



به نام خدا
دانشگاه تهران
دانشکده مهندسی
برق و کامپیوتر



درس شبکه‌های عصبی و یادگیری عمیق
تمرین امتیازی

نام و نام خانوادگی	بهراد موسایی شیرمحمد محمد جواد رنجبر کلهرودی
شماره دانشجویی	۸۱۰۱۰۱۱۷۳ ۸۱۰۱۰۱۲۷۸
تاریخ ارسال گزارش	۱۴۰۲.۱۰.۲۷

۵	LoRA ۲
۵	۲-۱. LoRA چگونه کار می کند؟
۵	فاین تون کردن کل پارامترها
۵	فاین تون کردن یک یا چند لایه از مدل
۶	Fine tune with LoRA
۸	۲۲. قرار است چه کاری بر روی داده ها صورت پذیرد؟
۹	توضیحات دیتاست
۱۱	۳۲. کد نویسی: آموزش مدل
۱۵	نتایج Roberta
۱۹	نتایج Freeze
۲۲	مقایسه کلی
۲۴	۴-۲. چرا LoRA ؟
۲۶	۳. تشخیص تقلب
۲۶	۳-۱. آشنایی با دیتاست
۲۸	۳-۲. پیاده سازی معماری مقاله
۳۲	۳-۳. نمونه برداری
۳۸	۳-۴. آموزش مدل با داده های جدید
۳۸	طبق پیاده سازی مقاله

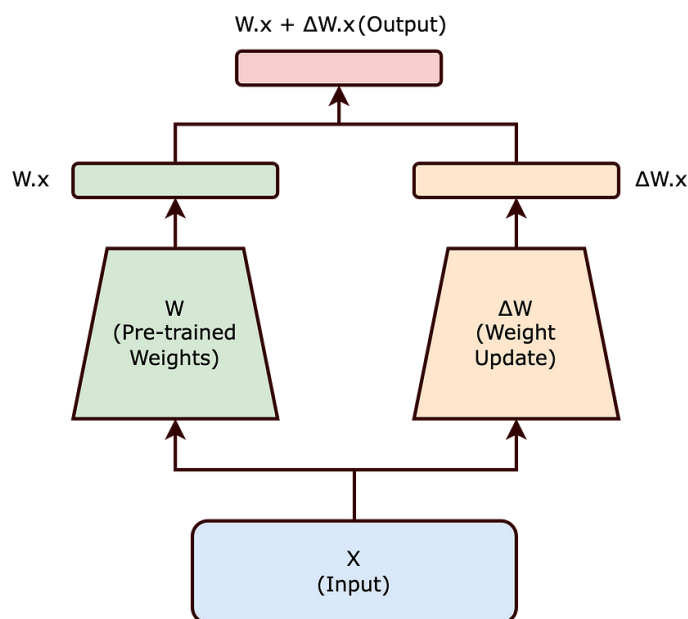
شکل ۱ روش سنتی تنظیم دقیق	۵
شکل ۲ نحوه بروزرسانی وزن‌های سنتی	۶
شکل ۳ نحوه تنظیم دقیق با LoRa	۷
شکل ۴ نحوه بروزرسانی وزن‌ها با LoRa	۷
شکل ۵ ضریب مقیاس در LoRa	۸
شکل ۶: روش بارگذاری داده‌ها	۹
شکل ۷ نمونه‌هایی از مجموعه داده MultiNLI	۹
شکل ۸: نمایی از دیتاست	۱۲
شکل ۹: مدل زبانی Roberta	۱۳
شکل ۱۰: پیاده سازی کلاس Callback	۱۵
شکل ۱۱: نتایج Roberta	۱۶
شکل ۱۲: ماتریس در هم آمیختگی Roberta	۱۶
شکل ۱۳: مدل Roberta	۱۸
شکل ۱۴: پیاده سازی کلاس Callback	۱۹
شکل ۱۵: نتایج Freeze	۱۹
شکل ۱۶: ماتریس در هم آمیختگی Freeze	۲۰
شکل ۱۷ تعداد پارامترهای قابل آموزش مدل قبل و بعد از LoRa	۲۰
شکل ۱۸: نتایج LoRa	۲۱
شکل ۱۹: ماتریس در هم آمیختگی LoRa	۲۱
شکل ۲۰ هیستوگرام داده‌ها	۲۷
شکل ۲۱ معماری مدل	۳۰
شکل ۲۲ هاپر پارامترهای مدل	۳۰
شکل ۲۳: نتایج پیاده سازی مدل	۳۱
شکل ۲۴: ماتریس در هم آمیختگی مدل	۳۱
شکل ۲۵: نتایج مدل	۳۲
شکل ۲۶ معیارهایی صحت و recall	۳۲
شکل ۲۷ تعداد کل داده‌های اقلیت	۳۴
شکل ۲۸ محاسبه ۲	۳۴
شکل ۲۹ محاسبه شعاع	۳۵

- شکل ۳۰: تعداد نمونه‌ی تصادفی ۳۵
- شکل ۳۱: تولید نمونه تصادفی ۳۶
- شکل ۳۲: هیستوگرام بعد از نمونه برداری ۳۸
- شکل ۳۳: نتایج با نمونه برداری ۴۰
- شکل ۳۴: ماتریس در هم ریختگی با نمونه برداری ۴۰
- شکل ۳۵: نتایج مربوط به نمونه برداری ۴۱
- شکل ۳۶: معیارهای صحت و recall با نمونه برداری ۴۱

۲-۱. LoRA چگونه کار می کند:

فاین تون کردن کل پارامترها

در فاین تون کردن کامل، مدل با وزن های پیش آموزش داده شده Φ_0 مقداردهی اولیه می شود و سپس به $\Phi_0 + \Delta\Phi$ به روزرسانی می شود. این به روزرسانی با دنبال کردن گرادیان برای حداکثر کردن هدف مدل سازی زبان شرطی انجام می گیرد. برای هر وظیفه جدید، مجموعه ای متفاوت از پارامترها $\Delta\Phi$ آموخته می شود که ابعاد آن برابر با ابعاد وزن های پیش آموزش داده شده $|\Phi_0|$ است. این می تواند یکی از معایب اصلی فاین تون کردن کامل باشد چرا که اگر مدل پیش آموزش داده شده بزرگ باشد (مانند GPT۳ با $|\Phi_0|$ تقریباً ۱۷۵ میلیارد)، ذخیره سازی و استقرار نمونه های مستقل و متعدد از مدل های فاین تون شده می تواند چالش برانگیز باشد، اگر اصلاً ممکن باشد.



شکل ۱ روش سنتی تنظیم دقیق

فاین تون کردن یک یا چند لایه از مدل

فاین تون کردن تنها برخی از لایه های مدل یک نوع ساده تر از فاین تون کردن کامل است که در آن تنها برخی از لایه ها به روزرسانی می شوند در حالی که بقیه لایه ها بدون تغییر باقی می مانند. به عنوان مثال، در یک مطالعه قبلی (Li & Liang ۲۰۲۱) درباره GPT۲، فقط دو لایه آخر مدل (FTT_{op2}) فاین تون شدند. این رویکرد اجازه می دهد تا با کاهش تعداد پارامترهایی که نیاز به به روزرسانی دارند، نیاز به ذخیره سازی و پردازش کمتری داشته باشیم. به عنوان مثال، ذخیره سازی ۱۰۰ مدل تطبیق داده شده با

این رویکرد فقط نیاز به ۳۵۴ گیگابایت فضا دارد، در حالی که با فاین تون کردن کامل برای هر مدل نیاز به ۳۵۰ گیگابایت فضا خواهیم داشت، که مجموعاً ۳۵ ترابایت فضا نیاز دارد.

Fine tune with LoRA

حال به صورت کامل داریم که LoRA , روش (LoRA(LowRank Adaptation) یک استراتژی فاین تون است که وزن های پیش آموزش داده شده ی مدل را ثابت نگه می دارد و ماتریس های قابل آموزش تجزیه رتبه پایین را به هر لایه از معماری Transformer تزریق می کند. این کار به طور چشمگیری تعداد پارامترهای قابل آموزش برای وظایف پایین دست را کاهش می دهد. به عنوان مثال، در مقایسه با GPT۳ که با آدام فاین تون شده، LoRA تعداد پارامترهای قابل آموزش را تا ۱۰۰۰۰ برابر و نیاز به حافظه GPU را تا سه برابر کاهش می دهد.

وقتی یک مدل زبان را به دقت تنظیم می کنیم، پارامترهای زیربنایی مدل را اصلاح می کنیم. برای ملموس تر کردن این ایده، می توانیم به روزرسانی پارامتر حاصل از تنظیم دقیق را همانطور که در معادله زیر نشان داده شده است، فرموله کنیم.

$$\overbrace{W_{ft}}^{\text{Finetuned Weights}} = \underbrace{W_{pt}}_{\text{Pretrained Weights}} + \overbrace{\Delta W}^{\text{Weight Update}}$$

شکل ۲ نحوه به روزرسانی وزن های سنتی

ایده اصلی پشت LoRA مدل سازی این به روز رسانی به پارامترهای مدل با تجزیه با رتبه پایین ۹ است که در عمل به عنوان یک جفت پیش بینی خطی اجرا می شود. LoRA لایه های از پیش آموزش دیده LLM را ثابت می گذارد و یک ماتریس تجزیه رتبه قابل آموزش را به هر لایه از مدل تزریق می کند. زیر را ببینید.

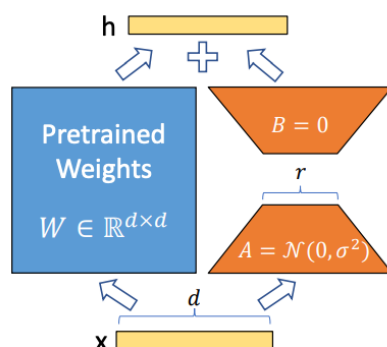


Figure 1: Our reparametrization. We only train A and B .

شکل ۳ نحوه تنظیم دقیق با LoRa

حال اجزای LoRa را توضیح می‌دهیم:

ماتریس تجزیه رتبه به زبان ساده، ماتریس تجزیه رتبه فقط دو طرح خطی است که ابعاد ورودی را کاهش داده و بازیابی می‌کند. خروجی این دو پیش بینی خطی به خروجی حاصل از وزن های از پیش آموزش دیده مدل اضافه می‌شود. لایه به روز شده ای که با افزودن این دو تبدیل موازی تشکیل می‌شود، مطابق شکل زیر فرموله شده است. همانطور که می‌بینیم، افزودن LoRa به یک لایه مستقیماً به روز رسانی وزن لایه زیرین را یاد می‌گیرد.

$$W_{ft} = W_{pt} + \underbrace{\Delta W}_{\text{Approximation}} = W_{pt} + \underbrace{AB}_{\text{Rank Decomposition Matrix}}$$

where $W_{ft}, W_{pt}, \Delta W, AB \in \mathbb{R}^{d \times d}$
 and $\underbrace{A \in \mathbb{R}^{d \times r}, B \in \mathbb{R}^{r \times d}}_{\text{Low Rank}}$

شکل ۴ نحوه بروزرسانی وزن ها با LoRa

محصول ماتریسی AB همان ابعاد به‌روزرسانی تنظیم دقیق را دارد. تجزیه به روز رسانی به عنوان محصولی از دو ماتریس کوچکتر تضمین می‌کند که به روز رسانی رتبه پایینی دارد و تعداد پارامترهایی را که باید آموزش دهیم به میزان قابل توجهی کاهش می‌دهد. به جای تنظیم دقیق پارامترها در لایه های LLM از پیش آموزش دیده، LoRa فقط ماتریس تجزیه رتبه را بهینه می‌کند و نتیجه ای به دست می‌دهد که به روز رسانی حاصل از تنظیم دقیق کامل را تقریبی می‌کند. ما A را با مقادیر تصادفی و کوچک

مقداردهی اولیه می‌کنیم، در حالی که B به صورت صفر مقداردهی می‌شود، و اطمینان حاصل می‌کنیم که فرآیند تنظیم دقیق را با وزن‌های اصلی و از پیش آموزش‌دیده‌شده مدل آغاز می‌کنیم.

افزایش r تقریب LORA را در مورد به‌روزرسانی کامل تنظیم دقیق بهبود می‌بخشد، اما مقادیر فوق‌العاده کوچک r در عمل کافی است و به ما این امکان را می‌دهد که هزینه‌های محاسباتی و حافظه را با کمترین تأثیر بر عملکرد به میزان قابل توجهی کاهش دهیم. به عنوان مثال، ما می‌توانیم از LORA برای تنظیم دقیق GPT-3 با استفاده از تنها ۰.۰۱٪ از کل پارامترها استفاده کنیم و همچنان عملکردی قابل مقایسه با تنظیم دقیق کامل داشته باشیم.

Scaling Factor

$$W_{ft} = W_{pt} + \frac{\overset{\alpha}{\underbrace{AB}}}{\underbrace{r}_{\text{Rank Decomposition Matrix}}}$$

شکل ۵ ضریب مقیاس در LoRA

ضریب پوسته‌پوسته شدن هنگامی که به‌روزرسانی رتبه پایین ماتریس وزن به دست آمد، قبل از اضافه کردن آن به وزن‌های از پیش آموزش‌دیده‌شده مدل، آن را با ضریب α مقیاس می‌کنیم. چنین رویکردی قانون انطباق نشان داده شده در بالا را به همراه دارد. مقدار پیش‌فرض ضریب مقیاس‌بندی یک است، به این معنی که وزن‌های از پیش تمرین‌شده و به‌روزرسانی وزن کم‌تر هنگام محاسبه پاس رو به جلو مدل، به طور مساوی وزن می‌شوند. با این حال، مقدار α را می‌توان تغییر داد تا اهمیت مدل از پیش آموزش‌دیده و انطباق ویژه کار جدید را متعادل کند. تجزیه و تحلیل تجربی اخیر نشان می‌دهد که مقادیر بزرگ‌تر α برای LoRA با رتبه بالاتر (یعنی $\alpha \propto r$) ضروری است.

۲۲. قرار است چه کاری بر روی داده‌ها صورت پذیرد؟

در این سوال قرار است از مدل زبانی بزرگ ROBERT نسخه Large استفاده کنید. این مدل را می‌توان به صورت آماده و با دستور زیر دریافت کرد. می‌توانید از هر روش دیگری نیز استفاده کنید ولی در این صورت از سالم بودن مدل دریافت شده اطمینان حاصل نمایید.


```
from transformers import AutoTokenizer, AutoModel

tokenizer = AutoTokenizer.from_pretrained("roberta-large")
model = AutoModel.from_pretrained("roberta-large")
```

شکل ۶: روش بارگذاری داده ها

با توجه به اطلاعاتی که داده شده، ابتدا به توضیحات کلی در مورد دیتاست پرداخته و سپس روش آموزش مدل با استفاده از LORA را شرح می‌دهیم:

توضیحات دیتاست

مجموعه داده چند سبک تداول زبان طبیعی (MultiNLI) یک مجموعه از ۴۳۳ هزار جفت جمله است که با اطلاعات تضمین متنی ارزیابی شده‌اند و به صورت گروهی به صورت برون‌سپاری جمع‌آوری شده‌اند. این مجموعه داده بر پایه مجموعه داده SNLI مدل شده است، اما با این تفاوت که گستره‌ای از سبک‌های متن‌های گفتاری و نوشتاری را پوشش می‌دهد و ارزیابی تعمیم از سبک متفاوت را پشتیبانی می‌کند. این مجموعه داده به عنوان پایه‌ای برای وظیفه مشترک کارگاه RepEval در کنفرانس ۲۰۱۷ EMNLP در کپنهاگ خدمت کرده است.

Premise	Label	Hypothesis
Fiction		
The Old One always comforted Ca'daan, except today.	<i>neutral</i>	Ca'daan knew the Old One very well.
Letters		
Your gift is appreciated by each and every student who will benefit from your generosity.	<i>neutral</i>	Hundreds of students will benefit from your generosity.
Telephone Speech		
yes now you know if if everybody like in August when everybody's on vacation or something we can dress a little more casual or	<i>contradiction</i>	August is a black out month for vacations in the company.
9/11 Report		
At the other end of Pennsylvania Avenue, people began to line up for a White House tour.	<i>entailment</i>	People formed a line at the end of Pennsylvania Avenue.

شکل ۷ نمونه‌هایی از مجموعه داده MultiNLI

مجموعه داده شامل جفت جملات از ده ژانر مختلف، مانند داستان، دولتی و مکالمات تلفنی است. هر جفت با یکی از سه دسته برچسب گذاری شده است: "مضاف"، "تضاد" یا "خنثی" که نشان دهنده رابطه بین دو جمله است.

MultiNLI به عنوان معیاری برای ارزیابی عملکرد مدل های مختلف NLP به کار گرفته شده است و نقش مهمی در پیشرفت تحقیقات در زمینه درک زبان طبیعی ایفا کرده است.

محققان و متخصصان اغلب از مجموعه داده MultiNLI برای آموزش و ارزیابی مدل ها استفاده می کنند و این به یک معیار استاندارد در جامعه NLP تبدیل شده است.

حال برای آموزش مدل باید مراحل زیر را طی کنیم:

مجموعه داده MultiNLI را بارگیری کنید:

مجموعه داده MultiNLI را که معمولاً در کتابخانه Hugging Face Datasets موجود است، بدست آورید.

توکن سازی:

از توکنایزر RoBERTa برای تبدیل داده های متن خام به فرمتی مناسب برای ورودی مدل استفاده کنید.

پیکربندی مدل:

مدل RoBERTa از قبل آموزش دیده را برای طبقه بندی توالی بارگذاری کنید و آن را برای کار خاص خود پیکربندی کنید.

تنظیم دقیق:

راه اندازی خط لوله تنظیم دقیق، از جمله تعریف پارامترهای آموزشی، بهینه ساز و زمان بندی نرخ یادگیری.

حلقه آموزشی:

حلقه آموزشی را پیاده سازی کنید، از طریق مجموعه داده آموزشی تکرار کنید، تلفات را محاسبه کنید، و وزن مدل را با استفاده از پس انتشار به روز کنید.

اعتبار سنجی:

به طور دوره ای مدل را بر روی مجموعه اعتبارسنجی ارزیابی کنید تا عملکرد آن در طول آموزش نظارت شود.

تنظیم فرایارامتر:

فرایارامترهایی مانند نرخ یادگیری، اندازه دسته ای و تعداد دوره های آموزشی را بر اساس عملکرد مدل در مجموعه اعتبار سنجی تنظیم کنید.

ارزیابی در مجموعه تست:

پس از آموزش، عملکرد تعمیم مدل را در مجموعه آزمون ارزیابی کنید.

ورودی مدل RoBERTa در مسائل پردازش زبان طبیعی (NLP) به شکل دنباله‌ای از توکن‌ها (واژه‌ها یا زیرکلمات) می‌باشد. این دنباله توکن‌ها به وسیله‌ی توکنایزر مربوط به مدل تولید می‌شود. برای مثال، اگر شما یک جمله دارید، این جمله توسط توکنایزر RoBERTa به یک دنباله از شناسه‌های عددی (توکن‌ها) تبدیل می‌شود که به مدل وارد می‌شود.

خروجی مدل RoBERTa در وظیفه‌های متفاوت می‌تواند متفاوت باشد. در وظیفه تشخیص دسته‌ای (مانند طبقه‌بندی دسته‌ها)، خروجی احتمالات اختصاص داده شده به هر دسته است. برای وظایف دیگر مانند ترجمه یا تولید متن، خروجی ممکن است شامل متن تولیدی باشد.

اگر مدل RoBERTa را به عنوان یک طبقه‌بندی‌کننده در MultiNLI داشته باشید، از آنجا که ما دارای ۳ کلاس هستیم نیاز به اضافه کردن یک لایه با سه نورون به انتهای مدل می‌باشیم.

برای اضافه کردن LoRA، به این صورت عمل می‌کنیم که LoRA لایه‌های از پیش آموزش دیده LLM را ثابت می‌گذارد و یک ماتریس تجزیه رتبه قابل آموزش را به هر لایه از مدل تزریق می‌کند.

۳۲. کد نویسی: آموزش مدل

حال شروع به پیاده سازی کد می‌کنیم:

در این بخش که با استفاده از کتابخانه Hugging Face Transformers و سایر بسته‌های مرتبط، یک محیط یادگیری ماشین را برای وظایف پردازش زبان طبیعی تنظیم و پیکربندی می‌کند. ابتدا کتابخانه‌های مختلف پایتون مانند PyTorch، Hugging Face Transformers و مجموعه داده‌ها را وارد می‌کند که برای توسعه و آموزش مدل‌های یادگیری ماشین ضروری هستند. بررسی می‌کند که آیا یک GPU CUDA در دسترس است یا نه و دستگاه را بر این اساس اختصاص می‌دهد. در این مورد خاص، اگر یک GPU در دسترس باشد، از آن استفاده خواهد کرد. در غیر این صورت، به استفاده از CPU برمی‌گردد. این انتخاب دستگاه برای بهینه سازی فرآیند آموزش مدل‌های یادگیری عمیق مهم است زیرا GPU ها می‌توانند محاسبات را به طور قابل توجهی تسریع کنند. کد اندازه دسته را روی ۳۲ تنظیم می‌کند که تعداد نمونه‌های ورودی پردازش شده در هر تکرار در طول آموزش را تعیین می‌کند. به طور کلی، این قطعه کد تنظیمات محیط ضروری را برای وظایف یادگیری ماشینی که شامل داده‌های متنی است، ایجاد می‌کند.

در ادامه کدهایی که پیاده سازی می کنیم به صورت زیر بیان می شود:

این بخش داده را برای کار MultiNLI MNLI از معیار ارزیابی درک عمومی زبان (GLUE) با استفاده از کتابخانه «مجموعه های داده» Hugging Face بارگیری و آماده می کند. مجموعه داده های MNLI را به سه بخش تقسیم می کند: یک مجموعه داده آموزشی ("train_data"), یک مجموعه داده اعتبارسنجی برای داده های منطبق ("val_m_data") و یک مجموعه داده اعتبارسنجی برای داده های ناهمخوان ("val_mm_data"). سپس اطلاعات مربوط به این مجموعه داده ها را چاپ می کند. علاوه بر این، یک متریک برای ارزیابی وظیفه MNLI به نام "metric" با استفاده از تابع "load_metric" از همان کتابخانه بارگیری می کند. این کد معمولاً برای راه اندازی و آماده سازی داده ها برای آموزش و ارزیابی مدل های یادگیری ماشین در تکلیف MNLI در معیار GLUE، که معمولاً برای درک زبان طبیعی و وظایف طبقه بندی متن استفاده می شود، استفاده می شود.

```
Train dataset:
Dataset({
  features: ['premise', 'hypothesis', 'label', 'idx'],
  num_rows: 392702
})
Validation matched dataset:
Dataset({
  features: ['premise', 'hypothesis', 'label', 'idx'],
  num_rows: 9815
})
Validation mismatched dataset:
Dataset({
  features: ['premise', 'hypothesis', 'label', 'idx'],
  num_rows: 9832
})
```

شکل ۸: نمایی از دیتاست

در ادامه پیاده سازی داریم:

برای پردازش داده های متنی، به ویژه برای یک کار پردازش زبان طبیعی با استفاده از یک مدل RoBERTa از پیش آموزش دیده استفاده می شود. در ابتدا، یک نشانه ساز ("RobertaTokenizer") مخصوص مدل "robertabase" بارگذاری می شود و متن به حروف کوچک تبدیل می شود. عملکرد "tokenize" برای توکن کردن دو فیلد متنی ("Premise" و "Hypothesis") در داده های

ورودی تعریف شده است، و اطمینان حاصل می‌کند که آنها به طول حداکثر ۱۲۸ توکن اضافه یا کوتاه شده‌اند. مجموعه داده‌های آموزشی ("train_data")، اعتبارسنجی منطبق ("val_m_data") و اعتبارسنجی نامتناسب ("val_mm_data") برای تکرارپذیری با ۴۲ عدد ترکیب می‌شوند و تنها ۱۰۰۰ نمونه اول از هر کدام انتخاب می‌شوند. سپس این داده‌ها با استفاده از `tokenize_function` به صورت دسته‌ای توکن می‌شوند. در نهایت، قالب این مجموعه داده‌ها روی تانسورهای `PyTorch` تنظیم می‌شود که فقط شامل ستون‌های ضروری است: «attention_mask, input_ids» و «label». این فرآیند مجموعه‌های داده را برای آموزش و اعتبارسنجی در یک مدل یادگیری ماشینی آماده می‌کند، و اطمینان حاصل می‌کند که در قالب و ساختار صحیح هستند.

```
RobertaForSequenceClassification(
  (roberta): RobertaModel(
    (embeddings): RobertaEmbeddings(
      (word_embeddings): Embedding(50265, 768, padding_idx=1)
      (position_embeddings): Embedding(514, 768, padding_idx=1)
      (token_type_embeddings): Embedding(1, 768)
      (LayerNorm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (encoder): RobertaEncoder(
      (layer): ModuleList(
        (0-11): 12 x RobertaLayer(
          (attention): RobertaAttention(
            (self): RobertaSelfAttention(
              (query): Linear(in_features=768, out_features=768, bias=True)
              (key): Linear(in_features=768, out_features=768, bias=True)
              (value): Linear(in_features=768, out_features=768, bias=True)
              (dropout): Dropout(p=0.1, inplace=False)
            )
            (output): RobertaSelfOutput(
              (dense): Linear(in_features=768, out_features=768, bias=True)
              (LayerNorm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
              (dropout): Dropout(p=0.1, inplace=False)
            )
          )
          (intermediate): RobertaIntermediate(
            (dense): Linear(in_features=768, out_features=3072, bias=True)
            (intermediate_act_fn): GELUActivation()
          )
          (output): RobertaOutput(
            (dense): Linear(in_features=3072, out_features=768, bias=True)
            (LayerNorm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
    )
  )
)
```

شکل ۹: مدل زبانی Roberta

در ادامه این بخش داریم:

آرگومان های آموزشی را برای تنظیم دقیق یک مدل RoBERTa از پیش آموزش دیده با استفاده از کتابخانه Transformers Hugging Face و تعریف تابعی برای محاسبه معیارها در حین ارزیابی تنظیم

می کند. کلاس `TrainingArguments` برای تعیین پارامترهای مختلف برای فرآیند آموزش استفاده می شود. این پارامترها شامل فهرست خروجی ("`Fine_tuned_Roberta`")، استراتژی های ارزیابی و ذخیره مدل ("دوران" برای هر دو)، نرخ یادگیری، اندازه دسته ای برای آموزش و ارزیابی، تعداد دوره های آموزشی، کاهش وزن برای منظم سازی، و دستورالعمل بارگذاری بهترین مدل در پایان آموزش بر اساس یک متریک مشخص ("دقت"). آرگومان «رزومه_از_نقطه_بازرسی» روی «درست» تنظیم شده است تا در صورت وجود، امکان ادامه آموزش از یک ایست بازرسی ذخیره شده را فراهم کند. تابع `compute_metrics` برای محاسبه عملکرد مدل در طول ارزیابی تعریف شده است. پیش‌بینی‌ها و برچسب‌ها را می‌گیرد، `argmax` پیش‌بینی‌ها را محاسبه می‌کند تا محتمل‌ترین پیش‌بینی‌های کلاس را به دست آورد، و سپس از یک «متریک» از پیش تعریف‌شده (مانند دقت) برای مقایسه این پیش‌بینی‌ها با برچسب‌های واقعی استفاده می‌کند و متریک محاسبه‌شده را برمی‌گرداند. این تنظیم برای آموزش یک مدل یادگیری ماشینی قوی و با عملکرد خوب ضروری است و امکان نظارت و بهینه سازی عملکرد آن را بر اساس دقت فراهم می‌کند.

حال در ادامه یک `callback` سفارشی برای کلاس `Hugging Face Trainer` تعریف می‌کند و یک محیط آموزشی برای یک مدل یادگیری ماشینی تنظیم می‌کند. کلاس «`CustomCallback`» یک زیر کلاس از «`TrainerCallback`» است و برای انجام اقدامات اضافی در پایان هر دوره در طول آموزش طراحی شده است. به طور خاص، برای ارزیابی مدل در مجموعه داده های آموزشی در پایان هر دوره، علاوه بر ارزیابی استاندارد در مجموعه داده های اعتبار سنجی تنظیم شده است. این با نادیده گرفتن روش «`on_epoch_end`» به دست می‌آید، جایی که بررسی می‌کند که آیا ارزیابی باید انجام شود («`control.should_evaluate`») و سپس ارزیابی روی مجموعه داده آموزشی («`self._trainer.train_dataset`») انجام می‌دهد. شی «`Trainer`» با پارامترهای مختلفی از جمله مدل، آرگومان‌های آموزشی («`args`»), مجموعه داده‌های آموزش و ارزیابی («`train_data`») و «`val_m_data`»، نشانه‌ساز و تابع «`compute_metrics`» سفارشی برای محاسبه متریک در طول ارزیابی: اگرچه پاسخ تماس سفارشی ("`CustomCallback`") تعریف شده است، اما نظر داده می‌شود و به مربی اضافه نمی‌شود، که نشان می‌دهد این یک جزء اختیاری است که در صورت نیاز می‌تواند فعال شود. متد «`trainer.train()`» در پایان فرآیند آموزش را شروع می‌کند. این تنظیم اجازه می‌دهد تا نظارت دقیقی بر عملکرد مدل، نه تنها بر روی داده های اعتبار سنجی دیده نشده، بلکه بر روی داده های آموزشی، ارائه بینشی در مورد پویایی یادگیری آن باشد.

[160/160 02:43, Epoch 5/5]			
Epoch	Training Loss	Validation Loss	Accuracy
1	No log	1.106683	0.328000
2	No log	1.090447	0.448000
3	No log	1.015701	0.488000
4	No log	1.002251	0.531000
5	No log	0.955415	0.557000

Checkpoint destination directory Fine_tuned_Roberta/checkpoint-32 already exists and is non-empty.Saving will proceed but saved results may be invalid.
 Checkpoint destination directory Fine_tuned_Roberta/checkpoint-64 already exists and is non-empty.Saving will proceed but saved results may be invalid.
 Checkpoint destination directory Fine_tuned_Roberta/checkpoint-96 already exists and is non-empty.Saving will proceed but saved results may be invalid.
 Checkpoint destination directory Fine_tuned_Roberta/checkpoint-128 already exists and is non-empty.Saving will proceed but saved results may be invalid.
 Checkpoint destination directory Fine_tuned_Roberta/checkpoint-160 already exists and is non-empty.Saving will proceed but saved results may be invalid.
 TrainOutput(global_step=160, training_loss=0.9854669570922852, metrics={'train_runtime': 164.0284, 'train_samples_per_second': 30.483, 'train_steps_per_second': 0.975, 'total_flos': 32889177216000.0, 'train_loss': 0.9854669570922852, 'epoch': 5.0})

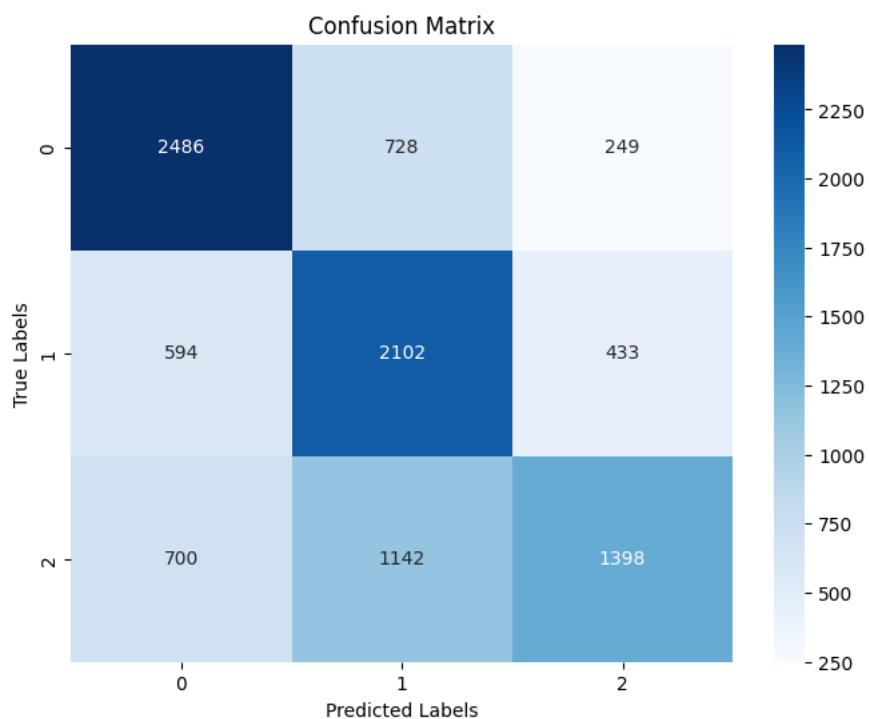
شکل ۱۰: پیاده سازی کلاس Callback

در ادامه پیاده سازی داریم قطعه کد برای ارزیابی عملکرد یک مدل یادگیری ماشینی آموزش دیده (احتمالاً یک مدل طبقه بندی متن با استفاده از Transformers Hugging Face) در یک مجموعه داده آزمایشی ("val_mm_data") است. ابتدا، پیش‌بینی‌های مدل روی مجموعه داده آزمایشی با استفاده از روش «trainer.predict()» به دست می‌آیند. سپس پیش‌بینی‌ها برای استخراج محتمل‌ترین برچسب‌های کلاس با استفاده از «np.argmax()» پردازش می‌شوند. این برچسب‌های پیش‌بینی شده با برچسب‌های واقعی از مجموعه داده آزمایشی برای محاسبه ماتریس سردرگمی مقایسه می‌شوند، ابزاری که به درک عملکرد مدل در کلاس‌های مختلف کمک می‌کند. سپس ماتریس سردرگمی چاپ می‌شود و تفکیک دقیقی از پیش‌بینی‌های صحیح و نادرست مدل برای هر کلاس ارائه می‌کند. علاوه بر این، یک گزارش طبقه‌بندی با استفاده از «classification_report()» تولید می‌شود که معیارهای مختلفی مانند دقت، یادآوری، و امتیاز F1 را برای هر کلاس و همچنین میانگین‌های کلی ارائه می‌کند. این معیارها یک نمای کلی جامع از عملکرد طبقه‌بندی مدل ارائه می‌دهند و نقاط قوت و ضعف آن را در پیش‌بینی کلاس‌های مختلف برجسته می‌کنند. این ارزیابی برای درک اثربخشی مدل و هدایت بهبودهای بیشتر بسیار مهم است.

نتایج Roberta

Classification Report:				
	precision	recall	f1-score	support
0	0.66	0.72	0.69	3463
1	0.53	0.67	0.59	3129
2	0.67	0.43	0.53	3240
accuracy			0.61	9832
macro avg	0.62	0.61	0.60	9832
weighted avg	0.62	0.61	0.60	9832

شکل ۱۱: نتایج Roberta



شکل ۱۲: ماتریس درهم آمیختگی Roberta

در ادامه این بخش به پیاده سازی مدل با Freeze می پردازیم:

مدل RoBERTa را برای یک کار طبقه بندی دنباله ای با تمرکز خاص بر فریز کردن لایه های خاص در طول آموزش، مقداردهی اولیه و پیکربندی می کند. مدل «RobertaForSequenceClassification» از مدل «robertabase» از پیش آموزش دیده Hugging Face بارگیری شده است که برای یک کار طبقه بندی با سه برچسب خروجی («تعداد_برچسبها=۳») طراحی شده است. هدف از انجماد لایه ها جلوگیری از به روزرسانی وزن آن لایه ها در طول آموزش است که می تواند برای تنظیم دقیق یک مدل در یک کار خاص با داده های محدود مفید باشد. در این قطعه، تمام پارامترها در لایه جاسازی مدل ثابت می شوند، همانطور که با تنظیم «requires_grad» روی «False» نشان داده می شود. به علاوه، شش لایه اول رمزگذار («تعداد_لایه های_به_فریز = ۶») نیز منجمد شده اند. یعنی در حین تمرین فقط وزنه های لایه های بالاتر و سر رده بندی به روز می شود. در نهایت، مدل برای آموزش کارآمد به یک دستگاه مشخص («DEVICE») منتقل می شود که می تواند یک CPU یا GPU باشد. این رویکرد انجماد جزئی یک مدل اغلب در سناریوهای یادگیری انتقال استفاده می شود تا از ویژگی های از پیش آموخته شده مدل پایه استفاده کند و در عین حال به مدل اجازه می دهد تا با ویژگی های کار جدید سازگار شود.

```

RobertaForSequenceClassification(
  (roberta): RobertaModel(
    (embeddings): RobertaEmbeddings(
      (word_embeddings): Embedding(50265, 768, padding_idx=1)
      (position_embeddings): Embedding(514, 768, padding_idx=1)
      (token_type_embeddings): Embedding(1, 768)
      (LayerNorm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (encoder): RobertaEncoder(
      (layer): ModuleList(
        (0-11): 12 x RobertaLayer(
          (attention): RobertaAttention(
            (self): RobertaSelfAttention(
              (query): Linear(in_features=768, out_features=768, bias=True)
              (key): Linear(in_features=768, out_features=768, bias=True)
              (value): Linear(in_features=768, out_features=768, bias=True)
              (dropout): Dropout(p=0.1, inplace=False)
            )
            (output): RobertaSelfOutput(
              (dense): Linear(in_features=768, out_features=768, bias=True)
              (LayerNorm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
              (dropout): Dropout(p=0.1, inplace=False)
            )
          )
          (intermediate): RobertaIntermediate(
            (dense): Linear(in_features=768, out_features=3072, bias=True)
            (intermediate_act_fn): GELUActivation()
          )
          (output): RobertaOutput(
            (dense): Linear(in_features=3072, out_features=768, bias=True)
            (LayerNorm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
    )
  )
)

```

شکل ۱۳: مدل Roberta

حال مشابه حالت قبل یک callback سفارشی را برای آموزش یک مدل یادگیری ماشینی با استفاده از کلاس Hugging Face Trainer تعریف می‌کند، به‌ویژه به منظور ارزیابی مدل در مجموعه داده‌های آموزشی در پایان هر دوره آموزشی. کلاس «CustomCallback» یک زیر کلاس از «TrainerCallback» است و برای انجام یک مرحله ارزیابی اضافی، روش «on_epoch_end» را لغو می‌کند. در این روش، اگر شرط «control.should_evaluate» درست باشد (که نشان دهنده پایان یک دوره است)، callback یک کپی از شی کنترل ایجاد می‌کند و سپس روش «ارزیابی» مربی را در مجموعه داده آموزشی فراخوانی می‌کند. (trainer.train_dataset_). این ارزیابی با «پیشوند_کلید_متریک» برچسب گذاری شده است، که آن را از ارزیابی استاندارد در مجموعه داده اعتبار سنجی متمایز می‌کند. خود شی «Trainer» با مدل، آرگومان‌های آموزشی، مجموعه داده‌های آموزش و ارزیابی، نشانه‌ساز و یک تابع محاسبات متریک سفارشی پیکربندی شده است. اگرچه

"CustomCallback" نمونه سازی شده است، اما نظر داده می شود و به مربی اضافه نمی شود، که نشان می دهد اختیاری است و در صورت نیاز می توان آن را فعال کرد. در نهایت، متد `trainer.train()` فرآیند آموزش را شروع می کند. این راه اندازی، با تماس سفارشی اختیاری، یک رژیم آموزشی جامع با توانایی نظارت بر عملکرد مدل نه تنها بر روی داده های اعتبارسنجی، بلکه بر روی داده های آموزشی برای تجزیه و تحلیل بیش از حد فراهم می کند.

[160/160 02:52, Epoch 5/5]

Epoch	Training Loss	Validation Loss	Accuracy
1	No log	1.105728	0.328000
2	No log	1.097146	0.366000
3	No log	1.063065	0.451000
4	No log	1.054658	0.450000
5	No log	1.011699	0.480000

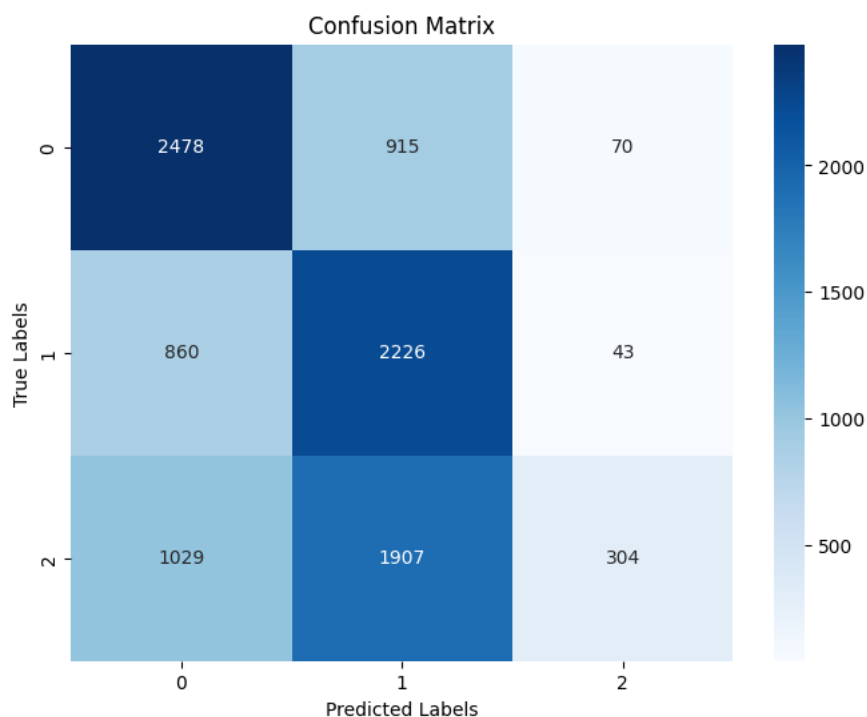
Checkpoint destination directory Fine_tuned_Roberta/checkpoint-32 already exists and is non-empty.Saving will proceed but saved results may be invalid.
 Checkpoint destination directory Fine_tuned_Roberta/checkpoint-64 already exists and is non-empty.Saving will proceed but saved results may be invalid.
 Checkpoint destination directory Fine_tuned_Roberta/checkpoint-96 already exists and is non-empty.Saving will proceed but saved results may be invalid.
 Checkpoint destination directory Fine_tuned_Roberta/checkpoint-128 already exists and is non-empty.Saving will proceed but saved results may be invalid.
 Checkpoint destination directory Fine_tuned_Roberta/checkpoint-160 already exists and is non-empty.Saving will proceed but saved results may be invalid.
 TrainOutput(global_step=160, training_loss=1.0496860504150392, metrics={'train_runtime': 173.2964, 'train_samples_per_second': 28.852, 'train_steps_per_second': 0.923, 'total_flos': 328891772160000.0, 'train_loss': 1.0496860504150392, 'epoch': 5.0})

شکل ۱۴: پیاده سازی کلاس Callback

نتایج Freeze

Classification Report:				
	precision	recall	f1-score	support
0	0.57	0.72	0.63	3463
1	0.44	0.71	0.54	3129
2	0.73	0.09	0.17	3240
accuracy			0.51	9832
macro avg	0.58	0.51	0.45	9832
weighted avg	0.58	0.51	0.45	9832

شکل ۱۵: نتایج Freeze



شکل ۱۶: ماتریس در هم آمیختگی Freeze

حال به توضیح بخش LoRA می‌پردازیم:

ابتدا با استفاده از کتابخانه‌ی peft بخش‌های LoRA را به کد اضافه می‌کنیم، این مدل LoRA برای وظیفه دسته‌بندی سری هست و می‌توانیم مقدار rank آن را مشخص کنیم که در تعداد پارامترهای یادگیری تاثیر دارد.

سپس LoRA را به مدل خود اضافه می‌کنیم، تعداد پارامترهای مدل قبل و بعد از اضافه کردن LoRA به شکل زیر می‌باشد:

```
trainable params: 355362819 || all params: 355362819 || trainable%: 100.0
trainable params: 4198403 || all params: 359561222 || trainable%: 1.167646215197255
```

شکل ۱۷ تعداد پارامترهای قابل آموزش مدل قبل و بعد از LoRA

حال مدل را finetune می‌کنیم و نتایج آن به صورت زیر می‌باشد:

```

Classification Report - Test Set:
              precision    recall  f1-score   support

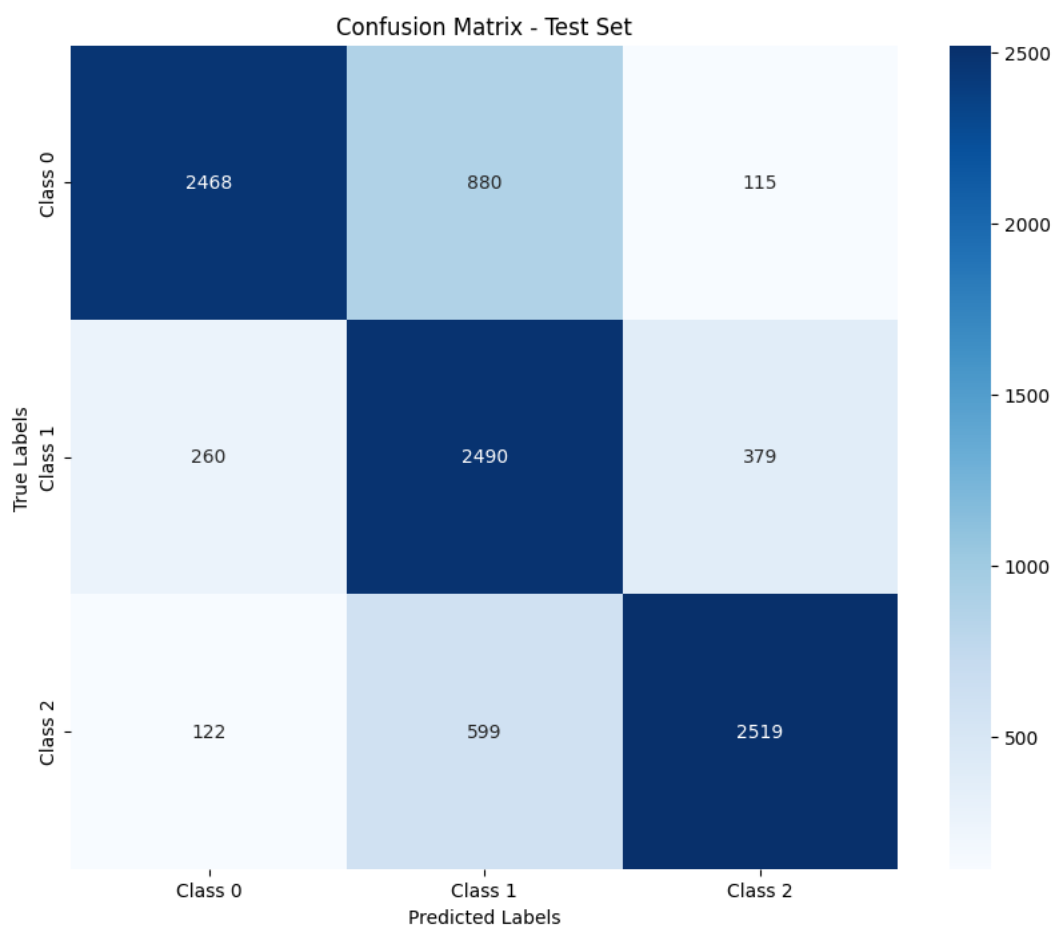
   Class 0       0.87        0.71        0.78       3463
   Class 1       0.63        0.80        0.70       3129
   Class 2       0.84        0.78        0.81       3240

 accuracy              0.76       9832
 macro avg              0.78       9832
 weighted avg           0.78       9832

Test Accuracy: 0.7605

```

شکل ۱۸: نتایج LoRA



شکل ۱۹: ماتریس درهم آمیختگی LoRA

مقایسه کلی

مفاهیم "LoRA" (انطباق با رتبه پایین) و لایه های "انجماد" در یک شبکه عصبی، به ویژه در زمینه مدل های ترانسفورماتور مانند BERT یا RoBERTa، نشان دهنده دو روش متمایز برای انطباق و تنظیم دقیق مدل های از پیش آموزش دیده برای خاص است. وظایف یا مجموعه داده ها

LoRA (انطباق با رتبه پایین)

مفهوم: LoRA شامل اصلاح یک مدل از پیش آموزش دیده با قرار دادن ماتریس های با رتبه پایین در لایه های آن است که در طول تنظیم دقیق یاد می شوند. به جای به روزرسانی وزن های اصلی مدل، LoRA ماتریس های کوچک و قابل آموزش را معرفی می کند که وزن های موجود را به روشی با رتبه پایین تطبیق می دهند. این رویکرد با هدف گرفتن سازگاری های خاص کار با پارامترهای آموزش پذیر کمتر است.

مزایا: LoRA امکان تطبیق کارآمد یک مدل را با افزایش نسبتاً کمی در تعداد پارامترها فراهم می کند. می تواند با ثبت نکات ظریف و ظریف و بدون نیاز به آموزش مجدد کل مدل، منجر به عملکرد بهتر در وظایف خاص شود.

سناریوهای استفاده: LoRA به ویژه در سناریوهایی مفید است که در آن شخص می خواهد یک مدل بزرگ از پیش آموزش دیده را با یک کار خاص بدون هزینه محاسباتی آموزش پارامترهای زیادی تطبیق دهد. همچنین زمانی مفید است که داده های مربوط به کار محدود باشد و تطبیق بیش از حد یک نگرانی باشد.

لایه های انجماد:

مفهوم: انجماد لایه ها در یک شبکه عصبی به معنای تنظیم وزن لایه های خاص غیرقابل آموزش است. این اغلب با مدل های از پیش آموزش دیده انجام می شود که در آن تنها چند لایه بالایی (یا فقط لایه طبقه بندی) برای یک کار خاص به خوبی تنظیم شده اند. لایه های پایین تر، که ویژگی های عمومی بیشتری را ثبت می کنند، منجمد می شوند.

مزایا: انجماد تعداد پارامترهایی را که نیاز به آموزش دارند کاهش می‌دهد، هزینه‌های محاسباتی و خطر بیش از حد برازش را کاهش می‌دهد. این به مدل اجازه می‌دهد تا دانش عمومی آموخته شده در طول آموزش را حفظ کند و فقط ویژگی‌های سطح بالاتر را برای کار خاص تطبیق دهد.

سناریوهای استفاده: فریز کردن معمولاً در یادگیری انتقال استفاده می‌شود، به‌ویژه زمانی که مدل‌های بزرگ از پیش آموزش دیده را روی یک مجموعه داده کوچک تنظیم می‌کنید. زمانی موثر است که مدل پایه از قبل روی یک مجموعه داده بزرگ و متنوع به خوبی آموزش داده شده باشد و تنها به انطباق‌های جزئی برای کار مورد نیاز نیاز است.

مقایسه:

هدف: هر دو روش با هدف تطبیق مدل‌های از پیش آموزش دیده با وظایف خاص به طور موثر انجام می‌شوند، اما آنها این کار را به روش‌های مختلف انجام می‌دهند. LoRA حداقل ظرفیت سازگاری را اضافه می‌کند، در حالی که انجماد سازگاری را به لایه‌های خاصی محدود می‌کند.

آموزش پارامتر: LoRA شامل آموزش ماتریس‌های با رتبه پایین اضافی است، در حالی که فریز کردن شامل آموزش تنها پارامترهای لایه‌های منجمد نشده است.

انعطاف پذیری: LoRA تعادلی بین حفظ دانش از پیش آموزش دیده و انطباق با وظایف جدید ارائه می‌دهد، در حالی که انجماد بیشتر در مورد استفاده از دانش موجود با حداقل سازگاری است.

کارایی محاسباتی: هر دو روش در مقایسه با تنظیم دقیق مدل کامل از نظر محاسباتی کارآمد هستند، اما کارایی خاص به تعداد پارامترهای ثابت شده یا رتبه ماتریس‌ها در LoRA بستگی دارد.

زمان محاسباتی: از آنجا که با استفاده از LoRa فقط تعداد کمی از پارامترهای مدل آموزش می‌بینند با استفاده از لورا توقع داریم در مدت زمان کمتری آموزش به اتمام برسد.

با توجه به نتایج مشخص است که مدل اصلی ۲۰ دقیقه در ۵ اپاک آموزش دیده است ولی مدل لورا فقط ۱۰ دقیقه طول کشیده است و نتایج به اندازه کافی مناسبی گرفته است.

مدل Lora با وجود مدت زمان کمتر آموزش و همچنین تعداد پارامتر کمتر، توانسته دقتی نزدیک دقت مدل اصلی (در دفعات مختلف ممکن بود تا ۵ درصد بهتر نیز بشود) گرفته است.

دقت مدل اصلی: ۷۵

دقت LoRA: ۷۳

یعنی مشخصاً Lora یک روش تنظیم دقیق بسیار خوب می‌باشد که با وجود مقداری کم شدن دقت می‌تواند تا حد زیادی عملکردی در حد تنظیم دقیق کل مدل داشته باشد و ما با منابع و زمان کمتر به این عملکرد می‌رسیم.

به طور خلاصه، LoRA و انجماد هر دو استراتژی‌های مؤثری برای تنظیم دقیق مدل‌های از پیش آموزش‌دیده هستند، اما نیازها و محدودیت‌های محاسباتی کمی متفاوت را برآورده می‌کنند. انتخاب بین آنها به نیازهای خاص کار، اندازه مجموعه داده و منابع محاسباتی موجود بستگی دارد.

۲-۴. چرا LoRA ؟

روش Lora به شدت efficient تر از finetune کردن کل مدل است و نتیجه تقریباً یکسانی با آن می‌گیرد، این روش با اضافه کردن تعداد کمی پارامتر به مدل در مدت زمان کمتر و با دقت و خطای مناسب همان نتایج مدل را به ما خواهد داد.

Finetune کردن کل مدل به شدت هزینه بر است و باید توانایی و داده‌های زیاد برای این کار داشته باشیم تا مدل مدل هم وظیفه جدید را یاد بگیرد و هم وظیفه قبلی را فراموش نکند، همین‌طور finetune کردن بخشی از مدل راه حلی منطقی تر و بهینه‌تر می‌باشد ولی با این حال بهترین نتایج را لزوماً نخواهد داد، روش Lora به طور کلی برتری‌های زیر را دارد.

- LoRA نسبت به تنظیم دقیق (fine-tuning) بسیار سریع‌تر است، زیرا فقط یک تعداد کوچک از پارامترها را به‌روزرسانی می‌کند به جای تمام پارامترهای مدل زبان بزرگ
- LoRA در مصرف حافظه و فضای ذخیره‌سازی به مراتب کارآمدتر از تنظیم دقیق است، زیرا تنها نیاز به ذخیره ماتریس‌های تجزیه رتبه برای هر وظیفه دارد به جای کل مدل تنظیم دقیق شده
- LoRA در کیفیت مدل در مقایسه با تنظیم دقیق بر روی مدل‌های زبان بزرگ مختلف مانند GPT-۲، DeBERTa، GPT-۳ بهتر یا مساوی عمل می‌کند، با این که تعداد کمتری پارامتر قابل آموزش دارد
- LoRA هیچ تأخیر استنتاجی ایجاد نمی‌کند، زیرا ماتریس‌های قابل آموزش می‌توانند در هنگام استقرار با وزن‌های یخ‌زده ترکیب شوند، به عنوان مقابل مانند آداپتورها که لایه‌های اضافی به مدل اضافه می‌کنند.

- LoRA یک روش ارتجاعی است که می‌تواند با تکنیک‌های دیگر تنظیم دقیق مانند تنظیم پیشوندی (prefix-tuning) ترکیب شود.

حال اگر قصد finetune کردن مدل برای چند وظیفه را داشته باشیم راه‌های زیر را داریم:

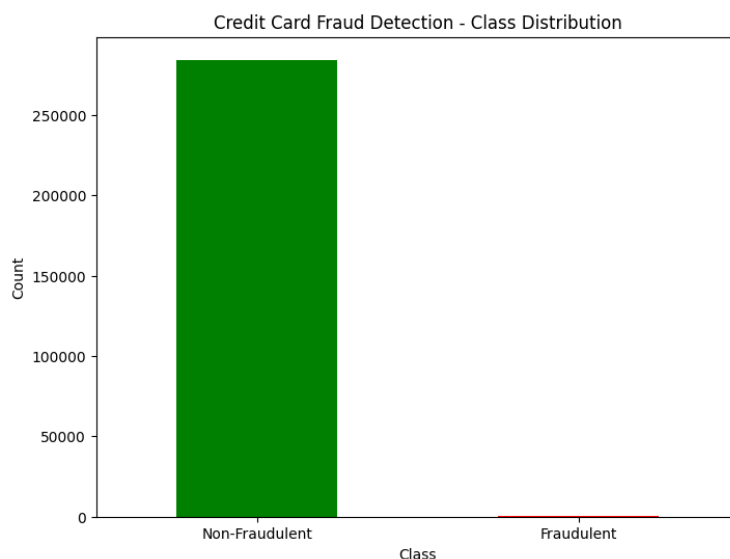
- مجموعه داده‌ای از این وظایف جدید و وظایفی که مدل قبلاً به آن آگاه بوده جمع کنیم و مدل را برای این هر دو آن‌ها آموزش بدهیم تا وظایف جدید را یاد بگیرد و وظایف قبلی را فراموش نکند.
- اگر مدل به اندازه کافی بزرگ باشد و قبلاً با داده‌های کافی آموزش دیده باشد، می‌توانیم با استفاده از روش‌های few shot این مساله را حل کنیم (که fine tuning نیست)
- می‌توانیم Lora را برای وظیفه اول آموزش بدهیم و آپدیت‌های آن را با وزن‌ها جمع کنیم، سپس Lora را برای وظیفه دوم نیز آموزش بدهیم و آپدیت‌ها را با وزن‌ها جمع کنیم. به طور کلی برای وظایف مختلف می‌توانیم آپدیت‌های lora را با هم جمع کرده و به وزن‌های مدل اضافه کنیم.

۳. تشخیص تقلب

۳-۱. آشنایی با دیتاست

این دیتاست مربوط به تراکنش‌های بانکی می‌باشد که آیا دارای کلاهبرداری بوده‌اند یا خیر. این مجموعه داده شامل تراکنش‌هایی با ویژگی‌های مختلفی است که از ۷۱ تا ۷۲۸ برچسب‌گذاری شده‌اند، که احتمالاً نتیجه تغییر تجزیه و تحلیل مؤلفه اصلی (PCA) به دلایل حفظ حریم خصوصی، و همچنین «زمان» و «مقدار» است. همچنین یک ستون 'Class' وجود دارد که متغیر هدف است و نشان می‌دهد که آیا تراکنش تقلبی است (۱) یا خیر (۰).

هیستوگرام توزیع کلاس در مجموعه داده، عدم تعادل قابل توجهی را بین دو کلاس نشان می‌دهد: غیر متقلبان (کلاس ۰) و متقلب (کلاس ۱). تعداد تراکنش‌های غیر متقلبان بسیار بیشتر از معاملات متقلبان است.



شکل ۲۰ هیستوگرام داده‌ها

این عدم تعادل یک چالش مهم برای مدل‌های یادگیری ماشینی است. اگر مدلی بر روی این مجموعه داده همانطور که هست آموزش داده شود، ممکن است نسبت به پیش‌بینی تراکنش‌ها به‌عنوان غیر جعلی تعصب داشته باشد، زیرا آن کلاس بر مجموعه داده تسلط دارد. در نتیجه، این مدل ممکن است در تشخیص تراکنش‌های جعلی، که کلاس اقلیت هستند، اما از منظر کشف تقلب، علاقه بالاتری دارند، عملکرد ضعیفی داشته باشد.

دقت مدل ممکن است بالا به نظر برسد، زیرا می‌تواند طبقه اکثریت را به درستی پیش‌بینی کند، اما توانایی آن در تشخیص طبقه اقلیت - که برای سیستم تشخیص تقلب بسیار مهم است - احتمالاً ناکافی است. به عبارت ساده‌تر دقت کلی مدل ممکن است ۹۹ درصد باشد اما روی داده‌های کلاس اکثریت دقت ۱۰۰ داشته باشیم و روی داده‌های کلاس اقلیت دقت ۰ که به این معنی است همه را کلاس اکثریت پیش‌بینی میکند و عملاً مدل چیزی یاد نگرفته است. این یک مشکل رایج در مجموعه‌های داده با کلاس‌های بسیار نامتعادل است و به استراتژی‌های خاصی برای رسیدگی نیاز دارد، مانند نمونه‌گیری مجدد از مجموعه داده، استفاده از تکنیک‌های تشخیص ناهنجاری، یا اعمال معیارهای مختلف برای ارزیابی مدل مانند فراخوان دقیق AUC-PR، F_1 ، امتیاز، یا استفاده از الگوریتم‌هایی که حساسیت کمتری نسبت به عدم تعادل کلاس دارند.

۲-۳. پیاده سازی معماری مقاله

در این بخش از کتابخانه TensorFlow و مدلی از نوع Sequential برای ساخت یک شبکه عصبی کانولوشنال D۱ استفاده می‌کند. این شبکه برای یادگیری از داده‌های ورودی با سه بعد $(X_train.shape[1], X_train.shape[2])$ استفاده می‌شود که در اینجا X_train به عنوان ورودی شبکه تعیین شده است.

در زیر توضیح داده‌ام که هر لایه از مدل چه کارهایی انجام می‌دهد:

۱. **Conv۱D Layer**: این لایه یک لایه کانولوشنال D۱ با ۳۲ فیلتر، هر کدام با یک کرنل به طول ۲ و تابع فعال‌سازی ReLU دارد. این لایه برای استخراج ویژگی‌های مختلف از داده‌های ورودی استفاده می‌شود.

۲. **BatchNormalization**: این لایه برای استانداردسازی میانگین و واریانس خروجی‌های لایه قبلی (Conv۱D Layer) استفاده می‌شود، که به بهبود آموزش مدل کمک می‌کند.

۳. **MaxPool۱D**: این لایه از عملگر MaxPooling برای کاهش ابعاد خروجی لایه قبلی استفاده می‌کند. این باعث کاهش تعداد ویژگی‌ها و افزایش مقیاس‌پذیری شبکه می‌شود.

۴. **Dropout**: این لایه برای جلوگیری از بر Overfitting در مدل استفاده می‌شود. با احتمال ۰.۲، تعدادی از واحدهای خروجی لایه قبلی را به صورت تصادفی خاموش می‌کند.

۵. **Conv۱D Layer**: این لایه یک لایه کانولوشنال دیگر با ۶۴ فیلتر و تابع فعال‌سازی ReLU است.

۶. BatchNormalization ۲: مانند ۱ BatchNormalization، این لایه برای استانداردسازی میانگین و واریانس خروجی‌های لایه قبلی (۲ Conv۱D Layer) استفاده می‌شود.

۷. MaxPool۱D ۲: مانند ۱ MaxPool۱D، این لایه از عملگر MaxPooling برای کاهش ابعاد خروجی لایه قبلی استفاده می‌شود.

۸. Dropout ۲: مانند ۱ Dropout، این لایه با احتمال ۰.۵ برای جلوگیری از Overfitting از طریق خاموش کردن واحدهای تصادفی استفاده می‌شود.

۹. Flatten: این لایه به داده‌های خروجی لایه قبلی اعمال می‌شود و آنها را به یک بردار یک بعدی تبدیل می‌کند.

۱۰. ۱ Dense: این لایه با ۶۴ واحد و تابع فعال‌سازی ReLU برای اعمال یادگیری ژرف بر روی ویژگی‌های استخراج شده توسط لایه‌های کانولوشنال استفاده می‌شود.

۱۱. ۳ Dropout: مانند ۱ Dropout و ۲، این لایه با احتمال ۰.۵ برای جلوگیری از Overfitting استفاده می‌شود.

۱۲. ۲ Dense: این لایه دیگر با ۶۴ واحد و تابع فعال‌سازی ReLU استفاده می‌شود.

۱۳. ۳ Dense: این لایه نهایی با یک واحد و تابع فعال‌سازی sigmoid برای انجام مسئله دسته‌بندی دودویی استفاده می‌شود.

این مدل به صورت یک لایه‌ی کانولوشنال و دو لایه‌ی میانی کاملاً متصل (Dense) به علاوه‌ی لایه‌های BatchNormalization، MaxPooling و Dropout برای ایجاد یک معماری عصبی کامل در

شبکه‌های عصبی عمیق استفاده شده است. این مدل به منظور یادگیری از داده‌های دوره‌ای D۱ مورد استفاده قرار می‌گیرد.

تصویر زیر تمام مدل را نشان می‌دهد که به صورت کامل از مقاله پیاده‌سازی شده است:

Layer	Description
1	Conv1D Layer (filters = 32, kernel_size = 2, activation = relu)
2	BatchNormalization
3	MaxPool1D (pool_size = 2)
4	Dropout (rate = 0.2)
5	Conv1D Layer (filters = 64, kernel_size = 2, activation = relu)
6	BatchNormalization
7	MaxPool1D (pool_size = 2)
8	Dropout (rate = 0.5)
9	Flatten
10	Dense (units = 64, activation = relu)
11	Dropout (rate = 0.5)
12	Dense (units = 64, activation = relu)

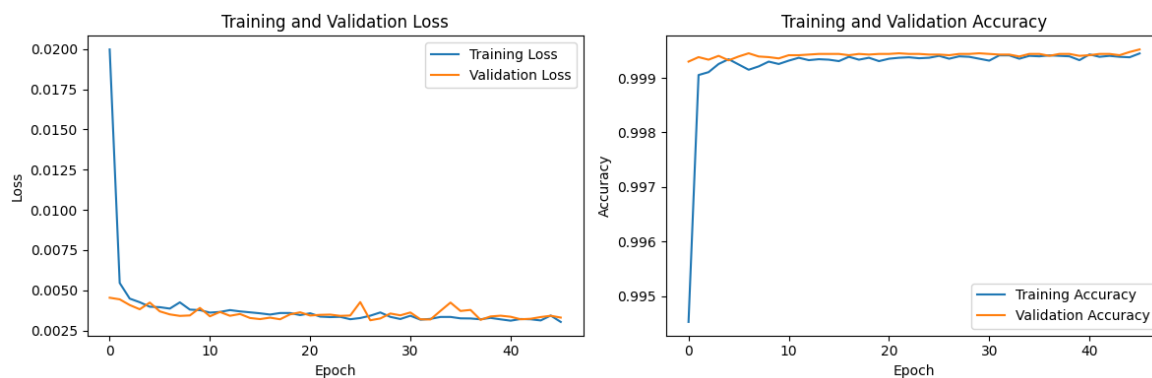
شکل ۲۱ معماری مدل

همچنین هایپرپارامترهای مدل به شرح زیر می‌باشند:

Hyperparameter	Value
epochs	46
optimizer	Adam
learning_rate	0.0001
loss	binary_crossentropy
metrics	accuracy

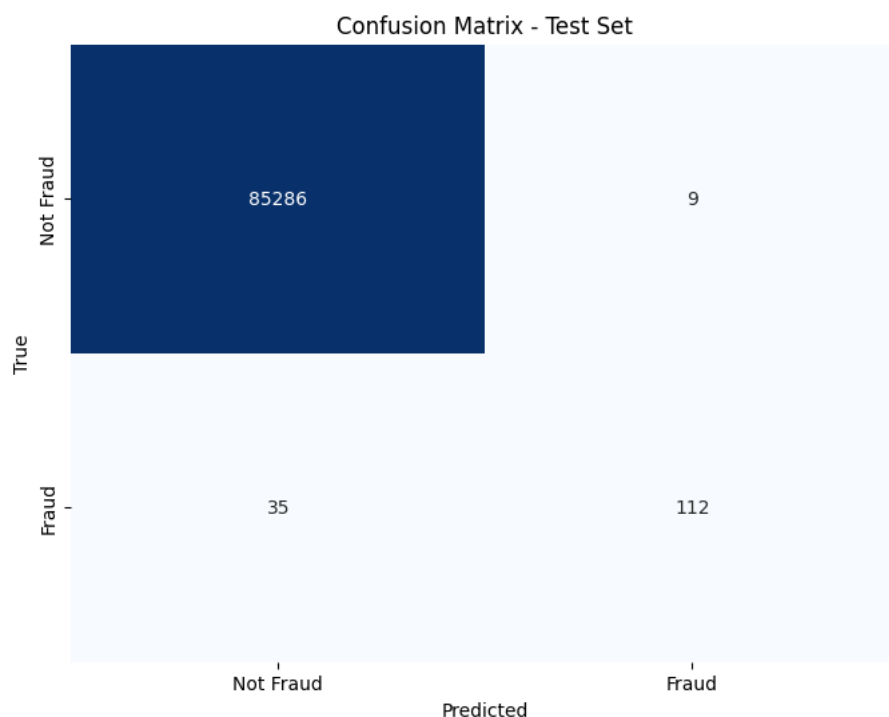
شکل ۲۲ هایپر پارامترهای مدل

حال مدل را برای ۴۶ تا ایپاک آموزش می‌دهیم که نتایج به صورت زیر شد:



شکل ۲۳: نتایج پیاده سازی مدل

با توجه به نمودار **loss** به نظر می‌آید مدل **overfit** نشده است و عملکرد روی داده‌های آموزش و **validation** به اندازه کافی خوب هست، با این حال دلیل این موضوع می‌تواند این باشد که داده‌های **validation** نیز شامل داده‌های کلاس اکثریت هست و مدل روی آن کلاس خوب عمل کرده، به طور کلی دقت برای نظر دادن راجع به این نوع از داده‌ها معیار مناسبی نیست. عملکرد مدل با استفاده از ماتریس درهم‌ریختگی به شکل زیر می‌باشد:



شکل ۲۴: ماتریس درهم آمیختگی مدل

Classification Report:				
	precision	recall	f1-score	support
0	1.00	1.00	1.00	85295
1	0.93	0.76	0.84	147
accuracy			1.00	85442
macro avg	0.96	0.88	0.92	85442
weighted avg	1.00	1.00	1.00	85442

شکل ۲۵: نتایج مدل

مشخص است که دقت مدل به صورت کلی خوب به نظر می‌رسد اما روی کلاس ۱ که داده‌های اقلیت هست عملکرد ضعیفی دارد. این موضوع را با استفاده از معیارهای دیگر می‌توان فهمید:

→ Precision: 0.8828125
Recall: 0.7687074829931972
F1-score: 0.8218181818181819

شکل ۲۶ معیارهایی صحت و recall

۳-۳. نمونه برداری

۱. توضیح روش ADASYN و مقایسه با روش‌های دیگر نمونه‌برداری

ADASYN یک روش نمونه‌برداری هوشمندانه است که برای مقابله با مشکلات عدم تعادل در داده‌های آموزشی در مسائل طبقه‌بندی ماشینی استفاده می‌شود. در مسائلی که تعداد نمونه‌های هر کلاس در داده‌ها به طور نسبی ناهمسان است، ممکن است مدل طبقه‌بندی با مشکل روبرو شود و به کلاس‌های کم‌نمونه‌ای توجه نکند. ADASYN با تولید نمونه‌های مصنوعی برای کلاس‌های کم‌نمونه‌ای و تعادل تعداد نمونه‌ها، به بهبود عملکرد مدل کمک می‌کند.

توضیح عملکرد ADASYN:

۱. تعیین وزن برای هر نمونه: این روش وزن مختص به هر نمونه را تعیین می‌کند. نمونه‌های کمیتی که برای تعادل در داده‌ها نیاز دارند، وزن بیشتری دریافت می‌کنند.

۲. محاسبه میزان نیاز به نمونه جدید: بر اساس اختلاف تعداد نمونه‌ها بین کلاس‌ها، میزان نیاز به نمونه جدید برای هر کلاس محاسبه می‌شود.

۳. تولید نمونه جدید: برای هر نمونه کمیتی که بر اساس محاسبات به آن نیاز است، نمونه‌های جدید تولید می‌شوند.

۴. افزودن نمونه‌های تولید شده به داده‌های آموزشی: نمونه‌های تولید شده به داده‌های آموزشی افزوده می‌شوند تا تعداد نمونه‌ها برای هر کلاس بهبود یابد.

مزایا:

- بهبود تعادل کلاس‌ها در داده‌های آموزشی.

- افزایش دقت مدل در کلاس‌های کم‌نمونه‌ای.

معایب:

- افزایش حجم داده‌های آموزشی به دلیل تولید نمونه‌های مصنوعی.

- حساسیت به نویز در داده‌ها.

- زمان‌بر بودن فرآیند تولید نمونه‌های مصنوعی.

- برای نمونه‌های اقلیت که به صورت پراکنده توزیع شده‌اند، هر محله ممکن است فقط دارای ۱ نمونه اقلیت باشد.

- دقت ADASYN ممکن است به دلیل ماهیت سازگاری آسیب ببیند.

نسبت به روش‌های دیگر نمونه‌برداری، ADASYN به دلیل توجه به نمونه‌های کم‌نمونه‌ای بهبود قابل توجهی می‌دهد. اما هر روشی مزایا و معایب خود را دارد و انتخاب بهترین روش بستگی به مشخصات خاص مسئله و داده‌های آموزشی دارد.

۳. روش پیاده‌سازی:

مراحل پیاده‌سازی به شرح زیر می‌باشد:

مرحله ۱

نسبت اقلیت به اکثریت را با استفاده از مثال‌های زیر محاسبه کنید:

$$d = \frac{m_s}{m_l}$$

که در آن m_s و m_l به ترتیب # نمونه های کلاس اقلیت و اکثریت هستند. اگر d کمتر از یک آستانه مشخص است، الگوریتم را مقداردهی اولیه کنید.

گام ۲

تعداد کل داده های اقلیت مصنوعی برای تولید را محاسبه کنید.

$$G = (m_l - m_s)\beta$$

شکل ۲۷ تعداد کل داده های اقلیت

در اینجا، G تعداد کل داده های اقلیت برای تولید است. β نسبت داده های اقلیت: اکثریت مورد نظر پس از ADASYN است. $\beta = 1$ به معنای مجموعه داده های کاملاً متعادل پس از ADASYN است.

مرحله ۳

k -نزدیکترین همسایگان هر یک از نمونه های اقلیت را پیدا کنید و مقدار r_i را محاسبه کنید. پس از این مرحله، هر نمونه اقلیت باید با یک محله متفاوت مرتبط شود.

$$r_i = \frac{\#majority}{k}$$

شکل ۲۸ محاسبه r

مقدار r_i نشان دهنده تسلط طبقه اکثریت در هر محله خاص است. محله های r_i بالاتر حاوی نمونه های کلاس اکثریت بیشتری هستند و یادگیری آنها دشوارتر است. برای تجسم این مرحله به زیر مراجعه کنید. در مثال، $K = 5$ (به دنبال ۵ همسایه نزدیک می گردد).

مرحله ۴

مقادیر r_i را عادی کنید تا مجموع همه مقادیر r_i برابر با ۱ باشد.

$$\hat{r}_i = \frac{r_i}{\sum r_i}$$

$$\sum \hat{r}_i = 1$$

شکل ۲۹ محاسبه شعاع

این مرحله عمدتاً مقدمه ای برای آسان کردن مرحله ۵ است.

مرحله ۵

مقدار نمونه های مصنوعی برای تولید در هر محله را محاسبه کنید.

$$G_i = G\hat{r}_i$$

شکل ۳۰ تعداد نمونه ی تصادفی

از آنجایی که r_i برای محله هایی که نمونه های طبقه اکثریت بر آنها غالب است، بالاتر است، نمونه های کلاس اقلیت مصنوعی بیشتری برای آن محله ها تولید می شود. از این رو، این به الگوریتم ADASYN ماهیت تطبیقی آن می دهد. داده های بیشتری برای محله های "سخت تر برای یادگیری" تولید می شود.

مرحله ۶

داده های G_i را برای هر مرحله ایجاد کنید. ابتدا مثال اقلیت را برای محله، X_i در نظر بگیرید. سپس، به طور تصادفی نمونه اقلیت دیگری را در آن محله، XZ_i انتخاب کنید. مثال مصنوعی جدید را می توان با استفاده از:

$$s_i = x_i + (x_{zi} - x_i)\lambda$$

شکل ۳۱ تولید نمونه تصادفی

در معادله بالا، λ یک عدد تصادفی بین ۰-۱ است، s_i نمونه مصنوعی جدید است، x_i و x_{zi} دو نمونه اقلیت در یک همسایگی هستند. تصویری از این مرحله در زیر ارائه شده است. به طور شهودی، نمونه های مصنوعی بر اساس ترکیب خطی x_i و x_{zi} ایجاد می شوند.

در کد یک کلاس به نام ADASYN` را تعریف می کنیم که الگوریتم ADASYN را برای افزایش نمونه ها در زمینه مسائل طبقه بندی نامتوازن پیاده سازی می کند. در زیر توضیحی از قسمت های اصلی کد آمده است:

در پیاده سازی این بخش داریم:

یک کلاس پایتون به نام ADASYN(Adaptive Synthetic Sampling) را تعریف می کند که مخفف Adaptive Synthetic Sampling است. الگوریتم ADASYN یک تکنیک نمونه برداری بیش از حد است که برای مقابله با عدم تعادل کلاس در مجموعه داده های یادگیری ماشین استفاده می شود. بر روی تولید نمونه های مصنوعی برای طبقه اقلیت بر اساس نقاط داده ای که یادگیری آنها دشوار است تمرکز دارد.

در اینجا جزئیاتی از اجزای اصلی کلاس ADASYN آورده شده است:

سازنده (__init__): این متد شی ADASYN را با چندین پارامتر مقداردهی اولیه می کند:

نسبت: نسبت مورد نظر تعداد نمونه در کلاس اقلیت پس از نمونه گیری مجدد به تعداد نمونه در کلاس اکثریت.

imb_threshold: آستانه ای برای تعیین اینکه آیا یک کلاس نامتعادل در نظر گرفته می شود.

k: تعداد نزدیک ترین همسایه ها برای استفاده در هنگام تولید نمونه های مصنوعی.

random_state: دانه ای برای مولد اعداد تصادفی که امکان تکرارپذیری را فراهم می کند.

verbose: پرچمی برای فعال یا غیرفعال کردن چاپ اطلاعات در حین پردازش.

fit: این روش ویژگی های X و هدف y را می پذیرد و اقدامات زیر را انجام می دهد:

آرایه های ورودی را بررسی و قالب بندی می کند.

یک مولد اعداد تصادفی را با دانه ارائه شده راه اندازی می کند.

طبقه اکثریت را شناسایی می کند و آمار کلاس را محاسبه می کند.

transform (تبدیل): این روش به طور کامل اجرا نمی شود زیرا X و y را نمی پذیرد و چیزی را بر نمی گرداند. به نظر می رسد در نظر گرفته شده است که نمونه برداری بیش از حد در مجموعه داده اعمال شود، اما ناقص است.

fit_transform: این روش مراحل تناسب و تبدیل را به یک مرحله ترکیب می کند. ابتدا نمونه **ADASYN** را با داده ها مطابقت می دهد و سپس نمونه برداری بیش از حد را انجام می دهد. سپس داده های بیش نمونه شده با داده های اصلی الحاق می شوند تا یک مجموعه داده جدید با توزیع کلاس متعادل ایجاد شود.

Generate_samples: این روش نمونه های مصنوعی را برای یک کلاس اقلیت معین تولید می کند. برای هر نقطه داده کلاس اقلیت، k نزدیکترین همسایه را پیدا می کند و با درون یابی بین نقطه داده و همسایگان آن، نمونه های جدیدی تولید می کند.

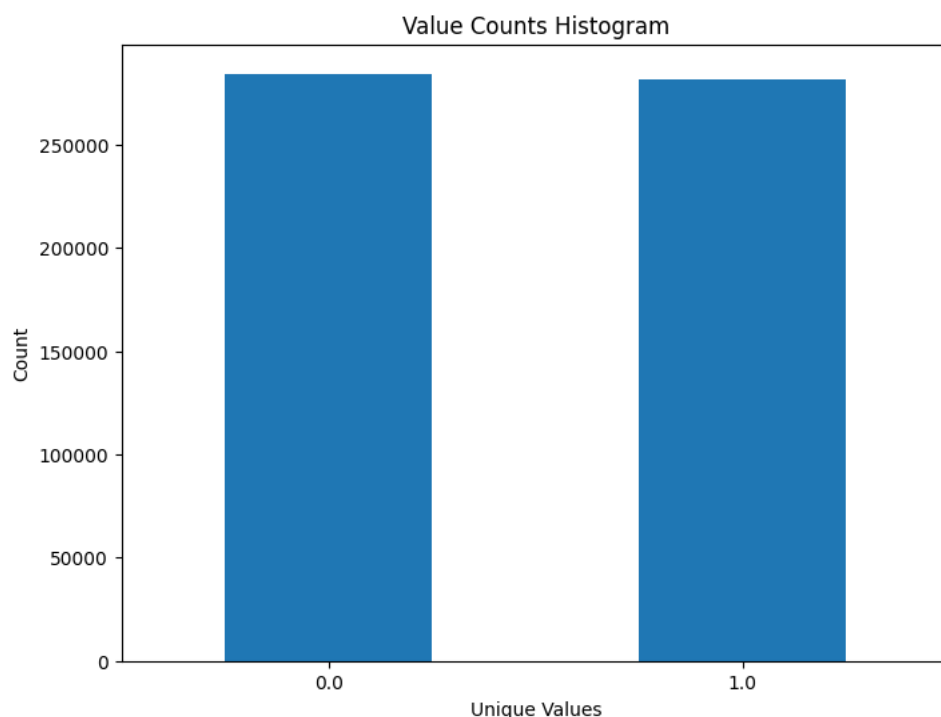
oversample: این روش نمونه برداری واقعی را انجام می دهد.

۴. زمان انجام نمونه برداری

نمونه برداری معمولاً باید قبل از تقسیم کردن داده ها به دسته های آموزش و تست انجام شود. این کار مهم است زیرا انجام نمونه برداری بعد از تقسیم بندی می تواند منجر به نشت اطلاعات بین دسته های آموزش و تست شود و در نتیجه ارزیابی نادرستی از عملکرد مدل را به دنبال داشته باشد.

۴. تحلیل داده های نمونه برداری شده

پس از انجام نمونه برداری، می توانید هیستوگرام هایی را برای مقایسه توزیع کلاس ها قبل و بعد از نمونه برداری ترسیم به صورت زیر می باشد.



شکل ۳۲: هیستوگرام بعد از نمونه برداری

۳-۴. آموزش مدل با داده های جدید

طبق پیاده سازی مقاله

حال برای پیاده سازی مدل عصبی با استفاده از TensorFlow و Keras ایجاد می کند که برای یادگیری از داده های دوره ای D1 طراحی شده است. این مدل شامل یک سری از لایه های مختلف است. دستورات مربوط به هر لایه به شرح زیر است:

۱. Conv1D Layer: این لایه یک لایه کانولوشنال D1 با ۳۲ فیلتر، هر کدام با یک کرنل به طول ۲ و تابع فعال سازی ReLU دارد. این لایه برای استخراج ویژگی های مختلف از داده های ورودی با ساختار $(X_train.shape[1], X_train.shape[2])$ استفاده می شود.

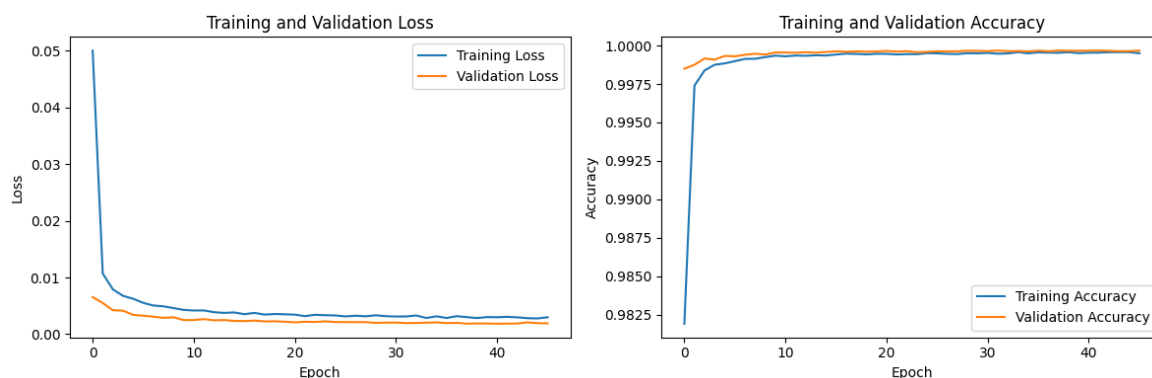
۲. BatchNormalization: این لایه برای استانداردسازی میانگین و واریانس خروجی های لایه قبلی (Conv1D Layer ۱) استفاده می شود. این کمک می کند تا آموزش مدل بهبود یابد و از مشکلات مربوط به مقیاس پذیری و واریانس داده جلوگیری شود.

۳. MaxPool1D: این لایه از عملگر MaxPooling برای کاهش ابعاد خروجی لایه قبلی استفاده می کند. این باعث کاهش تعداد ویژگی ها و افزایش مقیاس پذیری شبکه می شود.

۴. ۱ Dropout: این لایه برای جلوگیری از اورفیتینگ (Overfitting) در مدل استفاده می‌شود. با احتمال ۰.۲، تعدادی از واحدهای خروجی لایه قبلی را به صورت تصادفی خاموش می‌کند.
۵. ۲ Conv۱D Layer: این لایه یک لایه کانولوشنال دیگر با ۶۴ فیلتر و کرنل به طول ۲ و تابع فعال‌سازی ReLU است.
۶. ۲ BatchNormalization: مانند ۱ BatchNormalization، این لایه برای استانداردسازی میانگین و واریانس خروجی‌های لایه قبلی (۲ Conv۱D Layer) استفاده می‌شود.
۷. ۲ MaxPool۱D: مانند ۱ MaxPool۱D، این لایه از عملگر MaxPooling برای کاهش ابعاد خروجی لایه قبلی استفاده می‌کند.
۸. ۲ Dropout: این لایه برای جلوگیری از اورفیتینگ با احتمال ۰.۵ استفاده می‌شود. این باعث می‌شود تا تعداد بیشتری از واحدها در طی آموزش خاموش شوند و مدل بهبود پذیر شود.
۹. Flatten: این لایه به داده‌های خروجی لایه قبلی (۲ MaxPool۱D) اعمال می‌شود و آنها را به یک بردار یک بعدی تبدیل می‌کند.
۱۰. ۱ Dense: این لایه با ۶۴ واحد و تابع فعال‌سازی ReLU برای اعمال یادگیری ژرف بر روی ویژگی‌های استخراج شده توسط لایه‌های قبلی استفاده می‌شود.
۱۱. ۳ Dropout: این لایه برای جلوگیری از اورفیتینگ با احتمال ۰.۵ استفاده می‌شود.
۱۲. ۲ Dense: این لایه با ۶۴ واحد و تابع فعال‌سازی ReLU استفاده می‌شود.
۱۳. ۳ Dense: این لایه نهایی با یک واحد و تابع فعال‌سازی sigmoid برای انجام مسئله دسته‌بندی دودویی (binary classification) استفاده می‌شود.

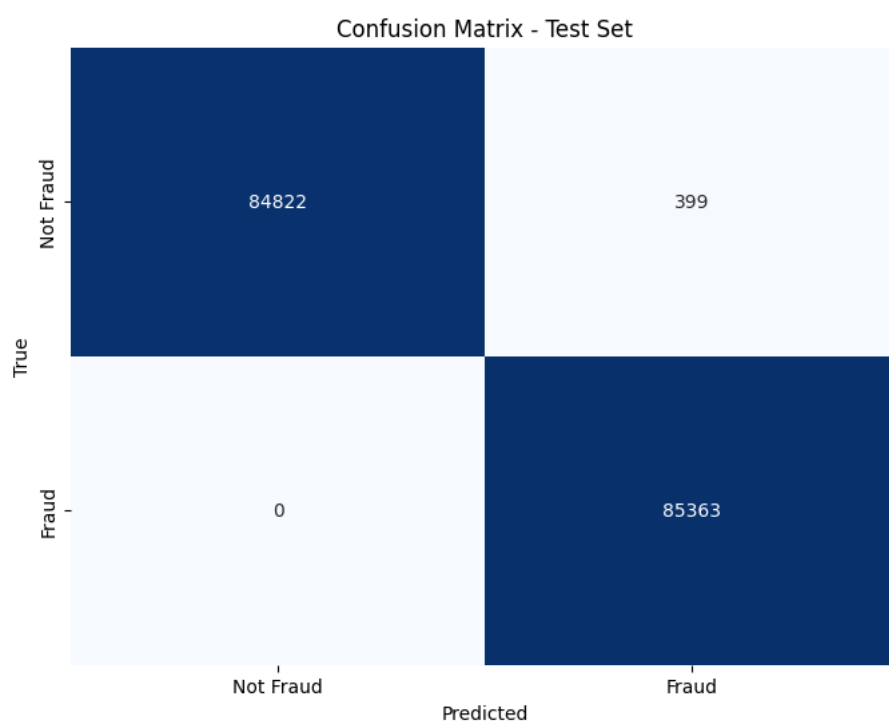
به این ترتیب، مدل نهایی دارای یک لایه ورودی (۱ Conv۱D Layer)، یک سری از لایه‌های میانی (BatchNormalization، MaxPooling، Dropout، ۲ Conv۱D Layer، ۲ BatchNormalization، ۲ MaxPool۱D، ۲ Dropout، Flatten، ۱ Dense، Dropout، ۳، ۲ Dense) و یک لایه خروجی (۳ Dense) است. این مدل مناسب برای مسائل دسته‌بندی دودویی با داده‌های دوره‌ای D۱ می‌باشد.

حال نتایج را به صورت زیر به دست می آوریم:



شکل ۳۳: نتایج با نمونه برداری

خطای هر هم داده‌های تست هم داده‌های ارزیابی در حال کم شدن است و تا حد زیادی مدل **overfit** نشده است و برای داده‌های خارج از آموزش نیز دارای عمومیت می‌باشد.



شکل ۳۴: ماتریس در هم ریختگی با نمونه برداری

Classification Report:				
	precision	recall	f1-score	support
0	1.00	1.00	1.00	85221
1	1.00	1.00	1.00	85363
accuracy			1.00	170584
macro avg	1.00	1.00	1.00	170584
weighted avg	1.00	1.00	1.00	170584

شکل ۳۵: نتایج مربوط به نمونه برداری

همچنین عملکرد این مدل برای معیارهایی غیر دقت به شرح زیر می‌باشد:

```
Precision: 0.9997390613324477
Recall: 0.9995138147752876
F1-score: 0.9996264253651248
```

شکل ۳۶ معیارهای صحت و **recall** با نمونه‌برداری

مشخص است که این مدل هم برای کلاس اقلیت هم برای کلاس اکثریت به بهترین نحو عمل کرده است و نوع توزیع هر دو کلاس را به درستی یاد گرفته است.