

(الف)

① با توجه به نمودار داده شده می توان دید که با شروع از مقدار  $h_0$  به سمت  $h_{\infty}$  می رسیم. همچنین صورت بردار آدرس  $h_1, h_2, \dots, h_{\infty}$  واضح است که  $h_{\infty}$  به مقدار  $0.18$  میل می کند. به عبارت دیگر  $h_0$  یکپیکر از  $0.18$  و پیکر آنگاه با وقت نمودار داده شده می توان دید که  $h_{\infty}$  به  $0.18$  میل می کند. همین روند برای وقتی که  $h_0$  دیگر از  $0.18$  بزرگتر یا کوچکتر نباشد.

ب) از آجایی که با انتخاب مقدار  $h_0$  در نهایت  $h_t$  به سمت  $0.18$  میل می کند. داریم

$$\frac{\partial h(t+1)}{\partial h(t)} = 2h(t+1)(1-h(t+1))$$

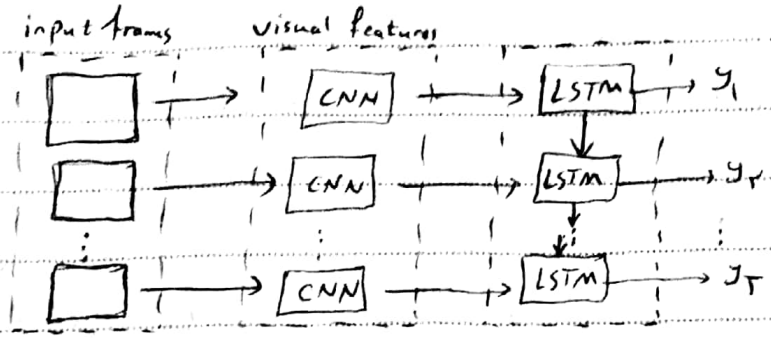
در این شبکه مشکل انتخاب یا حذف شدن گره ها در مخفی وجود دارد.

ج) همانگونه که توضیح داده شد مستقیماً می تواند که نقطه متادل  $h_t = 0.18$  می باشد. حال دایم اینکه مستقیماً کنیم این نقطه متادل Source است یا Sink باید مشخص دوم  $h_{t+1}$  نسبت به  $h_t$  را در نقطه  $0.18$  جای بگیریم. بعد از آن می بینیم مستقیماً می تواند که مستقیماً دوم در این نقطه متدل است. بعد از این نقطه  $0.18$  یک نقطه متادل Sink می باشد.

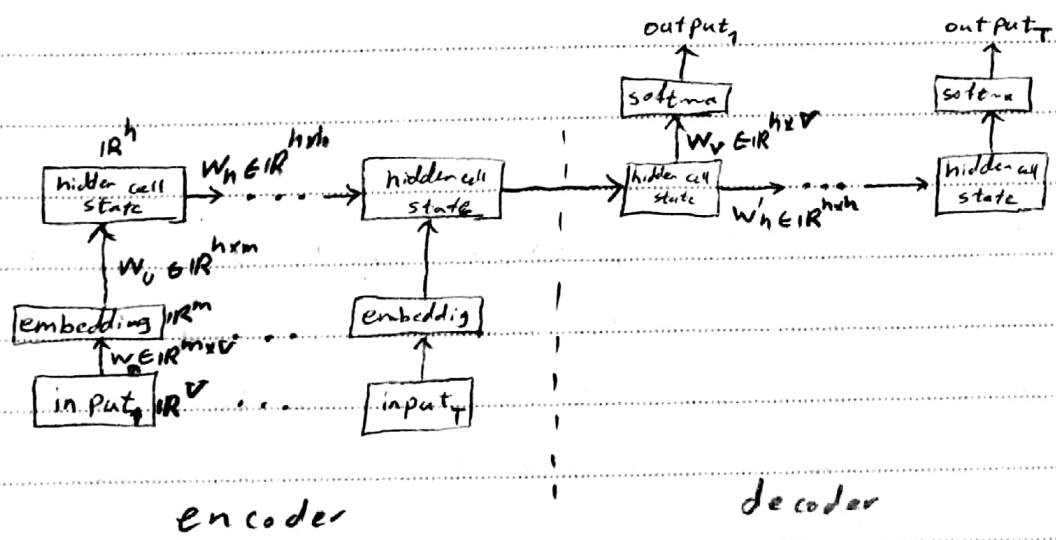
د) با توجه به رابطه  $h_{t+1} = \text{Relu}(2h_t - 1)$  مستقیماً است که اگر  $h_0$  بزرگتر از  $0.5$  باشد پیکر آنگاه با جای گذاری  $h_t$  به سمت  $\infty$  میل می کند و مشکل  $\text{explode}$  وجود دارد. حال اگر  $h_0$  کوچکتر از  $0.5$  باشد پیکر آنگاه با جای گذاری  $h_t$  به سمت  $0$  میل می کند و مشکل  $\text{Vanish}$  وجود دارد. به ازای  $h_0$  برای  $0.5$  مقدار  $h_t$  همواره برای  $0.5$  باقی می ماند.

$$h_t = \text{Relu}\left(2^t h_0 - \frac{2^t - 1}{2}\right)$$

② الف) وایان انجام شبکه در پردازش ویدیو می‌توان هر فریم از ویدیو را به یک شبکه CNN داد و سپس بردار بازخالی حاصل را به یک LSTM داد.



ب) کاری انجام این شبکه یک شبکه encoder decoder به صورت زیر می‌باشد. شبکه در ورودی را به encoder داده (البته ابتدا با استفاده از یک embedding شبکه در ورودی را به یک بردار تبدیل می‌کنیم). سپس <sup>cell state</sup> خط آفر این encoder را به عنوان ورودی شبکه decoder می‌دهیم. وظیفه این شبکه بازگشت تبدیل خودی از رشته حاصل زمان مشخص است. لازم به ذکر است که در شبکه های encoder - decoder گفته شده از طایفه های ~~RNN~~ استفاده می‌کنیم.



(i) برای محاسبه تعداد کل جوابات داریم:

$$\text{تعداد جوابات در سطح embedding} = m(V + V - 1) = m(2V - 1)$$

جمع      ضرب

تعداد جوابات در سطح hidden state

$$h_t = \sigma(W_U X + W_h h_{t-1} + \text{bias})$$

$\mathbb{R}^{hm} \rightarrow \mathbb{R}^m \rightarrow \mathbb{R}^{hv} \rightarrow \mathbb{R}^h \rightarrow \mathbb{R}^h$   
 $\downarrow \quad \quad \quad \downarrow \quad \quad \quad \downarrow$   
 $h(2m-1) \quad h(2h-1) \quad h$

Sigmoid تابع فرم =  $\frac{1}{1 + e^{-x}}$

$$\hookrightarrow 2mh + 2h^2 + 2h$$

بنابراین تعداد جوابات داده encoder عبارت است از:

$$[2mh + 2h^2 + 2h + 2mV - m]^T$$

برای محاسبه تعداد جوابات در سطح hidden state در decoder:

تابع نرمال سازی (softmax) داریم:

$$\text{softmax}, y_i = \text{softmax}(W_V h_i + \text{bias})$$

$\mathbb{R}^{hv} \rightarrow \mathbb{R}^h \rightarrow \mathbb{R}^h$   
 $\downarrow \quad \quad \quad \downarrow$   
 $V(2h-1) \quad V$

softmax تابع فرم =  $\frac{e^{x_i}}{\sum_j e^{x_j}}$

$$\hookrightarrow 2vh + 2v$$

بنابراین تعداد جوابات داده decoder عبارت است از:

$$[2h^2 + 2h + 2vh + 2v]^T$$

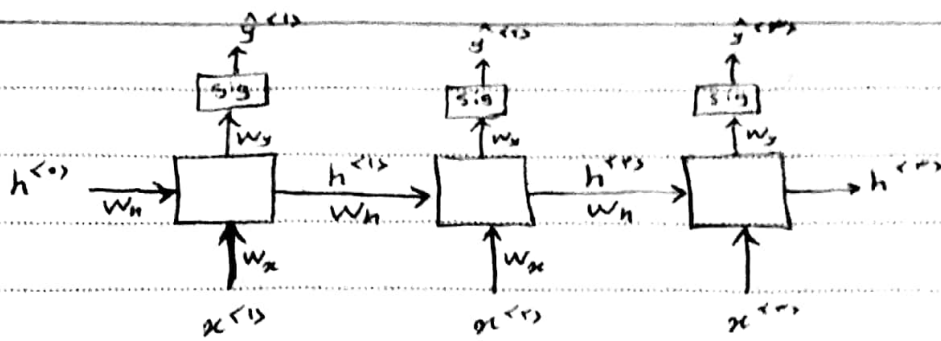
بنابراین تعداد کل جوابات عبارت است از:

$$[2h^2 + 2h + 2mh + 2mV + 2vh - m + 2v]^T$$

(ii) برای محاسبه تعداد بارها در سطح داریم:

<p>encoder بارها در سطح:</p> $m \times V + h \times m + h \times h + h$ <p style="text-align: center;"> <math>\underbrace{m \times V}_{W_m} \quad \underbrace{h \times m}_{W_U} \quad \underbrace{h \times h}_{W_h} \quad \underbrace{h}_{\text{bias}}</math> </p>	<p>decoder بارها در سطح:</p> $h \times h + h \times V + h + V$ <p style="text-align: center;"> <math>\underbrace{h \times h}_{W_h} \quad \underbrace{h \times V}_{W_V} \quad \underbrace{h}_{\text{bias}} \quad \underbrace{V}_{\text{hidden state}}</math> </p>
--	--

تعداد کل بارها =  $mV + mh + rh^2 + hv + rh + v$



(a) ③

$$h^{(t)} = w_h h^{(t-1)} + w_x x^{(t)} \quad \hat{y}^{(t)} = \sigma(w_y h^{(t)})$$

$$E^{(t)} = (y^{(t)} - \hat{y}^{(t)})^2$$

$$\left. \begin{aligned} E^{(r)} &= (y^{(r)} - \hat{y}^{(r)})^2 \\ \hat{y}^{(r)} &= \sigma(w_y h^{(r)}) \end{aligned} \right\} \Rightarrow E^{(r)} = (y^{(r)} - \sigma(w_y h^{(r)}))^2$$

از اینجا که  $w_y$  در تابع  $h^{(r)}$  نیستی، داریم:

$$\frac{\partial E^{(r)}}{\partial w_y} = -2(y^{(r)} - \sigma(w_y h^{(r)})) h^{(r)} \sigma(w_y h^{(r)}) (1 - \sigma(w_y h^{(r)}))$$

$$= -2(y^{(r)} - \hat{y}^{(r)}) \hat{y}^{(r)} (1 - \hat{y}^{(r)}) h^{(r)}$$

$$\frac{\partial E^{(r)}}{\partial w_h} = \frac{\partial E^{(r)}}{\partial \hat{y}^{(r)}} \frac{\partial \hat{y}^{(r)}}{\partial h^{(r)}} \left( \sum_{i=1}^r \frac{\partial h^{(i)}}{\partial h^{(r)}} \frac{\partial h^{(i)}}{\partial w_h} \right)$$

$$\frac{\partial h^{(r)}}{\partial w_h} = \sum_{i=1}^r \frac{\partial h^{(i)}}{\partial h^{(r)}} \frac{\partial h^{(i)}}{\partial w_h}$$

$$h^{(r)} = w_h h^{(r-1)} + w_x x^{(r)}$$

$$= w_h (w_h h^{(r-2)} + w_x x^{(r-1)}) + w_x x^{(r)}$$

$$= w_h^2 h^{(r-2)} + w_h w_x x^{(r-1)} + w_x x^{(r)}$$

$$= w_h^2 (w_h h^{(r-3)} + w_x x^{(r-2)}) + w_h w_x x^{(r-1)} + w_x x^{(r)}$$

$$= w_h^3 h^{(r-3)} + w_h^2 w_x x^{(r-2)} + w_h w_x x^{(r-1)} + w_x x^{(r)}$$



$$\frac{\partial h^{(r)}}{\partial h^{(i)}} = w_h^r, \quad \frac{\partial h^{(r)}}{\partial h^{(i)}} = w_h, \quad \frac{\partial h^{(r)}}{\partial h^{(i)}} = 1$$

$$\frac{\partial h^{(i)}}{\partial w_h} = h^{(i)}, \quad \frac{\partial h^{(i)}}{\partial w_h} = r w_h h^{(i)} + w_x x^{(i)}, \quad \frac{\partial h^{(i)}}{\partial w_h} = r w_h^r h^{(i)} + r w_h w_x x^{(i)} + w_x x^{(i)}$$

$$\Rightarrow \sum_{i=1}^r \frac{\partial h^{(i)}}{\partial h^{(i)}} \frac{\partial h^{(i)}}{\partial w_h} = w_h^r h^{(i)} + r w_h^r h^{(i)} + w_h x^{(i)} + r w_h^r h^{(i)} + r w_h w_x x^{(i)} + w_x x^{(i)}$$

$$= 4 w_h^r h^{(i)} + w_h (r w_x x^{(i)} + x^{(i)}) + w_x x^{(i)}$$

$$\frac{\partial E^{(r)}}{\partial \hat{y}^{(r)}} = r (\hat{y}^{(r)} - y^{(r)})$$

$$\frac{\partial \hat{y}^{(r)}}{\partial h^{(r)}} = w_y \hat{y}^{(r)} (1 - \hat{y}^{(r)})$$

$$\Rightarrow \frac{\partial E^{(r)}}{\partial w_h} = r w_y \hat{y}^{(r)} (\hat{y}^{(r)} - y^{(r)}) (1 - \hat{y}^{(r)}) (4 w_h^r h^{(i)} + w_h (r w_x x^{(i)} + x^{(i)}) + w_x x^{(i)})$$

$$\frac{\partial E^{(r)}}{\partial w_x} = \frac{\partial E^{(r)}}{\partial \hat{y}^{(r)}} \frac{\partial \hat{y}^{(r)}}{\partial h^{(r)}} \left( \sum_{i=1}^r \frac{\partial h^{(i)}}{\partial h^{(i)}} \frac{\partial h^{(i)}}{\partial w_x} \right)$$

$$\frac{\partial h^{(i)}}{\partial w_x} = x^{(i)}, \quad \frac{\partial h^{(i)}}{\partial w_x} = w_h x^{(i)} + x^{(i)}, \quad \frac{\partial h^{(i)}}{\partial w_x} = w_h^r x^{(i)} + w_h x^{(i)} + x^{(i)}$$

$$\Rightarrow \frac{\partial E^{(r)}}{\partial w_x} = r w_y \hat{y}^{(r)} (\hat{y}^{(r)} - y^{(r)}) (1 - \hat{y}^{(r)}) (3 w_h^r x^{(i)} + 2 w_h x^{(i)} + x^{(i)})$$

ب) مشاهده کردیم که در محاسبه گرادیان ها، در مرحله Backprop برای این گرادیان یک نقطه متغیر، ماتریس  $w_h$

در محاسبات ضرب ماتریسی در این زمان که گرادیان ها در محاسبات درگیر می شوند، در گرادیان حاصل یک عامل ضرب می شود.

دایم: در واقع:  $h^{(t)} \propto W_h^t h^{(0)}$  که می توان نوشت:  $h^{(t)} \propto Q^T \Lambda^t Q h^{(0)}$

حال اگر بزرگترین singular value ماتریس  $w_h$  بزرگتر از یک باشد  $\leftarrow$  explode و اگر کوچکتر از یک باشد  $\leftarrow$

Vanish

$$h^{(t)} = W^T h^{(t-1)}$$

$$W = \begin{bmatrix} 0.08 & 0.12 \\ 0.12 & 0.17 \end{bmatrix}$$

$$h^{(0)} = I_r \quad (f)$$

$$\det(W - \lambda I) = 0$$

(w)

$$(0.08 - \lambda)(0.17 - \lambda) - 0.12^2 = 0 \Rightarrow \lambda_1 = 0.2, \lambda_2 = -0.9$$

$$(W - \lambda_1 I) x_1 = 0 \Rightarrow \begin{bmatrix} 0.18 & 0.12 \\ 0.12 & 0.17 \end{bmatrix} x_1 = 0$$

$$\Rightarrow x_1 = \begin{bmatrix} -0.18 \\ 0.17 \end{bmatrix}$$

$$(W - \lambda_2 I) x_2 = 0 \Rightarrow \begin{bmatrix} -0.12 & 0.12 \\ 0.12 & -0.18 \end{bmatrix} x_2 = 0$$

$$\Rightarrow x_2 = \begin{bmatrix} 0.17 \\ 0.18 \end{bmatrix}$$

$$W = Q \Lambda Q^{-1}$$

$$Q = \begin{bmatrix} -0.18 & 0.17 \\ 0.17 & 0.18 \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} 0.2 & 0 \\ 0 & -0.9 \end{bmatrix}$$

$$\Rightarrow Q^{-1} = Q^T \Rightarrow W = Q \Lambda Q^T \Rightarrow W^T = Q^T \Lambda Q$$

$$h^{(t)} = W^T h^{(t-1)} \Rightarrow h^{(t)} = (Q^T \Lambda Q) h^{(t-1)}$$

$$\Rightarrow h^{(r)} = Q^T \Lambda^r Q = \begin{bmatrix} -0.18 & 0.17 \\ 0.17 & 0.18 \end{bmatrix} \begin{bmatrix} 0.2^r & 0 \\ 0 & -0.9^r \end{bmatrix} \begin{bmatrix} -0.18 & 0.17 \\ 0.17 & 0.18 \end{bmatrix}$$

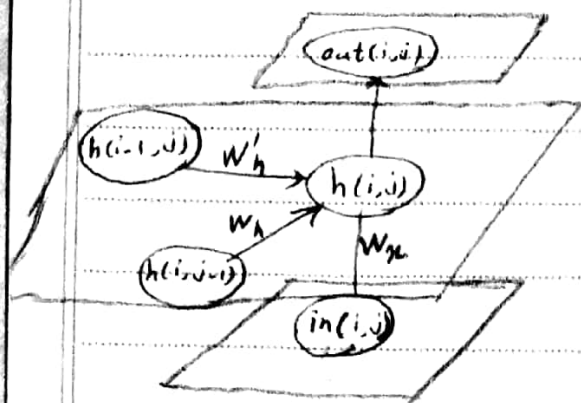
$$h^{<t>} = Q^T \Lambda^t Q$$

ب) چنانکه که مشاهده می شود داریم:

از آنجایی که داریم  $\Lambda = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0.4 & 0 \\ 0 & 0 & 0.19 \end{bmatrix}$  بنابراین سانس  $\Lambda^t$  با بزرگ شدن  $t$  به

~~سانس~~ سانس صفر میل می کند. بنابراین می توان گفت این نوع از شبکه بازگشتی با افزایش تعداد تکراری ورودی غنی نمائند خردی مناسب تولید کند. همچنین به علت مناسب در فاز Backward باعث Vanish گون گرادیان و عدم آپدیت وزن ها می شود.

② در شبکه های RNN یک بعدی، توانایی پردازش داده های نقطه یک بعدی مثل توالی از کلمات یا توالی از اعداد وجود داشت و امکان پردازش داده های چند بعدی مثل تصویری و ویدیو وجود نداشت. شبکه های MDRNN قابلیت پردازش داده های چند بعدی را دارند. ایده ی آن ها به این صورت است که یک اتصال بازگشتی واحد که در RNN ها وجود داشت را با  $n$  اتصال بازگشتی جایگزین کنند که  $n$  برای تعداد بعد ورودی است. در این نوع شبکه ها در زمان forward علاوه بر ورودی هر لحظه، خروجی های لایه های پنهان هر یک از بعد های دیگر نیز در فصلی بازنهاده می شوند. به عنوان مثال اگر داده ورودی ۲ بعدی باشد ( $n=2$ )، در هر استیپ مقدار hidden state در هر نقطه به صورت زیر است:

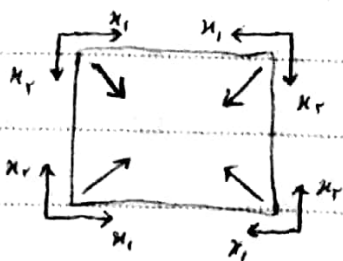


$$h(i,n) = w_x x(i,n) + \underbrace{w_h' h(i-1,n)}_{\text{if } i > 0} + \underbrace{w_h h(i,n-1)}_{\text{if } n > 0}$$

اگر بخواهیم مفهوم MDRNN را با مفهوم BLSTM ترکیب کنیم به مفهوم Multi Dimensional Bidirectional LSTM می رسیم.

BLSTM به این صورت است که در مرحله forward داده های توالی ورودی هم از اول به سمت آخر و هم از آخر به سمت اول پردازش می شوند. برای ساختن خروجی در هر لحظه از دو hidden state که با هم مرتبط به جهت  $\rightarrow$  و  $\leftarrow$  استفاده می کنیم.

حال ساختار Multi Dimensional LSTM برای صورتی است که در هر LSTM Unit برای یک جهت target به تعداد ابعاد ورودی از این نوع گیت داریم. حال اگر بخواهیم ساختار 2D BLSTM را شرح دهیم به این صورت است که در هر لایه  $2^n = 2^2 = 4$  hidden state وجود دارد که جهت های خاص به آن ها به صورت زیر است:





ج) برای محاسبه مرحله forward در MORNN، ابتدا وزنهای ورودی  $W_x$  را در هر یک از داده‌های ورودی ضرب می‌کنیم. سپس به ازای هر یک از این داده‌ها، ماتریس وزن  $W_h$  را در خروجی hidden state آن داده در بعد گسری ضرب می‌کنیم. به این نحوی که مجموع این ضربات  $hidden\ state$  برای داده در بعد فعلی محاسبه می‌شود.  $psedu\ code$  مربوط به صورت زیر است:

```

for  $i_1 = 0$  to  $I_1 - 1$ 
  for  $i_2 = 0$  to  $I_2 - 1$ 
    ...
    for  $i_n = 0$  to  $I_n - 1$ 
       $b \leftarrow W_x \times (i_1, i_2, \dots, i_n)$ 
      for  $j = 1$  to  $n$ 
        if  $i_j > 0$ 
           $b \leftarrow b + W_h^{ij} h(i_1, \dots, i_{j-1}, \dots, i_n)$ 
       $h(i_1, i_2, \dots, i_n) \leftarrow g(b)$ 
      تابع فعال‌ساز
  
```

\* از آنجایی که  $n$  بعد داریم و برای هر محاسبه  $hidden\ state$  در هر لحظه نیاز به پردازش تمامی  $hidden\ state$  ها، فرای برد قبل تمامی  $n$  بعد داریم. بنابراین وزنهایی که وجود دارند عبارت است از یک وزن  $W_x$  که عضو  $R^{in \times h}$  است،  $n$  وزن  $W_h$  نیاز که هر کدام عضو  $R^{h \times h}$  هستند و یک وزن  $W_y$  که عضو  $R^{h \times out}$  است. به صورت مشابه یک لایس  $h$  و  $n$  لایس  $h$  یک لایس  $out$  داریم. بنابراین تعداد کل پارامترها عبارت است از:

$$in \times h + n \times (h \times h) + h \times out + h + n \times h + out$$

حال اگر می‌خواهیم سوزی داشته باشیم می‌توانیم محاسبات مربوط به (اعداد مختلف را به یک سوزی متناظر و نگهداریم). تا این عملیات به صورت سوزی در همین انجام شوند. از طرفی می‌توانیم عملیات مربوط به پردازش ورودی در لحظات مختلف را نیز سوزی انجام داد. نقطه عملیات مربوط به محاسبه  $hidden\ state$  را می‌توان سوزی انجام داد زیرا به لحظه قبل بستگی دارد. بنابراین صورت زیرین اجرای کلی متناظر است با جدول دنباله ورودی