

به نام خدا



دانشگاه تهران



دانشکده مهندسی برق و کامپیوتر

**درس شبکه‌های عصبی و یادگیری عمیق**

**تمرین چهارم**

نام و نام خانوادگی	مریم دادخواه – جواد سراج
شماره دانشجویی	810101186 - 810101151
تاریخ ارسال گزارش	1402.03.03

7	پاسخ 1. توصیف عکس.....
7	1-1. پیش پردازش داده‌ها.....
8	1-2. معماری شبکه ( Encoder و Decoder).....
9	1-3. پیاده سازی شبکه Decoder-Encoder با وزن‌های فریز شده شبکه ResNet18.....
12	1-4. پیاده سازی شبکه Decoder-Encoder با شبکه ResNet18 با وزن‌های trainable.....
14	1-5. بحث و نتیجه‌گیری.....
16	پاسخ 2. تشخیص اندیشه.....
16	1-2. معماری LSTM و embedding.....
20	2-2. پیش پردازش دادگان.....
21	2-3. پیاده سازی طبقه بندی نیت.....
36	2-4. پیاده سازی مدل Responder.....

- شکل 1. نمونه تصویر و کپشن‌های آن در دیتاست.....7
- شکل 2. نمودار کاهش تابع هزینه بر روی داده‌های آموزشی و داده‌های اعتبارسنجی- وزن‌های فریز شده ResNet18.....10
- شکل 3. نمونه تصویر و کپشن تولید شده - freezed ResNet 18.....10
- شکل 4. نمونه تصویر و کپشن تولید شده - freezed ResNet 18.....11
- شکل 5. نمونه تصویر و کپشن تولید شده - freezed ResNet 18.....11
- شکل 6. نمودار کاهش تابع هزینه بر روی داده‌های آموزشی و داده‌های اعتبارسنجی- وزن‌های trainable برای ResNet18.....12
- شکل 7. نمونه تصویر و کپشن تولید شده - trainable ResNet 18.....13
- شکل 8. نمونه تصویر و کپشن تولید شده - trainable ResNet 18.....13
- شکل 9. نمونه تصویر و کپشن تولید شده - trainable ResNet 18.....14
- شکل 10. ساختمان مدل اول.....22
- شکل 11. نمودار loss و accuracy برای داده های train و validation برای state hidden 100.....23
- شکل 12. نمودار loss و accuracy برای داده های train و validation برای state hidden 25.....23
- شکل 13. نتایج پیاده سازی مدل اول بر داده های test برای state hidden 25.....24
- شکل 14. نتایج پیاده سازی مدل اول بر داده های test برای state hidden 100.....24
- شکل 15. ماتریس پراکندگی مدل اول با state hidden 25.....25
- شکل 16. ماتریس پراکندگی مدل اول با state hidden 100.....25
- شکل 17. ساختمان مدل دوم.....26
- شکل 18. نمودار loss و accuracy برای مدل دوم کتگوری class main روی داده های train و val ، با state hidden 25.....28
- شکل 19. نمودار loss و accuracy برای مدل دوم کتگوری class sub روی داده های train و val ، با state hidden 25.....28
- شکل 20. نتایج حاصل از مدل دوم بر روی داده های test برای class main با state hidden 25.....29
- شکل 21. نتایج حاصل از مدل دوم بر روی داده های test برای class sub با state hidden 25.....29
- شکل 22. ماتریس پراکندگی مدل دوم با state hidden 25 بر روی داده های test برای class main.....31

- شکل 23. ماتریس پراکندگی مدل دوم با 25 state hidden بر روی داده های test برای class sub ..... 31
- شکل 24. نمودار loss و accuracy برای مدل با 100 state hidden با داده های train و val برای class main ..... 34
- شکل 25. نمودار loss و accuracy برای مدل با 100 state hidden با داده های train و val برای class sub ..... 34
- شکل 26. نتایج حاصل از مدل دوم بر روی داده های test برای class main با 100 state hidden ..... 35
- شکل 27. نتایج حاصل از مدل دوم بر روی داده های test برای class sub با 100 state hidden .. ..... 35
- شکل 28. ماتریس پراکندگی مدل دوم با 100 state hidden بر روی داده های test برای class main ..... 36
- شکل 29. ماتریس پراکندگی مدل دوم با 100 state hidden بر روی داده های test برای class sub ..... 36
- شکل 30. ساختمان مدل responder ..... 37
- شکل 31. نحوه reshape کردن داده های جواب ..... 38
- شکل 32. کد مربوط به نحوه تبدیل خروجی شبکه به پاسخ های متنی ..... 39
- شکل 33. پاسخ های شبکه طراحی شده به سوال های مطرح شده ..... 40





## جدول‌ها

جدول 1. مقادیر rate learning برای lstm با مقادیر 2 states hidden ..... 27

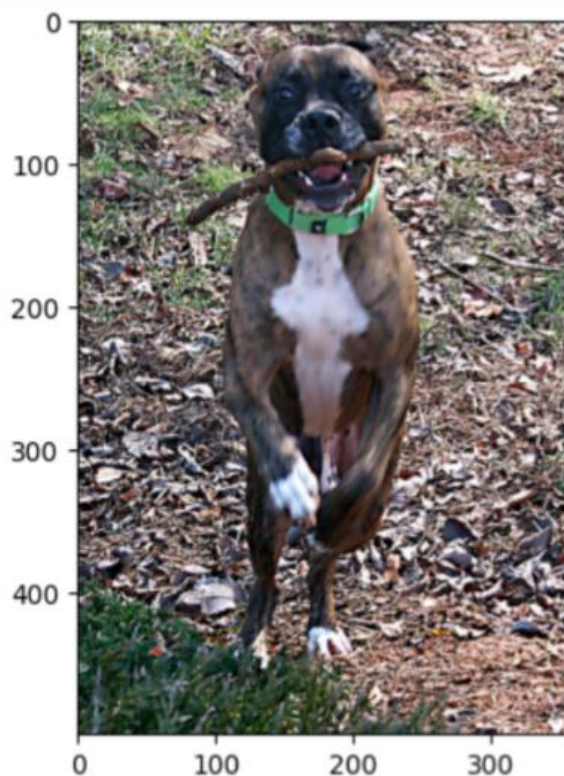
جدول 2. مقایسه عملکرد مدل دوم برای تعداد 25 و 100 state hidden برای لایه LSTM..... 33

## پاسخ 1. توصیف عکس

در این سوال به بررسی و آموزش شبکه‌ای برای تولید کپشن از تصاویر ورودی می‌پردازیم. این مدل از resnet 18 برای استخراج ویژگی از تصویر استفاده می‌کند و یک لایه LSTM برای تولید کپشن استفاده خواهد کرد. در ادامه به جزییات این شبکه و دیتاست مورد استفاده می‌پردازیم.

### ۱-۱. پیش پردازش داده‌ها

ما ابتدا داده‌ها را با استفاده از API کگل دانلود کردیم. سپس مسیرهای تصویر و کپشن‌ها را در یک دیتافریم pandas قرار دادیم. هر نمونه دیتاست شامل یک تصویر و 5 کپشن است که توسط افراد داده شده است. نمونه ای از تصویر و کپشن‌های مربوط به آن در شکل زیر داده شد است.



```
:caption - A brown and white dog is jumping up , with a stick in its mouth .  
:caption - A brown dog , holding a stick .  
:caption - A brown dog with green collar holding a stick in its mouth .  
:caption - A dog in a green collar fetches a stick from the dead leaves .  
:caption - The brown and white dog with a green collar is biting a stick .
```

شکل 1. نمونه تصویر و کپشن‌های آن در دیتاست



در ادامه به پردازش متن‌ها می‌پردازیم. برای پردازش داده‌ها اولین پردازش، ابتدا داده‌ها را توکنایز می‌کنیم و به کلمات تبدیل می‌کنیم. در ادامه، بعد کلمات را lower case می‌کنیم. سپس Vocabulary را می‌سازیم. همچنین علائم نگارشی موجود در متن نیز حذف شده‌اند. برای ساختن Vocabulary، دیکشنری می‌سازیم و به ازای هر کلمه تعداد تکرار آن محاسبه می‌کنیم و اگر تعداد تکرار کلمه‌ای از حد آستانه‌ای<sup>۱</sup> کمتر باشد آن داده را از Vocabulary حذف می‌کنیم. همچنین یک lookup table تشکیل می‌دهیم و به هر کلمه ایندکس نسبت می‌دهیم. در ابتدای کار حداستانه 5 در نظر گرفته شده است.

0:"<UNK>",3:"<EOS>",2:"<SOS>",1:"<PAD>" را نیز به Vocabulary اضافه می‌کنیم که به مشخص شدن ابتدا و انتهای جمله و مقابله با کلمات ناشناخته پیش‌بینی شد توسط مدل کمک می‌کند. همچنین به دلیل این که تعداد کلمات در جملات متفاوت است از padding برای یکسانسازی طول جملات استفاده می‌کنیم. برای padding، بر اساس بزرگ‌ترین جمله در یک Minibatch تمام جملات را pad می‌کنیم.

برای پردازش تصاویر ابتدا تصاویر ابعاد تمام تصاویر را به ابعاد  $224 \times 224$  تبدیل کردیم تا ورودی به استاندارد لازم برای شبکه ResNet برسد. در نهایت مجموع دیتاست ما شامل 40455 کامنت و 40455 تصویر است که تعداد یونیک تصاویر 8091 است. یعنی به ازای هر تصویر، 5 کامنت موجود داریم.

## 2-۱. معماری شبکه (Encoder و Decoder)

شبکه نامبرده در مقاله دارای ساختار Encoder-Decoder است. در بخش Encoder یک شبکه Resnet18 برای، تبدیل تصویر ورودی به بردار ویژگی استفاده شده است. این شبکه نقش استخراج بردار ویژگی برای ورود به بخش Decoder را بر عهده دارد. خروجی ResNet را از یک لایه Linear عبور می‌دهیم تا به ابعاد embedding متناسب با بخش Decoder برسانیم.

در بخش Decoder ابتدا یک لایه Embedding قرار می‌گیرد که ابعاد آن  $\text{embeddingsize} \times \text{vocabsize}$  است و وظیفه آن Embedding هر کلمه به بردار ویژگی آن کلمه است. سپس embedding کلمات درون کامنت را با بردار ویژگی تصویر کانکت می‌کنیم و این تنسور را به بخش LSTM می‌دهیم. LSTM مطابق با خواسته‌های صورت سوال ورودی embedding size دارد و تعداد  $\text{hidden size} = 256$  است و تعداد لایه‌های LSTM هم 1 در نظر گرفته شده است. سپس خروجی LSTM از یک لایه خطی (Fully connected) عبور داده می‌شود و خروجی این لایه به ابعاد Vocab size است تا بتوان کلمه موردنظر را پیش‌بینی کرد.

<sup>1</sup> Threshold

در بخش نهایی هدف تخمین یک کپشن مناسب است. برای این کار دقت می‌کنیم که فقط تصویر به عنوان ورودی داریم و هیچ متنی وجود ندارد و متن توسط شبکه تولید می‌شود. برای تولید متن توسط شبکه، بردار ویژگی embed شده تصویر را به لایه LSTM که در بخش قبل توضیح دادیم می‌دهیم و سپس LSTM یک بردار ویژگی تولید میکند که این بردار را به لایه fully connected می‌دهیم و یک بردار به ابعاد vocab size به ما خروجی میدهد که با argmax میتوان کلمه تخمینی با بیشترین احتمال را تخمین زد. سپس کلمه تخمین زده شده به عنوان ورودی برای تخمین کلمه بعدی استفاده می‌شود. اما توجه میکنیم در این مرحله ما ایندکس کلمه تخمینی را داریم و باید این ایندکس به لایه embedding داده شود تا بردار embedding آن کلمه استخراج شود. دو شرط مهم در این بخش رعایت می‌شود که در ادامه به آن می‌پردازیم.

- 1- ابتدا حلقه شرطی را طوری قرار میدهیم تا تعداد کلمات تخمینی از حدی بیشتر نشود. مثلاً حداکثر برای هر تصویر 20 کلمه predict کنیم. (ماکسیمم طول کپشن تولیدی)
- 2- شرطی قرار می‌دهیم که هر جا کلمه "<EOS>" تولید شد همان لحظه حلقه تولید کپشن را متوقف کنیم.

### پارامترهای شبکه و هایپرپارامترها

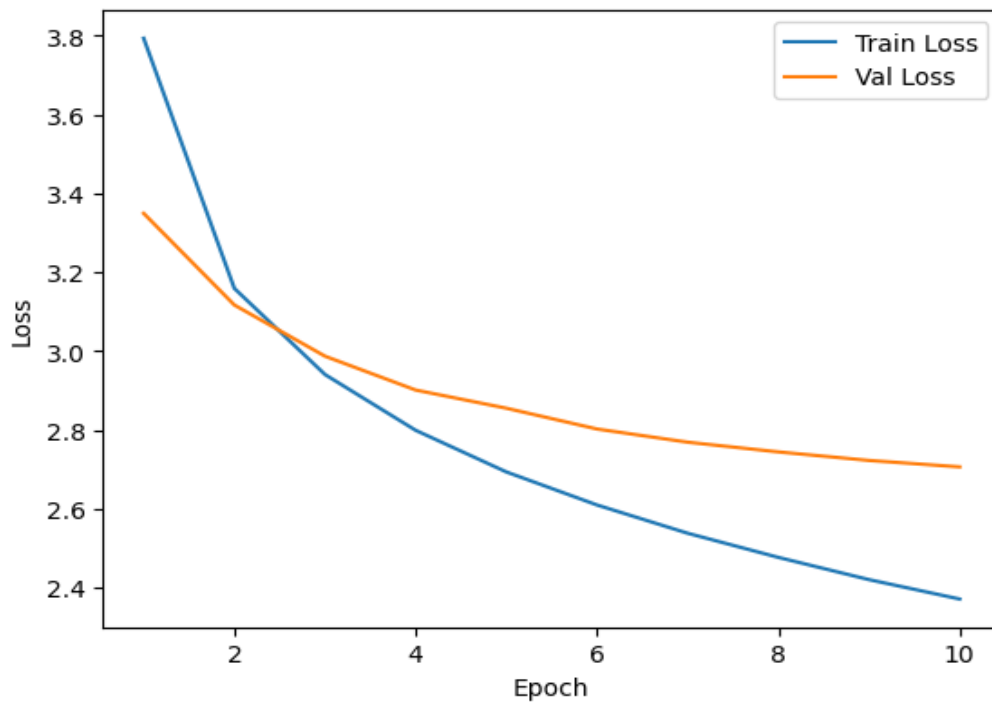
برای آموزش شبکه از تابع هزینه، cross entropy استفاده میکنیم. برای بهینه سازی، بهینه ساز Adam و با نرخ یادگیری 0.0001 را به کار می‌گیریم. Embedding Size برای هر کلمه را 200 در نظر گرفتیم و hidden size برای LSTM را 256 در نظر گرفتیم.

برای رسیدن به مدل مناسب، چند نرخ یادگیری متفاوت و embedding size های متفاوت مورد بررسی قرار گرفت تا بهینه‌ترین پارامترها برای ترین شبکه را انتخاب کنیم.

همچنین مطابق مقاله، ما برای لایه LSTM نرخ dropout = 0.30 در نظر گرفتیم. البته در مقاله dropout = 0.5 قرار داده شد که برای دیتاست ما مقدار 0.3 به نتایج بهتری رسید.

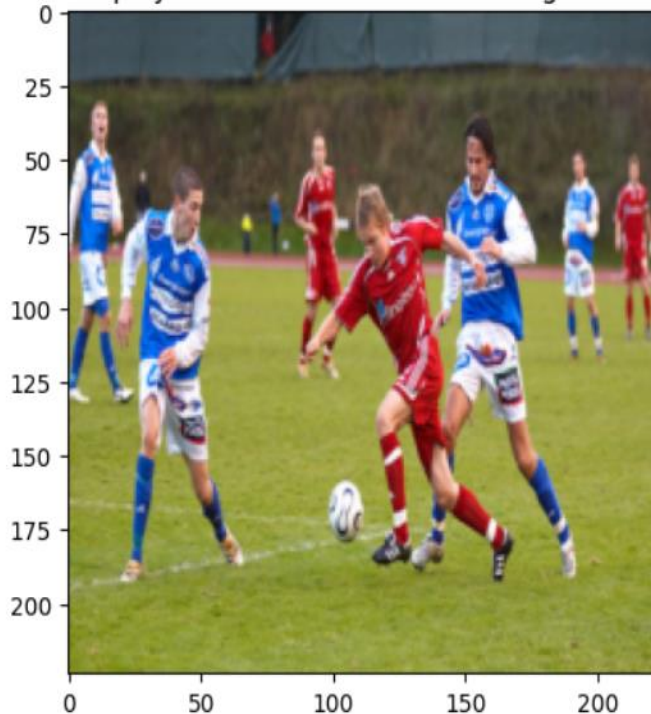
### 3-1. پیاده سازی شبکه Decoder-Encoder با وزن‌های فریز شده شبکه ResNet18

با توجه به خواسته‌های این پروژه، ابتدا وزن‌های Pretraibed شده مدل ResNet18 که به عنوان Encoder به کار گرفته شده است را فریز می‌کنیم. (به جز لایه آخر). سپس مدل را روی دیتاها ترین می‌کنیم. قطعه کد زیر برای فریز کردن وزن‌ها به کار گرفته می‌شود. نمودار تابع هزینه برای این حالت در زیر آورده شده است.



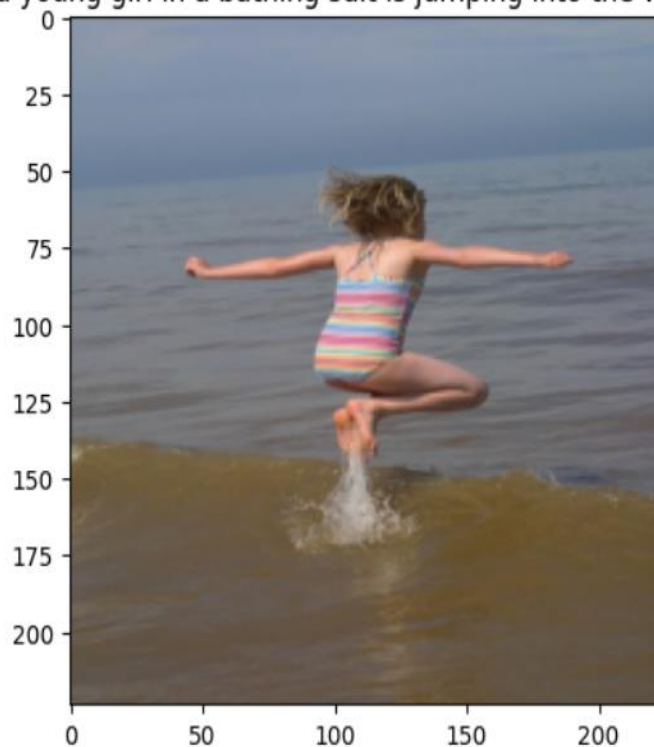
شکل 2. نمودار کاهش تابع هزینه بر روی داده‌های آموزشی و داده‌های اعتبارسنجی - وزن‌های فریز شده ResNet18

<SOS> a soccer player in a red uniform is running on the field . <EOS>



شکل 3. نمونه تصویر و کپشن تولید شده - freeze ResNet 18

<SOS> a young girl in a bathing suit is jumping into the water . <EOS>



شکل 4. نمونه تصویر و کپشن تولید شده - **frozen ResNet 18**

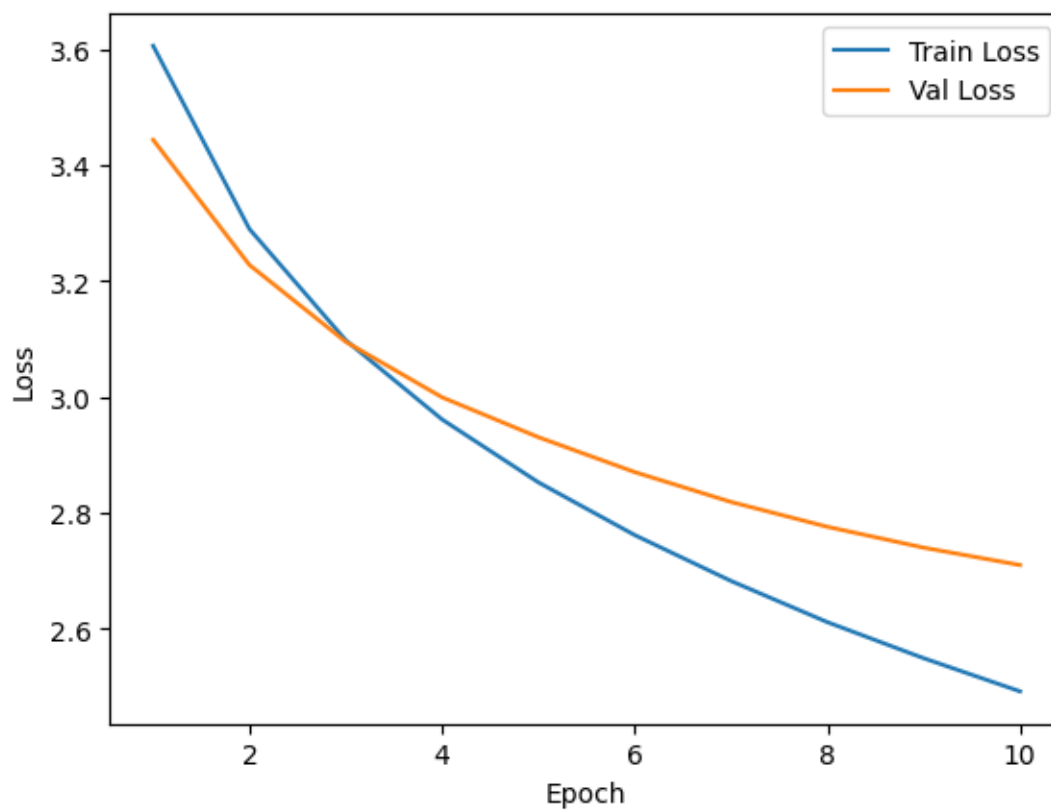
<SOS> a man in a black shirt and jeans is standing on a rock in front of a crowd .



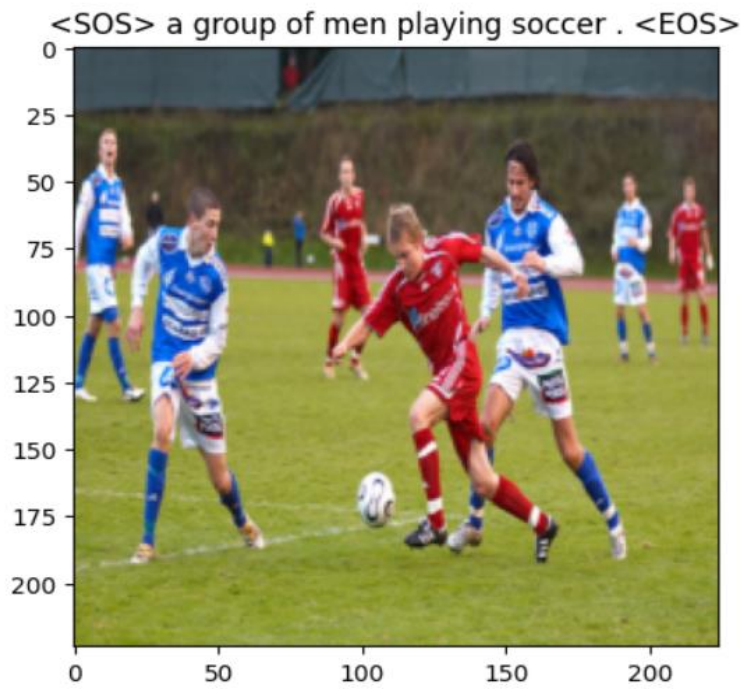
شکل 5. نمونه تصویر و کپشن تولید شده - **frozen ResNet 18**

#### 4-1. پیاده سازی شبکه Decoder-Encoder با شبکه ResNet18 با وزن های trainable

در این بخش مدل Resnet18 با قابلیت trainable بودن تمامی وزن ها برای آموزش روی دادگان استفاده می شود.



شکل 6. نمودار کاهش تابع هزینه بر روی داده های آموزشی و داده های اعتبارسنجی - وزن های trainable برای ResNet18

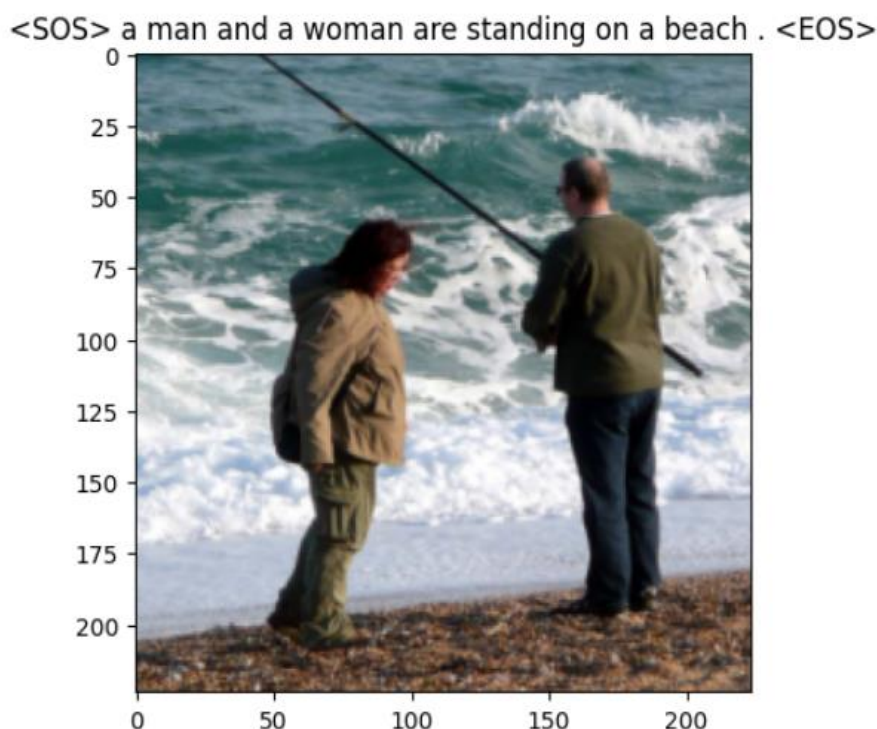


شکل 7. نمونه تصویر و کپشن تولید شده – trainable ResNet 18



شکل 8. نمونه تصویر و کپشن تولید شده – trainable ResNet 18





شکل 9. نمونه تصویر و کپشن تولید شده – trainable ResNet 18

## 5-1. بحث و نتیجه‌گیری

در این بخش به بحث و نتیجه‌گیری مدل‌های آموزش داده شده می‌پردازیم. ما ابتدا مدل را با وزن‌های pretrained از مدل resnet18 برای استخراج ویژگی استفاده نمودیم و در بخش 1-4 از مدل resnet18 که وزن‌ها قابل آموزش باشند استفاده کردیم. با توجه به نمودارهای کاهش هزینه می‌توان مشاهده نمود که کاهش هزینه در هر 2 مدل تقریباً یکسان بوده و در 10 اپیک آموزش loss به مقدار قابل توجهی کاهش پیدا کرده است. همان‌طور که مشاهده می‌شود مدل پتانسیل overfitting داشته است. ما برای مقابله با overfitting ابتدا نرخ یادگیری را کاهش دادیم. تکنیک دیگری که می‌توان برای جلوگیری از overfitting به کار گرفت early stopping است.

همان‌طور که در بخش نمایش تصاویر و کپشن‌های تولید شده توسط مدل مشاهده می‌شود نتایج و کپشن‌ها منطقی هستند و ارتباط معناداری بین کپشن‌ها و تصاویر برقرار است. البته مدل دوم بعضی جاها بهتر عمل کرده. برای مثال، در مقایسه کپشن‌های شکل 9 و شکل 5، می‌بینیم مدل دوم توانسته توضیحات مناسب‌تری بدهد. با کمی دقت در روی نمودارهای تابع هزینه درمیابیم که مدل دوم، اورفیت کمتری نسبت به مدل اول دارد و این می‌تواند دلیلی برای بهتر بودن مدل دوم نسبت به مدل اول بر روی داده‌های

test و validation باشد. از طرف دیگر مدل دوم تعداد 11 میلیون پارامتر بیشتر برای ترین کردن دارد و این باعث کندتر شدن روند در زمان آموزش شبکه می‌شود.



## ۱-۲. معماری LSTM و embedding

LSTM (حافظه کوتاه مدت بلندمدت) نوعی (RNN) است که بر محدودیت های خاص RNN های سنتی غلبه می کند. در اینجا برخی از مزایای LSTM نسبت به RNN ها آورده شده است:

### 1. Long-term dependencies: شبکه های LSTM برای حفظ وابستگی های طولانی

مدت در داده های sequential طراحی شده اند. RNN های معمولی برای یادگیری و حفظ اطلاعاتی که مربوط به time step های خیلی گذشته هستند، دچار مشکل می شوند لذا آنها برای کارهایی که نیاز به درک وابستگی های طولانی مدت دارند، مورد استفاده نیستند. معماری LSTM شامل سلول های حافظه (memory cells) و gate های پیچیده تر است که به آنها اجازه می دهد اطلاعات را در توالی های طولانی تر ضبط و منتشر کنند.

### 2. Handling vanishing gradients: RNN ها می توانند در طول زمان و در طی فزاینده

back propagation دچار مشکل vanishing gradient شوند. این امر مانع از فرآیند یادگیری، به خصوص برای توالی های طولانی می شود. LSTM با معرفی مفهوم سلول های حافظه و استفاده از gate های مختلف (input gate, forget gate, and output gate) برای کنترل جریان اطلاعات و گرادین ها در شبکه، این مشکل را کاهش می دهد. این به شبکه های LSTM کمک می کند تا اطلاعات را در توالی های طولانی تر به طور مؤثرتری یاد بگیرند و حفظ کنند.

### 3. Memory retention: سلول های حافظه (memory cells) در شبکه های LSTM به

آنها اجازه می دهند تا با استفاده از مکانیسم gate ها، اطلاعات را در هر مرحله زمانی به طور انتخابی حفظ یا فراموش کنند. این به LSTM ها امکان می دهد یاد بگیرند که کدام اطلاعات برای ذخیره کردن مهم هستند و کدامها را می توان دور انداخت.

4. **Flexibility and adaptability**: شبکه های LSTM بسیار انعطاف پذیر هستند و می توانند با انواع مختلف داده های متوالی سازگار شوند. آنها را می توان به آسانی با افزودن سلول های حافظه اضافی یا اصلاح مکانیسم های دروازه برای تطابق با وظایف خاص یا بهبود عملکرد، گسترش یا تغییر داد.

5. **Reduced training time**: در حالی که افزایش پیچیدگی شبکه های LSTM نیاز به محاسبات اضافی دارد، آنها اغلب در طول آموزش در مقایسه با RNN های پایه سریعتر همگرا می شوند. این به این دلیل است که شبکه های LSTM می توانند یاد بگیرند که بر روی اطلاعات مهم تمرکز کنند و ورودی های نامربوط را نادیده بگیرند و در نتیجه یادگیری کارآمدتری داشته باشند.

Word embedding تکنیکی است که در پردازش زبان طبیعی (NLP) برای نمایش کلمات یا عبارات به عنوان بردارهای متراکم و کم بعد در یک فضای برداری پیوسته استفاده می شود. هدف از word embedding، به دست آوردن روابط معنایی و نحوی بین کلمات است، و ماشین ها را قادر می سازد تا زبان انسان را بهتر درک و پردازش کنند.

دلایل اصلی استفاده از word embedding به شرح زیر است:

1. **بازنمایی معنایی (Semantic representation)**: word embedding به کلمات با معانی مشابه اجازه می دهد تا نمایش های برداری مشابهی داشته باشند. این تکنیک با به دست آوردن روابط معنایی، ماشین ها را قادر می سازد تا معنی و context کلمات را به روشی دقیق تر درک کنند. به عنوان مثال، در یک مدل word embedding به خوبی آموزش دیده، کلماتی مانند "شاه" و "ملکه" دارای بردارهایی هستند که در مقایسه با کلماتی مانند "سگ" یا "ماشین" به یکدیگر نزدیکتر هستند.

2. **کاهش ابعاد**: word embedding، بازنمایی کلمه های رمزگذاری شده به صورت one-hot با ابعاد بالا را به بردارهای متراکم با ابعاد پایین تر کاهش می دهند. این کاهش ابعاد، محاسبات را کمتر می کند و امکان مدیریت آسان تر واژگان بزرگ را فراهم می کند.

3. **Feature learning**: مدل‌های word embedding بر روی حجم زیادی از داده‌های متنی آموزش داده می‌شوند و در طول این فرآیند، مدل‌ها ویژگی‌های مفیدی را از متن یاد می‌گیرند. این ویژگی‌های آموخته شده را می‌توان برای کارهای downstream NLP مانند تجزیه و تحلیل احساسات، ترجمه ماشینی، پاسخگویی به سوالات و طبقه بندی متن مورد استفاده قرار داد.

روش‌های مختلفی برای تولید word embedding وجود دارد. در اینجا چند تکنیک متداول را به اختصار توضیح می‌دهیم:

1. **Word2Vec**: Word2Vec یک الگوریتم محبوب برای تولید word embedding است. این الگوریتم دو رویکرد زیر را ارائه می‌دهد: Continuous Bag-of-Words (CBOW) و Skip-gram. CBOW کلمه هدف را بر اساس کلمات در متن پیش بینی می‌کند، در حالی که Skip-gram کلمات متن را بر اساس یک کلمه هدف داده شده پیش بینی می‌کند. مدل‌های Word2Vec با استفاده از مجموعه بزرگی از داده‌های متنی آموزش داده می‌شوند و embedding حاصل، روابط بین کلمات را بر اساس الگوهای همزمانی آنها در داده‌های آموزشی ثبت می‌کنند.

2. **GloVe: Global Vectors for Word Representation**: یکی دیگر از روش‌های پرکاربرد برای word embedding است. این روش آمار جهانی از کاربرد همزمان (word co-occurrence) کلمات را با تکنیک فاکتورسازی برای تولید بردارهای کلمه ترکیب می‌کند. GloVe embedding بر روی داده‌های متنی بزرگ آموزش داده می‌شوند و روابط معنایی و نحوی بین کلمات را ثبت می‌کنند.

3. **BERT: Bidirectional Encoder Representations from Transformers**: یک مدل زبانی پیشرفته است که word embedding و همچنین بازنمایی کلمات متنی (contextualized) را تولید می‌کند. BERT روی مقادیر زیادی از داده‌های متنی بدون برچسب با استفاده از یک مدل زبانی از قبل آموزش داده شده تولید و سپس برای کارهای خاص fine-tuned می‌شود.

GloVe embeddings می‌توانند برای کلماتی که معانی متعددی دارند، مناسب باشند اما محدودیت‌های خاصی در گرفتن ظرایف معنایی ریز دارند.

GloVe embeddings بر روی مجموعه‌های بزرگی از متن آموزش داده می‌شوند، جایی که از آمار word co-occurrence برای تولید بردارهای کلمه استفاده می‌شود. این بدان معنی است که کلماتی که اغلب با هم در context مشابه ظاهر می‌شوند، بازنمایی‌های برداری مشابهی خواهند داشت. برای کلماتی با معانی چندگانه، تعبیه‌های GloVe تمایل دارند تا معنای متوسط یا کلی کلمه را بر اساس بافت کلی آن در داده‌های آموزشی نشان دهند.

در حالی که تعبیه‌های GloVe می‌توانند سطحی از چند معنایی را به تصویر بکشند، اما ممکن است بین معانی و senses مختلف یک کلمه به صورت دقیق تمایز قائل نشوند. این به این دلیل است که GloVe در طول فرایند آموزش، با تمام رخدادهای یک کلمه به عنوان یک موجودیت واحد رفتار می‌کند، و برای senses و حواس مختلف کلمه به طور جزئی تمایزی قائل نمی‌شود.

با این حال، توجه به این نکته مهم است که GloVe embeddings هنوز هم می‌توانند در بسیاری از کارهای NLP مفید باشند، به‌ویژه زمانی که معنی یا بافت کلی یک کلمه برای کار داده شده کافی باشد. به عنوان مثال، در تجزیه و تحلیل احساسات یا وظایف طبقه بندی متن، که در آن احساس کلی یا دسته بندی یک جمله مهم است، GloVe embeddings می‌توانند اطلاعات ارزشمندی را ارائه دهند.

اگر نیاز به تمایز دقیق‌تر از معنای کلمه باشد، تکنیک‌ها یا مدل‌های تخصصی تمایز حس کلمه که به صراحت برای گرفتن حواس کلمه طراحی شده‌اند، مانند WordNet یا جاسازی‌های حسی، ممکن است مناسب‌تر باشند.

## ۲-۲. پیش پردازش دادگان

در این بخش به اختصار اطلاعاتی کلی درباره دادگان مورد استفاده را ارائه می دهیم و فرایندهای پیش پردازش آنها را شرح خواهیم داد.

دیتاست مورد استفاده در این بخش شامل 5452 سطر و 3 ستون می باشد. هر سطر شامل یک سوال است که در ستون text قرار گرفته است. هر سوال نیز دارای دو برچسب است. یکی از آن دو تحت عنوان label-fine بوده که شامل مقادیر 0 تا 46 است و دیگری با نام label-coarse شامل مقادیر 0 تا 5 می باشد. برای فرایند پیش پردازش دادگان از دو کتابخانه tensorflow و NLTK استفاده کردیم. در زیر به شرح فرایندهای مورد استفاده در این مرحله خواهیم پرداخت.

1. Convert text to lowercase
  2. Word Tokenize: استفاده از nltk.word\_tokenize که یک جمله ورودی را به کلمات تشکیل دهنده آن split می کند.
  3. Remove Stop Words: حذف کلمات معروف به Stop Words مانند the, is, and و ...
  4. Stemming and Normalization: استخراج ریشه کلمات و حذف پیشوند، پسوند و غیره از آنها و سپس نرمالیزه کردن آنها با تبدیل بردار کلمات به جملات حاوی آنها.
  5. Tokenize with tensorflow: در این مرحله با استفاده از تابع Tokenizer از کتابخانه tensorflow، هر کلمه در جملات نرمالیزه شده در طی مراحل قبل را به مقداری عددی منحصر به فرد نگاشت می کند تا برای مدل ها قابل فهم شوند.
  6. Pad Text: هدف از این مرحله آن است که جملات را به اندازه های یکسان تبدیل کنیم تا برای استفاده در مدل آماده شوند. طول هر جمله را به اندازه طول بلندترین جمله و با استفاده از zero padding تبدیل می کنیم.
  7. Creating embedding matrix: در این مرحله ابتدا embedding پیشنهادی مقاله که در فایل glove.6B.300d.txt قرار دارد را داخل حافظه لود کرده و سپس ماتریس Embedding وزن مورد استفاده در مدل را با توجه به Embedding های glove.6B.300d و نیز کلمات موجود در متن دیتاست، آماده سازی می کنیم.
- glove.6B.300d embedding به یک مدل embedding از پیش پردازش شده با نام GloVe اشاره دارد که حاوی 6 میلیارد توکن و بردارهای با ابعاد 300 می باشد. (یعنی هر کلمه را به 300 بعد نگاشت می کند).

## 3-2. پیاده سازی طبقه بندی نیت

در این قسمت ابتدا به بررسی ساختار مدل ها پرداخته و سپس سعی می کنیم با استفاده از متریک های مختلف عملکرد هر کدام را بررسی کنیم.

```
Model: "First_model"
```

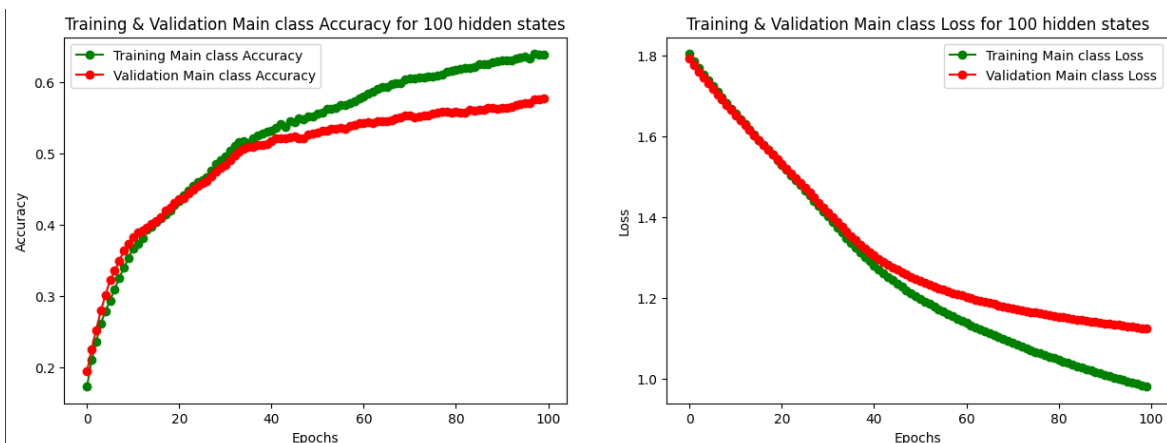
Layer (type)	Output Shape	Param #
input1 (InputLayer)	[(None, 18)]	0
embedding_17 (Embedding)	(None, 18, 300)	2072100
lstm_18 (LSTM)	(None, 25)	32600
flatten_24 (Flatten)	(None, 25)	0
mainclass (Dense)	(None, 6)	156

```
=====  
Total params: 2,104,856  
Trainable params: 32,756  
Non-trainable params: 2,072,100
```

شکل 10. ساختمان مدل اول

در شکل 10 ساختمان مدل اول نمایش داده شده است. در این ساختار ابتدا لایه input قرار دارد که اندازه بردار های ورودی آن به اندازه بلندترین جمله یعنی 18 است. سپس لایه embedding قرار دارد که همانطور که پیش از این نیز گفته شد حاوی 18 بردار 300 بعدی می باشد بدین معنی که هر کلمه ورودی را به یک فضای برداری 300 بعدی نگاشت می کند. سپس لایه LSTM قرار دارد که در این شکل تعداد hidden state های آن 25 است. نکته قابل توجه این است که در این مدل ما تنها به خروجی سلول آخر لایه LSTM نیاز داریم به همین علت خروجی این لایه یک بردار 25 بعدی است. سپس خروجی لایه LSTM را برای دادن به لایه بعدی یعنی لایه flat, linear می کنیم و در نهایت به یک لایه dense با 6 نورون (به اندازه تعداد main class) و با تابع فعالساز softmax می دهیم. تعداد کل پارامتر ها و نیز پارامتر های trainable و non-trainable نیز در شکل مشخص است. non-trainable params در واقع همان وزن های لایه embedding هستند که به صورت transform learning از وزن های glove.6B.300d استفاده کردیم.

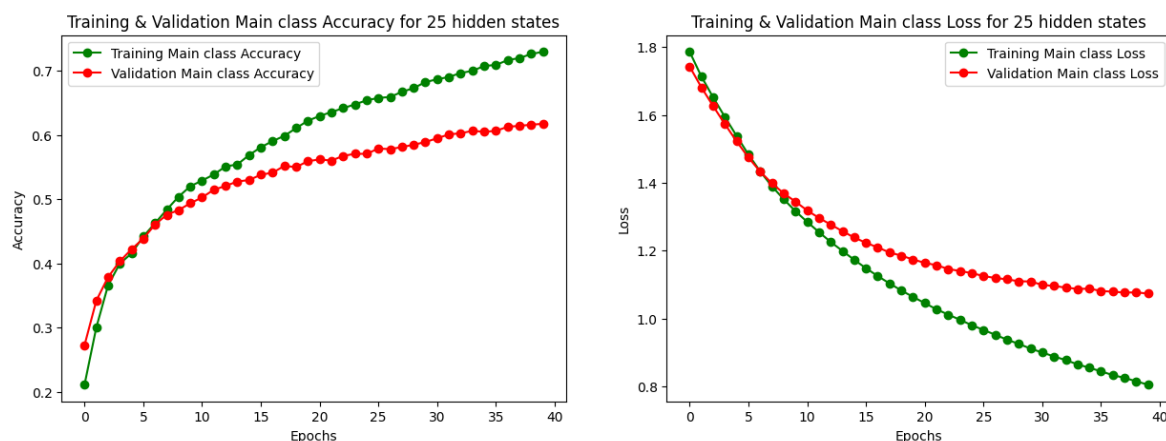
نهایتاً مدل را با تابع `loss`، `"sparse_categorical_crossentropy"` و `Adam optimizer` کامپایل کرده و در طی 100 اپیاک و با `validation_split = 0.4` آموزش دادیم. در ادامه نتایج کسب شده را بررسی خواهیم کرد.



شکل 11. نمودار `loss` و `accuracy` برای داده های `train` و `validation` برای 100 hidden state

شکل بالا، نمودار `loss` و `accuracy` را برای داده های `train` و `validation` در طی 100 اپیاک و برای تعداد 100 hidden states نشان می دهد. برای این مدل از `learning_rate = 1e-5` استفاده کردیم. همانطور که مشخص است در طی هر اپیاک دقت مدل افزایش و `loss` آن کاهش یافته است تا در نهایت در اپیاک آخر به مقادیر زیر رسیده است :

**loss: 0.9803 - accuracy: 0.6396 - val\_loss: 1.1248 - val\_accuracy: 0.5777**



شکل 12. نمودار `loss` و `accuracy` برای داده های `train` و `validation` برای 25 hidden state

در شکل 12 ، نمودار `loss` و `accuracy` را برای داده های `train` و `validation` در طی 40 اپاک و برای تعداد `hidden states` 25 نشان می دهد. برای این مدل از `learning_rate = 1e-4` استفاده کردیم. همانطور که مشخص است در این مدل نیز طی هر اپاک دقت مدل افزایش و `loss` آن کاهش یافته است تا در نهایت در اپاک آخر به مقادیر زیر رسیده است :

**loss: 0.8047 - accuracy: 0.7297 - val\_loss: 1.0734 - val\_accuracy: 0.6176**

	precision	recall	f1-score	support
0	0.55	0.95	0.70	138
1	0.69	0.48	0.57	94
2	0.00	0.00	0.00	9
3	0.65	0.55	0.60	65
4	0.88	0.50	0.64	113
5	0.62	0.59	0.61	81
accuracy			0.63	500
macro avg	0.57	0.51	0.52	500
weighted avg	0.67	0.63	0.62	500

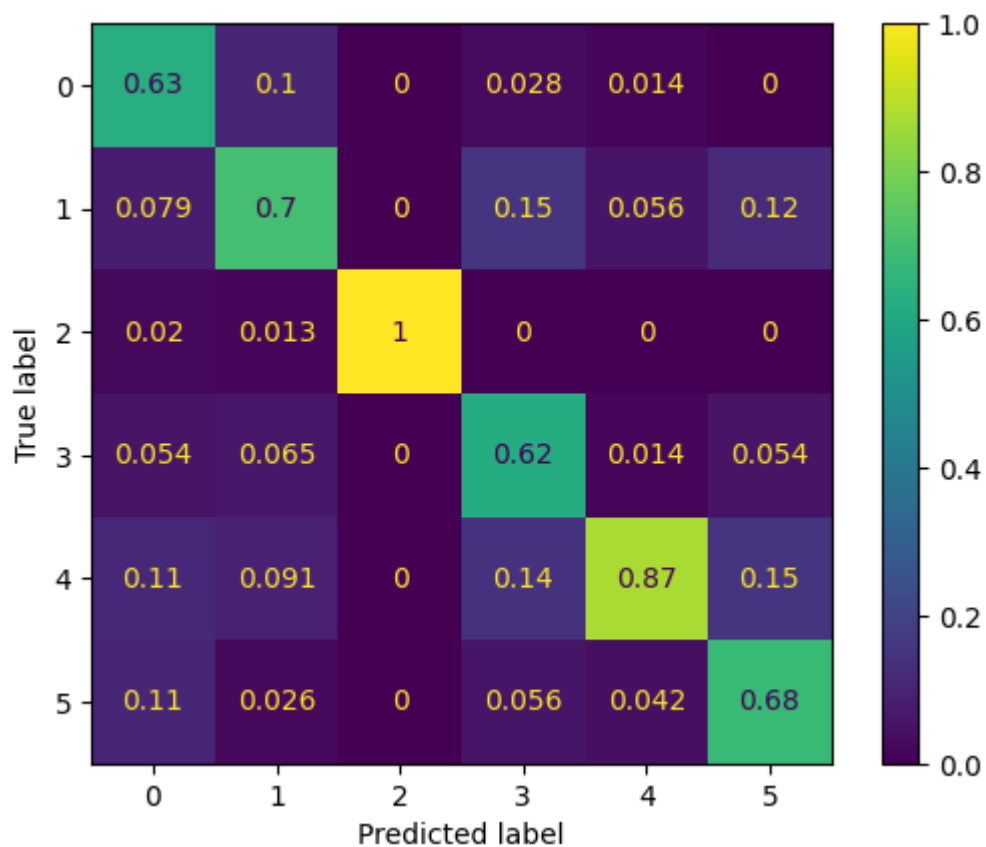
شکل 13. نتایج پیاده سازی مدل اول بر داده های **test** برای **hidden state 25**

	precision	recall	f1-score	support
0	0.52	0.96	0.67	138
1	0.59	0.40	0.48	94
2	0.00	0.00	0.00	9
3	0.56	0.43	0.49	65
4	0.92	0.43	0.59	113
5	0.59	0.56	0.57	81
accuracy			0.59	500
macro avg	0.53	0.46	0.47	500
weighted avg	0.63	0.59	0.57	500

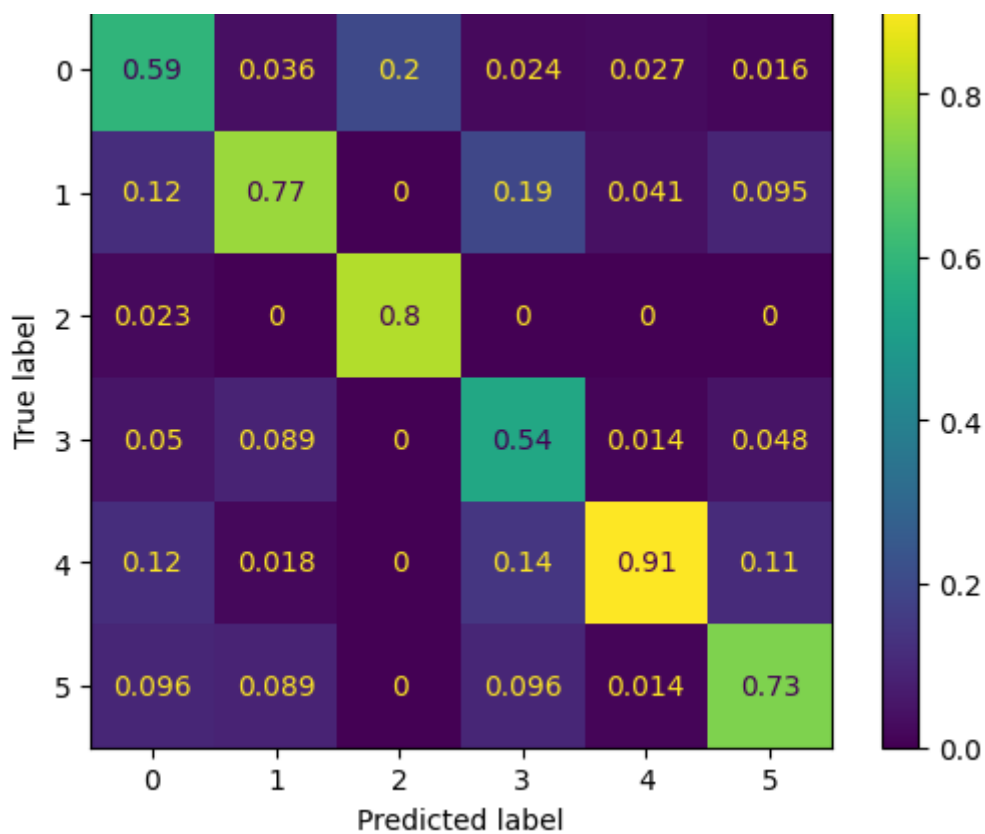
شکل 14. نتایج پیاده سازی مدل اول بر داده های **test** برای **100 hidden state**

شکل بالا نتایج حاصل از مدل با 25 `hidden states` و شکل پایین با 100 عدد است. در هر دو مورد دقت حاصل با دقت مقاله بسیار فاصله دارد. در کل مدل پیشنهادی با تعداد 25 `hidden states` توانسته به دقت های بهتری به نسبت دیگر مدل برسد.





شکل 15. ماتریس پراکندگی مدل اول با 25 hidden state



شکل 16. ماتریس پراکندگی مدل اول با 100 hidden state

همچنین ماتریس آشفتگی مربوط به مدل با هر دو تعداد LSTM hidden states در شکل های بالا نمایش داده شده است. در هر دو مورد قطر اصلی به نسبت سایریرین مقادیر بیشتری دارد که این نشانه خوبی است. مدل اول بهترین عملکرد را با جملات کلاس 2 داشته است و توانسته به طور کامل همه آنها را پیش بینی کند. مدل دوم بهترین عملکرد را برای کلاس 4 داشته است. مقایسه ماتریس پراکندگی دو مدل نشان می دهد که هر مدل در پیش بینی برخی کلاس ها از دیگری بهتر عمل کرده است.

```
#Building the base model
NUM_HIDDEN_STATES = 25
input_layer = Input(shape=(max_length+1,), name = 'input1')
embedding = Embedding(vocab_size, 300, weights=[embedding_matrix], input_length =max_length+1 ,trainable=False)(input_layer)
lstm = LSTM(NUM_HIDDEN_STATES , return_sequences=True)(embedding)

#subclass prediction
flat = Flatten()(lstm[:,-1])
first_predictions = Dense(num_sub_class,activation='softmax',name='subclass')(flat)
#mainclass prediction
flat_2 = Flatten()(lstm[:,-2])
second_predictions = Dense(num_main_class,activation='softmax',name='mainclass')(flat_2)

second_model = Model(inputs=input_layer, outputs=[first_predictions , second_predictions] , name="Second_model")
opt = optimizers.Adam(learning_rate=1e-4)
second_model.compile(optimizer=opt,
                    loss={
                        'subclass': 'sparse_categorical_crossentropy',
                        'mainclass': 'sparse_categorical_crossentropy'
                    },
                    metrics={
                        'subclass': 'accuracy',
                        'mainclass': 'accuracy'
                    })
second_model.summary()
```

Model: "Second\_model"

Layer (type)	Output Shape	Param #	Connected to
input1 (InputLayer)	(None, 19)	0	[]
embedding (Embedding)	(None, 19, 300)	2072100	['input1[0][0]']
lstm (LSTM)	(None, 19, 25)	32600	['embedding[0][0]']
tf.__operators__.getitem (SlicingOpLambda)	(None, 25)	0	['lstm[0][0]']
tf.__operators__.getitem_1 (SlicingOpLambda)	(None, 25)	0	['lstm[0][0]']
flatten (Flatten)	(None, 25)	0	['tf.__operators__.getitem[0][0]']
flatten_1 (Flatten)	(None, 25)	0	['tf.__operators__.getitem_1[0][0]']
subclass (Dense)	(None, 47)	1222	['flatten[0][0]']
mainclass (Dense)	(None, 6)	156	['flatten_1[0][0]']

```
=====
Total params: 2,106,078
Trainable params: 33,978
Non-trainable params: 2,072,100
```

شکل 17. ساختمان مدل دوم

شکل 17، نمایش دهنده ساختمان مدل دوم هست. تفاوت این مدل با مدل قبلی در این است که در اینجا ما دو خروجی main class و sub class داریم که برای آنها از خروجی آخر و یکی مانده به آخر لایه LSTM استفاده می کنیم. لذا ما به خروجی همه cell های lstm نیاز داریم و برای همین return\_sequence را برای لایه LSTM برابر True قرار می دهیم تا خروجی cell های داخلی را نیز داشته باشیم. سپس خروجی cell های مورد نظر را به لایه های dense با تعداد نوروں برابر با تعداد sub class و main class داده و از تابع softmax عبور می دهیم. مدل ما دو خروجی دارد که برای هر دو از "sparse\_categorical\_crossentropy" به عنوان تابع loss استفاده می کنیم. تعداد پارامترهای trainable و non-trainable مدل نیز در شکل نمایش داده شده است. نکته قابل توجه در این معماری آن است که برای ورودی شبکه، در آخر هر آرایه یک 0 اضافه می کنیم تا از خروجی LSTM آن برای طبقه بندی subclass استفاده کنیم. Hyperparameter های استفاده شده تعداد hidden state لایه LSTM و نیز learning rate برای Adam optimizer است که در جدول زیر مقادیر استفاده شده را آورده ایم:

جدول 1. مقادیر learning rate برای lstm با مقادیر 25,100 hidden states

# Hidden states for LSTM layer	Learning rate
25	3e-5
100	e-51

Hyperparameter های استفاده شده با آزمون و خطا به دست آمدند و آنهایی که بهترین دقت را برای مدل طراحی شده، داشتند به عنوان پارامترهای نهایی استفاده شدند. اگر چه نهایتاً نتوانستیم به دقت گزارش شده در مقاله دست پیدا کنیم و چون در مقاله درباره پارامترهای استفاده شده چیزی گفته نشده، نهایتاً مدل را با بهترین نتیجه ای که توانستیم به آن دست پیدا کنیم گزارش کردیم.

گزارش های مربوط به عملکرد مدل دوم برای 25 hidden states به شرح زیر می باشد:

**train\_subclass\_loss:** 2.2034

**train\_mainclass\_loss:** 0.9109

**train\_subclass\_accuracy:** 0.3925

**train\_mainclass\_accuracy:** 0.6796

**val\_subclass\_loss:** 2.4075

**val\_mainclass\_loss:** 1.0905

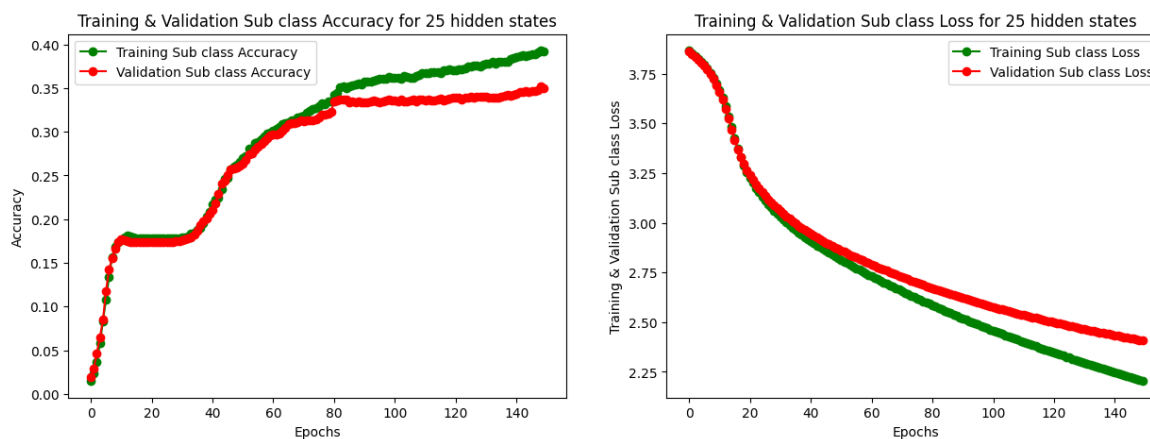
**val\_subclass\_accuracy:** 0.3494

**val\_mainclass\_accuracy:** 0.5924

همانطور که گفتیم دقت های به دست آمده از دقت های مقاله فاصله دارد ولی نکته آن است که مدل ما در این مرحله و با این دقت ها overfit نشده است و برای به دست آوردن دقت های بیشتر روی training data مدل overfit می شد.



شکل 18. نمودار loss و accuracy برای مدل دوم کتگوری main class روی داده های train و val ، با 25 hidden state



شکل 19. نمودار loss و accuracy برای مدل دوم کتگوری sub class روی داده های train و val ، با 25 hidden state

دو شکل بالا نمودار های accuracy و loss برای داده های train و val برای دو گروه subclass و main class برای تعداد 25 hidden states را نشان می دهد . همانطور که مشخص است در اجرای الگوریتم برای 150 اپیاک ، نمودار accuracy برای هر دو کلاس تقریباً صعودی و نمودار loss نزولی است . این نکته که مدل در نهایت overfit نشده است نیز در شکل مشخص است چرا که اختلاف بین نتایج train و val زیاد نیست.

	precision	recall	f1-score	support
0	0.52	0.93	0.67	138
1	0.60	0.39	0.47	94
2	0.00	0.00	0.00	9
3	0.59	0.52	0.55	65
4	0.87	0.36	0.51	113
5	0.57	0.62	0.60	81
accuracy			0.58	500
macro avg	0.53	0.47	0.47	500
weighted avg	0.62	0.58	0.56	500

شکل 20. نتایج حاصل از مدل دوم بر روی داده های test برای main class با 25 hidden state

accuracy			0.40	500
macro avg	0.07	0.08	0.06	500
weighted avg	0.29	0.40	0.30	500

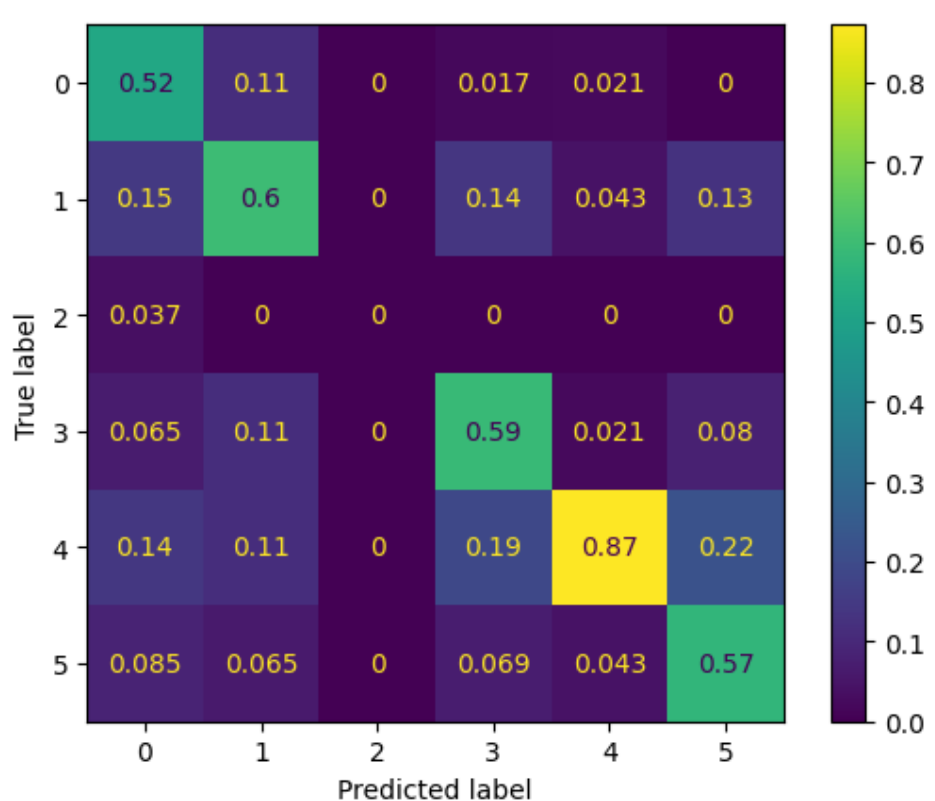
شکل 21. نتایج حاصل از مدل دوم بر روی داده های test برای sub class با 25 hidden state

سپس مدل را بر روی داده های test آزمایش کرده و از متریک هایی که در دوشکل بالا دیده می شود برای ارزیابی عملکرد مدل استفاده کردیم. بر طبق آنچه در نتایج حاصل دیده می شود ، عملکرد مدل در پیش بینی subclass به نسبت main class بدتر بوده است و یکی از دلایل آن میتوان تعداد زیاد کلاس ها (47) موجود در این کتگوری باشد که برای train شدن به دیتای بیشتری نیاز دارد . برای پیش بینی main class، مدل بهترین عملکرد را در پیش بینی کلاس 0 و بدترین عملکرد را در پیش بینی کلاس 2 داشته است . برای کلاس های موجود در main class به علت آنکه تعداد آنها کم بود مقایسه متریک های آن راحت تر است ولی برای sub class به علت تعداد زیاد کلاس ها این مقایسه راحت نیست و به همین علت فقط گزارش نهایی آن را آورده و کلاس به کلاس بررسی نکردیم. (نتایج حاصل از آن در بخش بعدی و در ماتریس آشفتگی بهتر قابل قیاس است.)

برای یادآوری !

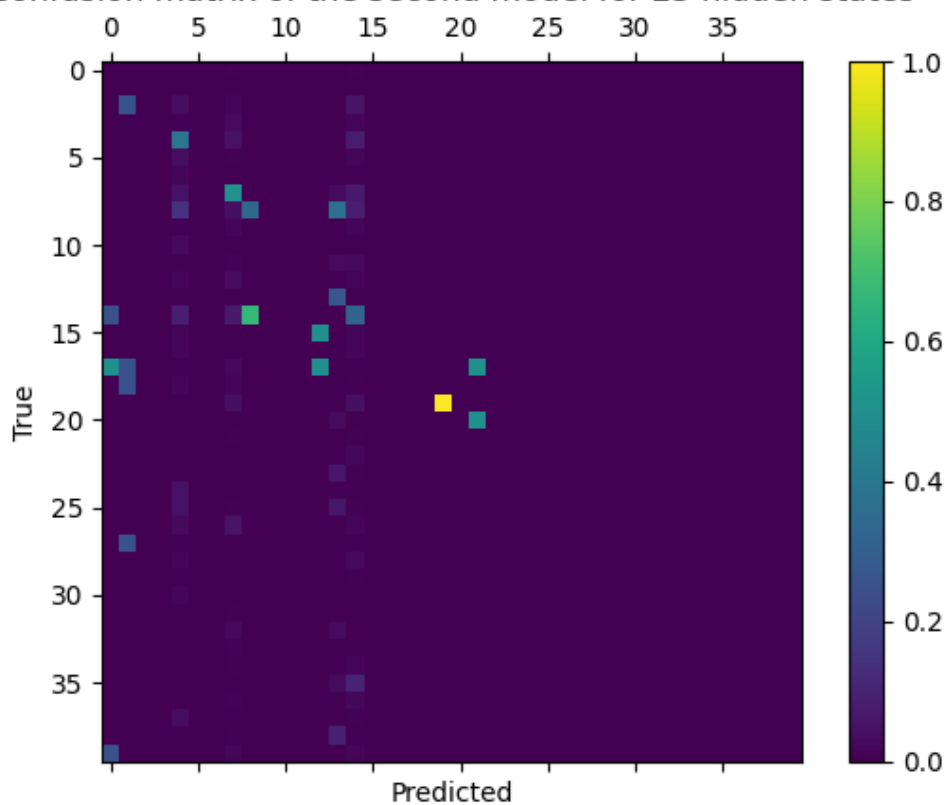
Precision: دقت پیش‌بینی‌های مثبت را اندازه‌گیری می‌کند،

Recall: کامل بودن پیش‌بینی‌های مثبت را اندازه‌گیری می‌کند.



شکل 22. ماتریس پراکندگی مدل دوم با 25 hidden state بر روی داده های test برای main class

Confusion matrix of the second model for 25 hidden states



شکل 23. ماتریس پراکندگی مدل دوم با 25 hidden state بر روی داده های test برای sub class

یکی دیگر از روش های مورد استفاده برای مقایسه عملکرد مدل ها استفاده از ماتریس آشفتگی آنهاست که در شکل 22,23 نمایش داده شده است. با نگاه به این ماتریس برای خروجی main class میتوان دید که مدل در پیش بینی کلاس 4 عملکرد به نسبت خوبی داشته است چرا که بیشترین تعداد پیش بینی درست آن برای این کلاس بوده است . یک نکته که در این ماتریس مشخص می شود و قبلا در سایر متریک ها متوجه آن نشدیم آن است که تعداد داده های با  $\text{main class} = 2$  به طور کلی در دیتاست کم بوده است و این غنی نبودن دیتا ست یکی از دلایل عملکرد ضعیف مدل است . البته این مسئله می تواند ناشی از تقسیم اشتباه داده ها به  $\text{train}$  و  $\text{test}$  نیز باشد که برای بررسی آن می توانیم از Splitting Strata استفاده کنیم تا داده های هر کلاس به نسبت مساوی تقسیم شوند . در خصوص ماتریس مربوط به sub class ، اول اینکه به علت تعداد بالای کلاس ها یکی از بهترین روش ها برای مقایسه عملکرد مدل همین استفاده از ماتریس آشفتگی است. ثانيا همانطور که دیده میشود تعداد پیش بینی های درست مدل زیاد نبوده و تنها برای یک کلاس توانسته به تعداد خوبی پیش بینی درست داشته باشد .

در ادامه همین نتایج را برای تعداد 100 hidden states بررسی خواهیم کرد.

**train\_subclass\_loss:** 1.8541

**train\_mainclass\_loss:** 0.8229

**train\_subclass\_accuracy:** 0.5139

**train\_mainclass\_accuracy:** 0.7083

**val\_subclass\_loss:** 2.1766

**val\_mainclass\_loss:** 1.0759

**val\_subclass\_accuracy:** 0.4269

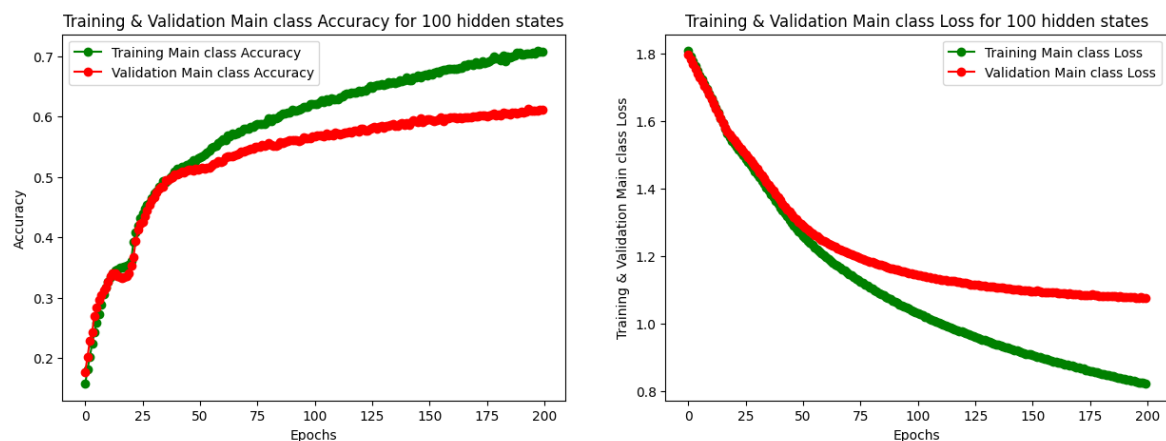
**val\_mainclass\_accuracy:** 0.6112



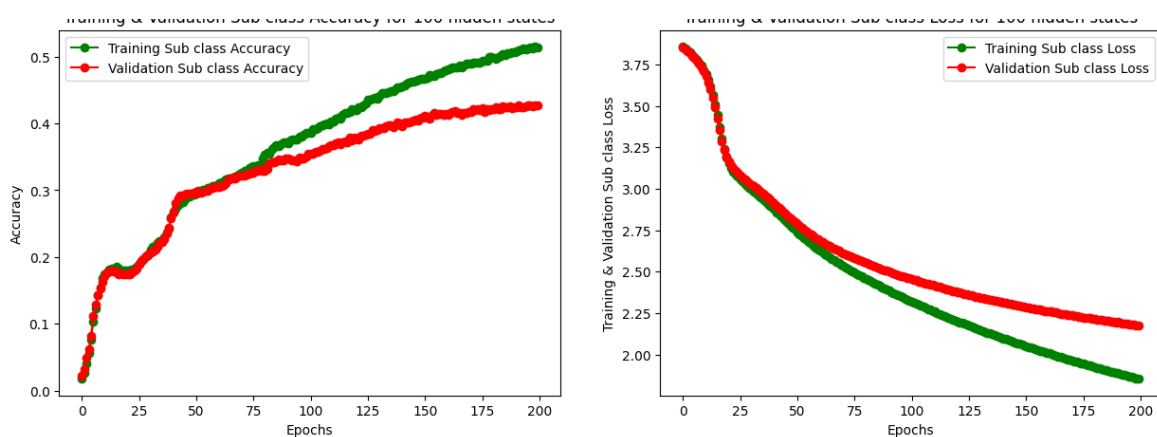
LSTM برای لایه 100 hidden state و 25 جدول 2. مقایسه عملکرد مدل دوم برای تعداد

	25 LSTM hidden states	100 LSTM hidden states
<b>Train subclass loss</b>	2.2034	<b>1.8541</b>
<b>Train mainclass loss</b>	0.9109	<b>0.8229</b>
<b>Train subclass accuracy</b>	0.3925	<b>0.5139</b>
<b>Train mainclass accuracy</b>	0.6796	<b>0.7083</b>
<b>Validation subclass loss</b>	2.4075	<b>2.1766</b>
<b>Validation mainclass loss</b>	1.0905	<b>1.0759</b>
<b>Validation subclass accuracy</b>	0.3494	<b>0.4269</b>
<b>Validation mainclass accuracy</b>	0.5924	<b>0.6112</b>

در جدول بالا مقایسه عملکرد مدل برای تعداد 25 و 100 hidden states برای لایه LSTM مشاهده می کنید. همانگونه که مشخص است عملکرد مدل برای تعداد 100 hidden states با اختلاف خوبی به نسبت دیگری بهتر است. مدل مجددا نتوانسته به دقت های مقاله برسد اما هم در train و هم در validation دقت مدل بهبود یافته است.



شکل 24. نمودار **loss** و **accuracy** برای مدل با **hidden state 100** با داده های **train** و **val** برای **main class**



شکل 25. نمودار **loss** و **accuracy** برای مدل با **hidden state 100** با داده های **train** و **val** برای **sub class**

دو نمودار بالا دقت و **loss** مدل با 100 hidden states را در طی 200 اپیاک برای داده های **train** و **test** نشان می دهد. در اینجا نیز در طی اپیاک ها دقت بهبود یافته و **loss** کم شده البته به نسبت مدل قبلی در نهایت اختلاف بین **train** و **validation** بیشتر است ولی به اندازه ای نیست که بتوان گفت مدل **overfit** شده است.

	precision	recall	f1-score	support
0	0.53	0.95	0.68	138
1	0.66	0.44	0.53	94
2	0.00	0.00	0.00	9
3	0.56	0.51	0.53	65
4	0.92	0.42	0.57	113
5	0.60	0.59	0.60	81
accuracy			0.60	500
macro avg	0.55	0.48	0.48	500
weighted avg	0.65	0.60	0.58	500

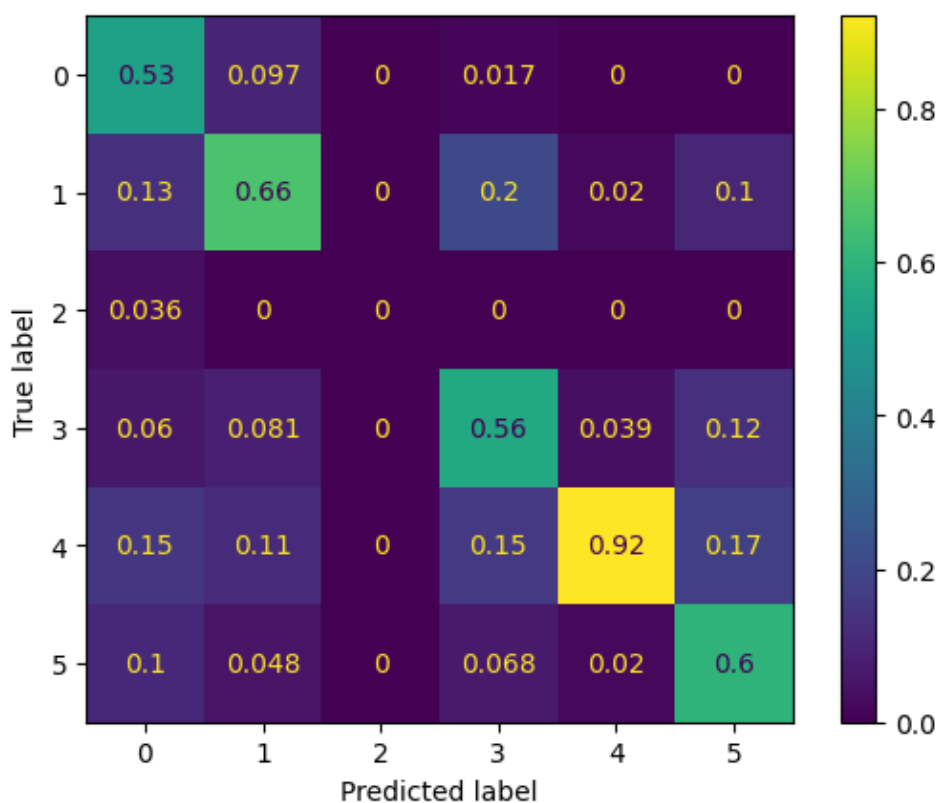
شکل 26. نتایج حاصل از مدل دوم بر روی داده های **test** برای **main class** با **hidden state 100**

accuracy			0.47	500
macro avg	0.15	0.13	0.13	500
weighted avg	0.38	0.47	0.39	500

شکل 27. نتایج حاصل از مدل دوم بر روی داده های **test** برای **sub class** با **hidden state 100**

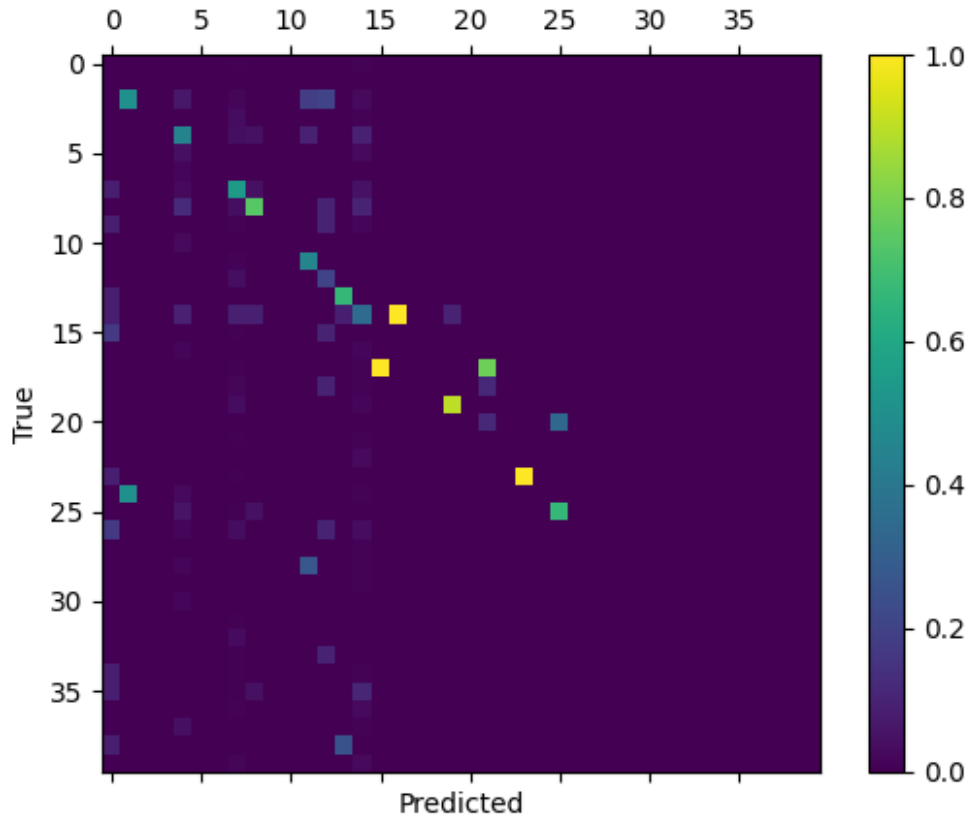
با مقایسه نتایج حاصل از اجرای مدل دوم بر روی داده های **test** که در اشکال بالا نمایش داده شده است ، می توان دید که برای **main class** مجددا کلاس 0 بهترین نتایج را داشته است و دوباره به طور کلی نتایج حاصل از اجرای مدل بر روی **subclass** نتایج بدتری به نسبت **main class** داشته است.

همچنین ماتریس آشفتگی ای مدل نیز در اشکال 28,29 نمایش داده شده است. برای داده های **main class** نتایج حاصل از این مدل مشابه نتایج حاصل از مدل قبلی می باشد ولی برای داده های **subclass** نتایج حاصل از این مدل به مراتب بهتر است چرا که تعداد داده های قرار گرفته در قطر اصلی بیشتر است اگرچه باید به این نکته نیز توجه شود که برخی از خانه های خارج از قطر اصلی نیز تعداد زیادی را به خود اختصاص داده اند که این نشان از ضعف مدل در پیش بینی داده های با آن لیبل ها است .



شکل 28. ماتریس پراکندگی مدل دوم با 100 **hidden state** بر روی داده های **test** برای **main class**

Confusion matrix of the second model for 100 hidden states



شکل 29. ماتریس پراکندگی مدل دوم با 100 **hidden state** بر روی داده های **test** برای **sub class**

## 4-2. پیاده سازی مدل Responder

در نهایت یک مدل responder بر روی مدل های طراحی شده سوار کرده تا نحوه استفاده از intent classification در پاسخگویی به سوالات را بررسی کنیم. بدین منظور از دیتاست QA\_data که در فایل تمرین قرار دارد استفاده کردیم. مراحل preprocessing اعمال شده بر روی این دیتاست نیز مشابه مراحل اعمال شده بر روی دیتاست قبلی که پیش تر توضیح دادیم است.

Model: "Responder"

Layer (type)	Output Shape	Param #	Connected to
input_layer (InputLayer)	[(None, 12)]	0	[]
embedding (Embedding)	(None, 12, 300)	287400	['input_layer[0][0]']
lstm (LSTM)	(None, 12, 100)	160400	['embedding[0][0]']
tf.__operators__.getitem_2 (SlicingOpLambda)	(None, 11, 300)	0	['embedding[0][0]']
tf.__operators__.getitem_1 (SlicingOpLambda)	(None, 100)	0	['lstm[0][0]']
tf.__operators__.getitem (SlicingOpLambda)	(None, 100)	0	['lstm[0][0]']
responder_lstm (Bidirectional)	(None, 11, 200)	320800	['tf.__operators__.getitem_2[0][0]', 'tf.__operators__.getitem_1[0][0]', 'tf.__operators__.getitem[0][0]', 'tf.__operators__.getitem_1[0][0]', 'tf.__operators__.getitem[0][0]']
Dense_layer (Dense)	(None, 11, 781)	156981	['responder_lstm[0][0]']

=====  
Total params: 925,581  
Trainable params: 638,181  
Non-trainable params: 287,400

شکل 30. ساختمان مدل responder

در شکل 30 ساختمان مدل responder را مشاهده می کنید. در لایه input جملات pad شده به عنوان ورودی شبکه داده می شود و سپس مشابه آنچه در مدل دوم داشته جریان داده های در شبکه ادامه میابد. نکته قابل توجه استفاده از یک لایه bidirectional lstm در شبکه طراحی شده است. خروجی cell آخر و یکی مانده به آخر لایه lstm که قبلا از آنها برای پیش بینی subclass و main class استفاده می کردیم، در این شبکه با یکدیگر ترکیب شده و به عنوان initial state به لایه bidirectional LSTM داده می شود. این همان استفاده از intent classification در پاسخ به سوالات است. توجه شود که تعداد hidden states های لایه bidirectional، 100 عدد است و چون خروجی همه cell های آنرا استخراج می کنیم نهایتا یک

خروجی (11\*200) خواهیم داشت. در آخر یک لایه dense قرار می‌دهیم با تعداد نوروں به اندازه تعداد کلمات خاص موجود در دیتاست. وظیفه این لایه این است که خروجی bidirectional LSTM را از 200 به 781 نگاشت کند تا بتوانیم از آن برای پیش‌بینی پاسخ استفاده کنیم.

تعداد پارامترهای trainable و non trainable این شبکه را نیز می‌توان مشاهده کرد.

```
train_ans_data = []
for sent in train_ans:
    temp = []
    for word in sent:
        b = np.zeros(vocab_size_ans)
        b[int(word)] = 1
        temp.append(b)
    train_ans_data.append(temp)

train_ans_data = np.array(train_ans_data)
print(train_ans_data.shape)

(500, 11, 781)
```

شکل 31. نحوه reshape کردن داده‌های جواب

با توجه به خروجی شبکه طراحی شده، نیاز داریم تا برای استفاده از gradient decent، داده‌های مربوط به جواب را نیز reshape کرده و با استفاده از روش one-hot encoding، آنها را به فضای 781 نگاشت کنیم. کد مربوط به این بخش نیز در شکل 31 مشاهده می‌شود.

نهایتاً شبکه را بر روی دیتاست داده شده train می‌کنیم و بعد سوال‌های داده شده در متن سوال را از آن می‌پرسیم. شکل 32 نحوه تبدیل خروجی شبکه به متن پاسخ را نشان می‌دهد. برای اینکار ابتدا یک reversed map از کلمات tokenize شده در قسمت قبل ایجاد می‌کنیم که در آن key‌ها همان اعداد یونیک اختصاص داده شده به هر کلمه و value‌ها کلمات مربوط به آن عدد هستند. سپس با استفاده از این map اعداد خروجی را به کلمات مرتبط نگاشت می‌کنیم.

```

# Creating a reverse dictionary
reverse_word_map = dict(map(reversed, tokenizer_ans.word_index.items()))

for idi , ans in enumerate(answer__):

    # Function takes a tokenized sentence and returns the words
    def sequence_to_text(list_of_indices):
        list_of_indices = list_of_indices[list_of_indices != 0]
        if list_of_indices.size == 0:
            return ""
        # Looking up words in dictionary
        words = [reverse_word_map.get(letter) for letter in list_of_indices]
        return(words[0])

    # Creating texts
    my_texts = list(map(sequence_to_text, ans))
    str_list = list(filter(None, my_texts))

    print("Question: " , test_questions[idi])
    if len(str_list) == 0:
        print("No Answer")
    else:
        print(str_list[0])

    print("*****")

```

شکل 32. کد مربوط به نحوه تبدیل خروجی شبکه به پاسخ های متنی

خروجی یا همان پاسخ های هر سوال در شکل 33 نمایش داده شده اند. همانطور که انتظار می رفت ، پاسخ ها کاملاً نامرتب با متن سوالات است که این با توجه به حجم بسیار پایین دیتاست برای چنین تسکی ، قابل توجیه است.

```
Question: How many people speak French?
caffeine
*****
Question: What day is today?
9
*****
Question: Who will win the war?
narrow
*****
Question: Who is Italian first minister?
blue
*****
Question: When World War II ended?
dollar
*****
Question: When Gandhi was assassinated?
metal
*****
```

شکل 33. پاسخ های شبکه طراحی شده به سوال های مطرح شده