

## Assignment 5

Vinta Redhu

ES18BTECH11028

- 1.
- a) Given that,  $t$  = sequence length  
 $l$  = number of layers  
 $n$  = number of neurons

Time complexity:

$$\text{Train time} = t n^2 l$$

$$\text{RNN - Test time} = t n^2 l$$

$$\text{Transformer - Train time} = t^2 n l$$

$$\text{Test time} = t^2 n l$$

Space complexity:

$$\text{RNN - Train time} = t n l$$

$$\text{Test time} = n l$$

$$\text{Transformer - Train time} = t n l$$

$$\text{Test time} = t n l$$

- b) • The time taken for computation is  $t^2 n l$  in transformers and  $t n^2 l$  in RNN.
- If the value of  $n$  is smaller than sequence length  $t$ , time taken by RNN will be less than time taken by transformer.
- Transformers in this case perform poorly because the computations at self attention layer are more than the normal feed forward network.
- c) • Self attention layer looking across the ~~all~~ tokens of a given input sequence is bottleneck for parallelism.
- Generally, the sequential operations in transformers are independent of sequence length. But these are expensive to decode.
- Parallel processing is what that makes transformer learn faster than RNN for longer sequences.

- d) • The feedforward network and layer norm do not look across the tokens.  
 • They only look at output content vector of self attention layer  
 • In this way parallelism is introduced as feedforward network work in parallel.

$$z = \sum_{i=1}^m (v_i \alpha_i)$$

$$\alpha_i = \frac{\exp(k_i^T q)}{\sum_{i=1}^m \exp(k_i^T q)}$$

a) Given,

$$z = v_j$$

which means that  $\alpha_j = 1$

and  $\alpha_i = 0 \forall i \neq j$

This means that  $k_j^T q \gg k_{i, i \neq j}^T q$

b) Given,  $\{k_1, k_2, \dots, k_m\} \rightarrow k_i \perp k_j \forall i \neq j$  and  $\|k_i\| = 1 \forall i$

$$\text{Let } q = \sum_{j=1}^m B_j k_j$$

$$\begin{aligned} \text{Let's evaluate } k_i^T q &= k_i^T \left( \sum_j B_j k_j \right) = \sum_j B_j k_i^T k_j \\ &= B_i \underbrace{\|k_i\|^2}_1 + \underbrace{\sum_{j \neq i} B_j k_i^T k_j}_0 \quad \text{as } k_i^T k_j = 0 \forall i \neq j \\ &= B_i \end{aligned}$$

$$\alpha_i = \frac{\exp(k_i^T q)}{\sum_{i=1}^m \exp(k_i^T q)} = \frac{\exp(B_i)}{\sum_{i=1}^m \exp(B_i)}$$

$$z \approx \frac{1}{2} (v_a + v_b) \Rightarrow \alpha_a = \alpha_b = \frac{1}{2}$$

and  $\alpha_i = 0 \forall i \neq a$  and  $i \neq b$

• we can obtain  $\alpha_a \approx \frac{1}{2}$ , by setting  $B_a = B_b \gg 0$  &  $B_i < 0 \forall i \neq a \neq b$

is the query vector  $q$   
and

$$\alpha_i \approx 0 \forall i \neq a \text{ and } i \neq b$$

3.

$$\begin{aligned}
 L(q) &= \int q\left(\frac{z}{n}\right) \log \frac{p(x, z)}{q\left(\frac{z}{n}\right)} dz \\
 &= \int q\left(\frac{z}{n}\right) \log \left[ \frac{p(x/z) p(z)}{q\left(\frac{z}{n}\right)} \right] dz \\
 &= E_{z \sim q\left(\frac{z}{n}\right)} \log \left[ \frac{p_0\left(\frac{x}{z}\right) p_0(z)}{q\left(\frac{z}{n}\right)} \right] \\
 &= E_{z \sim q\left(\frac{z}{n}\right)} \left[ -\log \frac{q\left(\frac{z}{n}\right)}{p_0(z)} + \log p_0\left(\frac{x}{z}\right) \right] \\
 &= E_{z \sim q\left(\frac{z}{n}\right)} \left[ \log p_0\left(\frac{x}{z}\right) \right] + E_{z \sim q\left(\frac{z}{n}\right)} \left[ -\log \frac{q\left(\frac{z}{n}\right)}{p_0(z)} \right] \\
 &= E_{z \sim q\left(\frac{z}{n}\right)} \left[ \log p_0\left(\frac{x}{z}\right) \right] - \underbrace{KL\left[q\left(\frac{z}{n}\right), p_0(z)\right]}_{\text{reconstruction loss}} \quad \text{regularisation term}
 \end{aligned}$$

- KLD measures similarity between two distributions
- reconstruction error here is the HLE of decoder.

4.

Given that,  $f(p, q) = pq$

$$\min_p \max_q f(p, q) = ?$$

a)

$K = (1, 1)$ ;  $n = 6$ ; step = 1

Minimizing w.r.t  $q_t$

$$\frac{\partial f}{\partial q_t} = \frac{\partial (p_t q_t)}{\partial q_t} = p_t$$

$$q_{t+1} = p_t + q_t \quad [\text{gradient update}] \quad \text{--- ①}$$

$$\text{Hence } f' = p_t q_{t+1}$$

Minimizing w.r.t  $p_t$  :-

$$\frac{\partial f'}{\partial p_t} = \frac{\partial (p_t q_{t+1})}{\partial p_t} = q_{t+1}$$

$$p_{t+1} = p_t - 1(q_{t+1}) \quad [\text{gradient update}]$$

↑  
substituting in Eq ①

4.

$$P_{t+1} = (q_{t+1} - q_t) - q_{t+1}$$

$$= -q_t$$

$$\therefore f_{t+1} = (-q_t)(P_t + q_t)$$

• Using the above to fill the table

$q_0$	$q_1$	$q_2$	$q_3$	$q_4$	$q_5$	$q_6$
1	2	1	-1	-2	-1	1
$P_0$	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$	$P_6$
1	-1	-2	-1	1	2	1

- b) • With the step=1, as we can see above there is no convergence. hence it is not possible to find the optimal value.
- To find the optimal value we can change the value of step size.
- c) • It attains equilibrium when the values remain same  
(f)

$$f_t = f_{t+1}$$

$$\left[ \begin{array}{l} P_t q_t = (-q_t)(P_t + q_t) \\ -P_t = P_t + q_t \end{array} \right] \quad \begin{array}{l} P_t q_t = (-q_t)(P_t + q_t) \\ q_t (2P_t + q_t) = 0 \end{array}$$

$$\therefore q_t = 0 \text{ or } P_t = -\frac{q_t}{2}$$

• Here at point 2,5  $P_t = -\frac{q_t}{2}$  which didnot lead to equilibrium.

• Hence  $q_t$  should be 0, which mean  $P_{t+1} = 0$ .

↓  
this or above condition should be satisfied to obtain equilibrium.