

emb-dim

emb-dim

center

$$|V| \begin{pmatrix} \text{Word_emb_1} \\ \vdots \\ \text{Word_emb_l} \end{pmatrix} \begin{matrix} v_1 \in \mathbb{R} \\ v_2 \\ \vdots \\ v_{|V|} \end{matrix}$$

$o_{-m} \dots o_1$	c	$o_1 \dots o_m$
$\underbrace{\hspace{4em}}_m$		$\underbrace{\hspace{4em}}_{m_1}$

$$\hat{v} = \frac{1}{2m} \sum_{j=1}^m (v_{o_{-j}} + v_{o_j})$$

emb-dim

emb-dim

$$|V| \begin{pmatrix} \text{word_emb_2} \\ \vdots \end{pmatrix} \begin{matrix} u_1 \\ u_1 \\ \vdots \\ u_{|V|} \end{matrix}$$

$$\Pr[c \mid o_{-m}, o_{-(m-1)}, \dots, o_1, o_1, \dots, o_m]$$

$$= \frac{\exp(u_c^T \cdot \hat{v})}{\sum_{j \in V} \exp(u_j^T \cdot \hat{v})}$$

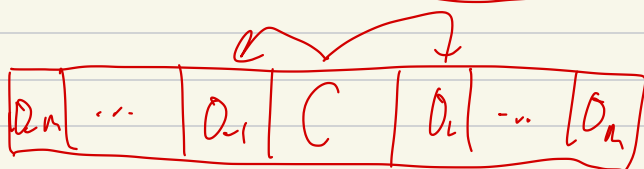
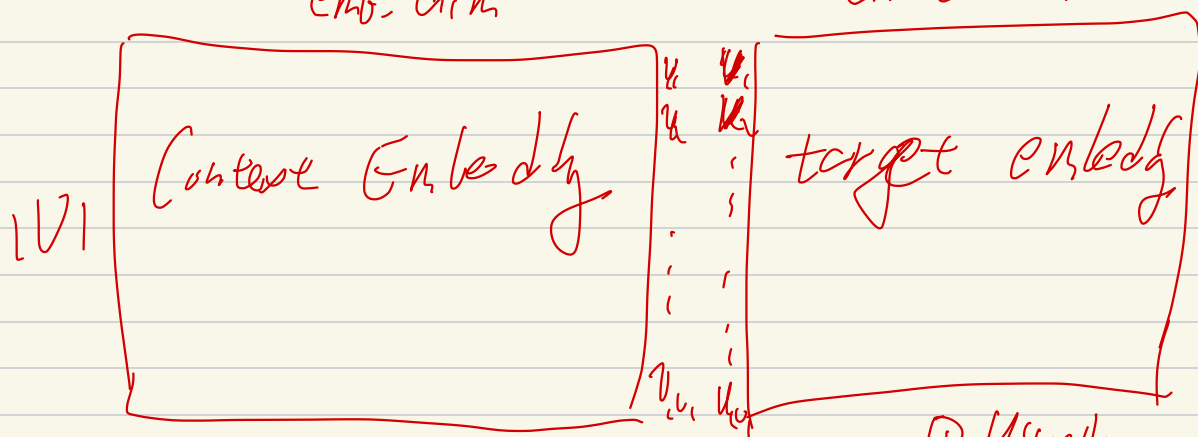
$$\underset{v, u}{\operatorname{argmin}} \sum_{\text{GED}} \log \left(\frac{\exp(u_c^T \cdot \vec{v})}{\sum_{j \in V} \exp(u_j^T \cdot \vec{v})} \right)$$

Output: v .

Skip gram

emb. dim

emb. dim



① Usually

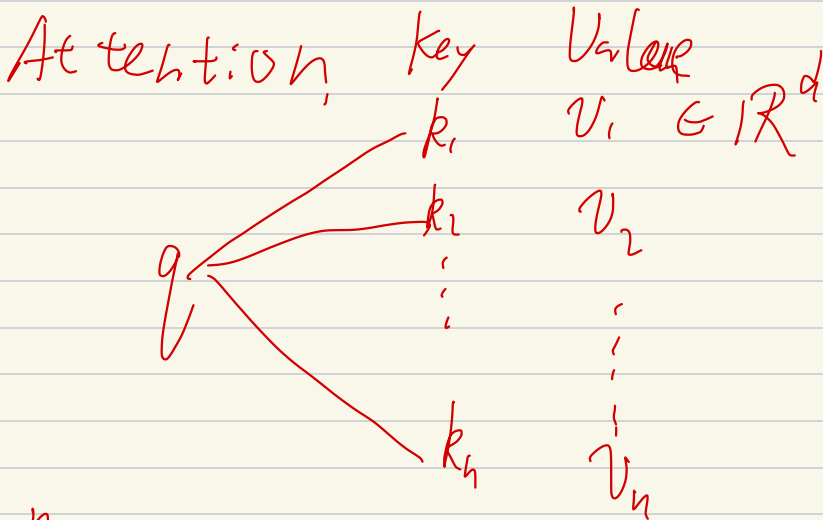
v

② $\frac{u+v}{2}$

$$\underset{\text{GED}}{\operatorname{min}} \sum \sum_{\substack{j=-m \\ j \neq 0}}^m \left(\log \frac{\exp(u_j^T \cdot v_c)}{\sum_{i \in V} \exp(u_i^T \cdot v_c)} \right) \quad \text{③ } (u, v)$$

$i \in V' \in V$

Negative Sampling,



$$\sum_{i=1}^n \frac{1}{2} f(q, k_i) \cdot v_i$$

Transformer, w_1 w_2 ... w_n

x_1 x_2 ... $x_n \in \mathbb{R}^d$

$$w_q, w_k, w_v \in \mathbb{R}^{d \times d}$$

$$q_i = w_q x_i \in \mathbb{R}^d, \quad k_i = w_k \cdot x_i, \quad v_i = w_v x_i$$

query

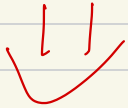
key

Value,

$$w'_{i,j} = \frac{q_i^T \cdot k_j}{\sqrt{d}} \xrightarrow{\text{softmax}} w_{ij} = \frac{\exp(w'_{ij})}{\sum_{j'=1}^n \exp(w'_{ij'})}$$

$$y_i = \sum_{j=1}^n w_{ij} \cdot v_j \in \mathbb{R}^d$$

$x_1 \ x_2 \dots x_n$



$y_1 \ y_2 \dots y_n$

① Better Parallelization \Rightarrow multi-head

② y is linear in v .

Add MLP $y \rightarrow \text{MLP}(y)$

③ Sequence Information.

+ Position Encoding

input = x + Position Embedding

④, No future information,
masked Attention.

⑤

Optimization ?

Skip-connection.

Layer Normalization