

Виды токенизации

1) Поделить по пробелам "мама ела ризу"

↓

"мама", "ела", "ризу"

2) BPE

byte-pair Encoding

"абаб", "оооо", "бвzg" ...
корпус

1) ["а", "б", "в", "г", "д"] ← словарь

2) "аа"

3) "аб"

4) "аа" + "а" = "ааа"

Находим пару токенов, которая встречается чаще всего. Добавляем её в словарь.

3) WordPiece

1) ["а", "б", "в", "г", "д"] ← словарь

2) $score = \frac{freq \text{ of pair}}{[freq w_1] \times [freq w_2]}$

["ног", "сказал"]

["ногс", "козал"]

3) добавляем пару с наибольшим score.

4) Unigram

1) словарь — все возможные токены. Например, BPE с макс. словарем.

["баи", "бум", "бах"] — корпус

"a", "d", "u", "y", "x", "da", "au", "dy", "yu", "ax",
 $\text{sum}(2 \ 3 \ 2 \ 1 \ 1 \ 2 \ 1 \ 1 \ 1 \ 1)$
 S таблицы

$$P("da", "u") = \frac{2}{S} \cdot \frac{2}{S}$$

$$P("d", "a", "u") = \frac{3}{S} \cdot \frac{2}{S} \cdot \frac{2}{S}$$

$$P("d", "au") = \frac{3}{S} \cdot \frac{1}{S}$$

2) Для каждого токена

считаем насколько уменьшается вероятность при его удалении.

Удаляем те, у которых вер-ть уменьшается меньше всего.

Языковые модели

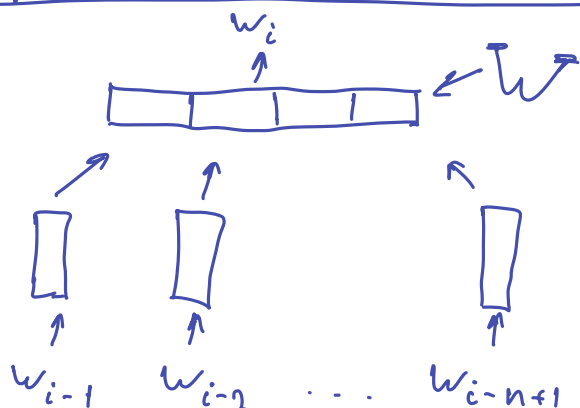
$$P(w_1, \dots, w_T) = p(w_1) \prod_{i=2}^T P(w_i | \underbrace{w_1, \dots, w_{i-1}})$$

N-gram language model

$$p(w_i | w_1, \dots, w_{i-1}) \approx p(w_i | w_{i-n}, \dots, w_{i-1})$$

$$p(w_i | w_{i-1}, \dots, w_{i-n+1}) = \frac{p(w_{i-n+1}, \dots, w_i)}{p(w_{i-n+1}, \dots, w_{i-1})} = \frac{\text{count}(w_{i-n+1}, \dots, w_i)}{\text{count}(w_{i-n+1}, \dots, w_{i-1})}$$

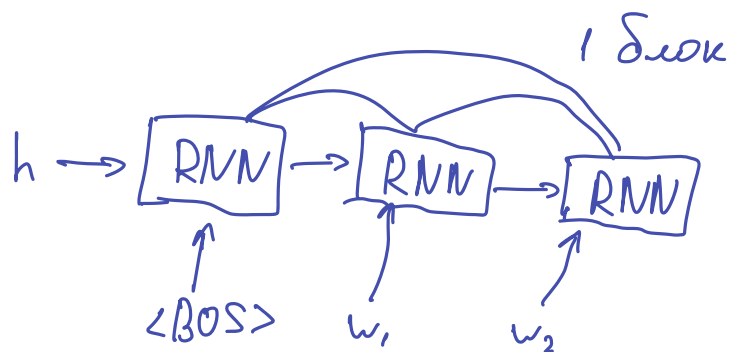
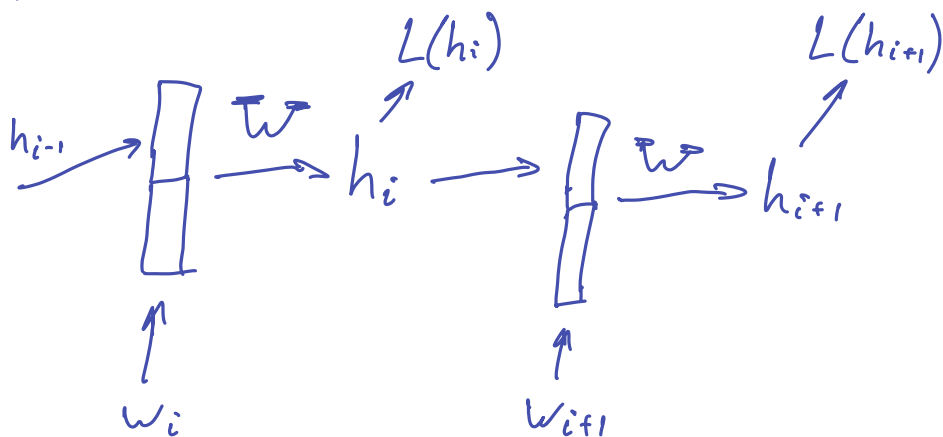
Нейросетевые языковые модели



+ Не надо хранить n-граммы

- увеличенное окно увеличивает W

Рекуррентная языковая модель (RNN)



Teacher forcing

$$\nabla_w L(h_i) = \frac{\partial L(h_i)}{\partial h_i} \cdot \frac{\partial h_i}{\partial h_{i-1}} \cdots \frac{\partial h_i}{\partial w} = \frac{\partial L(h_i)}{\partial h_i} \prod_{j=i}^1 \frac{\partial h_j}{\partial h_{j-1}} \frac{\partial h_i}{\partial w}$$