



- مهلت ارسال پاسخ تا ساعت ۵۹ : ۲۳ روز مشخص شده است.
- در طول ترم امکان ارسال با تاخیر پاسخ همه ی تمرین (به استثنای هفته ی امتحان میانترم) تا سقف هفت روز وجود دارد. پس از گذشت این مدت، پاسخ های ارسال شده پذیرفته نخواهند بود.
- همکاری و همفکری شما در انجام تمرین مانعی ندارد اما پاسخ ارسالی هر کس حتما باید توسط خود او نوشته شده باشد.
- در صورت همفکری و یا استفاده از هر منابع خارج درسی، نام همفکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
- لطفا تصویری واضح از پاسخ سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.

مسئله ی ۱. (۴۰ نمره)

در این سوال، شما با مدل های مولد مبتنی بر شبکه های عصبی پیچشی و به طور خاص مدل PixelCNN آشنا خواهید شد. هدف این سوال پیاده سازی مدل PixelCNN بر روی مجموعه داده MNIST است که یک مجموعه از تصاویر دست نوشته اعداد است. شما باید مدل PixelCNN را طوری طراحی کنید که بتواند تصاویر MNIST را به صورت مولد (یعنی پیکسل به پیکسل) تولید کند.

۱. آماده سازی داده ها

مجموعه داده MNIST را با استفاده از کتابخانه های رایج یا به صورت دستی بارگذاری کرده و داده ها را به شکل مناسب برای ورودی مدل PixelCNN تبدیل کنید.

۲. پیاده سازی مدل PixelCNN

یک شبکه عصبی پیچشی طراحی کنید که به صورت پیکسل به پیکسل کار کند. برای این کار از کانولوشن های ماسک دار استفاده کنید تا مدل فقط بتواند پیکسل های قبلی را برای پیش بینی پیکسل فعلی ببیند.

- برای پیاده سازی این شبکه عصبی، از دو نوع ماسک متفاوت A و B باید استفاده کنید. تفاوت این دو نوع را شرح دهید...

- برای پیاده سازی مدل می توانید از تنظیمات زیر استفاده کنید. استفاده از این تنظیمات اجباری نبوده و در صورتی که با هایپرپارامترهای دیگر جواب مناسبی می گیرید، میتوانید این مقادیر را تغییر دهید.

- به طور کلی از ۸ لایه پیچشی استفاده کنید به گونه ای که اندازه کرنل ۶ لایه اول ۷ باشد و اندازه کرنل دو لایه آخر ۱ باشد.

- تعداد کانال ها را ۶۴ در نظر بگیرید.

- بعد از هر لایه پیچشی از تابع فعال ساز ReLU استفاده کنید. بعد از آخرین لایه از تابع Sigmoid بهره ببرید.

۳. آموزش مدل

مدل را روی مجموعه داده آموزش دهید. هایپرپارامترها از قبیل نرخ یادگیری و تعداد دوره را به گونه ای تنظیم کنید که مدل آموزش ببیند. سپس نمودار خطای داده آموزش و تست را بر حسب دوره را رسم کرده و در گزارش خود نمایش دهید.

۴. تولید نمونه

پس از آموزش مدل ۱۰۰ نمونه تصویر تولید کرده و در گزارش نمایش دهید.

تمام قسمت ها باید توسط خودتان پیاده سازی شود. در صورت استفاده از کدهای آماده و موجود در اینترنت نمره ای به شما تعلق نمی گیرد. در صورتی که در قسمت های غیر اصلی از کدهای موجود در اینترنت استفاده می کنید، باید منبع خود را ذکر کنید.

مسئله‌ی ۲. (۶۰ نمره)

در این سوال، شما با مدل‌های مولد که از شبکه‌های RNN برای تولید متن استفاده می‌کنند، آشنا خواهید شد. هدف این سوال پیاده‌سازی یک مدل Decoder برای تولید متن فارسی است که با استفاده از داده‌های ویکی‌پدیا فارسی آموزش دیده است. همچنین، در این سوال با متد Transfer Learning و استفاده از مدل‌های از پیش آموزش دیده نیز آشنا خواهید شد.

۱. آماده سازی و پیش پردازش دادگان

در این تمرین، ابتدا لازم است داده‌های ویکی‌پدیا فارسی را از پلتفرم Kaggle بارگذاری کنید. برای این کار می‌توانید از کد زیر استفاده کنید:

```
۱
۲ mkdir ~/.kaggle
۳ cp ./kaggle.json ~/.kaggle/
۴ chmod 600 ~/.kaggle/kaggle.json
۵ kaggle datasets download miladfa7/persian-wikipedia-dataset -f Persian-WikiText-1.txt
۶
```

پس از بارگذاری داده‌ها، باید پیش پردازش‌های لازم را انجام دهید و آن‌ها را به شکلی مناسب برای ورودی مدل RNN تبدیل کنید. توجه داشته باشید که انجام صحیح این مرحله در عملکرد مدل بسیار اهمیت دارد. داده‌های ویکی‌پدیا معمولاً شامل نویز و اطلاعات غیرضروری هستند. بنابراین، کاراکترها، کلمات، یا هر قسمتی که باید از متن حذف شود را شناسایی کنید. به عنوان مثال، stop word ها معمولاً از داده‌های ورودی مدل حذف می‌شوند. توضیح دهید چرا بهتر است این کلمات از متن حذف شوند و همچنین هر مورد مشابه دیگر را با ذکر دلیل حذف کنید.

سپس با استفاده از کتابخانه hazm، متن را نرمال کنید. از متد lemmatize در این کتابخانه استفاده کرده و علت استفاده از آن را توضیح دهید.

در مرحله آخر، داده‌ها را توکنایز کنید. این یکی از مهم‌ترین مراحل پیش‌پردازش داده‌های متنی است. Byte-Pair Encoding (BPE) یکی از متدهای رایج برای توکنایز کردن داده‌های متنی است. در این بخش، درباره نحوه عملکرد این الگوریتم و مزایا و معایب آن توضیح دهید.

۲. آموزش مدل

در این تمرین، هدف ما آموزش یک مدل با استفاده از متد n-gram است. ابتدا در مورد نحوه عملکرد این متد تحقیق کنید.

سپس با استفاده از کتابخانه PyTorch، Dataset و Dataloader را به گونه‌ای پیاده‌سازی کنید که برای آموزش مدل به صورت n-gram مناسب باشد. یک مدل RNN طراحی کنید که برای تسک تولید توکن بعدی (Next Token Generator) مناسب باشد. شما می‌توانید از هر یک از انواع مدل‌های RNN استفاده کنید. تعیین معماری و هایپرپارامترهای مدل نیز بر عهده شماست. توجه داشته باشید که تعداد پارامترهای مدل باید با توجه به نوع تسک و حجم داده‌های آموزشی به درستی انتخاب شود.

سپس مدل طراحی شده را با استفاده از یک تابع loss مناسب آموزش دهید. نمودار کاهش loss در طول مراحل آموزش را رسم کرده و تمامی پارامترهای استفاده شده در مدل را در گزارش خود ذکر کنید.

۳. ارزیابی مدل

درباره معیار Perplexity برای ارزیابی مدل‌های زبانی تحقیق کرده و از آن برای ارزیابی مدل آموزشی خود استفاده کنید.

۴. تولید متن

در این بخش، از مدل آموزش دیده خود استفاده کنید. یک قسمت از یک متن را به عنوان ورودی به مدل بدهید و ادامه متن را با توجه به این ورودی تولید کنید.

۵. استفاده از ترنسفورمرها

در این بخش، به منظور تولید متن و ارتقاء کیفیت متن تولید شده در قسمت قبل با استفاده از دیتاست آماده شده در قسمت قبل، یک مدل ترنسفورمر از ابتدا (فراماسکرچ) طراحی کرده و خروجی های آن را با نتایج بخش قبلی مقایسه کنید.

۶. بهینه سازی محاسبات توجه در ترانسفورمرها

فرض کنید یک مدل ترانسفورمر با طول توالی ورودی $n = 4096$ و ابعاد فضای نمایشی $d = 512$ دارید. پیچیدگی محاسباتی مدل ترانسفورمر معمولی برابر $O(n^2 \cdot d)$ است که برای ورودی های طولانی بسیار پرهزینه می شود. اکنون به سوالات زیر با دقت و محاسبات کامل پاسخ دهید:

- (آ) تعداد کل عملیات محاسباتی در محاسبات توجه استاندارد را برای این مدل محاسبه کنید.
- (ب) فرض کنید از تکنیک KV Cache استفاده می شود که فقط نیاز به محاسبه توجه برای ۲۰٪ از توالی های جدید است. تعداد کل عملیات محاسباتی را با این فرض دقیقاً محاسبه کنید.
- (ج) حال فرض کنید از تکنیک Grouped Query Attention نیز استفاده می شود، به طوری که محاسبات برای این ۲۰٪ از توالی ها به ۶۴ گروه تقسیم شود و به صورت موازی انجام شود. تعداد عملیات محاسباتی را با دقت محاسبه کنید.
- (د) فرض کنید علاوه بر استفاده از KV Cache و Grouped Query Attention، به دلیل محدودیت های سخت افزاری و حافظه، هزینه سر بار ناشی از مدیریت حافظه پنهان و تقسیم بندی گروه ها باعث افزایش ۳۵ درصدی کل زمان محاسبات شود. تعداد کل عملیات نهایی (با در نظر گرفتن سر بار) را دقیقاً محاسبه کنید.
- (ه) اگر به دلیل تنگنای حافظه، پردازنده تنها بتواند ۷۵٪ از عملکرد معمول خود را ارائه دهد، چند عملیات اضافی ناشی از این محدودیت به محاسبات اضافه می شود؟
- (و) در نهایت، با فرض اینکه از روش بهینه سازی محاسباتی مانند Linformer استفاده می شود که پیچیدگی زمانی را به $O(n \cdot \log(n) \cdot d)$ کاهش می دهد، تعداد عملیات نهایی را برای شرایطی که از KV Cache، Grouped Query Attention و محدودیت های سخت افزاری استفاده می شود، محاسبه کنید. آیا این بهینه سازی در شرایط خاص (مانند داده های بسیار طولانی) موثرتر از روش های قبلی خواهد بود؟ توضیح دهید چرا یا چرا نه.

تمام قسمت ها باید توسط خودتان پیاده سازی شود. در صورت استفاده از کدهای آماده و موجود در اینترنت نمره ای به شما تعلق نمی گیرد. در صورتی که در قسمت های غیر اصلی از کدهای موجود در اینترنت استفاده می کنید، باید منبع خود را ذکر کنید.

موفق باشید (: