

# BERT: Bidirectional Encoder Representations from Transformer

## Pre-training of Deep Bidirectional Transformers for Language Understanding

- transformer를 활용한 language representation

- transfer-learning 방식으로 학습됨. 즉, 대량의 unlabeled data로 pre-training 후, 특정 task를 가지고 있는 labeled data로 fine-tuning.

- Comparing unidirectional 하거나, ELMo처럼 shallow bidirectional 한 모델과 달리, BERT는 deeply bidirectional 한 모델.

\* Unidirectional: 텍스트를 한 방향으로만 처리. (거의와 같은 모델은 텍스트를 왼쪽에서 오른쪽으로 읽음. 이러한 문법은 처리된 단어의 오른쪽 context(이후의 단어들)를 고려하지 않음이라, 극단 '방향성'인 task (ex. 텍스트생성)에 적합.

\* shallow bidirectional: 두 개의 방향의 레이어를 사용하여 텍스트를 양방향으로 처리. 하나의 레이어는 왼쪽에서 오른쪽으로, 다른 하나는 오른쪽에서 왼쪽으로 텍스트를 읽음. 그러나 이 두 레이어는 서로 독립적으로 작동하며, 각 레이어는 텍스트의 한 방향의 context만을 고려. ELMo는 이 두 레이어의 출력을 결합하여 더 풍부한 단어 표현을 생성.

\* deeply bidirectional: 모든 레이어에서 양방향의 context를 동시에 고려. 이는 각 단어가 주변 모든 단어(이전과 이후 모두)의 context를 바탕으로 이해될 수 있게 됨. BERT는 이러한 깊은 양방향 특성으로 인해 텍스트의 의미를 더 정확하게 파악하고, 다양한 NLP task에서 뛰어난 성능을 보임.

- BERT에는 pre-training 과정과 그 이후의 fine-tuning 과정이 있음.

\* pre-training에서 unlabeled 된 corpus를 transformer의 encoder architecture과 MLM과 NSP 과정을 통해서 feature가 나옴.

\* 이 feature를 task specific 용이 적용하여 fine-tuning을 통해 모델링을 하는 것.

- BERT는 모델의 크기에 따라 base 모델과 large 모델을 제공.

\* BERT-base:  $L=12$ ,  $H=768$ ,  $A=12$ , Total parameters = 110M

\* BERT-large:  $L=24$ ,  $H=1024$ ,  $A=16$ , Total parameters = 340M

( $L$ : transformer block의 layer 수,  $H$ : hidden size,  $A$ : self-attention heads 수)

### Input/Output Representations

- BERT의 input은 단어 embedding 값의 합으로 이루어져 있음.

- 모든 sentence의 첫번째 token은 [CLS] (special classification token)임.

\* 특정 context의 의미를 대표하는 역할.

\* classification을 위한 embedding.

\* classification task일 경우 cls token 위치에서 가장 높은 hidden unit output을 feature로 생각해서 classification 진행. 아니면 역시.

- 특정 구분

\* [SEP] token 사용.

\* segment embedding 사용: 앞의 문장을 sentence A embedding, 뒤의 문장을 sentence B embedding을 구분함. (모든 고정된 값)

### Pre-training

#### ① Masked Language Model (MLM)

- 무작위하게 몇 개의 토큰을 mask 시킴.

- transformer 구조에 들어가 주변 단어의 context만을 보고 mask 된 단어를 예측.

## ② Next Sentence prediction

- 두 문장을 pre-training 시에 같이 (공통) 두 문장이 이어지는 문장인지 아닌지 맞추는 것.

## fine tuning

- pre-training을 가진 모델을 가지고 downstream 한 방향으로 적용하여 finetuning 시함.

\* 이걸 가지고 세팅하고 pre-training 같은 hyper parameter 다 재설정 필요함.

(이제) batch size: 16, 32 / learning rate (Adam):  $5e^{-5}$ ,  $3e^{-5}$ ,  $2e^{-5}$  / # of epochs: 2, 3, 4