

Assignment 4

9/8/24

1. a. i) For  $\alpha_i$  to be  $\sim 1$ ,  $k_{i+j}^T q$  must be  $\sim 0$  & $k_i^T q$  must be  $\gg 0$ , i.e.:ii)  $q_{i+j} = -k_{i+j}$ ,  $q_i = k_i$ .c.  $\approx \sqrt{1/\alpha_i} \approx \sqrt{1}$ b.)  $q_{j \neq a,b} = -k_{j \neq a,b}$  &  $q_{j=a,b} = k_{j=a,b}$  s.t. $1 \leq j, a, b \leq n$ ,  $j \neq a \neq b$ ,normalize:  $\|q\| = \sqrt{(k_a + k_b)^T (k_a + k_b)}$ 

$$= \sqrt{k_a^T k_a + k_b^T k_b + k_a^T k_b + k_b^T k_a}$$

$$= \sqrt{2}$$

$$\Rightarrow q = \frac{k_a + k_b}{\sqrt{2}}$$

$$\text{check: } k_a^T q = \frac{1}{\sqrt{2}} k_a^T k_a + \frac{1}{\sqrt{2}} k_a^T k_b = \frac{1}{\sqrt{2}}$$

$$k_b^T q = \frac{1}{\sqrt{2}} k_b^T k_a + \frac{1}{\sqrt{2}} k_b^T k_b = \frac{1}{\sqrt{2}}$$

c. i.)  $q = \frac{\mu_a + \mu_b}{\sqrt{2}}$  b/c  $k_i \approx \mu_a$  given small cov.ii.) Since  $k_a$  could  $\approx \frac{1}{2} M a^2$ ,  $q$  could be much larger, which when not mult. w/  $k_a$  would cause  $\alpha_a \sim 1$  & var to be wdy. disp.d. i.) Since  $c \approx \frac{1}{2} (v_a + v_b) = \frac{1}{2} (c_1 + c_2)$ , pick  $v_a, v_b = c_1, c_2$  by setting  $q_1 = \frac{k_a + k_b}{2}$ ,  $q_2 = \frac{k_a - k_b}{2}$ .ii) Same var in  $c_1, c_2$  as before but now avg out, so not as much.

e) By avg. head. anticipation, there is reduced var despite high var.



coll. scores.

2. a. i.)  $Q_{perm} = X_{perm} \cdot W_Q = P X W_Q = P Q; K_{perm} = P K;$   
 $V_{perm} = P V; H_p = \text{softmax}\left(\frac{Q P K^T}{\sqrt{d}}\right) \cdot V_p$   
 $= \left(\frac{P Q (P K)^T}{\sqrt{d}}\right) P V$   
 $= \left(\frac{P Q K^T P^T}{\sqrt{d}}\right) P V$   
 $= P \text{softmax}\left(\frac{Q K^T}{\sqrt{d}}\right) P^T P V$   
 $= P H_p$

$Z_p = \text{ReLU}(H_p \cdot W_1 + 1 \cdot b_1) W_2 + 1 \cdot b_2$   
 $= \text{ReLU}(P \cdot H \cdot W_1 + 1 \cdot b_1) W_2 + 1 \cdot b_2$   
 $= P \text{ReLU}(H W_1 + 1 \cdot b_1) W_2 + 1 \cdot b_2$   
 $= P Z$

ii.) Weights don't use permutation in calculation, so everything after gets shifted by perm & text would be jumbled.  $\otimes$

b. i.) Yes, since there is delayed step increase in pos. emb., which correlates with actual text.  $\otimes$

ii.) Since denom so large (i.e.  $\gg T=1$ ), values never  $\gg 2\pi$ , so safe.  $\otimes$

3. a-c.) code

d) Dev acc.: 1.2%

Baseline: 5%  $\otimes$

e.) code

f) Dev acc.: 27.6%  $\otimes$





$$g.i.) \begin{pmatrix} \cos \theta_1 + i \sin \theta_1 \\ \vdots \\ \cos \theta_{d/2} + i \sin \theta_{d/2} \end{pmatrix} \odot \begin{pmatrix} x_t^{(1)} + i x_t^{(2)} \\ \vdots \\ x_t^{(d-1)} + i x_t^{(d)} \end{pmatrix}$$

$$= \begin{pmatrix} \cos \theta_1 x_t^{(1)} - \sin \theta_1 x_t^{(2)} + i (\sin \theta_1 x_t^{(1)} + \cos \theta_1 x_t^{(2)}) \\ \vdots \\ \cos \theta_{d/2} x_t^{(d-1)} - \sin \theta_{d/2} x_t^{(d)} + i (\sin \theta_{d/2} x_t^{(d-1)} + \cos \theta_{d/2} x_t^{(d)}) \end{pmatrix}$$

$$= \begin{pmatrix} (\cos \theta_1 x_t^{(1)} - \sin \theta_1 x_t^{(2)}) + i (\sin \theta_1 x_t^{(1)} + \cos \theta_1 x_t^{(2)}) \\ \vdots \\ (\cos \theta_{d/2} x_t^{(d-1)} - \sin \theta_{d/2} x_t^{(d)}) + i (\sin \theta_{d/2} x_t^{(d-1)} + \cos \theta_{d/2} x_t^{(d)}) \end{pmatrix}$$

$$ii) \langle \text{RoPE}(z_1, t_1), \text{RoPE}(z_2, t_2) \rangle$$

$$= \langle z_1 e^{it_1 \theta}, z_2 e^{it_2 \theta} \rangle$$

$$= \text{Re} \left( \overline{z_1 e^{it_1 \theta}} z_2 e^{it_2 \theta} \right)$$


$$= \text{Re} \left( z_1 \bar{z}_2 e^{i\theta(t_1 - t_2)} \right)$$


$$\langle \text{RoPE}(z_1, t_1 - t_2), \text{RoPE}(z_2, 0) \rangle$$


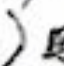
$$= \langle z_1 e^{i(t_1 - t_2)\theta}, z_2 e^{i0\theta} \rangle$$

$$= \text{Re} \left( z_1 e^{i(t_1 - t_2)\theta} \bar{z}_2 \right)$$

$$iii) \text{Dev acc.: } 19.2\%$$

4 a.) Via pretraining, the model better learned the facts. 

b.) 1: lose trust in systems (in over-reliance to knowledgeable), 

2: spread misinformation  (misinformation) 

c.) Remembering most similar vectors to the prompt, & return that, which can be utterly wrong. 