

به نام خالق رنگین کمان

ستاره باباجانی – 99521109

سوال 1: a) این معماری برای موارد (ب) و (ج) استفاده میشود.

این معماری برای وظایف دسته بندی احساسات و تشخیص جنسیت از گفتار مناسب است ولی برای تشخیص گفتار مناسب نیست زیرا آنها برای وظایفی طراحی شده اند که در آنها دنباله ای از ورودیها داریم و خروجی ما واحد است. در موارد (ب) و (ج) خروجی واحد ولی در مورد (آ) خروجی چندتا است.

حال به بررسی هر گزینه میپردازیم:

1. طبقه بندی احساسات:

- ورودی دنباله ای از کلمات است.
- خروجی یک برچسب احساسی (مثبت یا منفی) است.
- RNN می تواند با در نظر گرفتن ترتیب کلمات و روابط آنها نحوه آشکار شدن احساس در طول متن را تصویر کند.

2. تشخیص جنسیت از گفتار:

- ورودی دنباله ای از فریم های صوتی (کلیپ گفتار) است.
- خروجی یک برچسب تک جنسیتی (مرد یا زن) است.

- RNN می تواند چگونگی تکامل ویژگی های صوتی مانند آهنگ های زیر و بمی صدا و آهنگ را ردیابی کند.

3. چرا تشخیص گفتار مناسب نیست:

- تشخیص گفتار مستلزم نگاشت دنباله ای از فریم های صوتی به دنباله ای از متن است. (پس به معماری many-to-many احتیاج داریم).
- در اصل، RNN های چند به یک در وظایفی که درک معنی یا الگوی کلی در یک دنباله برای تولید یک طبقه بندی دقیق و واحد ضروری است، برتری می یابند.
- (b) گزینه (ج) یعنی یک طرفه (مقدار y_t تنها به x_1, \dots, x_t وابسته است و به x_{t+1}, \dots, x_{30} وابسته نیست) درست است زیرا اخلاق پنبه به آب و هوای فعلی و چند روز گذشته بستگی دارد و به آب و هوای آینده که هنوز اتفاق نیوفتاده است، بستگی ندارد.
- یک RNN یک طرفه، که توالی آب و هوا را به صورت زمانی پردازش می کند، می تواند به طور موثر این رابطه را مدل کند. استفاده از یک RNN دو جهته غیر ضروری خواهد بود و با در نظر گرفتن اطلاعات نامربوط آینده، به طور بالقوه نویز ایجاد می کند.

بنابراین، یک RNN یک طرفه مناسب‌ترین انتخاب برای این مشکل است، که به طور دقیق از وضعیت خلق و خوی پنبه بر اساس تاریخچه آب و هوای مربوطه عکس می‌گیرد.

(C) گزینه (ج) درست است زیرا در یک مدل زبان سعی میشود که گام بعدی بر اساس دانش تمامی مراحل قبلی پیش بینی شود.

- تخمین $P(y, y_1, \dots, y_{t-1})$ شامل پیش‌بینی کل دنباله تا $t-1$ است که تمرکز RNN در مرحله زمانی t نیست.
- تخمین $P(y_1)$ زمینه کلمات قبلی را نادیده می‌گیرد که برای مدل سازی زبان بسیار مهم است.
- تخمین $P(y_t | y, y_1, \dots, y_t)$ ممکن نیست زیرا RNN هنوز y_t را در مرحله زمانی t مشاهده نکرده است.

سوال 2:

$$\Rightarrow \frac{\partial \log(\hat{y}_{t,i})}{\partial \hat{y}_{t,i}} = \frac{1}{\hat{y}_{t,i}}, \quad \frac{\partial J_t}{\partial \hat{y}_t} = - \sum_{i=1}^V y_{t,i} \times \frac{\partial \log(\hat{y}_{t,i})}{\partial \hat{y}_{t,i}} = - \sum_{i=1}^V y_{t,i} \times \frac{1}{\hat{y}_{t,i}}$$

$$\frac{\partial J_t}{\partial \theta_t} = \frac{\partial J_t}{\partial \hat{g}_t} \times \frac{\partial \hat{g}_t}{\partial \theta_t}$$

$\hat{y}_t = \alpha(\phi_t)$ (توقع، تابع خطی)
 $\Rightarrow \frac{\partial \hat{y}_t}{\partial \phi_t} = \sum_{i=1}^p (y_{t,i} x_{t,i}) \times \alpha'(\phi_t)$
 فقط ϕ_t را تغییر دهیم

ف. حلقه‌ها را از یک منبع است RNN یک حلقه را هم پس مقایسه $\frac{\partial J_t}{\partial h_i}$ فقط برای t :
علاسه می شوند یعنی به آنکه واسطه تلفت :

$$\Rightarrow \frac{\partial J_t}{\partial n} = \frac{\partial J_t}{\partial g_t} \times \frac{\partial g_t}{\partial o_t} \times \left(\frac{\partial o_t}{\partial n} \times \left(\prod_{k=0}^{t-1} \frac{\partial n_{t-k}}{\partial n_{t-k-1}} \right) \right)$$

$g_{ot} = \text{die Chance zu leben}$ $\rightarrow O_t = w_{yn} \cdot h_t \Rightarrow \frac{\partial o_t}{\partial n} = w_{yn}$

folgt $h_t = \psi(z_t) = \psi(w_{nn} \cdot h_{t-1} + w_{nx} \cdot x_t)$

$$\Rightarrow \frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial h_t}{\partial z_t} \times \frac{\partial z_t}{\partial h_{t-1}} = \psi'(z_t) \times w_{hh}$$

$$\frac{\partial J_t}{\partial \theta_i} = g_{\theta_i} \times w_{yn} \times (w_{hh})^{t-i} \times \prod_{k=0}^{t-i-1} v'(z_k) \rightarrow = g_{h_{t_i}}$$

$$i=1 \rightarrow \frac{\partial J_t}{\partial n_i} = g_{0i} \times w_{yn} \times w_{nn}^{t-1} \times \prod_{k=0}^{t-1} v'(z_k), \quad t=1 \rightarrow (t-1)$$

$$\Rightarrow h_t = \psi(w_{hh} \cdot h_{t-1} + w_{hx} \cdot x_t) \quad (2)$$

$$\frac{\partial h_t}{\partial w_{hh}} = \frac{\partial h_t}{\partial z_t} \times \frac{\partial z_t}{\partial w_{hh}} = \psi'(z_t) \times h_{t-1}$$

$$\rightarrow \frac{\partial J_t}{\partial w_{hh}} = \sum_{i=1}^t \left(\frac{\partial J_t}{\partial h_i} \right) \times \frac{\partial h_i}{\partial w_{hh}} = \sum_{i=1}^t g_{h_{t_i}} \times \psi'(z_t) \times h_{i-1} = g_{w_{hh},t}$$

$g_{h_{t_i}} \leftarrow$

$$\Rightarrow \frac{\partial J}{\partial w_{hh}} = \sum_{t=1}^T \frac{\partial J_t}{\partial w_{hh}} = \sum_{t=1}^T g_{w_{hh},t} \quad (3)$$

سوال 3: الف) همان طور که میدانیم این مکانیزم مقداری را انتخاب میکند که بیشترین شباهت را به پرس و جو با استفاده از ضرب داخلی دارد.

حال برای محاسبه بیشترین شباهت (خروجی) مراحل زیر را طی میکنیم:

• محاسبه ضرب داخلی بین پرس و جو و کلیدها:

$$\Rightarrow q \cdot \text{Keys}[0] = 3 * 1 + (-1 * 2) + (-1 * 3) = -2$$

$$\Rightarrow q \cdot \text{Keys}[1] = 3 * 2 + (-1 * 2) + (-1 * 1) = 3$$

$$\Rightarrow q \cdot \text{Keys}[2] = 3 * 0 + (-1 * 1) + (-1 * -1) = 0$$

$$\Rightarrow q \cdot \text{Keys}[3] = 3 * 0 + (-1 * -2) + (-1 * -4) = 6$$

• اعمال argmax روی مقادیر بدست آمده برای انتخاب اندیس

مشابه ترین:

$$\Rightarrow \text{argmax}([-2, 3, 0, 6]) = 3$$

• محاسبه مقدار خروجی از طریق اندیس محاسبه شده:

$$\Rightarrow \text{Output} = \text{Values}[3] = [6, 1, 2]$$

پس خروجی لایه برای q داده شده با استفاده از این مکانیزم، [6,1,2] است.

ب) استفاده از argmax در زمینه مکانیسم های توجه به دلیل ماهیت غیر قابل تمایز آن با وجود اینکه از نظر محاسباتی ساده است، می تواند چالش هایی را در طول آموزش ایجاد کند:

- تابع argmax شاخص حداکثر مقدار را در مجموعه ای از مقادیر انتخاب می کند، به این معنی که گسسته است و به راحتی قابل بهینه سازی مبتنی بر گرادیان نیست: در مکانیسم های توجه، تابع softmax معمولاً برای استخراج وزن های توجه بر اساس نمرات شباهت (محاسبه شده با استفاده از پرس و جو و کلیدها) استفاده می شود که مجموع نمرات را 1 می کند و امکان توزیع یکنواخت و احتمالی را فراهم می کند که می توان آن را متمایز کرد و در پس انتشار برای آموزش استفاده کرد. وقتی از argmax به جای softmax استفاده می شود، یک ناپیوستگی در گرادیان ها ایجاد می کند. این ناپیوستگی محاسبه مستقیم گرادیان ها با توجه به پارامترهای کلیدها و پرس و جوها را غیرممکن می کند زیرا عملیات argmax خود دارای یک گرادیان کاملاً تعریف شده نیست. در نتیجه، این جریان موثر شیب ها را در طول شبکه در طول انتشار پس از آن مهار می کند.

- توجه argmax ممکن است برای کارهایی که به درک روابط پیچیده بین کلیدها و پرس و جوها دارند، بیش از حد ساده باشد بطوریکه تفاوت‌های کوچک یا شباهت در توزیع‌ها را در بر نگیرد.
- چون خروجی بر اساس شبیه‌ترین کلید بدون هیچ گونه تبدیل قابل به یادگیری است، مدل نمیتواند در طول فرآیند آموزش بهبود کلیدها را بیاموزد.
- به جای استفاده از argmax ، استفاده از ترفند-Gumbel Softmax، می تواند مفید باشد. این روش یک عنصر تصادفی را معرفی می کند که عملیات argmax گسسته را با یک تقریب پیوسته و قابل تمایز تقریب می زند. این به گرادیان‌ها اجازه می‌دهد تا در طول آموزش از طریق مکانیسم توجه جریان پیدا کنند و شبکه را قادر می‌سازد تا نمایش‌های بهتری از پرس و جوها و کلیدها را بیاموزد.

پایان