

Word Embeddings (LLM NOTES)

Mapping notes - P111
Positional Embeddings - P112
Vector DB's - P95

- in reality tokenization frequently divides words, not whole words every time but we will focus on it by thinking tokens are clean whole words

Truth: This process (known famously Lie: this process (known famously

token \leftarrow and same for output token it can be any token not always a full word

All words $\approx 50k = 50k$ tokens

- The model has a pre defined Vocabulary, some list of all possible tokens (words, symbols, spaces, number etc)

and the first matrix we will see is the embedding matrix has a

single column for each of these words these columns determine what vector each word get in that first step

We label W_E = Embedding Matrix.

1	2	3	4	5	6	7	8	9
1	2.1	5	7	...	2	5		
2	2	6	9	...	7	7		
3	3	7	9	...	6	8		
4	4	8	1	...	8	9		
5	5	9	2	...	7	9		

Embedding Matrix W_E

the sky was blue
 $\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} 5 \\ 6 \\ 7 \end{bmatrix}$ EM picks the vector for blue

* Values begin random and learned through data *

- We embed these words as tokens which means its a vector in some very high dimension space like

12,288 dimensions in GPT-3. So

we visualize in 3D, as it learn these embedding in training it settles on

embedding's whose direction have meaning

for Ex paris and France's embedding are

close as they're related but up and down for ex are opposite in meaning and space.

in the Ex 2 we see how Embedding are related and why using vectors is good

say we did not know what a Female leader is we could do King + (Woman - man) and finding the vector (embedding) closest to that. $E(\text{Queen}) \approx E(\text{King}) + E(\text{woman}) - E(\text{man})$

- this is a perfect Ex in reality they might be further apart and are impossible to visualize, they are learned in training.

- further more "Queen" in training data would not always be a female leader

Ex "Queen the boat" hence the 50 embedding for Queen won't be the gives formula but close still.