

(1)

جمله است ۱: $P(\text{ما امروز}) = \frac{452}{1872} = \frac{4(w_2 w_1)}{c(w_1)} \rightarrow \text{استفاده از جدول bigram}$
 $= \frac{452}{1872} = \frac{c(w_1)}{c(w_1)} \rightarrow \text{استفاده از جدول unigram}$

$P(\text{امروز کتاب}) = \frac{231}{1943} \approx 0,118 \approx 12\%$ $0,20 \approx \frac{1872}{9521} = P(\text{ما})$

$P(\text{کتاب خوانیم}) = \frac{320}{1245} \approx 0,257$

$0,24 \times 0,12 \times 0,257 = 0,007442$ = احتمال جملی (ملازمه کتاب خوانیم):

$0,20 \times 0,24 \times 0,257 =$ Unigram: احتمالی اول را نیز در نظر بگیریم:

جمله است ۲: ما دیروز داستان خوانیم.

$P(\text{ما}) = \text{احتمال unigram} = \frac{1872}{9521} = 0,196$

$P(\text{ما امروز}) = \frac{11}{1872} = 0,005875 \approx 0,22$

$P(\text{دیروز داستان}) = \frac{68}{2021} = 0,0336 \approx 0,034$

$P(\text{داستان خوانیم}) = \frac{345}{945} = 0,365$

$0,22 \times 0,034 \times 0,365 = 0,002733$ = احتمال جمله - حاصل ضرب احتمالات: بعد از در نظر گرفتن احتمال "ما"

$0,196 \times 0,22 \times 0,034 \times 0,365 = 0,000535$ (احتمال با در نظر گرفتن "ما"):

(۲) می دانیم: $P(B|A) = \frac{P(A, B)}{P(A)}$ و $P(A, B) = P(A)P(B|A)$ (chain rule)

variable مارا زیاد کنیم به طور کلی داریم:

$$P(x_1 \rightarrow x_n) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2) \dots P(x_n|x_1, \dots, x_{n-1})$$

= پس احتمال رخ دادن یک جمله به صورت دنباله ای از کلمات w_1 تا w_n را می توان بر (دستی) مشابه نوشتن داد:

$$P(w_1 \rightarrow w_n) = P(w_1)P(w_2|w_1) \dots P(w_n|w_1, \dots, w_{n-1})$$

(۳) برای به دست آوردن احتمال یک جمله از w_1 تا w_n طبق قانون زنجیره ای احتمال می داریم:

$$P(w_{1:n}) = P(w_1)P(w_2|w_1)P(w_3|w_{1:2}) \dots P(w_n|w_{1:n-1}) = \prod_{k=1}^n P(w_k|w_{1:k-1})$$

منظور از $w_{1:k}$: w_1, w_2, \dots, w_k است که w_k بیت سرشم است.
طبق این $w_{1:k}$ ما می توان گفت که کلمه هر کلمه بجای آنکه به کل کلمات جمله قبل از آن وابسته باشد می توان به صورت تدریجی وابسته به اصول محلی جملش در نظر گرفت. (برای bigram و)

بین ترتیب می توان گفت $P(w_n|w_{1:n-1}) \approx P(w_n|w_{1:k-1})$ Markov assumption
بین ترتیب با یکدیگر این احتمالات در معادله اول داریم:

$$P(w_{1:n}) \approx \prod_{k=1}^n P(w_k|w_{k-1})$$

(۳) در مرحله اول فقط ۲ گزینه داریم و beam search کردن آنها را می داریم:

log probability $(-1, w_2) \rightarrow$ "network" و $(-0.25, w_2) \rightarrow$ "neutral"

(۳) پس حالت نامحتمل را برای کرده و در نهایت ۲ تا می برتر نگه می داریم.

"neutral network" $(-1, w_2)$

log probability $(-1.32, w_2 - 0.25 = -0.75) \rightarrow$ "network neutral"

(۳) در مرحله دوم دنباله ای که با استفاده از دو عبارت برتر از مرحله اول در مرحله دوم می سازیم پس می باشد مراحل قبل عبارت برتر را انتخاب می کنیم که به صورت زیر است:

"Neutral network network" $(-1.65, w_2 - 0.25 + -0.25 = -1.75)$

"Network neutral network" $(-1.93, w_2 - 0.25 - 0.25 = -1.75)$

(۴) خرابی انتخاب نمی اندازیم زیرا برای که بیشترین احتمال را دارد "neutral neutral network" است با احتمال $-1.93 = -0.98 + -0.98 + -0.98 = -2.94$ در مرحله دوم از ترس می مولا توجه خارج می شود زیرا احتمال "neutral neutral" جزو عبارت اول به طرز آسانی.

(۳)

(۱) برای هر مرحله با توجه به اینکه تعداد دارگان m است، پیچیدگی از $O(m)$ می باشد.

(۲) برای انتخاب کردن k تا از بهترین ها ابتدا باید $\log m$ بار sort کنیم که از اردر $m \log m$ زمان می برد. تا اینجا می دانیم $O(m \log m) = O(m \log m + m)$

(۳) کار مرحله قبل را برای k تا از رتبه ها به انجام می دهیم $O(k m \log m)$

(۴) پس باید تعداد k رتبه ی تولید شده را مرتب کرد:

$$O(k^2 \log(k)) = O(k^2 \times \log(k))$$

یعنی تا اینجا داریم:

$$O(k m \log(m) + k^2 \log(k)) = O(k m \log m)$$

(۵) این کار را باید برای T مرحله انجام دهیم \Rightarrow داریم:

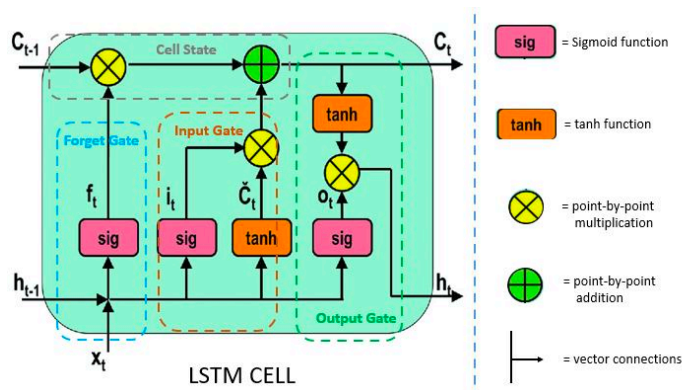
$$T \times O(k m \log m) = O(T k m \log m)$$

(A.4)

وقتی که فقط گیت Forget یک LSTM را حفظ کنیم و گیت های Input و Output را حذف کنیم، تغییرات زیادی در عملکرد شبکه اتفاق می افتد. برای درک بهتر این تغییرات، ابتدا به ساختار LSTM و نقش هر گیت در آن نگاهی می اندازیم.

LSTM یک نوع خاص از شبکه های عصبی بازگشتی (RNN) است که برای حفظ اطلاعات بلندمدت و جلوگیری از مشکل محو شدن گرادیان (Vanishing Gradient) در RNN طراحی شده است. ساختار LSTM شامل چهار بخش اصلی است:

1. گیت فراموشی (Forget Gate): این گیت تصمیم می گیرد که کدام اطلاعات از حافظه ی سلول (Cell Memory) قبلی باید فراموش شود و کدام اطلاعات باقی بمانند.
 2. گیت ورودی (Input Gate): این گیت تصمیم می گیرد که چه اطلاعات جدیدی باید به حافظه ی سلول اضافه شود.
 3. حافظه ی سلول (Cell Memory): این بخش حافظه ی میان مدت است که اطلاعات در طول زمان در آن ذخیره می شود.
 4. گیت خروجی (Output Gate): این گیت تصمیم می گیرد که اطلاعات چه قدر از حافظه ی سلول باید به عنوان خروجی شبکه استفاده شود.
- در تصویر زیر ساختار توابع استفاده شده در این گیت ها را مشاهده می کنیم:



اگر فقط گیت Forget را حفظ کنیم و گیت‌های Input و Output را حذف کنیم، این به این معنی است که شبکه دیگر نمی‌تواند تصمیم بگیرد که کدام اطلاعات باید به حافظه‌ی سلول اضافه شود و کدام اطلاعات باید برای تولید خروجی استفاده شوند.

بدین ترتیب حافظه‌ی سلول به طور مداوم تغییر نخواهد کرد، زیرا اطلاعات جدید به آن اضافه نمی‌شود و همچنین بدون گیت خروجی، شبکه قادر به انتقال اطلاعات از حافظه‌ی سلول به خروجی نخواهد بود و خروجی ثابت و بی‌تغییر می‌شود. پس ساختار شبکه بسیار ساده‌تر شده و نمی‌تواند به خوبی اطلاعات را جذب کند و خروجی‌های مختلفی تولید کند. این تغییرات باعث افت شدید عملکرد شبکه و توانایی آن در یادگیری و پردازش داده‌ها خواهد شد.

(B.4)

وقتی که مقدار گیت Forget را به صفر تنظیم می‌کنیم، شبکه LSTM دیگر هیچ اطلاعاتی را از حافظه‌ی سلول (Cell Memory) خود ذخیره نمی‌کند و شبیه RNN می‌شود. این تغییرات می‌تواند تأثیرات جالبی بر روی عملکرد شبکه داشته باشد:

- از دست دادن سازگاری: شبکه توانایی تنظیم پویا مقدار اطلاعات را بر اساس ورودی و زمینه فعلی از دست می‌دهد. این باعث کاهش انعطاف‌پذیری آن در رسیدگی به انواع مختلف توالی می‌شود.
- ناپدید شدن گرادینان: فراموش کردن مداوم تمام اطلاعات می‌تواند منجر به ناپدید شدن گرادینان در طول آموزش شود. این امر یادگیری وابستگی‌های طولانی مدت در داده‌ها را برای شبکه دشوار می‌کند.
- کاهش توانایی یادگیری: با فراموشی محدود، شبکه تلاش می‌کند تا بر مرتبط‌ترین اطلاعات تمرکز کند و وضعیت داخلی خود را به طور موثر به روز کند. این مانع از ظرفیت کلی یادگیری آن می‌شود.

(C.4)

مزایا:

- بهبود نمایش ویژگی‌های پیچیده:
- با افزایش تعداد لایه‌ها، شبکه قادر به یادگیری ویژگی‌های پیچیده‌تر و سطوح بالاتری از داده می‌شود. هر لایه LSTM به شبکه اجازه می‌دهد که ویژگی‌های جدیدی از داده را استخراج کرده و به لایه‌های بالاتر منتقل کند.
- افزایش توانایی یادگیری داده‌ها:
- شبکه LSTM با لایه‌های بیشتر می‌تواند الگوهای زمانی پیچیده‌تری را یاد بگیرد. این امر می‌تواند بهترین نتایج را در وظایفی مانند ترجمه ماشینی یا پیش‌بینی دنباله‌ای (sequence prediction) فراهم کند.
- کاهش مشکل محو شدن گرادینان:
- با افزایش عمق شبکه (تعداد لایه‌ها)، ممکن است مشکل محو شدن گرادینان که معمولاً در شبکه‌های عمیق به وجود می‌آید، کاهش یابد. این به این معنی است که شبکه قادر به بهترین استفاده از اطلاعات تاریخی در داده‌های توالی خواهد بود.

معایب:

- افزایش پیچیدگی و زمان آموزش:
- با افزایش تعداد لایه‌ها، پیچیدگی شبکه افزایش می‌یابد و زمان آموزش نیز افزایش می‌یابد. شبکه‌های عمیق‌تر نیاز به تعداد بیشتری پارامتر (وزن‌ها) دارند و به دنباله‌های طولانی‌تری نیاز دارند تا آموزش ببینند.
- مشکل Overfitting:
- شبکه‌های عمیق‌تر ممکن است به سرعت به داده‌های آموزشی بیش‌برازش کنند یعنی ویژگی‌های غیرضروری یا نویز در داده‌های آموزش را یاد بگیرند و عملکرد خوبی بر روی داده‌های تازه‌ای (که در فاز آموزش دیده نشده‌اند) ارائه ندهند.
- پیچیدگی بالا برای تنظیم پارامترها:
- با افزایش تعداد لایه‌ها، تنظیم پارامترهای شبکه (مانند نرخ یادگیری و نرخ Dropout) نیاز به دقت و دانش بیشتری دارد.

(I.5)

افزایش متغیر (استفاده از نقاط داده گذشته بیشتر):

مزایا:

- ثبت روندهای بلندمدت: ترکیب داده های بیشتر به LSTM اجازه می دهد تا روندهای بازار گسترده تری را ثبت کند و به طور بالقوه چرخه ها یا الگوهای بلندمدتی را شناسایی کند که ممکن است در بازه زمانی کوتاه تر قابل مشاهده نباشند.
- بهبود استحکام: یک مجموعه داده بزرگتر می تواند مدل را در برابر نویز و نوسانات در داده ها قوی تر کند و منجر به پیش بینی های پایدارتر و قابل اعتمادتر شود.

معایب:

- Overfitting بیش از حد: گنجاندن داده های زیاد می تواند منجر به overfitting شود، مدل الگوهای خاصی را در داده های آموزشی و نویز ها به خاطر می سپارد و در تعمیم به داده های دیده نشده عملکرد ضعیف خواهد داشت.
- افزایش زمان آموزش: آموزش LSTM با مجموعه داده بزرگتر به منابع محاسباتی بیشتری نیاز دارد و زمان بیشتری می برد.

کاهش متغیر (با استفاده از نقاط داده گذشته کمتر):

مزایا:

برعکس قسمت قبل کاهش متغیر می تواند باعث کاهش زمان آموزش - کاهش Overfitting شود.

معایب:

برعکس قسمت قبل می توان به از دست دادن روندهای بلندمدت و امکان پیش بینی های نادرست و افزایش حساسیت به نویز اشاره کرد.

Used links:

<https://aroussi.com/post/python-yahoo-finance>

https://matplotlib.org/stable/gallery/color/named_colors.html

<https://www.youtube.com/watch?v=b61DPVFX03I>

<https://medium.com/@CallMeTwitch/building-a-neural-network-zoo-from-scratch-the-long-short-term-memory-network-1cec5cf31b7>