

Q.1	a)	Time complexity	RNN	Transformer
		train	$t \ln^2$	$t^2 \ln$
		test	$t \ln^2$	$t^2 \ln$
b)		Space complexity	RNN	Transformer
		train	$t \ln$	$t \ln$
		test	$\ln$	$t \ln$

b) RNNs is sequential <sup>bottleneck</sup> because of sequence length. whereas for a transformer it only depends on number of layers. Hence when sequences are longer - small-layered transformers can learn much faster than vanilla RNNs

c) Yes, it is bottleneck. The sequential operations in transformers are independent of length but are quite complex/expensive to decode. It make it a tradeoff.

d) No, feed forward network and layer norm don't look across tokens. However, various operations/paths can be executed parallelly while flowing through FF and norm layers.

Q.2 a)  $z = v_j$  implies  $x_j = 1$  and  $x_{i: i \neq j} = 0$

This means that  $K_j^T q \ggg K_{i: i \neq j}^T q$

b) Assume,  
 $q = \sum_{j=1}^m \beta_j K_j$

$$\text{now } K_i^T q = K_i^T \left( \sum_j \beta_j K_j \right)$$

$$= \sum_j \beta_j K_i^T K_j$$

$$= \beta_i \|K_i\|^2 + 0 \quad \dots \quad \left\{ \begin{array}{l} \text{as } K_i \text{ are} \\ \text{orthogonal} \\ \text{max} \end{array} \right.$$

$$= \beta_i \times 1$$

$$\text{so } x_i = \frac{\exp(\beta_i)}{\sum_{j=1}^m \exp(\beta_j)}$$

$$Z \approx \frac{1}{2} (V_a + V_b) \Rightarrow x_a = x_b = \frac{1}{2}$$

$$\text{and } x_{i: i \neq a \text{ and } i \neq b} = 0$$

$$\text{setting } \beta_a = \beta_b \gg 0 \quad \text{and } \beta_{i: i \neq a \text{ and } i \neq b} \ll 0$$

in query vector  $q$  will achieve the same.  
 This will give  $x_a = x_b \approx \frac{1}{2}$

$$\text{and } x_{i: i \neq a \text{ and } i \neq b} \approx 0$$

Q.3

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

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

$$= \mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log \frac{p(z)}{q_\phi(z|x)} + \log p_\theta(x|z) \right]$$

Approximating  
 $q$  by  $q_\phi$   
 $p$  by  $p_\theta$

$$= \mathbb{E}_{z \sim q_\phi(z|x)} \left[ -\log \frac{q_\phi(z|x)}{p_\theta(z)} \right] + \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)]$$

$$= -KL(q_\phi(z|x) \parallel p_\theta(z)) + \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)]$$

$$= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - KL(q_\phi(z|x) \parallel p_\theta(z))$$

reconstruction term

regularization term,  
reconstruction

Here  $\mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)]$  is the regularization

term which aligns the forces the generated image to be similar to original one.

Here  $KL(q_\phi(z|x) \parallel p_\theta(z))$  is the regularization term which aligns the distribution  $q_\phi(z|x)$  with prior  $p_\theta(z)$



8.4  $f(p, q) = pq$

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

$\frac{d f(p, q)}{d q} = p$  will evaluate to 0  
as  $\max_q pq \geq 0$

hence only setting  $p=0$  will minimize  $\max_q pq$

a)

$q_0$	$q_1$	$q_2$	$q_3$	$q_4$	$q_5$	$q_6$
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 it is not possible to reach optimal value as after 6 iterations the new value of  $p, q$  are same as  $p_0, q_0$ . make it oscillatory sequence.

It is because the step size which is 1 is too large. Since min-max problems are unstable, small step size is preferable.

c) Nash equilibrium point is 0.

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

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

$$q_t(2p_t + q_t) = 0$$

$$\Rightarrow 2p_t = -q_t \quad \text{or} \quad q_t = 0$$

↓  
Not valid solution

(see part a)

hence  $q_t = 0 \Rightarrow p_{t+1} = 0$

and  $p, q$  from then  $= 0$

$$\left\{ \begin{array}{l} p_{t+1} = -q_t \\ q_{t+1} = p_t + q_t \end{array} \right.$$