

$\text{text1} \begin{pmatrix} 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$
нле не оч. попр.

$\begin{pmatrix} 1 \\ 2 \end{pmatrix}$
 ← класс 1, positive
 ← класс 2, negative

строки - объекты (тексты)
 столбцы - признаки

если есть информация для обучения, правильные ответы, то находите еще признаки

после обучения мы можем получить text3 , в виде набора признаков и получить предсказание класса (positive / negative)

Метод: будем определять вероятность класса по набору признаков

$$P\{C = \text{positive} \mid W_1 = w_1, W_2 = w_2, \dots, W_N = w_N\}$$

↑ класс ↑ feature, признак, слово

формула Байеса

Соправлено:

$$P\{\text{pos} \mid w_1, w_2, \dots, w_n\} = P\{w_1, w_2, \dots, w_n \mid \text{pos}\} \cdot \frac{P(\text{pos})}{P(w_1, \dots, w_n)}$$

$$P(A|B) \stackrel{\text{оп}}{=} \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) \stackrel{\text{оп}}{=} \frac{P(B \cap A)}{P(A)}$$



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

$$P\{\text{neg} \mid w_1, \dots, w_n\} = P\{w_1, \dots, w_n \mid \text{neg}\} \cdot \frac{P(\text{neg})}{P(w_1, \dots, w_n)}$$

$$P\{\text{pos}\} \cdot P\{\text{neg}\} = P(w_1, \dots, w_n)$$

- сопоставление

Уточню, что выдраны классы, гарантированно сравнимы

$$P\{w_1, \dots, w_n | \text{pos}\} = P\{\text{pos}\} \cdot \dots$$

$$P\{w_1, \dots, w_n | \text{neg}\} = P\{\text{neg}\} \cdot \dots$$

$$P(A|B) = P(A)P(B)$$

если A, B
независимы

- Переформулировка о независимости:

$$P\{w_1, \dots, w_n | C\} = P\{w_1 | C\} \cdot \dots \cdot P\{w_n | C\}$$

уточню, что определены классы, также сравнимы

$$P\{w_1 | \text{pos}\} \cdot \dots \cdot P\{w_n | \text{pos}\} \cdot P\{\text{pos}\} \quad \text{и}$$

$$P\{w_2 | \text{neg}\} \cdot \dots \cdot P\{w_n | \text{neg}\} \cdot P\{\text{neg}\}$$

каждый множитель можно оценить по требованию
каждому

$$P\{W_i = w_i | C = \begin{matrix} \text{pos, neg} \\ \downarrow \\ C \end{matrix} \} = \frac{C(\text{гоусметит класс } C, \text{ где } W_i \text{ есть или нет})}{1 + 0}$$

Пример

$$P\{C = C\} = \frac{\text{гок из } C}{\text{всего гок}} = \frac{C(\text{гоусметит класс } C)}{C(\text{класс } C)}$$

	мке	не	ор	нопр	
text1	1	0	0	1	pos
text2	1	1	1	0	neg
texts	0	0	1	1	pos

$$P\{\text{мке} = 1 | C = \text{pos}\} = \frac{1 - 1 \text{ pos гок со словом мке}}{2 - \text{pos гок}}$$

$$P\{\text{не} = 0 | C = \text{pos}\} = \frac{1}{2}$$

$$P\{\text{ноправнас} = 1 | C = \text{pos}\} = \frac{2}{2} \quad \begin{matrix} 2 \text{ гоусметит } \text{из} \text{ pos} \\ \text{со словом "ноправнас"} \end{matrix}$$

$$P\{x=1 \mid c = \text{neg } y\} = \frac{1}{1} - \text{гол} \cdot \text{со сработке не}$$

$$\frac{1}{1} - \text{всего neg гол}$$

Вероятности классов:

$$P\{C = \text{neg } y\} = \frac{1}{3}$$

$$P\{C = \text{pos } y\} = \frac{2}{3}$$

Как и в n-грамм вероятностях $P\{y \dots y\} = 0$

Эвент = 0. Также нужно спрашивать - что

все вероятности $\neq 0$.

На практике, вычисляя

$$P\{w_1 | c\} \dots P\{w_n | c\} P\{c\} \quad \nearrow \text{max}$$

Без $-\log$:

$$= \log P\{w_1 | c\} - \dots - \log P\{w_n | c\} - \log P\{c\}$$

это можно сделать меньше.

$\underbrace{\hspace{10em}}$
минусовать
там

$\underbrace{\hspace{10em}}$
минусовать
там

Что еще можно сделать после обучения:

можно найти самые большие значения

$P\{w_i | c\}$ - w_i сильно вероятны в классе c .