

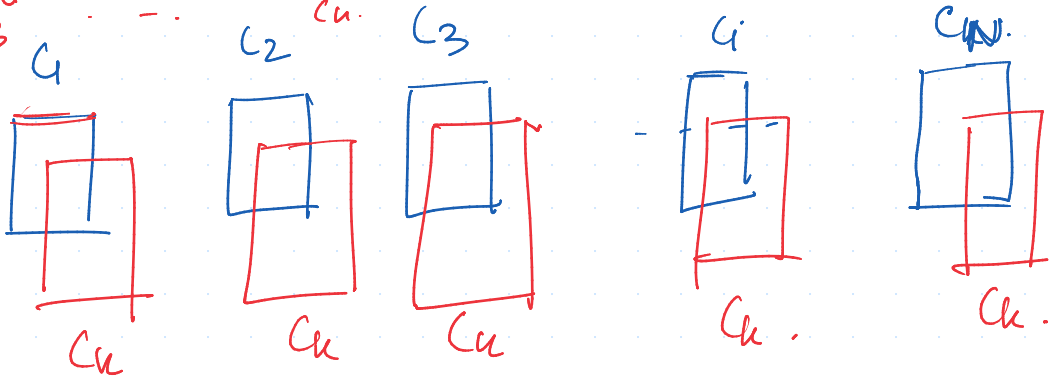
# Lec - 10 : Transformers

21/3/26

word  $\rightarrow$  tokenization  $\rightarrow$  embedding  $\rightarrow$   $C_k$   $\xleftrightarrow{\text{code}}$

the cat sat on the mat.

$k$ th embedding.



with the

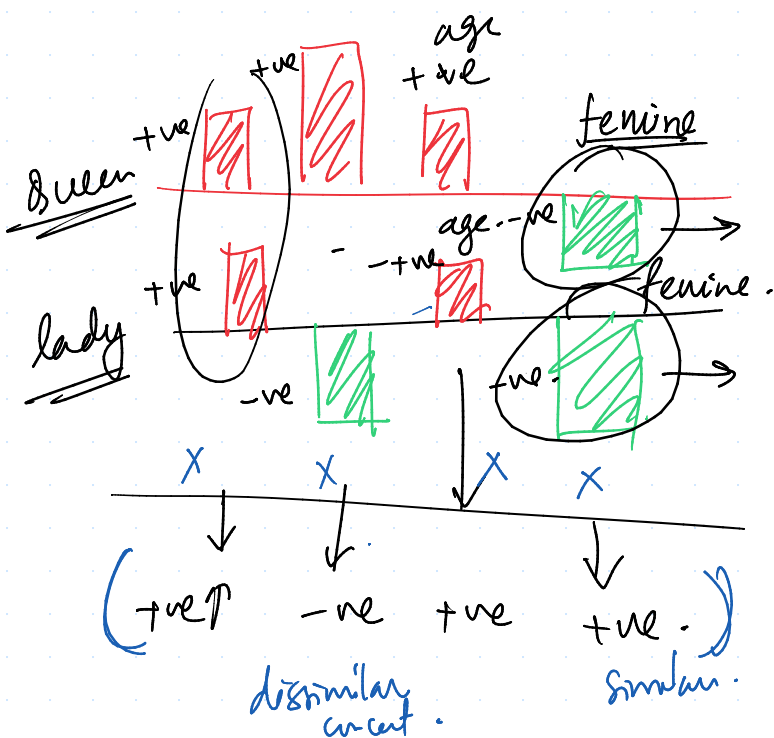
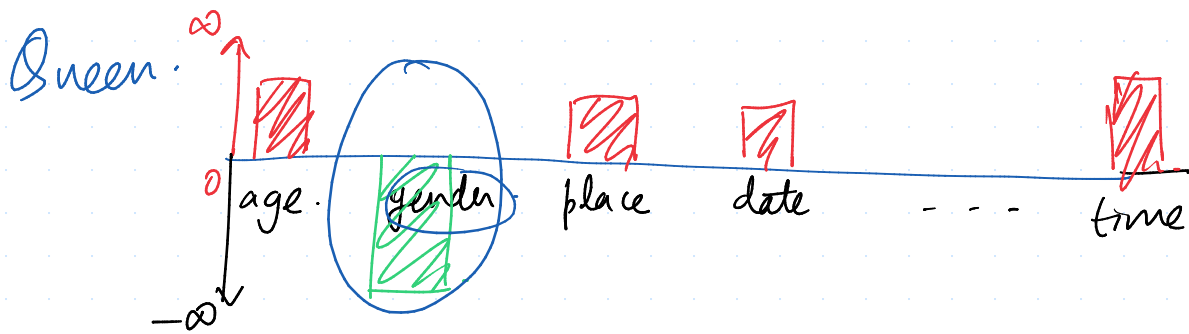
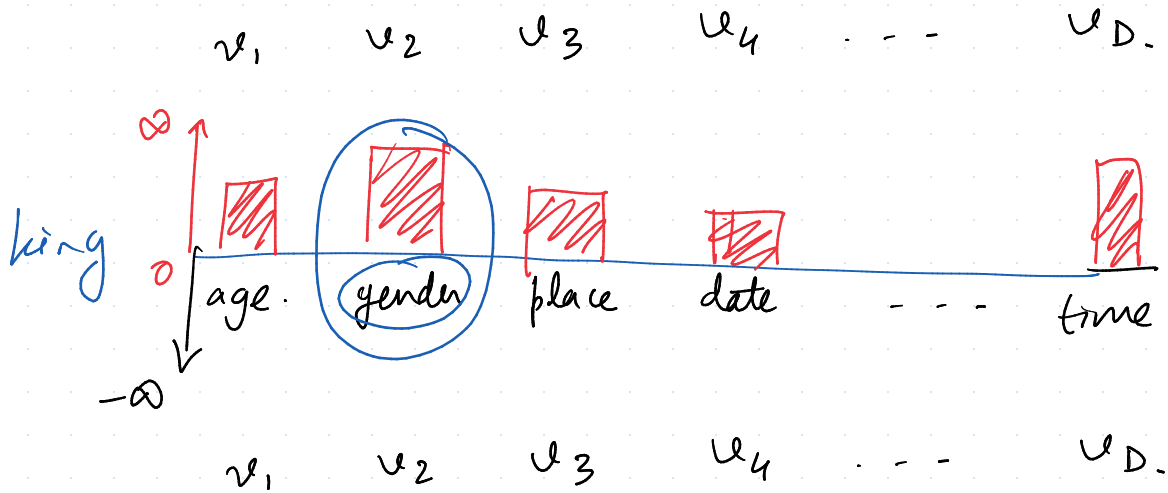
$k$ -th word, perform the inner product to each of the  $n$ th word.

embedding  $\left\{ \begin{array}{l} C_k = \langle v_1^{(k)}, v_2^{(k)}, v_3^{(k)}, \dots, v_D^{(k)} \rangle \rightarrow D\text{-dim vector.} \\ C_i = \langle v_1^{(i)}, v_2^{(i)}, v_3^{(i)}, \dots, v_D^{(i)} \rangle \end{array} \right.$

Dot product / inner product  $= \sum_{j=1}^D v_j^{(k)} v_j^{(i)} = \text{scalar.}$

$C_i$ 's  $\rightarrow$  are complex combination of different concepts of words learnt based on the expos. training data.

Let us assume that we have trained a model (word-2-vec) on a corpus, we have learnt the word embeddings for the corpus.



When concepts are similar, they have a similar interpretation in the feature space.

Similar concepts when multiplied gets a large +ve value & dissimilar concepts when multiplied gets a large -ve value.

$(-\infty, +\infty)$

→ Probability dist

• +ve values.

•  $\Sigma \Rightarrow 1$ .

-10	-2	10	0	-5	
-5	-7	-3	2	1	
↓	↓	↓	↓	↓	
50	14	<u>-30</u>	0	-5	-
		↓			
30	30	30	30	30	-
80	44	0	30	25	-

Adding the least value.

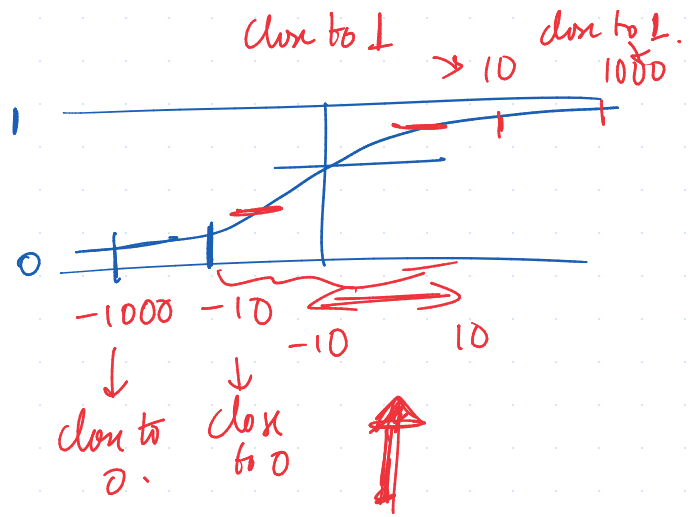
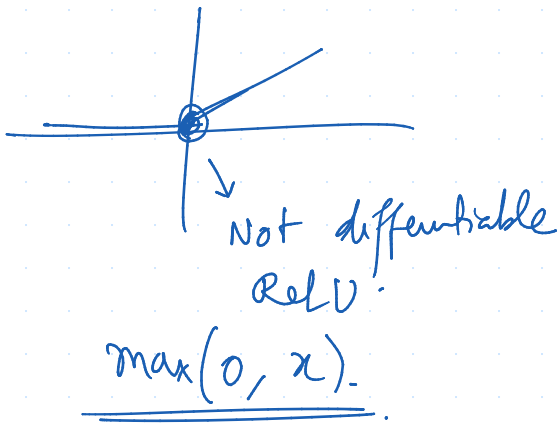
↓  
Normalize.

Dividing by the max value.

Not differentiable. Finding the least value is

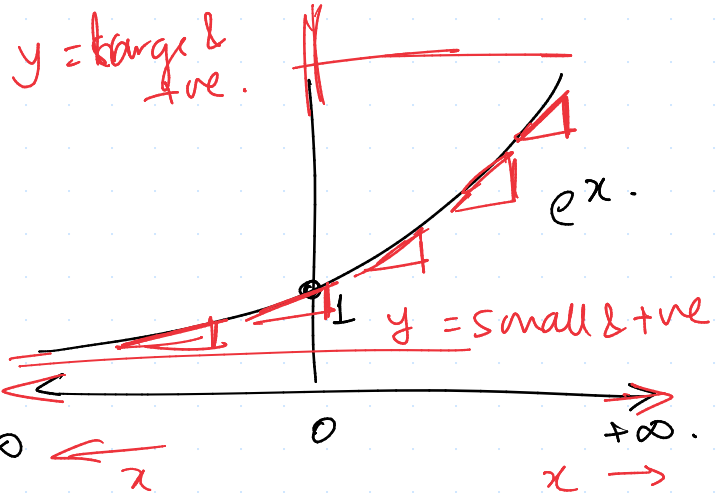
like sorting -- (min, max) → differentiable(x)

finding min-max is not differentiable.



monotonically increasing  
function.

$$\frac{d}{dx} e^x = e^x$$



$e$  [resultant vector]

$\Rightarrow$  +ve values.

$\rightarrow$  since  $e^x$  is a monotonically increasing function  
the small values scales up accordingly

$\exp(C_i \cdot C_u) \rightarrow$  +ve output.

preserve the meaning of the product.

$$\underbrace{\gamma_{k \rightarrow 1}, \gamma_{k \rightarrow 2}, \dots, \gamma_{k \rightarrow i}, \dots, \gamma_{k \rightarrow N}}_{\substack{\downarrow \\ c_k \cdot c_1 \quad c_k \cdot c_2 \quad \dots \quad c_k \cdot c_N}}$$

relative sensitivity of the  $k$ th word to the  $i$ th word

Normalization

$$\gamma_{k \rightarrow i} = \frac{\exp(c_k \cdot c_i)}{\exp(c_k \cdot c_1) + \exp(c_k \cdot c_2) + \dots + \exp(c_k \cdot c_N)}$$

$$= \frac{\exp(c_k \cdot c_i)}{\sum_{n=1}^N \exp(c_k \cdot c_n)} \geq 0.$$

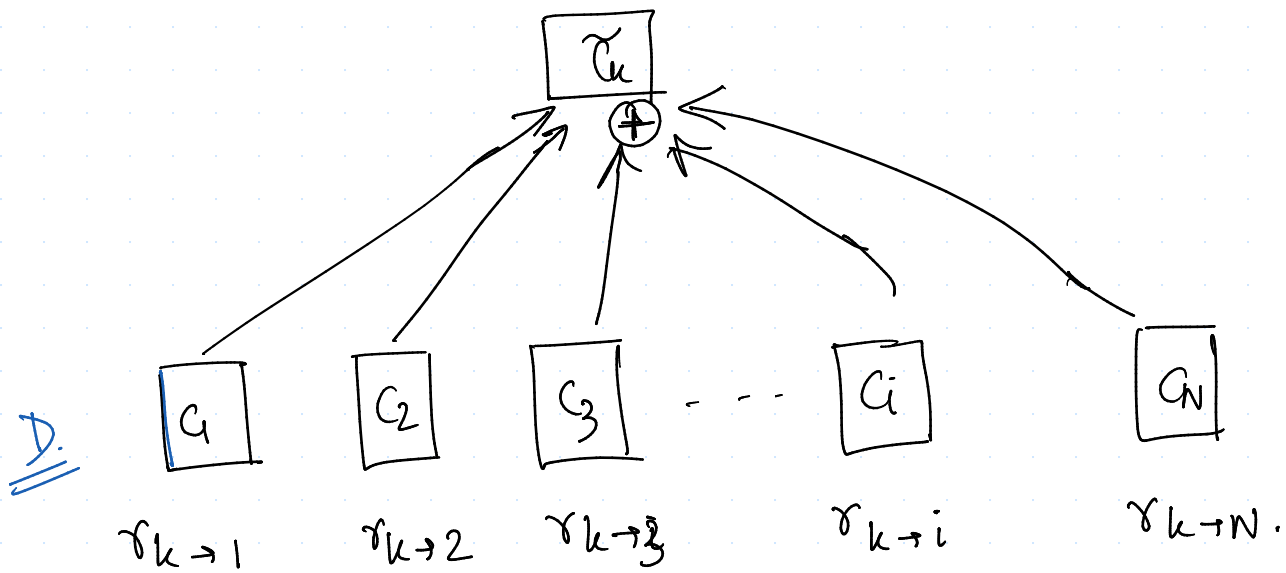
between  
(0-1)

relative relationship b/w 2 word vectors.

probability distribution.

-ve values  $\rightarrow$  small & +ve.

+ve values  $\rightarrow$  large & +ve.



$C_k \rightarrow$  not informed of the context / independent of the context.

$\begin{bmatrix} \vdots \end{bmatrix}$      $\begin{bmatrix} \vdots \end{bmatrix}$      $\begin{bmatrix} \vdots \end{bmatrix}$      $\begin{bmatrix} \vdots \end{bmatrix}$   
 $C_1 \cdot C_k$  ,  $C_2 \cdot C_k$  , ... ,  $C_i \cdot C_k$  , ...  $C_N \cdot C_k$ .  
 $\sigma_{k-1}$      $\sigma_{k-2}$      $\sigma_{k-i}$      $\sigma_{k-N}$

(relative relationship)

Multiply the relative degrees by corresponding codes & add them up  $\tilde{C}_k$ .

If  $C_k$  is highly related to a particular word, then the corresponding relative relationship quantified through  $\sigma$  will be larger.

$\alpha_{k \rightarrow 1}$  ; how much attention to pay to the corresponding word/code.

$$\tilde{C}_k = (\alpha_{k \rightarrow 1} \times C_1) + (\alpha_{k \rightarrow 2} \times C_2) + \dots + (\alpha_{k \rightarrow i} \times C_i) + \dots + (\alpha_{k \rightarrow N} \times C_N)$$

$$= \sum_{n=1}^N \alpha_{k \rightarrow n} \times C_n$$

$\tilde{C}_k \uparrow$   $\rightarrow$  pay high attention to  $C_k$

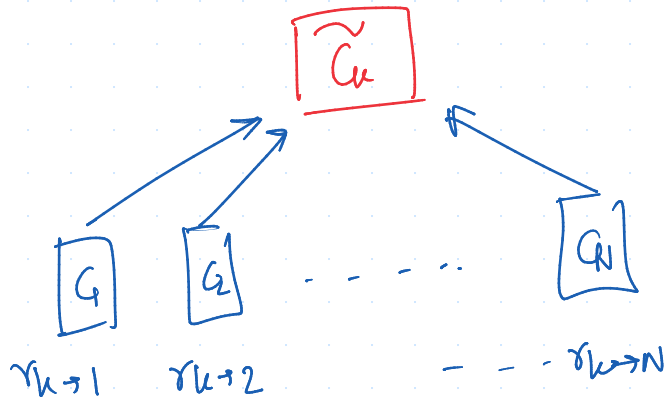
$\tilde{C}_k \downarrow$   $\rightarrow$  pay low attention to  $C_k$ .

$C_k$  to  $\tilde{C}_k$   $\rightarrow$  mapping / vector transformed with context.

takes into account the surrounding context.

Values.

Keys.



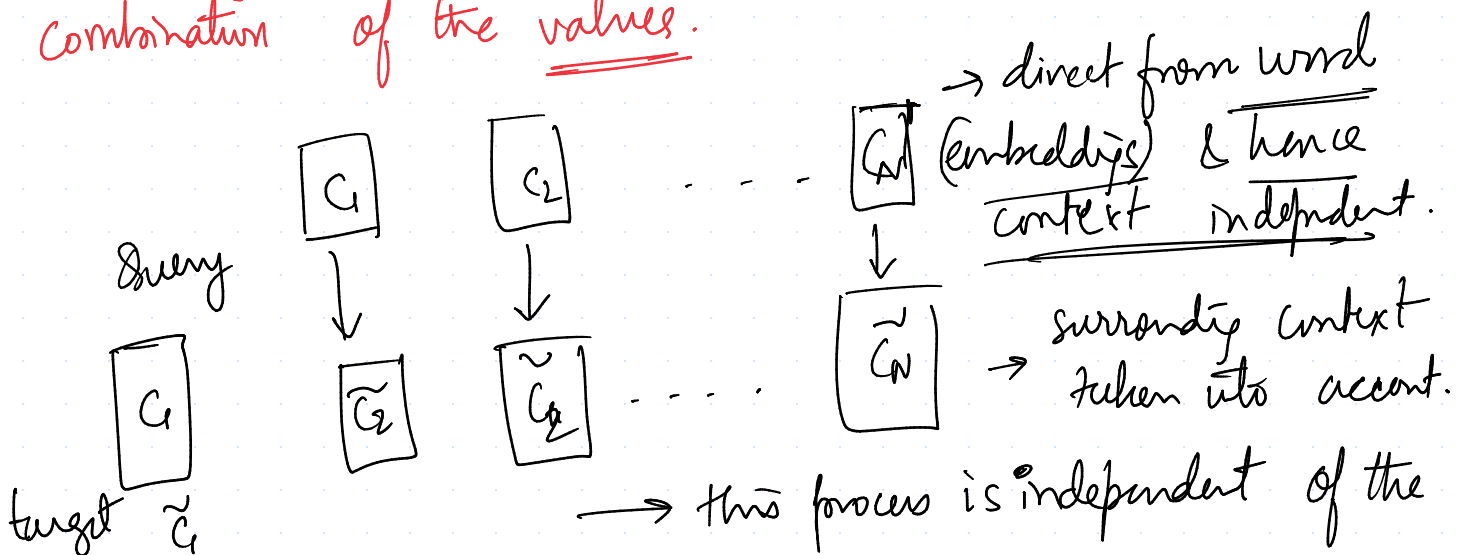
Query



Code  $C_k$  transformed to  $\tilde{C}_k$ ,

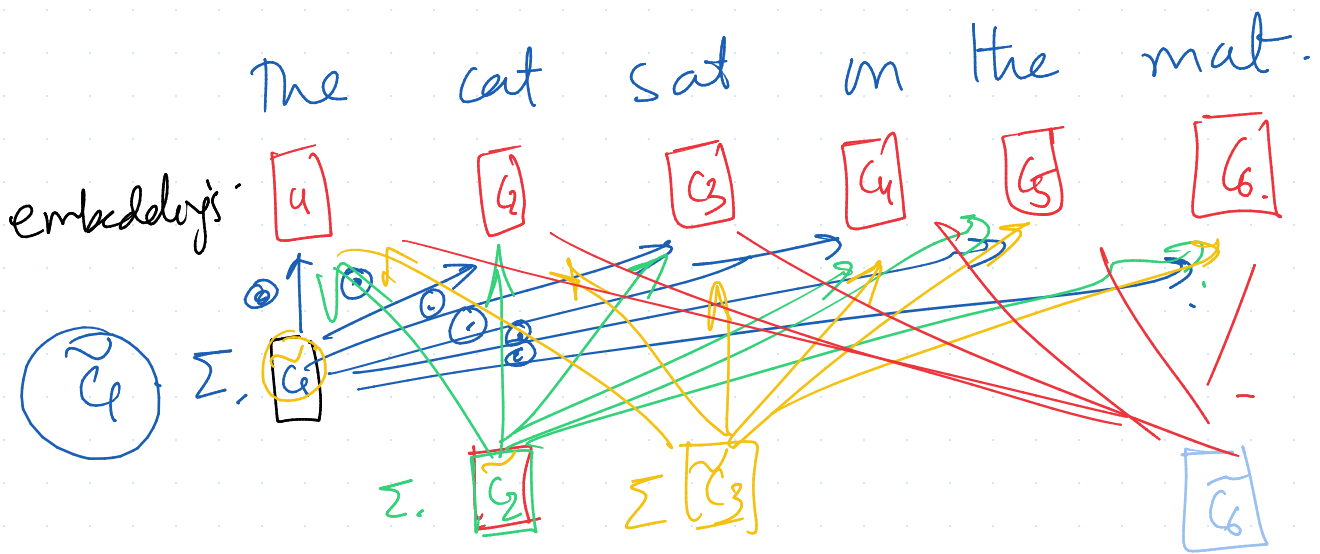
Relative weights  $\{r_{k-1}, r_{k-2}, \dots, r_{k-N}\}$  are always one & sum to 1.

Each query  $C_k$  mapped to  $\tilde{C}_k$  is a convex combination of the values.



word order.

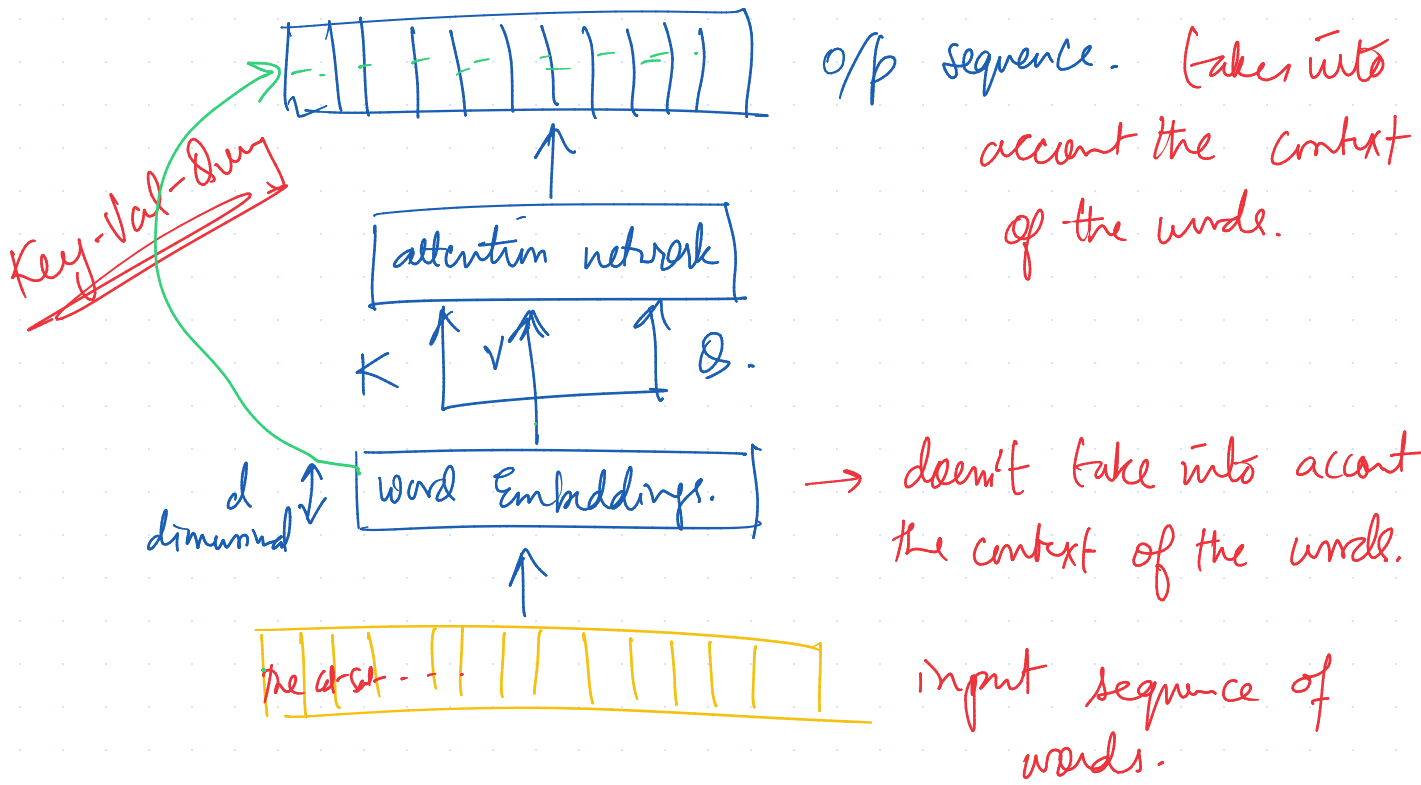
The word 'order' matters, if we want to capture the contextual meaning.



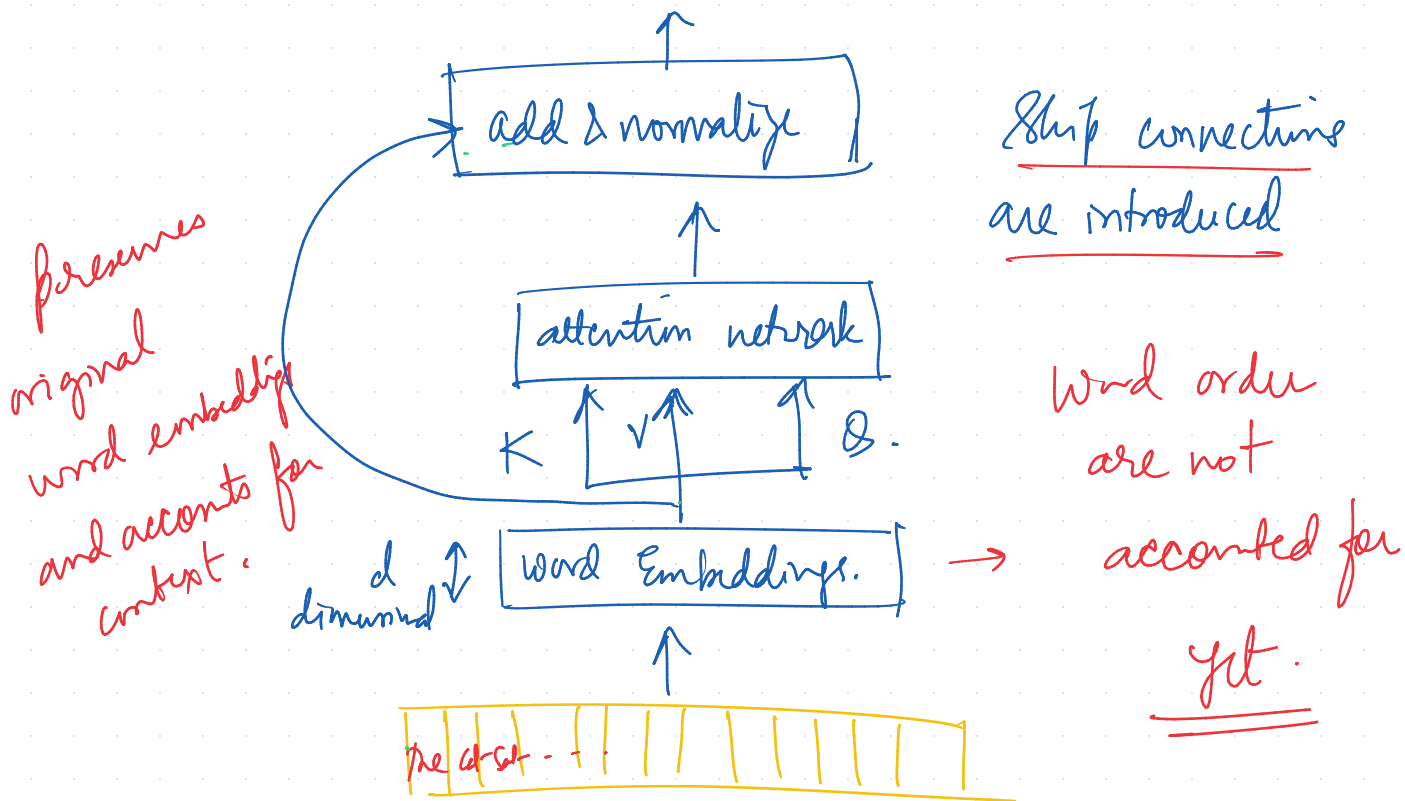
Semantic  
meaning is  
different.

The cat sat on the mat. ✓

The mat sat on the cat. ✗



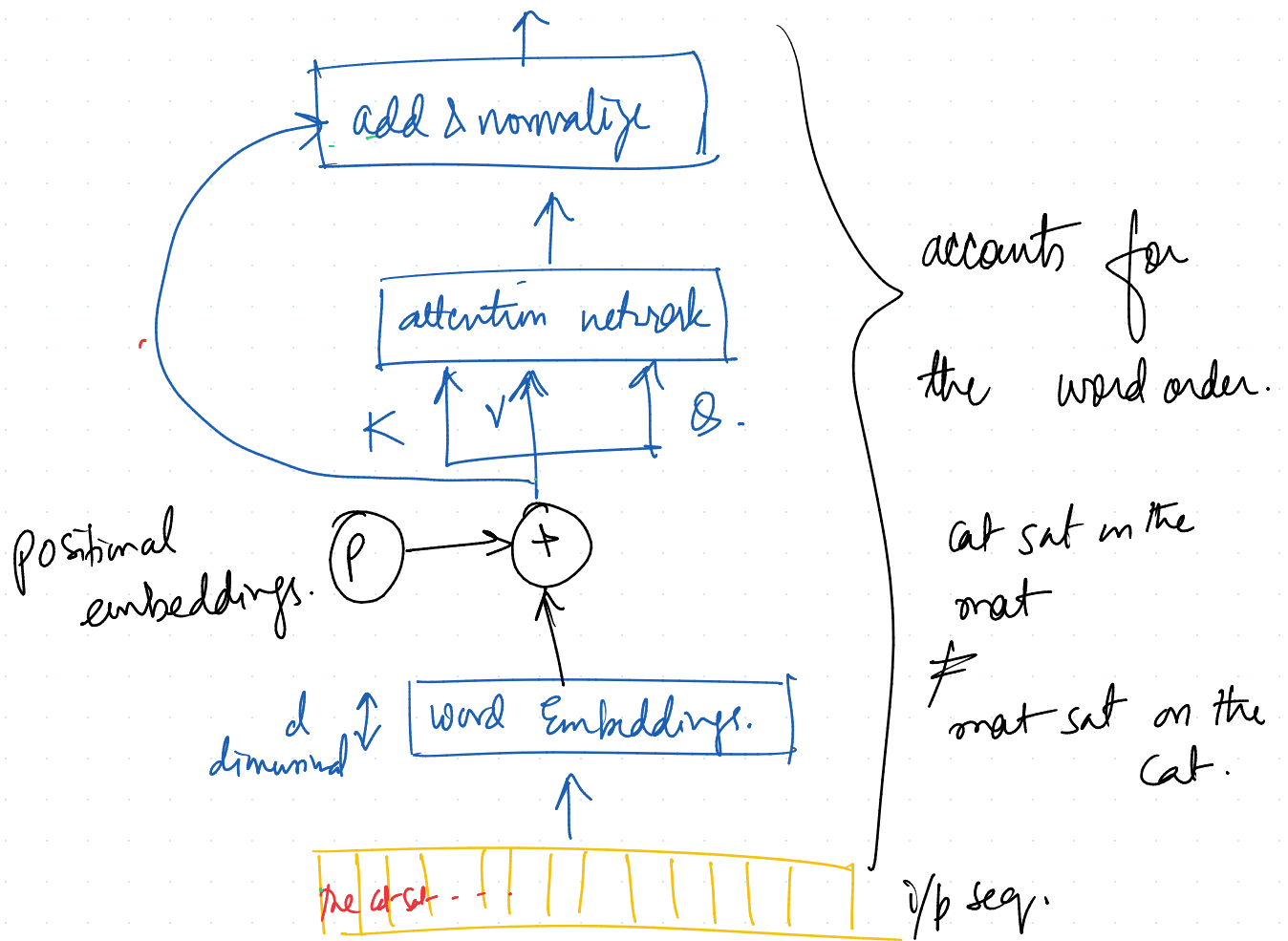
We join the word embeddings at the bottom, when we go from the bottom to the top.



skip connections  
are introduced

word order  
are not  
accounted for  
yet.

# Positional Embeddings



## positional embeddings.

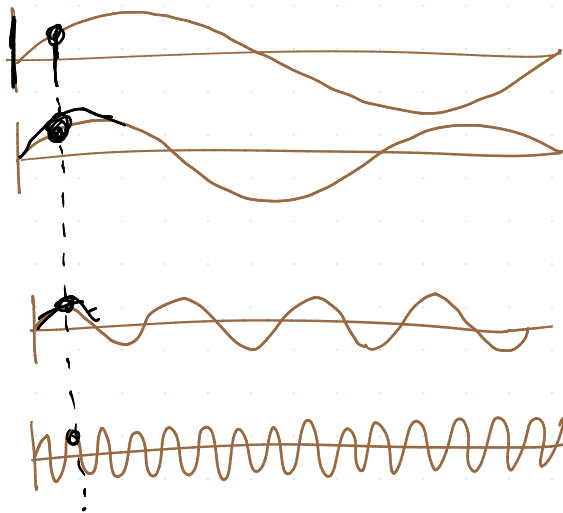
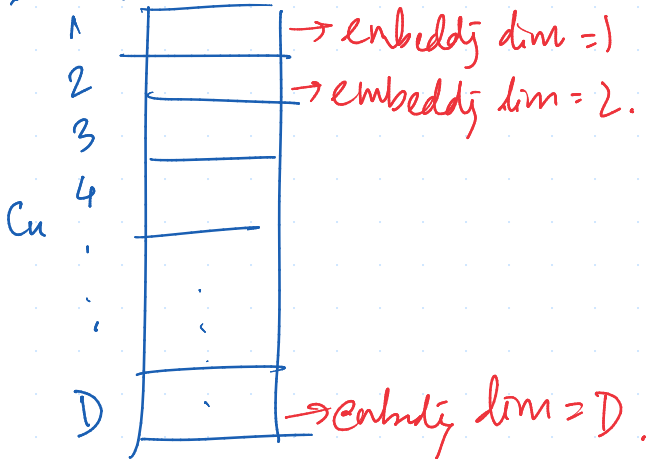
$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d}}\right)$$

$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$

} d-dimensional  
positional  
embeddings.

→ represented by sine & cosine waves.

embeddings



sine-wave →

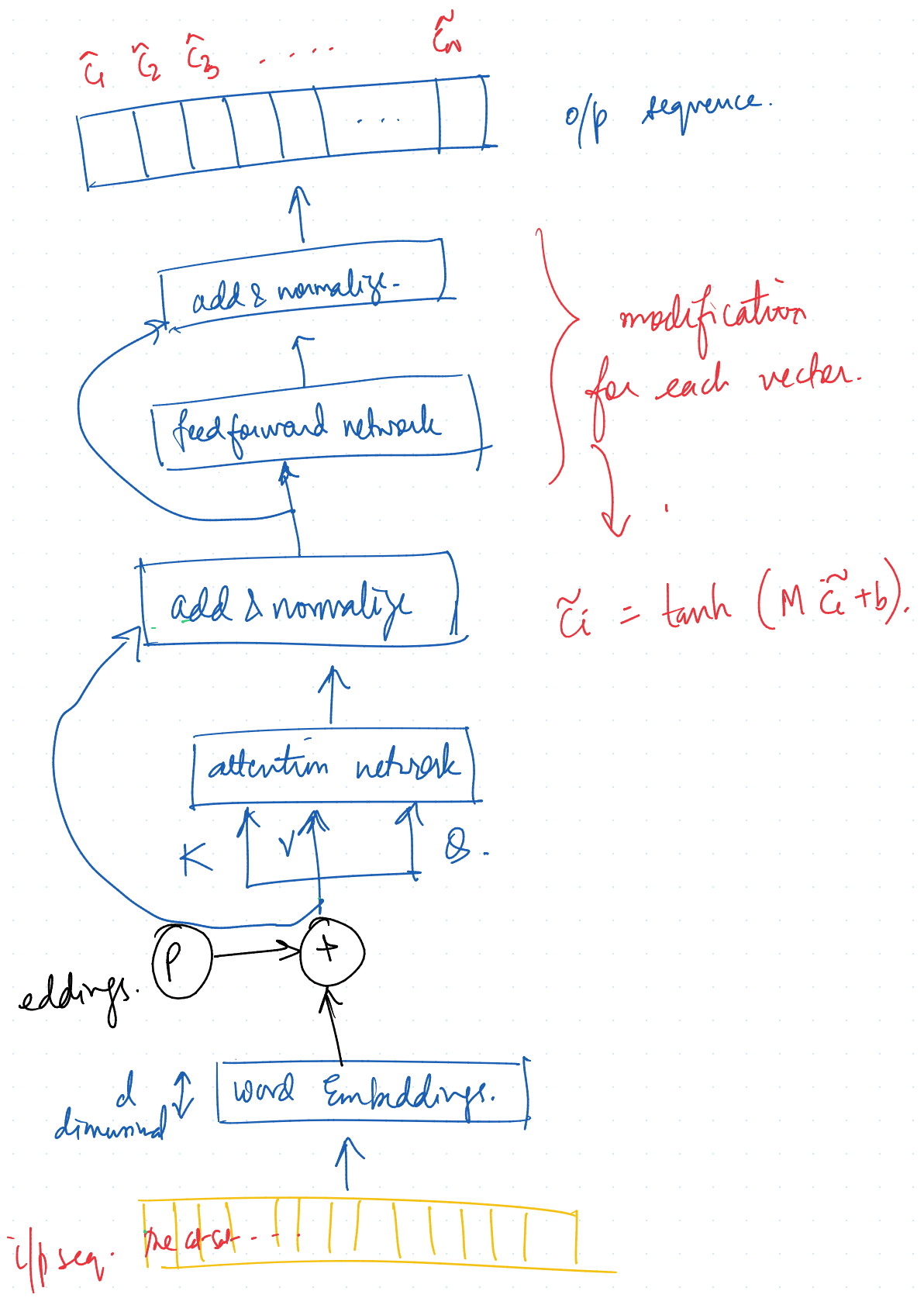
closer word will get more embedding context than the later words.

$$\sin\left(\frac{\text{pos}}{10000^{2i/d}}\right)$$

$$\begin{array}{c} \text{pos} \\ \hline \uparrow \frac{2i}{(10000)^{10}} \\ \downarrow \frac{\text{pos}}{(10000)^{\frac{2i}{10000}}} \end{array}$$

★ ★

check potential embeddings from blogs etc in the internet.



Words that are nearby positionally will have similar positional embeddings, therefore their inner products will be high.

Words that are far apart  $\rightarrow$  word positional embeddings will be different  $\rightarrow$  hence  $\rightarrow$  -ve inner products hence they will not be related.