

Language model interpolation

- unseen word sequence에 대처하는 방법
- 다른 Language Model을 linear하게 열정비율로 섞는 것

$$\tilde{P}(w_n | w_{n-k} \dots, w_{n-1}) = \lambda \underbrace{P_1(w_n | w_{n-k} \dots, w_{n-1})}_{\text{Language model}} + (1-\lambda) \underbrace{P_2(w_n | w_{n-k} \dots, w_{n-1})}_{\text{Language model}}$$

★ 그나마 모든 corpus를 합쳐서 만들지 않는 이유

Domain specific corpus의 양이 적어서 단어의 분포를 잘 나타낼 수 있는 가능성 ↑

적절한 interpolation ratio(λ)를 찾는 것 중요

Back off

n-gram의 n을 줄여 나가며 예측을 구하는 것 ex) Trigram \rightarrow Bi-gram \rightarrow unigram

$$\tilde{P}(w_n | w_{n-k} \dots, w_{n-1}) = \lambda_1 P(w_n | w_{n-k} \dots, w_{n-1}) \\ + \lambda_2 P(w_n | w_{n-k+1} \dots, w_{n-1}) \\ + \dots \\ + \lambda_k P(w_n),$$

$$\text{where } \sum_i \lambda_i = 1.$$

Summary

전통적인 NLP는 단어를 discrete symbol로 보기 때문에 문제가 발생

- Exact matching에 대해서만 count를 하며, 예측은 approximation

ex) 분홍과 빨강이 얼마나 유사한지 신경X 시크 아예 다른 것

Perplexity (PPL)

- Test set에 대해서 언어모델을 이용하여 확률(Likelihood)을 구하고 PPL 수식을 넣어 성능 측정

$$PPL(x_1, \dots, x_n; \theta) = P(x_1, \dots, x_n; \theta)^{-\frac{1}{n}} = \sqrt[n]{\frac{1}{P(x_1, \dots, x_n; \theta)}}$$

PPL이 작을수록 좋은 언어모델

Perplexity

- Timestep 별 평균 branch의 수
- 즉사위 PPL: 매 timestep마다 가능한 가짓수인 G

PPL이 ↓ 확률 분포가 sharp 하다.

PPL이 ↑ 확률 분포가 Flat 하다.

언어모델의 경우

- 낮을수록 Good
- 확률의 역수에 로그값이로 기하 평균
- 매 time step마다 평균적으로 헛참리크(no clue 완전히) 있는 단어의 수

N -gram

- 장점:
- 쉽게 대형 시스템에 적용 가능
 - N -gram 훈련 및 추론 방식이 굉장히 간편

단점

- unseen word sequence에 대해 미흡
- 멀리 있는 단어에 대해 대처 불가
ex) Markov Assumption
- N 이 커질수록 용량도 커짐

Auto Regressive Task

· NLG

- 이전 time step의 출력에 따라 현재時刻의 state가 바뀌는 것 (이전의 출력이 앞으로의 것에 영향을 끼치는 것)

So

- 학습시에는 이전 점답의 출력이 아닌 정답을 RNN에 넣어줌

Problem

1. inference를 위한 code 필요
2. 학습과 추론 방식의 차이가 발생하여 성능 저하
Transformer의 등장 이유?