

FASTFORMER

Nguyễn Trần Đăng Hoàng



Fastformer

1. Những hạn chế của Transformer

- Cốt lõi chính của Transformer là cơ chế self-attention, cho phép mô hình hoá context bên trong input sequence,
+ self-attention tính dot product các biểu diễn đầu vào theo từng cặp vị trí trong sequence \rightarrow do đó độ phức tạp tính toán sẽ là $O(N^2)$ với N là độ dài của input sequence \rightarrow xử lý không hiệu quả với long input sequence

- Có nhiều biến thể của Transformer để có thể xử lý với long input sequence (BigBird, Linformer, Longformer, Linear Transformer, ...) tuy nhiên vẫn chưa mô hình hoá đầy đủ được global context.

2. Transformer và cơ chế self-attention.

- Transformer được xây dựng dựa trên Multi-head self-attention để có thể mô hình hoá hiệu quả contexts bên trong sequence bằng cách capture tương tác giữa từng cặp vị trí trong input sequence. Một h-head self-attention có biểu thức sau:

$$\text{Multi-head}(Q, K, V) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h) \cdot W^O$$

Trong đó +, $Q, K, V \in \mathbb{R}^{N \times d}$ là ma trận đầu vào query, key, value.

Với N là độ dài sequence, d là hidden dimension tại mỗi attention head

+ $W^O \in \mathbb{R}^{hd \times d}$ là learnable parameters dùng để linear transformation.

- Mỗi attention head có biểu thức như sau:

$$\text{head}_i = \text{Attention}(Q W_i^Q, K W_i^K, V W_i^V) = \text{softmax}\left(\frac{Q W_i^Q \cdot (K W_i^K)^T}{\sqrt{d}}\right) \cdot V W_i^V$$

Trong đó $\omega, w_i^Q, w_i^K, w_i^V \in \mathbb{R}^{d \times d}$ là ma trận learnable parameter

→ Từ biểu thức trên ta có thể thấy độ phức tạp tính toán là $O(N^2)$, với N là độ dài của input sequence

3. Fastformer

- là 1 biến thể của Transformer, sử dụng cơ chế additive attention để mô hình hoá hiệu quả contexts sequence với độ phức tạp tuyến tính.

⊕ Kiến trúc của fastformer

- Input matrix $E \in \mathbb{R}^{N \times d}$ (với N là độ dài của input sequence, d là hidden dimension), các vector tương ứng sẽ là $[e_1, e_2, \dots, e_N]$
 - Thực hiện phép linear transformation

Với input matrix:

$$\begin{cases} Q = E \cdot W^Q \\ K = E \cdot W^K \\ V = E \cdot W^V \end{cases} \text{ với } W^Q, W^K, W^V \in \mathbb{R}^{d \times d}$$

là learnable parameter matrix

- Khi đó $Q, K, V \in \mathbb{R}^{N \times d}$

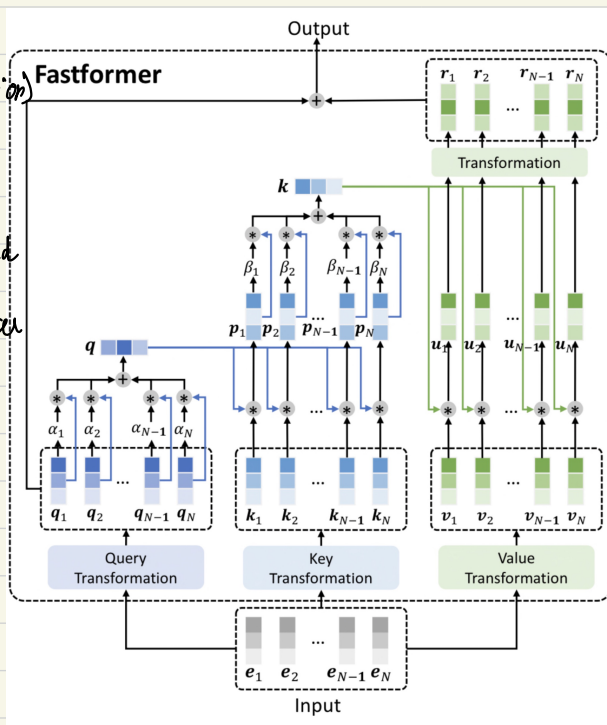
$$Q = [q_1, q_2, \dots, q_N]$$

$$K = [k_1, k_2, \dots, k_N]$$

$$V = [v_1, v_2, \dots, v_N]$$

- Mô hình hoá contexts của input sequence dựa trên tương tác của Query, Key, Value

→ Một cách hiệu quả để giảm độ phức tạp tính toán đó là summarize attention matrix trước khi mô hình hoá tương tác giữa chúng. → Additive attention có thể summarize một cách hiệu quả các thông tin quan trọng trong thời gian tuyến tính



→ Sử dụng additive attention để summarize query matrix thành global query vector; vector này sẽ có thông tin global contexts bên trong attention query. Attention weight α_i tại vector thứ i trong query matrix được tính bằng

$$\alpha_i = \frac{\exp\left(\frac{w_q^T \cdot q_i}{\sqrt{d}}\right)}{\sum_{j=1}^N \exp\left(\frac{w_q^T \cdot q_j}{\sqrt{d}}\right)}$$

Với $w_q \in \mathbb{R}^d$ là learnable parameter vector

→ Khi đó global attention query vector $q \in \mathbb{R}^d$ được tính bằng

$$q = \sum_{i=1}^N \alpha_i \cdot q_i$$

$$\alpha_1 \begin{bmatrix} q_{11} & q_{12} & \dots & q_{1d} \end{bmatrix}$$

+

$$\alpha_2 \begin{bmatrix} q_{21} & q_{22} & \dots & q_{2d} \end{bmatrix}$$

+

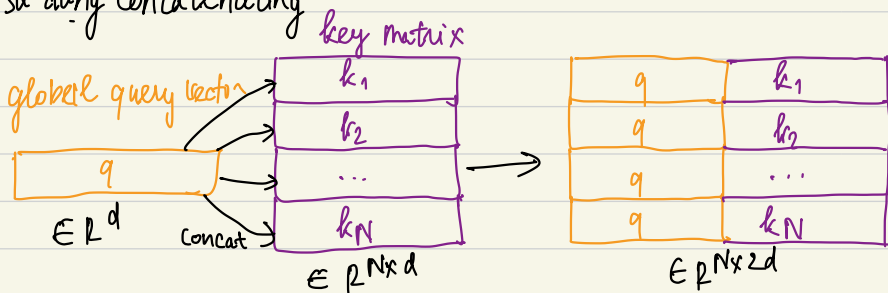
$$\vdots$$

$$\alpha_N \begin{bmatrix} q_{N1} & q_{N2} & \dots & q_{Nd} \end{bmatrix}$$

$$\rightarrow \begin{bmatrix} q \end{bmatrix} \in \mathbb{R}^d$$

— Mô hình hoá tương tác giữa global query vector với key matrix

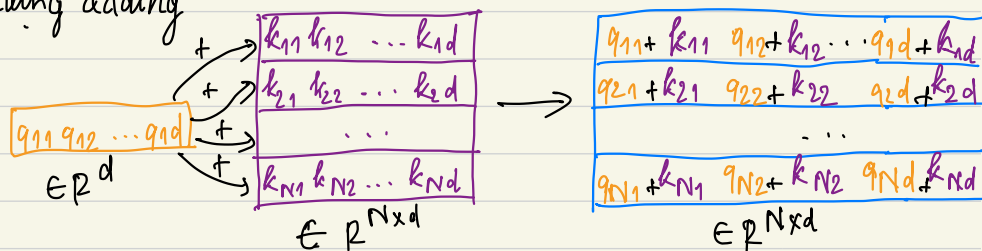
→ Sử dụng concatenating



→ Không mô hình hoá được context vì chỉ đơn giản là concat 2 vector

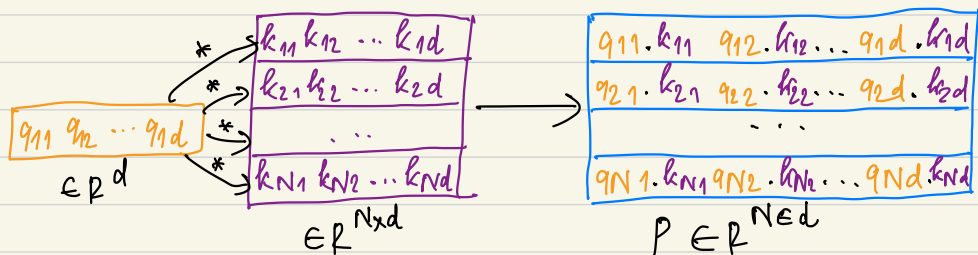
mà không xem xét tới tương tác giữa chúng

→ sử dụng adding



→ không là phương pháp tốt vì vì chỉ có thể mô hình hoá tương tác tuyến tính giữa 2 vector, khiến có thể không học chính xác biểu diễn contexts

→ sử dụng element-wise product $p_i = q * k_i$ $P \in \mathbb{R}^{N \times d}$



→ Mô hình hoá được mối quan hệ phi tuyến giữa 2 vector, giúp mô hình hoá được context phức tạp trong long sequence

— Tương tự sử dụng additive attention để summarize global context-aware key matrix P . Additive attention weigh p_i tại vector p_i được tính bởi

$$p_i = \frac{\exp\left(\frac{w_k^T \cdot p_i}{\sqrt{d}}\right)}{\sum_{j=1}^N \exp\left(\frac{w_k^T \cdot p_j}{\sqrt{d}}\right)}$$

$w_k \in \mathbb{R}^d$ là learnable parameter Vector

→ khi đó global key vector $k \in \mathbb{R}^d$ được tính bằng

$$k = \sum_{i=1}^N p_i \cdot p_i$$

- Tương tự để mô hình hoá tương tác giữa global key vector và attention value matrix, ta sử dụng element-wise product

$$u_i = k * v_i \quad u \in \mathbb{R}^{N \times d}$$

- Sử dụng linear transformation với ma trận tương tác key-value u để học những biểu diễn ẩn. Output matrix sẽ là $P = [p_1, p_2, \dots, p_N] \in \mathbb{R}^{N \times d}$

- Ma trận P sẽ được cộng với ma trận Query ban đầu → output của Fastformer

→ Mỗi global key và query vector sẽ tương tác với value và key vector để học biểu diễn contexts

→ Stacking nhiều Fastformer layers sẽ mô hình hoá đầy đủ thông tin contexts

→ Sử dụng kĩ thuật parameters sharing, chia sẻ transformation parameter trong value và query để giảm chi phí memory. Thêm vào đó sử dụng parameter sharing giữa các layer Fastformer khác nhau để giảm overfitting

⊗ độ phức tạp tính toán

- Additive attention để học global query và global key sẽ có chi phí tính toán và chi phí bộ nhớ đều là $O(N \cdot d)$ với N : là độ dài input sequence
 d : là hidden dimension

- Element-wise product cũng có chi phí tính toán và chi phí bộ nhớ là $O(N \cdot d)$

→ Vậy độ phức tạp sẽ là $O(N \cdot d)$, hiệu quả hơn rất nhiều so với Transformer với độ phức tạp $O(N^2 \cdot d)$

- Nếu kỹ thuật parameters sharing được sử dụng, tổng số lượng param tại mỗi layer Fastformer sẽ là $3hd^2 + 2hd \ll$ với ít nhất $4hd^2$ param trong Transformer

⊗ Những bài toán mà Fastformer cho thấy performance gần như là tốt nhất so với các biến thể khác của Transformer

- News recommendation task ← text modeling
← understanding user interest
- Text summarization task
- Sentiment and topic classification task

⊗ Ảnh hưởng của parameter sharing

- Query-Value sharing: có performance tốt hơn một chút với không sharing
- Query-Value + Head wise sharing: performance bị sụt giảm do các attention-head có xu hướng capture những mẫu contexts khác nhau phải sharing sẽ không có lợi trong việc mô hình hoá contexts
- Query-Value + layer wise sharing = tăng performance vì giúp giảm thiểu overfitting

→ Fastformer sử dụng query-value + layer wise sharing sẽ làm tăng performance cũng như giảm parameters size