

Transformer

20220307 수학과 고등기

No.

Attention \rightarrow linear combination \rightarrow skip gram

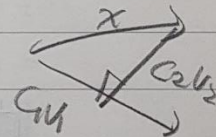
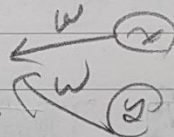
$$x = (x \cdot v_1) v_1 + (x \cdot v_2) v_2 + \dots + (x \cdot v_n) v_n$$

$$\underbrace{(w_a x)}_{a_1} \underbrace{(w_k v_1)}_{\tilde{v}_1} + \dots + \underbrace{(w_a x)}_{a_n} \underbrace{(w_k v_n)}_{\tilde{v}_n}$$

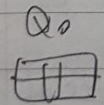
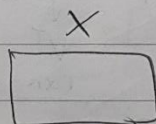
$$\tilde{x} = a_1 \tilde{v}_1 + a_2 \tilde{v}_2 + \dots + a_n \tilde{v}_n$$

$$\frac{e^{a_1}}{\sum_{i=1}^n e^{a_i}} \tilde{v}_1 + \frac{e^{a_2}}{\sum_{i=1}^n e^{a_i}} \tilde{v}_2 + \dots + \frac{e^{a_n}}{\sum_{i=1}^n e^{a_i}} \tilde{v}_n$$

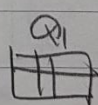
$$\frac{(w_a x) \cdot (w_k v_1)}{\sum_{i=1}^n (w_a x) \cdot (w_k v_i)}$$



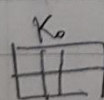
Multi-head attention



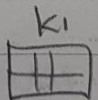
$w_0 Q$



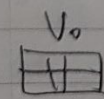
$w_1 Q$ (new)



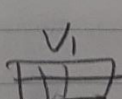
$w_0 K$



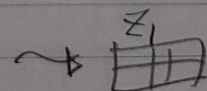
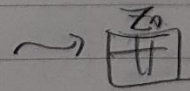
$w_1 K$ (new)



$w_0 V$



$w_1 V$ (new)



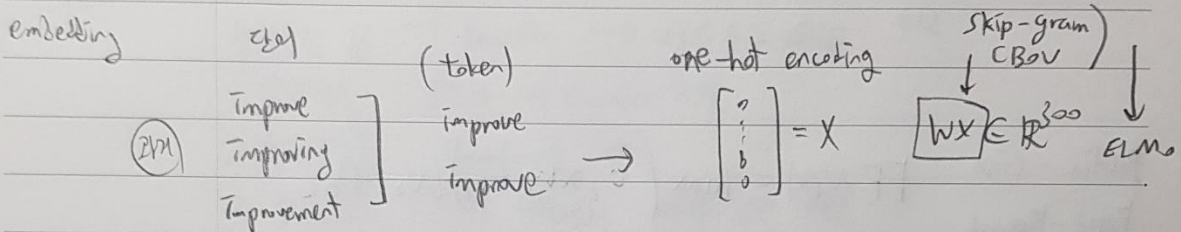
문장 \rightarrow subword Tokenizer \rightarrow 토큰 ID 시퀀스

Embedding layer \rightarrow 각 토큰의 벡터

+

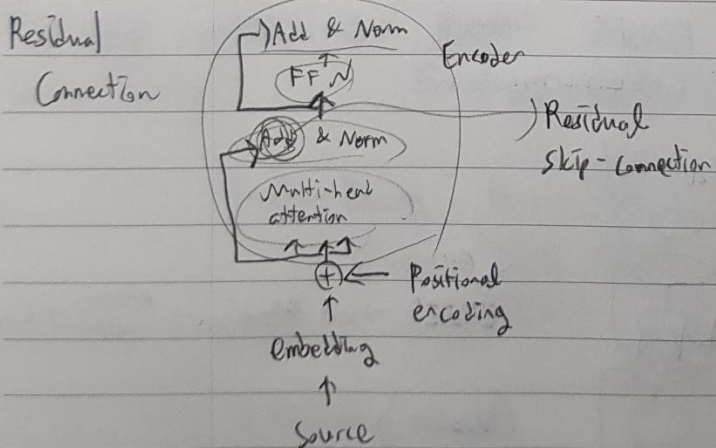
Positional Embedding \rightarrow 위치 정보를 더한 벡터 $\sin\left(\frac{pos}{10000^{2i/d}}\right)$

Transformer 입력



(B2) subword \rightarrow word piece, BPE

unaffordable = un + afford + able \langle subword \rangle
 $= [78, 178, 93] \leftarrow$ 정수 index



embedding-dim = num-heads \times head-dims

Masked
Attention

$$\text{Masked Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}} + \text{mask}\right) \cdot V$$

필요한 이유

! transformer는 입력 전체를 동시에 처리, 미래 단어를 봐야 함

! Mask로 자기자신보다 앞의 단어만 보게 제한

Training

입력문장 $X \rightarrow \text{Encoder} \rightarrow \text{context vector}$ 출력문장 $Y \rightarrow \text{Decoder} \rightarrow \text{다음 토큰 예측}$ 1. 입력문장을 토큰화 $\rightarrow \text{Encoder}$ 전달

2. 출력문장을 Decoder에 입력

3. Decoder는 각 위치에서 다음 토큰 예측

4. Cross Entropy Loss 계산

$$L = -\sum \log p(y_{t+1} | x, y_{<t})$$

5. Optimizer 통해 파라미터 업데이트

병렬화 학습,

병렬화 학습: RNN과 다름

Teacher Forcing

Encoder: 전체 시퀀스를 한번에 attention으로 처리

Teacher Forcing: Decoder 학습시

이전에 생성된 예측값 아닌

정답토큰을 다음 입력으로 사용

Transformer 예시

① Embedding Layer

학습하는 parameter

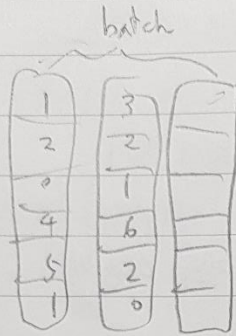
② Attention Weight

③ Position-wise FFNN

④ Layer Norm 파라미터

⑤ 최종 출력 Projection

Layer
Normalization



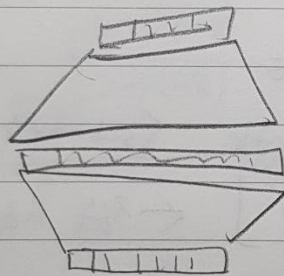
정규화 대상: 각 샘플의 특성 (dim)

mean 2 3 3

std 2 2 2

Position-wise
Feed-Forward
Network

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$



입력: 하나의 토큰 벡터

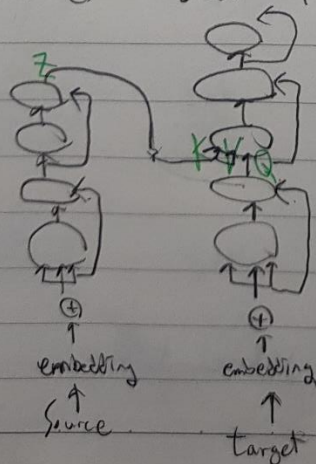
출력: 같은 차원의 벡터

특성: Position-wise

MoE Position-wise FFNN을 여러개의 전문가 네트워크로 확장한 구조
선택된 게이트 네트워크가 결정

$$MoE(x) = \sum g_i(x) \cdot E_i(x)$$

K, V의 역할
in Encoder → decoder



<Encoder>

입력 토큰들 → self Attention → encoder output → K, V

<Decoder>

출력 토큰들 → Masked self Attention → Q

Cross Attention (Q, K, V)