

سوال ①

می‌خواهیم در RNN داده شده خروجی y_t در آخرین لحظه زمانی را استخراج کنیم، همین است در هر است زمان به صورت زیر آپدیت می‌شود.

$$h_t = h_{t-1} + x_t$$

دری در لحظه t

 $h_0 = \text{initial hidden state}$

$$\rightarrow y_t = \sigma(10 \cdot h_t) \rightarrow \text{final time step}$$

$$\begin{aligned} \text{RNN recurrence: } h_1 &= -h_0 + x_1 \\ h_2 &= -(-h_0 + x_1) + x_2 = h_0 - x_1 + x_2 \\ h_3 &= -h_0 + x_1 - x_2 + x_3 \\ &\vdots \\ h_t &= (-1)^t h_0 + \sum_{i=1}^t (-1)^{t-i} x_i \end{aligned}$$

$$T \text{ خروجی در } T = y_T = \sigma(10 \cdot (h_0 + \sum_{i=1}^T (-1)^{T-i} x_i))$$

زوج \rightarrow

سوال ②

① ایجاد بالا و پراکنده گی: وکتورهای one-hot دارای ابعاد بسیار بالا هستند (طول وکتور برابر با تعداد واژگان است). این ویژگی باعث ناکارآمدی در محاسبات و محذوری می‌شود. (مخصوصاً با افزایش واژگان)

② هارد-کود: یعنی اضافه کردن یک واژه جدید نیازمند افزودن یک بعد جدید به وکتور است. این ویژگی باعث می‌شود one-hot برای واژگان بزرگ یا داینامیک انعطاف پذیری و مقیاس پذیری نداشته باشد.

③ ناهمبستگی: محاسبات را لحاظ نمی‌کنند. (مثل امپدینگ نیست)

سوال ③

$$h_t = \sigma(W h_{t-1} + U x_t) \rightarrow h_{t+1} = \sigma(W h_t + U x_{t+1})$$

$$\rightarrow \frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_t} \quad \text{زنجیره ای} \quad \rightarrow \frac{\partial h_{t+1}}{\partial h_t} = (\sigma'(W h_t + U x_{t+1})) \cdot W$$

$\hookrightarrow \sigma' = \sigma(1 - \sigma)$

$$\rightarrow \frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+1}} \bullet \left[(\sigma(W h_t + U x_{t+1})) \odot (1 - \sigma(W h_t + U x_{t+1})) \right] \cdot W$$

$\hookrightarrow \text{Point wise prod}$

$$\rightarrow \frac{\partial h_{t+1}}{\partial h_t} = \text{diag}(\sigma'(W h_t + U x_{t+1})) \cdot W$$

(ب) البته به بخیر الف:

$$\rightarrow \frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_T} \prod_{k=t}^{T-1} \frac{\partial h_{k+1}}{\partial h_k} \quad \rightarrow \frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_T} \prod_{k=t}^{T-1} \left[(\sigma'(W h_k + U x_{k+1})) \cdot W \right]$$

$$\frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_T} = \frac{\partial L}{\partial h_T} \cdot \prod_{k=t}^{T-1} (w \cdot \sigma'(wh_k + Ux_{k+1})) \quad \text{سوال 3} \quad (1)$$

$$\text{maximizing } \sigma'(z) = \sigma(z)(1-\sigma(z)) \xrightarrow[\sigma(z)=\frac{1}{2}]{\text{quadratic form}} \max(\sigma'(z)) = \frac{1}{2} \left(1 - \frac{1}{2}\right) = \frac{1}{4}$$

$$\rightarrow \frac{\partial L}{\partial h_t} \leq \left| \frac{\partial L}{\partial h_T} \right| \cdot \prod_{k=t}^{T-1} |w| \cdot \frac{1}{4} \Rightarrow \left| \frac{\partial L}{\partial h_t} \right| \leq \left| \frac{\partial L}{\partial h_T} \right| \cdot |w|^{T-t} \cdot \left(\frac{1}{4}\right)^{T-t}$$

$$\frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+2}} \cdot \frac{\partial h_{t+2}}{\partial h_t} + \frac{\partial L}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial h_t}$$

$$\frac{\partial h_{t+2}}{\partial h_t} = \sigma'(wh_{t+1} + Mh_t + Ux_{t+2}) \cdot M + \sigma'(wh_{t+1} + Mh_t + Ux_{t+2}) \cdot w \cdot \frac{\partial h_{t+1}}{\partial h_t}$$

$$\rightarrow \frac{\partial h_{t+1}}{\partial h_t} = \sigma'(wh_t + Mh_{t-1} + Ux_{t+1}) \cdot w$$

$$\rightarrow \frac{\partial h_{t+2}}{\partial h_t} = \sigma'(wh_{t+1} + Mh_t + Ux_{t+2}) \cdot M + \sigma'(wh_{t+1} + Mh_t + Ux_{t+2}) \cdot w \cdot \sigma'(wh_t + Mh_{t-1} + Ux_{t+1}) \cdot w$$

$$\rightarrow \frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+2}} \cdot \left(\sigma'(wh_{t+1} + Mh_t + Ux_{t+2}) \cdot M + \sigma'(wh_{t+1} + Mh_t + Ux_{t+2}) \cdot w \cdot \sigma'(wh_t + Mh_{t-1} + Ux_{t+1}) \cdot w \right) + \frac{\partial L}{\partial h_{t+1}} \cdot \sigma'(wh_t + Mh_{t-1} + Ux_{t+1}) \cdot w$$

ترج M باعث ایجاد یک کانکشن مستقیم (skip connection) می شود و مدارات تکرار ضرب های درونی که در میانه و بایک کاهش یک گره ها و اینست که چون ضرب گره ها قبل ندارد (گره ها این که از توابع فعال سازی ناشی می شوند) در کل با ایجاد جریان مستقیم گره ها باعث می شود کاهش شدت ناک باشد و باعث حفظ اطلاعات قبل به تسهیل یادگیری می شود. (بابت بهینه سازی و بایک توانایی های بلند تر و RNN های عمیق تر استفاده شود.

ج) نوعی ای گزاردان کسینک :

clipping by value : در این نوعی اگر هر صولفه گزاردان از یک ترسهولده مشخص بیشتر شود، مقدار آن به همان ترسهولده محدود می شود.

clipping by norm : در این نوعی اگر نرخ کل بردار گزاردان از یک مقدار آستانه بیشتر شود، بردار گزاردان نزال سازی می شود. یعنی اگر نرخ بردار گزاردان بزرگتر از ترسهولده باشد کل بردار مقیاس بندی شده تا نرخ آن برابر ترسهولده شود. ~~(در اینجا مقیاس بندی می شود)~~
حزایای by norm نسبت به by value :

۱) حفظ جهت گزاردان : در روش اول هر صولفه به صرح جداگانه تفسیر می کند که ممکن است جهت گزاردان را تغییر دهد. اما در روشی اول چون کل بردار گزاردان مقیاس بندی می شود جهت اصلی حفظ می شود.

2) همگرایی پایدار تری خواهیم داشت و از بایاس یا عدم تعادل ناشی از clipping انتخابی جلوگیری می شود.

آ) مشکل آموزش تولید توانی در این مدل که مدل خروجی خود را به عنوان ورودی مرحله بعد استفاده می کند. اگر مدل در مرحله اولیه خطایی کند این خطا در مراحل بعدی تشدید می شود و محاسبه دقتی مدل را کاهش می دهد.

برای حل این مشکل از روش teacher force استفاده می شود. به جای استفاده از پیش بینی های مدل برای مرحله بعد، از خروجی های درست داده که train استفاده می شود. (grand truth) این کار به مدل کمک می کند تا با استفاده از زمینه صحیح آموزش ببیند و خطا انباشته نشود.

ب) مشکل اصلی Teacher Forcing چیزی به نام Exposure bias است. در زمان آموزش مدل همیشه از خروجی های درست (ground truth)

به عنوان ورودی مرحله بعد استفاده می کند. اما در زمان تست، مدل به این خروجی های دسترس نمی ندارد و باید به پیش بینی های خود تکیه کند. این تفاوت بین محیط تست و تریین مایه های می شود که مدل اگر خطای کوچک در تست مرتکب شود نتواند آن را جبران کند، چون برای مدیریت خطای خود آموزش ندیده است. این امر به تجمع خطا در طول توانی منجر می شود.

ج) برای حل مشکل bias exposure می توان از تکنیکی به نام schedule sampling استفاده کرد. در این روش مدل به تدریج به پیش بینی های خودش به جای خروجی های درست تکیه می کند. به این صورت که با پیشرفت آموزش، احتمال استفاده از پیش بینی های مدل افزایش می یابد. این انتقال تدریجی به مدل کمک می کند تا با خطای خود کار کند و برای مرحله ی تست که باید کاملاً به پیش بینی های خودش متکی باشد آماده شود. این روش اختلاف بین محیط های train, test را کاهش داده و اثرات bias exposure را به حداقل می رساند.

$$\begin{cases} y_t = 1 & \text{if previous inputs including current one is 1} \\ y_t = 0 & \text{o.w.} \end{cases}$$

$$h_t = h_{t-1} \cdot x_t$$

$$h_t = \begin{cases} 1 \\ 0 \end{cases} \quad \text{with the same condition as } y_t.$$

the output defined $\Rightarrow y_t = h_t$, $h_0 = 1$

$$\text{update state: } h_t = 1 \Rightarrow h_t = f(w_2 h_{t-1} + w_1 x_t + b_1)$$

$$f(z) = \begin{cases} 1 & z \geq 0.5 \\ 0 & z < 0.5 \end{cases}$$

output

$$\rightarrow y_t = w_3 h_t + b_2$$

it is similar to and gate $\rightarrow \begin{cases} w_1 = 1 \\ w_2 = 1 \\ b_1 = -1.5 \end{cases}$ only if both are one

$w_2 = 1$ hidden state weight

$w_1 = 1$: input weight

$b_1 = -1.5$: for \otimes

$w_3 = 1 \rightarrow$ directly shows output
 $b_2 = 0 \rightarrow$ no need

$$x_t = 1110100$$

3 evaluate

time steps \Rightarrow	1	2	3	4	5	6	7	8	...
input	1	1	1	1	0	1	0	0	...
$h_t = h_{t-1} \cdot x_t$	1.0	1.0	1.0	1.0	1.0	0.0	0.0	0.0	...
$y_t = h_t$	1	1	1	1	0	0	0	0	...

خروجی جمع و جمع به یکدیگر نماند. بار آخر که صحیح اند با تریج به evaluation، منطبق.

وردی در توالی
باینری است. $\begin{cases} x_1^{(t)} \\ x_2^{(t)} \end{cases}$ value of two sequence at time

سوال 6

$$y^{(t)} = \begin{cases} 1 \\ 0 \end{cases} \quad \text{توالی یک و صفر است} \quad \text{o.w.}$$

$$\begin{aligned} h_1^{(t)} &= 1 & \text{اگر } x_1^{(t)} = 0 \text{ و } x_2^{(t)} = 0 &\rightarrow \text{o.w.} \Rightarrow 0 = h_1^{(t)} \\ h_2^{(t)} &= 1 & \text{اگر } x_1^{(t)} = 1 \text{ و } x_2^{(t)} = 1 &\rightarrow \text{o.w.} \Rightarrow 0 = h_2^{(t)} \end{aligned}$$

$$\text{تابع فعال: } h^{(t)} = \phi(w x^{(t)} + b)$$

$$x^{(t)} = \begin{bmatrix} x_1^{(t)} \\ x_2^{(t)} \end{bmatrix}, W^{2 \times 2} \rightarrow W = \begin{bmatrix} -1 & -1 \\ 1 & 1 \end{bmatrix}, b = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\text{For } h_1^{(t)} \Rightarrow \text{input } Wx^{(t)} + b = -x_1^{(t)} - x_2^{(t)} + 1 > 0 \Leftrightarrow x_1^{(t)} = x_2^{(t)} = 0 \text{ اگر}$$

$$\text{For } h_2^{(t)} \Rightarrow \text{input } Wx^{(t)} + b = x_1^{(t)} + x_2^{(t)} - 1 > 0 \Leftrightarrow x_1^{(t)} = 1 = x_2^{(t)} \text{ اگر}$$

$y^{(t)} \propto y^{(t-1)}$ & $h^{(t)}$ & activation func of current hidden layer

$$\text{so } \rightarrow y^{(t)} = \begin{cases} \phi(v^T h^{(t)} + r y^{(t-1)} + c) & t > 1 \\ \phi(v^T h^{(t)} + c_0) & \end{cases}$$

$$v = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \begin{matrix} r=1 \\ c=-1 \\ c_0=0 \end{matrix}$$

$$\text{hidden layer } h^{(t)} \rightarrow x_1^{(t)} = x_2^{(t)}$$

output $y^{(t)} \rightarrow$ propagates the result of the comparison across time steps.
 $\hookrightarrow t=1$ start

$$\phi(z) = u(z) \quad \text{step func}$$

$$W = \begin{bmatrix} -1 & -1 \\ 1 & 1 \end{bmatrix} \quad b = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \quad r = 1 \quad c = -1 \quad c_0 = 0$$

simple initialization

بازمانده می‌دهی تا آنکه

خروجی را به خونی اجرا کنه \rightarrow same sequence $\rightarrow 1$ and different sequence $\rightarrow 0$