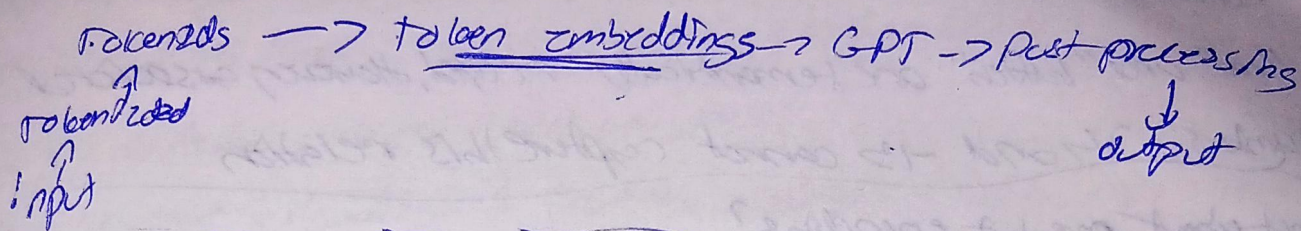


Token Embeddings



① conceptual understanding of why token embeddings are needed

② Demo

③ How are token embeddings created for LMs

↳ (a) Initialize embedding weights w/ random values

(b) This initialization serves as the starting point for the LM learning process

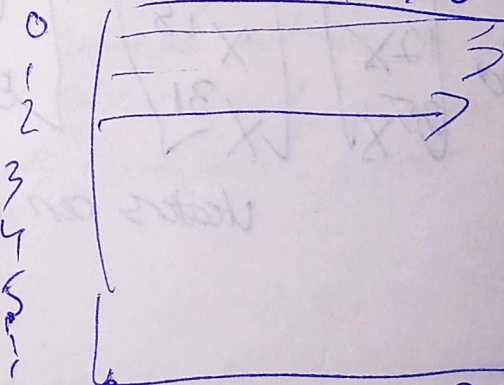
(c) The embedding weights are optimized as part of the LM training process

(d) Vocabulary (usually sorted alphabetically)

brown $\rightarrow 0$
dog $\rightarrow 1$
fox $\rightarrow 2$
jump $\rightarrow 3$
lazy $\rightarrow 4$
love $\rightarrow 5$
quack $\rightarrow 6$
she $\rightarrow 7$

Token IDs are
converted into
embedding
vectors

Token IDs vector Dimension
 $\rightarrow 7 \times 768$



Embedding layer
weight matrix

This is optimized during

(50257, 768)
Dimensions

What are token embeddings?

Often to represent words as #'s?

"cat" $\rightarrow 34$

"dog" $\rightarrow 2, 2$, etc

"boston" $\rightarrow -20$

"london" $\rightarrow -15$

"cat" and "london" are semantically related, however associated

Numbers 34 and -15 cannot capture this relation

What about one-hot encodings?

1 create dictionary of words

2 assign sequential one-hot encoding to each word

"dog" \rightarrow dict $\rightarrow [0, 0, 1, 0, 0, 0, \dots]$

"boston" \rightarrow dict $\rightarrow [0, 0, \dots, 1, 0, \dots]$

Also fails to capture semantic relationship.

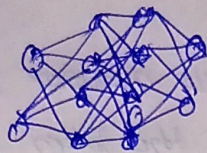
Semantically similar words should have similar vectors

	"dog"	"cat"	"apple"	"banana"	
boston	$\times 13$	$\times 31$	1	2	apple & banana dog & cat
edible	2	3	$\times 22$	$\times 38$	
4 legs?	$\times 18$	$\times 21$	0	0	
makes sound	$\times 12$	$\times 18$	0.5	0.2	
pet	$\times 35$	$\times 31$	5	7	

Vectors can capture semantic meaning.

Use an Artificial Neural Network to create vector embeddings

kernel
"cat"
"apple"
"banana"



$[23, 2, 19, 12, 35]$
 $[31, 3, 21, 18, 15]$
 $[1, 22, 0, 0.5, 5]$
 $[2, 38, 0, 0.2, 7]$

(6) Embedding layer is a lookup operation that retrieves rows from the embedding layer weight matrix using a token ID

(8) What an embedding layer actually does?

* Suppose we have following 3 training examples

idx = torch.tensor([2, 3, 1])

* Embedding dimension = 5

* Embedding = torch.embedding(4, 5)

* Embedding weights

$\begin{bmatrix} 0 \\ 2 \\ 3 \end{bmatrix}$ ----- (8)

* Embedding (7x8) (only 2, 3, 1)

one-hot $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$

\Rightarrow Some CS neural network linear layer

Each input (3) has 5 dimensions, so the

$\begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}$



$\Rightarrow \text{output} = X \cdot W^T =$

$\begin{bmatrix} 3 & 3 & 5 \end{bmatrix}$

Since as what we set when we need embeddings

(4x5) each column represents a neuron (layer)

$W^T = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$

linear

* Both embedding layer and linear layer need for same output

* Embedding layer is much more computationally efficient, since linear layer has many unnecessary multiplications $\frac{1}{2}$ of the size of one-hot embeddings