

Word Embedding (p-3)

- A LLM can only process a fixed number of tokens at a time this is called context size. GPT = 2048 max

Ex: The smartest person is 4 vectors/tokens: Context size = 4

- This means as you go on and on with a conversation with a LLM it slowly starts to lose the context and thread of the conversation if you continue too long as it has only so much context from before (prev text)

* - An embedding

- at the end we have a probability distribution of all possible next tokens. We use the last vector for this

Ex The smartest person is $\begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}$ \rightarrow $\begin{cases} \text{Einstein 76.1} \\ \text{The 3rd} \\ \text{World 20.1} \\ \vdots \end{cases}$

- This involves 2 steps

- start at random
- learns in training

2) We take an unembedding matrix W that maps the very last vector in the context to a list of 50k values one for each token in vocab

$$2) \text{ [unembedding Matrix } W_u \cdot \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 50 \\ 6.4 \\ -30 \\ 54 \\ \vdots \\ -20 \end{bmatrix} \begin{matrix} 1 \\ \text{about} \\ \text{blue} \\ \text{boy} \\ \vdots \end{matrix}$$

2) Then we use a softmax function to normalize it for all probabilities $a_i \geq 0$ upto 1 (probability distribution)

2) \rightarrow Softmax \Rightarrow $\begin{cases} \text{cinstar} & 0.78 \\ \text{the} & 0.30 \\ \text{would} & 0.20 \\ \vdots & \end{cases}$

? What about the other vectors in the final layer why only the last one?
They all have context and meaning so why not use them?

- This is because in the training process the $\begin{bmatrix} \text{smartest } 5.1 \\ \text{one } 2.1 \\ \text{etc } 0 \end{bmatrix}$ the smartest $\begin{bmatrix} \text{person } 2.1 \\ \text{one } 5.1 \end{bmatrix} \dots$
 it's more efficient if you use each vector
 in the final layer to simultaneously make
 a prediction for what would come
 immediately after it. $\begin{bmatrix} \text{the} \\ \begin{bmatrix} : \\ : \\ : \end{bmatrix} \end{bmatrix} \begin{bmatrix} \text{smartest} \\ \begin{bmatrix} : \\ : \\ : \end{bmatrix} \end{bmatrix} \dots \begin{bmatrix} \text{is} \\ \begin{bmatrix} : \\ : \\ : \end{bmatrix} \end{bmatrix}$