

① اگر ماتریس های $W_1 \in \mathbb{R}^{K \times m}$ و $W_2 \in \mathbb{R}^{m \times K}$ را که SVD تجزیه کنیم داریم:

$$W_1 = U_1 \Sigma_1 V_1^T, \quad W_2 = U_2 \Sigma_2 V_2^T$$

$$L(W_1, W_2; X) = \frac{1}{n} \|X - W_2 W_1 X\|_F^2 + \lambda \|W_1\|_F^2 + \lambda \|W_2\|_F^2$$
 تابع هدف به صورت زیر است:

اگر نرم های مربوط به تنظیم سازی درجه نه داشته باشند، میهم کردن تابع هزینه منجر به این می شود که $W_2 W_1$ مانند یک ماتریس نگاشت دمج عمل کند که X را به K نزدیکترین Singular vector های آن نگاشت کند. بنابراین می توان گفت که $W_2 W_1$ تقریباً یک ماتریس همانی است. بنابراین W_2 و W_1 تقریباً متکدر می هستند.

حال حدس آن نرم های تنظیم ساز را به صورت زیر می نهد:

$$\lambda (\|W_1\|_F^2 + \|W_2\|_F^2) = \lambda (\|\Sigma_1\|_F^2 + \|\Sigma_2\|_F^2) = \lambda \sum_{i=1}^K \left(\sigma_i^2 + \frac{1}{\sigma_i^2} \right)$$

$$\text{میهم کردن عبارت بالا به صورت } 0 = \frac{-2}{\sigma^3} + 2\sigma = 0 \rightarrow \left(\sigma^2 + \frac{1}{\sigma^2} \right)' = 0 \text{ نتیجه می دهد که } \sigma = 1$$

بنابراین می توان نتیجه گرفت که W_2 و W_1 pseudo-inverse متکدر و همچنین singular value های غیر صفر آن ها نیز نزدیک به یک است. این بدان معناست که ردیف های W_1 تقریباً متعامد و همچنین ستون های W_2 نیز تقریباً متعامد هستند. در نهایت می توان گفت که با داشتن شرایط صورت بالا یک inductive bias برای W_2 داریم که این است که تقریباً pseudo-inverse ماتریس W_1 است.

② تابع هزینه بار سازی $J = \frac{1}{2} \|x' - x\|^2$ هر دو بردار ورودی x و x' را در نظر می گیریم و از

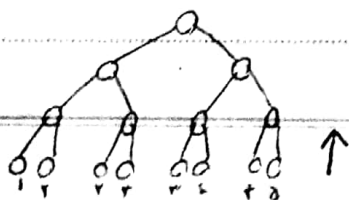
اینکه شبکه به یکی از آن ها تقدم ندهد جلوگیری می کند. بنابراین اگر تمام بردارهای ورودی را در نظر بگیریم برای ساخت بردار P .

$$\begin{aligned} W_1 &\in \mathbb{R}^{D_p \times 2 D_x}, \quad b_1 \in \mathbb{R}^{D_p \times 1}, \quad W_2 \in \mathbb{R}^{2 D_x \times D_p}, \quad b_2 \in \mathbb{R}^{2 D_x \times 1} \\ W_3 &\in \mathbb{R}^{D_c \times D_p}, \quad b_3 \in \mathbb{R}^{D_c \times 1} \end{aligned}$$

تعداد پارامترها $= 4 D_p D_x + D_p + 2 D_x + D_c D_p + D_c$

ج. می توان جمله ورودی را به جهت کلمات مجاور تقسیم کرد پس بردار P هر یک از این جهت ها را به دست آورده پس به صورت بازگشتی همین کار را روی بردارهای P به دست آمده انجام می دهیم. این کار را ادامه می دهیم تا به یک بردار P واحد برسیم.

در واقع به صورت درختی متقابل نگاری کنیم.



② (ا) اگر ابعاد ورودی یک پیکر آکام تعدادی decoder برای مدل کردن در روشی های general و tanh layer یکسان است اما در ابعاد دیگر tanh layer تعدادی مدل کردن بهتر دارد. general و tanh layer از نظر پیچیدگی حسابی تقریباً با هم برابرند. در روش dot با ابعاد افقی وجود ندارد اما در general تعداد ابعاد برای w_a و $l_{d \times l_e}$ یکسان است. در روش tanh layer تعداد ابعاد برای V_a^T و D_{w_a} و تعداد ابعاد برای w_a یکسان است! $D_{w_a} \times D_{D_{w_a}}$ از لحاظ هزینه حسابی dot نسبت به دو روش دیگر کمتر است اما تعدادی مدل کردن گسری دارد. همچنین در مرحله Backprop نیز هزینه گسری دارد. در روش tanh layer هزینه score را tanh اعمال شده از حالت صاف خارج می کند و هزینه در بازه [0, 1] نگاشت می دهد. بنابراین می تواند از explode شدن گره های پنهان جلوگیری کند و همچنین مشکل Vanish شدن گره های را حل کند. در شبکه های seq2seq مدل general از دو روش دیگر بهتر است چون نسبت به dot تعدادی مدل کردن بهتری دارد و همچنین نسبت به tanh layer هزینه حسابی گسری دارد.

(ب) می گوییم مدل سیر عمیق و پیکر برای t دارد که مطابق زیر است:

$$h_d^{t-1} \rightarrow \alpha^t \rightarrow c^t \rightarrow h_d^t$$

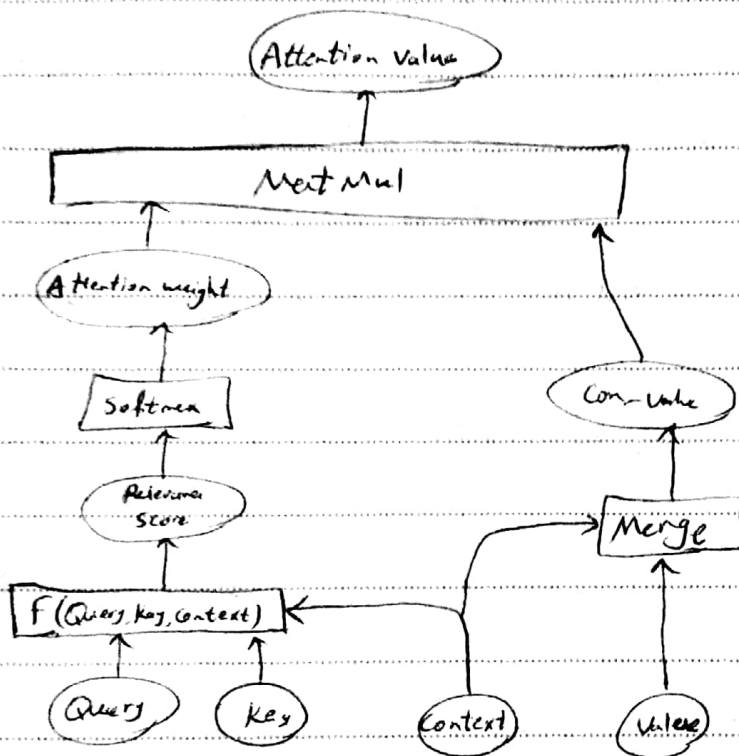
اما می گوییم دوم سیر یکم هزینه های دارد که به صورت زیر است:

$$h_d^t \rightarrow \alpha^t \rightarrow c^t \rightarrow \tilde{h}_d^t$$

می گوییم اول در encoder از روش bidirectional استفاده شده است که باعث دوم هزینه در مدل یادگیری می شود. اما در می گوییم دوم از روش یک طرفه استفاده شده است.

در می گوییم اول جای مناسب score از روش concat استفاده شده اما در می گوییم دوم از چند روش از جمله dot، general، concat و location استفاده شده که location نتیجه بهتری می دهد. به علاوه یکی از تدابیر گشت می گوییم دوم نسبت به روش اول برتری دارد.

ج) روشی برای دستیابی به این هدف، برای هر یک از این توکن‌ها، یک وزن توجه (attention weight) محاسبه می‌شود. این وزن به هر یک از توکن‌های Context بستگی دارد. برای محاسبه این وزن، ابتدا Query و Key را در یک ضرب داخلی (dot product) قرار می‌دهیم. سپس نتیجه را با نرمال‌سازی (softmax) تبدیل به یک توزیع احتمالی می‌کنیم. این توزیع، وزن توجه (attention weight) را می‌دهد. در نهایت، این وزن را با Value ضرب می‌کنیم تا به Contextual Value برسیم. در نهایت، همه Contextual Value‌ها را با هم جمع می‌کنیم تا به Attention Value برسیم.



$N \rightarrow$ Sequence length
 $D \rightarrow$ dim. of Q and K
 $M \rightarrow$ dim. of V

Computation Cost: $O(N \max(D, M))$

Memory: $O(NM + ND) + O(N^2) + O(N^2M) + O(NM) = O(N^2M + ND)$

\downarrow \downarrow \downarrow \downarrow
 K, Q, V $Q_i^T K_j$ $\exp(\frac{Q_i^T K_j}{\sqrt{D}})$ V_i'

$$\text{sim}(q, k) = (q^T k + 1)^2$$

$$\phi(x) = [1, x_1, x_2, \dots, x_D, x_1^2, x_2^2, \dots, x_D^2, \sqrt{2} x_1 x_2, \dots, \sqrt{2} x_{D-1} x_D]^T$$

$$V_i' = \frac{\sum_{j=1}^N \phi(Q_i)^T \phi(K_j)}{\sum_{j=1}^N \phi(Q_i)^T \phi(K_j)} V_j$$

هزینه محاسباتی برای یک کرنل چند بعدی درجه ۲ و برای $O(N^2 M)$ است. همین حافظه نیز
 نیاز برای $O(NM \max(D, M))$ است. برآورد آن گفت که اگر $N > D^2$ است. استاندارد از این رابطه
 است از رابطه قبلی برآورد