

# April 19, Transformers

Tokens - Patches of sound / text / picture

↓

Vectors → Pass through Attention. ①

- talk to other vectors.
- update their values

e.g. Model - Language Model  
Fusion Model

update it's meaning (vector)

## ② feedforward

- don't talk to each other (same operation in parallel)
- kinda like asking question, each vector, updating them based on answer to those question of the prev token

③ repeats. → until end, essential meanings had been baked into the "last token" of the sequence

sound natural?

apple, banana, and bread are ???

↓ we kept doing, new token as seed to repeat calculation

system prompt

↳ You're an apple.

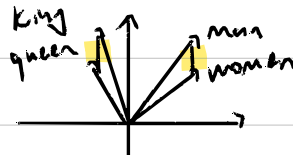
Prediction

↳ User: Give me a bunnn

↳ Apple: ???

## Tokenizing (vector)

- train tends to settle vector in a space, where direction have a semantic meaning

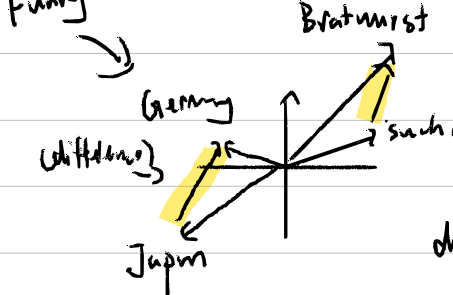


$$E(\text{queen}) - E(\text{king}) = E(\text{woman}) - E(\text{man})$$

$$E(\text{queen}) \approx E(\text{king}) + E(\text{woman}) - E(\text{man})$$

e.g. demonstrate: demonstrate → garden

$$E(\text{Hilary}) + E(\text{Germany}) \approx E(\text{Mussolini}) + E(\text{Italy})$$



dot product:

$$\langle \text{plur w/ dogs} \rangle \cdot \text{dot prod}$$

$$\langle \text{plur w/ cats} \rangle \cdot \text{dot prod}$$

✓ (larger)

$$\langle \text{plur} \rangle \cdot \text{w/ dogs} \cdot \text{dot prod.}$$

$$\langle \text{plur} \rangle \cdot \text{w/ cats} \cdot \text{dot prod.}$$

Fun: You can find how plural a word is

\* embeddings[token] = vector

N words/tokens

$$\begin{bmatrix} \text{apple} & \text{banana} & \text{cats} & \text{dogs} \\ +8.2 & -4.2 & .2 & .8 & -8.6 \\ +7.2 & +6.4 & -1 & 7 & -9.2 \\ -8.1 & +3.8 & -4 & 4 & +1.8 \\ -8.7 & -1.7 & +1.2 & 8 & 8 \end{bmatrix} = W_E$$

only for Embedding)

Parameter = dimension

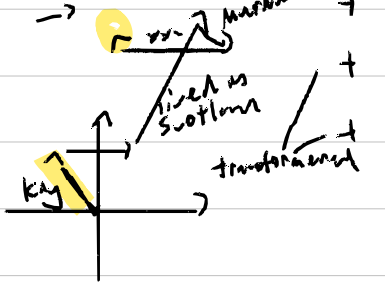
N words

(Embedding matrix)

\* vector are updated based on context

\* In transformer, these aren't just representing

new king



+ "meaning" of these word

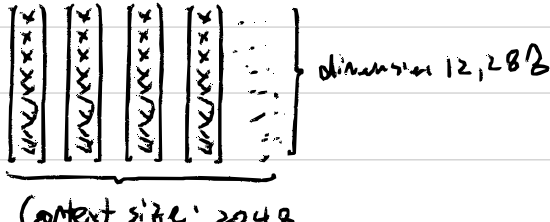
+ positions of that word

+ context

• the direction changes by other "blocks" of the sequence, the resulting pointing into a more specific, nuanced direction

\* think how you shaped a word. in a sentence.

Context: how many token can it process (how much can it incorporate, when predict or next word)



What happens in very end? : Unembedding

Harry potter, least favorite, Professor → How

①: for the last vector

- use a matrix, match the last vector, to a list, of 50,000 vocabulary

→ normalize it with softmax

why only use last embedding to predict?

②: why other vectors are just sitting there, with it's own context rich meaning?

• more efficient to use each other embeddings, to predict what goes right after it.

• last token is the summarization of previous tokens.

$$\begin{bmatrix} x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \\ x & x & x & x & x \end{bmatrix} \begin{bmatrix} x \\ x \end{bmatrix} \rightarrow \begin{bmatrix} 200 \\ -8 \\ 1 \\ 6 \\ -9 \end{bmatrix} \begin{matrix} \text{good} \\ \text{apple} \\ \text{cat} \\ \text{new} \\ \text{cut} \end{matrix} \rightarrow \text{softmax}$$

Unembedding Matrix

$W_u$

(vectors)

$W_u = \text{word} \cdot \text{dimension}$

(summed of embedding)

\* unembedding is a classifier.

use the query to find dot product, which is most similar.

• the output is a semantic (mean, you need to find the best fit for the meaning)

• Vector is meaning,  $W_u$  is a mapper to possibilities, to tokens (words)

## Softmax

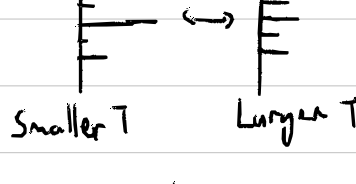
• Normalizing it, add up to one

• Temperature: throwing a T in the Softmax process

→ it's like temperature, because it resembles temperature in thermodynamic equations;

• T is larger, more weight to lower value

→ distribution is more uniform



e.g. • lower temperature: always most predictable word

• higher temperature: chance to choose less likely words

True derivative of Goldilocks

Logits - Softmax → Probabilities

$$+6.0 \quad e^{x/b / \sum_{n=0}^{N-1} e^{x_n/T}}$$

0.19

$$-5.0 \quad e^{x/b / \sum_{n=0}^{N-1} e^{x_n/T}}$$

0.04

$$+4.0 \quad e^{x/b / \sum_{n=0}^{N-1} e^{x_n/T}}$$

0.14

$$+1.5 \quad e^{x/b / \sum_{n=0}^{N-1} e^{x_n/T}}$$

0.10

$$+9.9 \quad e^{x/b / \sum_{n=0}^{N-1} e^{x_n/T}}$$

0.35

• logits: raw, unnormalized output



April 24th.

Attention in Transformers, DLB

VECTORS: Semantics in space

Transformers: through adjusting the embeddings.

don't merely encode each individual words (vector)

but take them in some more richer contextual

Meanings

eg.

- American shrew mole
- One mole of carbon dioxide
- bispy of a mole

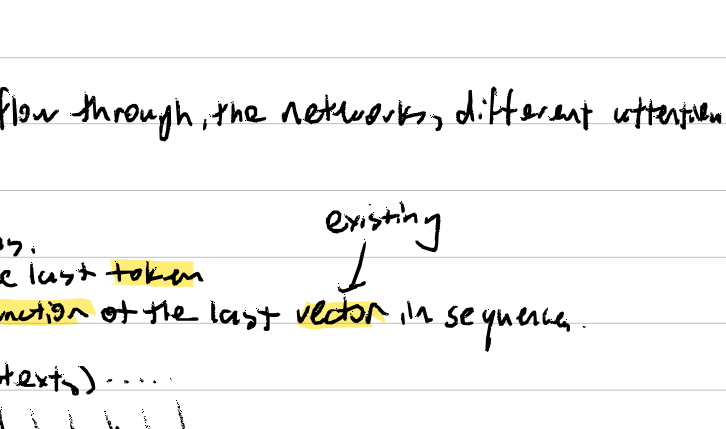
1. After first step of transformers, the "mole" would mean the **sum**.

Because first step of token-embeddings is just a lookup table

2. It's the next step of Transformers, the surrounding

embedding have the chances to pass the contextual

information to "mole"



Other than just retaining a meaning of a word, the attention block allows to move information of one embedding, to another

- information ferrying
- information richer than just a word

After all the vectors flow through, the networks, different attention block:

- the computations: predict the last token
- ... a function of the last vector in a sequence.

(Context): ...

Therefore, the modeler was ???

with all previous information

this last token through attention computation, will include all the information before it, that's reason to predict the (Contextual) next word

"attention head" Example

"a fluffy blue creature roamed the verdant forest"

Embeddings for each word:  $E_1, E_2, E_3, E_4, E_5, E_6, E_7, E_8$

Color, position of the word is also updating the vectors

Updated "creature"

1:  $E_4$  (creature) ask, any word there any adj before me?

$(E_2, E_3) \rightarrow 2m$

Query

if you learn it's "creature", what to ask?

Computing for this Query, matrix multiplying to embedding

$$W_Q \cdot E_4 = Q_4 \leftarrow \text{Querying query}$$

$$Q_4 = \begin{bmatrix} 3 & 2 & 2 \\ 2 & 2 & 2 \\ 2 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ 2 \end{bmatrix}$$

Any adjectives before me?

if you learn it's "creature", what to ask?

"creature"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"

"a"

"fluffy"

"blue"

"creature"

"banned"