



به نام خدا
دانشگاه تهران
دانشکده مهندسی
برق و کامپیوتر



درس شبکه‌های عصبی و یادگیری عمیق
تمرین چهارم

محمد سپهری	نام دستیار طراح	پرسش ۱
msepehri898@gmail.com	رایانامه	
پرهام بیچرانلو	نام دستیار طراح	پرسش ۲
Parhambicharanlu1378@gmail.com	رایانامه	
۱۴۰۲.۰۲.۲۷	مهلت ارسال پاسخ	

قوانین.....	ت
پرسش ۱. توصیف عکس.....	۱
۱-۱. مجموعه دادگان و پیش پردازش آنها.....	۲
۲-۱. مدل شبکه.....	۳
۳-۱. پیش‌بینی شبکه.....	۳
۴-۱. پرسش‌ها.....	۵
پرسش ۲. تشخیص اندیشه.....	۶
۱-۲. معماری LSTM و embedding.....	۶
۲-۲. پیش پردازش دادگان.....	۷
۳-۲. پیاده سازی طبقه بندی نیت.....	۷
۴-۲. پیاده سازی مدل Responder	۷

شکل‌ها

- شکل ۱. خروجی یک مدل آموزش دیده برای Image Captioning ۱
- شکل ۲. تصویر مدل مورد بررسی در سوال اول ۳
- شکل ۳. نحوه استفاده از مدل در زمان تست جهت تولید جمله ۴
- شکل ۴. الگوریتم بازگو کننده شبکه شکل ۳ جهت تولید جمله ۴

قبل از پاسخ دادن به پرسش‌ها، موارد زیر را با دقت مطالعه نمایید:

- از پاسخ‌های خود یک گزارش در قالبی که در صفحه‌ی درس در سامانه‌ی Elearn با نام **REPORTS_TEMPLATE.docx** قرار داده شده تهیه نمایید.
- پیشنهاد می‌شود تمرین‌ها را در قالب گروه‌های دو نفره انجام دهید. (بیش از دو نفر مجاز نیست و تحویل تک نفره نیز نمره‌ی اضافی ندارد) توجه نمایید الزامی در یکسان ماندن اعضای گروه تا انتهای ترم وجود ندارد. (یعنی، می‌توانید تمرین اول را با شخص A و تمرین دوم را با شخص B و ... انجام دهید)
- **کیفیت گزارش شما در فرآیند تصحیح از اهمیت ویژه‌ای برخوردار است؛** بنابراین، لطفاً تمامی نکات و فرض‌هایی را که در پیاده‌سازی‌ها و محاسبات خود در نظر می‌گیرید در گزارش ذکر کنید.
- در گزارش خود مطابق با آنچه در قالب نمونه قرار داده شده، برای شکل‌ها زیرنویس و برای جدول‌ها بالانویس در نظر بگیرید.
- الزامی به ارائه توضیح جزئیات کد در گزارش نیست، اما باید نتایج بدست آمده از آن را گزارش و تحلیل کنید.
- **تحلیل نتایج الزامی می‌باشد، حتی اگر در صورت پرسش اشاره‌ای به آن نشده باشد.**
- **دستیاران آموزشی ملزم به اجرا کردن کدهای شما نیستند؛** بنابراین، هرگونه نتیجه و یا تحلیلی که در صورت پرسش از شما خواسته شده را به طور واضح و کامل در گزارش بیاورید. در صورت عدم رعایت این مورد، بدیهی است که از نمره تمرین کسر می‌شود.
- **کدها حتماً باید در قالب نوت‌بوک با پسوند .ipynb تهیه شوند، در پایان کار، تمامی کد اجرا شود و خروجی هر سلول حتماً در این فایل ارسالی شما ذخیره شده باشد.** بنابراین برای مثال اگر خروجی سلولی یک نمودار است که در گزارش آورده‌اید، این نمودار باید هم در گزارش هم در نوت‌بوک کدها وجود داشته باشد.
- **در صورت مشاهده‌ی تقلب امتیاز تمامی افراد شرکت‌کننده در آن، 100- لحاظ می‌شود.**
- تنها زبان برنامه نویسی مجاز **Python** است.
- **استفاده از کدهای آماده برای تمرین‌ها به هیچ وجه مجاز نیست.**
- نحوه محاسبه تاخیر به این شکل است: پس از پایان رسیدن مهلت ارسال گزارش، حداکثر تا یک هفته امکان ارسال با تاخیر (به ازای هر روز 5 درصد کسر نمره) وجود دارد، پس از این یک هفته نمره آن تکلیف برای شما صفر خواهد شد.

- لطفا گزارش، کدها و سایر ضمایم را به در یک پوشه با نام زیر قرار داده و آن را فشرده سازید، سپس در سامانه‌ی Elearn بارگذاری نمایید:

HW[Number] _[Lastname]_[StudentNumber]_[Lastname]_[StudentNumber].zip

(مثال: HW1_Ahmadi_810199101_Bagheri_810199102.zip)

- برای گروه‌های دو نفره، بارگذاری تمرین از جانب یکی از اعضا کافی است ولی پیشنهاد می‌شود هر دو نفر بارگذاری نمایند.

پرسش ۱. توصیف عکس^۱

یکی از حوزه‌های جذاب در یادگیری ماشین، توصیف یک عکس با یک جمله است. در واقع هدف ایجاد و آموزش مدلی است که بتواند یک تصویر را به عنوان ورودی بگیرد و در نهایت یک جمله در توصیف آن عکس در خروجی خود تولید کند. تصویر زیر نمونه‌ای از خروجی این شبکه را نشان می‌دهد.



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."

شکل ۱. خروجی یک مدل آموزش دیده برای **Image Captioning**

حال در این تمرین قصد داریم یک مدل برای رسیدن به این هدف پیاده‌سازی نماییم. ساختار کلی این مدل‌ها به این صورت است که یک شبکه CNN جهت تولید ویژگی‌های تصاویر وجود دارد و در کنار آن روش‌های مختلفی برای Embedding جملات موجود است که در نهایت بردار ویژگی تصاویر و متن در کنار هم قرار گرفته و به عنوان ورودی یک شبکه بازگشتی اعمال می‌شود تا در نهایت جمله نهایی را تولید نماید. در ادامه بیشتر با بخش‌های مختلف آن آشنا خواهید شد. مقاله‌ای که شما در این بخش از تمرین می‌توانید به آن رجوع کنید مقاله [Image Captioning](#) است که به پیوست هم برای شما قرار داده شده است.

¹ Image Captioning

۱-۱. مجموعه دادگان و پیش پردازش آنها

با مطالعه مقاله اشاره شده متوجه می‌شوید که سه مجموعه داده معرفی شده است. مجموعه دادگانی که باید شما در این تمرین استفاده کنید flickr8k است که مجموعه دادگانی با سایز کوچکتر در مقاله اشاره شده است. این مجموعه داده را می‌توانید از پیوند زیر دریافت کنید:

<https://www.kaggle.com/datasets/adityajn105/flickr8k>

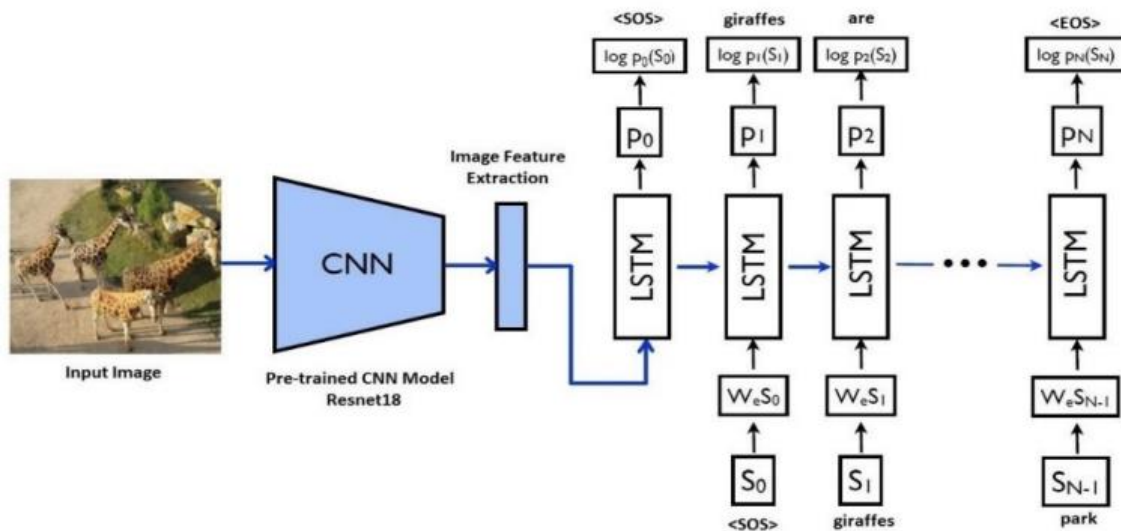
این مجموعه از دو بخش به نام Image و Caption.txt تشکیل شده است که پوشه Image شامل ۸۰۹۱ تصویر و Caption.txt شامل ۴۰۴۵۵ جمله است که برای هر تصویر ۵ جمله مختلف توسط افراد مختلف جمع آوری شده است. در کنار هر جمله نام تصویر مورد نظر نیز آورده شده است. با آماده سازی تصاویر برای اعمال به شبکه‌های کانولوشنی پیش‌تر آشنا شدید. در اینجا جملات نیز باید پیش‌پردازش شوند تا به بردارهایی از اعداد تبدیل شوند. ما در اینجا برای سادگی پیشنهاد می‌کنیم که از لایه Embedding در پایتورچ استفاده کنید که نحوه کار با این لایه را در پیوند زیر مشاهده می‌کنید: (شما می‌توانید از سایر روش‌ها هم به انتخاب خودتان بهره ببرید که نیاز هست که در گزارش خودتان به آن اشاره کنید)

<https://pytorch.org/docs/stable/generated/torch.nn.Embedding.html>

پارامتری به نام Embedding_dim در آن وجود دارد که می‌توانید آن را ۳۰۰ در نظر بگیرید که البته انتخاب آن در اختیار شما می‌باشد که در واقع این عدد مشخص می‌کند که برای هر کلمه یک بردار عددی با طول ۳۰۰ در نظر بگیرید. نکته که مهمی که در پیش‌پردازش داده‌ها باید توجه نمایید، این است که باید برای هر جمله از توکن‌های شروع و پایان جمله <SOS> و <EOS> استفاده نماییم. که توکن‌های خاصی می‌باشد که توسط خود شما تعریف می‌شوند. همچنین باید مجموعه لغات موجود در مجموعه دادگان خود را پردازش و به هر کدام از آنها یک Index نسبت دهید. بهتر است علامت‌های نگارشی از جملات حذف شوند. همچنین از آنجایی که جملات Caption‌ها طول‌های متفاوتی دارند باید طول آن‌ها باهم یکسان شوند، که این کار را با Padding مناسب می‌توانید انجام دهید که می‌توان یک طول مشخص ثابت را در نظر گرفت یا یکسان‌سازی را در هر mini batch انجام داد.

۲-۱. مدل شبکه

در شکل شماره ۲ مدل کلی مد نظر را مشاهده می‌کنید. همان‌طور که مشاهده می‌کنید، بخشی از مدل جهت استخراج ویژگی تصاویر مورد استفاده قرار می‌گیرد. در این مسئله ما قصد داریم از یک مدل از پیش آموزش داده شده Resnet18 استفاده نماییم. این مدل در کتابخانه پایتورچ قابل دسترس می‌باشد و از آخرین لایه شبکه کانولوشنی آن ویژگی‌های تصویر استخراج می‌شود که در نهایت نیاز است به یک لایه خطی جهت استخراج ویژگی‌های مورد نظر با ابعاد مناسب جهت ورود به شبکه بازگشتی، استفاده نمود.

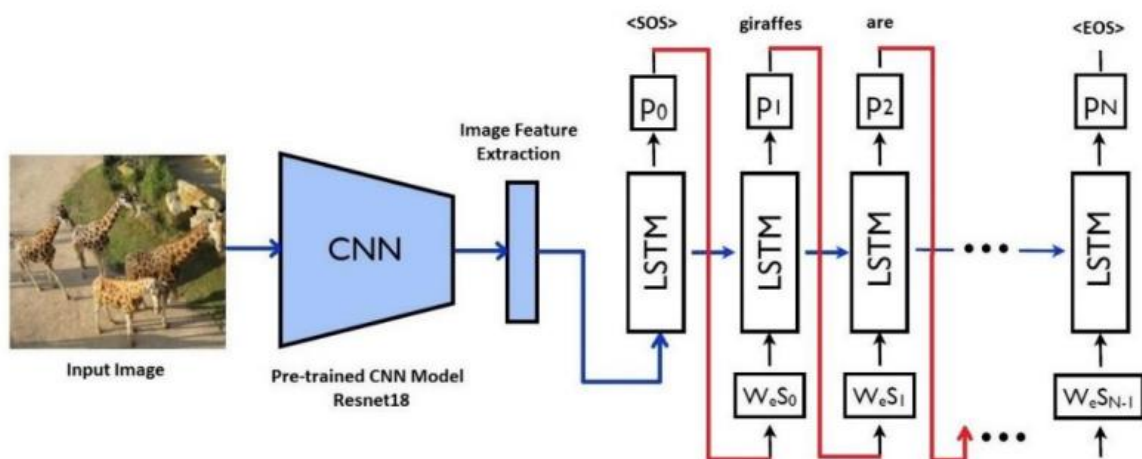


شکل ۲. تصویر مدل مورد بررسی در سوال اول

در این قسمت از یک لایه شبکه LSTM با تعداد ۲۵۶ لایه پنهان استفاده می‌نماییم و بردارهای Embed شده جملات در کنار بردار تصویر به آن داده شده و خروجی آن به یک لایه خطی به سبب ورودی Hidden State و سبب خروجی تعداد کلمات موجود در مجموعه دادگان اعمال می‌شود و به این ترتیب به محاسبه خطا و پیش‌بینی مدل می‌پردازیم.

۳-۱. پیش‌بینی شبکه

بعد از آموزش شبکه، نیاز دارید تا شبکه را ارزیابی نمایید. جهت ارزیابی شبکه باید به صورتی که در شکل ۳ نشان داده شده از شبکه استفاده نماییم.



شکل ۳. نحوه استفاده از مدل در زمان تست جهت تولید جمله

همان‌طور که می‌دانیم در زمان تست شبکه آموزش داده شده، Caption وجود ندارد و ما باید برای یک تصویر Caption تولید نماییم. برای این منظور روش‌های مختلفی وجود دارد ولی ما در اینجا مدل بالا را پیشنهاد می‌دهیم. در یک تابع به عنوان ورودی، تصویر تست و مدل آموزش داده شده را جهت پیش‌بینی کلمات اعمال می‌کنیم. قطعه کد زیر الگوریتم این شبکه را نمایش داده‌است.

```
input_data = Trained_Model.CNN(image)
states = None #(Hn, Cn)

for _ in range(max_length):
    hiddens, states = Trained_Model.lstm(input_data, states)
    output = Trained_Model.linear(hiddens)
    predicted_index = output.argmax()
    input_data = Trained_Model.Embedding(predicted_index)
    caption_prediction.append(predicted_index)

    if predicted_index.item() == "<EOS>":
        break
```

شکل ۴. الگوریتم بازگو کننده شبکه شکل ۳ جهت تولید جمله

در نهایت caption_prediction مجموعه index های کلمات می باشد که در نهایت به کمک دایره لغات موجود در مجموعه دادگان قابل تبدیل به کلمات می باشد. توجه داشته باشید که الگوریتم فوق فقط مراحل کار را نشان داده است و نیاز به بازنویسی درست، رعایت ابعاد تنسورها و غیره دارد که بر عهده شما می باشد. البته استفاده از هر شیوه دیگری جهت تست و تولید جملات بلامانع است.

۴-۱. پرسش ها

در این بخش به پرسش های زیر با توجه به بخش های پیش پردازش، مدل شبکه و پیش بینی شبکه برای هر پرسش پاسخ دهید:

۱. از یک مدل از پیش آموزش داده شده Resnet18 به عنوان شبکه CNN استفاده نمایید و به جز لایه خطی آخر تمامی لایه های آن را Freeze نمایید تا در عملیات بروزرسانی وزن ها شرکت نداشته باشند. سپس خروجی آن را در کنار بردارهای Embed شده جملات به یک لایه شبکه LSTM یک طرفه اعمال کرده و نمودار خطای آموزش و تست را در طول یادگیری گزارش نمایید. از تابع خطای CrossEntropy و تابع بهینه ساز Adam می توانید استفاده نمایید. بعد از فرآیند آموزش، ۳ عدد عکس از دادگان تست را جهت پیش بینی مدل، به آن اعمال کرده و خروجی آن را در گزارش کار خود ذکر نمایید. (۵۰ نمره)

(جزئیات بارم: پیش پردازش: ۱۰ نمره، مدل شبکه: ۱۰ نمره، پیش بینی شبکه: خروجی خطا: ۱۵ نمره و خروجی تصویر: ۱۵ نمره)

۲. با حفظ موارد گفته شده سؤال قبل تمامی لایه های شبکه Resnet18 را Unfreeze نمایید و مجدداً موارد خواسته شده در سؤال قبل را بررسی نمایید و نتایج بدست آمده را با سؤال قبل مقایسه کنید. (۵۰ نمره)

(جزئیات بارم: پیش پردازش: ۱۰ نمره، مدل شبکه: ۱۰ نمره، پیش بینی شبکه: خروجی خطا: ۱۵ نمره و خروجی تصویر: ۱۵ نمره)

پرسش ۲. تشخیص اندیشه^۱

مسئله‌ی پرسش و پاسخ یکی از مهم‌ترین و پیچیده‌ترین مسائل شناخته شده در حوزه پردازش زبان‌های طبیعی و بازیابی اطلاعات است. در مسئله‌ی پرسش و پاسخ انتظار داریم مدل به پرسش‌های داده شده به طور خودکار پاسخ مرتبط و درست بدهد. یکی از مراحل مهم برای انجام این وظیفه تشخیص اندیشه‌ی پرسش است. برای همین طبقه بندی اندیشه‌ی پرسش‌ها به منظور درک هدف درخواست‌های کاربر برای پاسخ سریع و دقیق می‌تواند بسیار کمک کننده باشد.

دقت کنید که در اینجا طبقه بندی اندیشه به معنای طبقه بندی سوالات با توجه به معنی آن‌ها نیست، بلکه هدف اینکه سوالات را با توجه به نوع پاسخ آن‌ها طبقه بندی کنیم.

هدف این تمرین این است که با توجه به جزئیات [مقاله‌ی Intent Classification in Question-Answering Using LSTM Architectures](#) یک طبقه بند اندیشه پیاده سازی کنید و در ادامه با کمک آن، مدل شما بتواند به سوالات پاسخ مرتبط بدهد. برای پیاده سازی درست حتما مقاله با دقت خوانده شود.

برای این پرسش از مجموعه دادگان **TREC** استفاده می‌کنید. همچنین برای قسمت آخر پرسش از دادگان QA_data استفاده کنید. همه‌ی این دیتاست‌ها در فایل تمرین پیوست شده‌اند. دیتاست TREC شامل فایل train با 5500 پرسش و فایل test نیز با 500 پرسش همراه با دسته‌های اندیشه و زیر دسته‌های آن‌ها است. فایل QA_data همان فایل test است که پاسخ پرسش‌ها هم در یک ستون مجزا کنار بقیه ستون‌ها قرار دارد.

۲-۱. معماری LSTM و embedding

(۱۵ نمره)

توضیح دهید که مزیت معماری LSTM به معماری RNN چیست. سپس علت استفاده از word embedding را بیان کنید و روش‌های تولید آن را شرح دهید. برای کلماتی که چند معنی متفاوت دارند آیا GloVe embedding مناسب است و چرا؟

¹ Intent Classification

۲-۲. پیش پردازش دادگان

(۱۵ نمره)

ابتدا دادگان را با استفاده از پیش پردازش‌های مورد نیاز مثل Tokenization، Normalization و سایر موارد آماده استفاده کنید.

۲-۳. پیاده سازی طبقه بندی نیت

(۵۰ نمره)

در این مقاله برای طبقه بندی دسته‌های اندیشه دو مدل پیشنهاد داده شده است که مدل اول فقط دسته‌ی سطح اول را پیش بینی می‌کند و مدل دوم هر دو سطح دسته را پیش‌بینی می‌کند. هر دو مدل را برای اندازه‌های متفاوت hidden state (25 و 100) پیاده سازی کنید و با مجموعه دادگان آموزش دهید. نمودارهای دقت و خطا را در طول زمان یادگیری رسم کنید. همچنین Confusion matrix، F1-Score، Recall، Precision، Accuracy را گزارش کنید. در آخر نتایج را تحلیل کنید.

۲-۴. پیاده سازی مدل Responder

(۲۰ نمره)

برای اینکه اهمیت طبقه بند اندیشه را درک کنید، مدل پیشنهادی مقاله با نام مدل Responder را که به پرسش‌های پرسیده شده جواب مرتبط(نه الزاما درست) می‌دهد را پیاده سازی کنید و آن را روی دادگان فایل QA_data که شامل 500 پرسش و پاسخ است آموزش دهید. سپس پرسش‌های زیر را به عنوان نمونه به مدل بدهید و پاسخی که مدل می‌دهد را گزارش کنید.

- How many people speak French?
- What day is today?
- Who will win the war?
- Who is Italian first minister?
- When World War II ended?
- When Gandhi was assassinated?