

هدف از تمرین ۳، پیاده‌سازی شبکه RNN به‌منظور حل مسئله sentiment classification می‌باشد. از داده ست Sports and Outdoors کامنت‌های آمازون به این منظور استفاده شده‌است. ۵ کلاس وجود دارد که کلاس ۵ به منظور بیشترین رضایت‌مندی است. ابتدا شبکه‌هایی با معماری‌های متفاوت بر روی دیتاست ذکرشده آموزش می‌بینند و سپس تاثیر بازنمایی‌های از قبل آموزش دیده شده و داده‌های ناموجود در دیتاست آموزش بررسی می‌شود. در انتها عملکرد مدل آموزش دیده بر روی دیتاست دیگری به نام دیتاست imdb سنجیده می‌شود.

بخش اول

۱- پیش‌پردازش دادگان

ابتدا دیتاست Sports and Outdoors توسط wget دانلود می‌شود. این دیتاست به صورت فشرده دانلود می‌شود به همین علت ابتدا باید از حالت فشرده خارج شود. این دیتاست از چند فیلد تشکیل شده‌است که برای تسک این تمرین تنها نیاز به قسمت overview که متن کامنت می‌باشد و overall که عددی بین ۱ تا ۵ است می‌باشد. هر خط از داده پردازش می‌شود و قسمت overview در کلیدی به همین نام در دیتافریم قرار می‌گیرد. Overall هم به همین صورت. در انتها قسمت overview به عنوان x دیتاست Overall هم به عنوان y دیتاست خروجی داده می‌شود. قبل از خروجی اما پیش پردازش‌هایی مانند حذف تگ‌ها، تبدیل جمله به آرایه ای از کلمه‌ها و تبدیل هر کلمه به ریشه آن صورت می‌گیرد. از این دیتاست به تعداد ۲.۴ میلیون داده لود می‌شود تا برای آموزش و تست مورد استفاده قرار بگیرد.

در هر مدل که به داده‌های متنی کار می‌کند باید به نحوی متن را به صورت مجموعه‌ای از اعداد به مدل ورودی داد. به این منظور در ابتدا از Tokenizer موجود در tf استفاده شده است. Tokenizer بر روی قسمت x داده فیت می‌شود و هر کلمه جمله را به صورت آرایه ای از اعداد در می‌آورد که هر عدد درواقع عدد متناظر به کلمه می‌باشد. به عبارتی Tokenizer با استفاده از داده‌های داده شده به آن یک دیکشنری تشکیل می‌دهد که هر کلمه را به یک عدد نسبت می‌دهد. سپس هر جمله را تبدیل به دنباله ای از اعداد متناظر با کلمه‌های موجود در آن جمله می‌کند.

نکته مورد توجه دیگر این می‌باشد که ورودی‌های شبکه باید طول برابر داشته باشند و به همین دلیل ابتدا با یک تابع کمکی میانگین طول جمله‌های موجود در دیتاست محاسبه می‌شود و سپس همه جمله‌ها به آن سائز تبدیل می‌شوند. همچنین دقت شد که y‌های مربوط به آموزش باید به فرمت hot-۱ دربیایند.

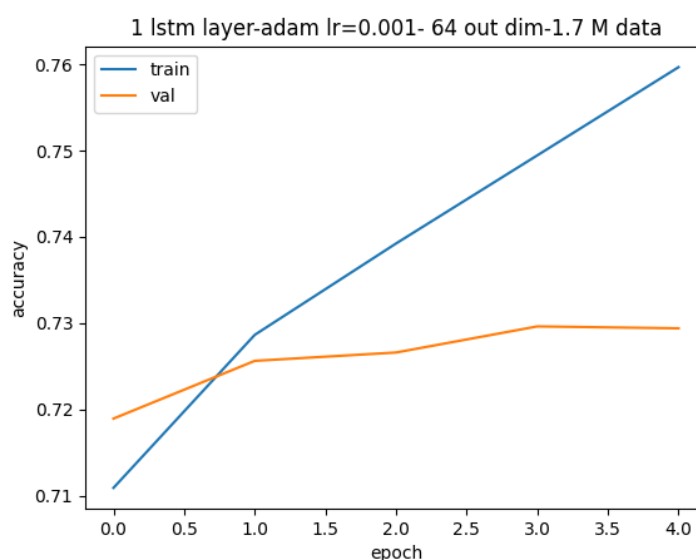
۲- پیاده‌سازی شبکه

برای معماری شبکه ابتدا از یک لایه embedding استفاده می‌شود. از آنجایی که با استفاده از Tokenizer یک دیکشنری ایجاد شد بردار اولیه برای هر کلمه اندازه ای برابر با سائز دیکشنری دارد و به همین علت، ورودی این لایه نیز سائز دیکشنری است. این لایه مسئول یادگیری بازنمایی کلمات می‌باشد. این لایه بازنمایی‌های ابتدایی را دریافت می‌کند و پس از آموزش مدل به عنوان خروجی بازنمایی‌های آموزش یافته را خروجی می‌دهد که سائز این بازنمایی خروجی باید به عنوان ورودی به این لایه داده شود. بعد از این لایه

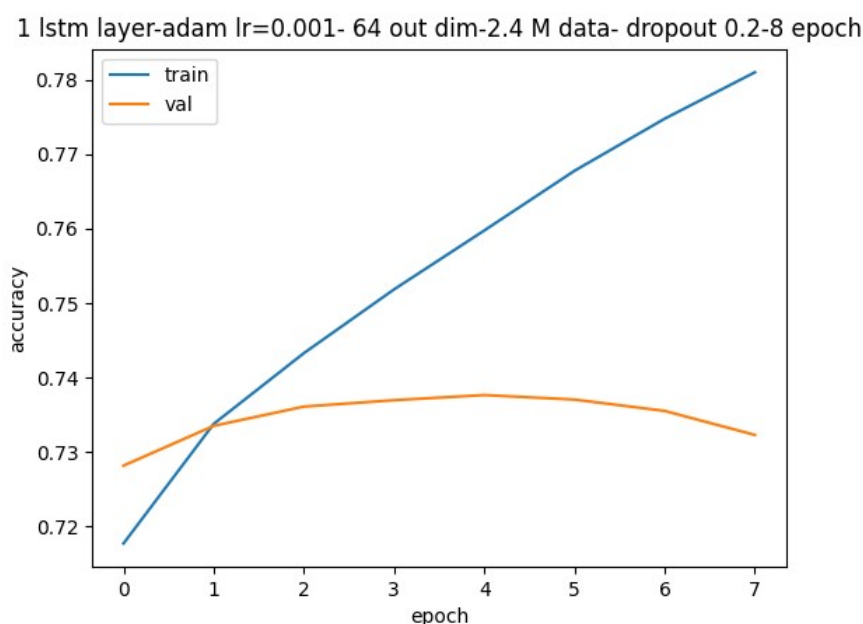
ابتدا با تعداد یک لایه، لایه های مختلف مانند LSTM، GRU، Bidirectional و SimpleRNN تست شده اند. در ادامه مشاهده می شود که این مدل ها با $overfit$ روبه رو می شوند و به این منظور از لایه های $dropout$ و $Batch Normalization$ نیز استفاده شده است. یک لایه $Dense$ با ۵ نورون برای ۵ کلاس با تابع فعال سازی سافت مکس در انتها قرار می گیرد. برای بهینه سازی با دو بهینه ساز $adam$ و $RMSprop$ تست انجام شده است. متریک انتخاب شده اکیورسی و تابع لاس هم $CategoricalCE$ انتخاب شده است.

۳- انجام آزمایش ها

در ابتدا مدل زیر با پارمترهای مشخص شده ران می شود و این مدل دچار مشکل اورفیت است. در ابتدا تعداد داده ۱.۷ میلیون استفاده می شد.

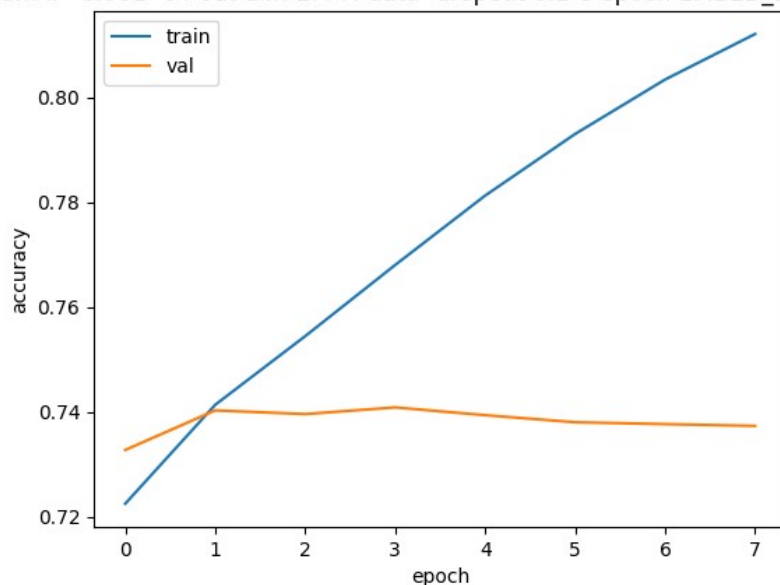


برای آدرس دادن به این مشکل تعداد داده ها به ۲.۴ میلیون افزایش داده شد و لایه $dropout$ با نرخ ۰.۲ اضافه شد و مدل زیر ران شد.



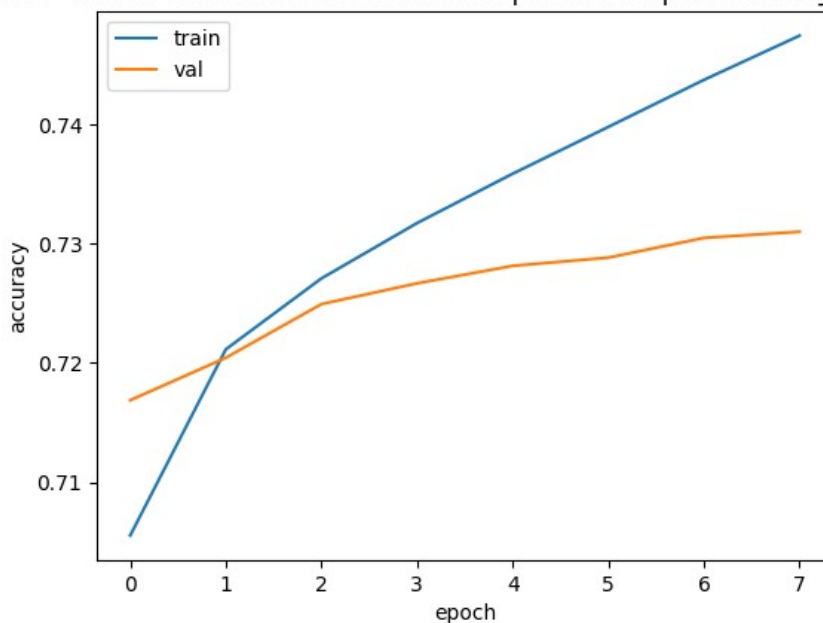
مشاهده می‌شود که همچنان مشکل اورفیت وجود دارد و حل نشده است و این در حالی است که تعداد داده‌های زیادی به شبکه داده شده است و از دراپ‌اوت هم استفاده شده است. تصمیم بعدی برای حل مشکل اورفیت این است که خروجی لایه بازنمایی از ۶۴ به ۱۲۸ افزایش پیدا کند.

adam lr=0.001- 64 out dim-2.4 M data- dropout 0.2-8 epoch-EMBED_DIM 128-



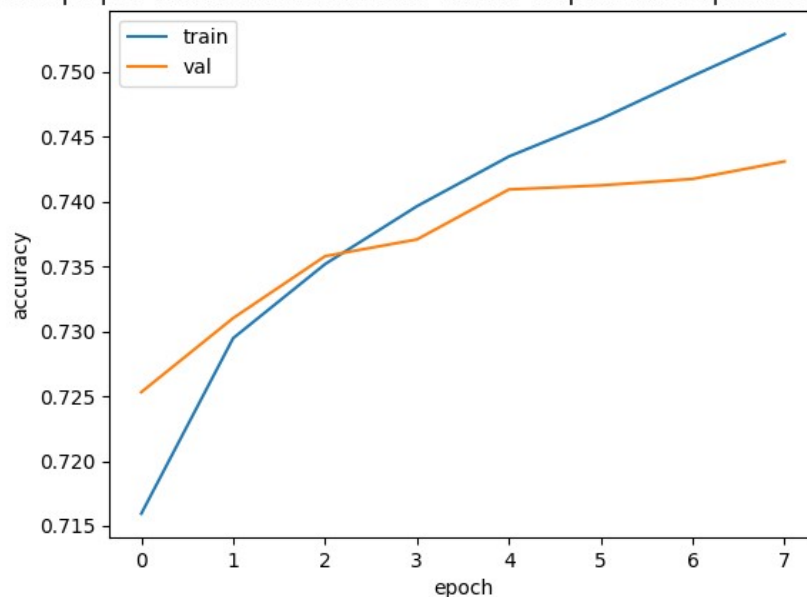
مشاهده می‌شود که تغییر خاصی رخ نداده است. (مقداری اکپورسی‌ها در هر دو ست افزایش پیدا کرده‌اند اما باز هم اورفیت وجود دارد). همچنین باید توجه داشت که این مدل‌ها با ۸ ایپاک هم به همگرایی نرسیدند و به علت زمان تخمینی حدود ۱۰ دقیقه برای ران هر مدل امکان استفاده از ایپاک‌های بیشتر وجود نداشت. آزمایش بعدی به منظور حل اورفیت کاهش Learning rate از ۰.۰۰۱ به ۰.۰۰۰۱ بود.

adam lr=0.0001- 128 out dim-2.4 M data-dropout 0.2-8 epoch-EMBED_DIM 128-



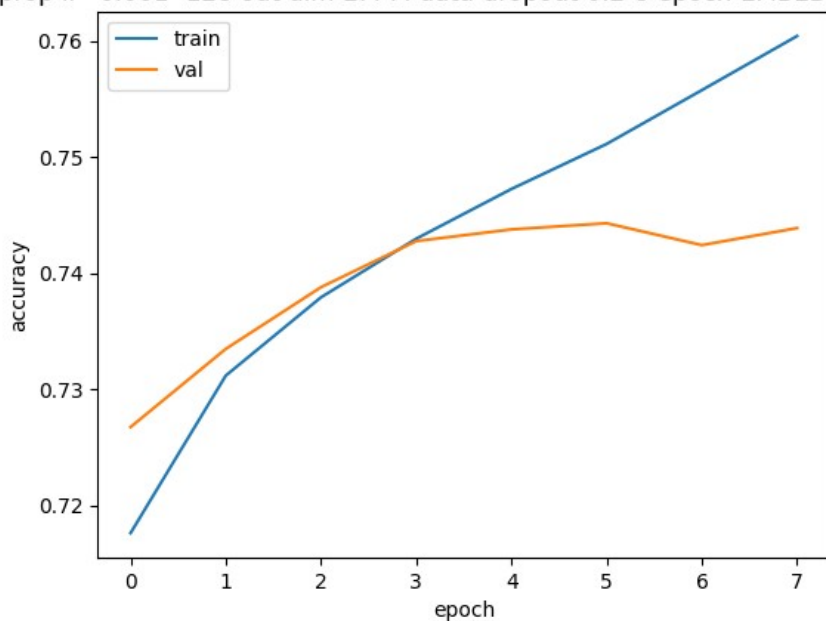
مشاهده می‌شود که اکپورسی ها مقداری کاهش پیدا کرده اند اما اورفیت کمتر شده است. تا اینجا همه مدل ها دارای یک لایه lstm بودند. در ادامه برای مقایسه بین انواع لایه‌های RNN مدل های زیر با پارمترهای مشخص شده ران شدند.
یک لایه Bidirectional:

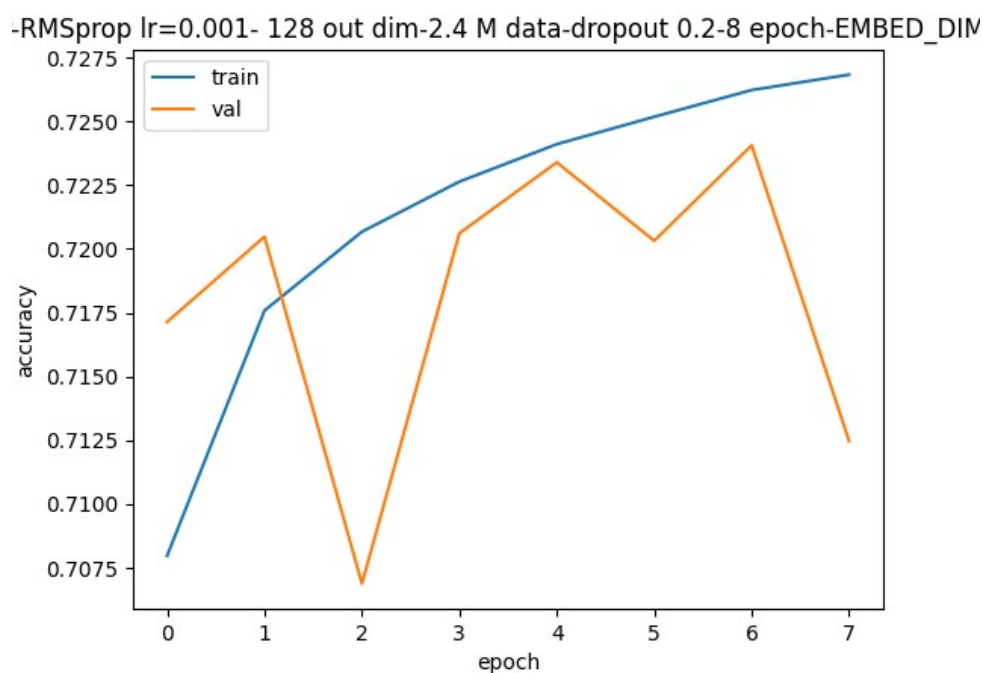
r-RMSprop lr=0.001- 128 out dim-2.4 M data-dropout 0.2-8 epoch-EMBED_DIM



یک لایه GRU:

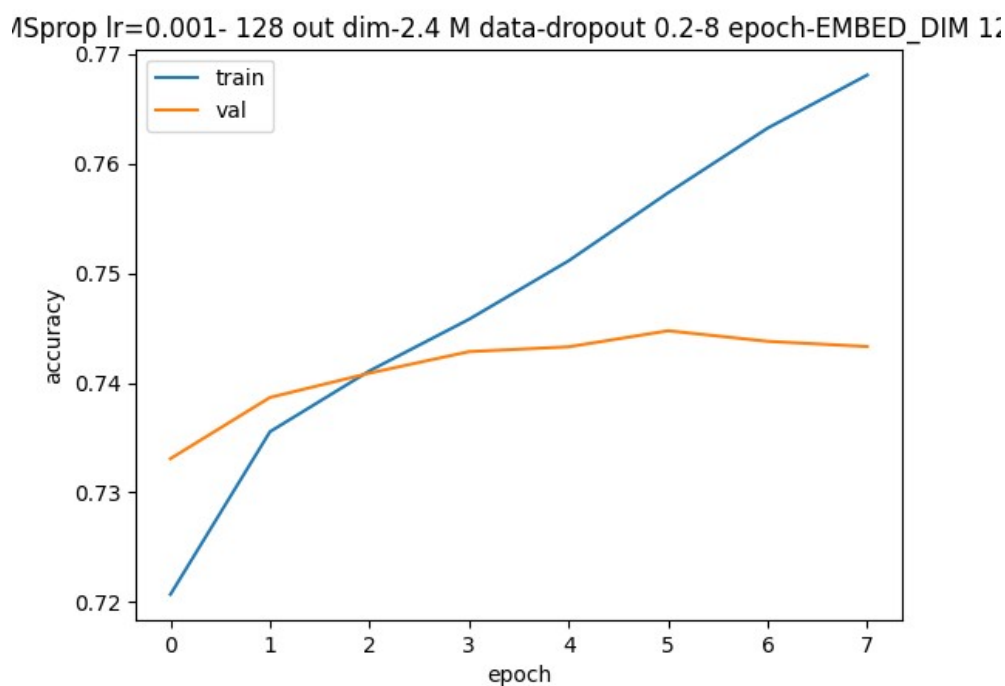
lSprop lr=0.001- 128 out dim-2.4 M data-dropout 0.2-8 epoch-EMBED_DIM 12





مشاهده می‌شود که بیشترین اکیورسی متعلق به مدل یک لایه GRU است و کمترین متعلق به SimpleRNN. یک لایه Bidirectional مقدار کمی از یک لایه lstm بهتر عمل کرده است اما هر دو کمتر از یک لایه GRU هستند. (مقایسه ها بر اساس جدول ۱ انجام شده است).

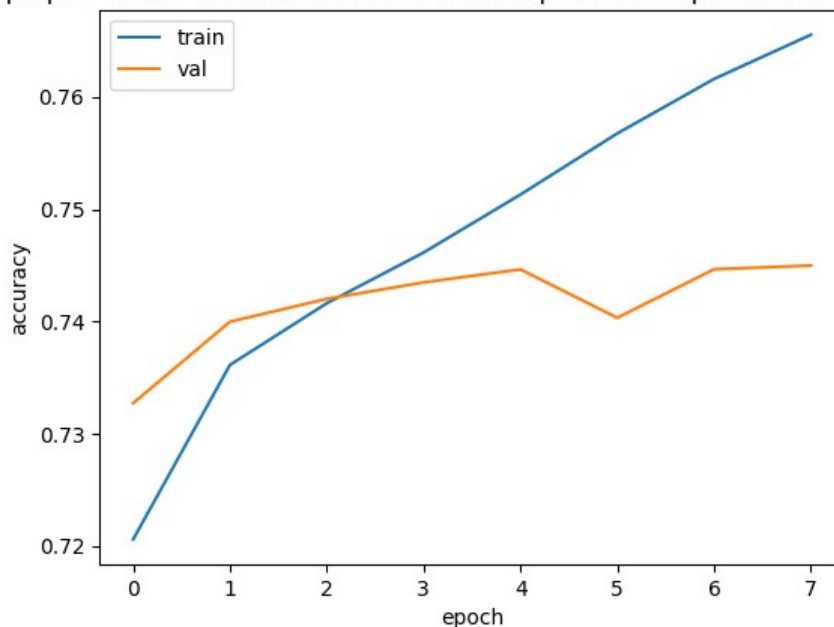
قدم بعدی آزمایش ها مربوط به بررسی تعداد لایه است. ابتدا مدل دو لایه lstm تست می‌شود.



مشاهده میشود که مقدار کمی در حد هزارم اکیورسی از یک لایه lstm بهبود داشته است. (بر اساس جدول ۱)

سه لایه lstm:

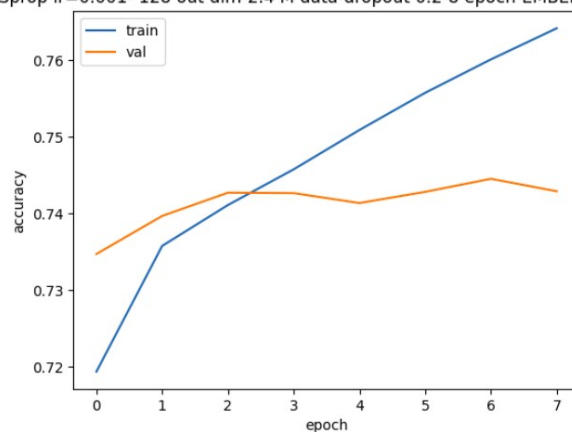
4Sprop lr=0.001- 128 out dim-2.4 M data-dropout 0.2-8 epoch-EMBED_DIM 1:



دقت شود که مقایسه این مدل ها با هم از روی این نمودار ها انجام نشده است و از روی اکیورسی که روی داده تست با تابع predict صورت گرفته است، انجام شده ایت و جدول حاوی اکیورسی عملکرد همه مدل ها بر روی داده تست در ادامه داخل جدول ۱ قرار گرفته است.

مدل ۳ لایه lstm بر اساس جدول ۱ از ۲ لایه بهتر عمل کرده است بر روی داده تست عمل کرده است و به طور کلی هم به عنوان بهترین مدل (مدل پایه را آزمایش های بخش های بعد) برای تست w2v و آزمایش های بعدی بر روی آن انتخاب می شود. مدل ۴ لایه lstm:

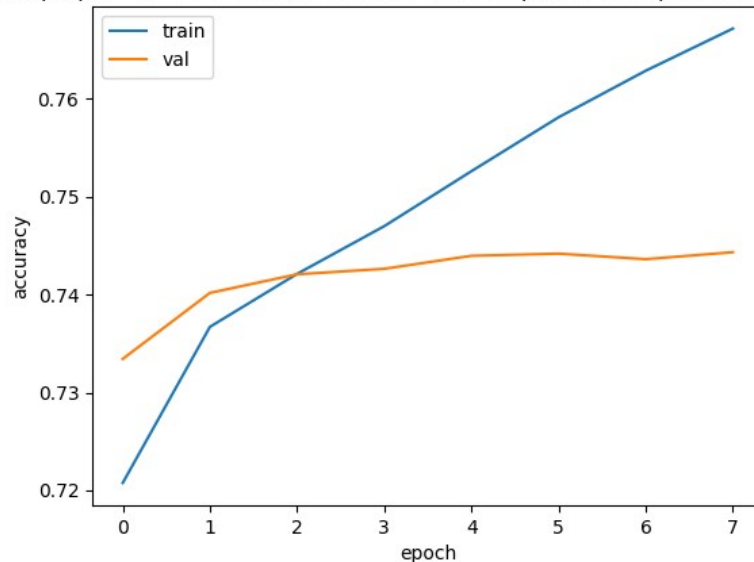
4 LSTM layer-RMSprop lr=0.001- 128 out dim-2.4 M data-dropout 0.2-8 epoch-EMBED_DIM 128-LSTM_OUT 128



مدل ۴ لایه lstm نسبت به ۳ لایه بهبود نداشته است.

در انتهای این مرحله نیز مدل ۳ لایه که ۲ لایه اول آن lstm و لایه سوم آن gru است آزمایش شده است.

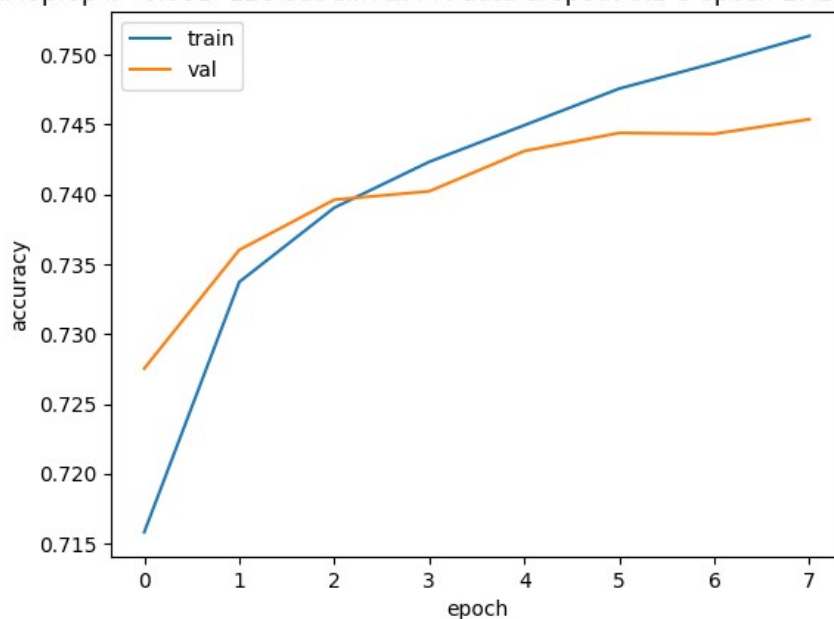
-RMSprop lr=0.001- 128 out dim-2.4 M data-dropout 0.2-8 epoch-EMBED_DIM



طبق جدول ۱ این مدل نسبت به مدل ۳ لایه و ۴ لایه lstm کاهش اکيورسی داشته است.

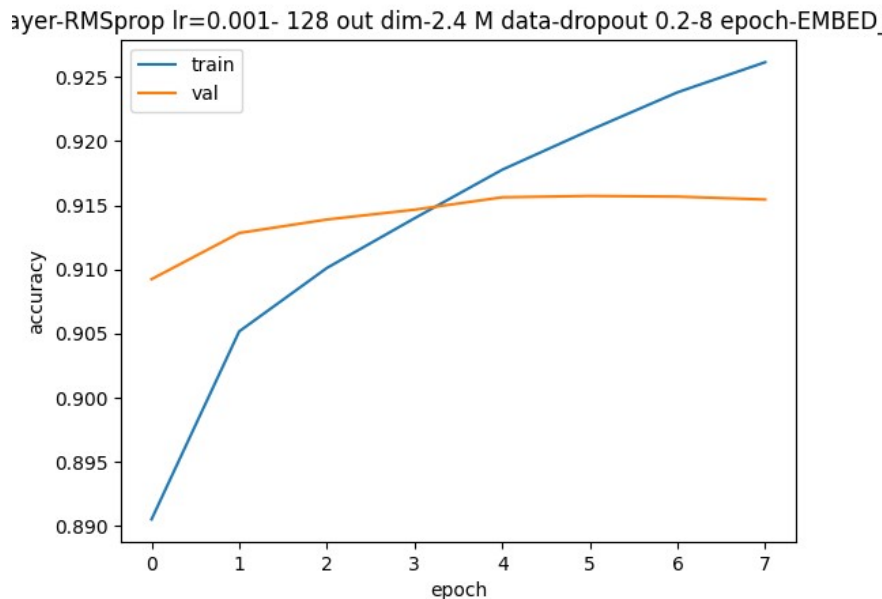
برای مرحله دوم آزمایش ها از بازنمایی های از پیش آموزش داده شده (w2v) Word2Vec استفاده می شود. از کتابخانه gensim این مدل از پیش آموزش داده بر روی دیتاست ویکی پدیا را با نام glove-wiki-gigaword-100 لود می شود این مدل برای هر کلمه بازنمایی به طول ۱۰۰ ایجاد می کند. برای استفاده از این بازنمایی ها کافیست تا در ابتدا برای همه لغات موجود در دیکشنری ساخته شده توسط Tokenizer بازنمایی w2v آن ها را از مدل لود شده دریافت کرده و در ماتریسی به نام ماتریس امبدینگ قرار دهیم. سپس در لایه Embedding این ماتریس را به عنوان وزن ورودی دهیم. در اینجا اضافه کردن بازنمایی های w2v به مدل ۳ لایه lstm انجام شده است.

-RMSprop lr=0.001- 128 out dim-2.4 M data-dropout 0.2-8 epoch-EMBED_DIM



با توجه به جدول ۱ عملکرد این مدل نسبت به مدل ۳ لایه lstm کاهش جزئی داشته است.

برای مرحله سوم آزمایش‌ها که آموزش دادن مدل روی داده‌های ۳ کلاس ۱ و ۳ و ۵ و تست آن روی داده‌های کلاس‌های ۲ و ۴ است، یک تابع نوشته شده است که داده‌های سه کلاس ۱، ۳، ۵ را بعنوان داده آموزش و داده‌های ۲ و ۴ را به عنوان داده تست خروجی می‌دهد. سپس بر روی این داده آموزش جدید نیز پیش‌پردازش‌های قبلی به وسیله Tokenizer انجام می‌شود و ورودی به مدل (بدون بازنمایی w2v) داده می‌شود. پس از آموزش این مدل پلات آن به صورت زیر است:



مشاهده می‌شود که اکيورسی آموزش و validation افزایش پیدا کرده است و به بالای ۹۰ رسیده است اما با توجه به جدول ۱ اکيورسی این مدل برای داده‌های تست برابر با ۰ است. داده‌های تست از کلاس‌هایی بودند که تا حالا مدل آن‌ها را ندیده است و مدل سعی می‌کند تا نزدیک‌ترین کلاسی را که می‌شناسد به آن‌ها نسبت دهد. به عنوان مثال به کلاس ۲ کلاس‌های ۱ و ۳ و به کلاس ۴ کلاس‌های ۳ و ۵ را نسبت دهد. اما باز هم چون این کلاس‌ها، کلاس‌های درست نیستند اکيورسی ۰ می‌شود.

در بخش سوم آزمایش‌ها عملکرد مدل آموزش دیده قبلی (بر روی داده‌های ۳ کلاس ۱ و ۳ و ۵) آزمایش می‌شود. به این منظور ابتدا دیتا ست imdb با wget دانلود می‌شود و توسط یک تابع مشابه تابع قبلی یک دیتا فریم یاخته می‌شود و پیش‌پردازش‌های مانند حذف تگ‌ها و تبدیل به حروف کوچک و تبدیل به ریشه مردن صورت می‌گیرد سپس این داده‌های به عنوان داده تست به مدل داده می‌شوند با استفاده از predict و اکيورسی برای این پاسخ‌ها شنجیده می‌شود که برابر با مقداری نزدیک به ۰ و در اینجا ۰.۰۲ است.

دلیل این موضوع این می‌باشد که داده‌های imdb دارای دو کلاس مثبت و منفی هستند که چون توسط مدل در حین آموزش دیده نشده‌اند، مدل داده‌های مثبت را به کلاس ۵ و داده‌های منفی را به کلاس ۱ نسبت می‌دهد و چون این کلاس‌ها، کلاس‌های درست نیستند اکيورسی مقدار پایینی به دست می‌آید.

نام مدل	اکيورسی روی داده تست
LSTM ۱ لایه	۰.۷۴۴۷۰۶۵۶۸۸۵۶۲۸۵۴
Bidirectional ۱ لایه	۰.۷۴۷۹۵۷۶۵۲۵۵۰۸۵۰۳
GRU ۱ لایه	۰.۷۴۹۴۵۸۱۵۲۷۱۷۵۷۲۵
SimpleRNN ۱ لایه	۰.۷۱۱۴۰۳۸۰۱۲۶۷۰۸۹
LSTM ۲ لایه	۰.۷۴۸۹۱۶۳۰۵۴۳۵۱۴۵۱
LSTM ۳ لایه	۰.۷۵۱۲۰۸۷۳۶۲۴۵۴۱۵۲
LSTM ۳ لایه	۰.۷۵۱۰۰۰۳۳۳۴۴۴۴۸۱۵
GRU ۱ و LSTM ۲ لایه	۰.۷۴۷۶۶۵۸۸۸۶۲۹۵۴۳۲
LSTM ۳ لایه با بازنمایی w2v	۰.۷۴۹۶۲۴۸۷۴۹۵۸۳۱۹۴