

language model · 언어모델

- 문장 자체의 출현확률을 예측
- 이전 단어들에 주어졌을 때 다음 단어를 예측하기 위한 모델

In Korean

① 단어의 순서가 중요하지 않다. ex)

② 단어의 생략이 가능

나는 학교에 갑니다. 버스를 타고
나는 버스를 타고 학교에 갑니다.
버스를 타고 나는 학교에 갑니다.
(나는) 버스를 타고 학교에 갑니다.

· 확률이 표시는 현상 발생

· 타고 → ., 학교에, 나는 등 하나 이상이 가능

~~~~~ Korean ~~~~~

$$D = \{x_i\}_{i=1}^N$$

*(문장)*

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^N \log P(x_{1:n} | \theta) \quad \text{where } x_{1:n} = \{x_1, x_2, \dots, x_n\}$$

chain rule

$$\begin{aligned} P(A, B, C, D) &= P(D | A, B, C) P(A, B, C) \\ &= P(D | A, B, C) P(C | A, B) P(A, B) \\ &= P(D | A, B, C) P(C | A, B) P(B | A) P(A) \end{aligned}$$

$$P(x_{1:n}) = P(x_1, \dots, x_n)$$

$$= \underbrace{P(x_n | x_1, \dots, x_{n-1}) \dots P(x_2 | x_1) P(x_1)}$$

$$= \prod_{i=1}^N P(x_i | x_{<i})$$

*( $x_1 \sim x_{n-1}$ 까지의 단어가 주어졌을 때  $x_n$ 의 확률)*

$$\log P(x_{1:n}) = \sum_{i=1}^N \log P(x_i | x_{<i})$$

$$\begin{aligned}
 P(\langle \text{BOS} \rangle, I, \text{love}, \text{to}, \text{play}, \langle \text{EOS} \rangle) &= P(\langle \text{EOS} \rangle | \langle \text{BOS} \rangle, I, \text{love}, \text{to}, \text{play}) P(\langle \text{BOS} \rangle, I, \text{love}, \text{to}, \text{play}) \\
 &= P(\langle \text{EOS} \rangle | \langle \text{BOS} \rangle, I, \text{love}, \text{to}, \text{play}) P(\text{play} | \langle \text{BOS} \rangle, I, \text{love}, \text{to}) P(\langle \text{BOS} \rangle, I, \text{love}, \text{to}) \\
 &= P(\langle \text{EOS} \rangle | \langle \text{BOS} \rangle, I, \text{love}, \text{to}, \text{play}) P(\text{play} | \langle \text{BOS} \rangle, I, \text{love}, \text{to}) P(\text{to} | \langle \text{BOS} \rangle, I, \text{love}) P(\langle \text{BOS} \rangle, I, \text{love}) \\
 &= P(\langle \text{EOS} \rangle | \langle \text{BOS} \rangle, I, \text{love}, \text{to}, \text{play}) P(\text{play} | \langle \text{BOS} \rangle, I, \text{love}, \text{to}) P(\text{to} | \langle \text{BOS} \rangle, I, \text{love}) P(\text{love} | \langle \text{BOS} \rangle, I) P(\langle \text{BOS} \rangle, I) \\
 &= P(\langle \text{EOS} \rangle | \langle \text{BOS} \rangle, I, \text{love}, \text{to}, \text{play}) P(\text{play} | \langle \text{BOS} \rangle, I, \text{love}, \text{to}) P(\text{to} | \langle \text{BOS} \rangle, I, \text{love}) P(\text{love} | \langle \text{BOS} \rangle, I) P(I | \langle \text{BOS} \rangle) P(\langle \text{BOS} \rangle)
 \end{aligned}$$

이제 최대가 되도록 파라미터 조정

# Count Based Approximation

$$\frac{P(\langle \text{EOS} \rangle | \langle \text{BOS} \rangle, I, \text{love}, \text{to}, \text{play})}{\text{이러한 확률을}} \approx \frac{\text{count}(\langle \text{Bos}, I, \text{love}, \text{to}, \text{play}, \langle \text{EOS} \rangle)}{\text{count}(\langle \text{BOS} \rangle, I, \text{love}, \text{to}, \text{play})}$$

등장 확률을 통해 계산

$$P(x_n | x_{<n}) \approx \frac{\text{count}(x_1, \dots, x_n)}{\text{count}(x_1, \dots, x_{n-1})}$$

## 문제점

- ① 해당 단어/문장에 하나도 없으면?  $\rightarrow P(\text{sentence}) = 0$
- 다른 chain rule 에 끼치는 영향이 매우 큼 한 단어만 없어도 문장의 확률 0

이를 개선하기 위해  
Markov Assumption 을 Apply

Markov Assumption  $\rightarrow$  학습 corpus에서 보지 못한 문장에 대한 확률값도 풀  
 . 이전 단어를 통해 현재 단어를 예측할려 할때 모든 이전 단어를 보는 대신  
 지정한 k step 전까지의 단어만을 활용해 현재 단어 예측

$$P(x | x_{<n}) \approx P(x_n | x_{n-1} \dots x_{n-k}) \approx \frac{\text{count}(x_{n-k} \dots, x_n)}{\text{count}(x_{n-k} \dots, x_{n-1})}$$

## n-gram Language Model (Markov Assumption)

| k | n-gram | 명칭       |
|---|--------|----------|
| 0 | 1-gram | uni-gram |
| 1 | 2-gram | bi-gram  |
| 2 | 3-gram | tri-gram |

$$n = k + 1$$

4-gram 부터는 four-gram

- $n$ 이 커질수록 오히려 확률이 정확하게 표현하는데 어려움  
 $n$ 이 ↑ 단어가 없거나 외래어 가능성 ↑  $n$ 이 ↓ 너무 짧게 봐서 외래어 가능성 ↑  
 적절한  $n$  필요

· 보통 3-gram 을 많이 사용

· corpus가 많은 경우 4-gram 사용하기도

## Smoothing, Discounting (여전히 unseen word sequence에 대해 미흡)

- Markov Assumption을 도입 하였지만 unseen word sequence의 확률 = 0이라는 문제가 아직 존재

↓ 그래서 0이 되는 것을 방지하기 위한 기술 등장

Add one smoothing

$$P(w_t | w_{<t}) \approx \frac{C(w_{1:t})}{C(w_{1:t-1})} \approx \frac{C(w_{1:t}) + 1}{C(w_{1:t-1}) + |\text{vocab size}|}$$

0이 되는 것을 방지

$$\hookrightarrow P(w_t | w_{<t}) \approx \frac{C(w_{1:t}) + m \times P(w_t)}{C(w_{1:t-1}) + m} \quad \text{즉가적으로}$$

Kneser-Ney Discounting

핵심 아이디어

다양한 단어가 뒤를 이을수록 unseen word sequence에 등장할 확률이 높지 않을까?

↳ 많이 등장한 단어의 종류가 다양할수록 해당 확률이 높다.

Modified Kneser-Ney Discounting 보통 많이 사용