

Классификация текстов.

Задача. есть конечное мн-во классиф.

- 1) {спам/не спам}
- 2) {полож/отриц} - анализ тональности
- 3) {Толстой, Достоевский, Пушкин, ...} - опр. автора
- 4) {спорт, политика, общество, ...} - для новостей

не тексты надо определить класс

Для обучающей Корпус - набор текстов, для которых классы уже известны. (это обучение с учителем).

На основе каких данных определять класс?

Features - признаки. - информация, на основе которой происходит классификация.

1 Тип признаков: мешок слов (bag of words)

f_i - i -ый признак

$f_i(\text{text})$ - число, вектор, эл-т конечного мн-ва.

Если словарь всех возможных слов $V = \{w_1, w_2, \dots, w_n\}$

1) $f_i(\text{text}) = \begin{cases} 1, & \text{если слово } w_i \text{ встречается} \\ 0, & \text{если нет} \end{cases}$

2) $f_i(\text{text}) =$ сколько раз слово w_i встретилось в тексте

$\text{text} = "я хочу спать"$. $V = \begin{matrix} я & кот & кровать & шкаф & кухня \\ w_1 & w_2 & w_3 & w_4 & w_5 \end{matrix}$

В векторе 1)

$f_1(\text{text}) = 1$ $f_2(\text{text}) = 0$ $f_3(\text{text}) = 0$ $f_4(\text{text}) = 1$ $f_5(\text{text}) = 1$

или $f(\text{text}) = \begin{pmatrix} 1 & 0 & 0 & 1 & 1 \end{pmatrix}$

- вектор признаков

В начале 2

$$f("я хочу чай, чай, чай") = (1 \ 0 \ 0 \ 3 \ 1)$$

3) tf-idf.

$$f_i(\text{text}) = \text{tf} \cdot \text{idf}, \text{ где}$$

tf = term frequency = сколько раз
выпало слово w_i

idf = inverse document frequency

$$= \frac{1}{\log(1 + \text{кол-во документов со словом } w_i)}$$

$$\text{tf-idf слова "он"} = \frac{\text{tf(он)}}{\log(1 + \text{всего doc})} \quad \left. \vphantom{\frac{\text{tf(он)}}{\log(1 + \text{всего doc})}} \right\} \text{ много}$$

это bag-of-words признаки, т.к. не учитывается порядок слов

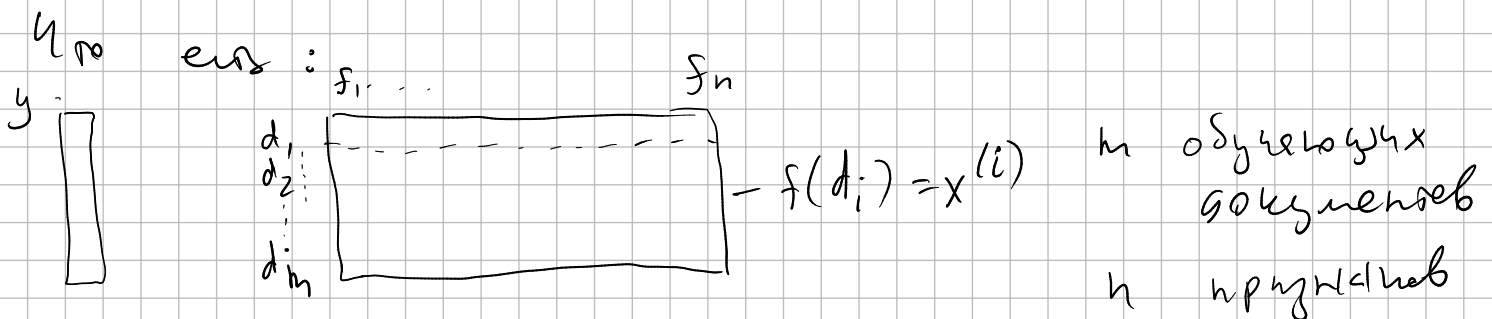
n-грамм признаки:

$f(\text{text})$ - сколько раз в тексте встречается опреде.
n-грамм.

Итого, \exists f-вектор признаков, q-ия считает сразу
все признаки. $f(\text{text}) = (f_1(\text{text}), f_2(\text{text}), \dots)$

Вектор признаков обычно очень длинный. Может быть
10 000-к признаков.

- сделать вектор короче: задать уменьшенный размер-
ность.



$x_j^{(i)}$ — это j -ый признак i -го документа.

$y^{(i)}$ — это правильные классы для документа d_i ($x^{(i)}$)

на первом этапе обучения — линейные переделы

↑ — предсказание.

Алгоритм определения класса $\hat{y}^{(i)}$ для $x^{(i)}$

для класса y_j

$$\theta_1^{(j)} x_1^{(i)} + \theta_2^{(j)} x_2^{(i)} + \theta_3^{(j)} x_3^{(i)} + \dots + \theta_n^{(j)} x_n^{(i)} = x^{(i)} \cdot \theta^{(j)}$$

так $x^{(i)}$

вычисляем для каждого класса (j) это значение,
выбираем (j), где значение максимальное.

Пример. Спам / Не спам

$\theta^{(0)}$

$\theta^{(1)}$

$V = \{ \text{купить, привет, заказ, спам} \}$

$$\theta^{(0)} = \begin{pmatrix} 2 \\ -0.1 \\ -2 \\ 1 \end{pmatrix}$$

$$\theta^{(1)} = \begin{pmatrix} -2 \\ 0.1 \\ 2 \\ -0.5 \end{pmatrix}$$

← \approx оспаривает ли или не оспаривает.

$\text{text} = \text{"привет, это спам"}$. f_i — i -ый признак

$$f(\text{text}) = (0, 1, 0, 1)$$

для класса 0: $f(\text{text}) \cdot \theta^{(0)} = 0 \cdot 2 - 1 \cdot 0.1 - 0 \cdot 2 + 1 \cdot 1 = 0.9$

для класса 1: $f(\text{text}) \cdot \theta^{(1)} = 0 \cdot (-2) + 1 \cdot 0.1 + 0 \cdot 2 - 1 \cdot 0.5 = -0.4$

$0.9 > -0.4 \Rightarrow$ класс 0 \Rightarrow спам.

Получите, что для определения класса нужно знать

2. Обучение - это поиск θ .

разные методы обучения - разные способы поиска θ .

Замечание увидим нейросеть в том, что мы

обсуждаем.

Текст попадает на входной слой

если слово 1 $\rightarrow 0$

1 $\rightarrow 0$

если слова нет

0 $\rightarrow 0$

0

0

0

0

0

входные

нейроны

(для признаков)

нейрон = слово

слои классов

$x \cdot \theta$

1

2

- аналогично

max

Нейрон каждой входной сети умножает на вес \rightarrow выбирает все значения с весом

Наивный Байесовский классификатор
Naive

Считаем, что текст является случайным словом, т.е. это результат случайного процесса

$p(y | \text{text}) = ?$

- это задача классификации.

Процесс:

Мат выбирает случайно один из классов в соответствии с априорным распределением.

Например, считаем, что спам: 5%

не спам: 95%

(мы будем искать эти θ -ы, важно, что мы считаем, что они есть)

Условие. Для каждого класса есть представление слов. Например,

в классе "чем" слова представлены

так:

привет	кухня	он	чем
30%	20%	30%	60%

≠ 100%

в классе "не чем"

привет	кухня	он	чем
60%	20%	30%	10%

для каждого слова нужно определить, есть оно или нет.

Условие "текст" сгенерирован. Словом не имеет значения какого класса:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad \left. \begin{array}{l} p(y^{(0)} | \text{text}) \\ p(y^{(1)} | \text{text}) \end{array} \right\}$$

то, что должно
— быть.

← same because same

$$p(y^{(0)} | \text{text}) = \frac{p(\text{text} | y^{(0)}) \cdot p(y^{(0)})}{p(\text{text})}$$

$$p(y^{(1)} | \text{text}) = \frac{p(\text{text} | y^{(1)}) \cdot p(y^{(1)})}{p(\text{text})}$$

где $p(\text{text})$ — max number he generates