\mathbb{E} \# \{w \in d \in C_t\} = \sum_{d \in D} z_{dt} \cdot n_w d$$

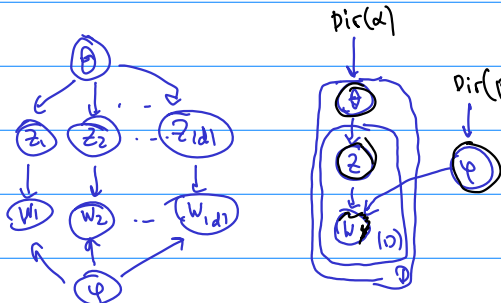
2) Topic Modeling

d - проект. бер. на макс $\theta_{dt} = p(t|d)$

$$\varphi_{tw} = p(w|t)$$

PLSA

$$p(d|\Theta, \Phi) = \prod_{w \in d} p(w|\Theta, \Phi) = \prod_{w \in d} \sum_{t \in T} \theta_{dt} \varphi_{tw}$$



$$d \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \approx d \begin{pmatrix} \Theta \\ \vdots \\ \Phi \end{pmatrix} \begin{pmatrix} w \\ \vdots \\ w \end{pmatrix}$$

$d \times w$ $d \times t$ $t \times w$

LDA - Latent Dirichlet allocation