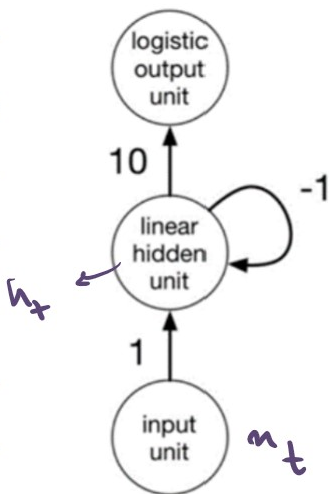


سوال ۱:



$$h_t = n_t - h_{t-1}$$

$$h_{t-1} = n_{t-1} - h_{t-2}$$

⋮

$$h_1 = n_1 - h_0$$

$$\Rightarrow h_t = n_t - n_{t-1} + h_{t-2} = n_t - n_{t-1} + n_{t-2} - h_{t-3}$$

$$\Rightarrow h_t = n_t - n_{t-1} + n_{t-2} - \dots - n_1 + h_0$$

$$h_1 = n_1 - h_0 \quad h_2 = n_2 - n_1 + h_0 \quad \Rightarrow h_T = \sum_{i=1}^T (-1)^i n_i + h_0$$

از نوع

$$y = \sigma(10h_T) \Rightarrow y_T = \sigma\left(10 \sum_{i=1}^T (-1)^i n_i\right) + 10h_0$$

سوال ۲: (۱) ابعاد بسیار بالا: هر کلمه باید بردار به طول هم حکایت حسیان داده شود پس حافظه زیار

استیج است و سرعت محاسبات پایین می آید.

(۲) آنکه یک کلمه اضافه شود باید ابعاد تمام بردارها را تغییر داد.

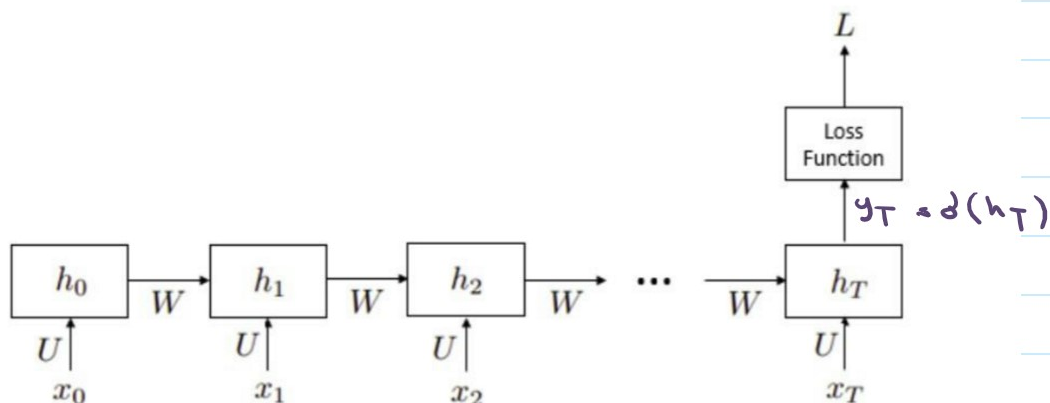
(۳) در one-hot ارتباط معنای کلمات در نظر گرفته می شود مثلا کلمه queen با king

معنای فاملا متفاوتی خواهند داشت. به عنوانی که شبیه مناسب است که کلمات در فضای نزدیک

فاصله کمی داشته باشند اما در one-hot رعایت نمی شود

(۴) sparse است و مار کردن با داده های sparse در شبیه سازی است.

الف)



$$h_{t+1} = \sigma(z) \quad z = Ux_{t+1} + Wh_t$$

$$\frac{\partial \mathcal{L}}{\partial h_t} = \frac{\partial \mathcal{L}}{\partial h_{t+1}} \times \frac{\partial h_{t+1}}{\partial h_t} \quad \frac{\partial h_{t+1}}{\partial h_t} = \frac{\partial h_{t+1}}{\partial z_{t+1}} \times \frac{\partial z_{t+1}}{\partial h_t}$$

$$\Rightarrow \frac{\partial h_{t+1}}{\partial h_t} = \sigma'(Ux_{t+1} + Wh_t) (1 - \sigma'(Ux_{t+1} + Wh_t)) \times W$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial h_t} = \frac{\partial \mathcal{L}}{\partial h_{t+1}} \times \underbrace{\sigma'(Ux_{t+1} + Wh_t)}_{z_{t+1}} \underbrace{(1 - \sigma'(Ux_{t+1} + Wh_t))}_{z_{t+1}} \times W$$

ب.

$$\frac{\partial \mathcal{L}}{\partial h_0} = \frac{\partial \mathcal{L}}{\partial h_1} \times \sigma'(z_1) (1 - \sigma'(z_1)) \times W$$

$$= \frac{\partial \mathcal{L}}{\partial h_2} \times \sigma'(z_2) (1 - \sigma'(z_2)) \sigma'(z_1) (1 - \sigma'(z_1)) W^2$$

$$= \frac{\partial \mathcal{L}}{\partial h_T} \times \prod_{i=1}^T \sigma'(z_i) (1 - \sigma'(z_i)) W^T$$

$$= \frac{\partial \mathcal{L}}{\partial h_T} \times \prod_{i=1}^T \sigma'(Ux_i + Wh_{i-1}) (1 - \sigma'(Ux_i + Wh_{i-1})) W^T$$

(آ) برای به دست آوردن حد بالا برای  $\frac{\partial \mathcal{L}}{\partial h_0}$

$$\max \left| \frac{\partial \mathcal{L}}{\partial h_0} \right|$$

$$\max \sigma(z) (1 - \sigma(z))$$

$$0 \leq \sigma(z) \leq 1$$

$$\sigma'(z) (1 - \sigma(z)) + \sigma(z) - \sigma(z)'$$

$$\sigma'(z) - 2\sigma'(z)\sigma(z) = 0 \Rightarrow 1 - 2\sigma(z) = 0 \Rightarrow \sigma(z) = \frac{1}{2}$$

$$\Rightarrow \max \sigma(z) (1 - \sigma(z)) = \frac{1}{4}$$

$$\left| \frac{\partial \mathcal{L}}{\partial h_0} \right| \leq \left| \frac{\partial \mathcal{L}}{\partial h_T} \right| \times \frac{1}{4}^T |W|^T = \frac{\partial \mathcal{L}}{\partial h_T} \left( \frac{|W|}{4} \right)^T$$

gradient exploding  $\rightarrow \left| \frac{W}{4} \right| < 1 \Rightarrow |W| \leq 4$

(ب)

$$z_{t+2} = Ux_{t+2} + Wh_{t+1} + Mh_t \quad h_{t+2} = \sigma(z_{t+2})$$

$$z_{t+1} = Ux_{t+1} + Wh_t + Mh_{t-1} \quad h_{t+1} = \sigma(z_{t+1})$$

$$\frac{\partial \mathcal{L}}{\partial h_t} = \frac{\partial \mathcal{L}}{\partial z_{t+2}} \times \frac{\partial h_{t+2}}{\partial h_t} + \frac{\partial \mathcal{L}}{\partial z_{t+1}} \times \frac{\partial h_{t+1}}{\partial h_t}$$

$$\frac{\partial h_{t+2}}{\partial h_t} = \sigma(z_{t+2})(1 - \sigma(z_{t+2}))M \quad \frac{\partial h_{t+1}}{\partial h_t} = \sigma(z_{t+1})(1 - \sigma(z_{t+1}))W$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial h_t} = \frac{\partial \mathcal{L}}{\partial z_{t+2}} \times \sigma(z_{t+2})(1 - \sigma(z_{t+2}))M + \frac{\partial \mathcal{L}}{\partial z_{t+1}} \times \sigma(z_{t+1})(1 - \sigma(z_{t+1}))W$$

$$\frac{\partial \mathcal{L}}{\partial h_{t+1}} = \frac{\partial \mathcal{L}}{\partial z_{t+2}} \times \sigma(z_{t+2})(1 - \sigma(z_{t+2})) \times W$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial h_t} = \frac{\partial \mathcal{L}}{\partial z_{t+2}} \sigma(z_{t+2})(1 - \sigma(z_{t+2}))M + \frac{\partial \mathcal{L}}{\partial z_{t+2}} \sigma(z_{t+2})(1 - \sigma(z_{t+2}))W^2$$

gradient ممکن است رخ داده اما  $\sigma(z_{t+1})(1 - \sigma(z_{t+1}))W^2$   $\downarrow$   
 vanishing  $\downarrow$   
 به مسیر جابجایی وجود دارد.

با اضافه کردن skip-connection  $\leftarrow$  می‌توان مسیر دیگری به تیرادین اضافه کرد.

تیرادین به جایی مسیر موازی  $W^N$  به مسیر جابجایی دارد که توان مسیر تیرادین را کاهش می‌دهد و این موضوع

از محاسبه گر جابجایی می‌کند. جریان تیرادین بهتر خواهد بود. احتمال مستقیم  $h_{t+2}$  به  $h_t$  می‌تواند تغییر کند که

در مسیر ممکن است باعث محاسبه گر شود را کم کند.

$$g = \frac{os \cdot os}{s \cdot w}$$

clipping by value : if  $g > \text{threshold} \Rightarrow g = \text{threshold}$

مزایا : پیاده سازی راحت و همپوشانی از بزرگ شدن مقدار ترددیان به صوبت از انقباض ترددیان  
معایب : مسیرهای ترددیان در نظر گرفته نمی شود به ممکن است فرآیند جبهه سازی را مختل کند

clipping by norm : if  $\|g\| > \text{threshold} \Rightarrow g = \frac{\text{threshold}}{\|g\|} g$

مزایا : مسیر ترددیان را حفظ می کند و فقط اندازه را تغییر می دهد  
برای زمانی که ابعاد بالاست مفید است (تغییر اندازه کلی را در نظر می گیرد)

چرا by clipping بهتر است ؟ ← همانطور که گفته شد صحت ترددیان فقط در مورد norm  
و به نسبت اندازه آن ، ترددیان تغییر می کنند به فرآیند جبهه سازی را مختل نمی کند.

از آنجا که اندازه ی کلی را در نظر می گیرد نه اندازه هر جزء (component) بولی

ترددیان ها با ابعاد زیاد می تواند استفاده شود.



انهم مشكلاتی که این پروگرام می‌کند:

اثر تکرار توکن با افزایش توکن‌سور (در زمان  $t$ ) این خطا به تمام توکن‌های که در timestep های

بعدی راه می‌یابد و کیفیت sequence را تضعیف می‌کند.  $\hookrightarrow$  Error Accumulation

همچنین کیفیت sequence مخفی تولید شده بسیار بر درسم توکن‌های اولیه وابسته است که خوب نیست

چون سبب تمام اولیه دقیق نیست.

در روش آموزش Teacher Forcing به جای token تولید شده در  $t$ ، مقدار

ground truth token به عنوان ورودی timestep بعدی استفاده می‌شود

این موضوع باعث می‌شود که مشکل Error Accumulation که توضیح داده شد،

در فرآیند آموزش رخ ندهد. همچنین استفاده از ground truth باعث می‌شود که ارتباط بین

token ها بهتر ترسفت می‌شود و همگامی سری‌های رخ می‌دهد. با استفاده از ground truth token

فرآیند یادگیری نیز آسان‌تر می‌شود.

ب) exposure bias مشکلی است که به در نتیجه آموزش و inference به وجود می آید.

در هنگام آموزش با ground truth token ها به عنوان ورودی استفاده می کنیم اما در

هنگام inference مدل از توکن تولید شده به عنوان ورودی استفاده می کند و ما target را خارج

و مدل بر خودش مشق است. از آنجایی که حرفه آموختن مدل یاد گرفته است که زمانی که ورودی

predicted توکن هست، چگونه باید عمل کند؛ عکس در آن در زمان inference خوب نخواهد بود

یعنی تواند ورودی های غلط یا noisy را حمل کند. over-reliance (اعتماد بیش از حد) به ground truth

و نتایج این روش با کاربردهای واقعی که نیاز به حمل کردن ورودی های غلط دارد، مشکلی است که وجود دارد.

ج)

در روش scheduled sampling، در هر مرحله باید احتمال مشخص ground truth

یا token پیشین شده استفاده می شود. به این صورت مدل یاد می گیرد که با token های پیشین

شده هم کار کند و خطاها را حمل کند و اینکه عکس در مدل بهتر می شود.

یک تابع scheduling وجود دارد که احتمال استفاده از توکن ground truth یا predicted

را در هر مرحله مشخص می کند. (این تابع در timestep های اولیه در ground truth با احتمال بیشتری

استفاده می کند و به تدریج احتمال استفاده از توکن های پیشین شده را بیشتر می کند، این روند باعث می شود که

در ابتدا به خوبی انجام شود و در ادامه حمل کردن ورودی های مشابه هم یاد می گیرد) و با توجه به احتمال تدریج

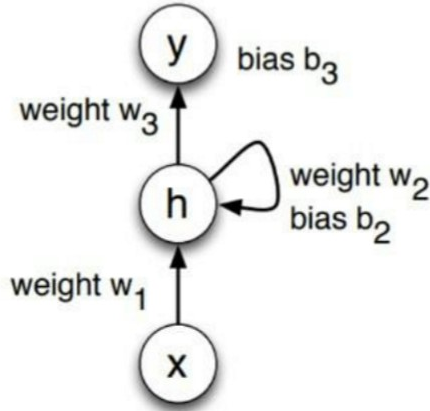
شده اینکه از چه ورودی استفاده کنیم قطع می شود.

این کار حاصله از train و inference را هم می کند، مدل را adapt می کند و از آنجایی

که به تدریج استفاده از توکن های پیشین شده بیشتر می شود، مدل به تدریج و در طول زمان در کنار یادگیری

مقادیر درست به ورودی های نادرست هم adapt می شود.

سوال ۵ :



$$n_t = g(w_1 n_t + w_2 h_{t-1} + b_2)$$

$$y = f(w_3 h_t + b_3)$$

بنابراین  $h_t$  اطلاعات از  $n$  های قبل از  $n_t$  را دربردارد.

حال ما می خواهیم آن ورودی ها را این لفظ یک برد خروجی یک

باشد را آن ورودی منفی شده خروجی منفی شده - اگر در نظر بگیریم

$$h_t = n_t \wedge h_{t-1} \quad \text{با فرض } h_0 = 1$$

$$w_1 = w_2 = 1 \quad b_2 = -\frac{3}{2}$$

تا وقتی که ورودی هایت باشد خروجی یک است :

$$f(n) = \begin{cases} 0 & n < 0 \\ 1 & n \geq 0 \end{cases} = \text{step}(n) \quad w_3 = 1 \quad b_3 = 0$$

$$g(n) = n \rightarrow \text{تابع هنجاری}$$

در خلاصه  $h_1$  سه ورودی دارد که آن‌ها عبارتند از  $m_1$  و  $m_2$  و شرطی که با این

$$\begin{matrix} m_1 & m_2 & h_1 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \end{matrix}$$

$$h(t) = \begin{bmatrix} h_1(t) \\ h_2(t) \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \begin{bmatrix} m_1^{(t)} \\ m_2^{(t)} \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

$$\Rightarrow h_1(t) = \phi(w_{11}h_1^{(t)} + w_{12}h_2^{(t)} + b_1) \quad w_{11} = 1 \quad w_{12} = -1 \quad b_1 = \frac{3}{2}$$

$$\begin{matrix} m_1 & m_2 & h_2 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \end{matrix}$$

$h_2$  هر ۲ تک بودن را تشخیص می‌دهد:

$$\left. \begin{matrix} \rightarrow \text{and} \rightarrow w_{21} = 1 \quad w_{22} = 1 \\ b_2 = -\frac{3}{2} \end{matrix} \right\}$$

$$\Rightarrow W = \begin{bmatrix} -1 & -1 \\ 1 & 1 \end{bmatrix} \quad b = \begin{bmatrix} \frac{3}{2} \\ -\frac{3}{2} \end{bmatrix}$$

$$g^{(t)} = \phi(v^T h^{(t)} + r y^{(t-1)} + c) \quad t \geq 1$$

$$\begin{matrix} h_1^{(t)} & h_2^{(t)} \\ 0 & 1 \\ 1 & 0 \\ 0 & 0 \end{matrix} \left. \begin{matrix} \rightarrow y^{(t-1)} = 1 \Rightarrow y^{(t)} = 1 \\ \rightarrow y^{(t-1)} = 0 \Rightarrow y^{(t)} = 0 \\ \rightarrow y^{(t)} = 0 \end{matrix} \right\}$$

$$v_1 = 1 \quad v_2 = 1 \quad r = 1 \quad c = -\frac{3}{2}$$

$$\Rightarrow y^{(t)} = \phi\left( \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} h_1^{(t)} \\ h_2^{(t)} \end{bmatrix} + y^{(t-1)} + \frac{-3}{2} \right)$$

$$y' = \phi(v^T h' + c) \Rightarrow \left. \begin{matrix} v^T h' = 1 \Rightarrow y' = 1 \\ v^T h' = 0 \Rightarrow y' = 0 \end{matrix} \right\} \Rightarrow c = -\frac{1}{2}$$