

Idea: Can we transform RNN to purely attention-based model?

→ but, how do we overcome temporal dependencies?

essentially a layer.

everything processed w.r.t no assumed order.

Self-Attention: For each x_i , encode h_i , then $v_i = v(h_i)$,
For every q_i and k_j , compute attention scores: (v_{ij}) and $k_i = k(h_i)$ and $q_i = q(h_i)$
and respective softmax $a_i = \frac{\sum_j a_{ij} v_j}{\sum_j a_{ij}}$ ← multiply by each value.

learned (W_v)

learned (W_k)

learned (W_q)

"value"

"key"

"query"

activation, this can be inputted into another s-a layer!

Transformer: involves self-attention, positional encoding (addresses sequence info), multi-headed attention (query multiple positions), adding nonlinearities, masked decoding (prevent attention lookups in the future).

→ position of words carries information, so add some info to x_t to indicate that.

$h_t = f(x_t, t)$ relative encoding → absolute encoding.

Frequency-based representations for positional encoding

$$PE(pos, 2i) = \sin(pos/10000^{2i/d})$$

$$PE(pos, 2i+1) = \cos(pos/10000^{2i/d})$$

from paper

oscillate btw -1/1

OR, we could learn this $d \times T$ matrix. More flexible/optimal. More complex.

Input to self-attention (from paper) = $emb(x_t) + p_t$.

Multi-Head Attention: ideally we'd want to incorporate more than 1 time step.
→ have multiple $k/q/v$ for every time step. Called heads, learned independently.

Naive self-attention lacks "expression" of picked information. Fully linear.

→ alternate self-attention and nonlinearity

$$e.g. h_t^l = \sigma(W^l a_t^l + b^l) \leftarrow \text{non-linear, learned. Basically NN layer. output from prev. self-attention layer.}$$

Big issue: self-attention can "see" the future. Not inherently sequential.

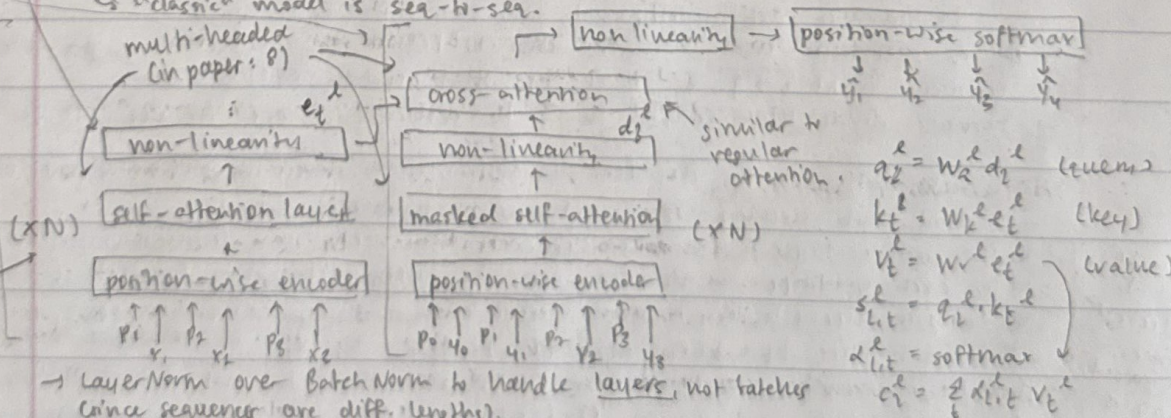
→ Masked attention → does not allow lookups past t .

$$s_{ij} = \begin{cases} q_i \cdot k_j & \text{if } i \leq j \\ -\infty & \text{otherwise (becomes 0 in softmax).} \end{cases}$$

"Transformer" → transforms one sequence into another at each layer

→ "classic" model is seq-to-seq.

decoder block (processors) y_i shifted right
encoder block (processors) x
in paper, $N=6$.



→ LayerNorm over BatchNorm to handle layers, not batches (since sequences are diff. lengths).

μ, σ are scalars, aggregate over different dimensions of a , not all a .
→ straightforward to compute!

→ implement LN + skip connection

$$\text{input: } h_t^{l-1} \rightarrow \bar{a}_t^l = \text{LayerNorm}(h_t^{l-1} + a_t^l)$$

"modify", not change, maintain well-behaved gradients

Benefits of Transformer:

- long-range connections
- much easier to parallelize
- adaptive, can be made deeper.
- practical involvement