

(1)

StatsQuesA - Transformer NNs

transformer is a type of nn

nn's need numbers as inputs so start w/ word embedding
 $= \text{number} \rightarrow \text{for word}$

embed networks are the same for each word

Positional encoding \rightarrow retain order context

- take word embedding and apply transformation to get positional encode
 \hookrightarrow usually cosine etc

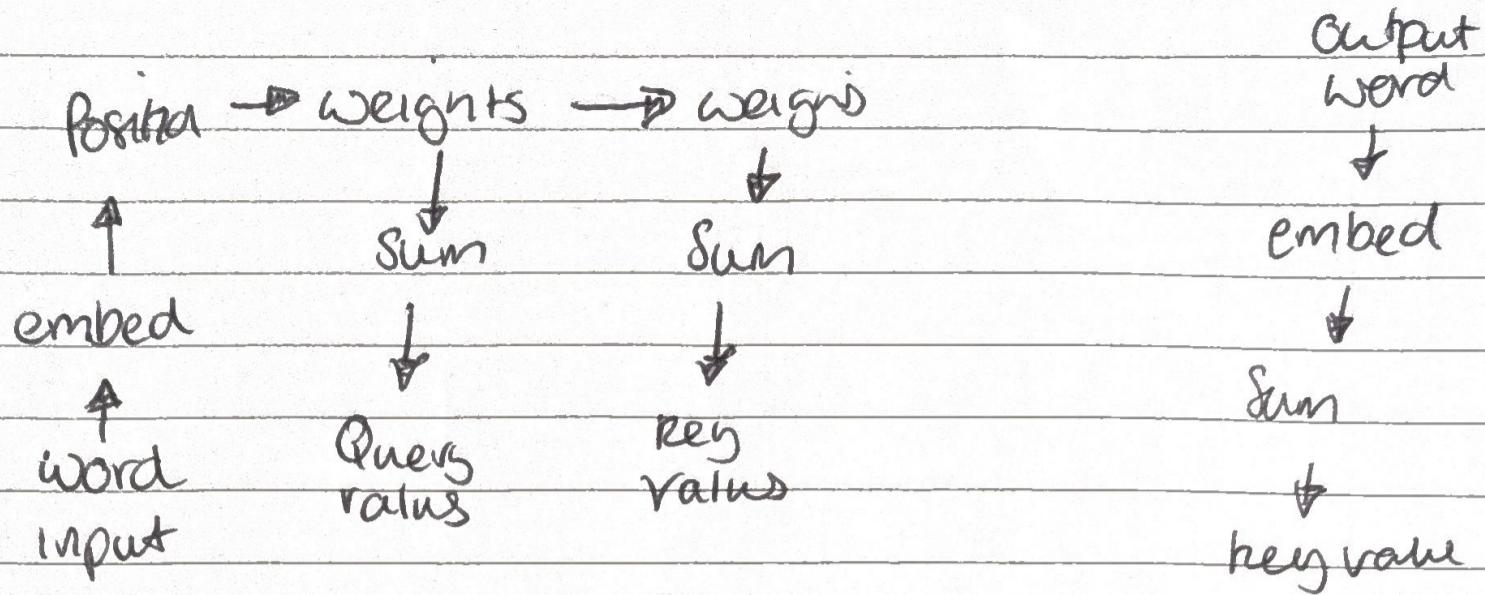
this is how a transformer keeps track of word positions

How do transformers keep track of the relationships between words?

Self-Attention

- \rightarrow keeps track of a word's similarity to all other words in a sentence including itself

(2)



use Key values to calc similarities w/ Query
 (have use dot prod for sim)

$$① \text{Query} \times \text{Input key} = \text{Sim}$$

$$② \text{Query} \times \text{Output key} = \text{Sim}$$

large = more similar

→ has more influence on encoding
 using softmax (0, 1) to harness word value

Decide the words similar to a word to influence its encoding

After a word in sentence has been compared to all other words

$$\text{again} (\text{Pos rals} \times \text{weigh}) \times \text{sum} = \text{Value}$$

value * softmax

do for input side & output & scale

(3)

Result = Self-Attention values for words s_i from sentence

Repeat Self-Attention Scores for each word

can re-use Query & value weights

can calc all words @ once, does need to be iterative

word \rightarrow embed \rightarrow Position \rightarrow Self-Attention

\hookrightarrow these are the progression of values that represent a word

Self-Attention values improve as they contain input from all other words

\hookrightarrow this enables words to have context

this helps establish how each word in the input is related to the others

\rightarrow Self Attention Cell =

- think of are w/ weights calc'ng the Queries, keys & values of cell

layer \rightarrow stack cells to capture more relationships in complicated inputs

- cells start w/ positional encoded values

in the first transformer paper they stacked 8 self-Attention cells

↳ multi-head attention

in the word value flow take the:

- Pos encoded values
- Self-Attention values
- & sum together

these are called residual connections

Allows self-Attention to focus on relationships

→ Position encodes are added back

the residual values are the encoded inputs for the transformer

4 stages - embed, Position, Self-Attn, residual

- encode into nums
- encode positions
- encode relationships
- easily & quickly train in parallel

→ there are lots more complex things that can be added to a transf

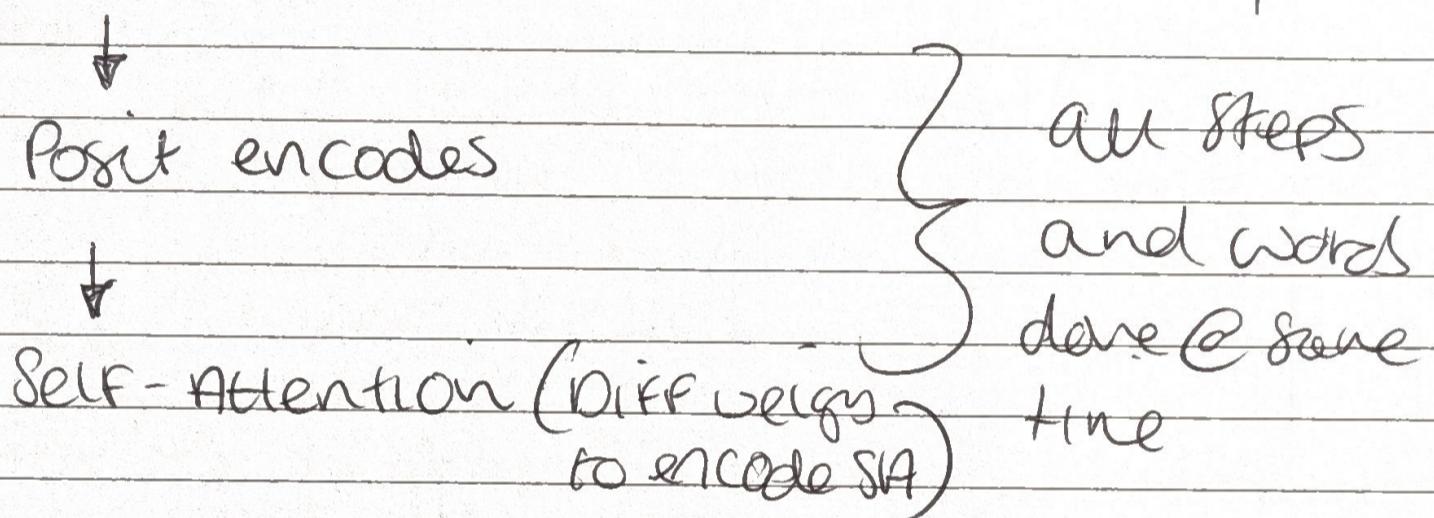
(5)

transformers are an encoder-decoder architecture

Previous has all been encoders into the resid connections

Decoder

Starts w/ input embeds - Start w/ eos



Previous looked @ self-att relations between sentence words

But also need to track between input & output (Attention)

Decoder needs to keep track of significant words in the input
(not lose any context)

"encoder-Decoder Attention"

- take word from decode & take Query val
- Calc key values from all words in encoder
- Calc all similarities (dot prod)
- Softmax to decide weight influence

Output = encoder-Decoder Attention value
for a output word

E-D Attention has diff weights to
self att

then calc a resid value from
the self-Attention + E-D Attention

Resid)

- E-D Attention - input/output relation
- Self-Att - intra sentence relation
- Positional encode
- word embed

finally resids go into a fully
connected layer to predict the
decoded word from vocab

repeat until Decoder outputs EOS

TRANSFORMERS

- ① Embed
- ② Position
- ③ Self-Attention
- ④ Encode-Decode Attention
- ⑤ Residual to chain together