

# Задание 9 от 19.10.20 (Вероятностные модели)

Валентин Александров, 620 группа

## Постановка задачи

- Запрос к поисковой системе  $q$ : "a b"
- В коллекции имеются следующие документы ( $d_1, d_2, d_3, d_4$ ):
  - a b c d
  - a a a
  - b b c
  - a b b c
- Применить языковую модель к этой коллекции
- Сравнить  $\lambda_1 = 0.5$  и  $\lambda_2 = 0.9$
- Как упорядочатся документы при этих значениях лямбда? Какая выдача кажется более правильной?

## Решение

$$P(w|d) = \lambda P_{mle}(w|M_d) + (1 - \lambda) P_{mle}(w|M_c)$$

$$\lambda_1 = 0.5$$

$$P(q|d_1) = (0.5 \frac{1}{4} + 0.5 \frac{7}{14}) \cdot (0.5 \frac{1}{4} + 0.5 \frac{5}{14}) = \frac{3}{8} \cdot 0.303571 = 0.113839$$

$$P(q|d_2) = (0.5 + 0.5 \frac{7}{14}) \cdot 0.5 \frac{5}{14} = 0.75 \cdot 0.178571 = 0.133928$$

$$P(q|d_3) = 0.5 \frac{7}{14} \cdot (0.5 \frac{2}{3} + 0.5 \frac{5}{14}) = 0.25 \cdot 0.178571 = 0.044642$$

$$P(q|d_4) = (0.5 \frac{1}{4} + 0.5 \frac{7}{14})(0.5 \frac{2}{4} + 0.5 \frac{5}{14}) = 0.375 \cdot 0.428571 = 0.160714$$

$$\lambda_2 = 0.9$$

$$P(q|d_1) = 0.071696$$

$$P(q|d_2) = 0.033928$$

$$P(q|d_3) = 0.031785$$

$$P(q|d_4) = 0.133571$$

	$\lambda_1 = 0.5$	$\lambda_2 = 0.9$
$P(q \mid d_1)$	0.113839	0.071696
$P(q \mid d_2)$	0.133928	0.033928
$P(q \mid d_3)$	0.044642	0.031785
$P(q \mid d_4)$	0.160714	0.133571

$\lambda_1 : d_4 > d_2 > d_1 > d_3$

$\lambda_2 : d_4 > d_1 > d_2 > d_3$

Ранжирование второй лямбды выглядит разумней (" $a b c d$ " > " $a a a$ "). Более того,  $d_2$  и  $d_3$  имеют низкие вероятности, что тоже выглядит разумно, так как эти документы не имеют оба слова  $a$  и  $b$  одновременно.