

## ① Stat Quest - Attention for NNs

Problem w/ basic en-dec is that LSTM turns sentence into a single context vec

this is fine for short phrases

but is input vocab of 1000+ words & possibility of longer, more complex sentence

Here, early words in the LSTM can be forgotten

early words might be v important for the context of the sentence

Recall, RNNs had a similar issue w/ remembering things long back

Main idea of LSTM was to solve through separate path for long & short term memories

but even w/ separate paths, if the paths are long words can get lost

Attention adds a bunch of new paths from the encoder to the decoder

decoder can directly access steps from the encoder

(2)

this video covers a combination of  
 Seq2Seq combined w/ Attention

Encoder-Decoder

↳ this is an a precursor to  
 transformers

how does attention connect the inputs  
 to each step of the decoder

this route is not simple

there are no rules about how  
 Attention should be added to  
 an encoder-Decoder model

↳ only conventions

Possible Steps

- ① Determin how similar the outputs  
 from the encoder LSTM @ each step  
 to the outputs of the decoder LSTM

Similarity Score between "words" (Seq)

Attention Algs have DIFF ways of  
 doing this but Cosine Similarity  
 is common

Cosine Similarity =

Calcs Sim between two seqs of nums  
between  $-1 \leq \text{Sim} \leq 1$

Encode

LSTM<sub>1</sub> → LSTM<sub>2</sub>

Decode

LSTM<sub>1</sub>

Sim Score

Sim score

LSTM<sub>1</sub> output from Decode gets a sim score w/ each step of the Encode LSTM<sub>2</sub>

→ Plugged into cosine formula

Cosine Sim can be replaced w/ a dot product to improve computation speed w/ similar affect

Qn how are the scores used?

Higher score = More influence on Decoder Step

Run score through softmax

→ recall soft max =  $\frac{\exp(x_i)}{\sum \exp(x_j)}$

All add upto 1

→ Result = % of encode word to use in influence

(4)

⇒ Scale input words by their softmax

Sum all scaled inputs into single val

↳ these are the Attention values  
for the decode LSTM, input +  
token/word

to determine the first pred word

⇒ Plug all of the Attention values  
from each layer into a Fully  
Connected layer w/ softmax

repeat until EOS is predicted

- ① unroll LSTM
- ② Comp Attention
- ③ Fully connect

⇒ Attention in encoder - Decoder

↳ the encoder is much the same

but decoder, each step has access  
to individual encodings for each  
input word

→ use softmax & similar score to determine  
what % of each encoded word should  
be used to help predict the next  
output word

(5)

A finding of attention is that we no longer need the LSTMs

↳ but this is an advancement shown in transformers