

به نام خدا



دانشگاه تهران



دانشکده مهندسی برق و کامپیوتر

درس شبکه‌های عصبی و یادگیری عمیق

تمرین پنجم

610399205

امیرعباس رضا سلطانی

610399199

نیما نیرومند

پرسش 1. تشخیص اخبار جعلی مبتنی بر مدل‌های ترنسفورمر

۱-۱. آشنایی با BERT و CT-BERT

در این سوال از روش یادگیری انتقالی استفاده می‌کنیم. در این روش می‌توانیم از ذخیره و انتقال اطلاعات بدست آمده از انجام وظیفه خاص برای انجام وظیفه جدید استفاده کنیم. از این روش زمانی استفاده می‌کنیم که داده‌های کمی برای یادگیری و ارزیابی مدل موجود است.

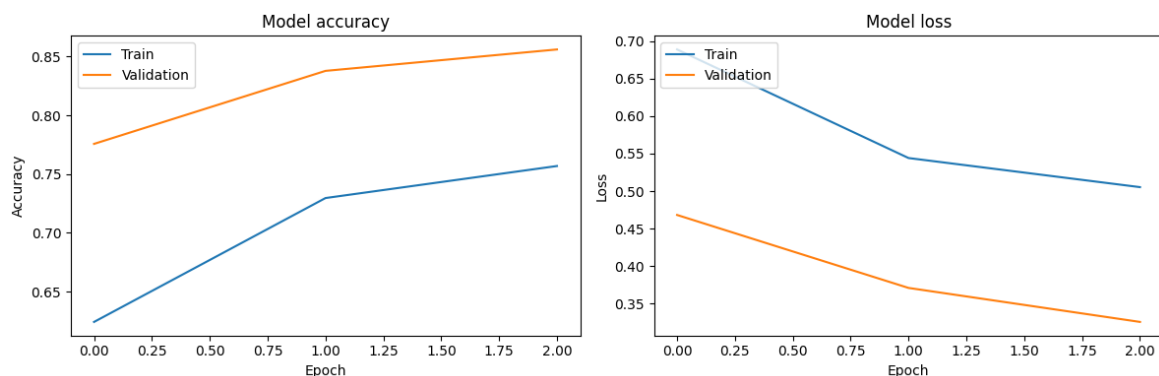
در پردازش زبان طبیعی، مدلی را روی مقدار زیادی داده متنی به صورت self supervised آموزش می‌دهند (pretraining). از خروجی این مدل که فهم contextual از ورودی دارد، با دو روش برای عمل یادگیری انتقالی برای مساله downstream استفاده می‌شود. روش Fine Tuning و روش Feature Based. در روش Feature Based، از نگاشت ویژگی‌های بدست آمده، به عنوان ورودی مدل دیگر (downstream) برای تصمیم‌گیری (طبقه‌بندی یا رگرسیون) استفاده می‌شود و صرفاً مدل دوم برای وظیفه آموزش داده می‌شود و وزن‌های مدل اول ثابت است.

اما در روش Fine Tuning، پس از آموزش مدل اول روی داده‌های زیاد، در مرحله یادگیری برای وظیفه downstream، وزن‌های شبکه اول (feature extractor) نیز در کنار وزن‌های شبکه دوم آموزش داده می‌شود. فایده روش دوم، فهم عمیق‌تر مدل نسبت به مفهوم عبارات به کار رفته در دیتاست است.

1-2. پیاده‌سازی مدل با رویکرد FineTuning

داده‌های ورودی مربوط به اخبار شایع در توییتر در مورد کوئید می‌باشد. هدف آموزش شبکه‌ای است که توانایی تشخیص درست یا غلط بودن خبر را دارد.

- در مدل اول، از شبکه BERT برای Word Embedding استفاده می‌کنیم. این مدل دارای 12 لایه، 12 attention head و 110 میلیون پارامتر است و روی متن‌های انگلیسی با 2500 میلیون توکن پیش‌آموزش دیده است. ابتدا جملات را به صورت رشته‌های 128 تایی tokenize می‌کنیم. خروجی شبکه، ماتریسی است که عضو اول آن (CLS) طول 768 تایی دارد. این رشته دارای context از کل عبارت ورودی است. خروجی را به عنوان ورودی به لایه fully connected با یک نود می‌دهیم و برای آن از تابع فعالساز sigmoid استفاده می‌کنیم و خروجی آن را به عنوان احتمال تعلق خبر به دسته واقعی تعبیر می‌کنیم.



تصویر ۱. نتایج مدل اول Feature Base

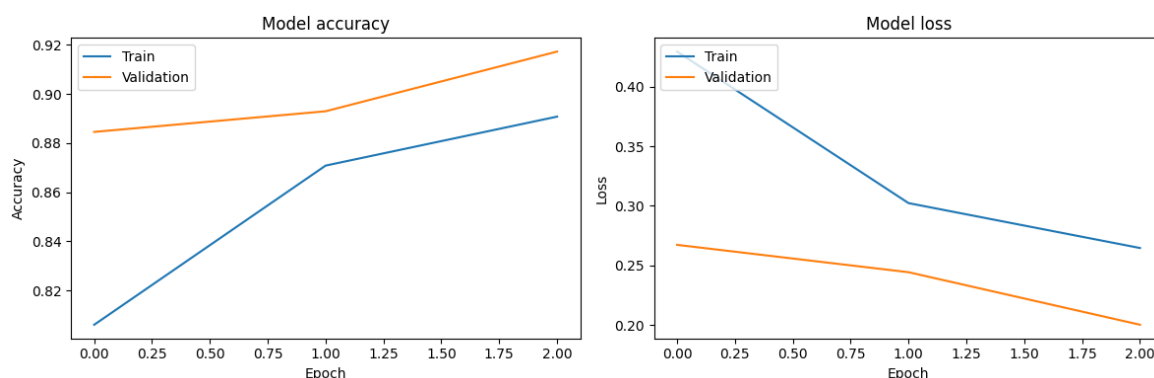
جدول ۱. نتایج مدل اول در رویکرد Feature Base

F1	loss	دقت
0.87	0.33	0.86

جدول ۲. درهم‌ریختگی مدل اول Feature Base

	1 Predicted	0 Predicted
1 Label	762	258
0 Label	52	1068

- در مدل دوم، شبیه مدل اول عمل می‌کنیم با این تفاوت که خروجی BERT را به صورت ماتریسی با طول 128 تایی در نظر می‌گیریم. این ماتریس قبل از ورود به لایه fully connected، از یک لایه GRU Bidirectional عبور می‌کند. خروجی آن وارد لایه fully connected با یک نود و تابع فعالساز sigmoid برای عمل طبقه‌بندی می‌شود.



تصویر ۲. نتایج مدل دوم Feature Base

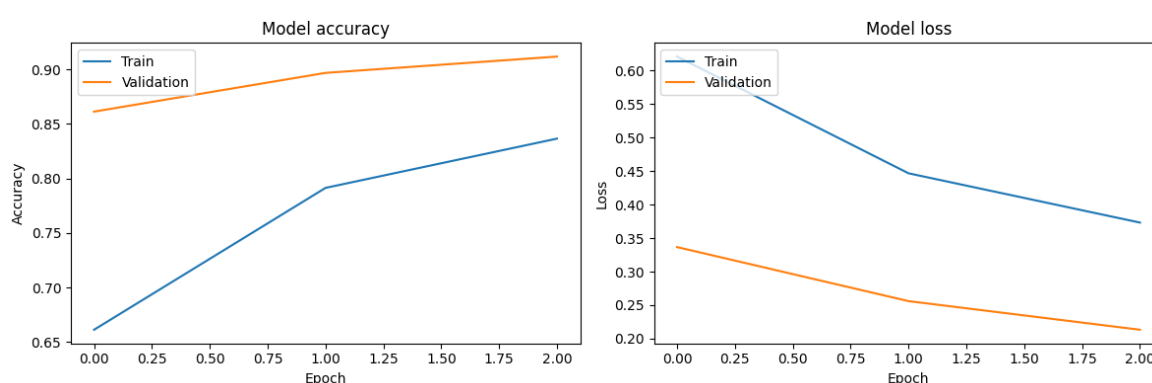
جدول ۳. نتایج مدل دوم در رویکرد Feature Base

دقت	loss	F1
0.92	0.21	0.92

جدول ۴. درهم‌ریختگی مدل دوم Feature Base

	1 Predicted	0 Predicted
1 Label	884	136
0 Label	39	1081

- مدل سوم شبیه مدل دوم است. با این تفاوت که BERT را پس از پیش پردازش روی 2500 میلیون توکن عمومی روی داده‌های با مفهوم کویید پیش پردازش می‌دهیم تا درک عمیق‌تری نسبت به اخبار ورودی داشته باشد. این مدل CT-BERT می‌نامیم.



تصویر ۳. نتایج مدل سوم Feature Base

جدول ۵. نتایج مدل دوم در رویکرد Feature Base

دقت	loss	F1
0.91	0.21	0.92

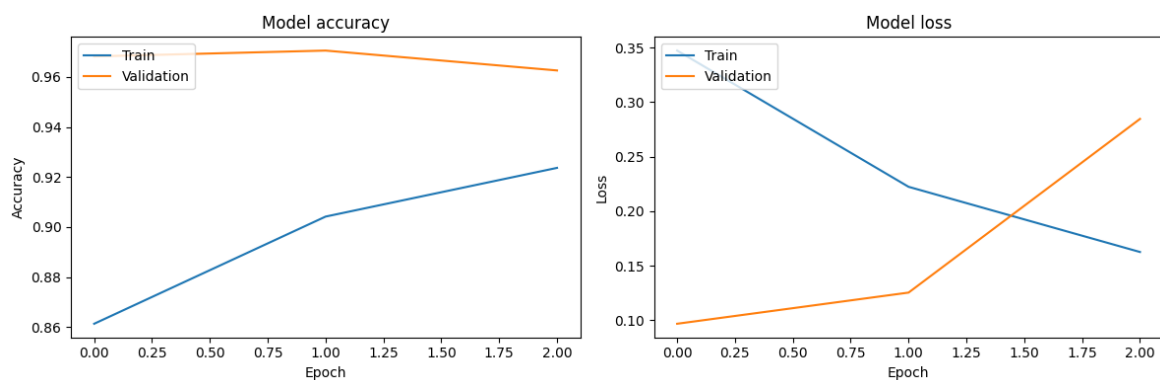
جدول ۶. درهم‌ریختگی مدل سوم Feature Base

	1 Predicted	0 Predicted
1 Label	893	127

0 Label	63	1057
---------	----	------

3-1. پیاده سازی مدل با رویکرد Feature-Based

- مدل اول:



تصویر ۴. نتایج مدل اول Fine Tuning

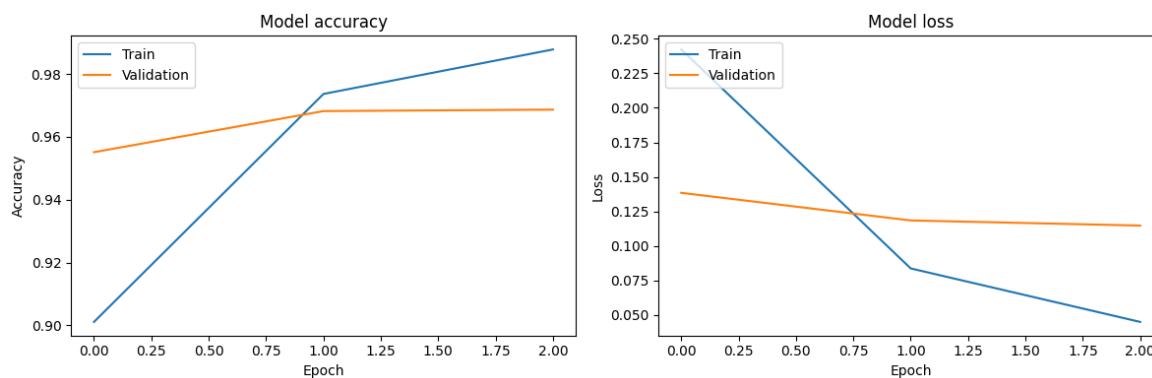
جدول ۷. نتایج مدل اول در رویکرد Fine Tuning

F1	loss	دقت
0.96	0.31	0.95

جدول ۸. درهم‌ریختگی مدل اول Feature Base

	1 Predicted	0 Predicted
1 Label	968	52
0 Label	17	1103

- مدل دوم:



تصویر ۵. نتایج مدل دوم Fine Tuning

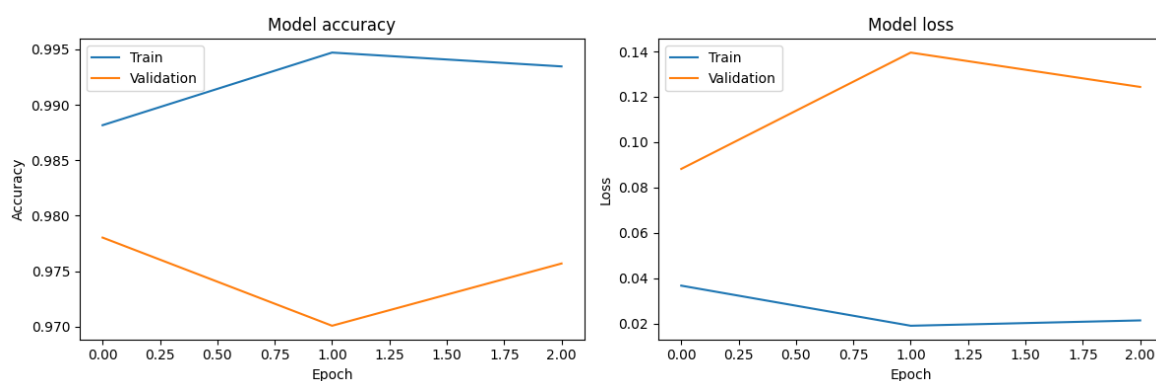
جدول ۱۰. نتایج مدل دوم در رویکرد Fine Tuning

F1	loss	دقت
0.97	0.11	0.96

جدول ۱۱. درهم‌ریختگی مدل دوم Feature Base

	1 Predicted	0 Predicted
1 Label	925	95
0 Label	3	1117

• مدل سوم:



تصویر ۶. نتایج مدل سوم Fine Tuning

جدول ۱۲. نتایج مدل سوم در رویکرد Fine Tuning

F1	loss	دقت
0.98	0.14	0.97

جدول ۱۳. درهم‌ریختگی مدل سوم Fine Tuning

	1 Predicted	0 Predicted
1 Label	975	45
0 Label	12	1108

4-1. تحلیل نتایج

در هر رویکرد یادگیری انتقالی، مدل سوم با استفاده از CT-BERT و مدل دوم با استفاده از BERT دقت تقریباً مشابه و بهتر از مدل دیگر کسب کردند. همچنین F1-Score آن‌ها نیز تقریباً برابر است. در رویکرد Fine Tuning، مدل سوم با استفاده از CT-BERT، دقت بهتری از دو مدل دیگر (که از BERT استفاده کردند) کسب کرده است. همچنین F1-Score بالاتری نسبت به دو مدل دیگر دارد. بنابراین می‌توان نتیجه گرفت که استفاده از مدل BERT با پیش‌پردازش بر روی داده‌های مرتبط با داده‌های ورودی برای وظیفه خاص (آموزش مدل CT-BERT) و سپس آموزش آن برای انجام یک وظیفه خاص، نتایج بهتری دارد. زیرا درک بهتر (جدایی پذیری بیشتر برای مفاهیم دورتر و نزدیکی بهتر برای مفاهیم های نزدیک تر) از جملات مربوط به کوئید دارد و با آموزش مجدد روی داده‌های مربوط به یک وظیفه خاص مربوط به کوئید، جدایی پذیری بالاتری در خروجی ایجاد می‌کند و عمل پیش‌بینی برای بخش‌های بعدی شبکه ساده‌تر می‌شود.

در رویکرد Feature-Based، مدل دوم (BERT+GRU+FC) دقت و F1-Score بهتری از مدل اول (BERT+FC) دارد. علت آن استفاده از GRU است که قدرت حافظه را به شبکه اضافه می‌کند و به عنوان یک لایه feature extractor عمل می‌کند. GRU الگوهایی که در توالی کلمات خروجی از BERT است را درک کرده و به لایه بعدی می‌دهد. این درک، قابلیت جدایی پذیری بالاتری به داده‌های دو کلاس می‌دهد (داده‌های هم کلاس نزدیک‌تر و داده‌ها از کلاس‌های متفاوت دورتر).

هر مدل، در رویکرد FineTuning نتایج بهتری نسبت به مدل متناظر در FeatureBase گرفته است. علت آن، افزایش اطلاعات لایه Embedding در مورد متن‌های ورودی مساله با توجه به هدف مساله است. در واقع لایه Embedding در کنار اطلاعات قبلی پیش‌آموزش دیده، اطلاعات جدیدی در مورد متون یاد می‌گیرد که برای تشخیص اخبار در لایه‌های بعدی مفید خواهد بود.

جدول ۱۴. نمونه تشخیص های غلط

مدل	پیش بینی	برچسب	خبر
اول	real	fake	An audio file by an alleged worker at a health institution in Rio de Janeiro. She says that healthcare workers on public institutions in Rio are forced to state whether a patient has COVID-19 or not even before he sees a doctor. This was allegedly being done to artificially inflate the number of cases
اول	real	fake	The Chinese government announced that "garlic is a ".preventative food for the the novel coronavirus
دوم	fake	real	On the 15/03 NCDC directly contacted a Twitter user who mentioned his friend who returned from UK had runny nose but could not reach authorities for testing. Within 12 hours of communication with us via DM a sample was collected. We're committed to doing our best https://t.co/fccdGij3uG
دوم	fake	real	IndiaFightsCorona #COVID19 Recoveries exceed active # cases by more than 2.1 million. There has been more than 4 times jump in the average weekly recoveries from the first week of July to last week of August. https://t.co/HRyyaWhdJR
سوم	real	fake	RT @factchecknet: In the absence of clarity and a rash of #misinformation the #COVID19 pandemic has created a ...breeding ground for #prejudi
سوم	fake	real	An even better piece of news: states reported fewer than 500 deaths. That hadn't happened since March. Yes there will probably be a larger number of deaths reported tomorrow as lagging weekend data gets posted but it is a significant pandemic milestone. https://t.co/Dflu3I5agZ

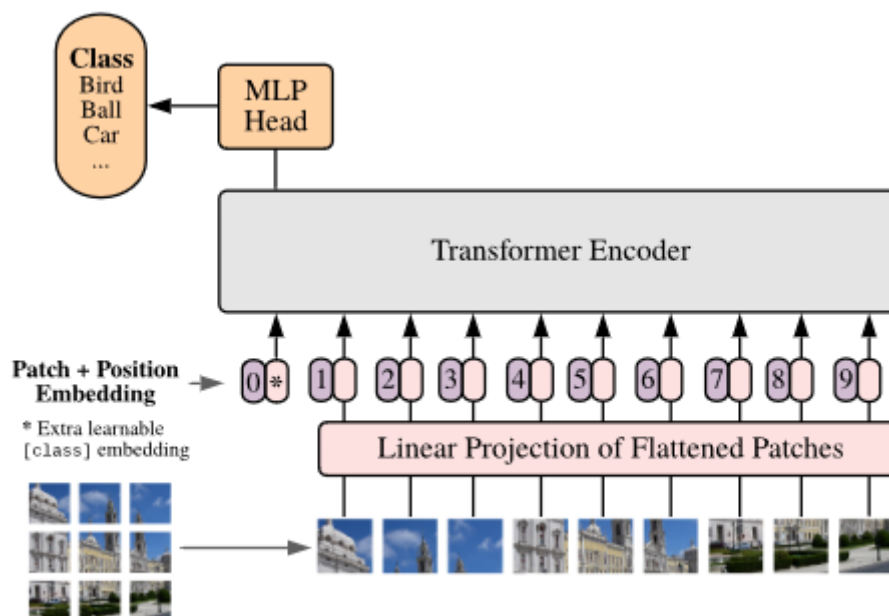
این طور که پیداست، می توان حدس زد که اگر خبر از طرف منبعی بیان شده باشد و آن منبع به صورت "some resource says:" یا "A news:" در ابتدای خبر بیان شود، مدل ها آن را درست تشخیص می دهد. همچنین در صورتی که لینک خبر در انتها ظاهر شود، خبر اشتباه تشخیص داده می شود.

پرسش ۲ - به کارگیری مدل‌های ترنسفرمر در طبقه بندی تصاویر

۱-۲. آشنایی با ترنسفرمرهای تصویر

الف) ViT یا همون vision transformer یک مدل است که برای پردازش تصویر استفاده می‌شود و به جای CNN ها سعی می‌شود از لایه‌های transformer استفاده کند. از آنجایی که در ابتدا ترنسفرمرها را برای داده‌های متنی تعریف کردیم و به صورت دنباله‌ای آن را توکن و سپس تعبیه (embedd) می‌کردیم حال پس برای عکس‌ها باید دنبال توکن کردن باشیم که همچنین مفهوم دنباله‌ای را نیز داشته باشد و برای این کار مثلا می‌توانیم فرض کنیم که قطعات کوچک عکس را به قطعات کوچک (مثلا ۱۶ در ۱۶ تقسیم می‌کنیم). و همانطور که گفتیم در ادامه این قطعات را به یک بردار تبدیل می‌کنیم و سپس مانند ترسفر برای متن‌ها موقعیت مکانی قطعات را نیز لحاظ می‌کنیم و سپس این بردارها را به شبکه ترسفرمر می‌دهیم و دیگر تغییر خاصی نسبت به حالت متنی نداریم.

(ب)



شکل ۱. ساختار ViT

Patch Embedding : تصاویر را به قطعه‌های کوچک‌تر تقسیم می‌کند و سپس هر قطعه را به یک

بردار تبدیل می‌کند.

Position Embedding: اطلاعات موقعیت مکانی را به بردارهای قطعات اضافه می‌کند.

Linear Projection of Flattened Patches: قطعات تصویر به صورت مسطح تبدیل می‌کند و به یک

بردار ویژگی تبدیل می‌شوند.

Transformer Encoder: شامل چندین لایه ترنسفرمر است.

MLP Head: در نهایت از ویژگی‌های استخراج شده چندین لایه فولی کانکت می‌زنیم تا ویژگی‌ها

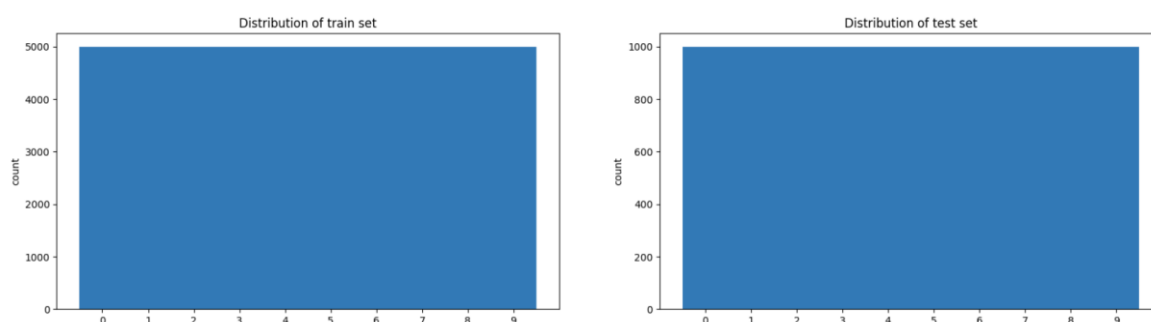
استخراج شوند.

ج) از ایرادات می‌توان این را اشاره کرد که CNN ها موقعیت مکانی هم در نظر می‌گرفتند اما در این مدل تا حدی این اتفاق نمی‌افتد و داده‌ها چون به صورت دنباله‌ای داده می‌شوند سخت می‌شود این را لحاظ کرد همچنین از نظر زمانی نیز طولانی‌تر و به داده‌های بیشتری نیاز دارند.

پس با توجه به این اشکالات احتمالا ایده‌ای که بتونه یکی کمک کننده باشد این است که به صورت ترکیبی از هردوی CNN و ViT بهره ببریم تا از خوبی‌ها هر کدام بتونیم استفاده کنیم.

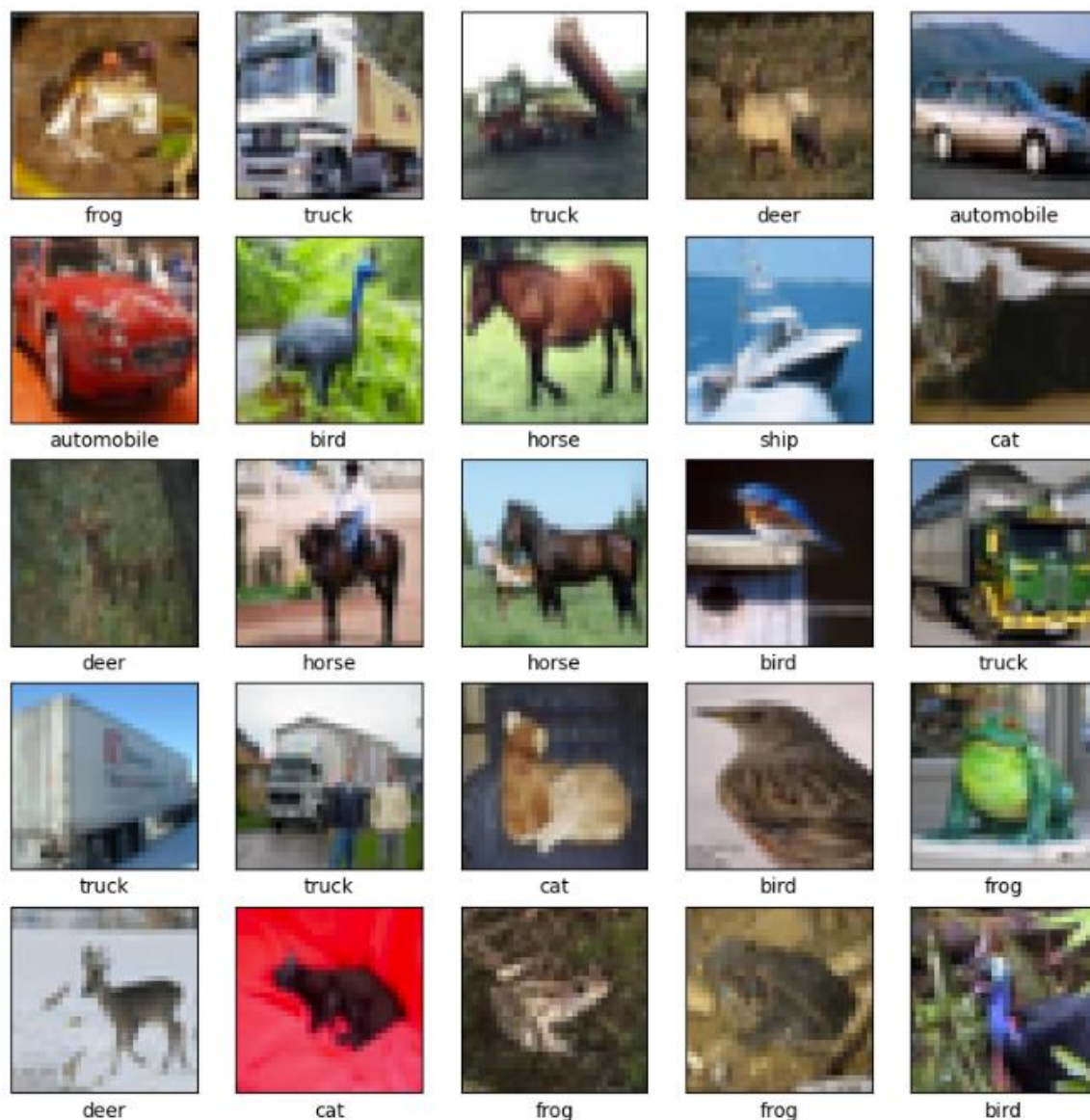
۲-۲. لود و پیش پردازش دیتاست

در این قسمت از مجموعه داده‌ی آماده CFAR10 استفاده می‌کنیم که شامل ۵۰۰۰ عکس برای مجموعه آموزشی و ۱۰۰۰ عکس برای مجموعه آزمایشی است همچنین عکسها ۳۲ در ۳۲ و رنگی (۳*۳۲*۳۲) هستند که دارای ۱۰ برچسب هستند و در ۱۰ کلاس طبقه‌بندی شده‌اند حال در هردو داده آموزشی و تست توزیع کلاس‌ها را بررسی می‌کنیم.



شکل ۱. توزیع آماری دسته‌ها در مجموعه آموزش و آزمایش

و همانطور که می‌بینیم در هردوی حالت آزمایش و تست داده‌ها به صورت یکنواخت و یکسان تقسیم شده‌اند (در آموزش از هر کلاس ۵۰۰ داده و در تست از هر داده ۱۰۰ تا) پس توزیع داده‌ها متوازن است. حال برای آشنایی بیشتر با مجموعه داده‌ها تصویر تعدادی را خروجی می‌دهیم. (توجه شود که کلاس‌های ما برچسب ۰ تا ۹ دارند و در عکسها به صورت دستی نام لیبل واقعی را نوشتیم).

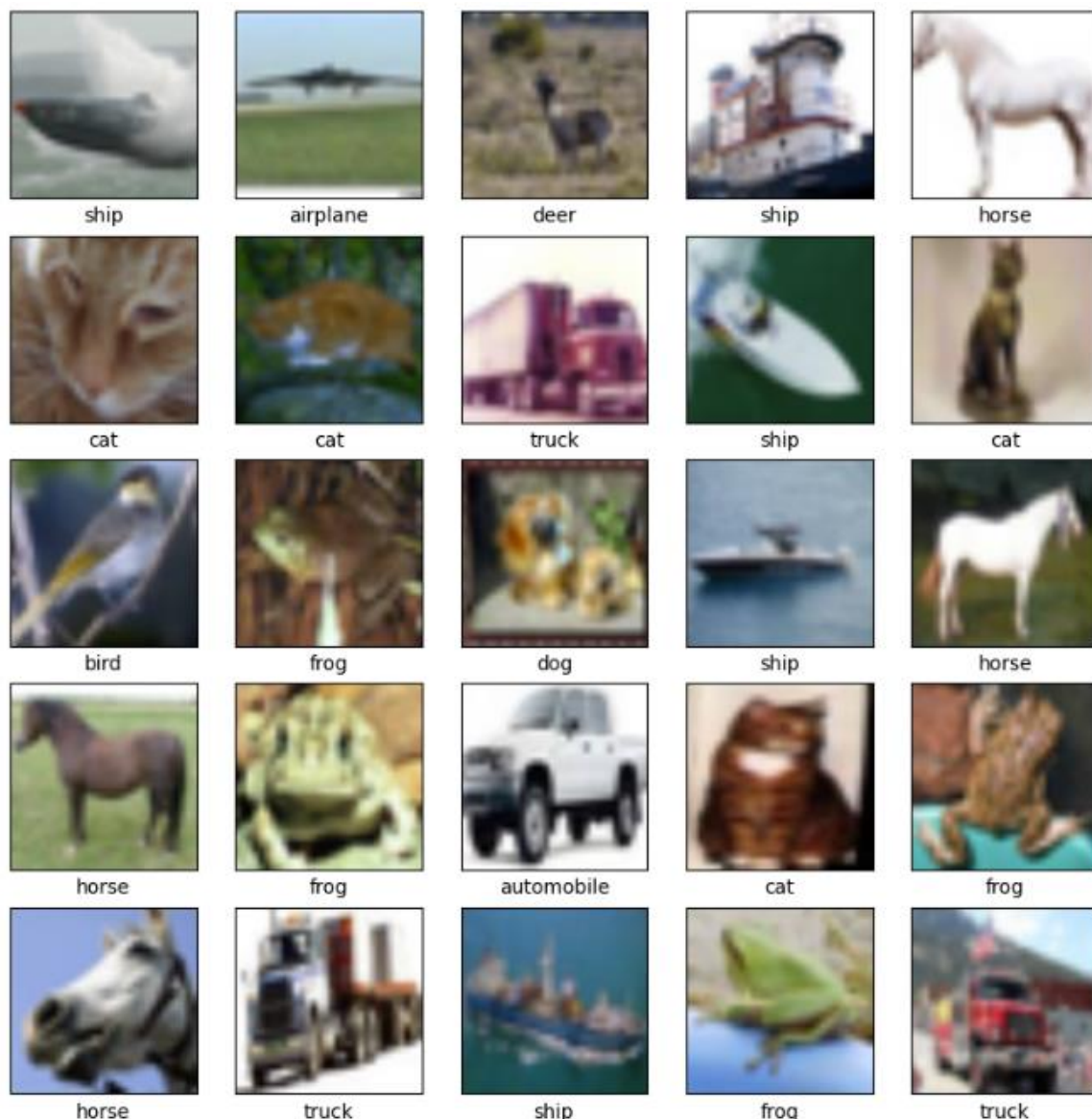


شکل ۲. نمونه‌ای از تصاویر دیتاست قبل از پیش پردازش

حال از آنجایی که مقادیر درایه‌های داخل آرایه عکسهای ما ۰ تا ۲۵۵ هستند برای اینکه نرمال کنیم و اعداد را در بازه ۰ تا ۱ داشته باشیم از min-max normalization استفاده می‌کنیم و همه مقادیر آرایه‌ها را بر ۲۵۵ تقسیم می‌کنیم. حال از آنجایی که transfer learning استفاده می‌کنیم و از مدل‌های آموزش دیده‌ای

که استفاده می‌کنیم ورودیشون با سایز ۲۲۴ در ۲۲۴ است پس سایز عکس‌هامون را تغییر می‌دهیم و با استفاده از روش bilinear آن‌ها را تغییر سایز می‌دهیم.

حال تعدادی از عکس‌ها را بعد از این تغییرات خروجی می‌دهیم



شکل ۳. نمونه‌ای از تصاویر دیتاست قبل از پیش پردازش

۲-۳. fine-tuning شبکه کانولوشنی

الف) از مدل پیش آموخته شده VGG19 استفاده می‌کنیم که بر روی مجموعه Imagenet1k

یادگیری شده حال قسمت فولی کانکت را حذف می‌کنیم و می‌خواهیم از ویژگی‌های استخراج شده توسط این مدل را استفاده کنیم. حال برای اینکه بعضی از لایه‌های مدل را دوباره آموزش بدیم طبق مقاله لایه‌های block5_conv1 به بعد را unfreeze می‌کنیم تا دوباره مدل را آموزش دهیم و fine-tuning کنیم حال برای تسک خودمان طبق گفته مقاله ابتدا لایه‌ها فلت می‌کنیم تا سائز داده‌ها یک بعدی شه حال از بردارهای به دست آمده را به یک لایه فولی کانکت ۲۵۶ با نورون وصل می‌کنیم. حال برای جلوگیری از بیش بردازش Dropout با مقدار احتمال ۵۰ درصد اضافه کنم حال برای خروجی از آنجایی که ۱۰ کلاس داریم یک لایه با ۱۰ نورون و با Softmax به عنوان تابع فعال‌ساز می‌گذاریم.

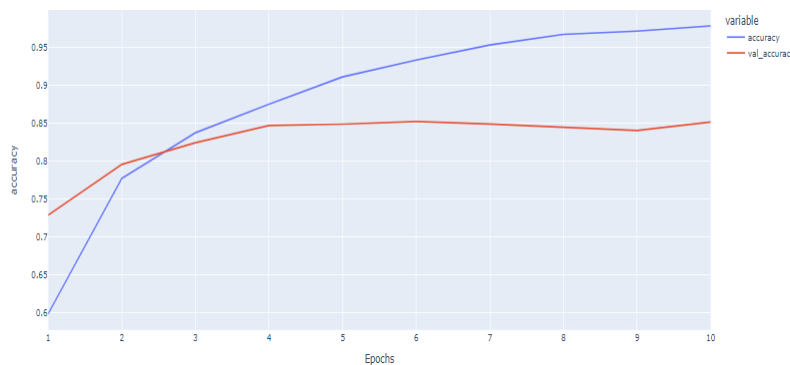
ب) تعداد پارامترهای trainable ۹۵۷۳۱۳۰، است از آنجایی که خروجی مدل VGG19 بعد از فلت کردن یک بردار ۵۱۲ تاست و به یک لایه ۲۵۶ نورونی وصل است پس $۱۳۱۳۲۸ = (۱+۵۱۲)*۲۵۶$ و بعد این ۲۵۶ را به یک لایه ۱۰ تایی وصل می‌کنیم و $۲۵۷۰ = (۱+۲۵۶)*۱۰$ و باقی مربوط به لایه‌های آخر unfreeze شده مدل VGG19 است.

ج) حال مدل را با بهینه‌ساز آدام و لرنینگ ریت ۰.۰۰۰۱ و تابع loss categorical_crossentropy و با batchsize ۶۴، را به تعداد ۱۰ اپیک ران می‌کنیم.

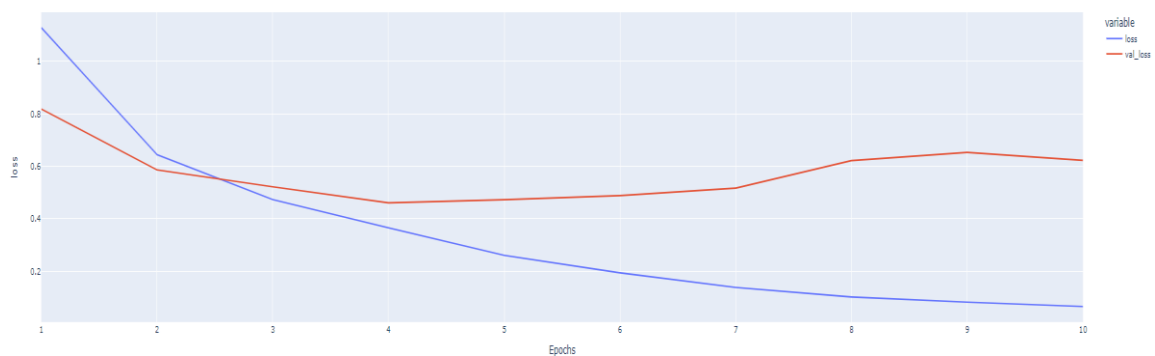
د) در پایان اپیک دهم دقت و مقدار تابع هزینه روی داده‌های آموزشی و ارزیابی به شکل زیر می‌شوند

دقت	
۹۷.۸	داده آموزشی
۸۵.۲	داده ارزیابی
تابع هزینه	
۰.۰۶	داده آموزشی
۰.۶۲	داده ارزیابی

همچنین نمودار تغییراتشان را نیز در زیر نشان دادیم



شکل ۵. نمودار تغییرات دقت در داده‌های آموزشی و ارزیابی



شکل ۶. نمودار تغییرات تابع هزینه در داده‌های آموزشی و ارزیابی

(۵)

میانگین زمان (بر حسب ثانیه)

۲۹۸.۶۶

۲۵.۲۱

زمان آموزش

زمان ارزیابی

۲-۴. fine-tuning شبکه ترنسفرمر

الف) از مدل پیش آموخته شده DeiTBaseDistilled استفاده می‌کنیم که بر روی مجموعه Imagenet1k یادگیری شده می‌خواهیم از ویژگی‌های استخراج شده توسط این مدل را استفاده کنیم. حال برای اینکه بعضی از لایه‌های مدل را دوباره آموزش بدیم طبق مقاله لایه‌های دوازدهم به بعد را از قسمت ترنسفرمر شبکه را unfreeze می‌کنیم تا دوباره مدل را آموزش دهیم و fine-tuning کنیم حال برای

تسک خودمان طبق گفته مقاله ابتدا لایه ها فلت می کنیم تا سایز داده ها یک بعدی شه حال از بردارهای به دست آمده را به یک لایه فولی کانکت ۲۵۶ با نورون وصل می کنیم. حال برای جلوگیری از بیش بردازش Dropout با مقدار احتمال ۵۰ درصد اضافه کنیم حال برای خروجی از آنجایی که ۱۰ کلاس داریم یک لایه با ۱۰ نورون و با Softmax به عنوان تابع فعال ساز می گذاریم.

ب) تعداد پارامترهای trainable ۹۶۳۰۴۲۶، است از آنجایی که خروجی مدل یک بردار ۱۰۰۰ تاست و به یک لایه ۲۵۶ نورونی وصل است پس $۲۵۶ * (۱ + ۱۰۰۰) = ۲۵۶۲۵۶$ و بعد این ۲۵۶ را به یک لایه ۱۰ تایی وصل می کنیم و $۱۰ * (۱ + ۲۵۶) = ۲۵۷۰$ و باقی مربوط به لایه های آخر unfreeze شده مدل است.

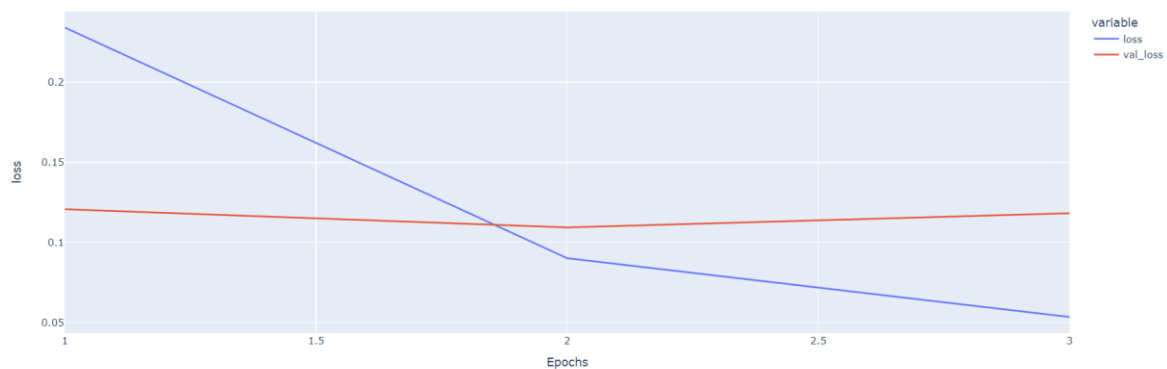
ج) حال مدل را با بهینه ساز آدام و لرنینگ ریت ۰.۰۰۰۱ و تابع loss، categorical_crossentropy و با batchsize ۶۴ را به تعداد ۳ اپاک (به علت طولانی بودن و رسیدن به جواب خوب در همین تعداد نسبت به حالت قبل تعداد اپاک ها را کم کردیم) ران می کنیم.

د) در پایان اپاک دهم دقت و مقدار تابع هزینه روی داده های آموزشی و ارزیابی به شکل زیر می شوند

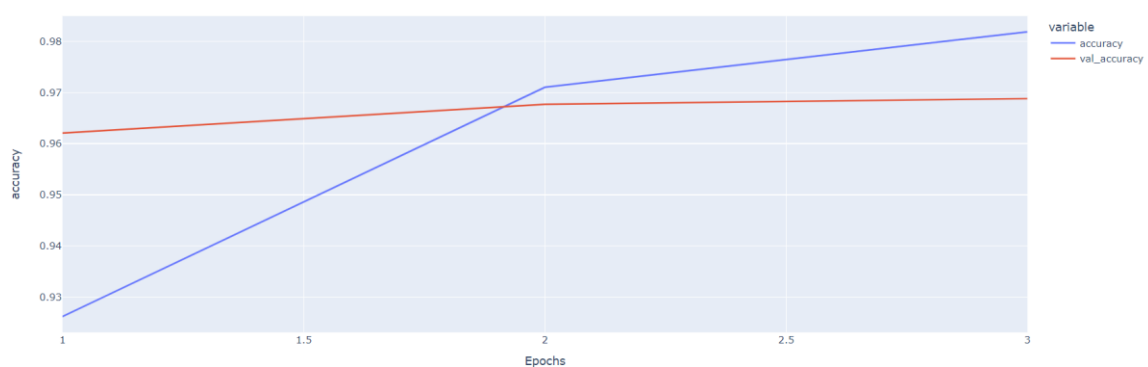
دقت	
۹۸.۲	داده آموزشی
۹۶.۷	داده ارزیابی

تابع هزینه	
۰.۰۵۳	داده آموزشی
۰.۱۱۸	داده ارزیابی

همچنین نمودار تغییراتشان را نیز در زیر نشان دادیم



شکل ۵. نمودار تغییرات تابع هزینه در داده‌های آموزشی و ارزیابی



شکل ۶. نمودار تغییرات دقت در داده‌های آموزشی و ارزیابی

(۵)

میانگین زمان (برحسب
ثانیه)

۱۳۴۷

۱۳۷

زمان آموزش

زمان ارزیابی

۲-۵. مقایسه نتایج

خب همینطور که در بخش‌های قبل گفتیم دقت ما روی مجموعه ارزیابی در ViT ، ۹۶.۷ و در CNN ، ۸۵.۲ شده است و همانطور که انتظار داشتیم مدل ترنسفرمری ما عملکرد بهتری نسبت به CNN داشته است که در مقدار تابه هزینه و مقایسه بینشون مشخص هست

و همانطور که از نمودار هایشان مشخص است مدل ترنسفری در همان ایپاک‌های اول به نتایج خوبی می‌رسد و از کانولوشنی بهتر عمل کرده است که می‌توان تاثیر معماری ترنسفرمر و مفهوم attention در بهبود دقت را مشاهده کرد

همچنین دو مدلی انتخابی تعداد پارامترهای قابل یادگیریشون تفاوتی زیادی را دارا نمی‌باشند اما خب از سویی معماری ترنسفرمر باعث می‌شود که زمان یادگیری آن خیلی بیشتر باشد و تقریبا ۴ برابر نسبت به مدل کانولوشنی شده است که این را می‌توان از معایب مدل کانولوشنی دانست .

در مقایسه با مقاله هم دقت ما روی مجموعه اعتبارسنجی در مدل ترسفرمر تقریبا همان شده است اما در مدل کانولوشنی حدود ۵ درصد کاهش داشته است که یکی از دلایل این تغییر می‌تواند تعداد ایپاک‌ها باشد اما با بررسی نمودار نظر میرسد که دقت همگرا شده بوده و بعید است تاثیر گذار بوده باشد و باقی شرایط به مانند مقاله بوده اما تفاوتی که ممکن است باشد این که داده افزایی بیشتر و مانند مقاله می‌توانست دقت را بالاتر ببرد و مدل نمونه‌های بیشتری را دارا باشد

