

N-грамм модели

Модель языка: p : предложение $\rightarrow [0, 1]$

Каждому предложению сопоставляется "вероятность"

$$\sum p(s) = 1.$$

Словарь

- выбор самого вероятного w_i из нескольких
- генерация текста, выбор вероятных предложений

$$p(s \text{ из } w) > p(s \text{ из } e_i) > 0$$

$$p(s) = p(w_1 w_2 w_3 \dots w_N) =$$

и слов в предл.

$$= p(w_N | w_1 w_2 \dots w_{N-1}) \cdot p(w_1 \dots w_{N-1})$$

$$p(AB) = p(B|A) \cdot p(A)$$

$$= p(w_N | w_1 w_2 \dots w_{N-1}) \cdot$$

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

$$p(w_{N-1} | w_1 w_2 \dots w_{N-2}) \cdot$$

$$p(w_2 | w_1) \cdot$$

$$p(w_1)$$

*** я вижу кот

$$\begin{aligned} p(\text{я вижу кот}) &= \\ &= p(\text{кот} | \text{я вижу}) \cdot \\ &\quad \cdot p(\text{я вижу}) \end{aligned}$$

= в начале предложения есть слово *

≈ оценить w_i как $p(w_i | w_1 w_2 \dots w_{i-1})$ как $p(w_i | \underbrace{w_{i-n+1} \dots w_{i-1}}_{n \text{ слов}})$

$$\left. \begin{aligned} &= p(w_N | w_{N-n+1} \dots w_{N-1}) \cdot \\ & p(w_{N-1} | \dots w_{N-2}) \cdot \\ & p(w_1 | \underbrace{*** \dots *}_{n-1} \text{ " "}) \end{aligned} \right\} (1)$$

н-граммная зависимость

$$p(W_N = w_N | \underbrace{W_{N-1} = w_{N-1}}_{\text{конкр. значение}} \underbrace{W_{N-2} = w_{N-2}}_{\text{сл. вел}})$$

если $n=2$.

$$p(w_1, w_2, \dots, w_n) = p(w_1 | *) \cdot p(w_2 | w_1) \cdot p(w_3 | w_2) \cdot \dots \cdot p(w_n | w_{n-1})$$

если $n=1$

$$p(w_1, w_2, \dots, w_n) = p(w_1) \cdot p(w_2) \cdot \dots \cdot p(w_n)$$

Как оценить $p(w_i | w_j)$?

нужен корпус текстов для обучения модели.
считаем встречаемость n -грамм, т.е. последовательностей из n слов. $c(\dots)$ = кол-во раз (count)

$\begin{cases} \text{я} & \text{из} & \text{глагол}, & \text{я} & \text{вых} & \text{кого}, & \text{я} & \text{вых} & \text{сон.} \end{cases}$

$$c(\text{я из}) = 1 \quad c(\text{я вых}) = 2$$

$$p(w_n | w_1, w_2, \dots, w_{n-1}) = \frac{c(w_1, \dots, w_{n-1}, w_n)}{\sum_{w \in V} c(w_1, \dots, w_{n-1}, w)}$$

В-т увидит w_n при условии, что перед этим прочитали w_1, \dots, w_{n-1}

$\sum_{w \in V}$
V-словарь

$$p(\text{вых} | \text{я}) = \frac{c(\text{я вых})}{c(\text{я})} = \frac{2}{3}$$

$$= \frac{c(w_1, \dots, w_{n-1}, w_n)}{c(w_1, \dots, w_{n-1})} \cdot (2)$$

Всё, модель построена!!!

$$p(w_1, \dots, w_n) = (1) = \text{оценивается через (2)}.$$

При генерации текста:

given $n-1$ слово w_1, w_2, \dots, w_{n-1}

given w_n выбирается случайно не распреде-

либо $p(\cdot | w_1, \dots, w_{n-1})$ потом w_{n-1} по
распределению $p(\cdot | w_2, w_3, \dots, w_n)$.

Сглаживание

Проблема 1. $p(\text{я вижу тире и звездочку}) = p(\text{тире и звездочка} | \dots)$
- не было слова.

$p(\text{тире} | \dots) = 0 \Rightarrow$ результат тоже 0.

Обычно все неизвестные слова заменяют на $\langle \text{UNK} \rangle$
(unknown)

Считают, что словарь - это все встречающиеся слова
+ слово UNK.

$p(w_n | w_1, \dots, w_{n-1}) = 0$, если слово w_n не встречалось
никогда раньше w_1, \dots, w_{n-1} . "я вижу слит"

\Rightarrow p всего предположить тоже
будет 0.

в корпусе есть все
эти слова, но нет этой
фразы!

при $n=1$

Сглаживание дает все
вероятности $\neq 0$.

$$p(\text{слит} | \text{я вижу}) = \frac{0}{\dots} = 0$$

Варианты сглаживания.

1. Лаплас (Lid...)

$$p(w_1 | w_1, \dots, w_{n-1}) = \frac{c(w_1, w_2, \dots, \overset{\text{всегда берем}}{\downarrow} \bar{w}_n) + \alpha}{c(w_1, \dots, w_{n-1}) + \alpha |V|}$$

V - словарь

$|V|$ - размер словаря

$p(\cdot | w_1, \dots, w_{n-1})$ вер-ть

Лаплас: $\alpha = 1$ (можно $\alpha = 0,1$ и т.п.)

2. Интерполяция, где откуда

$$p(w_n | w_1, \dots, w_{n-1}) = \frac{c(w_1, w_2, \dots, w_n)}{c(w_1, w_2, \dots, w_{n-1})}$$

← n грамм

$$\frac{c(w_2, \dots, w_n)}{c(w_2, \dots, w_{n-1})}$$

← n-1 грамм

$$\frac{c(w_3, \dots, w_n)}{c(w_3, \dots, w_{n-1})}$$

← n-2 грамм

т.е. $\tilde{p}(\text{cheir} | \text{я чгг}) = \frac{p(\text{cheir} | \text{я чгг})}{\frac{c \dots}{c \dots}} = \frac{p(\text{cheir} | \text{чгг})}{\frac{c \dots}{c \dots}}$

↙ среза*

$$p(\text{cheir}) = \frac{c \dots}{c \dots}$$

Как использовать оценки меньшей размерности?
 интерполяция. сложить с коэффициентами.
 } n=3 где пример

$$\tilde{p}(w_3 | w_1, w_2) = \alpha_3 p(w_3 | w_1, w_2) + \alpha_2 p(w_3 | w_2) + \alpha_1 p(w_3)$$

$$\alpha_3 + \alpha_2 + \alpha_1 = 1 \quad \alpha_i \geq 0$$

α_i - его коэффициент.

против откат. Если $p(w_3 | w_1, w_2) = 0$, взять $p(w_3 | w_2)$ Если тоже 0 $p(w_3)$

что $p(\cdot | w_1, w_2)$ - такая вероятность, генерируемая

$$\tilde{p}(w_3 | w_1, w_2) = \begin{cases} 0.6 p(w_3 | w_1, w_2), & \text{если } \neq 0 \\ 0.4 p(w_2 | w_2) \cdot 0.6 & \text{иначе, если } \neq 0 \\ 0.4 p(w_3) \cdot 0.6 & \text{иначе, если } \neq 0 \\ 0.4 \cdot 1 & \end{cases}$$

последний метод
Kneser - Ney

где $h=2$

$$\tilde{p}(w_2 | u_1) = \frac{c(w, w_2) - \delta}{c(w_1)} + \frac{1}{c(w_2 | w_2)} \frac{|\{ \bar{w} : (\bar{w}, w_2) \in u_1 \}|}{c(w_2 | w_2)}$$

одно
если $c(w, w_2) \neq 0$
то
оценка различных
возможных слов через
 v_2 .
(new York)
где York
базис
слов.

Как оценить качество модели?

1. Матрица качества
р {текста}.

если текст реален, \Rightarrow
его бит потока δ и $\uparrow \uparrow$
max.

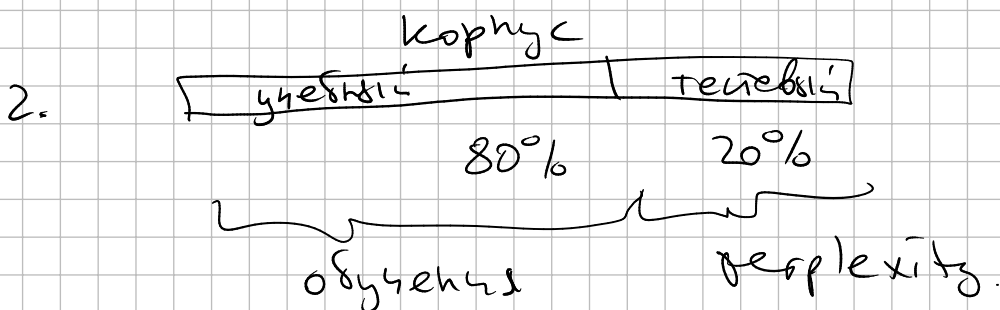
$$p \{ \text{текста} \} = p(w_1, v_2, \dots, w_n) = (1).$$

невозможно посчитать из-за underflow (очень мало)

$$-\log_2 p \{ \text{текста} \} = -\log_2 p(w_1 | \ast \ast \ast) - \log_2 p(w_2 | \ast \dots \ast w_1) -$$

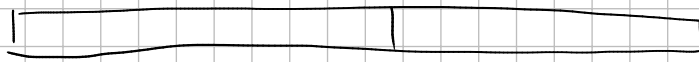
\uparrow будет мало.

это называется perplexity. Должно быть как
маленькое число.



3. По какому критерию выбирать языковую модель.

- Лангес ($\alpha=1, \alpha=0.1$)
 - Кнейсер-Нег δ
 - Simple Backoff
- } принцип выбора



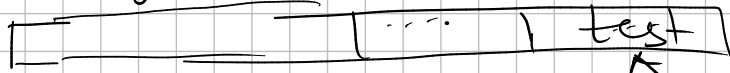
модель 1
модель 2

perplexity модель
1/2

← можно
выбирать
языковую
модель.

Можно

обучение. Настройка



60%

20%

test

оценить perplexity
языковой модели.

↑

можно
выбирать языковую модель

Можно

- язык корпуса
- настройка языковой модели
- выбирать языковую и оценивать её perplexity
- непрерывно обучать.