

تمرین سری چهارم

درس یادگیری عمیق

در این تمرین هدف پیاده‌سازی تسک پردازش گفتار با یک مدل end to end شبکه عصبی می‌باشد. مجموعه دادگان مورد استفاده در لینک زیر قرار دارند جهت آشنایی به لینک زیر مراجعه کنید.

[voxforge dataset](#)

یک فایل برای دانلود و پیش‌پردازش داده‌ها آپلود شده است. کتابخانه مورد نیاز این تمرین pytorch می‌باشد. که ورژن‌های مورد نیاز در فایل توضیحات جداگانه آمده است. **توصیه:** اجرای تسک‌های مربوط به speech نیازمند زمان آموزش و حافظه قابل توجهی هستند. بنابراین توصیه می‌شود از google colab استفاده شود.

مرحله 0

دانلود و پیش‌پردازش داده‌ها:

با استفاده از script موجود در فایل‌های دانه‌دوی با اجرای کد download_prepare_data.sh داده‌های شما دانلود و پیش‌پردازش خواهند شد. برای جزییات بیشتر به فایل readme.md مراجعه کنید.

مرحله 1

معرفی ساختار شبکه

همانطور که می‌دانیم شبکه‌های موسوم به transformer networks با حذف لایه‌های بازگشتی و استفاده از لایه‌های feedforward و یک مدل attention جدید بسیار مورد توجه قرار گرفت. در این تمرین می‌خواهیم تسک پردازش گفتار را با استفاده از یک شبکه transformer به صورت end to end پیاده‌سازی کنیم (برای سادگی قسمت decoder مربوط به شبکه transformer را در نظر نمی‌گیریم). برای پیاده‌سازی این قسمت دقت داشته باشید که ابتدا دانلود و پیش‌پردازش را روی داده‌ها، توسط کد داده شده انجام می‌دهیم. سپس برای استفاده از مدل transformer داده‌های صوتی ابتدا وارد قسمت input embedding و سپس وارد Encoder می‌شوند. و داده‌های متنی متناسب با داده صوتی داده شده به قسمت output embedding بدون وارد شدن به قسمت decoder وارد تابع هزینه ctc می‌شود). ساختار شبکه transformer به صورت زیر می‌باشد: (برای درک بهتر و جزییات بیشتر این [مقاله](#) مطالعه شود).

Transformer Network Modules:

- ❖ Input Embedding
- ❖ Output Embedding (پیاده‌سازی لازم نیست)
- ❖ Positional Encoding
- ❖ Encoder
- ❖ Decoder (پیاده‌سازی لازم نیست)

برای embed کردن داده‌های ورودی صوتی از ساختار شبکه زیر استفاده می‌کنیم.

Input Embedding (emb_cnn feature embedding):

- (0): Conv2d(1,32, kernel_size=(41, 11), stride=(2, 2), padding=(0, 10))
- (1): BatchNorm2d(32, eps=1e-05, momentum=0.1)

تابع بهینه‌ساز (optimizer): برای بهینه‌سازی از تابع بهینه‌ساز Adam با پارامتر های زیر استفاده کنید

$$adam\ optimizer : \beta_1 = 0.9, \beta_2 = 0.98, \varepsilon = 10^{-9}$$

و طبق فرمول زیر learning rate را در طول زمان آموزش تغییر دهید:

$$lrate = d_{model}^{-0.5} \cdot \min(step_num^{-0.5}, step_num \cdot warmup_steps^{-1.5})$$

که step_num شماره مرحله آموزش است. Warmup_steps = 4000

مرحله ی 2

ساختار کلی شبکه به این صورت می‌باشد که برای سادگی کار در شبکه transformer قسمت decoder را در نظر نمی‌گیریم. فرض کنیم داده‌های ورودی (صوتی) ابتدا با کد داده شده تبدیل به spectrogram میشوند و بعد از آن وارد یک شبکه برای embedding میشوند و پس از ترکیب با قسمت positional encoding وارد یک شبکه encoder با ساختار توضیح داده شده در بالا می‌شود و با عبور از لایه encoder به همراه خروجی متنی وارد تابع هزینه etc می‌شود و مقدار هزینه محاسبه می‌شود.

الف:

سه معیار WER و CER و Loss را در طول آموزش محاسبه و در یک نمودار برای داده های آموزشی (train) و داده‌های اعتبارسنجی (validation) رسم کنید.

ب:

برای چند داده ورودی (۱۰ داده ورودی صوتی) دنباله خروجی متنی حاصل از دنباله ورودی تولید شده توسط مدل را در هر 50 گام محاسبه و نمایش دهید.
برای دیکد کردن دنباله از روش best path decoding استفاده کنید.

پارامترهای مدل را به صورت زیر در نظر بگیرید.

```
Sample_rate = 16000,
Window_size = 0.02
hidden_size = int(math.floor( (sample_rate * window_size) / 2) + 1)
hidden_size = int(math.floor(hidden_size - 41) / 2 + 1)
hidden_size = int(math.floor(hidden_size - 21) / 2 + 1)
dim_input = hidden_size * 32
num_heads= 8
dim_model=256
dim_key=64
dim-emb= 256
```

dim_value= 64
dim-inner= 1024
num-layers_decoder = 4
Learning rate= 1e-4
batch-size= 12

نکات کلی تمرین:

- در صورت مشاهده مشابهت کد بین هر دو دانشجو، نمره تمرین هر دو نفر صفر لحاظ خواهد شد.
- در صورت مشاهده هرگونه مشابهت کد با کدهای موجود در اینترنت، نمره تمرین صفر لحاظ خواهد شد. اگر بخشی از کد که از قسمت‌های اصلی تمرین نمی‌باشد را کدهای آماده استفاده می‌کنید حتما لینک و منبع آن را اعلام کنید.
- لازم به ذکر است نیمی از نمره تمرین مربوط به گزارش می‌باشد. بنابراین رعایت اصول نگارشی حائز اهمیت می‌باشد.
- در نوشتن گزارش، لحاظ جزئیات گزارش الزامی است مانند موارد زیر
 - ارجاع دادن به مطالب و اشکالی که از مقاله و وبسایت ها گرفته شده اند
 - توضیح اشکال و جداول در caption
 - نوشتن فرمول و قرار ندادن عکس‌های فرمول در متن
 - نوشتن نتایج شبیه‌سازی ها به صورت جدول و شکل
 - درست بودن متن از لحاظ قواعد دستور زبانی و نگارشی
 - موارد تکمیلی در فایل template آمده است.
- گزارش تمرین را به صورت فایل pdf و در کنار کدهای تمرین در سایت آپلود نمایید.
- نحوه نامگذاری به صورت studentnumber_Homeworknumber.pdf باشد.
- برای هرگونه پرسش پیرامون این تمرین را با ایمیل‌های aliparchekan@gmail.com, esmaeilfarhang@gmail.com مطرح نمایید.