b Zipf

log freq



$$\underbrace{\log freq}_{y} = C - \underbrace{\log rank}_{x}$$

log rank

1) попробуйте "нормализовать" слова.



Кот
Кота
Котом
... → Кот

1 слово

?? → стало (сущ.)
стал → стал (гл.)
(сущ → ИП, ед.ч)
(гл → инфин.)
...

→ лемма
норм. форма
→ стем (отрезать
окончание)

playing → play
player → play
dancing → danc (e)
...

Англ. Porter Stemmer
— алг, правила отрезания
окончаний.

Дом. задание I

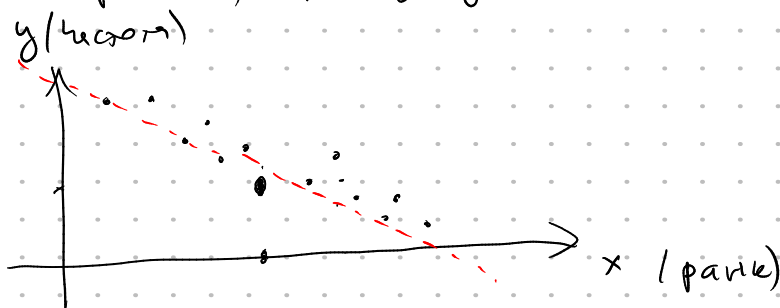Нормализуем через морф. анализатор

— берём первый в-т анализа.

pymorphy 2
— для слова выдаёт
варианты нач. форм,
если не знает —
угадывает.

Дом. задание II

Попытаемся подобрать прямую, проходящую ближе всего к
нашим точкам.



y (частота)

x (ранг)

Задача Давайте пытаемся предсказывать частоту по рангу

Есть примеры (обучающие)

ранк частота

много {
| 100 | 10 |
| 20 | 50 |
| 120 | ? | — предсказать

Задача предсказания непрерывной величины — задача [регрессии]
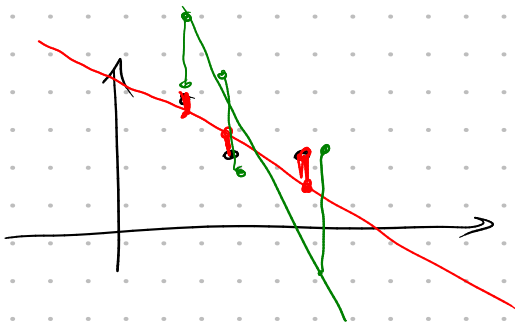
(см. классификацию)

Линейная регрессия, предполагает.

$$Y = \beta + k\,X + \boxed{Err} \sim N(0, \sigma)$$

$\underbrace{Y}$ — надо предсказать

$\underbrace{X}$ — дано

Научиться предсказывать = подобрать $\beta$ и $k$.
при подборе делаем так, чтоб $\sum\limits_{i=1}^{N} err_i^2 \longrightarrow min$

сумма по примерам.

ош.        ош.
зел $\gg$ красн

Как подобрать прямую?

I. питру

Данные с числами лучше хранить в питру-матрицах.

питру — создание и работа с матрицами из чисел.
значительно эффективнее и по памяти, и по времени,
по сравнению со списками и циклами.

```
import питру as np

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

y = np.array([[10,20], [30,40]])
```
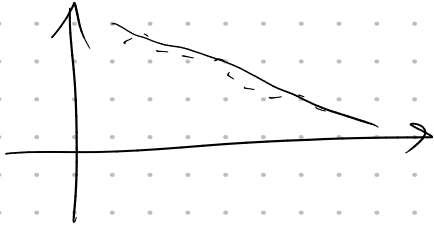
| 10 | 20 | 30 |

| 10 | 20 |
| 30 | 40 |

объект модели

$$\begin{array}{cc} X & y \end{array}$$

$$\begin{array}{cc} 0 & 7 \\ 1 & 2 \\ 1.5 & 4 \\ 2 & 3 \end{array}$$

в Zipf



проверьте, что получилось $\approx -1$

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

$$\|$$

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

$$k = -1$$

$$freq = \frac{C}{rank^{gen+1}}$$

$N$-грам модели.

Строим модель языка.

$\exists$ язык — это множество предложений.

Например, в русском

"я иду спать" $\in$ Русский

"я иду насос" $\notin$ Русский

"мы стали более лучше одеваться"
$\notin$ Русский

Лучше иначе:

$$P(\text{предложение}) \longrightarrow [0,1]$$

Число от 0 до 1

вероятность встретить предлож.

$P(\text{"я иду спать"}) >>> P(\text{"я иду насос"})$

Для вероятностной модели важно $\sum\limits_{\text{предл}} P(\text{предл}) = 1$

Это вероятностная модель языке.

Модель языке $(P: \text{предл} \to \text{числа})$

зачем

— при генерации текста, выбор лучшего варианта.

I sit by the table $\to$ Я сижу у слона (1)
Я сижу у таблицы (2)

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

Выберем 1ый вариант.

— поиск ошибок. Я гляжу кита $^{(1)}$ — было
Я гляху кота $^{(2)}$ — и.б.

$$P(1) < P(2) \implies \text{и.б. Вы имели в виду "кота".}$$

У нас будет N-грам модель.

$$P(\underbrace{w_1}_{\text{слова}}, w_2, w_3, w_4 \ldots w_n) = \begin{array}{l}\text{Гипотеза: след слово}\\ \text{зависит от } N-1\\ \text{предыдущих.}\end{array}$$

N=2

$$= P(\underset{\text{начало}}{<s>}\; w_1, w_2\, w_3 \ldots w_n\; \underset{\text{конец}}{</s>}) :=$$

$$= \underbrace{P(w_1 | <s>)}_{} \cdot P(w_2 | w_1) \cdot P(w_3 | w_2) \cdot \ldots \cdot (P(</s> | w_n) \cdot$$

вероятность слова
$w_2$ после $w_1$

P-ть слева $w_1$ после $<s>$

$$P(A|B) := \frac{P(AB)}{P(B)}$$
усл. в-ть

кубик
$$P(\text{чет} | \geqslant 3) = \frac{P(\text{чет} \geqslant 3)}{P(\geqslant 3)} = \frac{2/6}{4/6} = 1/2$$

$$1\,2\,\boxed{3\,4\,5\,6}$$
$$\tfrac{1}{6}\,\tfrac{1}{6}\,\tfrac{1}{6}\,\tfrac{1}{6}\,\tfrac{1}{6}\,\tfrac{1}{6}$$

$$P(\text{нечёт} | \geqslant 4) = \frac{P(\text{чёт} \geqslant 4)}{P(\geqslant 4)} = \frac{1/6}{3/6} = \frac{1}{3}$$

$$P(\text{нечёт}) = \frac{3}{6} = \frac{1}{2}$$

$N = 3$

$$P\left( <s><s> \; w_1, \; w_2 \; w_3 \dots \; w_n </s> \right) =$$

$$= P\left( w_1 | <s><s> \right) \cdot P\left( w_2 | <s> w_1 \right) \cdot P\left( w_3 | w_1 w_2 \right) \dots$$

$$\dots \cdot P\left( w_i | w_{i-2} w_{i-1} \right) \dots P\left( </s> | w_{n-1} w_n \right).$$

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

Метод максимального правдоподобия (MLE)

m
a
x
i
k
e
l
y
h
o
od
$+$
s
t
i
m
a
t
i
o
n

(Дан) корпус, где много текстов/предложений/слов.

можно посчитать   $\underline{P(\text{корпуса})} = P(\text{предл 1}) \cdot P(\text{предл 2}) \cdot \dots$

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

$$= P(w_1 | <s>) \dots \dots \dots \dots \longrightarrow \max.$$

Max достигается при

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

← сколько раз в корпусе было написано слово w, потом u

↑ сколько слов w в корпусе.

Модель построения.

Дан корпус $\longrightarrow$ считаем $\longrightarrow$ [ N-гр. модель $P(u|w)$ ]  предложение

$P(u|w)$   →  $P(предложения)$

$$P(\text{я иду спать}) = P(\text{я}|<s>) \cdot P(\text{иду}|\text{я}) \cdot P(\text{спать}|\text{иду}) \cdot P(</s>|\text{спать})$$

— $P(u|w) = 0$, если в корпусе нет подряд $wu$.

$P(..... wu ...) = 0$. Нежелательно иметь 0 вероятность, они будут часто встречаться.

Сглаживание — $P(u|w) =$ чуть исправит значение, чтобы всегда $\neq 0$.

$P$ должна остаться вероятностью.

$$\sum_{u \in V} P(u|w) = 1 \qquad V - весь словарь$$

например

$$P(\text{я}|<s>) + P(\text{он}|<s>) + P(\text{они}|<s>) +$$
$$P(\text{стол}|<s>) + ... + P(\text{ящер}|<s>) = 1$$
все слова.

— как сравнить качество разных сглаживаний? в маш. обучении — подбор гиперпараметров.

— оценить качество модели

перерыв до 10:45

# Пример:   Корпус      N=2

<s> Я вижу стол </s>
<s> Я вижу стул </s>
<s> Я вижу сыр </s>
<s> Я ем сыр </s>
<s> Я стол </s>

Результат обучения

$$\frac{C(\text{<s>}\,Я)}{C(\text{<s>})} = \frac{C(Я\ \text{вижу})=3}{C(Я)=5}$$

$\Sigma = 1 \longrightarrow$

$$\frac{C(\text{вижу стол})}{C(\text{вижу})} = \frac{1}{3}$$

| w \ u | <s> | Я | вижу | ем | стол | стул | сыр | </s> |
|---|---|---|---|---|---|---|---|---|
| <s> | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Я | 0 | 0 | 3/5 | 1/5 | 1/5 | 0 | 0 | 0 |
| вижу | 0 | 0 | 0 | 0 | 1/3 | 1/3 | 1/3 | 0 |
| ем | 0 | | | | | | 1 | |
| стол | 0 | | | | | | | 1 |
| стул | 0 | | | | | | | 1 |
| сыр | 0 | | | | | | | 1 |
| </s> | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

$$P(Я\ \text{вижу}) = P(Я|\text{<s>}) \cdot P(\text{вижу}|Я) \cdot P(\text{</s>}|\text{вижу}) =$$

$$= 1 \cdot \frac{3}{5} \cdot 0 = 0.$$

Избавление от $P(u|w)=0$.

гиперпараметр.
$= 1\quad 0.1\quad 0.01$

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

в его словаре

$\boxed{\lambda = 1}$

| w \ u | <s> | Я | вижу | ем | стол | стул | сыр | </s> |
|---|---|---|---|---|---|---|---|---|
| <s> | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Я | 1/13 | 1/13 | 4/13 | 2/13 | 2/13 | 1/13 | 1/13 | 1/13 |
| вижу | 0 | 0 | 0 | 0 | 1/3 | 1/3 | 1/3 | 0 |
| ем | 0 | | | | | | 1 | |
| стол | 0 | | | | | | | 1 |
| стул | 0 | | | | | | | 1 |
| сыр | 0 | | | | | | | 1 |
| </s> | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

$\Sigma = 1$

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

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

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

При стэхивании все неизвестные слова (не из корпуса) заменяются на <u>UNKN</u>.

<u>Качество модели.</u> Оценивает на корпусе.

обычно делим исходный корпус на части

корпус



80% обучение

20% оценка

test

$P(\text{тест. корпуса}) =$

$= P(\text{предл 1}) \cdot P(\text{предл 2}) \cdot$

$P(\text{предл 3}) \cdot \ldots\ldots \longrightarrow max$

Perplexity $\longrightarrow$ примерно как $P$, но $\longrightarrow min$