

به نام خالق رنگین کمان

ستاره باباجانی - گزارش تمرین سری 4

سوال 1: به محاسبه احتمال جملات خواسته شده میپردازیم: (مجموع کلمات پیکره 9521 میباشد).

محاسبه احتمال یونیگرم ها به شرح زیر است:

$$p(\text{ما}) = \frac{1872}{9521}, p(\text{خواندیم}) = \frac{1495}{9521}, p(\text{دیروز}) = \frac{2021}{9521},$$
$$p(\text{امروز}) = \frac{1943}{9521}, p(\text{داستان}) = \frac{945}{9521}, p(\text{کتاب}) = \frac{1245}{9521}$$

حال اگر در جدول بایگرم ها فرض کنیم کلمه اول در ستون سمت راست و کلمه بعدی آن در ردیف بالای آن قرار داشته باشد، خواهیم داشت:

• احتمال رخداد جمله تست 1:

$$p(\text{کتاب} | \text{خواندیم}) \cdot p(\text{امروز} | \text{کتاب}) \cdot p(\text{ما} | \text{امروز}) = p(\text{ما امروز کتاب خواندیم})$$
$$= \frac{452}{1872} * \frac{231}{1943} * \frac{320}{1245} = 0.0074$$

• احتمال رخداد جمله تست 2:

$$p(\text{ما دیروز داستان خواندیم}) =$$

$$p(\text{داستان} | \text{خواندیم}) \cdot p(\text{دیروز} | \text{داستان}) \cdot p(\text{ما} | \text{دیروز}) = \frac{411}{1872} * \frac{68}{2021} * \frac{345}{945} = 0.0027$$

سوال 2:

سوال ۷: برای اثبات ساده:

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

داریم:

$$P(B|A) = \frac{P(A,B)}{P(A)} \rightarrow P(A,B) = P(A) \cdot P(B|A)$$

حالت عمومی:

$$P(A,B,C) = P(A) \cdot P(B|A) \cdot P(C|A,B)$$

برای هر کلمه w_1 احتمال آن $P(w_1)$ است.

برای هر کلمه w_2 احتمال آن $P(w_2 | w_1)$ است.

برای هر کلمه w_3 احتمال آن $P(w_3 | w_1^2)$ است.

برای هر کلمه w_n احتمال آن $P(w_n | w_1^{n-1})$ است.

حالت عمومی را می‌توانیم به صورت زیر بنویسیم:

$$P(w_1^n) = \prod_{k=1}^n P(w_k | w_1^{k-1})$$

سوال 3:

(a) در مرحله اول با توجه به $k = 2$ ، هر دو کلمه را نگه می‌داریم:

- "neural" with score = -0.65
- "network" with score = -0.73

(b) در این مرحله خواهیم داشت:

- For "neural" with the score of -0.65:
 - Expand to "neural neural": total score will be -1.45
 - Expand to "neural network": total score will be -1.25
- For "network" with the score of -0.73:
 - Expand to "network neural": total score will be -1.33
 - Expand to "network network": total score will be -1.53

حال با توجه به $k = 2$ دوتا از بهترین نتایج مرحله قبل را حفظ میکنیم:

- “neural network” with score = -1.25
- “network neural” with score = -1.33

(c) حال دوباره هر یک از نود های مرحله قبل را گسترش میدهیم:

- For “neural network” with the score of -1.25:
 - Expand to “neural network neural”: total score will be -2.05
 - Expand to “neural network network”: total score will be -1.85
- For “network neural” with the score of -1.33:
 - Expand to “network neural neural”: total score will be -2.13
 - Expand to “network neural network”: total score will be -1.93

حال با توجه به $k = 2$ دوتا از بهترین نتایج مرحله قبل را حفظ میکنیم:

- “neural network network” with score = -1.85
- “network neural network” with score = -1.93

(d) خیر. همان طور که میدانیم در این مثال بهترین دنباله ها “neural neural” و “network network network” به ترتیب با امتیاز -1.46 , -1.54 هستند که beam search نتوانسته آنها را پیدا کند.

این الگوریتم لزوما دنباله کلی با بیشترین احتمال را برنمیگرداند بلکه بهترین را در میان تعداد محدودی از دنباله هایی که ارزیابی میکند، پیدا میکند. همچنین این الگوریتم با استفاده از پارامتر k تعادلی بین کارایی محاسباتی و کیفیت جواب نهایی برقرار میسازد. (هرچه k بیشتر، راه حل دقیق تر ولی کندتر و بالعکس.)

(e) همان طور که میدانیم در هر مرحله زمانی t برای یک دنباله k دنباله هایی که در حال حاضر ذخیره شده اند، RNN امتیازات را برای همه m کلمه بعدی محاسبه کرده ($k*m$ فراخوانی) و سپس بهترین k دنباله از این $k*m$ تا برای مرحله محل انتخاب میشوند. پس از آنجایی که این فرآیند برای هر یک از t مرحله دنباله تکرار میشود، پیچیدگی زمانی اجرای کلی آن بصورت $O(t*k*m)$ خواهد بود.

سوال 4:

(a) اگر فقط گیت forget داشته باشیم، LSTM فقط می تواند تصمیم بگیرد که کدام بخش از وضعیت سلول قبلی را حفظ کند، اما نمی تواند کنترل کند که چه اطلاعات جدیدی به وضعیت سلول اضافه می شود.

- گیت input: گیت ورودی کنترل می کند که چه مقدار از اطلاعات جدید (از ورودی فعلی و حالت پنهان قبلی) به حالت سلول اضافه می شود. بدون آن، LSTM توانایی خود را برای افزودن اطلاعات جدید به حالت سلول از دست می دهد. این بدان معنی است که وضعیت سلول فقط می تواند در طول زمان کاهش یابد یا ثابت بماند، زیرا اطلاعات فقط فراموش می شوند، هرگز به روز نمی شوند یا با ورودی های جدید تکمیل نمی شوند.

- گیت output: گیت خروجی کنترل می کند که چه مقدار از حالت سلول برای ایجاد حالت پنهان فعلی (خروجی) استفاده می شود. بدون آن، کل حالت سلول همیشه به عنوان خروجی در معرض دید قرار می گیرد، که

می‌تواند منجر به مسائل نفوذ بیش از حد از حالت سلول و عدم کنترل بر روی اینکه چه بخشی از حالت داخلی در معرض لایه‌ها یا خروجی‌های بعدی قرار می‌گیرد، شود.

- تغییرات خروجی: LSTM احتمالاً در طول زمان توانایی کمتری در یادگیری الگوهای پیچیده پیدا می‌کند زیرا نمی‌تواند اطلاعات جدید را در ورودی‌ها جمع‌آوری کند.

(b) گیت forget در یک LSTM میزان حفظ حالت سلول قبلی (از آخرین مرحله زمانی) را در وضعیت سلول فعلی کنترل می‌کند. این کار را با اعمال یک ضرب گیت (بین 0 و 1) در حالت سلول قبلی انجام می‌دهد.

- بدون گیت forget: وقتی گیت فراموشی روی صفر تنظیم می‌شود، به طور موثری از حفظ هرگونه اطلاعات از وضعیت سلول قبلی در وضعیت سلول فعلی جلوگیری می‌کند. این بدان معناست که LSTM هر آنچه را که قبلاً می‌دانست در هر مرحله زمانی فراموش می‌کند و از انتقال هر گونه اطلاعات به گام‌های زمانی جلوگیری می‌کند.

- تاثیر بر یادگیری و پیش‌بینی: باعث از دست دادن وابستگی‌های طولانی مدت می‌شود. یکی از نقاط قوت اولیه LSTM ها توانایی آنها در به خاطر سپردن اطلاعات در دوره‌های طولانی است. صفر کردن گیت فراموشی این قابلیت را حذف می‌کند و در نتیجه توانایی LSTM را برای استفاده از بافت تاریخی یا وابستگی‌هایی که فراتر از یک مرحله زمانی هستند فلج می‌کند.

(C)

- قابلیت های یادگیری پیشرفته: افزودن لایه های LSTM بیشتر به شبکه اجازه می دهد تا بازنمایی های پیچیده تری را بیاموزد. هر لایه می تواند سطح متفاوتی از انتزاع را بیاموزد. لایه های بیشتر به طور کلی ظرفیت شبکه را افزایش می دهد و به آن اجازه می دهد تا روابط و وابستگی های پیچیده تری را در داده ها مدل کند.
- بهبود درک متنی: شبکه های عمیق تر LSTM به طور بالقوه می توانند وابستگی های طولانی مدت را بهتر از شبکه های کم عمق تر جذب کنند، زیرا اطلاعات دارای مراحل بیشتری برای پردازش و پالایش است که هر یک سطحی از درک زمینه را اضافه می کند. چندین لایه LSTM می توانند اطلاعات را به صورت سلسله مراتبی پردازش کنند. لایه های پایین تر می توانند وابستگی های کوتاه مدت را مدیریت کنند، در حالی که لایه های بالاتر می توانند این وابستگی ها را در بازه های طولانی تر ادغام کنند، بنابراین درک کلی توالی را بهبود می بخشند.
- پس بطور خلاصه، افزایش تعداد لایه های LSTM می تواند توانایی مدل را برای یادگیری الگوهای پیچیده و بهبود عملکرد در وظایفی که نیاز به درک وابستگی های بلندمدت دارند، افزایش دهد. با این حال، این امر با افزایش هزینه محاسباتی، پیچیدگی آموزش، و پتانسیل بیش از حد برازش همراه است.

سوال 5: قسمت های خواسته شده برای کد، در فایل ضمیمه شده، قرار دارد. حال برای قسمت i خواهیم داشت:

همانطور که میدانیم متغیر timestep تعریف شده، میزان توجه به گذشته را شرح میدهد. حال در اینجا به بررسی اثر افزایش و یا کاهش آن میپردازیم:

• افزایش timestep:

- اطلاعات تاریخی (گذشته) بیشتر: با افزایش گام زمانی، مدل پنجره بزرگتری از قیمت های تاریخی را برای هر پیش بینی در نظر می گیرد. این بدان معناست که هر دنباله ورودی مدل حاوی نقاط داده گذشته بیشتری است.
- پتانسیل برای درک بهتر روندها: با زمینه تاریخی بیشتر، مدل ممکن است روندهای بلندمدت و الگوهای فصلی قیمت بیت کوین را بهتر به تصویر بکشد، که می تواند برای پیش بینی قیمت های آینده بسیار مهم باشد.
- افزایش پیچیدگی مدل و هزینه محاسباتی آن: مراحل زمانی بزرگتر می تواند به مدل های پیچیده تری منجر شود، زیرا مدل نیاز به یادگیری از زمینه وسیع تری دارد. این می تواند زمان آموزش و منابع محاسباتی مورد نیاز را افزایش دهد.
- خطر تطبیق بیش از حد: اگر گام زمانی خیلی بزرگ باشد، این خطر وجود دارد که مدل با داده های تاریخی بیش از حد مطابقت داشته باشد، به خصوص اگر داده های آموزشی متنوع کافی در دسترس نباشد. برازش بیش از حد می تواند باعث شود که مدل در داده های آموزشی عملکرد خوبی داشته باشد اما در داده های دیده نشده ضعیف باشد.

• کاهش timestep:

- داده های تاریخی (گذشته) کمتر: کاهش مرحله زمانی به این معنی است که مدل از نقاط داده تاریخی کمتری برای هر پیش بینی استفاده می کند. این مدل را بر روندهای اخیر متمرکز می کند.

▪ آموزش سریعتر مدل: با داده های کمتری برای پردازش برای هر پیش بینی، ممکن است مدل سریعتر آموزش ببیند و به منابع محاسباتی کمتری نیاز داشته باشد.

▪ پتانسیل برای عدم تناسب: اگر گام زمانی خیلی کوچک باشد، ممکن است مدل به اندازه کافی زمینه تاریخی را برای پیش بینی های دقیق ثبت نکند و به طور بالقوه منجر به عدم تناسب شود. این بدان معناست که مدل ممکن است برای ثبت الگوهای اساسی در داده ها بسیار ساده باشد.

انتخاب گام زمانی یک اقدام متعادل کننده است. ما باید مقداری را انتخاب کنیم که زمینه تاریخی کافی برای ثبت الگوهای مرتبط را بدون تحت تأثیر قرار دادن مدل یا تمرکز آن بر روی داده های نامربوط فراهم کند.