

# Домашнее задание №8

Аспанов Андрей, ИУ9-11М

Даны документы и их классы

$$D_1 = (x_1, x_2, x_3) \quad C_1$$

$$D_2 = (x_1, x_2, x_4) \quad C_1$$

$$D_3 = (x_4, x_5, x_6) \quad C_2$$

Определить класс документа на основе Байесовского классификатора

$$D_4 = (x_1, x_4, x_5)$$

Использовать 2 способа: Multinomial и Bernoulli

## Multinomial

В многомерной модели наивного Байесовского классификатора документ  $D_i$  - это вектор признаков  $x_i$ , показывающих, встречалось ли в том или ином документе слово.

$$C_{NB} = \underset{C_j \in C}{\operatorname{argmax}} \left\{ P(C_j) \prod_{i=1}^n P(x_i, C_j) \right\} \leftarrow \text{общая формула для Байесовского классификатора}$$

$$C_{NB} = \underset{C_j \in C}{\operatorname{argmax}} \left\{ P(D_4|C_1) P(C_1), P(D_4|C_2) P(C_2) \right\}$$

$$C_{NB} = \underset{C_j \in C}{\operatorname{argmax}} \left\{ \underbrace{P(x_1|C_1) P(x_4|C_1) P(x_5|C_1)}_{P(D_4|C_1)} P(C_1), \underbrace{P(x_1|C_2) P(x_4|C_2) P(x_5|C_2)}_{P(D_4|C_2)} P(C_2) \right\}$$

$$P(x_i, C_j) = \frac{N(x_i, C_j) + 1}{N(C_i) + k}$$

Общее число уникальных признаков, принадлежащих классу  $C_i$

Сколько раз встречается признак  $x_i$  при условии класса  $C_j$

суммарное число признаков обучающей выборки

← с учётом сглаживания по Лапласу

# Multinomial

Признак / (слово) документа	$P(x_i   C_1)$	$P(x_i   C_2)$
$x_1$	$\frac{2+1}{6+6} = \frac{1}{4}$	$\frac{0+1}{3+6} = \frac{1}{9}$
$x_2$	$\frac{2+1}{6+6} = \frac{1}{4}$	$\frac{0+1}{3+6} = \frac{1}{9}$
$x_3$	$\frac{1+1}{6+6} = \frac{1}{6}$	$\frac{0+1}{3+6} = \frac{1}{9}$
$x_4$	$\frac{1+1}{6+6} = \frac{1}{6}$	$\frac{1+1}{3+6} = \frac{2}{9}$
$x_5$	$\frac{0+1}{6+6} = \frac{1}{12}$	$\frac{1+1}{3+6} = \frac{2}{9}$
$x_6$	$\frac{0+1}{6+6} = \frac{1}{12}$	$\frac{1+1}{3+6} = \frac{2}{9}$

$$C_{NB} = \operatorname{argmax} \left\{ \frac{1}{4} \cdot \frac{1}{6} \cdot \frac{1}{12} \cdot \frac{2}{3}, \frac{1}{9} \cdot \frac{2}{9} \cdot \frac{2}{9} \cdot \frac{1}{3} \right\}$$

$$C_{NB} = \operatorname{argmax} \left\{ \underset{\substack{\text{"} \\ P(\text{дч-это} \\ \text{класс } C_1)}}{0,00231}, \underset{\substack{\text{"} \\ P(\text{дч-это} \\ \text{класс } C_2)}}{0,001828} \right\} \Rightarrow C_1.$$

Bernulli:

$N(X_i, C_j) + 1$   
 $N(C_i) + 2$   
 Потому что по Th. Байеса предполагаем, что мы как бы сделали +2 лишних испытания, одно из которых удачно

Признак (слово) документа  
 общее число признаков, прин-к классу  $C_i$   
 $P(X_i|C_1)$

сколько раз встречается признак  $X_i$  при условии класса  $C_j$   
 $P(X_i|C_2)$

Признак (слово) документа	$P(X_i C_1)$	$P(X_i C_2)$
X1	$\frac{2+1}{6+2} = \frac{3}{8}$	$\frac{0+1}{3+2} = \frac{1}{5}$
X2	$\frac{2+1}{6+2} = \frac{3}{8}$	$\frac{0+1}{3+2} = \frac{1}{5}$
X3	$\frac{1+1}{6+2} = \frac{1}{4}$	$\frac{0+1}{3+2} = \frac{1}{5}$
X4	$\frac{1+1}{6+2} = \frac{1}{4}$	$\frac{1+1}{3+2} = \frac{2}{5}$
X5	$\frac{0+1}{6+2} = \frac{1}{8}$	$\frac{1+1}{3+2} = \frac{2}{5}$
X6	$\frac{0+1}{6+2} = \frac{1}{8}$	$\frac{1+1}{3+2} = \frac{2}{5}$

$$C_{NB} = \operatorname{argmax} \{ P(X_1|C_1) P(X_4|C_1) P(X_5|C_1) \cdot P(C_1), P(X_1|C_2) P(X_4|C_2) P(X_5|C_2) P(C_2) \}$$

$$P_{NB} = \operatorname{argmax} \left\{ \frac{2}{3} \cdot \left( \frac{3}{8} \cdot \left(1 - \frac{3}{8}\right) \left(1 - \frac{1}{4}\right) \cdot \frac{1}{4} \cdot \frac{1}{8} \cdot \left(1 - \frac{1}{8}\right) \right), \right.$$

$$\left. \frac{1}{3} \cdot \left( \frac{1}{5} \cdot \left(1 - \frac{1}{5}\right) \left(1 - \frac{1}{5}\right) \cdot \frac{2}{5} \cdot \frac{2}{5} \cdot \left(1 - \frac{2}{5}\right) \right) \right\}$$

$$C_{NB} = \operatorname{argmax} \left\{ \frac{2}{3} \cdot \frac{3}{8} \cdot \frac{5}{8} \cdot \frac{3}{4} \cdot \frac{1}{4} \cdot \frac{1}{8} \cdot \frac{7}{8}, \frac{1}{3} \cdot \frac{1}{5} \cdot \frac{4}{5} \cdot \frac{4}{5} \cdot \frac{2}{5} \cdot \frac{2}{5} \cdot \frac{3}{5} \right\}$$

$$C_{NB} = \operatorname{argmax} \left\{ \frac{630}{196608}, \frac{192}{46875} \right\}$$

$$C_{NB} = \operatorname{argmax} \{ 0,003204; 0,004096 \} \Rightarrow \text{более вероятен класс } C_2$$

Вывод: по многокомпонентному Байесовскому классификатору более вероятен класс  $C_2$  по методу Бернулли - класс  $C_2$ .