

Integrating knowledge in Language Model

Language Model $\begin{cases} \text{Masked LM} & \text{ex. } \text{went} \text{ I [MASK] to the } \text{store} \text{ [MASK].} \\ \text{Standard LM} & \text{The students opened their } \text{books}. \end{cases}$

- predictions generally make sense, but are not all factually correct.

Due to:

- 1) Unseen facts
- 2) Rare facts
- 3) Model sensitivity

개선 방향

1) Querying traditional knowledge bases

2) Querying LMs as knowledge bases

$\Rightarrow T5$

갖춰진 dataset

· manual annotation 필요

· Complex NLP pipeline 필요

· More flexible 하지만

hard to interpret, trust, modify

Knowledge를 LM에 주입하는 방법

- 1) Add pretrained entity embeddings
- 2) Use an external memory
- 3) Modify the training data

ERNIE - entity embedding

- 사전 학습된 BERT를 확장

ex. Bob Dylan wrote Blowin' in the Wind in 1962.

에 대한 wiki data 사전 학습에 기존 것과 합성

T-Encoder \Rightarrow K-Encoder

KGLM - KG 형태에 의존적

KNN-LM

- 의미가 유사한 text를 KNN로 찾아내고 단어 예측하는 방법들로 LM 학습
- KNN과 LM prob 을 함께 결합해 사용 가능

$$P(y|x) = \lambda P_{KNV}(y|x) + (1-\lambda)P_{LM}(y|x)$$

- pretraining 과정에서 자연스럽게 knowledge 주입

WKLM - 모델이 true knowledge와 false knowledge 구분하기
- Weakly supervised Knowledge Pretrained LM

BERT 이후 masking strategy 변경한 다양한 모델 제안
그 중에서도 span 단위로 masking 수행하는 모델 등장

Salient span masking

- Named entity를 이용해 salient span 생성
- 간단하지만 entity 정보 효과적이 학습
- QA task에서 두각

REALM : Retrieval-Augmented LM-Training

- masking 방법을 학습해 QA task 성능 개선

LAMA Probe

- 동일 setting에서 학습해 어떤 LM이 가장 많은 정보 포함하는지 비교