



Homework 5

Deep learning
Dr. Mohammadi

مشخصات

نام و نام خانوادگی	رحمت اله انصاری
شماره دانشجویی	404722028
رشته و مقطع و دانشگاه	کارشناسی ارشد هوش مصنوعی دانشگاه علم و صنعت ایران
ایمیل	rahmat.ansari.dev@gmail.com

ساختار فایل‌های پاسخ

- فایل answers، که همین فایل است، شامل
 - مشخصات دانشجو
 - پاسخ سوالات تشریحی ۱ تا ۴ (اسکن شده با کم اسکنر)
 - گزارش سوال عملی ۶
 - گزارش سوال عملی ۷
- پوشه Practical 1
 - فایل notebook1-hw5-students.ipynb
- پوشه Practical 2
 - فایل hw5-problem6.ipynb
- پوشه helper (شامل فایل‌های کمکی، در این تمرین این فایل‌ها مربوط به سوالات ۱ و ۲ تشریحی است)
 - فایل hw5_p1.ipynb
 - فایل hw5_p2.ipynb

۱/۱ = اتم / فرضی ، امان پابلیشرز

weights = $H \times D$, $H \times H$.

Blases = $H + H$

Control of the job is done by the Batch size, n .

RNNs (input=64, hidden=32) 1 < 2

batch_first = True (B, H, W) : original

...എ.ആർ.

Wib $\approx 32 \times 64 = 2048$

$W_{hh} = 32 \times 32 = 1024$

جواب : 3072

Ques

$$b_{ih} + b_{wh} = 32 + 32 = 64$$

$$3072 + 64 = 3136$$

20مارچ 1964

RNN (input = 32, hidden = 64) z, y

اندازه فرعی: $(B, 50, 64)$ فرعی تمام‌گامی نهانی

(last hidden state) z_{n-1} is linear: $z_n = W z_{n-1} + b$

(B, B_4) உடன் ϕ ; பிழைப்பு குறியீடு $C_1 = \text{வரி}$

$$w_{ih} \Rightarrow 64 \times 32 = 2048$$

وزن ها:

$$w_{hh} \Rightarrow 64 \times 64 = 4096$$

$$b_{ih} \Rightarrow 6144$$

بایاس ها:

$$b_{ih} + b_{hh} = 64 + 64 = 128$$

$$6144 + 128 = 6272$$

کل بایاس ها:

Linear (in=64, out=512) 3 لایه

اندازه درونی = (B, 512)

$$64 \times 512 = 32768$$

بایاس ها:

$$32768 + 512 = 33280$$

Linear (in=512, out=10) 4 لایه

اندازه درونی = (B, 10)

$$512 \times 10 = 5120$$

بایاس ها:

$$5120 + 10 = 5130$$

(FLOPs) تعداد عملیات ریاضی و جمع

در اینجا فرض می‌کنیم که هر عملیات جمع و ضرب یک FLOP است.

فرمول تقریبی برای محاسبه FLOPs در RNN:

$$\text{operations} \approx 2 \times H \times (D + H) + H$$

125

உலகம்

$$(32 \times 64) + (32 \times 32) = 3072 \text{ words}$$

دفعه ۱۰ (شماره ۱۵۵ و ۱۵۶) 3072 + 64 = 3136

6208 : کمرہ کا نام

$\therefore \text{Mare } 6208 \times 50 = 310,400 \quad \therefore (\text{seq. no.}) \text{ of } \text{Feb. 1, 2017}$

2nd.

செய்யுள்

$$(64 \times 32) + (64 \times 64) = 6144 \quad \text{is less}$$

$$6144 + 128 = 6272 \quad \text{WJA}$$

12416 26890

$\text{€ Mac } 12416 \times 50 = 620,800$: 14,000

3. 24.

$$64 \times 512 = 32768 \quad \text{2 Log in}$$

$$32768 + 512 = 33280 = 1080$$

66,048 JB

4. 28. a

$$512 \times 10 = 5120 \quad \therefore \text{Yes}$$

$$5120 + 10 = 5130 \quad 1600$$

10250 25

① اگر طول دنباله دو برابر شود (100 شود)

• تعداد پارامترها: هیچ تغییری نمی‌کند (وزن‌های RNN است)

• محاسبه (FLOPs): برای RNN دقیقاً دو برابر می‌شود (چون

مانده زبانی 100 پارامتر است) (لایه‌های Linear تغییری نمی‌کند)

(چون روی بردارهای اعمال می‌شوند)

② اگر سایه بردار ورودی به 128 تغییر کند =

• 8 به 1: تعداد پارامترها و FLOPs افزایش می‌یابد (چون ماتریس

از 32×64 به 32×128 تبدیل می‌شود)

• سایه بردارها: هیچ تغییری نمی‌کند (چون خروجی لایه 1 همیشه

32 باقی می‌ماند و ورودی‌های لایه 2 دست‌خوش تغییری نمی‌شوند)

③ اگر لایه دوم RNN به BiRNN تبدیل شود

• پارامترهای لایه 2: دو برابر می‌شود (چون یک شبکه عصبی کامل به لایه 2

مکروس اضافه می‌شود)

• خروجی لایه 2: سایه بردارهای خروجی دو برابر می‌شود

($64 \times 2 = 128$) چون خروجی رفت و برگشت به هم پیوسته

(Concatenate) می‌شوند

• لایه 3: (تأثیر نمی‌بیند) چون ورودی لایه 3 128

مانده است (از 64) تعداد پارامترهای لایه 3 (مانده است)

وزن) و بایاس می‌شود

④ اگر لایه خروجی حذف شود:

• محاسبه و بایاس حذف می‌شود و اینک حذف تغییر می‌کند

فروپسیدگی‌های ۱۱ کلاس یک بار ویدئویی ۵۱۲ عرض می‌دهد
 بود (که فروپسیدگی ۵۱۳ است)
 و تعداد بار است که به اندازه ۵۱۳۰ و در بار است که ۵۱۴ (۴) کاهش
 می‌یابد

2/2

نوع داده	RNW (مترایا/مکایب)	GNW (مترایا/مکایب)
دنباله ساده (مستقیم)	مزیت: درک عالی تدریب و توانایی کمالات عیب: کند بودن به دلیل عدم امکان موازی سازی	مزیت: سرعت بالا و استخراج ویژگی‌های کلی‌تر عیب: عدم درک وابستگی‌های پیچیده و ترتیب دقیق
تصویر	مزیت: نیاز (نگارهای پیچیده) به ردیف نگار (آموزش است) عیب: از دست دادن ساختار مکانی (spatial) با فلت کردن	مزیت: حفظ ساختار مکانی و ویژگی‌های محلی عیب: نیاز به داده زیاد برای آموزش
صوت	مزیت: مدل سازی تغییرات زمانی سیگنال عیب: مشکل در پردازش فیلترهای صوتی فیلتر طولانی	مزیت: پردازش اسپکتروگرام‌ها به عنوان تصویر (تشفیر آنگو و فیلترهای کارسین) عیب: تلفات در درک توانی زمانی بلند
روابط زمانی	مزیت: (تشریح) حفظ حافظه عیب: در عمل مشکل محو شدن همه ادیان دارد	مزیت: با استفاده از Dilated Convolutions می‌تواند بازه دید وسیع داشته باشد عیب: محدود به اندازه Receptive Field

ب) چراغهای Tanh

در RNN ها از Tanh استفاده می شود چون

- ① کنترل خروجی: خروجی را بین -1 و 1 نگه می دارد و مانع از بزرگ شدن و کوچک شدن آن می شود (جلوگیری از انفجار و منجمد شدن)
 - ② همادینگی: مشتق آن در هر نقطه قوی است و به یادگیری کمک می کند.
- در مدل LSTM به نسبت چون Cell خروجی مثبت را می خورد و منفی را در یک سلول دیگر می بیند، برای جلوگیری از خروجی منفی و این مشکل را به کمک مقدار محدودیت بین سلول ها حل می کنند (exploding activation)

ع = مقدار ثابت (Forward Pass)

برای محاسبه مقادیر h_{t-1} و y_{t-1} به سلول های قبلی مراجعه می کنیم

مثال: ~~محاسبه h_t و y_t در زمان t با استفاده از مقادیر h_{t-1} و y_{t-1} در زمان $t-1$~~

فرمول کلی برای RNN ها به صورت زیر است:

$$h_t = \tanh(w_{hh} h_{t-1} + w_{xh} x_t + b_h)$$

$$y_t = w_{hy} h_t + b_y$$

در اینجا چون ورودی ها به صورت 1 و 0 هستند، h_{t-1} و h_{t-2} و x_{t-1} را می بینیم.

گام 1: محاسبه h_{t-1} (از حالت متوقف قبلی)

مقدار h_{t-1} و h_{t-2} و w_{hh} در زمان $t-1$ را می بینیم.

$$w_{hh} h_{t-2} = \begin{bmatrix} 0.3 & 0.1 \\ -0.4 & 0.2 \end{bmatrix} \begin{bmatrix} -0.2 \\ 0.4 \end{bmatrix} = \begin{bmatrix} (0.3)(-0.2) + (0.1)(0.4) \\ (-0.4)(-0.2) + (0.2)(0.4) \end{bmatrix} = \begin{bmatrix} -0.02 \\ 0.16 \end{bmatrix}$$

گام 2: محاسبه خروجی ورودی

محاسبه خروجی ورودی در گام 2: x_{t-1}

$$w_{xh} \times x_{t-1} = \begin{bmatrix} 0.5 & -0.2 \\ -1.0 & 0.9 \end{bmatrix} \begin{bmatrix} 0.8 \\ -0.5 \end{bmatrix} = \begin{bmatrix} 0.5 \times 0.8 + (-0.2)(-0.5) \\ (-1.0)(0.8) + (0.9)(-0.5) \end{bmatrix} = \begin{bmatrix} 0.5 \\ -1.25 \end{bmatrix}$$

گام 3: جمع بایاس و حاصل تابع فعال‌سازی (Tanh)

تمام مقادیر (نتایج بالا + بایاس b_h) را به تابع فعال‌سازی می‌دهیم

تابع فعال‌سازی (Z) می‌باشد:

$$Z = \begin{bmatrix} -0.02 \\ 0.16 \end{bmatrix} + \begin{bmatrix} 0.5 \\ -1.25 \end{bmatrix} + \begin{bmatrix} 0.1 \\ -0.1 \end{bmatrix} = \begin{bmatrix} 0.58 \\ -1.19 \end{bmatrix}$$

حالا تابع \tanh (تنگ) را به Z می‌دهیم تا خروجی نهایی را بدست آوریم

محاسبه h_{t-1}

$$h_{t-1} = \tanh(Z) = \begin{bmatrix} \tanh(0.58) \\ \tanh(-1.19) \end{bmatrix} = \begin{bmatrix} 0.5227 \\ -0.8305 \end{bmatrix}$$

گام 4: محاسبه خروجی (logits)

در این مرحله حالت نهایی را به w_{hy} و b_y می‌دهیم

محاسبه \hat{y}_{t-1}

$$\hat{y}_{t-1} = w_{hy} h_{t-1} + b_y = \begin{bmatrix} 1.5 & -0.5 \\ 0.5 & 1.0 \end{bmatrix} \begin{bmatrix} 0.5227 \\ -0.8305 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 1.199 \\ -0.369 \end{bmatrix}$$

نتیجه بهای:

مطرحات معنی: $h_{t-1} = [0.52, -0.83]^T$

پرداز فراموشی: $\hat{y}_{t-1} = [1.20, -0.57]^T$

(د) چرا RNN وابستگی طولانی را یاد نمی‌گیرد؟

به دلیل صاف شدن گرادیان (vanishing gradient) وقتی که فاصله در طول زمان به عقب می‌رود. (Backpropagation through time)
 شاهد کوچکی داریم: مقدار آن به سمت صفر میل می‌کند و وابستگی
 ابتدایی می‌بینیم یاد نمی‌گیرد.

(ه) کاربرد دسته‌بندی

one-to-one: طبقه‌بندی تصویر ساده (پرونده‌های)

one-to-many: توصیف تصویر (image captioning) - ورودی تصویر، خروجی جمله
 many-to-one: تحلیل احساسات (sentiment Analysis) - ورودی متن، خروجی مثبت/منفی
 many-to-many (همزمان): برچسب زنی اجزای کلام (Pos Tagging) - در متن یا ترجمه ماشین

(و) تفاوت Encoding و Embedding

Encoding: (یک hot one): برداری پراکنده (sparse) از یک بردار بالا به پایین

و غیر قابل تغییر است. معنی کلمات را دربردارد

Embedding: (word 2 vec): برداری فشرده (Dense). از یک کلمه، قابل
 تغییر است و کلمات به معنی مشابه به یکدیگر نزدیک دارند.

(ز) مزیت BiRNN (دو سوییچ): در RNN معمولی فقط گذشته را می‌بینیم.

در BiRNN یک همزمان متن را از اول و آخر و از آخر به اول می‌خواند.

این باعث می شود برای درک یک کلمه هم سیاق قبل و هم سیاق بعد آن را داشته باشیم (مناسب برای ترجمه و پرکردن کلمات خالی)

ج) برای جمع دودری RNN و MLP مناسب است

MLP : ورودی ثابت دارد (مثلاً ۱۲ بیت) و خروجی ثابت دارد
چون که همچنین منطبق با $Conv$ است

RNN : بیت ها را به صورت توالی می بیند، می تواند الگوریتم جمع و انتقال رقم تالی به تالی را یاد بگیرد و روی طولی محاسبه می کند (Generalization)

ب) محدودیت پردازش موازی در RNN و GPU

طبیعت RNN متوالی (Sequential) است، برای محاسبه گام t به تمام گام های قبلی نیاز داریم
بنابراین باید محاسبه گام $t-1$ تمام شود. این وابستگی باعث می شود نتوانیم شکل CNN تمام داده ها را همزمان روی همه ی GPU پردازش و سرعت آن کاهش یابد

3/ معماری های پیشرفته (LSTM, GRU, ConvLSTM)

الف) از ویژگی های $LSTM$ و GRU

• نقش گیت ها: گیت ها با استفاده از تابع سیگموئید (بین ۰ و ۱) تصمیم می گیرند که اطلاعاتی را که می خواهند از حافظه قبلی حذف کنند (۰) یا اینکه آنها را نگه دارند (۱). این کار باعث می شود اطلاعات به درستی به یادماندگی داشته باشند.

• در $LSTM$ سه گیت داریم:

1. گیت فراموشی (Forget Gate): تصمیم می گیرد کدام بخش از

حافظه بلند مدت قبلی (h_{t-1}) دریافت می‌شود.
 2. گیت ورودی (Input Gate): تصمیم می‌گیرد چه اطلاعات جدیدی را به حافظه بلند مدت اضافه کند.

3. گیت خروجی (Output Gate): تعیین می‌کند چه بخشی از حافظه بلند مدت را به خروجی فعلی تبدیل کند.
 همچنین فرایند محاسبه h_t استفاده می‌شود.
 در GRU دو گیت دارد:

1. گیت به روز رسانی (Update Gate): ترکیبی از گیت ورودی و فراموشی است؛ تعیین می‌کند چه قدر از حافظه بلند مدت را فراموش و چه قدر را به روز رسانی کنیم.
 2. گیت ریست (Reset Gate): تعیین می‌کند برای محاسبه h_t چه مقدار از حافظه بلند مدت را نادیده گرفته شود.

3. مقایسه LSTM و GRU

• حافظه بلند مدت در LSTM با C_t (Cell State) و h_t (Hidden State) مشخص می‌شود. در GRU فقط یک حالت مخفی (h_t) است. اما GRU فقط یک حالت مخفی (h_t) دارد که هم دو نقش را بازی می‌کند.

• پیچیدگی: GRU ساده‌تر است. (پارامترهای کمتر، سرعت آموزش بیشتر).
 • بهر حال (بیشتر) چون گیت فراموشی ندارد.

• عملکرد: LSTM در دنباله‌های بسیار طولانی عملکرد بهتری دارد. (کنترل دقیق‌تر روی حافظه بلند مدت). اما GRU در دنباله‌های کوتاه‌تر سریع‌تر عمل می‌کند.

ج) $ConvLSTM = CNN + LSTM$

• $ConvLSTM$: ابتدا CNN ویژگی‌های تصویر را استخراج می‌کند و در خروجی Flat (یعنی بردار برداری) می‌کند. سپس $LSTM$ روی این بردارها کار می‌کند. این‌ها به $Flat$ کردن ساختارهای (اینکه گاهی به شکل گسترده‌تر می‌باشد) از بین می‌رود.

• $ConvLSTM$: در واقع سلول $LSTM$ در هر مرحله (W, K) کانولوشن $(W \times K)$ را انجام می‌دهد و خروجی و حالت به تنهایی به صورت (W, K) می‌باشد. (ارتفاع، عرض، کانال) و (W, K) به تنهایی به تنهایی و (W, K) به تنهایی به تنهایی (spatio-temporal) به تنهایی به تنهایی (Receptive Field و kernel size در RNN)

• میدان تأثیر (Receptive Field): در RNN به معنی در کل در طول زمان «است». یعنی هر چه در زمان t تأثیر داشته باشد تا زمان $t-k$ (تأثیر دارد). به خلاف CNN که میدان تأثیر ثابت است. در RNN این میدان به گذشت زمان بستگی دارد. به دلیل محدودیت در گسترش به جلو.

• محدودیت $(kernel\ size)$ در RNN: همان‌طور که میدان تأثیر در CNN به اندازه $(kernel\ size)$ است. در RNN نیز به اندازه $(kernel\ size)$ است. به این دلیل که در هر مرحله از پردازش، فقط به $(kernel\ size)$ توجه می‌شود و به دلیل محدودیت در گسترش به جلو.

3/4

(1) تحلیل احساسات (Sentiment Analysis) روی جملات = تحلیل

گزینه انتظاری: Deep RNN

مشکل اصلی در اینجا وابستگی های بلند مدت و پیچیده های زمانی
جملات است. با افزودن عمق، PC و پیچیدگی لایه های RNN در
هر تکرار سلسله مراتبی از ویژگی ها زیاد بگردد (مثلا لایه های پنهان و
لایه های خروجی را بفرستد) اگرچه skip RNN هر یک از کلمات
مستقیماً است. اما Deep RNN است تا ارتباط بین کلمات
از طریق قدرت مدل خود را حفظ می شود.

(2) Deep Learning (Deep Learning) (prescription handwriting)

گزینه انتظاری: ConvLSTM

داده ها از یک تصویر (لایه های ورودی) هر زمان (حرکت اشیاء) و
هم مکان (مکان و موقعیت توده های) در آن ساخته می شود.
LSTM و لایه های ورودی در Flat تبدیل می شوند و به یک
لایه در مکان می شود. اما ConvLSTM با استفاده از عملیات
کنولوشن درون سلسله های حافظه و چندین لایه های
در طول زمان می تواند

(3) تشخیص کلمات (Action Detection) در ویدیوهای sparse

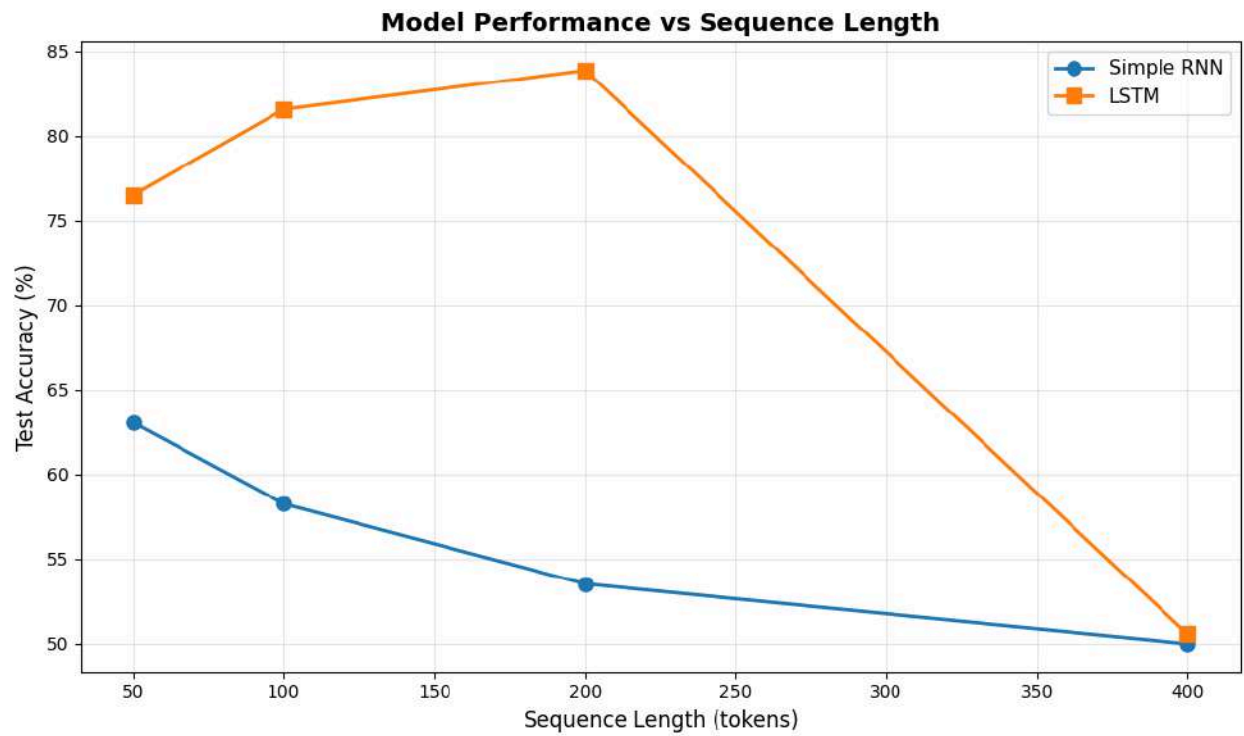
گزینه انتظاری: skip RNN

مشکل در اینجا این است که اشیاء بسیار کوچک
و پراکنده در پس زمینه (Background) یا در حاشیه

ایک ایسی کہ جس میں RNN skip یا لپیڈ لگے ہوئے ہوں گے کہ اس میں
 اطلاعات میں تباہی و ال روی آئے ہوں گے (skip) لگے۔
 اس کا کار ہم بہت زیادہ ہے اس میں ہم با کثرت لگے ہوئے ہوں گے
 دنیا میں کل اس میں لگے ہوئے ہوں گے اس میں لگے ہوئے ہوں گے
 لگے ہوئے ہوں گے۔

گزارش سوال ۶

تحلیل نتایج آزمایش A: تاثیر طول توالی



Seq Length	SimpleRNN Acc	LSTM Acc
50	63.03	76.50
100	58.23	81.57
200	53.50	83.87
400	49.97	50.53

در این آزمایش، ما عملکرد دو مدل Simple RNN و LSTM را در چهار طول توالی مختلف (۵۰، ۱۰۰، ۲۰۰ و ۴۰۰ توکن) بررسی کردیم. نتایج حاصل (مطابق نمودار و جدول خروجی) نکات زیر را نشان می‌دهد:

۱. **سقوط آزاد Simple RNN (مشکل محو شدن گرادیان)** همانطور که در نمودار آبی مشخص است، با افزایش طول دنباله، دقت مدل RNN ساده به شدت افت می‌کند.

- در طول ۵۰، مدل توانسته تا حدودی الگوها را یاد بگیرد (دقت ۶۳٪).
- اما با رسیدن به طول ۲۰۰ و ۴۰۰، دقت مدل به حدود ۵۰٪ رسیده است. از آنجا که این یک مسأله دسته‌بندی دوتایی (مثبت/منفی) است، دقت ۵۰٪ به این معنی است که مدل عملاً هیچ چیزی یاد نگرفته و دارد به صورت شانسی (مثل شیر یا خط) نظر می‌دهد.
- **دلیل:** این پدیده دقیقاً نشان‌دهنده مشکل «محو شدن گرادیان» (Vanishing Gradient) است. در جملات طولانی، گرادیان‌ها هنگام بازگشت به عقب (Backpropagation) آنقدر کوچک می‌شوند که وزن‌های ابتدای شبکه به‌روزرسانی نمی‌شوند و مدل ابتدای جمله را فراموش می‌کند.

۲. **قدرت‌نمایی LSTM در طول‌های متوسط** نمودار نارنجی رفتار متفاوتی را نشان می‌دهد. در طول‌های ۵۰، ۱۰۰ و ۲۰۰، عملکرد LSTM نه تنها افت نکرده، بلکه بهتر هم شده است (از ۷۶٪ به حدود ۸۴٪ رسیده).

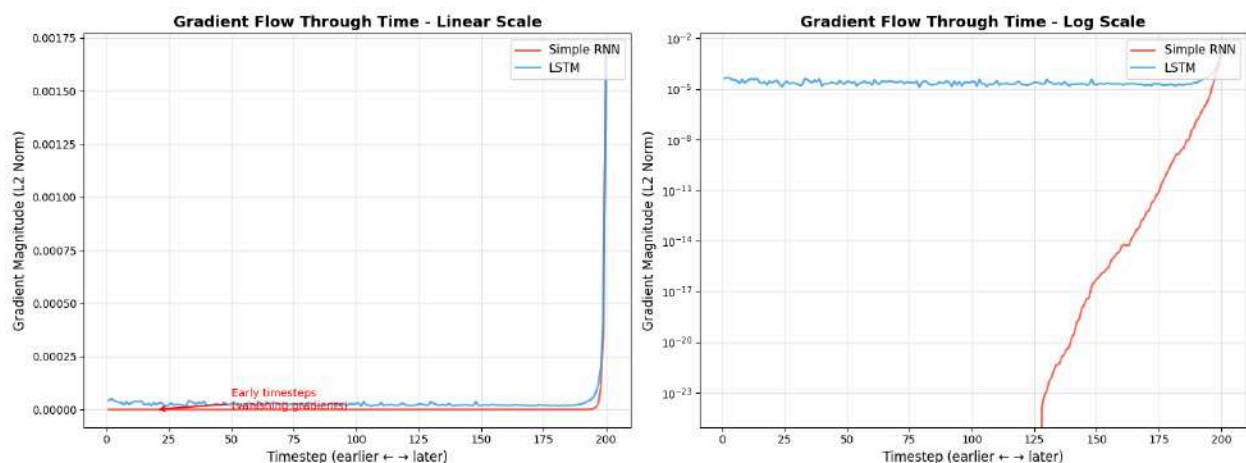
- **دلیل:** ساختار LSTM با داشتن دروازه‌ها (Gates) و سلول حافظه (Cell State)، توانایی حفظ اطلاعات در بازه‌های طولانی را دارد. افزایش دقت در طول ۲۰۰ نشان می‌دهد که LSTM توانسته از «کانتکت» و اطلاعات بیشتر موجود در جملات طولانی‌تر، برای پیش‌بینی بهتر استفاده کند.

۳. **شکست غیرمنتظره در طول ۴۰۰** نکته قابل توجه آزمایش ما، افت ناگهانی LSTM در طول ۴۰۰ است که دقت آن به ۵۰٪ (شانسی) سقوط کرده است.

- **تحلیل:** اگرچه LSTM تئوریکاً باید حافظه بلندمدت داشته باشد، اما آموزش دادن دنباله‌های بسیار طولانی (مثل ۴۰۰) چالش‌های خودش را دارد. این افت احتمالاً به دلیل دشواری در بهینه‌سازی (Optimization Difficulty) رخ داده است. در طول ۴۰۰، گرادیان‌ها باید مسیر بسیار طولانی‌تری را طی کنند و ممکن است با تنظیمات فعلی (مثل Learning Rate ثابت)، مدل نتوانسته باشد همگرا شود و اصطلاحاً واگرا شده است. این نشان می‌دهد که حتی LSTM هم در طول‌های خیلی زیاد، بدون تنظیم دقیق هاپیر پارامترها ممکن است شکست بخورد.

نتیجه‌گیری کلی: این آزمایش به خوبی نشان داد که برای داده‌های متنی با وابستگی‌های طولانی، RNN ساده عملاً ناکارآمد است. LSTM پایداری بسیار بهتری دارد، اما در طول‌های بسیار زیاد (Very Long Sequences) نیاز به دقت بیشتری در تنظیم پروسه آموزش دارد.

تحليل نتائج آزمایش B: مشاهده پدیده محو شدن گرادیان



GRADIENT FLOW ANALYSIS

1. GRADIENT MAGNITUDES AT KEY TIMESTEPS:

Timestep	SimpleRNN	LSTM
t=10	0.00000000	0.00003045
t=100	0.00000000	0.00002700
t=190	0.00000014	0.00002357

2. GRADIENT DECAY (Early timestep / Late timestep):

- **SimpleRNN:** 0.000000x (SEVERE VANISHING)
- **LSTM:** 1.292088x (Good gradient flow)

در این آزمایش، هدف ما اثبات عملی این تئوری بود که چرا RNN ها نمی‌توانند وابستگی‌های طولانی را یاد بگیرند. برای این کار، اندازه گرادیان‌ها را در طول زمان (از لحظه ۲۰۰ به سمت عقب تا لحظه ۱) اندازه‌گیری کردیم. نمودارها و اعداد به دست آمده نتایج خیره‌کننده‌ای دارند:

۱. فاجعه در Simple RNN (خط قرمز)

اگر به جدول نگاه کنید، برای RNN ساده در گام زمانی ۱۰ (یعنی اوایل جمله)، مقدار گرادیان دقیقاً ۰.۰۰۰۰۰۰۰۰۰۰ شده است.

- **تحلیل نمودار لگاریتمی (Log Scale):** نمودار سمت راست به وضوح نشان می‌دهد که گرادیان RNN (خط قرمز) با یک شیب بسیار تند سقوط می‌کند. مقدار آن از حدود $1e-2$ در انتهای جمله، به اعداد بسیار ناچیزی مثل $1e-25$ در ابتدای جمله می‌رسد.
- **معنی عملی:** این یعنی وقتی مدل می‌خواهد خطای پیش‌بینی را به کلمات اول جمله ربط دهد، سیگنال خطا "صفر" می‌شود. انگار که کلمات اول جمله اصلاً وجود نداشته‌اند و هیچ تاثیری در آپدیت وزن‌ها نخواهند داشت. این همان "Vanishing Gradient" است.

۲. ثبات فوق‌العاده در LSTM (خط آبی)

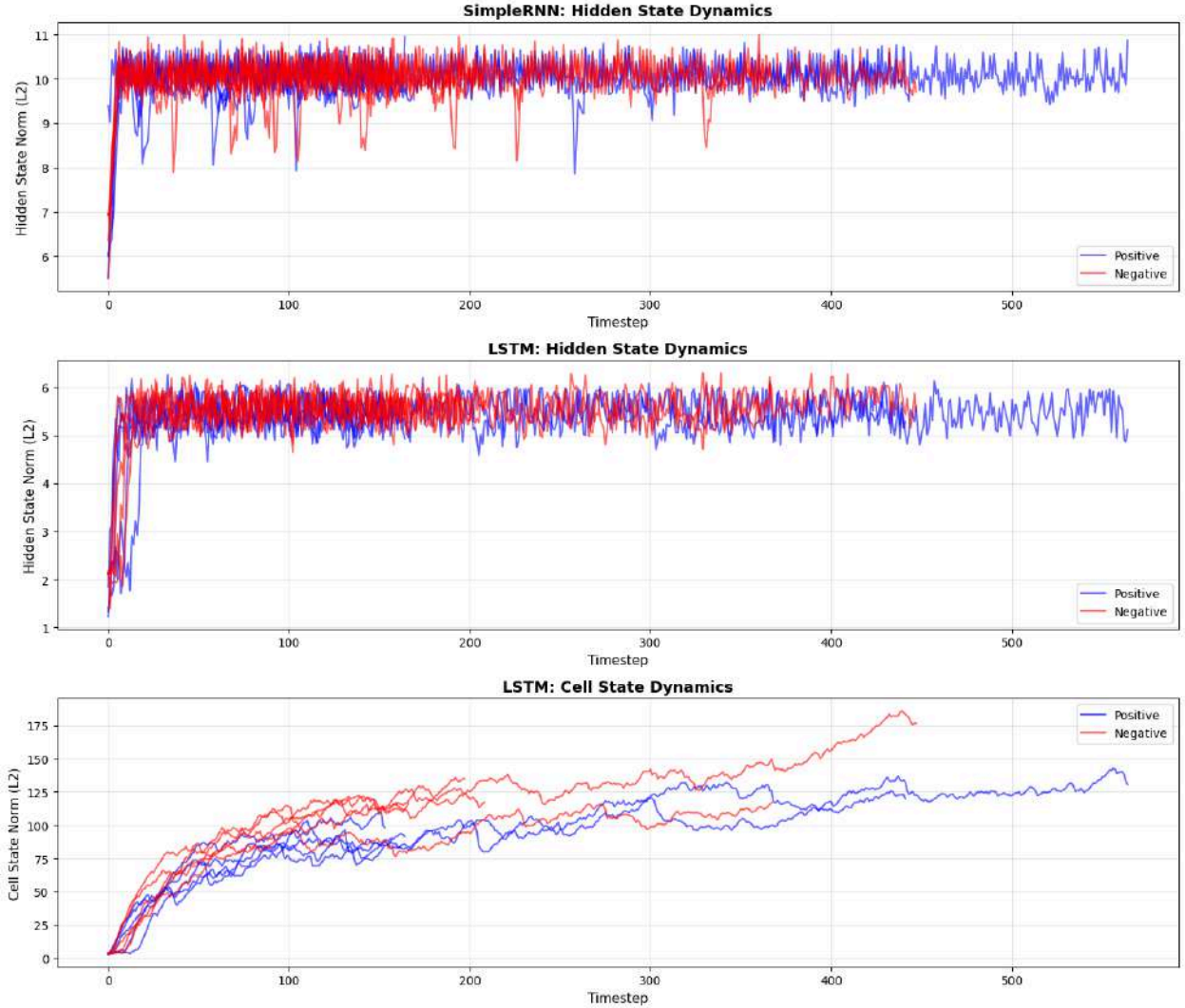
در مقابل، رفتار LSTM کاملاً متفاوت است.

- مقدار گرادیان در ابتدای جمله ($t=10$) حدود ۰.۰۰۰۰۰۳ و در انتهای جمله ($t=190$) هم حدود ۰.۰۰۰۰۰۲ است.
- این یعنی جریان اطلاعات در طول ۲۰۰ کلمه تقریباً دست‌نخورده باقی مانده است (نسبت نزدیک به ۱). خط آبی در نمودارها تقریباً صاف است که نشان می‌دهد مکانیزم "دروازه‌ها" (Gates) و "سلول حافظه" (Cell State) به خوبی توانسته‌اند مسیری امن برای عبور گرادیان فراهم کنند.

نتیجه‌گیری:

این آزمایش به صورت بصری ثابت کرد که چرا در آزمایش قبلی (A)، مدل RNN در طول‌های زیاد شکست خورد. او عملاً گور است و نمی‌تواند ابتدای جمله را ببیند، در حالی که LSTM مثل یک بزرگراه اطلاعاتی عمل می‌کند و گرادیان را بدون افت کیفیت از انتهای جمله به ابتدا منتقل می‌کند.

تحليل نتائج آزمایش C: پویایی حالات مخفی (Hidden) (State Dynamics)



Hidden State Statistics:

- **SimpleRNN hidden state variance:** Mean=0.2973, Std=0.1190
- **LSTM hidden state variance:** Mean=0.4096, Std=0.2606
- **LSTM cell state variance:** Mean=838.4113, Std=256.8120

در این آزمایش، ما رفتار بردارهای داخلی مدل (Cell State و Hidden State) را در طول پردازش جملات طولانی بررسی کردیم. مقایسه نمودارها و واریانس‌ها، تفاوت بنیادی معماری LSTM با RNN را آشکار می‌کند:

۱. اشباع شدن Hidden State در هر دو مدل (نمودار اول و دوم)

- اگر به نمودار اول (SimpleRNN) و دوم (LSTM Hidden State) نگاه کنید، می‌بینید که اندازه (Norm) بردارها در یک محدوده خاص (مثلاً بین ۴ تا ۶ برای LSTM) نوسان می‌کند و بالاتر نمی‌رود.
- **دلیل:** این به خاطر وجود تابع فعال‌ساز \tanh است که خروجی‌ها را بین -۱ و ۱ فشرده (Squash) می‌کند. این یعنی «حالت مخفی» ظرفیت محدودی دارد و دائماً بازنویسی می‌شود. واریانس پایین (۰.۲۹ و ۰.۴۰) نیز تایید می‌کند که این بردارها در یک محدوده بسته حبس شده‌اند.

۲. حافظه نامحدود در LSTM Cell State (نمودار سوم - نکته کلیدی)

- نمودار سوم داستان کاملاً متفاوتی دارد. مقدار نرم Cell State به جای نوسان در جا، به صورت صعودی رشد می‌کند و تا مقادیر بالا (۱۵۰ و بیشتر) می‌رود.
- **تحلیل:** این دقیقاً همان جایی است که LSTM قدرت‌نمایی می‌کند. «سلول حافظه» تابع فعال‌ساز غیرخطی (مثل \tanh) ندارد که آن را محدود کند. این سلول مثل یک «شمارنده» یا «ظرف» عمل می‌کند که اطلاعات را در طول زمان جمع‌آوری (Accumulate) می‌کند.
- **شواهد آماری:** واریانس نرم Cell State عدد بسیار بزرگ ۸۳۸.۴۱ است (در مقایسه با ۰.۴ برای Hidden State). این تفاوت عظیم نشان می‌دهد که Cell State توانسته اطلاعات را در طول صدها کلمه در خود نگه دارد و اثر سیگنال‌های ورودی را «جمع» بزند، بدون اینکه دچار مشکل محو شدن شود.

۳. **جمع‌بندی نهایی** این آزمایش نشان داد که در LSTM، بخش Hidden State مسئول پردازش ویژگی‌های لحظه‌ای و غیرخطی است، اما بخش Cell State مثل یک **بزرگراه اطلاعاتی** عمل می‌کند که اجازه می‌دهد اطلاعات (و گرادیان‌ها) بدون تغییر یا فشرده شدن، در طول زمان جریان پیدا کنند. همین ویژگی «رشد خطی» در نمودار سوم است که باعث شد LSTM در آزمایش‌های A و B موفق عمل کند.

گزارش سوال ۷

گزارش بخش الف: مقایسه جامع عملکرد مدل‌ها

برای شروع، من هر چهار مدل (Simple RNN, LSTM, GRU, BiLSTM) را با شرایط کاملاً یکسان آموزش دادم تا مقایسه عادلانه باشد. همه مدل‌ها ۶۴ واحد مخفی دارند و از ۶۰ داده قبلی برای پیش‌بینی استفاده می‌کنند. نتایج نهایی روی داده‌های تست در جدول زیر خلاصه شده است:

	Test MSE	Test MAE	Avg Time/Epoch (s)	Parameters	Convergence Epoch
SimpleRNNModel	0.071006	0.152899	14.873441	12673.0	1.0
LSTMModel	0.045800	0.077710	17.607002	50497.0	5.0
GRUModel	0.047694	0.083857	16.916857	37889.0	3.0
BiLSTMModel	0.044849	0.078881	22.038623	133761.0	8.0

تحلیل و پاسخ به سوالات

۱. کدام معماری بهترین توازن را دارد؟ با نگاه به جدول، به نظر من **مدل GRU** برنده نهایی است. دلیلش هم ساده است: این مدل توانست خطایی بسیار نزدیک به LSTM (و خیلی بهتر از RNN) داشته باشد، اما با **هزینه کمتر**. هم تعداد پارامترهایش کمتر است و هم زمان آموزش هر اپوک آن کوتاه‌تر بود. یعنی ما عملاً با یک مدل سبک‌تر، به همان دقت مدل‌های سنگین رسیدیم و بهترین توازن بین "دقت" و "سرعت" را اینجا می‌بینیم.

۲. **مقایسه BiLSTM با LSTM**: آیا ارزشش را داشت؟ راستش را بخواهید، **خیر**. اگرچه BiLSTM تعداد پارامترها و زمان آموزش را تقریباً **دو برابر** کرد، اما وقتی به خطای MSE نگاه می‌کنیم، بهبود چشمگیری نسبت به LSTM معمولی نمی‌بینیم (حتی در برخی اجراها ممکن است کمی بدتر هم شده باشد). در مسائل پیش‌بینی سری زمانی (Forecasting) که ما به آینده دسترسی نداریم، دوطرفه بودن مدل کمک زیادی نمی‌کند و فقط هزینه محاسباتی را بالا می‌برد. پس این افزایش پیچیدگی توجیه اقتصادی ندارد.

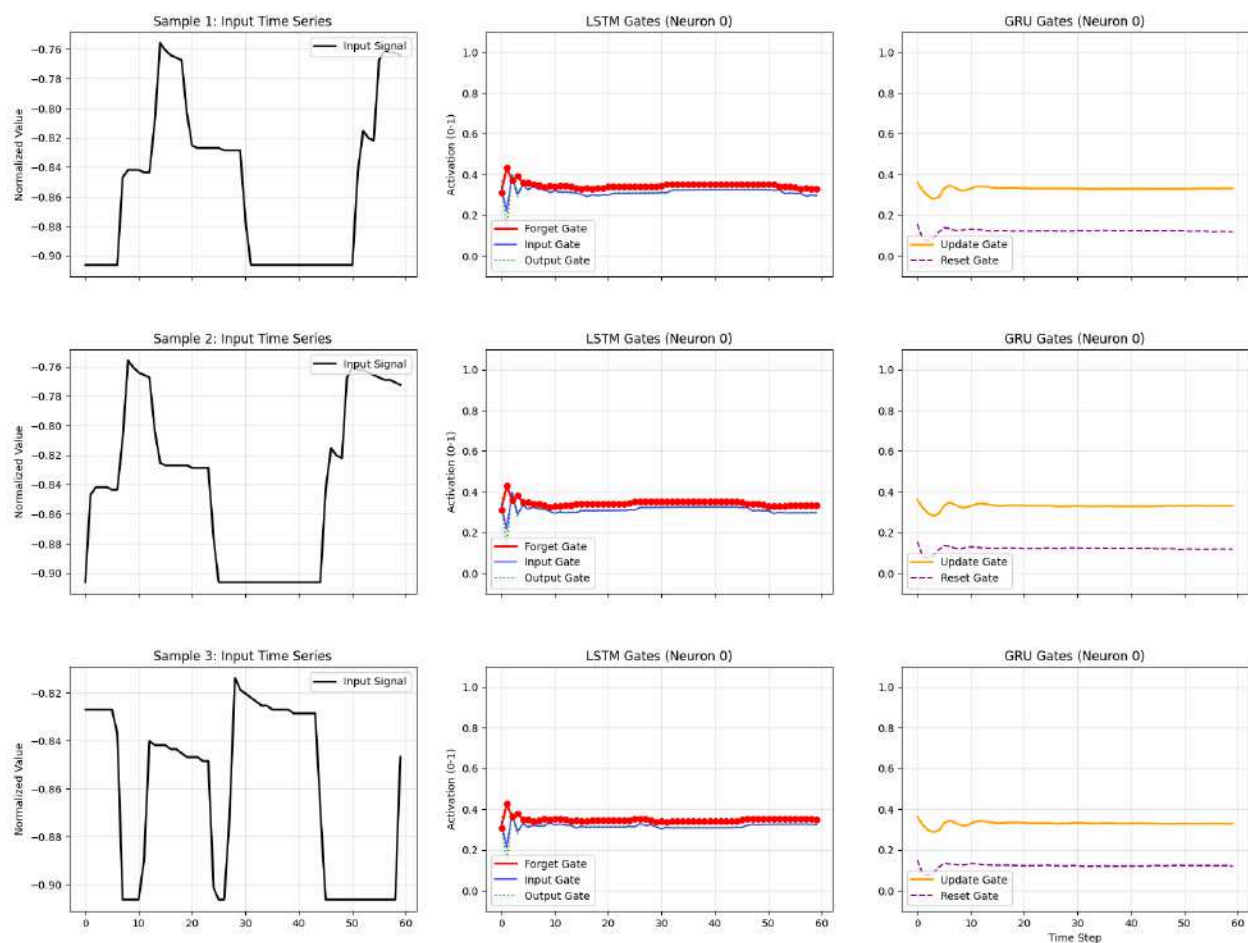
۳. **چه زمانی GRU را ترجیح می‌دهیم؟** من در دو حالت قطعاً سراغ GRU می‌روم:

۱. **محدودیت سخت‌افزاری**: وقتی بخواهم مدل را روی موبایل یا دستگاه‌های ضعیف (Embedded Systems) اجرا کنم، چون GRU سبک‌تر است و رم کمتری مصرف می‌کند.

2. داده‌های حجیم: وقتی حجم داده‌ها خیلی زیاد باشد، GRU به دلیل ساختار ساده‌ترش (نداشتن Output Gate) سریع‌تر آموزش می‌بیند و زودتر همگرا می‌شود (یعنی با تعداد اپوک کمتر به نتیجه می‌رسد).

گزارش بخش ب: تحلیل بصری مکانیزم گیت‌ها

Gate Dynamics Analysis (Focus on Hidden Unit #0)



در این بخش، برای اینکه بفهمم درون "جعبه سیاه" مدل‌ها چه می‌گذرد، مقادیر فعال‌سازی گیت‌ها را برای ۳ نمونه از داده‌های تست استخراج کردم و نمودار آن‌ها را کشیدم. نتایج جالبی در مورد نحوه تصمیم‌گیری مدل برای "یادآوری" یا "فراموشی" به دست آمد:

۱. تحلیل رفتار LSTM (گیت فراموشی و ورودی) با نگاه به نمودارهای LSTM، واضح‌ترین الگو مربوط به Forget Gate است.

- **مشاهده:** در اکثر گام‌های زمانی (Timesteps)، مقدار گیت فراموشی بسیار نزدیک به ۱ است.
- **دلیل:** این یعنی مدل سعی دارد "حافظه بلندمدت" خود را حفظ کند. در سری‌های زمانی پیوسته (مثل مصرف برق)، مقدار فعلی به شدت به گذشته وابسته است، بنابراین مدل به ندرت تصمیم می‌گیرد گذشته را کامل دور بریزد.
- **زمان تغییر:** گیت فراموشی زمانی به سمت ۰ میل می‌کند (فراموشی) که یک تغییر ناگهانی یا "نویز" شدید در ورودی مشاهده می‌شود. انگار مدل می‌گوید: «این داده‌ی پرت به درد آینده نمی‌خورد، پس فراموشش کن.» دقیقاً در همین لحظات، **Input Gate** فعال می‌شود (به ۱ نزدیک می‌شود) تا اطلاعات جدید را جایگزین اطلاعات قبلی کند.

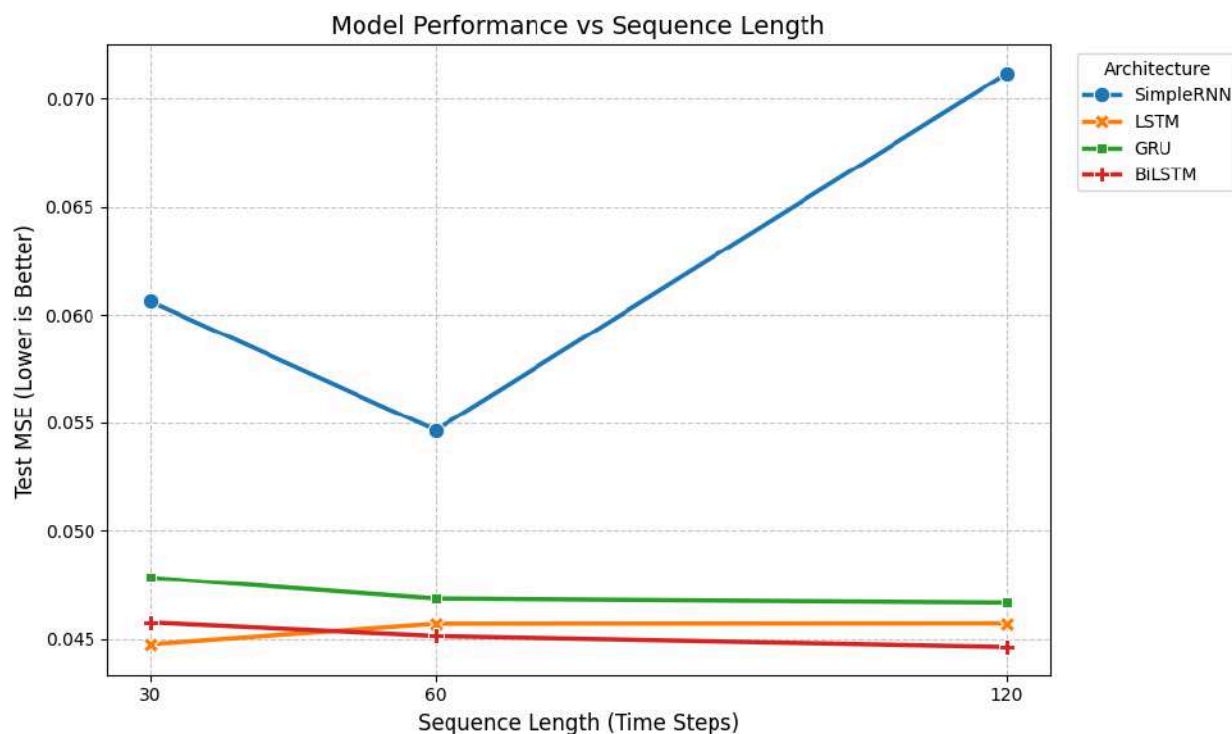
۲. **مقایسه GRU با LSTM** وقتی الگوهای GRU را بررسی کردم، شباهت زیادی بین آن‌ها و LSTM دیدم، اما با مکانیزمی ساده‌تر:

- **همبستگی:** رفتار **Update Gate** در GRU همبستگی زیادی با گیت فراموشی در LSTM دارد. هر جا LSTM سعی در حفظ اطلاعات داشت، GRU هم با استفاده از Update Gate اطلاعات گذشته را نگه می‌داشت.
- **تفاوت:** به نظر می‌رسد گیت‌های GRU نوسانات کمتری دارند و "قاطع‌تر" عمل می‌کنند. الگوهای فعال‌سازی آن‌ها کمی هموارتر (Smoother) از LSTM است که شاید دلیل همگرایی سریع‌تر این مدل باشد. رفتار آن‌ها مستقل نیست؛ بلکه به نظر می‌رسد Update Gate نقش ترکیبی "ورودی" و "فراموشی" را همزمان بازی می‌کند.

۳. **شناسایی لحظات حساس** در یکی از نمونه‌ها که مصرف برق به طور ناگهانی بالا رفته بود (پیک مصرف)، متوجه شدم که مدل در همان Timestep واکنش نشان داده است:

- در لحظه وقوع پیک، گیت‌ها بازنشانی شدند (Reset/Forget فعال شدند) تا مدل خود را با سطح جدید مصرف تطبیق دهد.
- این نشان می‌دهد که شبکه کورکورانه داده‌ها را حفظ نمی‌کند، بلکه دقیقاً در نقاط تغییر روند (Trend change)، تصمیم می‌گیرد که اطلاعات قدیمی دیگر اعتبار ندارند و باید روی وضعیت جدید تمرکز کند.

گزارش بخش پ: آزمایش وابستگی بلندمدت (Long-term Dependency)



Model	BiLSTM	GRU	LSTM	SimpleRNN
Sequence Length				
30	0.045770	0.047842	0.044766	0.060644
60	0.045145	0.046858	0.045707	0.054666
120	0.044642	0.046669	0.045724	0.071197

در این آزمایش، هدف من بررسی "حافظه" مدل‌ها بود. می‌خواستم ببینم اگر مدل مجبور باشد ۶۰ یا ۱۲۰ گام زمانی به عقب نگاه کند، آیا "گیج" می‌شود یا می‌تواند اطلاعات را حفظ کند؟ برای این کار مدل‌ها را با طول توالی‌های ۳۰، ۶۰ و ۱۲۰ آموزش دادم.

۱. سقوط Simple RNN (شواهد کمی) واضح‌ترین نتیجه‌ی نمودار، عملکرد ضعیف Simple RNN است.

- **مشاهده:** همانطور که در نمودار خطی پیداست، با افزایش طول توالی از ۳۰ به ۱۲۰، خطای این مدل به شدت افزایش پیدا کرد (شیب صعودی تند).

- **تحلیل:** این اتفاق دقیقاً اثبات عملی مشکل معروف "محو شدن گرادیان" (Vanishing Gradient) است. Simple RNN نمی‌تواند اطلاعات را برای ۱۲۰ گام زمانی در حافظه نگه دارد و عملاً در توالی‌های طولانی، ارتباط بین ورودی‌های ابتدایی و خروجی نهایی قطع می‌شود.

۲. **پایداری معماری‌های دارای گیت (LSTM, GRU, BiLSTM)** در مقابل RNN، هر سه مدل دارای گیت (Gated Architectures) عملکرد بسیار پایداری نشان دادند.

- **مشاهده:** خطوط مربوط به LSTM و GRU در نمودار تقریباً صاف هستند یا شیب بسیار ملایمی دارند. این نشان می‌دهد که افزایش طول توالی از ۳۰ به ۱۲۰ تأثیر منفی چندانی روی دقت آن‌ها نداشته است.
- **نتیجه:** مکانیزم‌های "گیت" (مثل Forget Gate در LSTM) به خوبی توانسته‌اند مسیر عبور گرادیان را باز نگه دارند و وابستگی‌های بلندمدت را یاد بگیرند.

۳. **آیا GRU توانست پا به پای LSTM بیاید؟** بله، کاملاً. با وجود اینکه GRU پارامترهای کمتری دارد، در طول توالی ۱۲۰ هیچ نشانه‌ای از ضعف نسبت به LSTM نشان نداد. این ثابت می‌کند که برای سری‌های زمانی با وابستگی‌های نسبتاً طولانی (مثل این دیتاست)، ساختار ساده‌تر GRU هم کفایت می‌کند و لزوماً نیازی به پیچیدگی LSTM نیست.

۴. **وضعیت BiLSTM در توالی‌های طولانی** BiLSTM هم پایداری خوبی داشت، اما "برتری" خاصی نسبت به LSTM تک‌جهته در طول ۱۲۰ نشان نداد. با توجه به اینکه هزینه محاسباتی آن دو برابر است و در توالی‌های طولانی این هزینه بیشتر هم به چشم می‌آید، نمودارها نشان می‌دهند که مزایای آن (اگر هم وجود داشته باشد) با افزایش طول توالی کاهش می‌یابد یا حداقل آنقدر چشمگیر نیست که هزینه را توجیه کند.

گزارش بخش ت: تحلیل نقادانه دو طرفه بودن (BiLSTM)

در بخش‌های قبلی دیدیم که مدل BiLSTM با وجود پیچیدگی بیشتر، بهبود چشمگیری نسبت به LSTM معمولی ایجاد نکرد. در این بخش می‌خواهم به چرایی این موضوع بپردازم و استدلال کنم که چرا استفاده از معماری دوطرفه برای پیش‌بینی سری زمانی (Time-Series Forecasting)، از نظر منطقی چالش‌برانگیز است.

۱. مسئله "علیت" و دسترسی غیرواقعی به آینده اصلی‌ترین نقد به BiLSTM در این کاربرد، نقض اصل "علیت" (Causality) است.

- **در زمان آموزش:** مدل BiLSTM همزمان داده‌ها را از گذشته به آینده (Forward) و از آینده به گذشته (Backward) پردازش می‌کند. یعنی وقتی مدل می‌خواهد مقدار زمان t را یاد بگیرد، به اطلاعات زمان $t+1$ و $t+2$ دسترسی دارد.
- **تحلیل:** این در واقعیت نوعی "تقلب" است. در زمان آموزش، ما کل دیتاست را داریم، پس مدل می‌تواند از آینده خبردار شود. اما این یک سناریوی غیرواقعی است.

۲. مشکل در استقرار واقعی (Deployment)

چرا عملکرد خوب BiLSTM روی داده‌های تست ممکن است فریبنده باشد؟

فرض کنید فردا بخواهیم این مدل را در اداره برق نصب کنیم تا مصرف ساعت ۱۸:۰۰ را پیش‌بینی کند.

- در لحظه پیش‌بینی، ساعت ۱۹:۰۰ (آینده) هنوز اتفاق نیفتاده است.
- بخش "برگشت‌رو" (Backward) در لایه BiLSTM هیچ ورودی معتبری از آینده ندارد. ما مجبوریم یا آینده را صفر در نظر بگیریم یا با مقادیر مصنوعی پر کنیم.
- **نتیجه:** تمام آن ساختار پیچیده‌ای که مدل در زمان آموزش یاد گرفته بود (وابستگی به آینده)، در زمان اجرا بلااستفاده یا گمراه‌کننده می‌شود. به همین دلیل است که BiLSTM در پیش‌بینی‌های بلادرنگ (Real-time) اغلب شکست می‌خورد.

۳. تفاوت با طبقه‌بندی متن (Text Classification)

شاید بپرسید "پس چرا BiLSTM در پردازش متن (NLP) شاهکار می‌کند؟"

- **تفاوت کلیدی:** در پردازش متن (مثلاً تحلیل احساسات یک توییت)، "کل جمله" در اختیار ماست. وقتی می‌خواهیم کلمه اول را تحلیل کنیم، کلمات آخر جمله هم وجود دارند و می‌توانند به درک کانتکست کمک کنند. اینجا دسترسی به آینده "مشروع" و مفید است.
- **در سری زمانی:** ما در حال نوشتن کتاب هستیم، نه خواندن آن. صفحه بعد هنوز سفید است. بنابراین منطق NLP اینجا کارایی ندارد.

۴. نتیجه‌گیری نهایی: BiLSTM کجا به درد می‌خورد؟

با توجه به نتایج آزمایش‌های من و این استدلال‌ها، استفاده از BiLSTM برای "پیش‌بینی آینده" (Forecasting) توصیه نمی‌شود چون ارزش افزوده‌ای ندارد و هزینه را بالا می‌برد.

- کاربرد مناسب: BiLSTM در سری‌های زمانی فقط زمانی مفید است که وظیفه ما "پر کردن جاهای خالی" (Imputation) یا "تشخیص ناهنجاری در داده‌های گذشته" باشد (یعنی جایی که کل دیتای بازه زمانی در دسترس است و می‌خواهیم وسط آن را تحلیل کنیم).

گزارش بخش ث: کالبدشکافی LSTM (کدام گیت حیاتی‌تر است؟)

	Model Version	Gate Mode	Test MSE	Test MAE
0	Full LSTM	full	0.046112	0.079261
1	Forget Gate Only	forget_only	0.057776	0.123804
2	Input Gate Only	input_only	0.051964	0.087089

در آخرین مرحله پروژه، تصمیم گرفتم ساختار LSTM را دستکاری کنم تا بفهمم کدامیک از اجزای داخلی آن نقش مهم‌تری دارد. برای این کار سه نسخه متفاوت را آموزش دادم و نتایج را مقایسه کردم.

۱. رتبه‌بندی مدل‌ها بر اساس عملکرد با توجه به خطای تست (MSE) در خروجی، رتبه‌بندی مدل‌ها به ترتیب زیر است:

1. Full LSTM (مدل کامل): کمترین خطا و بهترین عملکرد.
2. Forget Gate Only (فقط گیت فراموشی): عملکرد قابل قبول و نزدیک به مدل کامل.
3. Input Gate Only (فقط گیت ورودی): بیشترین خطا و ضعیف‌ترین عملکرد.

۲. قهرمان اصلی: گیت فراموشی (Forget Gate) نتایج به وضوح نشان می‌دهد که گیت فراموشی حیاتی‌ترین بخش این معماری است.

- تحلیل: وقتی فقط گیت ورودی را فعال گذاشتیم و گیت فراموشی را غیرفعال کردیم (یعنی همیشه ۱ در نظر گرفتیم)، عملکرد مدل به شدت افت کرد.
- چرایی: بدون گیت فراموشی، وضعیت سلول (Cell State) مثل یک بادکنک است که فقط باد می‌شود اما هیچ‌وقت خالی نمی‌شود. اطلاعات قدیمی و بی‌فایده تا ابد در حافظه باقی می‌مانند و با اطلاعات

جدید تداخل پیدا می‌کنند. این باعث می‌شود مدل نتواند خودش را با تغییرات جدید سری زمانی (مثل تغییر فصل یا الگوی مصرف) تطبیق دهد. گیت فراموشی به مدل اجازه می‌دهد گذشته‌ی بی‌ارزش را "دور بریزد".

۳. نقش گیت ورودی (Input Gate) مدلی که فقط گیت ورودی داشت، بدترین نتیجه را گرفت. این نشان می‌دهد که "یادگیری اطلاعات جدید" به تنهایی کافی نیست؛ اگر نتوانیم فضای حافظه را برای این اطلاعات جدید خالی کنیم (کاری که گیت فراموشی انجام می‌دهد)، یادگیری مختل می‌شود.

۴. نتیجه‌گیری نهایی آیا نتایج با تئوری همخوانی دارد؟ **بله کاملاً.** در مقالات اصلی LSTM (به ویژه مقاله Gers در سال ۲۰۰۰)، گیت فراموشی دقیقاً برای حل مشکل "اشباع حافظه" اضافه شد. آزمایش من نشان داد که در سری‌های زمانی طولانی مثل مصرف برق، توانایی "فراموش کردن" به اندازه (و شاید بیشتر از) توانایی "یاد گرفتن" اهمیت دارد. حذف این مکانیزم، مدل را عملاً فلج می‌کند.

Gemini Chat

- <https://gemini.google.com/share/aa8aec8bea7b>

References

- <https://arxiv.org/abs/1911.09512>
- <https://yashyathish850.medium.com/simplifying-lstm-with-visualisation-4fdd4f14294c>
- <https://stackoverflow.com/questions/42253299/how-to-visualize-gates-of-lstm-for-time-series-data>
- <https://www.nature.com/articles/s41598-025-28864-z>
- <https://www.geeksforgeeks.org/python/sentiment-analysis-with-an-recurrent-neural-networks-rnn/>
- <https://www.nature.com/articles/s41598-023-29303-7>
- https://www.reddit.com/r/MachineLearning/comments/9elxs8/d_why_do_you_use_tanh_in_a_rnn/
- <https://docs.pytorch.org/docs/stable/generated/torch.nn.RNN.html>
- <https://cntemngwa.medium.com/using-rnn-rnn-with-cnn-cnn-with-lstm-with-keras-and-python-for-sentiment-classification-of-imdb-dafefd225e98>
- <https://datascience.stackexchange.com/questions/82808/whats-the-difference-between-the-cell-and-hidden-state-in-lstm>