

# Классификация текстов.

Задача. есть конечное мн-во классов:

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

не зная того, 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, или конечное мн-во.

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 версии 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

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

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

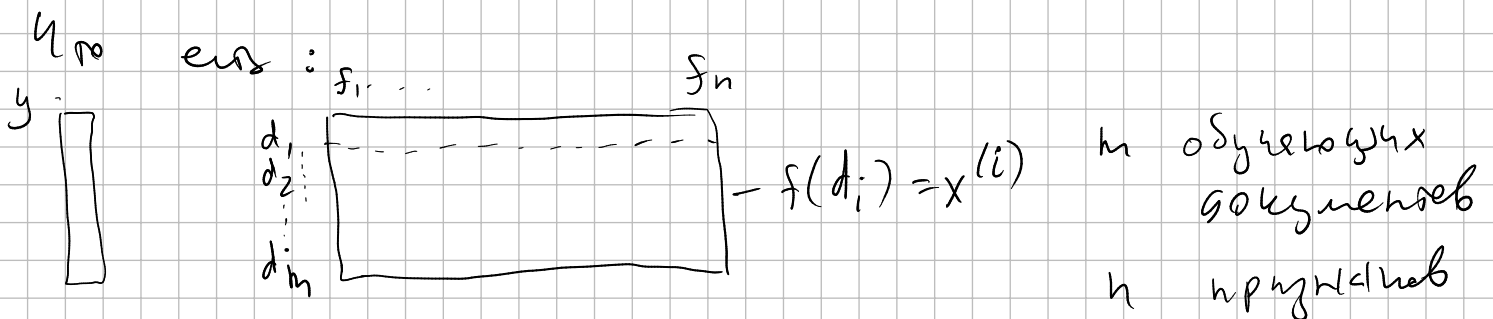
n-грамм признаки:

$f(\text{text})$  - 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

$\theta^{(0)}$

$\theta^{(1)}$

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

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

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

←  $\approx$  оспаривает ли или не оспаривает.

$\text{Text} = \text{"приват, это спам"} \quad f_i - 1/\text{признак}$

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

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

$0.9 > -0.4 \Rightarrow \text{класс } 0 \Rightarrow \text{спам}$

Получите, что для определения класса нужно знать

2. Обучение - это поиск  $\theta$ .

разные методы обучения - разные способы поиска  $\theta$ .

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: 50%  
не спам: 50%

(мы будем искать эти  $\theta$ -ы, важно, что мы считаем, что они есть)

Услов. Для каждого класса есть представление слов. Например,

в классе "чем" слова представлены

так:

привет	кухня	он	чем
30%	20%	10%	40%

= 100%

в классе "не чем"

привет	кухня	он	чем
60%	5%	30%	5%

= 100%

для каждого слова нужно определить, есть оно или нет.

Услов "текст" стандартизован. Строим не текст, а разреженные классы:

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

$$\left\{ \begin{array}{l} p(y^{(0)} | \text{text}) \\ p(y^{(1)} | \text{text}) \end{array} \right\}$$

то, что должно  
— ответ.

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

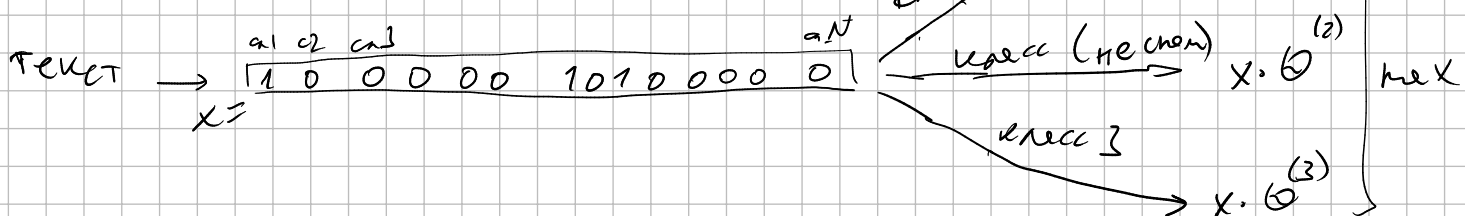
← same because same

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

→ next

где  $p(\text{text})$  — max number he gets to  $p(\text{text})$ .

Например, Митинский класс.



can text  $w_1, w_2, \dots, w_N$

$$p(y^{(0)} | \text{text}) \sim p(\text{text} | y^{(0)}) \cdot p(y^{(0)}) \stackrel{\text{наблюдать}}{=} p(w_1 | y^{(0)}) \cdot p(w_2 | y^{(0)}) \cdot \dots \cdot p(w_N | y^{(0)}) \cdot p(y^{(0)})$$

вероятности  
оценить.

$$p(w | y^{(i)}) \approx p(y^{(i)}) \quad \text{можно}$$

$$p(w | y^{(i)}) = \frac{\text{число раз встретился слово } w \text{ в классе } i}{\text{всего слов в тесте из класса } i}$$

$$p(y^{(i)}) = \frac{\text{число тестовых классов } i}{\text{число всего тестов.}}$$

1) проблема вероятности много - underflow.

$$-\log p(y^{(0)} | \text{text}) = -\log(w_1 | y^{(0)}) - \log(w_2 | y^{(0)}) - \dots - \log(w_N | y^{(0)}) - \log p(y^{(0)}).$$

i.e.  $\Rightarrow$   $\begin{bmatrix} -\log p(w_1 | y^{(0)}) & -\log p(w_2 | y^{(0)}) & \dots \end{bmatrix}$

слова  $w_1, \dots, w_N$   $\rightarrow$   $N$  слов в классе  $y^{(0)}$   
 $-\log p(y^{(0)})$

2. regularization.

$$p(w_i | y^{(i)}) = \frac{c(w_i | y^{(i)}) + \alpha}{\text{всего слов в классе } i + \alpha \cdot N}$$

$\alpha = 0.1$   
 $\downarrow$  можно раз  
 $\uparrow$  слово в классе  $y^{(i)}$   
 $\uparrow$  слово в классе  $y^{(i)}$

Параметр  $\alpha$  - можно не считать

$\alpha = 0.01$   
 $0.1$   
 $0.05$   
 $0.2$

что оптимизировать: точность =  $\frac{\text{кол-во прав ответов}}{\text{общее кол-во вопросов}}$

все примеры  $\rightarrow$   $\text{train}$   
 $\text{validation}$   
 $\text{test}$

K-fold verification.

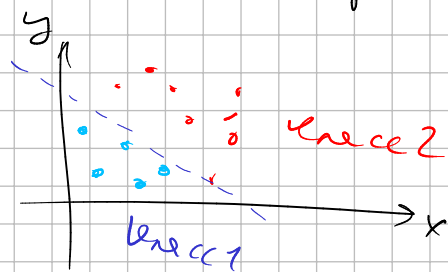


3. подаём пример  $x^{(i)}$  на сеть, получаем класс  $\hat{y}$ .
4. если  $y^{(i)} = \hat{y}$ , увеличиваем счётчик.
5. иначе уменьшаем веса на активированных вершинах, увеличиваем — на активированных соседях вершин.
6. Повторяем, пока не достигнем примера  $x^{(i)}$ .

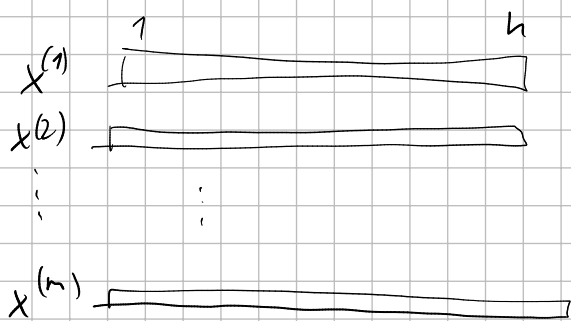
— Особенности алгоритма.

1) online-режим. Классифицируем по мере получения данных. Т.е. можно пользоваться до конца обучения.

2) Алгоритм сходится, если исх. мн-во примеров линейно разделимо



### Логистическая регрессия



и протестируем  $\rightarrow y^{(i)}$  — прав класс

$\rightarrow y^{(m)}$

все классы  $y_1, y_2, \dots, y_k$

пример  $x$ , класс  $y$

$p(y|x)$  — вер-ть класса  $y$  при условии текста  $x$

$$p(y|x) \approx \frac{e^{+x \cdot \theta^{(y)}}}{e^{+x \cdot \theta^{(y_1)}} + e^{+x \cdot \theta^{(y_2)}} + \dots + e^{+x \cdot \theta^{(y_k)}}} \in [0, 1]$$

$$\sum p(y_i|x) = 1$$



Чго хотим оптимизировать @:

$$p(\text{все классы правильно}) = p(y^{(1)}|x^{(1)}) \cdot p(y^{(2)}|x^{(2)}) \dots p(y^{(n)}|x^{(n)})$$

→ max