



به نام خدا



دانشگاه تهران  
دانشکده مهندسی برق و کامپیوتر  
یادگیری عمیق با کاربردها

تمرین شماره چهار

نام و نام خانوادگی	امید نائیج نژاد
شماره دانشجویی	610301189
تاریخ ارسال گزارش	1404/3/20

## فهرست مطالب

5 .....	سوال 1
7 .....	سوال 2
11.....	سوال 3
14.....	مراجع

## فهرست شکل ها

- شکل 1 - زیربخش های هد چندگانه توجه.....5
- شکل 2 - زیربخش های هر هد توجه.....6
- شکل 3 - نمودار تابع هزینه حین آموزش.....13

## فهرست جدول ها

## سوال 1

ابعاد سیگنال قبل از ورود به هر بلاک و پس از خروج از آن در بخش رمزنگار به صورت زیر می باشد:

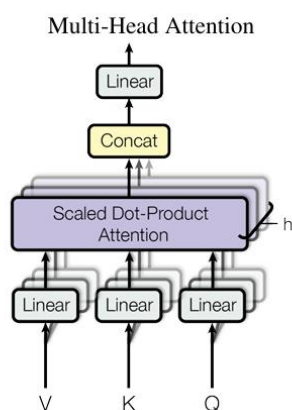
در ابتدا سیگنال داده ای به شکل یک ماتریس با ابعاد  $(32, 4096)$   $\rightarrow$  (Batch, Seq Length) است. در این مرحله هر توکن به صورت یک رشته فعلا ذخیره شده است.

پیش از ورود داده به ساختار مدل های ترنسفورمری، باید با کمک روش های امبدینگ ورودی (تعبیه)، نمایش وان-هات توکن ها را با یک تبدیل خطی به یک فضای معنادار از لحاظ فاصله بردارها منتقل کنیم. (مشابه روشی که در مدل های رمزنگار-رمزگشا مبتنی بر شبکه های عصبی بازگشتی برای کاربرد توصیف تصاویر در تمرین قبلی انجام دادیم). بنابراین ابعاد سیگنال به صورت  $(32, 4096, 1024)$  تبدیل خواهد شد. برای نمایش هر توکن از یک بردار به اندازه امبدینگ توکن ورودی یعنی 1024 استفاده کرده ایم.

حالا به بلاک های انکودری (رمزنگار) میرسیم.

- در بلوک هد چندگانه توجه، بردار نهان به تعداد هدها تقسیم میشود تا ورودی و خروجی هر هد به دست بیاید:

✓ هر هد توجه، مطابق با شکل 1، در ابتدا با استفاده از یک تبدیل خطی برای محاسبه  $Q, V, K$  ابعاد سیگنال را از  $(32, 4096, 1024)$  به  $(32, 4096, 768)$  تبدیل میکند و متناسب با تعداد هد ها آنها را جدا میکند:  $(32, 6, 4096, 128) \rightarrow 768 \div 6 = 128$



شکل 1 - زیربخش های هد چندگانه توجه

✓ حالا هر هد توجه به صورت جداگانه مطابق با شکل 2، و رابطه 1 خروجی محاسبه میشود:

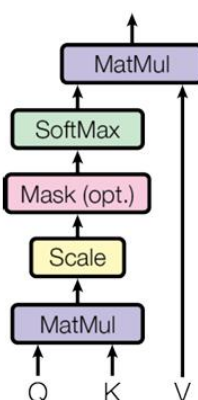
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

1.  $QK^T = (32, 6, 4096, 4096)$
2. تقسیم کردن ماتریس بر یک اسکالر و تابع سافت مکس واضحاً ابعاد را تغییر نمیدهند.
3. بعد از ضرب کردن ماتریس بدست آمده در بردار  $V$ ، مجدداً سیگنال به صورت  $(32, 6, 4096, 128)$  خواهد بود.

✓ حالا پس از پردازش های هر هد توجه باید خروجی های بدست آمده را به همدیگر بچسبانیم (شکل 1) پس سیگنال دوباره بصورت  $(32, 4096, 768)$  درخواهد آمد.

- ✓ در انتهای بلاک هد چندگانه توجه، تبدیل خطی خروجی انجام میشود تا سیگنال به حالت پیش از ورود به این بلاک یعنی (32, 4096, 1024) تبدیل شود.
- بدیها در بخش نرمال سازی لایه و ارتباط باقیمانده، ابعاد تغییر نمیکند.
- در لایه پرسرو بلاک رمزگشا:
- ✓ در لایه پرسرو اول، ابعاد ورودی نصف میشود: (32, 4096, 512)
- ✓ در لایه پرسرو دوم، ابعاد ورودی تغییر نمیکند: (32, 4096, 512)
- ✓ در لایه پرسرو سوم، ابعاد ورودی دوباربر میشود: (32, 4096, 1024)
- مجددا در بخش نرمال سازی لایه و ارتباط باقیمانده، ابعاد تغییر نمیکند.

### Scaled Dot-Product Attention



شکل 2 - زیربخش های هر هد توجه

### پارامترهای مدل:

- هر بلوک رمزنگار:  $3,149,056 + 2048 + 1,312,768 + 2048 = 4,465,920$
- هد چندگانه توجه:  $2,361,600 + 787,456 = 3,149,056$
- تبدیل خطی Q/K/V: هر کدام یک لایه پرسرو از 1024 به 128 هستند پس تعداد پارامترهای این بخش با توجه به اینکه 6 هد توجه داریم برابر است با:  $6 \times 3 \times (1024 \times 128 + 128) = 2,361,600$
- تبدیل خطی خروجی: یک لایه پرسرو از 768 به 1024 است پس تعداد پارامترهای این بخش برابر است با:  $768 \times 1024 + 1024 = 787,456$
- نرمال سازی لایه: 2 بردار قابل آموزش به اندازه نرون های آن لایه یعنی  $2 \times 1024 = 2048$
- لایه های پرسرو:  $524,800 + 262,656 + 525,312 = 1,312,768$
- در لایه پرسرو اول، ابعاد ورودی نصف میشود:  $1024 \times 512 + 512 = 524,800$
- در لایه پرسرو دوم، ابعاد ورودی تغییر نمیکند:  $512 \times 512 + 512 = 262,656$
- در لایه پرسرو سوم، ابعاد ورودی دوباربر میشود:  $512 \times 1024 + 1024 = 525,312$
- نرمال سازی لایه: 2 بردار قابل آموزش به اندازه نرون های آن لایه یعنی  $2 \times 1024 = 2048$

- هر بلوک رمزگشا:  $3,149,056 + 2048 + 3,149,056 + 2048 + 1,312,768 + 2048 = 7,617,192$
- هد چندگانه ماسک شده توجه:  $2,361,600 + 787,456 = 3,149,056$
- تبدیل خطی Q/K/V: هر کدام یک لایه پسرو از 1024 به 128 هستند پس تعداد پارامترهای این بخش با توجه به اینکه 6 هد توجه داریم برابر است با:  $6 \times 3 \times (1024 \times 128 + 128) = 2,361,600$
- تبدیل خطی خروجی: یک لایه پسرو از 768 به 1024 است پس تعداد پارامترهای این بخش برابر است با:  $768 \times 1024 + 1024 = 787,456$
- نرمالسازی لایه: 2 بردار قابل آموزش به اندازه نرون های آن لایه یعنی  $2 \times 1024 = 2048$
- توجه رمزنگار-رمزگشا: تعداد پارامترها کاملاً مشابه با هد چندگانه ماسک شده توجه است:  $2,361,600 + 787,456 = 3,149,056$
- نرمالسازی لایه: 2 بردار قابل آموزش به اندازه نرون های آن لایه یعنی  $2 \times 1024 = 2048$
- لایه های پسرو:  $524,800 + 262,656 + 525,312 = 1,312,768$
- در لایه پسرو اول، ابعاد ورودی نصف میشود:  $1024 \times 512 + 512 = 524,800$
- در لایه پسرو دوم، ابعاد ورودی تغییر نمیکند:  $512 \times 512 + 512 = 262,656$
- در لایه پسرو سوم، ابعاد ورودی دوبرابر میشود:  $512 \times 1024 + 1024 = 525,312$
- نرمالسازی لایه: 2 بردار قابل آموزش به اندازه نرون های آن لایه یعنی  $2 \times 1024 = 2048$

مدل ترنسفورمری ما شامل 6 بلاک انکودر و 8 بلاک دیکدر است پس تعداد کل پارامترهای مدل برابر است با:

$$4,465,920 \times 6 + 7,617,192 \times 8 = 26,795,520 + 60,936,192 = \underline{\underline{87,731,712}}$$

## سوال 2

مدل های زبانی بزرگ مانند GPT، LLaMA و Claude معمولاً در سه مرحله آموزش داده می شوند:

1. پیش آموزش (Pre-training)
2. تنظیم دقیق با نظارت (Supervised Fine-Tuning - SFT)
3. همراستاسازی با بازخورد انسانی (Alignment via Reinforcement Learning from Human Feedback - RLHF)

هر مرحله نقش خاصی در ساخت مدل های قابل استفاده و منطبق با ارزش های انسانی دارد.

### ۱. مرحله پیش آموزش

الف) این مرحله در چه داده هایی و با چه هدفی انجام می شود؟

در این مرحله، مدل با استفاده از داده های متنی خام وبدون برچسب مانند کتابها، مقالات ومحتوای وب، آموزش می بیند. هدف از این بخش، یادگیری الگوهای زبانی عمومی مانند دستور زبان، معناشناسی و دانش عمومی است.

- هدف: یادگیری ساختار زبان و دانش عمومی
- نوع داده: متون خام و بدون برچسب از منابع مختلف (مثل دیالوگ های سریال Friends در این تمرین)
- روش آموزش: با پیش بینی توکن بعدی دریک دنباله، مدل را آموزش می دهیم.

## ب) چرا مدل در این مرحله هنوز آمادگی استفاده مستقیم در کاربردهای انسانی نیست؟

مدل های پیش آموزش شده توانایی درک یا پیروی از دستورات خاص را ندارند. ممکن است پاسخ های نامناسب از نظر اخلاقی، مغرضانه یا نادرست تولید کنند. چون مدل ترجیحات انسانها را برای سوال(پرامپت)هایی که میپرسند نمیداند و ارزش منفی کلمات بی ادبانه را نمیفهمد؛ پس امکان تولید خروجی توهین آمیز و نامناسب وجود دارد. بنابراین، نیاز به مراحل بعدی برای تنظیم دقیق و هم راستاسازی دارد.

## ۲. مرحله آموزش با نظارت

### الف) تفاوت این مرحله با پیش آموزش چیست؟

همانطور که گفته شد در مرحله پیش آموزش (Pre-training) :

- مدل زبانی برای اولین بار آموزش می بیند و با ساختار جملات آشنا میشود.
- در این مرحله، از داده های بدون برچسب استفاده می شود (مثل تمام محتوای اینترنت یا کتابها).
- هدف یادگیری ساختار زبان، روابط بین کلمات و جملات و بدست آوردن یک دانش عمومی و کلی است.
- نوع یادگیری اغلب به صورت self-supervised است، مثلاً مدل سعی می کند کلمات حذف شده را حدس بزند یا جمله بعدی را پیش بینی کند.

اما در مرحله آموزش با نظارت (Supervised Fine-Tuning) که بعد از پیش آموزش انجام می شود:

- از داده های برچسب خورده استفاده می شود؛ یعنی برای هر ورودی، پاسخ درست توسط انسان مشخص شده.
- هدف از این مرحله یادگیری نحوه انجام تسک های متفاوت (مثل پاسخ به سوال، ترجمه، طبقه بندی، چت کردن، خلاصه سازی) است.
- این مرحله باعث می شود مدل قابلیت تعامل واقعی با انسانها را به دست بیاورد.



ب) نوع داده‌هایی که در این مرحله استفاده میشوند را شرح دهید و یک مثال بزنید؟

- داده‌ها باید شامل جفت‌های (ورودی، خروجی) باشند که توسط انسان تهیه شده‌اند.
- این داده‌ها اغلب در قالب مجموعه‌هایی از سوال و جواب، دستور و پاسخ، جمله و ترجمه، یا متن و برچسب ارائه می‌شوند. پس در این مرحله به این labelled data یا داده‌های برچسب‌خورده نیاز داریم. واضح است که جمع‌آوری داده مورد نیاز این مرحله بسیار هزینه برتر از مرحله اول است. در مرحله قبلی با نوشتن یک web crawler به راحتی کل متن‌های موجود در وب سایت‌ها را میتوانیم دانلود کنیم اما در این مرحله برای آماده‌سازی دیتا به انسان نیاز داریم.
- مثال‌هایی از داده‌های مورد نیاز برای تسک‌های مختلفی که از یک مدل زبانی انتظار میرود بتواند انجام دهد:

1. ترجمه:

○ ورودی: «Translate 'I love you' to Spanish»

○ خروجی: «Te quiero»

2. پاسخ به سوال:

○ ورودی :

«What is the boiling point of water?»

○ خروجی:

«100 °C at standard atmospheric pressure»

3. طبقه‌بندی احساسات:

○ ورودی: «I had a terrible day.»

○ خروجی: «Negative»

4. خلاصه‌سازی متن:

○ ورودی: یک مقاله یا متن بلند

○ خروجی: یک جمله‌ی خلاصه‌شده شامل مطالب مهم آن

ج) چرا این مرحله نقش مهمی در تعامل انسانی-مدل دارد؟

✓ انسان‌محور شدن مدل:

- مدل بدون آموزش با نظارت، فقط یک زبان‌فهم است، اما نمی‌داند چطور در تعامل واقعی با انسان‌ها پاسخ دقیق و سودمند بدهد.
- آموزش با نظارت به مدل کمک می‌کند تا سبک پاسخ‌دهی مناسب، دقیق و اخلاقی را بیاموزد.

✓ کاهش پاسخ‌های نادرست یا بی‌ربط:

- مدل‌هایی که فقط با پیش‌آموزش تربیت شده‌اند، ممکن است پاسخ‌هایی بدهند که گمراه‌کننده، بی‌معنا یا حتی خطرناک باشد.
- داده‌های نظارت‌شده مثل راهنماهایی هستند که به مدل نشان می‌دهند: «چه چیزی قابل قبول است و چه چیزی نیست.»

✓ افزایش دقت و قابلیت اطمینان:

- در کاربردهای حساس مثل پزشکی، حقوق، یا آموزش، مدل باید دقیق باشد و علاوه بر یادگیری ساختار زبان، باید با اصطلاحات تخصصی و بار حقوقی بعضی کلمات نیز آشنا شود.
- آموزش با نظارت باعث می‌شود مدل به خروجی‌های استاندارد و درست نزدیک‌تر شود.

✓ پایه‌ای برای آموزش مبتنی بر بازخورد انسانی (RLHF):

- این مرحله پیش‌نیاز مرحله بعدی است که در آن مدل از طریق بازخورد انسانی بهینه‌تر می‌شود. اگر این مرحله انجام نشود، RLHF نمی‌تواند مؤثر باشد.

### ۳. هم راستاسازی

#### الف) هدف اصلی از این مرحله چیست؟

هدف هم راستاسازی (Alignment) این است که مدل زبانی را طوری تنظیم کنیم که با ارزش‌ها و انتظارات انسانی هماهنگ شود.

یعنی حتی اگر مدل از نظر فنی درست کار کند و جواب درستی از نظر علمی تولید کند، باید یاد بگیرد که چگونه پاسخ‌هایی مفید، محترمانه، بی‌طرف و ایمن ارائه دهد تا انسان‌ها بتوانند به راحتی با آن کار کنند و برای کودکان بدآموزی نداشته باشد.

#### ب) یکی از روش‌های رایج برای انجام هم راستاسازی را نام ببرید و نحوه‌ی کار آن را به طور خلاصه توضیح دهید.

یکی از رایج‌ترین روش‌ها یادگیری تقویتی با بازخورد انسانی (RLHF = Reinforcement Learning from Human Feedback) نام دارد. نحوه‌ی کار این روش به صورت زیر است:

1. ابتدا مدل روی یک مجموعه داده‌ی با نظارت آموزش می‌بیند (مرحله قبلی).
2. سپس انسان‌ها چند خروجی مختلف مدل را برای یک ورودی مشخص رتبه‌بندی می‌کنند (مثلاً پاسخ بهتر، محترمانه‌تر یا دقیق‌تر را انتخاب می‌کنند).
3. با استفاده از این رتبه‌بندی‌ها، یک مدل پاداش (reward model) آموزش داده می‌شود.
4. در نهایت، با استفاده از الگوریتم‌های یادگیری تقویتی (مثل PPO)، مدل زبانی به گونه‌ای بهینه می‌شود که پاسخ‌هایی تولید کند که مطابق با پاداش انسانی باشند.

### سوال 3

#### دانلود و آماده سازی داده

از لینک داده شده، دیتاست را دانلود کرده و با استفاده از ستون text آن، corpus را میسازیم سپس آن را به دنباله هایی با طول یکسان و برابر با *block\_size* تقسیم میکنیم تا برای مدل ترنسفورمری ما قابل پردازش باشد.

#### تنظیم هایپر پارامترها

کاملاً مشابه با (جدول ۲: ابرپارامترهای پیشنهادی برای آموزش مدل) موجود در فایل توضیح تمرین، تنظیمات انجام شد.

#### پیاده سازی مائول های مدل

در ابتدا در سطح کاراکتر توکنایزیشن را انجام میدهیم. در ادامه مقایسه ای بین روشهای مختلف توکنایزیشن را ارائه میدهیم:

##### 1. توکن سازی در سطح کاراکتر

###### • مزایا:

- تعداد واژگان بسیار کم (در حد مثلاً 100-200 کاراکتر)
- عدم وجود مشکل در مدیریت کلمات ناشناخته (Out of Vocab)

###### • معایب:

- دارای سرعت پردازش و تولید بسیار کمتر نسبت به حالت های بعدی
- عدم درک ساختارهای معنایی سطح بالا مانند کلمات، عبارات و قواعد نحوی.
- مدل باید وابستگی های طولانی مدت را در سطح کاراکتر یاد بگیرد، که بدیهتاً دشوارتر از یادگیری وابستگی های طولانی مدت در سطح کلمه است.

##### 2. توکن سازی در سطح کلمه

###### • مزایا

- توالی های کوتاه تر و پردازش سریعتر
- درک بهتر از ساختار معنایی جملات

###### • معایب

- واژگان بسیار بزرگ (صدها هزار کلمه) میشود که برای مدل های چند زبانه مشکل ساز خواهد بود.
- وجود مشکل در مدیریت کلمات ناشناخته (Out of Vocab) برعکس روش در سطح کاراکتر

##### 3. توکن سازی در سطح زیرکلمه (Subword-Level)

###### • مزایا

- ترکیبی از مزایای دو روش بالا را دارد.

- واژگان با اندازه متوسط (مثلاً 30-50 هزار توکن)
- توانایی مدیریت کلمات خارج از واژگان با تقسیم آن ها به زیرکلمات

#### • معایب

- پیاده سازی پیچیده تر نسبت به توکن سازی در سطح کاراکتر

بنابراین:

- ✓ برای مدل های کوچک یا پروژه های آموزشی، توکن سازی در سطح کاراکتر مناسب است.
- ✓ برای مدل های بزرگ و کاربردهای واقعی، توکن سازی در سطح زیرکلمه (مانند BPE یا SentencePiece) ترجیح داده می شود، زیرا تعادل خوبی بین اندازه واژگان و توانایی مدل در درک ساختار زبان فراهم می کند.

### آموزش مدل

تابع `train_model` پیاده سازی شده است و با ورودی گرفتن `optimizer, scheduler, tokenizer, max_iter, eval_interval, batch_size` مدل را آموزش می دهد. باید توجه کرد که برای آموزش مدل های زبانی بزرگ برخلاف شبکه های عصبی کانولوشنی و بازگشتی که در هر اپیک روی کل داده آموزشی گرادیان تابع هزینه را محاسبه می کنیم و با استفاده از الگوریتم پس انتشار خطا وزن ها شبکه را آپدیت می کنیم، اینجا در هر اپیک یک نمونه تصادفی از کل داده آموزشی را انتخاب می کنیم و گرادیان تابع هزینه را صرفاً روی یک بسته محاسبه می کنیم. چون در مدل های زبانی با یک دنباله سر و کار داریم که میتواند هر چیزی باشد (برخلاف تسک های بینایی ماشین و پردازش زبان طبیعی که مثلاً صرفاً یک طبقه بندی ساده است) و با این نمونه گیری تصادفی به نوعی کوثری کاربر را داریم شبیه سازی می کنیم. به طور مشابه نحوه عملکرد مدل روی داده های ارزیابی نیز با همین منطق روی صرفاً یک بسته اندازه گیری و گزارش شده است.

برای آموزش مدل های زبانی، بهینه ساز `AdamW` و `CosineAnnealingLR` برای `Learning Rate Scheduler` پیشنهاد شده است. در زیر به طور خلاصه توضیحاتی درباره این دو روش نوشته ام:

**AdamW** یک الگوریتم بهینه سازی است که از بهینه ساز `Adam` مشتق شده است. تفاوت اصلی بین `Adam` و `AdamW` در نحوه مدیریت کاهش وزن (تنظیم کننده `L2`) است.

- در بهینه ساز اصلی `Adam`، تنظیم کننده `L2` (کاهش وزن) با افزودن یک عبارت متناسب با وزن ها مستقیماً به گرادیان ها اعمال می شود. این کار به طور مؤثری کاهش وزن را با نرخ های یادگیری تطبیقی ترکیب می کند که منجر به عملکرد نامطلوب یا دشوارتر شدن تنظیم مؤثر کاهش وزن می شود. اعمال کاهش وزن مستقیماً به گرادیان ها به این معنی است که این مقیاس بندی بر عبارت کاهش وزن نیز تأثیر می گذارد، که ممکن است مطلوب نباشد.
- `AdamW` که توسط لوشچیلوف و هوتر معرفی شد، کاهش وزن را از فرآیند به روزرسانی گرادیان جدا می کند و به جای افزودن عبارت کاهش وزن به گرادیان ها، آن را مستقیماً پس از به روزرسانی گرادیان به وزن ها اعمال می کند. این جداسازی تضمین می کند که نرخ های یادگیری تطبیقی `Adam` با کاهش وزن تداخل پیدا نمی کنند و امکان کنترل دقیق تر بر تنظیم کننده را فراهم می کند.

چرا AdamW برای آموزش LLM ها مفید است:

- بهبود تعمیم‌پذیری LLM : الگوریتم AdamW با جدا کردن کاهش وزن، یک مکانیزم تنظیم‌کننده مؤثرتر را فراهم می‌کند. این به جلوگیری از حفظ داده‌های آموزشی توسط مدل کمک می‌کند و در عوض آن را به یادگیری ویژگی‌های تعمیم‌پذیرتر ترغیب می‌کند که منجر به عملکرد بهتر در داده‌های دیده نشده می‌شود.
- عملکرد بهتر: شواهد تجربی در جامعه توسعه دهندگان LLM نشان می‌دهد که AdamW اغلب منجر به همگرایی برتر و عملکرد نهایی بهتر در مقایسه با Adam اصلی می‌شود، به‌ویژه برای معماری‌های مبتنی بر ترانسفورمر در مقیاس بزرگ.

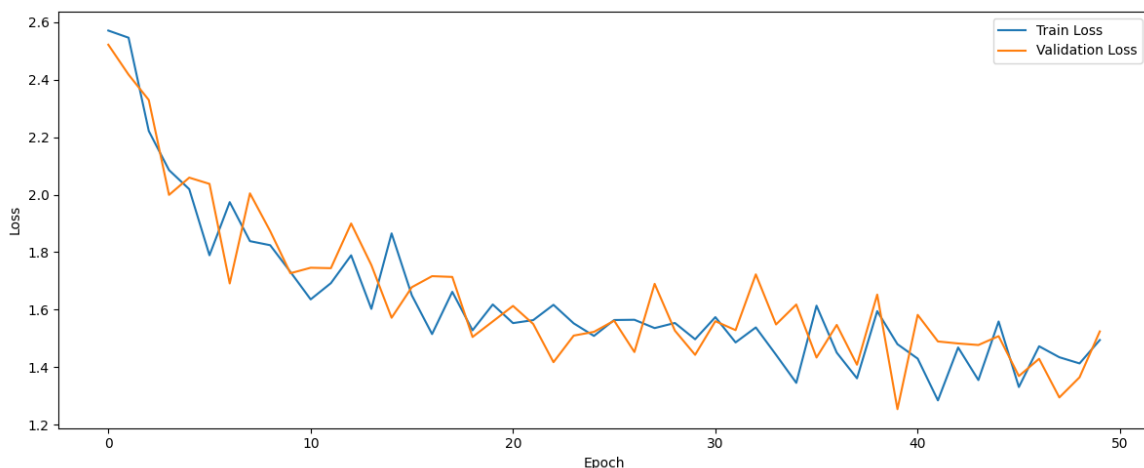
CosineAnnealingLR یک برنامه‌ریز نرخ یادگیری است که نرخ یادگیری را در طول آموزش به دنبال یک منحنی کسینوسی تنظیم می‌کند:

$$\eta_t = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min}) \left(1 + \cos\left(\frac{T_{Cur}}{T_{max}}\pi\right)\right)$$

که در آن:

- $\eta_{min}$  حداقل نرخ یادگیری است.
- $\eta_{max}$  نرخ یادگیری اولیه (حداکثر) است. در این تمرین 0.001.
- $T_{Cur}$  تعداد تکرارها/ایوک‌ها از ابتدای آموزش است.
- $T_{max}$  تعداد کل تکرارها/ایوک‌ها برای چرخه کاهنده کسینوسی است. در این تمرین 5000.

مقدار اولیه بالای نرخ یادگیری، باعث میشود تا مدل در ابتدا با سرعت خوبی آموزش ببیند و در ادامه با کاهش تدریجی نرخ یادگیری به همگرا شدن مدل کمک میکند.



شکل 3- نمودار تابع هزینه حین آموزش

تولید متن

در این مثال، خروجی مدل را با استفاده از (top\_k decoding (top\_k = 5 و temperature = 0.8 میبینید:

ورودی اولیه:

“Joey: “

خروجی مدل:

“their fored.

Yeah, three and star”

در این مثال، خروجی مدل را با استفاده از (Beam Search Decoding (beam\_width=3 و temperature = 1 میبینید:

ورودی اولیه:

“Joey: “

خروجی مدل:

“their fored.

Yeah, three and star”

باتوجه به نحوه پیاده سازی تابع generate\_text که لاجیت ها را بر مقدار temperature تقسیم میکند:

- اگر مقدار temperature کمتر از یک باشد، خروجی مدل پیش بینی پذیرتر و دقیقتر و بدون خلاقیت خواهد بود.
- اگر مقدار temperature بیشتر از یک باشد، مدل خلاق تر و متنوع تر عمل خواهد کرد.

## مراجع

1. A. Vaswani et al., “Attention is all you need,” Adv. Neural Inf. Process. Syst., vol. 30, 2017.
2. [\[2504.12501v1\] Reinforcement Learning from Human Feedback](#)
3. [Decoupled Weight Decay Regularization](#)