



DSBA Transformer survey paper study

A Survey of Transformers

#3: Attention

arXiv preprint



고려대학교 산업경영공학과

Data Science & Business Analytics Lab
이유경, 김명섭, 윤훈상, 김지나, 허재혁, 김수빈

발표자 : 김지나

1. Linearized Attention

1. Feature map
2. Aggregation Rule

2. Prototype and Memory Compression

1. Attention with Prototype Queries
2. Attention with Compressed Key-Value Memory

3. Low-rank Self-Attention

1. Low-rank Parameterization
2. Low-rank Approximation

Linearization을 통해 attention의 computational complexity $O(T^2) \rightarrow O(T)$ 줄임

- 기존 Attention: $Q, K, V \in R^{TxD}$ 에 대한 attention matrix를 위한 $\text{softmax}(QK^T)V$ 연산
 - QK^T 연산은 Quadratic, computational complexity $O(T^2)$
- Linearized Attention: $\text{softmax}(QK^T)$ 연산을 위해 QK^T 를 $Q'K'^T$ 로 disentangle
 - Computational complexity $O(T)$
 - K'^TV 연산 먼저 수행 후, Q' 와 연산
 - $Q'K'^TV \Rightarrow Q'(K'^TV)$

Linearized Attention

- Un-normalized attention matrix

$$\hat{\mathbf{A}} = \exp(\mathbf{Q}\mathbf{K}^\top)$$

- $\exp(\cdot)$ is applied element-wise
- 기존 softmax 취한 score에 따른 attention matirx에서 normalization을 위한 denominator 생략
- Regular Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D_k}}\right)\mathbf{V} = \mathbf{AV}$$

$$\mathbf{Z} = \mathbf{D}^{-1}\hat{\mathbf{A}}\mathbf{V} \quad \text{where} \quad \mathbf{D} = \text{diag}(\hat{\mathbf{A}}\mathbf{1}_T^\top)$$

$\mathbf{1}_T^\top$: the all-ones column vector of length T

Linearized Attention

- Approximate or replace the unnormalized attention matrix $\exp(QK^T)$ with $\phi(Q)\phi(K)^T$

$$\hat{A} = \exp(QK^T) \rightarrow \phi(Q)\phi(K)^T$$

- ϕ : is a feature map that is applied in row-wise manner

$$z_i = \sum_j \frac{\text{sim}(q_i, k_j)}{\sum_{j'} \text{sim}(q_i, k_{j'})} v_j, \quad \rightarrow$$

Regular Attention $\triangleq \text{sim}(\cdot, \cdot)$

: the exponential of inner product $\exp(\langle \cdot, \cdot \rangle)$

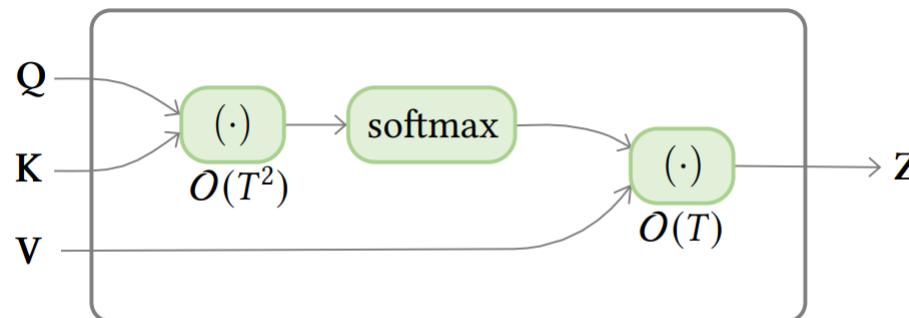
$$\begin{aligned} z_i &= \sum_j \frac{\phi(q_i)\phi(k_j)^T}{\sum_{j'} \phi(q_i)\phi(k_{j'})^T} v_j \\ &= \frac{\phi(q_i) \sum_j \phi(k_j) \otimes v_j}{\phi(q_i) \sum_{j'} \phi(k_{j'})^T}, \end{aligned}$$

$\text{sim}(\cdot, \cdot)$: a kernel function $K(x, y) = \phi(x)\phi(y)^T$

\otimes : outer product

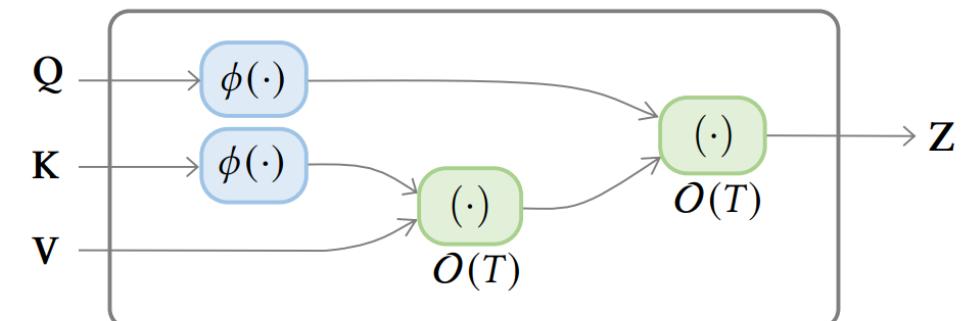
Linearized Attention

$$\mathbf{z}_i = \sum_j \frac{\text{sim}(\mathbf{q}_i, \mathbf{k}_j)}{\sum_{j'} \text{sim}(\mathbf{q}_i, \mathbf{k}_{j'})} \mathbf{v}_j,$$

Regular Attention $\ominus| sim(\cdot, \cdot)$: the exponential of inner product $\exp(\langle \cdot, \cdot \rangle)$ 

(a) standard self-attention

$$\begin{aligned}\mathbf{z}_i &= \sum_j \frac{\phi(\mathbf{q}_i)\phi(\mathbf{k}_j)^\top}{\sum_{j'} \phi(\mathbf{q}_i)\phi(\mathbf{k}_{j'})^\top} \mathbf{v}_j \\ &= \frac{\phi(\mathbf{q}_i) \sum_j \phi(\mathbf{k}_j) \otimes \mathbf{v}_j}{\phi(\mathbf{q}_i) \sum_{j'} \phi(\mathbf{k}_{j'})^\top},\end{aligned}$$

 $sim(\cdot, \cdot)$: a kernel function $K(x, y) = \phi(x)\phi(y)^\top$ \otimes : outer product

(b) linearized self-attention

Linearized Attention

$$\mathbf{z}_i = \sum_j \frac{\text{sim}(\mathbf{q}_i, \mathbf{k}_j)}{\sum_{j'} \text{sim}(\mathbf{q}_i, \mathbf{k}_{j'})} \mathbf{v}_j,$$


Regular Attention의 $\text{sim}(\cdot, \cdot)$

: the exponential of inner product $\exp(\langle \cdot, \cdot \rangle)$

$$\begin{aligned}\mathbf{z}_i &= \sum_j \frac{\phi(\mathbf{q}_i)\phi(\mathbf{k}_j)^\top}{\sum_{j'} \phi(\mathbf{q}_i)\phi(\mathbf{k}_{j'})^\top} \mathbf{v}_j \\ &= \frac{\phi(\mathbf{q}_i) \sum_j \phi(\mathbf{k}_j) \otimes \mathbf{v}_j}{\phi(\mathbf{q}_i) \sum_{j'} \phi(\mathbf{k}_{j'})^\top},\end{aligned}$$

Vector 연산의 summation

$\text{sim}(\cdot, \cdot)$: a kernel function $K(x, y) = \phi(x)\phi(y)^\top$

\otimes : outer product

Attention can be **linearized** by first computing the highlighted terms

연산량 매우 줄어듦

Linearized Attention

$$\begin{aligned} \mathbf{z}_i &= \sum_j \frac{\phi(\mathbf{q}_i)\phi(\mathbf{k}_j)^\top}{\sum_{j'} \phi(\mathbf{q}_i)\phi(\mathbf{k}_{j'})^\top} \mathbf{v}_j \\ &= \frac{\phi(\mathbf{q}_i) \sum_j \phi(\mathbf{k}_j) \otimes \mathbf{v}_j}{\phi(\mathbf{q}_i) \sum_{j'} \phi(\mathbf{k}_{j'})^\top}, \end{aligned}$$

- Memory matrix

$$\sum_j \phi(\mathbf{k}_j) \otimes \mathbf{v}_j$$

- maintains a memory matrix by aggregating associations represented by outer products of (feature mapped) keys and values

$$\phi(\mathbf{q}_i) \sum_j \phi(\mathbf{k}_j) \otimes \mathbf{v}_j$$

- retrieve a value by multiplying the memory matrix with feature mapped query with proper normalization.

(1) Feature map $\phi(\cdot)$ (2) Aggregation rule

Feature Maps

- **Linear Transformer** (ICML 2020, 110회 인용)
 - Simple feature map

$$\phi_i(\mathbf{x}) = \text{relu}(x_i) + 1$$

기존의 dot product attention을 approximate하는 것을 목표로 하지 않고, 비슷한 수준의 성능을 내는 것을 목표로 하여, standard transformer에 준하는 성능 달성

Feature Maps

- **Performer – first version** (arXiv 2020, 17회 인용)

기존의 dot product attention을 approximate하는 것을 목표로 함

- Random feature map

$$\phi(\mathbf{x}) = \frac{h(\mathbf{x})}{\sqrt{m}} [f_1(\omega_1^\top \mathbf{x}), \dots, f_m(\omega_m^\top \mathbf{x}), \dots, f_l(\omega_1^\top \mathbf{x}), \dots, f_l(\omega_m^\top \mathbf{x})],$$

$$f_1, \dots, f_l : \mathbb{R} \rightarrow \mathbb{R} \text{ and } h : \mathbb{R}^D \rightarrow \mathbb{R}.$$

where $\omega_1, \dots, \omega_m \stackrel{\text{iid}}{\sim} \mathcal{D}$ are drawn from some distribution $\mathcal{D} \in \mathcal{P}(\mathbb{R}^D)$

Gaussian kernel을 approximate하기 위해 아래와 같은 kernel 함수를 사용

$$h(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x}\|^2}{2}\right), l = 2, f_1 = \sin, f_2 = \cos.$$

Feature Maps

- Random Feature Attention (ICLR 2021, 21호 인용)

Performer(ver. 1)와 유사

- Random feature map

$$\phi(\mathbf{x}) = \frac{h(\mathbf{x})}{\sqrt{m}} [f_1(\omega_1^\top \mathbf{x}), \dots, f_m(\omega_m^\top \mathbf{x}), \dots, f_l(\omega_1^\top \mathbf{x}), \dots, f_l(\omega_m^\top \mathbf{x})],$$

$$f_1, \dots, f_l : \mathbb{R} \rightarrow \mathbb{R} \text{ and } h : \mathbb{R}^D \rightarrow \mathbb{R}.$$

where $\omega_1, \dots, \omega_m \stackrel{\text{iid}}{\sim} \mathcal{D}$ are drawn from some distribution $\mathcal{D} \in \mathcal{P}(\mathbb{R}^D)$

query, key를 feature space에 보내기 전 **\mathbf{l}_2 -normalization** 하기 때문에, $h(\mathbf{x}) = \mathbf{1}$

$$h(\mathbf{x}) = \exp\left(\frac{\|\mathbf{x}\|^2}{2}\right), l = 2, f_1 = \sin, f_2 = \cos.$$

Feature Maps

- Random Feature Attention (ICLR 2021, 21회 인용)

Performer(ver. 1)와 유사

- Random feature map

The trigonometric random feature map leads to an unbiased approximation,
it **does not guarantee non-negative attention scores**.

$$f_1, \dots, f_l : \mathbb{R} \rightarrow \mathbb{R} \text{ and } h : \mathbb{R}^D \rightarrow \mathbb{R}$$

where $\omega_1, \dots, \omega_m \stackrel{\text{iid}}{\sim} \mathcal{D}$ are drawn from some distribution $\mathcal{D} \in \mathcal{P}(\mathbb{R}^D)$

query, key를 feature space에 보내기 전 \mathbf{l}_2 -normalization 하기 때문에, $h(x) = 1$

$$h(\mathbf{x}) = \exp\left(\frac{\|\mathbf{x}\|^2}{2}\right), l = 2, f_1 = \sin, f_2 = \cos.$$

Feature Maps

- **Performer – second version** (ICLR 2021, 117회 인용)
 - Positive random feature map

$$\phi(\mathbf{x}) = \frac{h(\mathbf{x})}{\sqrt{m}} [f_1(\omega_1^\top \mathbf{x}), \dots, f_m(\omega_m^\top \mathbf{x}), \dots, f_l(\omega_1^\top \mathbf{x}), \dots, f_l(\omega_m^\top \mathbf{x})],$$

$f_1, \dots, f_l : \mathbb{R} \rightarrow \mathbb{R}$ and $h : \mathbb{R}^D \rightarrow \mathbb{R}$.

where $\omega_1, \dots, \omega_m \stackrel{\text{iid}}{\sim} \mathcal{D}$ are drawn from some distribution $\mathcal{D} \in \mathcal{P}(\mathbb{R}^D)$

$$h(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x}\|^2}{2}\right), l = 1, f_1 = \exp$$

**Guarantees unbiased and non-negative approximation
of dot-product attention**

⇒ Performer(ver. 1) 보다 stable하고,
더 좋은 approximation 결과 보임

Feature Maps

- Performer – second version (ICLR 2021, 117호 인용)
 - Positive random feature map

$$\phi(\mathbf{x}) = \frac{h(\mathbf{x})}{\sqrt{m}} [f_1(\omega_1^\top \mathbf{x}), \dots, f_m(\omega_m^\top \mathbf{x}), \dots, f_l(\omega_1^\top \mathbf{x}), \dots, f_l(\omega_m^\top \mathbf{x})],$$
$$f_1, \dots, f_l : \mathbb{R} \rightarrow \mathbb{R} \text{ and } h : \mathbb{R}^D \rightarrow \mathbb{R}.$$

where $\omega_1, \dots, \omega_m \stackrel{\text{iid}}{\sim} \mathcal{D}$ are drawn from some distribution $\mathcal{D} \in \mathcal{P}(\mathbb{R}^D)$

$$h(\mathbf{x}) = 1, l = 1, f_1 = \text{ReLU}.$$

Effective in various tasks including machine translation
and protein sequence modeling.

Feature Maps

- **Linear Transformers Are Secretly Fast Weight Memory Systems** (arXiv 2021, 5회 인용)

feature space 상에서의 orthogonality를 이용할 수 있는 feature map 제안

$$\phi_{i+2(j-1)D}(\mathbf{x}) = \text{ReLU}([\mathbf{x}, -\mathbf{x}])_i \text{ReLU}([\mathbf{x}, -\mathbf{x}])_{i+j}$$

for $i = 1, \dots, 2D, j = 1, \dots, v.$

*input $x \in R^D$
the feature map $\phi : R^D \rightarrow R^{2vD}$*

Aggregation Rule

The associations $\{\phi(k)_j \otimes v_j\}$ are aggregated into the memory matrix by simple summation

⇒ 새로운 association을 memory network S 에 추가할 때, 선택적으로 association을 drop하는 것이 효과적

Aggregation Rule

- Random Feature Attention (ICLR 2021, 21회 인용)
 - Gating mechanism

$$\begin{aligned} g_t &= \text{sigmoid}(\mathbf{w}_g \cdot \mathbf{x}_t + b_g), \\ \mathbf{S}_t &= g_t \mathbf{S}_{t-1} + (1 - g_t) \phi(\mathbf{k}_t) \otimes \mathbf{v}_t, \\ \mathbf{z}_t &= g_t \mathbf{z}_{t-1} + (1 - g_t) \phi(\mathbf{k}_t). \end{aligned}$$

- \mathbf{w}_g and b_g are learned parameters, and \mathbf{x}_t is the input representation at timestep t .
- By multiplying the learned scalar gates $0 < g_t < 1$ against the hidden state $(\mathbf{S}_t, \mathbf{z}_t)$,
history is exponentially decayed, favoring more recent context.

Aggregation Rule

- **Linear Transformers Are Secretly Fast Weight Memory Systems** (arXiv 2021, 5회 인용)

Association의 단순 합으로 memory matrix를 update하는 것은, memory matrix의 capacity를 제한하는 것, 따라서 write-and-remove를 통해 capacity를 확장하는 방식을 제안

- Write-and-remove update

$$\mathbf{k}^{(i)}, \mathbf{v}^{(i)}, \mathbf{q}^{(i)} = \mathbf{W}_k \mathbf{x}^{(i)}, \mathbf{W}_v \mathbf{x}^{(i)}, \mathbf{W}_q \mathbf{x}^{(i)}$$

$$\bar{\mathbf{v}}^{(i)} = \mathbf{W}^{(i-1)} \phi(\mathbf{k}^{(i)})$$

$$\beta^{(i)} = \sigma(\mathbf{W}_\beta \mathbf{x}^{(i)})$$

$$\mathbf{v}_{\text{new}}^{(i)} = \beta^{(i)} \mathbf{v}^{(i)} + (1 - \beta^{(i)}) \bar{\mathbf{v}}^{(i)}$$

$$\begin{aligned} \mathbf{W}^{(i)} &= \mathbf{W}^{(i-1)} + \underbrace{\mathbf{v}_{\text{new}}^{(i)} \otimes \phi(\mathbf{k}^{(i)})}_{\text{write}} - \underbrace{\bar{\mathbf{v}}^{(i)} \otimes \phi(\mathbf{k}^{(i)})}_{\text{remove}} \\ &= \mathbf{W}^{(i-1)} + \beta^{(i)} (\mathbf{v}^{(i)} - \bar{\mathbf{v}}^{(i)}) \otimes \phi(\mathbf{k}^{(i)}) \end{aligned}$$

$$\mathbf{y}^{(i)} = \mathbf{W}^{(i)} \phi(\mathbf{q}^{(i)})$$

$\beta(i)$ is the “write-strength. only depends on input $x(i)$

새로운 input key-value pair가 들어오면,

1) k^i 와 직전 memory matrix W^{i-1} 의 association \bar{v}^i 을 구하고

2) 현재 v^i 와 \bar{v}^i 의 convex combination을 memory matrix에 update함

목적

Reduce the complexity of attention by

(1) Query prototyping: reducing the number of **queries**

(2) Memory compression: reducing the number of **key-value pairs**

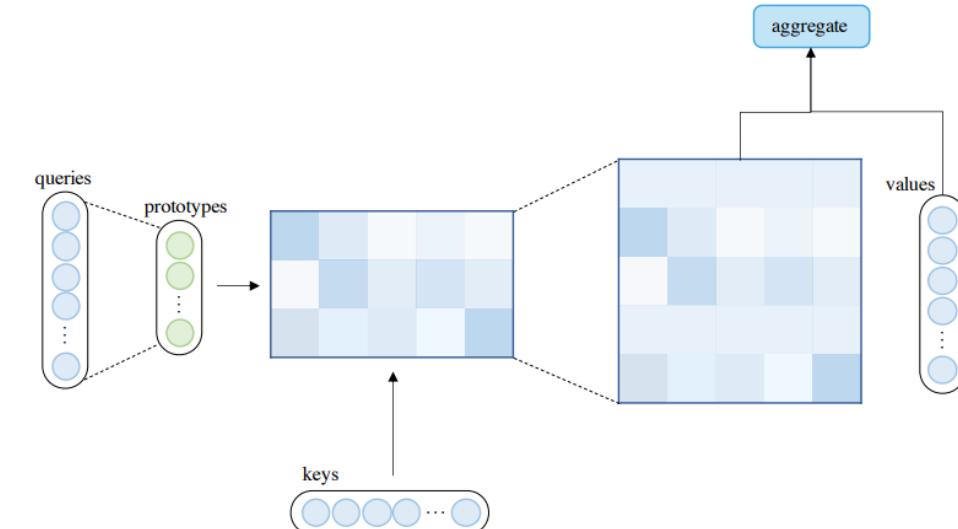
The key-value pairs are often referred to as a **key-value memory**

(1) Attention with Prototype Queries

Attention with Prototype Queries

Several prototypes of queries serve as the main source to compute attention distributions
query representation 중,

- 특정 position의 query distribution을 copy
- discrete uniform distribution으로 채움



(a) Query prototyping

Attention with Prototype Queries

- **Clustered Attention** (arXiv 2020, 20회 인용)
 - groups queries into several clusters and then computes attention distributions for cluster centroids.

① Centroid 구하기

$$Q_j^c = \frac{\sum_{i=1}^N S_{ij} Q_i}{\sum_{i=1}^N S_{ij}}$$

$S_{ij} = 1$, if the i -th query Q_i belongs to the j -th cluster and 0 otherwise.

② Centroid에 대한 attention score 계산 ③ Centroid에 의한 새로운 value 계산

$$A^c = \text{softmax} \left(\frac{Q^c K^T}{\sqrt{D_k}} \right)$$

$Q^c \in \mathbb{R}^{C \times D_k}$ as the centroid matrix

$$\hat{V}^c = A^c V.$$

④ 가장 가까운 centroid에 대한 attention value 도출

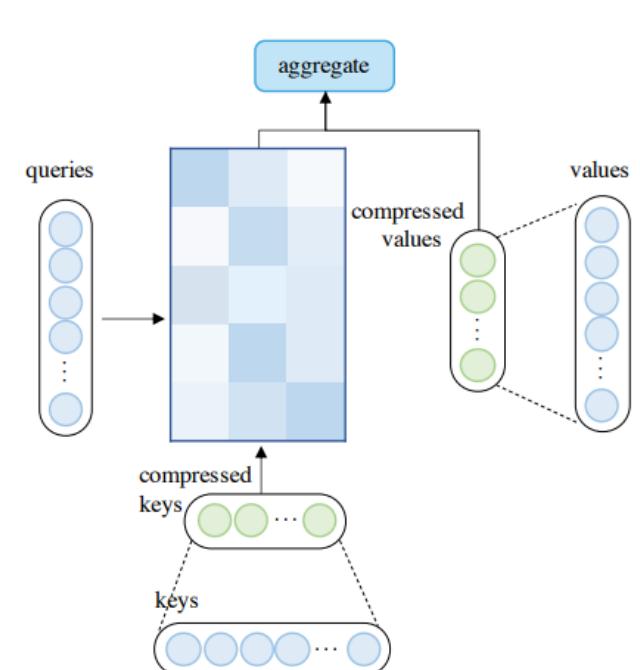
$$\hat{V}_i = \sum_{j=1}^C S_{ij} \hat{V}_j^c.$$

Attention with Prototype Queries

- **Informer** (AAAI 2021, 17회 인용)
 - Query sparsity measurement를 제안하여, 상위 u 개의 query만을 가지고 attention distribution 계산
 - 나머지 query에는 discrete uniform distribution 부여
 - **Query sparsity measurement**
 - Query의 attention distribution과 the discrete uniform distribution 사이의 Kullback-Leibler divergence 값을 기반으로 정의
 - Attention distribution이 the discrete uniform distribution과의 차이가 클수록 몇몇 key에 dominant한 attention을 주는 query라고 할 수 있기 때문에, KLD가 클수록 중요한 query로 봄

Attention with Compressed Key-Value Memory

reduce the complexity by reducing the number of the key-value pairs before applying the attention mechanism



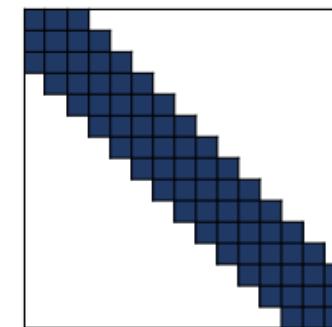
(b) Memory compression

Attention with Compressed Key-Value Memory

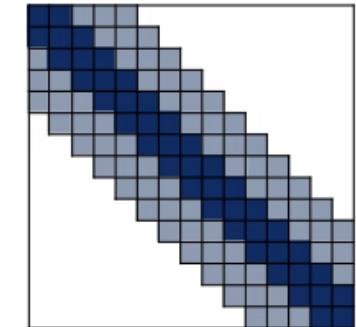
- **Memory Compressed Attention**(ICLR 2018, 379회 인용)
 - Strided convolution을 사용하여 key와 value의 개수를 줄임
 - Kernel size k에 따라 줄이고자 하는 key, value 개수를 조절하며, attention 연산량이 줄어들기 때문에, 같은 시간 안에 vanilla transformer보다 훨씬 긴 sequence를 처리할 수 있다.
- **Set Transformer** (In Proceedings of ICML 2019, 229회 인용),
- **Luna** (arXive 2021, 1회 인용)
 - Trainable global node를 정의하여, input으로부터의 정보를 취합하게 하며, 취합된 정보는 $\text{input}_0 | \text{attend}$ 할 memory로 사용된다.

Attention with Compressed Key-Value Memory

- **Linformer**(arXiv 2020, 132회 인용)
 - Key, value에 대해 linear projection을 적용하여 length n 에서 더 짧은 n_k 의 sequence로 사영
 - 단점: autoregressive attention에 사용할 수 없음
- **Poolingformer** (ICML 2021, 1회 인용)
 - 여러 개의 pooling operation을 사용하여, key와 value를 줄이면서도 receptive field를 키우고자 함



(a) Single-level local attention



(b) Two-level pooling attention

Low-rank Self-Attention

$T \times T$ 의 self-attention matrix가 low-rank = $\text{rank}(\text{self-attention matrix}) << T$

- (1) The low-rank property could be explicitly modeled with parameterization
- (2) The self-attention matrix could be replaced by a low-rank approximation

Low-rank Self-Attention

(1) Low-rank Parameterization.

Low-rank Parameterization

$T \times T$ 의 self-attention matrix가 low-rank = rank(self-attention matrix) << T

특히, input의 짧을 때, $D_k > T$ 로 설정하는 것은, over-parameterization이기 때문에 overfitting을 초래한다.

- **Low-Rank and Locality Constrained Self- Attention for Sequence Modeling**

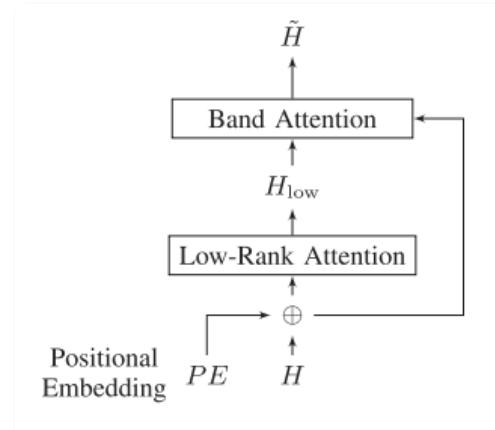
(IEEE/ACM Trans. Audio, Speech and Lang 2019, 3회 인용)

- To capture long-range non-local interaction

Self-attention matrix를 decompose하여 만든 D_k 차원의 low-rank attention

- To capture local dependency

Band attention



Low-rank Approximation

low-rank matrix approximation을 통해 self-attention의 computational complexity 줄이고자 함

- **Nyström method** (AAAI 2021, 14회 인용)
 - Length T의 input에 대한 strided average pooling을 통해 m개의 landmark node만을 선택
 - 선택된 landmark query, key matrix를 \tilde{Q} , \tilde{K} 라 했을 때, 아래의 식에 의해 attention을 추정할 수 있다.

$$\tilde{\mathbf{A}} = \text{softmax}(\tilde{\mathbf{Q}}\tilde{\mathbf{K}}^\top) \left(\text{softmax}(\tilde{\mathbf{Q}}\tilde{\mathbf{K}}^\top) \right)^{-1} \text{softmax}(\tilde{\mathbf{Q}}\mathbf{K}^\top)$$

- $\mathbf{M}^{-1} = \left(\text{softmax}(\tilde{\mathbf{Q}}\tilde{\mathbf{K}}^\top) \right)^{-1}$ 라고 할 때, M^{-1} 이 존재하지 않을 수 있기 때문에 CSALR(IJCAL, 2020, 1회 인용)에서는 M에 identity matrix를 더해 M^{-1} 의 존재를 보장

Q&A

감사합니다