

StatsQuest: Decoder Only Transformations

Masked self-Attention words contain input from all of the other words that came earlier \rightarrow Context

{ masked Attention cells }
Applied to position encode value

\hookrightarrow into residual connections

these residual numbers represent ^{a word} ~~words~~

We want to use these numbers to predict the next word

Decode-only transformer

We want to build something that can both encode the prompt + generate output

word \rightarrow embed \rightarrow Mask self-att \rightarrow residual \rightarrow MLP \rightarrow ~~word~~

\hookrightarrow MLP predicts the next word

- Possible outputs could be anything in the vocab

\triangleright in the original paper, the MLP is the word embed network & weights in reverse (not all do this way)

Diff between transformer & Decode only

Decode only uses the same components to encode & generate output

Same cells, weights & biases

Masked Self-Attention is calculate to each current word & everything preceding it

Self-Attention is applied equally to the input prompt & output generates

Mask Self-Att a Decoder-only to determine how words in the prompt are related

& keep track of important input words when generating the output

Basic transformers use:

- encoder to learn prompt
- Decoder to generate output
- uses SELF-ATT, not masked SELF-ATT
 - ↳ learns how all words in prompt are related
 - ↳ Masked just uses words before

Normal transformer uses masked in the Decoder only