

# Moving Towards GPT

## Generatively Pre-trained Transformer.

### Self Attention?

What does self-attention do?

We already have a NN model which can predict words or tokens ahead.

But, the current embedding of words into the embedding space isn't satisfactory.

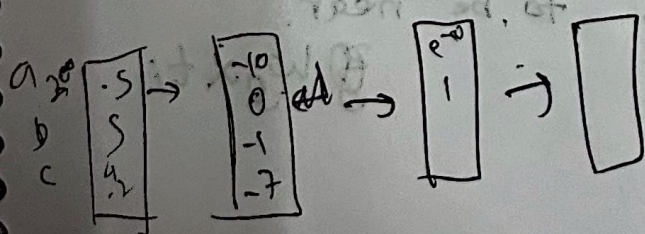
We can't ~~to develop~~ the mapping from words to their embeddings to have richer context hidden in it.

So now, for eg. we want

$$E(\text{man}) - E(\text{woman}) = E(\text{nephew}) - E(\text{neice})$$

Basically adding a vector  $x$  on current

words to their ~~feminine~~ <sup>female</sup> versions etc.





We want the model to understand the context.

So we want the previous words to send information to the last current word in the sequence, and then the prediction will be based on the last in the sequence.

Eg: have adjectives adjust meanings of corresponding nouns.

initial embedding  $\longrightarrow$  converts words to a high dimensional space [here 12000] and their position.

Now, first we need to identify which ~~words~~ preceding words (tokens) affect the current word embedding, for this we use a  $K, Q$

Key Query system.

$\downarrow$   $\downarrow$   
 $W_K$   $W_Q$   $\longrightarrow$  dim = smaller emb space  $\times$  Oq emb spa  
 $128 \times 12000$

let  $A_i$  be the set of tokens preceding a token  $t_i$  that affect  $t_i$  more

$\therefore$  We want  $W_K \times A_{ij}$  to be near.

in the smaller space for

$W_Q \times t_i$

Basically, for all words in the current context window:

<del>the</del>	thiefy	blee	crean	rounded	the	werd	for
	$\downarrow$ $w_{\text{word}} \rightarrow \vec{g}_1$						
$\vec{w}_i \rightarrow \vec{k}_i$	$\vec{k}_i \cdot \vec{g}_1$	$\vec{k}_i \cdot \vec{g}_2$					
ph							
crea							
roan							
the							
wa							
fore							

$\therefore \vec{k}_i \cdot \vec{g}_j \Rightarrow$  How much should word/token  $i$  affect word/token  $j$ .

also we divide

$$\frac{\vec{k}_i \cdot \vec{g}_j}{\sqrt{d}}$$

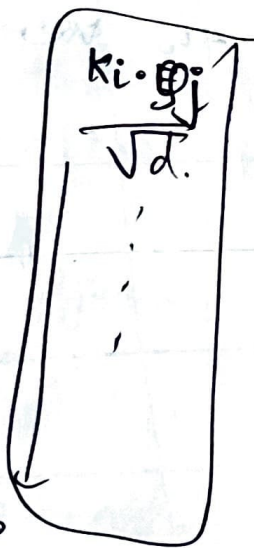
$\rightarrow$  dim of sentence emb space



Then we apply softmax on each col

$\theta_j$

$k_1$   
 $k_2$   
 $\vdots$

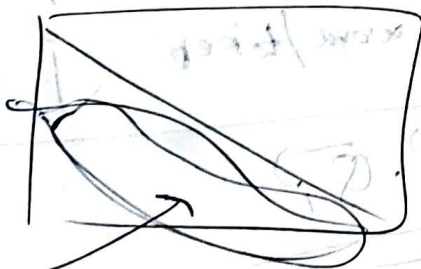


Softmax to get prob of how much that word affects word.  $j$

At prediction time our model only has access to words prior to the one being predicted so we can only want words before current one to affect them :)

$\therefore$  we set

$$k_i \cdot \theta_j = -\infty \quad \forall i > j$$



so their prob after softmax becomes 0.

called MASKING.

in bidirectional models, where the goal is to understand the meaning of a sentence rather than predict. then we can use that info too.

Seeing this square matrix we can see how context size is a huge bottleneck for models.

Now that we know how much a word affects another, how to update that value.

Value vector that lives in the emb space

$$\begin{aligned} \therefore \vec{E}_j &= \vec{E}_j + \Delta \vec{E}_j \\ &= \vec{E}_j + \sum \left( \frac{\vec{K}_i \cdot \vec{Q}_j}{\sqrt{d}} \right) \cdot \vec{V}_i \end{aligned}$$

$$\vec{E}_j = \vec{E}_j + \sum \text{softmax} \left( \frac{(\vec{W}_K \times \vec{E}_i) \cdot (\vec{W}_Q \times \vec{E}_j)}{\sqrt{d}} \right) \cdot \vec{W}_V \times \vec{E}_i$$

No. of paramet

$$Q \rightarrow 12 \times 10^3 \times 120 \sim 1.5 \times 10^6$$

$$K \rightarrow 1.5 \times 10^6$$

$$V \rightarrow (1.2 \times 10^4)^2 \rightarrow 150 \times 10^6$$



Value  $\rightarrow$  too large

$\therefore$  try to keep  $|value| \sim |Queue|$  order.

$\therefore$  we try to break down Value.

$V_d$  and  $V_r$ .

large  $V_d$   $\rightarrow$  small  $V_r$   $\rightarrow$  large emb space

2.  $V_r \cdot E_i$

$V_r \times (V_d \times E_i)$

$1.2 \times 10^4 \times 1.2 \times 10^2 \sim 1.5 \times 10^6$

$\therefore (w_{v_l} + w_{v_r} + w_k + w_g) = 6 \times 10^6$  params.

Basically.

$V_r \left( \sum Prob \cdot (V_d \times E_i) \right)$

$\uparrow$   
calc in smaller space.

cross attention:  $K$  and  $Q$  act on 2 diff datasets eg:  $\mathbb{Z}$  translation  $\uparrow$  lang to another.  
here we again do not meet as we are trying to understand the rel<sup>n</sup> b/w 2 languages.

here  $K$  &  $Q$  would tell pretty which word 'i' in lang 1 corresponds to word 'j' in lang 2.

### \* Multiple Attention heads:

We run this calculation to find  $\Delta E_i$  MANY times independently  $n$  in parallel.  
 $\Delta E_i = \Delta E_{i1} + \Delta E_{i2} + \dots$   
each head.

Diff heads can specialise on diff aspects of self attention:

- 1  $\rightarrow$  Syntactic roles (noun-verb-align)
- 2  $\rightarrow$  local word collaboration ("New York")
- 3  $\rightarrow$  - - -

How to encourage specialisation in diff heads?

① Random initialization

② Loss func is common so if 2/more heads do the same thing  $\rightarrow$  Redundant  $\rightarrow$  loss fn pushes them to diff things.



in papers/stuff, the  $V_{ij}$  matrices (from all heads) are often all concatenated into one matrix

$$\begin{matrix} 128 \times \text{no. of heads} \\ \hline 12000 \end{matrix}$$

$V_{ij}$

GPT-3

weights devoted to Attention!!!

n-layers  $\times$  n-heads  $\times$   $\left[ \begin{matrix} d_{\text{emb}} * d_{\text{query}} \leftarrow k \\ d_{\text{emb}} * d_{\text{query}} \leftarrow q \\ d_{\text{value}} * d_{\text{emb}} \leftarrow v \\ d_{\text{emb}} * d_{\text{value}} \leftarrow v \end{matrix} \right] \approx$

$\uparrow$

$\uparrow$

$\uparrow$

$10^2$

$10^2$

$4 \times 1.5 \times 10^6$

$\approx$

$6 \times 10^{10}$

$\approx 60$  billion

$\sim \frac{1}{3} \times 180$  billion

total ~~per~~ perennae



How do LLMs store facts?

Each MLP layer adds information to the words.

Johnson - Linderstrauss lemma:

No. of vectors that you can cram into an  $n$ -dim  
space  $\approx$   $2^n$  exponentially with  
n.

not intuitive but it allows LLMs to  
have independent ideas - ALMOST  $2^n$  to  
even other are have exponential increase.  
thus allowing so many more ide

$\therefore$  models scale so well with size

this is one reason

$\therefore$  diff rows  $n$  columns in MLP weight  
matrices correspond to ideas.