

Lecture 23: Transformers

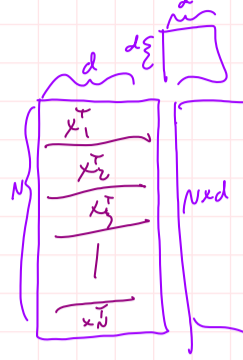
Announcement:

* project info up tonight.

Goal: network which can handle sequence data

3 answers: conv, RNN, attention

Convolutional (with pre-extracted patches)

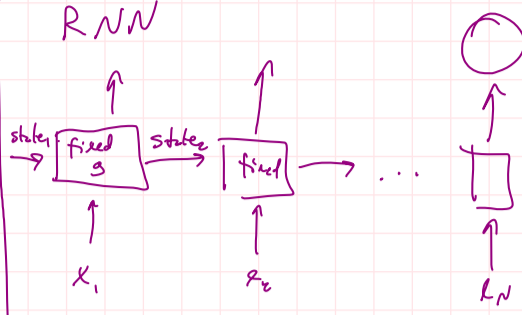


Comparison against linear layer with Nd

- * "less powerful" than linear
- * few params (seem to suffice)
- * can learn/gen to other sequence lengths

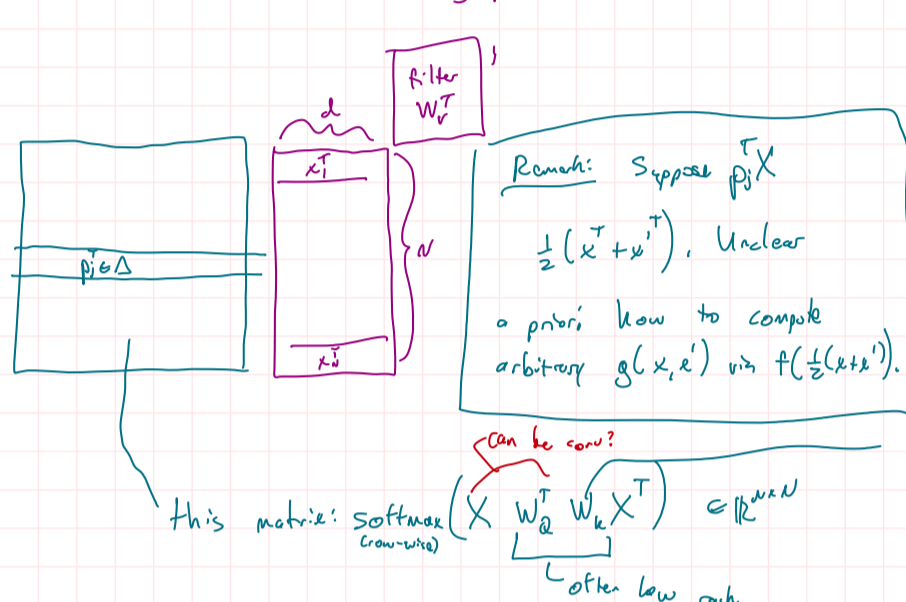
* computing a function with interacting patches requires multiple layers & overlapping patches (e.g., 14x14 go board, 3x3 conv filter).

RNN



- * if state = ϕ , similar to convolution
- * unlike conv, can model patch interaction, but state limits how far back interactions may be calculated ("fixed-size memory")
- * persevering state also makes training finicky. (e.g., vanishing/exploding gradients.)

Attention (Vaswani et al. '17)
convolution + fix to allow interacting patches



Remark (naming):
 $W_{value}, W_{key}, W_{query}$ names are from NPM (neural Turing machine).

Remark: Suppose $p_i^T X$
 $\frac{1}{2}(x_i^T + x_i'^T)$, unclear a priori how to compute arbitrary $g(x_i, x_i')$ via $f(\frac{1}{2}(x_i + x_i'))$.

can be conv? $(X W_q^T W_k X^T) \in \mathbb{R}^{N \times N}$ $W_k \in \mathbb{R}^{d \times d}$
often low rank

Remark (computation). Naive implementation does dN^2 scalar multiplications (applied community fights over this.)

Bells & whistles:

- * masking: elementwise product of $\text{softmax}(\dots)$ with a fixed $\sum_{0,1,2,3,4,5,6,7,8,9}$ ("causal")
- * "positional encoding": note that attention is permutation invariant wrt sequence structure.

Example (Vaswani et al.)

$$\begin{bmatrix} x_{i-1} \\ | \\ x_{i+1} \end{bmatrix} \mapsto \begin{bmatrix} x_{i-1} + \cos\left(\frac{i}{10000}\right) \\ x_i + \sin\left(\frac{i}{10000}\right) \\ x_{i+1} + \cos\left(\frac{i+2}{10000}\right) \\ \vdots \end{bmatrix}$$

AlphaFold paper

$$x_i \mapsto x_i + te_i \text{ assuming } d \geq N$$

- * self-attention vs attention
- * "multihead": multiple parallel attention, concatenate outputs (get a higher-order tensor), apply ReLU network.

Theorems

* (Wei-Ma-et al.) Can simulate T-step Turing machine with depth T and $W_k, W_v \in \mathbb{R}^{\log T \times \log T}$

* (Yao-Peng-Papadimitriou-Narasimhan)

Define Dyck_{k,0} language: k types of parentheses, D nestings
e.g. $k=3, D=4$ $\frac{1}{2} \left(\begin{matrix} 3 \\ 2 \\ 1 \end{matrix} \right) \left(\begin{matrix} 3 \\ 2 \\ 1 \end{matrix} \right) \left(\begin{matrix} 3 \\ 2 \\ 1 \end{matrix} \right) \left(\begin{matrix} 3 \\ 2 \\ 1 \end{matrix} \right)$ (abstracts parse tree)

"Theorem" Can identify & generate Dyck_{k,0} with depth D and all inner dimensions poly(k).

Proof idea - layer i identifies matching delimiters at depth D-i+1
IH: either higher depths are not in this language OR all higher depth delimiters are marked.

Inductive step: Scan right from each delimiter, skipping marked symbols, & identify matching bracket OR determine not in language

Negative result: Dyck_k is not possible

Large language models

Best { predict occluded words given two sentences, does one follow the other

GPT-3 = next word prediction

Question: why are these valuable "upstream tasks"