

□

GPT(Generative Pre-trained Transformer)

1. GPT vs. BERT

GPT

어제 카페 갔었어 거기 사람 많더라

BERT

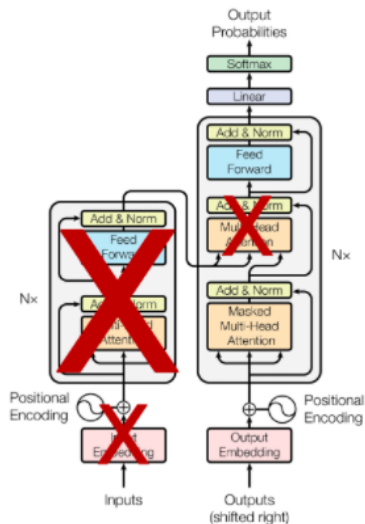
어제 카페 갔었어 □ 사람 많더라

	GPT	BERT
모델	Pre-trained Language Model	Pre-trained Masked Language Model
구조	트랜스포머에서 디코더만 이용	트랜스포머에서 인코더만 이용
pre-training 방법	이전 단어들이 주어졌을 때 다음 단어를 예측하도록 pre-training	문장 중간에 빈칸을 만들고 해당 빈칸에 어떤 단어가 적절한지 예측하도록 pre-training
방향성	단방향(unidirectional)	양방향(bidirectional)

2. GPT의 구조

- 트랜스포머에서 인코더를 제외하고 디코더만 사용
- 인코더 쪽에서 보내오는 정보를 받는 모듈(Multi-Head Attention)도 제거

GPT는 엄밀히 말하면 트랜스포머의 디코더 구조를 이용한다기 보다는 **Masked Multi-Head Attention**을 이용한다고 볼 수 있다.



2. GPT의 Embedding

1) 워드 임베딩

- GPT 역시 BERT와 마찬가지로 워드 임베딩을 위해 **서브워드 토큰나이저**를 사용한다.
- GPT의 경우 길이가 50257인 단어 집합을 사용한다.

2) Position Embedding

- BERT와 마찬가지로 각 단어의 위치 정보를 포함해주기 위한 Position Embedding을 사용한다.
- 워드 임베딩과 더해서 사용하기 위해 임베딩 크기는 워드 임베딩과 같게 맞춰준다.

3. GPT의 Pre-training

1) Input & Label

- GPT의 Pre-training은 다음에 올 토큰을 예측하는 방식으로 학습되어 별다른 라벨이 필요 없는 Unsupervised Learning이다.
- 따라서 input은 하나의 sequence가 되고 label은 각 단어의 다음에 올 단어들의 sequence가 된다.

Fetch window_size + 1 tokens

Dataset :	23	38	1251	123 56 7 33 9 12	12623	78	5678
Input Ids :	123	56	7	33	9	
Label Ids :		56	7	33	9	12

<https://ainote.tistory.com/17>

2) Future Mask (Attention Mask)

-
- | | seq_length | | | | | |
|------------|------------|---|---|-----|-----|---|
| seq_length | F | T | T | ... | T | T |
| | F | F | T | T | ... | T |
| | F | F | F | T | ... | T |
| | ⋮ | ⋮ | ⋮ | ⋮ | | ⋮ |
| | ⋮ | ⋮ | ⋮ | ⋮ | | ⋮ |
| | F | F | F | F | ... | F |
| | F | F | F | F | ... | F |
- Future Mask

- ### 3) Forwarding 과정

- Query Key
-
- 0101 $\xrightarrow{1,0}$ 0101 \rightarrow 카피
- 카피 카피
- $\frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2}$ $\frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2}$
- 7H1 7H1 masking
- 키값 키값
- 암호문 암호문

