

Assignment - 05

K. Surya Prakash
EE18BT01ECH11026

- Q1) t : seq. length
 l : no. of layers
@ n : no. of neurons

For training:

Arch.	Time Complexity	Space Complexity
RNN	$O(t \cdot n^2 \cdot l)$	$O(t \cdot l \cdot n)$
Trans former	$O(t^2 \cdot n \cdot l)$	$O(t \cdot l \cdot n)$

For testing:

Arch	Time Complexity	Space Complexity
RNN	$t \cdot n^2 \cdot l$	$l \cdot n$
Transformer	$t^2 \cdot n \cdot l$	$t \cdot n \cdot l$

1b)

If $n < t$

\Rightarrow Complexity for RNN: $O(tn^2l)$

Complexity for transformer: $O(t^2nl)$

\Rightarrow In this case:

RNN performs better than a transformer

This is because:

Since $t > n \Rightarrow O(tn^2l) < O(t^2nl)$

\Rightarrow More computation will be taken in self-attention layer itself in Transformer.

1c) Yes.

\rightarrow Self-attention layer considers all the input-sequence tokens & produces a matrix that learns the relation b/w each token.

\Rightarrow Although this might rule-out the sequential modeling. But, the current hidden layer depends on the

past hidden layer of every token, which makes it expensive in decoder.

→ Because, decoder uses masked-self attention. Thus computation in self-attention increases as sequence size increases.

⇒ Hence, this is a bottleneck for parallelism. But at the same time, this layer helps us do parallel computation in the later layers, since all the attention has been learnt earlier.

(d) No. Feed forward & layer norm. are ~~not~~ computed ~~only~~ based on the o/p of the self-attention layers (which look through tokens).

⇒ Thus these blocks become independent across time (token), and can be computed in parallel. ⇒ Parallelism

2Q)

$$a) \quad z = v_j \Rightarrow \alpha_i = 1 \quad \text{if } i=j \\ 0 \quad \text{else} \\ \Rightarrow \alpha_j = 1 / \sum_i \alpha_i$$

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

$$\Rightarrow \exp(k_j^T q) \approx \sum_{i=1}^m \exp(k_i^T q)$$

$$\text{ie; } \exp(k_j^T q) \gg \exp(\underbrace{k_x^T q}_{x \in \{1, \dots, m\} - \{j\}})$$

This tells us that:

query vector 'q' is aligned with

$$k_j \Rightarrow k_j^T q \gg k_x^T q \quad ; \quad x \in \{1, \dots, m\} - \{j\}$$

$$\Rightarrow \boxed{q \parallel k_j}$$

(b) given: $\{k_1, \dots, k_m\}$ are ~~orthogonal~~ orthonormal

Sol.

$$\Rightarrow k_i^T k_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

$$\Rightarrow \text{Given: } z = \frac{1}{2}(v_a + v_b)$$

$$z = \sum_{i=1}^m d_i v_i = \frac{1}{2}(v_a + v_b)$$

$$\Rightarrow \boxed{d_a = d_b = \frac{1}{2}} : d_i = 0 \quad \forall i \in \{1, \dots, m\} - \{a, b\}$$

This can happen if:

$$\text{Let } q = \sum_{i=1}^m c_i k_i$$

$$\Rightarrow k_i^T q = \underbrace{c_i (k_i^T k_i)}_{c_i(1)} + \underbrace{\sum_{\substack{j=1 \\ j \neq i}}^m c_j (k_i^T k_j)}_0$$

$$\therefore d_a = \frac{\exp(k_a^T q)}{\sum_{i=1}^m \exp(k_i^T q)} = \frac{\exp(c_a)}{\sum_{i=1}^m \exp(c_i)} = \frac{1}{2}$$

Similarly

$$d_b = \frac{\exp(c_b)}{\sum \exp(c_i)} = \frac{1}{2}$$

Hence: $\therefore \alpha = \beta = 1/2$

This can happen if

$$\boxed{\begin{array}{l} c_a = c_b \\ (c_a, c_b) \gg c_i \end{array}} \quad \forall i \in \{1, \dots, m\} - \{a, b\}$$

* A simple case is: if

$$\boxed{q = c(k_a + k_b)} \Rightarrow \begin{array}{l} c_a = c_b = c \\ c_i = 0; \forall i \notin \{a, b\} \end{array}$$

Q3)

$$= \mathcal{L}(q) = \int q(z|x) \log \left(\frac{p(x, z)}{q(z|x)} \right) dz$$

$$= \int q(z|x) \cdot \log \left(\frac{p(x|z) \cdot p(z)}{q(z|x)} \right) dz$$

$$= \int q(z|x) \log(p(x|z)) dz - \int q(z|x) \log \left(\frac{q(z|x)}{p(z)} \right) dz$$

$$= \mathbb{E}_{z \sim q(z|x)} (\log(p(x|z))) - \text{KL}(q(z|x) \parallel p(z))$$

$$\therefore KL(P||Q) = E_{x \sim P}(\log(P/Q)) = \int P(x) \log(P/Q(x)) dx$$

$$\mathcal{L}(q) = \underbrace{E_{z \sim q(z|x)}(\log(P(x|z)))}_{\text{Reconstruction}} - \underbrace{KL(q(z|x) || P(z))}_{\text{Regularisation}}$$

* Reconstruction: (Main loss term)

$E(\log(P(x|z)))$ amounts for MLE
 $z \sim q(z|x)$

\Rightarrow This is analogous with mse loss b/w
 actual x (img) & reconstructed $\text{img}(\tilde{x})$

* Regularisation:

\Rightarrow This pushes $q(z|x)$ to be as close as possible
 to $p(z)$ distribution which is a
 assumed to be gaussian with var: 1

Q4)

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

$$\text{Obj} := \min_p \max_q f$$

① Update steps (given p_t, q_t) $\Rightarrow f = p_t q_t$

(i) Obj maximises w.r.t q

\Rightarrow Require gradient ascent

$$\Rightarrow q_{t+1} = q_t + \frac{\partial f}{\partial q_t} = q_t + p_t$$

$$\boxed{q_{t+1} = q_t + p_t} \Rightarrow \boxed{f = p_t q_{t+1}}$$

(ii) Obj minimises w.r.t p

\Rightarrow Requires grad. descent

$$p_{t+1} = p_t - \frac{\partial f}{\partial p_t} = p_t - q_{t+1}$$

$$\boxed{p_{t+1} = p_t - q_{t+1}}$$

$$\boxed{f = p_{t+1} q_{t+1}}$$

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) No. In this setting, it's impossible to find a optimal value for f .

Reason:

From the table above, we can find that p, q values oscillate periodically, which do not lead to convergence.

(c) Equilibrium point:

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

$$p_{t+1} q_{t+1} = p_t q_t$$

$$\Rightarrow (p_t - q_{t+1})(q_{t+1}) = p_t q_t$$

$$\Rightarrow (p_t - q_t - p_t) \overset{(q_t + p_t)}{(-q_t)} = p_t q_t$$

$$= -q_t (q_t + p_t) = p_t q_t$$

$$\Rightarrow q_t + p_t = -p_t \text{ (or) } q_t = 0$$

$$\Rightarrow \boxed{q_t = -2p_t}; \boxed{q_t = 0}$$

\Rightarrow Hence if $q_t = 0$ (which does not occur any time in this case)

$$q_t = 0 \Rightarrow p_t = 0 \left(\because p_t = -\frac{q_t}{2} \right) \Rightarrow \begin{aligned} q_{t+1} &= 0 + 0 = 0 \\ p_{t+1} &= 0 - 0 = 0 \end{aligned}$$

\therefore Therefore, if one of them reaches '0'

\Rightarrow The fun. reaches an eq. @ 0

$$\boxed{f = 0}$$

—————X—————