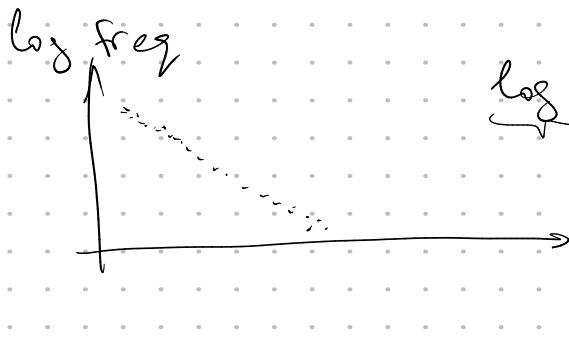


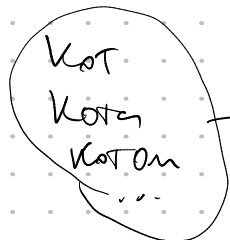
Zipf



$$\log \text{freq} = C - \log \text{rank}$$

$\underbrace{\log \text{freq}}_y = C - \underbrace{\log \text{rank}}_x$

1) попробуйте "нормализовать" слова.



кот
1 слово

морф. форма

?? став(цы)
 став → став(ка)
 (цы → ип, ег)
 га → шфын.

→ лена (офенд
 окончание)

playing → play
 player → ply
 dancing → dance (e)

Анн. Porter Stemmer
 - алг, правил отрезания
 окончаний.

Дан язык I
 Нормализуем через морф. анализатор

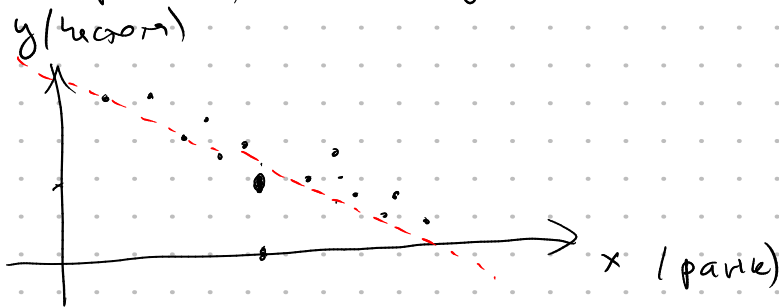
- берем первый в-т анализа.

morphology 2

- для слова берет
 варианты из-д форм,
 если не знает -
 угадывает.

Дан язык II

Попытаемся подобрать прямую, проходящую ближе всего к
 нашим точкам.



Задача. Давайте попытаемся предсказать частоту по рангу

Есть примеры (обучающие)

	ранки	весы
микро	100	10
	20	50
	120	?

— прогнозы

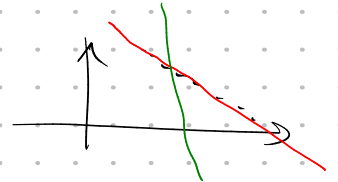
Задача регрессии —
 непрерывной величины —
 задача регрессии

(см. классификация)

Линейная регрессия, предполагается.

$$Y = \theta + k \cdot X + \text{Err} \sim N(0, \sigma)$$

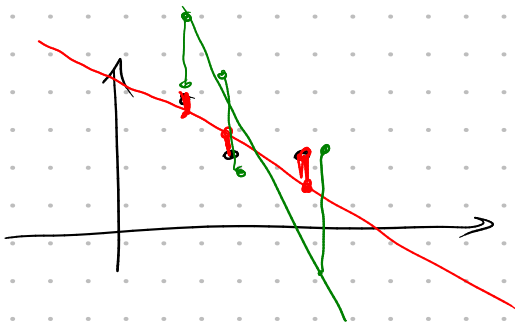
где Y — прогнозы, X — дано



Найти коэффициенты θ и k = подобрать

при выборе данных так, чтоб $\sum_{i=1}^N \text{err}_i^2 \rightarrow \min$

Сумма по примерам.



оч. зел >> оч. красн

Как подобрать прямо?

I. numpy с помощью функции хранит в numpy-матрицах.

numpy — создание и работа с матрицами из чисел.

значительно эффективнее и по памяти, и по времени,
 по сравнению со списками и словарями.

```
import numpy as np
```

```
x = np.array([10, 20, 30])
```

```
y = np.array([[10, 20], [20, 40]])
```

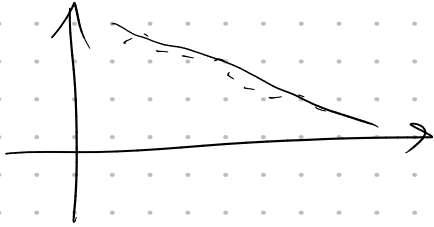
10	20	30
----	----	----

10	20
20	40

одушка модель

X_1	y
0	7
1	2
1.5	4
2	3

$\log f \sim f$



$$\log \text{freq} = C - \log \text{rank}$$

проверить, что получается ≈ -1

$$\log \text{freq} = b + k \cdot \log \text{rank}$$

$k = -1$

$$\text{freq} = \frac{C}{\text{rank}^{|k|+1}}$$

N -грамм модель.

Строки модели строки.

] язык \rightarrow это можно предложить.

Например, в русском

"я могу спать" \in Русский

"я могу насос" \notin Русский

"мы стали более лучше одеваюсь" \notin Русский

лучше иначе:

$$P(\text{предложение}) \rightarrow \text{член от } 0 \text{ до } 1$$

$$[0, 1]$$

вероятность встретить предлож.

$$P(\text{"я могу спать"}) \gg P(\text{"я могу насос"})$$

Для вероятности могут быть $\sum_{\omega \in \Omega} P(\omega) = 1$

→ то вероятности могут быть

Могут быть ($P: \Omega \rightarrow \mathbb{R}$)

знач.

— при генерации текста, выбор следующего символа.

I sit by the table → \downarrow ω \neq ω (1)
 \downarrow ω \neq ω (2)

$$P(1) > P(2)$$

Видею эти вероят.

— более ошибок. \downarrow ω \neq ω (1) — обн.
 \downarrow ω \neq ω (2) — н.с.

$P(1) < P(2) \Rightarrow$ н.с. во время в беге "коп"

У нас будет N -правая модель.

$P(\omega_1, \omega_2, \omega_3, \omega_4, \dots, \omega_n) =$ гипотеза: нег. способ
 зависит от $N-1$
 предыдущих.

$N=2$

$$= P(\langle s \rangle \omega_1, \omega_2, \omega_3, \dots, \omega_n \langle /s \rangle) :=$$

$$= P(\omega_1 | \langle s \rangle) \cdot P(\omega_2 | \omega_1) \cdot P(\omega_3 | \omega_2) \cdot \dots \cdot (P(\langle /s \rangle | \omega_n))$$

← вероятности симв. ω_2 после ω_1

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

ген. б-то

когда $P(\text{чет} \geq 3) = \frac{P(\text{чет} \geq 3)}{P(\geq 3)} = \frac{2/6}{4/6} = 1/2$

$\frac{1}{6} \frac{1}{6} \frac{1}{6} \frac{1}{6} \frac{1}{6} \frac{1}{6}$

$$P(\text{нечет} | \geq 4) = \frac{P(\text{нечет} \geq 4)}{P(\geq 4)} = \frac{1/6 \cdot 12352}{3/6} = 1/3$$

$$P(\text{нечет}) = 3/6 = 1/2$$

$N=3$

$$P(\langle s \rangle \langle s \rangle w_1 w_2 w_3 \dots w_n \langle s \rangle) =$$

$$= P(w_1 | \langle s \rangle \langle s \rangle) \cdot P(w_2 | \langle s \rangle w_1) \cdot P(w_3 | w_1 w_2) \cdot \dots$$

$$\cdot P(w_i | w_{i-2} w_{i-1}) \cdot \dots \cdot P(\langle s \rangle | w_n w_n).$$

Как оценить $P(u | w)$? слово u после слова w ?

Метод максимального правдоподобия (MLE)

maximum likelihood estimation

Дан корпус, где много слов / предложений / слов.

можно считать $P(\text{корпус}) = P(\text{предл} 1) \cdot P(\text{предл} 2) \cdot \dots$

правдоподобие $\cdot P(\text{предл} n) =$

$$= P(w_1 | \langle s \rangle) \cdot \dots$$

$\rightarrow \max.$

Max достигается при

$$P(u | w) = \frac{c(wu)}{c(w)}$$

число раз в корпусе это число слов w , перед u

число слов w в корпусе.

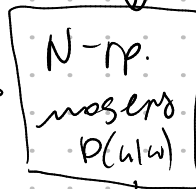
Модель построена.

предположение

Дан корпус

→ читаем $P(u|w)$

→



↓
 $P(\text{предположение})$

$$P(s \text{ из } \text{модель}) = P(s|<s>) \cdot P(\text{из } s) \cdot P(\text{модель} | \text{из } s) \cdot P(|<s> \text{модель})$$

- $P(u|w) = 0$, если в корпусе нет слова w .

$P(\dots w \dots) = 0$. возможно иметь 0 вероятностей, они будут часто встречаться.

Смраживание - $P(u|w)$ = что использовать значение, чтобы всегда $\neq 0$.

P гоним остальные вероятности.

$$\sum_{u \in V} P(u|w) = 1$$

V - весь словарь.

Например

$$P(s|<s>) + P(ок|<s>) + P(оки|<s>) + P(спл|<s>) + \dots + P(спер|<s>) = \int_{\text{все слова}}$$

- как сравнить качества разных смраживаний?

в маш. обучении - подбор гиперпараметров

- оценить качество модели

перерыв до 10:45

Пример: Копус N=2

- <5> г вуху сра <1>
- <5> г вуху суга <1>
- <5> г вуху сир <1>
- <5> г ем сир <1>
- <5> г сра <1>

результат обучения

$$\frac{C(\langle 5 \rangle \langle 1 \rangle)}{C(\langle 5 \rangle)} = \frac{C(\langle 1 \rangle \text{ вуху})}{C(\langle 1 \rangle)} = 3$$

$$\Sigma = 1 \rightarrow$$

$$\frac{C(\text{вуху сра})}{C(\text{вуху})} = \frac{1}{3}$$

w \ u	<5>	г	вуху	ем	сра	суга	сир	<1>
<5>	0	1	0	0	0	0	0	0
г	0	0	3/5	1/5	1/5	0	0	0
вуху	0	0	0	0	1/3	1/3	1/3	0
ем	0	0	0	0	0	0	0	0
сра	0	0	0	0	0	0	0	1
суга	0	0	0	0	0	0	0	1
сир	0	0	0	0	0	0	0	1
<1>	0	0	0	0	0	0	0	0

$$P(\langle 1 \rangle \text{ вуху}) = P(\langle 1 \rangle | \langle 5 \rangle) \cdot P(\text{вуху} | \langle 1 \rangle) \cdot P(\langle 5 \rangle | \text{вуху}) = 1 \cdot \frac{3}{5} \cdot 0 = 0$$

Узависене от $P(u|w) = 0$ ранее премир = 1 0.1 0.01

$$P(u|w) = \frac{C(wu) + \lambda}{C(w) + \lambda |V|}$$

где $\lambda = 1$

w \ u	<5>	г	вуху	ем	сра	суга	сир	<1>
<5>	0	1	0	0	0	0	0	0
г	1/13	1/13	4/13	2/13	2/13	1/13	1/13	1/13
вуху	0	0	0	0	1/3	1/3	1/3	0
ем	0	0	0	0	0	0	0	0
сра	0	0	0	0	0	0	0	1
суга	0	0	0	0	0	0	0	1
сир	0	0	0	0	0	0	0	1
<1>	0	0	0	0	0	0	0	0

$$\frac{0}{5} \rightarrow \frac{1}{5+1.8}$$

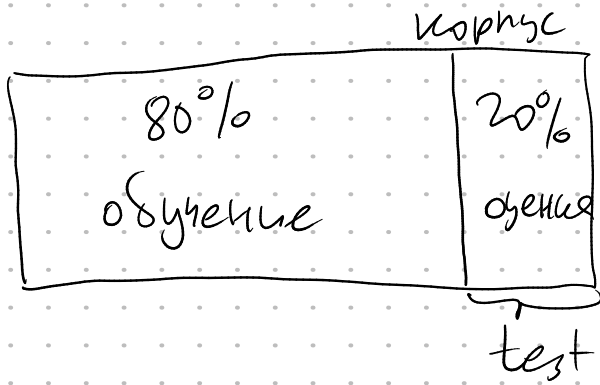
$$\frac{1}{5} \rightarrow \frac{2}{5+1.8}$$

$$\frac{3}{5} \rightarrow \frac{4}{5+1.8}$$

При стриживании все неизвестные слова (не из корпуса) заменяются на UNK.

Качество модели оценивается на корпусе.

Обычно для оценки корпуса не используют



$$P(\text{тест. корпус}) = \\ = P(\text{слова 1}) \cdot P(\text{слова 2}) \cdot \\ P(\text{слова 3}) \cdot \dots \rightarrow \text{max}$$

Perplexity \rightarrow насколько нет D , но \rightarrow with