

2022학년도 1학기 컴퓨터언어학

22강 셀프 어텐션 계층과 트랜스포머 소개

박수지

서울대학교 인문대학 언어학과

2022년 5월 25일 수요일

오늘의 목표

- 1 Attention 점수 계산을 일반화하여 query, key, value에 관한 식으로 표현할 수 있다.
- 2 셀프 어텐션 계층(self-attention layer)에서 입력 시퀀스를 처리하는 과정을 설명할 수 있다.
- 3 잔차 연결(residual connection)이 무엇인지 설명할 수 있다.
- 4 레이어 정규화(layer normalization)의 목적과 방법을 설명할 수 있다.

	he	hit	me	with	a	pie
il	α_{11}	α_{21}	α_{31}	α_{41}	α_{51}	α_{61}
a	α_{12}	α_{22}	α_{32}	α_{42}	α_{52}	α_{62}
m'	α_{13}	α_{23}	α_{33}	α_{43}	α_{53}	α_{63}
entarté	α_{14}	α_{24}	α_{34}	α_{44}	α_{54}	α_{64}

복습

부호화기-복호화기 모형에서의 attention

도착어 문장의 현재 단계에서 출발어 문장의 각 토큰을 얼마나 주목할 것인가?

복호화기의 직전 단계 은닉 상태 $\vec{h}_{i-1}^d \in \mathbb{R}^{d_h}$

부호화기의 모든 단계 은닉 상태 $\vec{h}_j^e \in \mathbb{R}^{d_h} (j = 1, 2, \dots, n)$

Attention 점수 $a_{ij} = \text{score}(\vec{h}_{i-1}^d, \vec{h}_j^e) = \vec{h}_{i-1}^d \cdot \vec{h}_j^e \in \mathbb{R}$

Attention 가중치 $\alpha_{ij} = \text{softmax}(a_{ij}) \in]0, 1[(j = 1, 2, \dots, n)$

문맥 벡터 $\vec{c}_i = \sum_{j=1}^n \alpha_{ij} \vec{h}_j^e \in \mathbb{R}^{d_h}$

Causal (or backward looking) self-attention layer

$$(\vec{x}_1, \dots, \vec{x}_n) \mapsto (\vec{y}_1, \dots, \vec{y}_n)$$

- 1 입력 시퀀스의 각 토큰을 처리할 때 현재 토큰 및 이전의 모든 토큰에 접근한다.
(이후의 토큰은 고려하지 않는다.)
⇒ 자기회귀적인(autoregressive) 방법으로 언어 모델을 만들 수 있다.
- 2 각 토큰에서 수행하는 계산은 다른 토큰의 계산과 독립적으로 이루어진다. (토큰을 순차적으로 처리하지 않는다.)
⇒ 순전파와 훈련을 병렬적으로 진행할 수 있다.

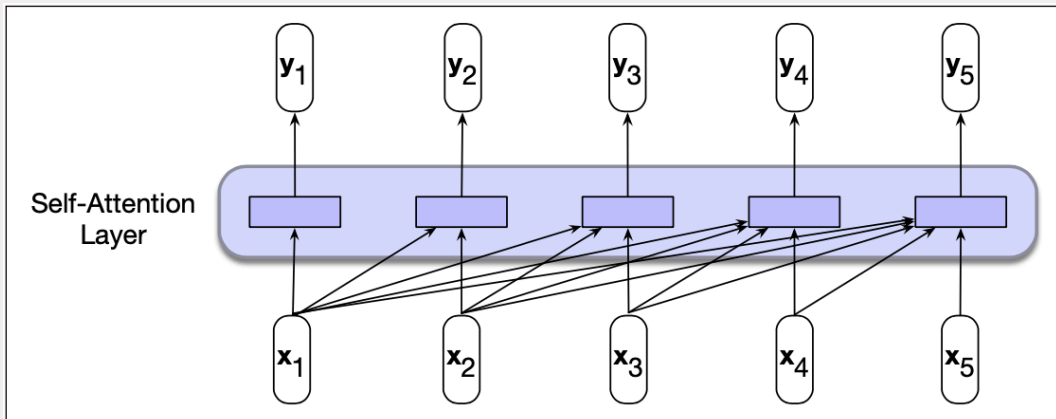


Figure 9.15 Information flow in a causal (or masked) self-attention model. In processing each element of the sequence, the model attends to all the inputs up to, and including, the current one. Unlike RNNs, the computations at each time step are independent of all the other steps and therefore can be performed in parallel.

“Causal” self-attention

입력 $\vec{x}_i \in \mathbb{R}^d$

선행 요소 $\vec{x}_j \in \mathbb{R}^d \quad \forall j \leq i$

Attention scores $a_{ij} = \text{score}(\vec{x}_i, \vec{x}_j) = \vec{x}_i \cdot \vec{x}_j \in \mathbb{R}$

Attention weights $\alpha_{ij} = \text{softmax}(\text{score}(\vec{x}_i, \vec{x}_j)) \in]0, 1[\quad \forall j \leq i$

출력 $\vec{y}_i = \sum_{j=1}^i \alpha_{ij} \vec{x}_j \in \mathbb{R}^d$

예시

입력 $\vec{x}_3 \in \mathbb{R}^d$

선행 요소 $\vec{x}_1, \vec{x}_2, \vec{x}_3 \in \mathbb{R}^d$

Attention scores $a_{3j} = \text{score}(\vec{x}_3, \vec{x}_j) = \vec{x}_3 \cdot \vec{x}_j \in \mathbb{R} \quad (j = 1, 2, 3)$

Attention weights $\alpha_{3j} = \text{softmax}(a_{3j}) = \frac{\exp(a_{3j})}{\exp(a_{31}) + \exp(a_{32}) + \exp(a_{33})}$

출력 $\vec{y}_3 = \alpha_{31}\vec{x}_1 + \alpha_{32}\vec{x}_2 + \alpha_{33}\vec{x}_3 \in \mathbb{R}^d$

셀프-어텐션 계층에서 일어나는 일

입력 \vec{x}_i 에 대응하는 출력 \vec{y}_i 를 얻는 방법

$$\vec{y}_i = \sum_{j=1}^i \alpha_{ij} \vec{x}_j = \sum_{j=1}^i \text{softmax}(\text{score}(\vec{x}_i, \vec{x}_j)) \vec{x}_j = \sum_{j=1}^i \text{softmax}(\vec{x}_i \cdot \vec{x}_j) \vec{x}_j$$

관찰

우변에서 입력 토큰들의 임베딩 벡터 \vec{x}_i 가 세 가지 다른 역할로 사용된다.

$$\text{output}_i = \sum_{j=1}^i \text{softmax}(\text{query}_i \cdot \text{key}_j) \text{value}_j$$

Attention score 수정 (1)

query, key, value 구별

Transformer의 접근 방법: 입력 벡터의 세 가지 표상

$$\vec{y}_i = \sum_{j=1}^i \text{softmax}(\vec{q}_i \cdot \vec{k}_j) \vec{v}_j$$

세 가지 역할을 각기 다른 벡터로 나타내기 위해 가중치 매개변수 행렬 \mathbf{W}^Q , \mathbf{W}^K , \mathbf{W}^V 를 도입한다.

query $\vec{q}_i = \vec{x}_i^T \mathbf{W}^Q \in \mathbb{R}^{d_k}$...“current focus of attention”

key $\vec{k}_j = \vec{x}_j^T \mathbf{W}^K \in \mathbb{R}^{d_k}$...“preceding input(s)”

value $\vec{v}_j = \vec{x}_j^T \mathbf{W}^V \in \mathbb{R}^{d_v}$

예시

입력 $\vec{x}_3 \in \mathbb{R}^d$ (열벡터)

Query $\vec{q}_3 = \vec{x}_3^\top \mathbf{W}^Q \in \mathbb{R}^{d_k}$ (행벡터)

선행 요소 $\vec{x}_1, \vec{x}_2, \vec{x}_3 \in \mathbb{R}^d$ (열벡터)

Keys $\vec{k}_1 = \vec{x}_1^\top \mathbf{W}^K, \vec{k}_2 = \vec{x}_2^\top \mathbf{W}^K, \vec{k}_3 = \vec{x}_3^\top \mathbf{W}^K \in \mathbb{R}^{d_k}$ (행벡터)

Attention scores $a_{3j} = \text{score}(\vec{x}_3, \vec{x}_j) = \vec{q}_3 \cdot \vec{k}_j \in \mathbb{R} \quad (j = 1, 2, 3)$

Attention weights $\alpha_{3j} = \text{softmax}(a_{3j}) = \frac{\exp(a_{3j})}{\exp(a_{31}) + \exp(a_{32}) + \exp(a_{33})}$

Values $\vec{v}_1 = \vec{x}_1^\top \mathbf{W}^V, \vec{v}_2 = \vec{x}_2^\top \mathbf{W}^V, \vec{v}_3 = \vec{x}_3^\top \mathbf{W}^V \in \mathbb{R}^{d_v}$ (행벡터)

출력 $\vec{y}_3 = \alpha_{31}\vec{v}_1 + \alpha_{32}\vec{v}_2 + \alpha_{33}\vec{v}_3 \in \mathbb{R}^{d_v}$

참고 원 논문(Vaswani et al., 2017에서는 $d_k = d_v = 1024$ 로 설정했음.

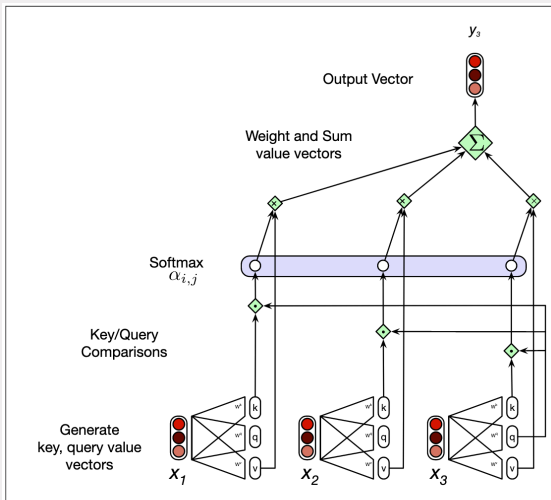


Figure 9.16 Calculating the value of y_3 , the third element of a sequence using causal (left-to-right) self-attention.

관찰

$i = 3$ 일 때(3번째 토큰을 살펴볼 때)
사용하는 벡터

Query \vec{q}_3

Keys $\vec{k}_1, \vec{k}_2, \vec{k}_3$

Values $\vec{v}_1, \vec{v}_2, \vec{v}_3$

주의

내적을 계산하기 위해서는 \vec{q}_i 와 \vec{k}_j 의
차원이 같아야 한다. (\vec{v}_j 의 차원은
달라도 됨)

내적 계산의 문제

$$\vec{q}_i \cdot \vec{k}_j = q_{i1}k_{j1} + q_{i2}k_{j2} + \cdots + q_{id_k}k_{jd_k} \in \mathbb{R}$$

- 1024차원 벡터의 내적은 1024개의 곱을 더하는 것이다.
 \Rightarrow 내적값이 커지면 지수함수(exp)를 취했을 때 overflow가 일어날 수 있다.

Overflow

```
>>> exp(709)
8.218407461554972e+307
>>> exp(710)
<stdin>:1: RuntimeWarning: overflow encountered in exp
inf
```

Attention score 수정 (2)

scaled dot-product

해결: scaled dot-product attention

내적값을 차원 d_k 의 제곱근으로 나누어 준다.

$$\text{score}(\vec{x}_i, \vec{x}_j) = \frac{\vec{q}_i \cdot \vec{k}_j}{\sqrt{d_k}}$$

남은 단계

입력 시퀀스 전체 $(\vec{x}_1, \dots, \vec{x}_N)$ 에 한꺼번에 적용하기

데이터

■ N개의 토큰 시퀀스를...

...N개의 d차원 임베딩 벡터(행벡터)를 쌓은 행렬 \mathbf{X} 로 표기하자.

$$\mathbf{X} = [\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N]^T = \begin{bmatrix} \vec{x}_1^T \\ \vec{x}_2^T \\ \vdots \\ \vec{x}_N^T \end{bmatrix} \in \mathbb{R}^{N \times d}$$

가중치 매개변수

■ Queries, keys, values

$$\mathbf{W}^Q \in \mathbb{R}^{d \times d_k} \Rightarrow \mathbf{Q} = \mathbf{XW}^Q \in \mathbb{R}^{N \times d_k}$$

$$\mathbf{W}^K \in \mathbb{R}^{d \times d_k} \Rightarrow \mathbf{K} = \mathbf{XW}^K \in \mathbb{R}^{N \times d_k}$$

$$\mathbf{W}^V \in \mathbb{R}^{d \times d_v} \Rightarrow \mathbf{V} = \mathbf{XW}^V \in \mathbb{R}^{N \times d_v}$$

사실

시퀀스 길이가 N 이면 $\text{score}(\vec{x}_i, \vec{x}_j)$ 의 모든 가능한 조합은 $(N \times N)$ 행렬로 표현 가능하다.

$$\mathbf{Q} = \begin{bmatrix} \vec{q}_1 \\ \vec{q}_2 \\ \vdots \\ \vec{q}_N \end{bmatrix}, \mathbf{K} = \begin{bmatrix} \vec{k}_1 \\ \vec{k}_2 \\ \vdots \\ \vec{k}_N \end{bmatrix} \in \mathbb{R}^{N \times d_k} \text{이면 } \mathbf{K}^T = [\vec{k}_1^T \ \vec{k}_2^T \ \dots \ \vec{k}_N^T] \in \mathbb{R}^{d_k \times N} \text{ 이므로}$$

$$\mathbf{Q}\mathbf{K}^T = \begin{bmatrix} \vec{q}_1 \vec{k}_1^T & \vec{q}_1 \vec{k}_2^T & \dots & \vec{q}_1 \vec{k}_N^T \\ \vec{q}_2 \vec{k}_1^T & \vec{q}_2 \vec{k}_2^T & \dots & \vec{q}_2 \vec{k}_N^T \\ \vdots & \vdots & \ddots & \vdots \\ \vec{q}_N \vec{k}_1^T & \vec{q}_N \vec{k}_2^T & \dots & \vec{q}_N \vec{k}_N^T \end{bmatrix} = \begin{bmatrix} \vec{q}_1 \cdot \vec{k}_1 & \vec{q}_1 \cdot \vec{k}_2 & \dots & \vec{q}_1 \cdot \vec{k}_N \\ \vec{q}_2 \cdot \vec{k}_1 & \vec{q}_2 \cdot \vec{k}_2 & \dots & \vec{q}_2 \cdot \vec{k}_N \\ \vdots & \vdots & \ddots & \vdots \\ \vec{q}_N \cdot \vec{k}_1 & \vec{q}_N \cdot \vec{k}_2 & \dots & \vec{q}_N \cdot \vec{k}_N \end{bmatrix} \in \mathbb{R}^{N \times N}$$

Attention weights $[\alpha_{ij}] = \left[\text{softmax}_j \left(\frac{\vec{q}_i \cdot \vec{k}_j}{\sqrt{d_k}} \right) \right] = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \in \mathbb{R}^{N \times N}$

최종

$$\text{SelfAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V} \in \mathbb{R}^{N \times d_v}$$

문제

현재 보고 있는 토큰 이후의 토큰은 고려하지 않는다는 조건을 반영해야 한다.

해결

$i < j$ 인 경우 $\vec{q}_i \cdot \vec{k}_j = -\infty$ 로 설정한다.

N

$q_1 \cdot k_1$	$-\infty$	$-\infty$	$-\infty$	$-\infty$
$q_2 \cdot k_1$	$q_2 \cdot k_2$	$-\infty$	$-\infty$	$-\infty$
$q_3 \cdot k_1$	$q_3 \cdot k_2$	$q_3 \cdot k_3$	$-\infty$	$-\infty$
$q_4 \cdot k_1$	$q_4 \cdot k_2$	$q_4 \cdot k_3$	$q_4 \cdot k_4$	$-\infty$
$q_5 \cdot k_1$	$q_5 \cdot k_2$	$q_5 \cdot k_3$	$q_5 \cdot k_4$	$q_5 \cdot k_5$

N

Figure 9.17 The $N \times N$ $Q\mathbf{T}^T$ matrix showing the $q_i \cdot k_j$ values, with the upper-triangle portion of the comparisons matrix zeroed out (set to $-\infty$, which the softmax will turn to zero).

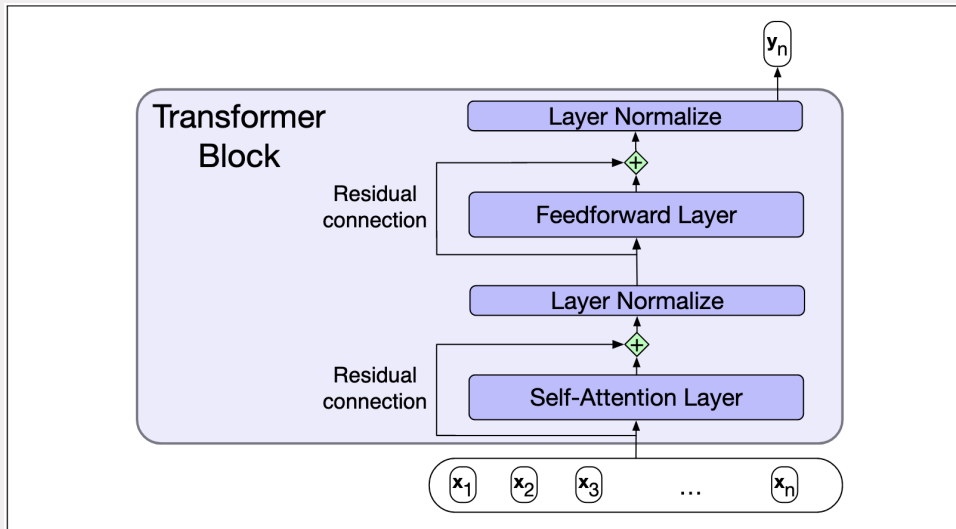


Figure 9.18 A transformer block showing all the layers.

$$\vec{z} = \text{LayerNorm}(\vec{x} + \text{SelfAttention}(\vec{x}))$$

$$\vec{y} = \text{LayerNorm}(\vec{z} + \text{FFNN}(\vec{z}))$$

트랜스포머 블록의 구성요소

- 셀프-어텐션 계층 ...방금 배웠음
- 순방향 계층 ...전에 배웠음
- 잔차 연결(Residual connection) ...아직 모름
- 계층 정규화(Layer normalization) ...아직 모름

$$\vec{z} = \text{LayerNorm} (\vec{x} + \text{SelfAttention} (\vec{x}))$$

$$\vec{y} = \text{LayerNorm} (\vec{z} + \text{FFNN} (\vec{z}))$$

잔차 연결

정보를 낮은 계층에서 중간 계층을 건너뛰고 높은 계층으로 직접 보내는 연결

$$\vec{z} = \text{LayerNorm} (\vec{x} + \text{SelfAttention} (\vec{x}))$$
$$\vec{y} = \text{LayerNorm} (\vec{z} + \text{FFNN} (\vec{z}))$$

계층 정규화

한 계층에 속하는 값들을 정규화(평균이 0, 표준편차가 1이 되도록 만드는 변환)하는 것.

$$\hat{\vec{x}} = \frac{\vec{x} - \mu}{\sigma}$$

$$\text{LayerNorm} (\vec{x}) = \gamma \hat{\vec{x}} + \beta$$

μ \vec{x} 의 평균

σ \vec{x} 의 표준편차

γ, β 매개변수

오늘 배운 것

- 셀프 어텐션 계층: query, key, value, scaled dot product
- 트랜스포머 블록: 잔차 연결, 레이어 정규화

남은 것

- Multihead attention, Positional embeddings (SLP3 9.7 나머지)
- Contextual embeddings: BERT and GPT-2,3 (SLP3 11장)