ⓐ

### for RNN

#### Time complexity

time complexity in RNN is directly proportional to square of number of layers. also it is proportional to number of neurons & sequence length. this will be the case in both training & testing time.

$$RNN_{Time} \propto l^2 \times m \times t$$

both train & test

#### Space Complexity

during testing time it doesnot depend on length of sentences (sequence length).

$$RNN_{Space\ train} \propto t \times n \times l$$

$$RNN_{Space\ test} \propto n \times l$$

### for Transformers.

#### Time complexity

each embedding in encoder maps to each embedding in decoder. therefore time complexity is proportional to square of sequence length

$$Transformers_{time} \propto t^2 \times n \times l$$

#### Space complexity
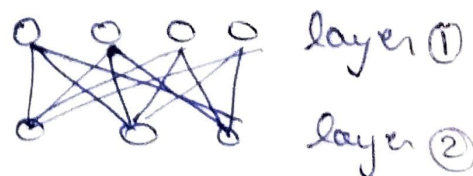
$$Transformers_{Space} \propto t \times n \times l$$

1.b) If the number of neurons in each layer is less than sequence length then learning will be effected badly. as same neuron will have to capture the weights that performs for more than 2 word embeddings of sequence

2c) Self attention in encoder layer is not a bottleneck ↑for parallelism as encoder has all the tokens in sequence to compute attention scores.

But for decoder, self attention acts as a bottleneck ↑for parallelism as it depends on the output of previous time step as well.

2d) Yes the feed forward network ~~support~~ look across the tokens.

layer ①
layer ②

This support parallelism as each layer in feed forward network can be ~~not~~ trained on different computational units.

Solution - ② (a)

The attention vector and weights are defined as :-

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

$$\alpha_i = \sum_{t=1}$$

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

⊙   $q_i$ = query vector

    $q_i = W_q x_i$

⊙   $k_i = W_k x_i$    (key vector)

⊙   $v_i = W_v x_i$    (value vector)

$$w_{ij} = q_i^T k_j$$
$$w_{ij} = softmax(w_{ij}) = \alpha_i$$

$$\boxed{z_i = \sum_{j} w_{ij} v_j}$$

This means that all the $j$ from $1$ to $m$
contributes towards the calculation of $z_i$.

Solution 2.b

$$\alpha'_i = \frac{\exp k_i^T q}{\sum_{l=1}^{m} \exp k_i^T q}$$

●    $k_i^T q \approx 0$    for   $i \notin \{a, b\}$     ─ ①

$$z = \frac{v_1 \exp(k_1^T q) + v_2 \exp(k_2^T q) \dots v_a \exp(k_a^T q) + v_b \exp(k_b^T q) \dots}{\exp(k_1^T q) \dots \dots \exp(k_n^T q)}$$

$$\boxed{z \approx \frac{v_a \exp(k_a^T q) + v_b \exp(k_b^T q)}{\exp(k_a^T q) + \exp(k_b^T q)}} \quad \left( if \ n \gg 0 \right)$$

\# given in question $z \approx \dfrac{v_a + v_b}{2}$

$$\frac{v_a \exp(k_a^T q) + v_b \exp(k_b^T q)}{\exp(k_a^T q) + \exp(k_b^T q)} = \frac{v_a + v_b}{2}$$

☉ let $\exp(k_a^T q) = a$    & $\exp(k_b^T q) = b$

$$\frac{a \cdot v_a + b \cdot v_b}{a + b} = \frac{v_a + v_b}{2}$$

$$\frac{a v_a}{2} + \frac{b v_b}{2} = \frac{b v_a}{2} + \frac{a v_b}{2} \quad ─ ②$$

from equation 2 we get.

$$a = b$$

$$\exp(K_a^T q) = \exp(K_b^T q)$$

$$\boxed{K_a^T q = K_b^T q} \longrightarrow \text{③}$$

\# Solving this analytically we get.

$$\boxed{q = (K_a + K_b).} \quad \text{✗✗}$$

we can check this by putting in eq$^n$ ③

$$K_a^T (K_a + K_b) = K_b^T (K_a + K_b)$$

$$K_a^T K_a + K_a^T K_b = K_b^T K_a + K_b^T K_b \quad \text{—④}$$

\# in question it is given $\underline{K_i \perp K_j}$ also

$\|K_i\| = 1$, therefore. eq$^n$ ④ can be solved as

$$\underset{1}{\overset{\not{1}}{\cancel{K_a^T K_a}}} + \underset{0}{\cancel{K_a^T K_b}} = \underset{0}{\cancel{K_b^T K_a}} + \underset{1}{\cancel{K_b^T K_b}}$$

$$\boxed{1 = 1}$$

\# Hence answer to the value of $\boxed{q = K_a + K_b}$

## Solution -3

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

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

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

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

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

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

$$L(q) = -D_{KL}\left(q(z|x) \,\|\, p(z)\right)$$

$$\qquad + E_{q(z|x)}\left(\log p(x|z)\right)$$

Solution - 4

to optimise   we   $\min_p \max_q f(p,v)$

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

$$\frac{d}{dq} f(q,p) = \frac{d(p.q)}{dq} = p.$$

$$\frac{d}{dp} f(p,q) = \frac{d\, p.q}{dp} = q.$$

Gradient descent.

$\rightarrow$  $q_{t+1} = q_t + \frac{df}{dq}$   $\{$ for maximising $\}$

$\rightarrow$  $p_{t+1} = p_t - \frac{df}{dp}$   $\{$ for minimising $\}$

$\rightarrow$ $q_{t+1} := q_t + p_t$

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

| $t \rightarrow$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| $q$ | 1 | 2 | 1 | -1 | -2 | -1 | 1 |
| $p$ | 1 | -1 | -2 | -1 | 1 | 2 | 1 |

$(p_t + q_t)$  $\leftarrow$

$(p_t - q_{t+1})$  $\leftarrow$

(b) By using above approach it is not possible to reach the optimal for all the parameters as we see in above example at $t = 6$ the values of $q$ & $p$ both reaches back to where it started.

(c) In GAN, we do not simultaneously reach the optimal for all parameters. instead we achieve Nash equillibrium where each parameter can not reduce their cost without altering the cost of other parameters.