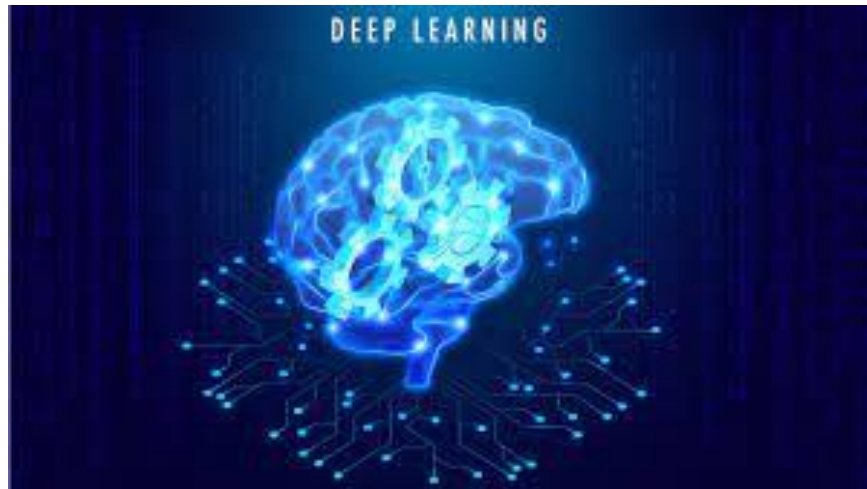


بسم الله الرحمن الرحيم



محمد عرفان زارع زردینی

تمرین سری پنجم درس یادگیری عمیق

۹۸۱۴۱۱۴۳۲

مدرس درس: استاد داوود آبادی

پاییز 1402

سوال 1)

(a)

معماری چند به یک RNN برای کارهایی که ورودی های متوالی (گام های زمانی متعدد) برای تولید یک خروجی واحد باید پردازش شود، مناسب است. حال به بررسی گزینه ها می پردازم:

الف) تشخیص گفتار، تبدیل یک کلیپ صوتی (صوت) به متن می باشد. با توجه به موقعیت هایی که ماهیت متوالی داریم در داده های صوتی، در طول زمان با معماری چند به یک RNN مطابقت دارد. پس برای این تسک مناسب است.

ب) این تسک با قطعه متن سروکار داریم. پس لزوماً به تحلیل متوالی متن در طول زمان برای تعیین احساسات نیست. برای این تسک بهتر است از CNN استفاده شود. پس برای این تسک RNN مناسب نیست.

ج) این تسک همچون تسک قسمت الف بوده و تبدیل صوت به متن است. در این فرمت ، ماهیت پردازش داده های صوتی متوالی برای تعیین جنسیت است. در طول زمان با معماری RNN چند به یک تطابق دارد. پس مناسب است.

(B)

این سناریو پیرامون پیش بینی رفتار گربه براساس داده های آب و هوای فعلی و داده های چند روز گذشته است. در این مورد ، با توجه به قسمت های مختلف، گزینه ها بررسی شده اند:

الف) این گزینه استفاده از یک RNN دو طرفه را پیشنهاد می کند زیرا هم اطلاعات گذشته و هم اطلاعات آینده را برای پیش بینی در نظر می گیرد. با این حال، در این سناریو، رفتار گربه احتمالاً بیشتر به داده های آب و هوای گذشته و آب و هوای فعلی بستگی دارد تا داده های آب و هوای

آینده. بنابراین، در نظر گرفتن داده های آب و هوای آینده ممکن است برای پیش بینی دقیق رفتار گربه ضروری یا مرتبط نباشد. پس این گزینه مناسب نیست.

ب) RNN های دو جهته ممکن است زمینه بیشتری را فراهم کنند، اما دقت گرادیان ها در Backpropagation ممکن است لزوماً چالش اصلی در اینجا نباشد. در حالی که RNN های دو طرفه به شبکه اجازه می دهند هم زمینه های گذشته و هم آینده را در نظر بگیرند، این امر دقت بهتری را برای کار خاص پیش بینی رفتار گربه بر اساس داده های آب و هوای تاریخی و فعلی تضمین نمی کند. پس این گزینه مناسب نیست.

ج) این گزینه، استفاده از یک RNN یک طرفه را پیشنهاد می کند، که در آن پیش بینی رفتار گربه در روز t (y_t) تنها به داده های آب و هوای تاریخی تا روز x_1, \dots, x_t متکی است. از آنجایی که رفتار گربه ممکن است لزوماً به داده های آب و هوای آینده بستگی نداشته باشد، یک RNN یک طرفه با بررسی اطلاعات گذشته و فعلی می تواند انتخاب مناسبی باشد.

د) مشابه گزینه قبلی، استفاده از RNN یک طرفه (یک طرفه) را پیشنهاد می کند. با این حال، با بیان اینکه رفتار گربه فقط به داده های آب و هوای روز جاری (x) بستگی دارد، مسئله را بیش از حد ساده می کند. این ممکن است وابستگی های زمانی در داده ها را نشان ندهد، زیرا رفتار ممکن است تحت تأثیر الگوهای آب و هوایی تاریخ قبل نیز قرار گیرد.

(C)

با توجه به خواسته مسئله ، گزینه هارو بررسی میکنیم:

الف) این گزینه بدان معناست که در مرحله زمانی t در RNN ، احتمال کل توالی توکن هارا از y_1 تا y_{t-1} تخمین میزند. با این حال، در RNN در در مرحله زمانی، معمولاً قسمت بعدی دنباله را به جای توزیع احتمال کل دنباله تا آن نقطه پیش بینی می کند.

ب) این گزینه بدان معناست که RNN احتمال اولین توکن که y_1 هست را تخمین می زند. با این حال، در یک مدل زبان RNN، در مرحله زمانی t ، مدل معمولاً توزیع احتمال نشانه بعدی را با توجه به دنباله تا مرحله زمانی $t-1$ پیش بینی می کند.

ج) این انتخاب نشان دهنده احتمال توکن در مرحله زمانی t با توجه به کل دنباله تا آن نقطه (از y_1 تا y_t) است. شامل زمینه دنباله تا مرحله زمانی t است، و با پیش بینی معمولی که توسط یک مدل زبان RNN در مرحله زمانی t انجام می شود، بیشتر همسو می شود. با این حال، به اشتباه شرایط را به عنوان نشانه در مرحله زمانی t مدنظر می گیرند، که اینطور نیست و مدل معمولاً نشانه بعدی را در دنباله پیش بینی می کند.

(د)

این گزینه احتمال توکن در مرحله زمانی t با توجه به دنباله تا مرحله زمانی $t-1$ است. در یک مدل زبان RNN، مدل توزیع احتمال نشانه بعدی را بر اساس دنباله تا مرحله زمانی قبلی $t-1$ پیش بینی می کند.

حال با توجه به بررسی های انجام شده، گزینه د) مناسب است. در مرحله زمانی t ، یک مدل زبانی RNN توزیع احتمال نشانه بعدی که y_t است را با توجه به دنباله تا مرحله زمانی قبلی $t-1$ تخمین میزند. این پیش بینی به تولید توکن بعدی در دنباله بر اساس زمینه ای که از توکن های قبلی آموخته شده است کمک می کند.

$$\frac{\partial J_t}{\partial \sigma_t} = - \sum_{i=1}^2 y_{t,i} \log(\hat{y}_{t,i}) \Rightarrow \hat{y}_{t,i} = \sigma(\sigma_{t,i}) \quad (ع) \quad 4$$

$$\frac{\partial J_t}{\partial \sigma_t} = - \sum_{i=1}^2 \frac{y_{t,i}}{\hat{y}_{t,i}} \cdot \frac{\partial \hat{y}_{t,i}}{\partial \sigma_{t,i}} \Rightarrow \frac{\partial \hat{y}_{t,i}}{\partial \sigma_{t,i}} = \hat{y}_{t,i} (1 - \hat{y}_{t,i}) \quad 5$$

$$\frac{\partial J_t}{\partial \sigma_t} = - \sum_{i=1}^2 \frac{y_{t,i}}{\hat{y}_{t,i}} \cdot \hat{y}_{t,i} (1 - \hat{y}_{t,i}) \Rightarrow \frac{\partial J_t}{\partial \sigma_t} = - \sum_{i=1}^2 (y_{t,i} - \hat{y}_{t,i}) \quad 6$$

$$h_t = \psi(z_t), z_t = w_{hh} \cdot h_{t-1} + w_{hx} \cdot x_t \Rightarrow J_t = - \sum_{i=1}^2 y_{t,i} \log(\hat{y}_{t,i}) \quad (د) \quad 7$$

$$\sigma_t = w_{yh} \cdot h_t \Rightarrow \frac{\partial J_t}{\partial h_i} = \sum_{j=1}^3 \frac{\partial J_t}{\partial \sigma_{t,j}} \cdot \frac{\partial \sigma_{t,j}}{\partial h_i} \Rightarrow \frac{\partial J_t}{\partial h_i} = \sum_{j=1}^3 g_{\sigma_{t,j}} \cdot \frac{\partial \sigma_{t,j}}{\partial h_i} \quad 8$$

$$\Rightarrow \sigma_t \cdot w_{yh} \cdot h_t \Rightarrow \frac{\partial J_t}{\partial h_i} = \sum_{j=1}^3 g_{\sigma_{t,j}} \cdot w_{y h i, j} \quad 9$$

$$\frac{\partial J_t}{\partial w_{hh}} = \sum_{i=1}^3 g_{h_{t,i}} \cdot \frac{\partial h_i}{\partial w_{hh}} \Rightarrow h_i = \psi(z_i) \Rightarrow z_i = w_{hh} \cdot h_{i-1} + w_{hx} \cdot x_i \quad (ج) \quad 10$$

$$\frac{\partial h_i}{\partial w_{hh}} = \frac{\partial \psi(z_i)}{\partial w_{hh}} + \frac{\partial \psi(z_i)}{\partial h_{i-1}} \cdot \frac{\partial h_{i-1}}{\partial w_{hh}} \Rightarrow \frac{\partial J_t}{\partial w_{hh}} = \sum_{i=1}^3 g_{h_{t,i}} \left(\psi'(z_i) + \frac{\partial \psi(z_i)}{\partial h_{i-1}} \cdot \frac{\partial h_{i-1}}{\partial w_{hh}} \right) \quad 11$$

$$\psi'(z_i) = \psi(z_i) \cdot (1 - \psi(z_i)) \Rightarrow \frac{\partial \psi(z_i)}{\partial h_{i-1}} = \psi'(z_i) \cdot w_{hh} \quad (ب) \quad 12$$

$$\Rightarrow \frac{\partial h_{i-1}}{\partial w_{hh}} = h_{i-2} \Rightarrow \frac{\partial J_t}{\partial w_{hh}} = \sum_{i=1}^3 g_{w_{hh} h_{t,i}} \left(\psi(z_i) + \psi(z_i) w_{hh} \cdot h_{i-2} \right) \quad 13$$

سوال 3)

منبع: (gpt3)

(الف)

سوال 3 الف)

برای توجه (attention) ، شباهت (similarity) از ضرب داخلی برداری کم + سینه نرمی (softmax) برای توجه (attention) استفاده می شود. پس به حاصل ضرب شباهت (similarity) و برداری کم + سینه نرمی (softmax) دسترسی داریم.

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ -1 \\ -1 \end{bmatrix} = -2 \quad \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ -1 \\ -1 \end{bmatrix} = 3$$
$$\begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ -1 \\ -1 \end{bmatrix} = 0 \quad \begin{bmatrix} 0 \\ 2 \\ -1 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ -1 \\ -1 \end{bmatrix} = 6$$

یکه چهارم، بیشترین شباهت را دارد. \Rightarrow similarity $[9, 2, 0, 6]$

مقادیر (values) مربوط به این یکبار برای یک $\text{output} = \text{values}[3]$ مربوط به یکبار $\text{output} = \text{values}[3]$

بنابراین ما استفاده از توجه (attention) برای پرسش خود داده شده، صرفی لایه توجه $\begin{bmatrix} 9 \\ 2 \\ 0 \\ 6 \end{bmatrix}$

(ب)

انتخاب طراحی بر مدل های آموزشی با استفاده از مکانیسم های توجه اثر گذار است که در ذیل توضیح دادم. وقتی نحوه انتقال گرادین ها از آخرین لایه به لایه اول شبکه بررسی شود، مهم است که تأثیر استفاده از توجه argmax را در نظر گرفته شود. حال به بررسی آن می پردازیم:

جریان شیب (Gradient Flow): در شبکه های عصبی، در حین backpropagation ، شیب ها از طریق لایه ها به سمت عقب پخش می شوند تا

وزن ها را به روز کنند. البته، با توجه argmax ، که شامل عملیات deterministic و غیر قابل تمایز است (از آنجایی که حداکثر شباهت را بدون در نظر گرفتن توابع قابل تمایز مانند softmax انتخاب می کند)، گرادیان ها ممکن است به آرامی در مکانیسم توجه جریان نداشته باشند.

چالش های آموزشی : ماهیت deterministic ، argmax می تواند منجر به ناپیوستگی گرادیان و مشکلات در انتشار گرادیان های معنی دار از طریق شبکه شود. این می تواند منجر به چالش هایی در یادگیری به طور مؤثر شود، به خصوص هنگام برخورد با داده های پیچیده یا نویزی.

برای بررسی این مسائل و بهبود روند آموزش داریم:

پارامترهای قابل یادگیری برای پرس و جوها و کلیدها: معرفی پارامترهای قابل یادگیری برای پرس و جوها و کلیدها به مدل اجازه می دهد تا به طور تطبیقی نمایش هایی را تنظیم کند و یاد بگیرد که روابط بین پرس و جوها و کلیدها را بهتر به تصویر بکشد. این پارامترهای قابل یادگیری را می توان در طول آموزش به روز کرد و امکان یادگیری مؤثرتر را فراهم می شود.

- منظم سازی و اصلاحات (Regularization and Modifications) : تکنیک هایی مانند روش های منظم سازی (مانند dropout) یا اصلاح مکانیسم توجه می تواند به کاهش ماهیت جبری argmax کمک کند. به عنوان مثال، معرفی یک مکانیسم نرم کننده برای smooth کردن و متمایز کردن توزیع توجه می تواند جریان گرادیان بهتری را در طول تمرین تسهیل کند.

بهبود نمایش داده ها یا کلیدها در طول آموزش با استفاده از آنها به عنوان پارامترهای قابل یادگیری امکان پذیر است. با اجازه دادن به این پارامترها برای به روز رسانی بر

اساس داده های آموزشی، مدل به طور بالقوه می تواند بازنمایی های معنی دار بیشتری بیاموزد و روابط بین پرس و جوها و کلیدها را بهتر نمایش دهد.

بدین سان باوجود اینکه روش هایی برای بهبود فرآیند آموزش توجه argmax هست ولی ممکن است چالش هایی را در گرفتن روابط پیچیده ایجاد نماید و نیاز به مدیریت دقیق در بهینه سازی دارد تا یادگیری موثر داشته باشیم.

سوال 4) گفته شد نیاز به توضیحات ندارد