

a)

$$f = \sigma(w_{xf}x_t + w_{hf}h_{t-1} + b_f) = \sigma([w_{hf}, w_{xf}] \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} + b_f)$$

$$\left. \begin{aligned} i &= \sigma(w_{xi}x_t + w_{hi}h_{t-1} + b_i) \\ g &= \tanh(w_{xg}x_t + w_{hg}h_{t-1} + b_g) \end{aligned} \right\} \rightarrow \text{input_gate} = i \odot g$$

$$o = \sigma(w_{xo}x_t + w_{ho}h_{t-1} + b_o)$$

$$c_t = c_{t-1} \odot f + \text{input_gate}$$

$$h_t = \tanh(c_t) \cdot o$$

b)

$$\frac{dE}{dh_t} = E_{\text{delta}} \quad , \quad \frac{dE}{dc_t} = \frac{dE}{dh_t} \times \frac{dh_t}{dc_t} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t))$$

$$\frac{dE}{do} = \frac{dE}{dh_t} \cdot \frac{dh_t}{do} = E_{\text{delta}} \cdot \tanh(c_t) \quad , \quad \frac{dE}{di_{\text{input}}} = \frac{dE}{dc_t} \cdot \frac{dc_t}{di_{\text{input}}} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t)) \cdot 1$$

$$\frac{dE}{dg} = \frac{dE}{di_{\text{input}}} \cdot \frac{di_{\text{input}}}{dg} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t)) \cdot i \quad , \quad \frac{dE}{di} = \frac{dE}{di_{\text{input}}} \cdot \frac{di_{\text{input}}}{di} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t)) \cdot g$$

$$\frac{dE}{df} = \frac{dE}{dc_t} \cdot \frac{dc_t}{df} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t)) \cdot c_{t-1}$$

$$\frac{dE}{dw_{xo}} = \frac{dE}{do} \cdot \frac{do}{dw_{xo}} = E_{\text{delta}} \cdot \tanh(c_t) \cdot o \cdot (1 - o) \cdot x_t$$

$$\frac{dE}{dw_{ho}} = \frac{dE}{do} \cdot \frac{do}{dw_{ho}} = E_{\text{delta}} \cdot \tanh(c_t) \cdot o \cdot (1 - o) \cdot h_{t-1}$$

$$\frac{dE}{dw_{hg}} = \frac{dE}{dg} \cdot \frac{dg}{dw_{hg}} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t)) \cdot i \cdot (1 - g^2) \cdot h_{t-1}$$

$$\frac{dE}{dw_{xg}} = \frac{dE}{dg} \cdot \frac{dg}{dw_{xg}} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t)) \cdot i \cdot (1 - g^2) \cdot x_t$$

$$\frac{dE}{dw_{xi}} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t)) \cdot g \cdot i \cdot (1 - i) \cdot x_t \quad , \quad \frac{dE}{dw_{hi}} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t)) \cdot g \cdot i \cdot (1 - i) \cdot h_{t-1}$$

$$\frac{dE}{dw_{xf}} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t)) \cdot f \cdot (1 - f) \cdot x_t \quad , \quad \frac{dE}{dw_{hf}} = E_{\text{delta}} \cdot (1 - \tanh^2(c_t)) \cdot f \cdot (1 - f) \cdot h_{t-1}$$

$$\frac{dE}{dbf} = E_{\text{delta}} \cdot (1 - \tanh^2(ct)) \cdot f(1-f) \quad , \quad \frac{dE}{db_i} = E_{\text{delta}} \cdot (1 - \tanh^2(ct)) \cdot g \cdot (i(1-i))$$

$$\frac{dE}{db_g} = E_{\text{delta}} \cdot (1 - \tanh^2(ct)) \cdot i \cdot (1 - g^2) \quad , \quad \frac{dE}{b_o} = E_{\text{delta}} \cdot \tanh(ct) \cdot o(1-o)$$

$$\frac{dE}{dx_t} = E_{\text{delta}} \cdot \left(\tanh(ct) \cdot (o(1-o)) \cdot w_{xo} + (1 - \tanh^2(ct)) \cdot i \cdot (1 - g^2) \cdot w_{xg} \right. \\ \left. + (1 - \tanh^2(ct)) \cdot g \cdot (i(1-i)) \cdot w_{xi} + (1 - \tanh^2(ct)) \cdot f(1-f) \cdot w_{xf} \right)$$

$$\frac{dE}{dh_{t-1}} = E_{\text{delta}} \cdot \left(\tanh(ct) \cdot (o(1-o)) \cdot w_{ho} + (1 - \tanh^2(ct)) \cdot i \cdot (1 - g^2) \cdot w_{hg} \right. \\ \left. + (1 - \tanh^2(ct)) \cdot g \cdot (i(1-i)) \cdot w_{hi} + (1 - \tanh^2(ct)) \cdot f(1-f) \cdot w_{hf} \right)$$

2 اگر فرض کنیم q, k ماتریس‌های تصادفی با ابعاد $d_k * d_k$ باشند و ورودی‌ها از توزیع $N(0,1)$ پیروی می‌کنند و iid باشند داریم:

اینکده
متغیر داخلی
 q_i در k_i
بزرگ است ب

$$\langle q_i, k_i \rangle = \sum_j q_{ij} k_{ij} \sim N(0, d_k)$$

این متغیر داخلی از توزیع نرمال با میانگین صفر و واریانس d_k پیروی می‌کند. بزرگی جکواردر از بزرگ شدن متغیر داخلی scale کردن مقادیر در $\frac{1}{\sqrt{d_k}}$ به توزیع $N(0,1)$ خواهیم رسید و از Gradient vanishing جلوگیری می‌شود.

3
(a)

1- با توجه به اینکه $\sum_i \alpha_i = 1$ است هر توان از α به عنوان توزیع احتمال استفاده کرد.

2- در صورت صحبت بسیار زیاد query با یکی از کلمات عدم صحبت با سایر کلمات، اکثر وزن روی α مربوط به آن کلمه خواهد بود.

3- اگر شرایط قسمت 2 باشد، اینکده مقدار α بیشتر مربوط به α ای می‌باشد که کلمات صحبت زیاد با query داشته باشد.

4- وقتی query با یکی از کلمات شبیه باشد Attention، بیشتر به Value مربوط به آن کلمه توجه می‌کند و خردی‌ها بیشتر تأثیر یافته از Value آنهاست.

$$v_a = c_1 a_1 + c_2 a_2 + \dots + c_m a_m \rightarrow a_j^T v_a = \underbrace{c_1 a_j^T a_1}_0 + \dots + \underbrace{c_j a_j^T a_j}_1 + \dots + \underbrace{c_m a_j^T a_m}_0 = c_j \quad -1 \quad (b)$$

$$v_b = d_1 b_1 + d_2 b_2 + \dots + d_p b_p \rightarrow a_j^T v_b = \underbrace{d_1 a_j^T b_1}_0 + \dots + \underbrace{d_p a_j^T b_p}_0 = 0$$

$$\Rightarrow \text{if } M = \begin{bmatrix} a_1^T \\ \vdots \\ a_m^T \end{bmatrix} \Rightarrow M(v_a + v_b) = M v_a + M v_b = v_a$$

$$c \approx \frac{1}{2} (v_a + v_b) = \frac{1}{2} v_a + \frac{1}{2} v_b \quad -2$$

$$\alpha_a \approx \alpha_b, \alpha_{i \neq a, b} \approx 0 \quad \rightarrow q = \beta(K_a + K_b), \quad \beta \text{ بزرگ}$$

$$\beta \gg 0, q = \beta K_a + \beta K_b \Leftrightarrow k_i \approx \mu_i \quad -1 \quad (c)$$

۱- از آنجایی که α ها بسیار کوچکند پس داریم

$$\alpha \approx 0 \rightarrow K_a \approx \gamma^M_a, \gamma \sim N(1, 1/2) \quad -2$$

$$\begin{matrix} k_i^T q \approx 0, i \neq \{a, b\} \\ K_b^T q = \beta \end{matrix} \Rightarrow c = \frac{e^{\gamma\beta}}{e^{\gamma\beta} + e^{\beta}} v_a + \frac{e^{\beta}}{e^{\beta} + e^{\gamma\beta}} v_b$$

با فرضی که مقدار c به سمت v_a است؛
کوچک کردن مقدار c به سمت v_b خواص رفت و اگر γ بزرگ باشد، مانند موارد قبلی خواص بود.

$$\beta \gg 0, q_1 = q_2 = \beta(K_a + K_b) \quad -1 \quad (d)$$

۲- در صورت استفاده از multihead-Attention، هر چند تعداد head ها بیشتر شود

مقدار c به سمت قبلی به میانه میل می کند، یعنی ۱ میل می کند، و $c \sim \frac{1}{2}(v_a + v_b)$ خواص بود

در Vision Transformers، تصویر ورودی به تعدادی Patch بدون overlap تقسیم می‌شوند و سپس هر Patch، Flatten می‌شود و به یک بردار $(N \times (P^2 \times c))$ که N تعداد Patch ها و P ساید هر Patch (قبل از Flat شدن) و c تعداد چنل‌ها می‌باشد. هر Patch به یک توکن در sequence است و اگر ساید تصویر ورودی تغییر کند تعداد و ساید Patch ها تغییر می‌کند. با Pre-train کردن مدل در ابعاد پایین تقسیم تعداد بزرگترها، در نتیجه سرعت محاسبات، حافظه بهبود می‌یابد و Fine-tune کردن آن نیاز به ساید تصویر بالاتر به generalization بهتر مدل کمک می‌کند. در مدیریت تغییر ابعاد تصویرها، ممکن است به معیاری برای جدید سازگار با ابعاد جدید نیاز باشد. روش پیشنهادی: Crop & Pad تصویر ورودی برای رسیدن به ابعاد دگوا. استفاده از شبکه‌های CNN برای رسیدن به ابعاد دگوا ورودی.

$$e_t = \begin{bmatrix} \sin\left(\frac{t}{f_1}\right) \\ \cos\left(\frac{t}{f_1}\right) \\ \sin\left(\frac{t}{f_2}\right) \\ \cos\left(\frac{t}{f_2}\right) \\ \vdots \\ \sin\left(\frac{t}{\frac{d_{model}}{2}}\right) \\ \cos\left(\frac{t}{\frac{d_{model}}{2}}\right) \end{bmatrix}, \quad f_m = \frac{1}{\lambda_m} := 10000 \frac{2m}{d_{model}}, \quad T^{(k)} E_{t,:} = E_{t+k} \quad (1)$$

مصفی‌ها

$$T^{(k)} = \begin{bmatrix} \Phi_1^{(k)} & & \\ & \Phi_2^{(k)} & \\ & & \ddots \\ & & & \Phi_{\frac{d_{model}}{2}}^{(k)} \end{bmatrix} \quad \Phi_m^{(k)} = \begin{bmatrix} \cos(r_m k) & -\sin(r_m k) \\ \sin(r_m k) & \cos(r_m k) \end{bmatrix}^T$$

و ساید درایه‌ها ماتریس‌ها 2x2 می‌باشند.

$$\begin{aligned} \sin(\lambda_m(t+k)) &= \cos r_m k \sin \lambda_m t + \sin r_m k \cos \lambda_m t \\ \cos(\lambda_m(t+k)) &= -\sin r_m k \sin \lambda_m t + \cos r_m k \cos \lambda_m t \end{aligned}$$

$$\Rightarrow \begin{bmatrix} \cos(r_m k) & \sin(r_m k) \\ -\sin(r_m k) & \cos(r_m k) \end{bmatrix} \begin{bmatrix} \sin \lambda_m t \\ \cos \lambda_m t \end{bmatrix} = \begin{bmatrix} \sin \lambda_m(t+k) \\ \cos \lambda_m(t+k) \end{bmatrix}$$

ماتریس جمع \Rightarrow اگر $\lambda = r$ رابطه قبلی صحیح خواهد بود.

$$\Phi_m^{(k)} = \begin{bmatrix} \cos \lambda_m k & \sin \lambda_m k \\ -\sin \lambda_m k & \cos \lambda_m k \end{bmatrix}^T, \quad \lambda_m = 10000 \frac{2m}{d_{model}}$$