

4장

4. Language Model

* Language Model

- 문서로부터 query 생성.

- 각 문서 d_i 로 probabilistic language model M_d 생성.

- $p(q_i | M_d)$ 기반으로 document ranking.

- probabilistic \leftrightarrow deterministic (automata)

	w	$p(w q_1)$
the		0.2
a		0.1
:		

표에 의해 확률 계산됨.

- "generative" model.

- MT, spell correction, speech recognition 등에서 활용.

* Formulating LM

- Chain rule

$$p(w_1, w_2, \dots, w_n) = \prod_i p(w_i | w_1, w_2, \dots, w_{i-1})$$

- 단점) 너무 많은 가능한 sentence가 있다. 모두 커버 X.

$$\text{Complexity} = O(V^{n^*}). n^*: \text{maximum sentence length}.$$

- Markov assumption

- 각 단어 발생은 서로 independent 한다.

$$p(w_1, \dots, w_n) = p(w_1)p(w_2) \dots p(w_n) : \text{Unigram Model}$$

- 단단, popular한 방법. 문법, 순서등을 고려 X.

- N-gram language models

- 이전 ($n-1$) 개의 단어에 conditioned.

$$\text{ex) bigram : } p(w_1, \dots, w_n) = p(w_1 | w_0)p(w_2 | w_1) \dots p(w_n | w_{n-1})$$

- 노이즈 높아질수록 spatial, time complexity가 너무 증가함.

- 정확해지긴 하지만, unigram을 보통 씀.

- 확률로 다음 단어 뽑기

- Cdf (누적분포함수)

* Language models for IR

- query 만들 확률이 가장 높은 document 찾기.

- 각 문서는 language model의 basis가 됨.

$$\cdot P(d|q) = \frac{P(q|d) P(d)}{P(q)}$$

• $P(q)$: constant.

• $P(d)$: prior. same for all d or PageRank

• $P(q|d)$ 와 $P(d|q)$ 를 등급으로 측정함.

* MLE와 MAP (best & 찾기)

- MLE (Maximum Likelihood Estimation)

$$\hat{\theta} = \arg\max_{\theta} P(X|\theta) \text{ 찾기}$$

• Small sample size 일 때, 어려움.

- MAP (Maximum A Posterior)

$$\hat{\theta} = \arg\max_{\theta} P(X|\theta) = \arg\max_{\theta} P(\theta|X) P(X)$$

• prior 정의 문제

* Query Likelihood Model

- $P(q|d)$ 은 ranking하기.

- 가정: 각 query 단어는 서로 독립적으로 생성된다.

$$\cdot P(q|d) = P(t_1, \dots, t_{|q|} | d) = \prod_{i=1}^{|q|} P(t_i | d)$$

- $q = w_1, w_2, \dots, w_n$ 이라고 하면

$$\cdot f(q, d) = \log p(q|d) = \sum_{w \in q} \log p(w|d) = \sum_{w \in q} c(w, q) \log p(w|d)$$

$$\left(p(w|d) = \frac{c(w, d)}{|D|} \right)$$

document LM

• 문제) 문법, 순서 고려 X

d에 있는 단어에 대해서 $p(w|d) = 0 \rightarrow f(q, d) = 0$. \therefore smoothing 필요.

- Smoothing : 본 단어에 대한 확률을 낮추고, 다른 단어들에게 확률 할당.

* Smoothing LM

- Reference LM 사용.

$$\cdot p(w|d) = \begin{cases} p_{seen}(w|d) & \text{if } w \text{ is in } d \\ \alpha_d p(w|c) & \text{otherwise} \end{cases}$$

collection based LM

$$\cdot f(q, d) = \log p(q|d) = \sum_{w \in q, c(w, d) > 0} c(w, q) \log p_{seen}(w|d) +$$

$$\underbrace{\sum_{w \in q, c(w, d) = 0} c(w, q) \log \alpha_d p(w|c)}$$

$$= \underbrace{\sum_{w \in q} c(w, d) \log \alpha_d p(w|c)} - \underbrace{\sum_{w \in q, c(w, d) > 0} c(w, d) \log \alpha_d p(w|c)}$$

$$\cdot f(q, d) = \sum_{w \in q, c(w, d) > 0} c(w, q) \log \frac{p_{seen}(w|d)}{\alpha_d p(w|c)} + \sum_{w \in q} c(w, q) \log \alpha_d p(w|c)$$

$$|q| \log \alpha_d \quad \downarrow \quad \sum_{w \in q} c(w, d) \log p(w|c)$$

$$\cdot \log p(q|d) = \sum_{w \in q} c(w, q) \log \frac{p_{seen}(w|d)}{\alpha_d p(w|c)} + |q| \log \alpha_d + \sum_{w \in q} c(w, q) \log p(w|c)$$

DL

IDF

DL normalization.

ignore

* Smoothing

- Additive smoothing

$$p(w|d) = \frac{C(w,d) + \alpha}{|d| + \alpha |V|}$$

(
d 깊이
d에서 w 본 개수
Vocab 크기)

- 대하기: Laplace smoothing.

- 문제: 각 단어들의 중요성이 다른데, 이를 무시함.

- Bayesian Smoothing

$$p(w|d) = \frac{C(w,d) + \alpha p(w|M_c)}{|d| + \alpha} \rightarrow \text{모든 collection} \text{의 term frequency (not document).}$$

- Known as dirichlet prior smoothing.

$$\alpha \in [0, +\infty]$$

$$p(w|d) = \frac{C(w,d) + \alpha p(w|c)}{|d| + \alpha} = \frac{|d|}{|d| + \alpha} \cdot \frac{C(w,d)}{|d|} + \frac{\alpha}{|d| + \alpha} p(w|c)$$

- JM Smoothing

$$\lambda \in [0, 1]$$

$$p(w|d) = (1 - \lambda) p(w|d) + \lambda p(w|c)$$

- λ 가 작으면 query를 모두 포함한 문서 위주.

λ 가 크면 긴 query에 적합. (긴 query는 $p(w|d)$ 가 0이 될 가능성 높음)

- Known as JM Smoothing

- p59 예제 풀기.

* Linear Interpolation

- N-1 gram 확률로 N gram 확률 smoothing 하기.

$$p(w_t | w_{t-1}, \dots, w_{t-n+1}) = \lambda_1 p(w_t | w_{t-1}, \dots, w_{t-n+1}) + \lambda_2 p(w_t | w_{t-1}, \dots, w_{t-n+2}) + \dots + \lambda_n p(w_t | w_t)$$
$$\sum_{i=1}^n \lambda_i = 1$$

* Bayesian smoothing 적용

$$p(w|d) = \frac{|d|}{|d|+\alpha} \cdot \frac{c(w,d)}{|d|} + \underbrace{\frac{\alpha}{|d|+\alpha}}_{\alpha d} p(w|c)$$

$$\alpha d p(w|c) = \frac{\alpha}{|d|+\alpha} p(w|c)$$

$$\Rightarrow \frac{p_{seen}(w|d)}{d p(w|c)} = \frac{\frac{c(w,d)+\alpha p(w|c)}{|d|+\alpha}}{\frac{\alpha p(w|c)}{|d|+\alpha}} = 1 + \frac{c(w,d)}{\alpha \cdot p(w|c)}$$

$$\begin{aligned} - \log p(q|d) &= \sum_{\substack{w \in q \\ w \in d}} c(w,d) \log \left(1 + \frac{c(w,d)}{\alpha p(w|c)} \right) + |q| \log \frac{\alpha}{|d|+\alpha} \\ &= \sum_{\substack{w \in q \\ w \in d}} c(w,d) \log \left(1 + \frac{c(w,d)}{\alpha} \cdot \frac{1}{p(w|c)} \right) + |q| \log \frac{\alpha}{|d|+\alpha} \\ &\quad \text{TF} \quad \text{IDF} \end{aligned}$$

- JM Smoothing은 똑같음. ($\frac{1+k}{k}$ term 추가 : doc length normalization)

→ 기초적인 것은 choose all correct 문제로.

- 계산보다, 이를 문제 많이나옴 (꼼꼼)

- 좋은 것 모두 고르시오.

- 하나만 고르면 복불장수 있음.

- 3개 이상은 거의 없음. (~27%)

- 속수의 범위, 이러한 속수에 대한 수명 고르시오.

- 확률 문제 2개 틀림.

- Sample of 10 balls : 10개 = 1 sample

- 10개 . 각 10점.