



پردازش زبان‌های طبیعی

تمرین سوم

محمد جواد رنجبر

۸۱۰۱۰۱۱۷۳

بهار ۱۴۰۲

## Contents

٤	..... توضیح فایل‌ها
٥	..... بخش ١
٥	..... One Hot encoding
٦	..... Glove
٧	..... Word2vec
٨	..... بخش ٢
٨	..... LSTM
١٠	..... GRU
١٢	..... مقایسه دو مدل

۵	شکل ۱ توزیع داده‌ها.....
۶	شکل ۲ عملکرد rnn برای one hot .....
۶	شکل ۳ ماتریس درهم‌ریختگی برای مدل rnn و one hot .....
۷	شکل ۴ عملکرد مدل rnn برای glove .....
۷	شکل ۵ ماتریس درهم‌ریختگی مدل RNN با glove .....
۸	شکل ۶ عملکرد مدل rnn برای word2vec .....
۸	شکل ۷ ماتریس درهم‌ریختگی مدل rnn برای word2vec .....
۹	شکل ۸ عملکرد lstm برای one hot .....
۹	شکل ۹ ماتریس درهم‌ریختگی برای مدل lstm و one hot .....
۱۰	شکل ۱۰ عملکرد lstm برای glove .....
۱۰	شکل ۱۱ ماتریس درهم‌ریختگی LSMT با glove .....
۱۱	شکل ۱۲ عملکرد مدل gru برای one hot .....
۱۱	شکل ۱۳ ماتریس درهم‌ریختگی مدل gru و one hot .....
۱۲	شکل ۱۴ عملکرد gru برای glove .....
۱۲	شکل ۱۵ ماتریس درهم‌ریختگی gru با glove .....

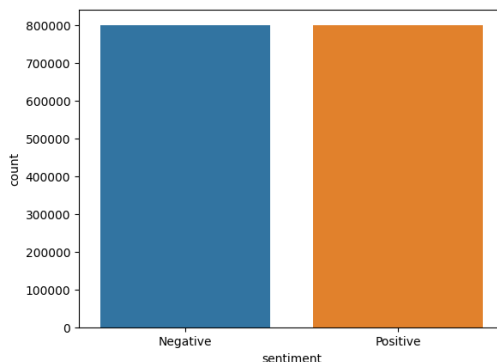
## توضیح فایل‌ها:

به واسطه سنگین شدن فایل colab تمرین در دو فایل انجام شده است که در فایل CA3\_Q\_one\_hot.ipynb بخش‌های وابسته به one hot encoding و در فایل CA2\_glove\_w2v.ipynb بقیه بخش‌های تمرین قرار داده شده است.

## بخش ۱:

ابتدا پیش پردازش های مناسب برای این توییت ها را انجام می دهیم که شامل موارد زیر می باشد:

- حذف stop word
  - حذف لینک ها و علامت های اضافی مانند # و @ (user ها و هشتگ نگه داشته می شوند).
  - از آنجا که emoji ها برای وظیفه، sentiment می توانند کاربردی باشند. آن ها را با اسامی emoji جایگزین می کنیم.
  - Lemmetazation
  - تبدیل label ها از ۰ و ۴ به ۰ و ۱ برای آسانتر شدن پردازش
  - برای حل مشکل هم اندازه بودن توییت ها دو روش وجود دارد:
    - استفاده از padding برای توییت ها که در کد پیاده سازی شده است اما سرعت آموزش را به شدت پایین می آورد.
    - کنار هم گذاشتن توییت های هم اندازه در هر batch
- همچنین نمودار توزیع داده های مختلف به شکل زیر می باشد، که تعداد داده ها در هر دو کلاس تقریباً برابر می باشد.

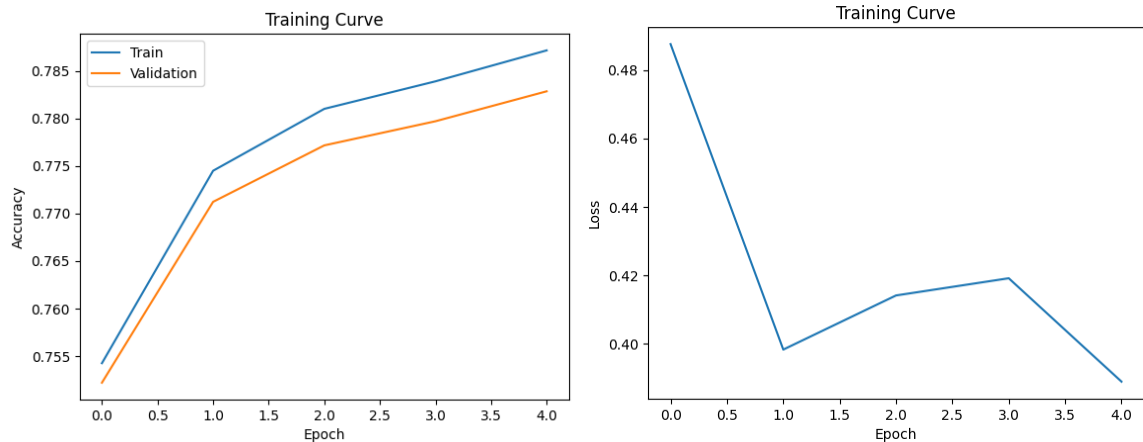


شکل ۱ توزیع داده ها

برای اینکه بتوانیم از داده ها استفاده کنیم و دچار مشکلات حافظه نشویم از ۱/۶ میلیون توییت تعداد ۶۰۰۰۰ را به صورت تصادفی انتخاب می کنیم. حال با استفاده از white space کلمات را tokenize می کنیم و با استفاده از embedding های مختلف، مدل RNN را آموزش می دهیم.

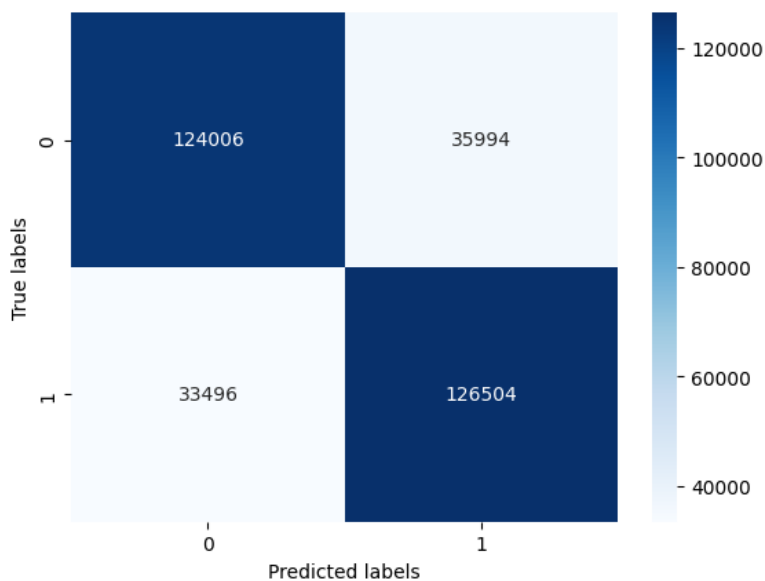
### One Hot encoding:

این نوع embedding بسیار سنگین بوده و اصلاً مناسب نیست، با این حال نتایج مدل به شکل زیر خواهد بود:



شکل ۲ عملکرد *rnn* برای *one hot*

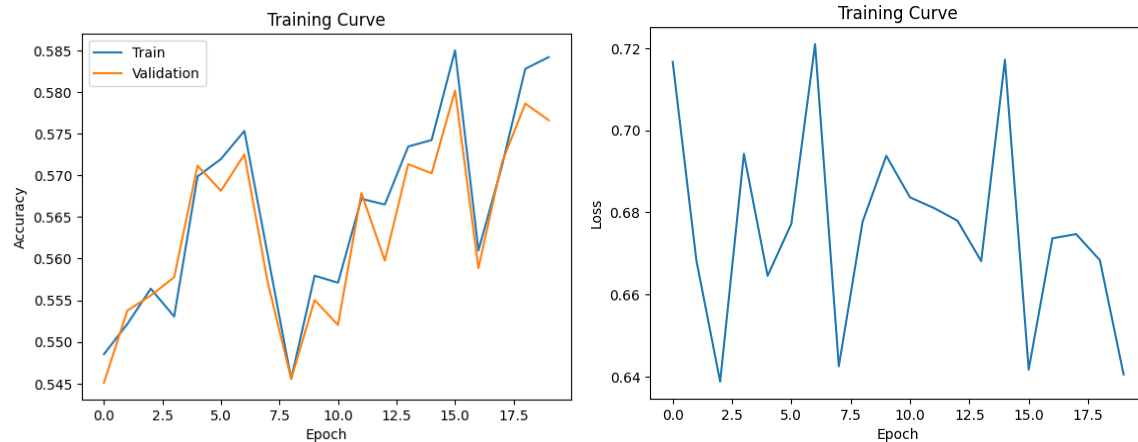
همچنین ماتریس درهم‌ریختگی به شکل زیر می‌باشد که بیشتر داده‌ها به درستی دسته‌بندی شده‌اند و مدل اشتباه خاصی نکرده است.



شکل ۳ ماتریس درهم‌ریختگی برای مدل *rnn* و *one hot*

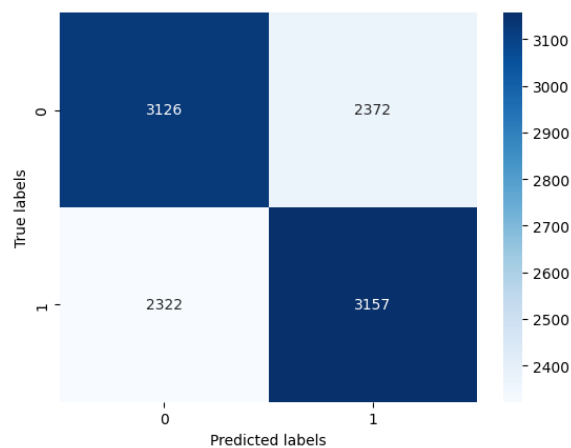
## Glove:

نتایج آموزش برای مدل RNN با این embedding به صورت زیر است.



شکل ۴ عملکرد مدل rnn برای glove

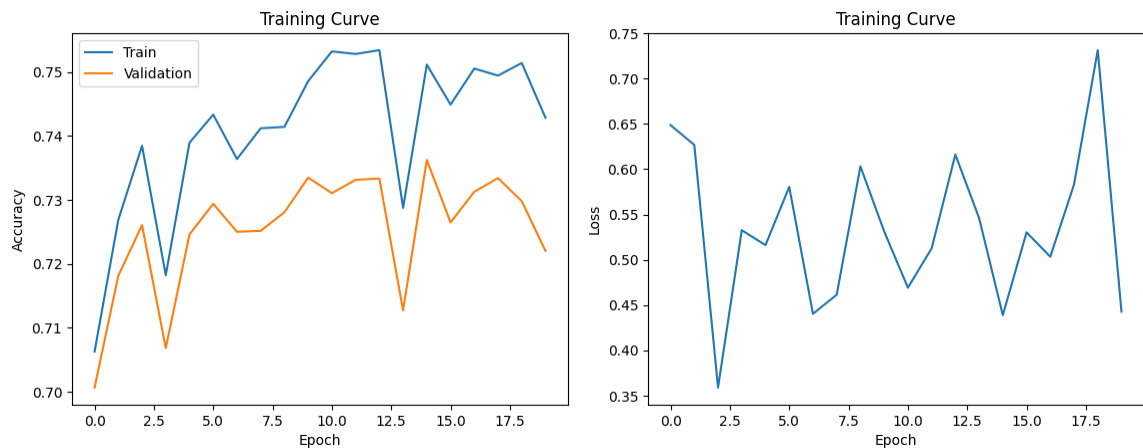
همانطور که مشخص است مدل دقت آنچنان خوبی ندارد و می‌تواند تعداد دفعات بیشتری نیز آموزش یابد. همچنین ماتریس درهم‌ریختگی این مدل به شکل زیر می‌باشد:



شکل ۵ ماتریس درهم‌ریختگی مدل RNN با glove

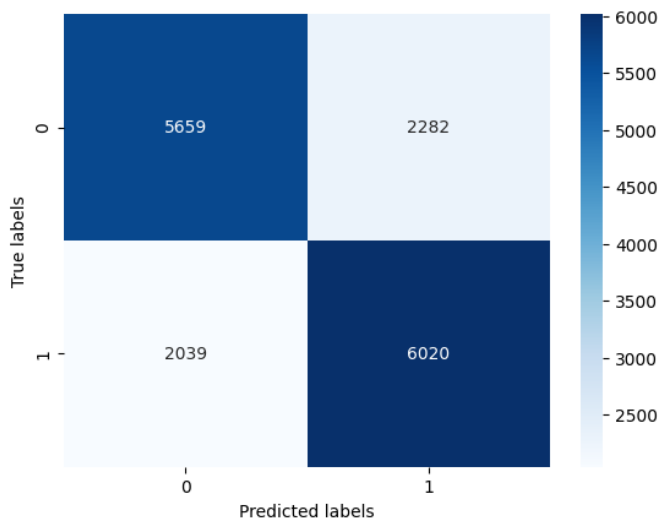
## Word2vec:

این embedding نسبت به دو embedding دیگر بهتر عمل می‌کند و نتایج آن به صورت زیر است.



شکل ۶ عملکرد مدل *rnn* برای *word2vec*

همچنین ماتریس درهم‌ریختگی به شکل زیر می‌باشد:



شکل ۷ ماتریس درهم‌ریختگی مدل *rnn* برای *word2vec*

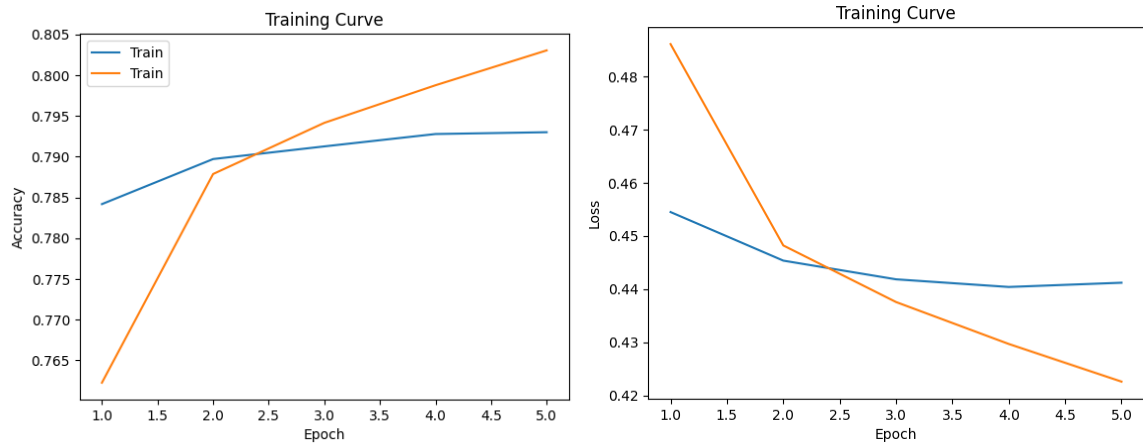
بخش ۲:

**LSTM:**

One Hot encoding:

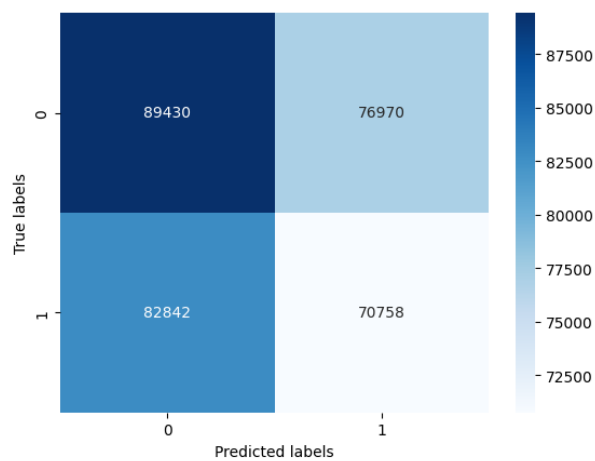
نتایج آموزش برای این مدل به شکل زیر می‌باشد:





شکل ۸ عملکرد LSTM برای one hot

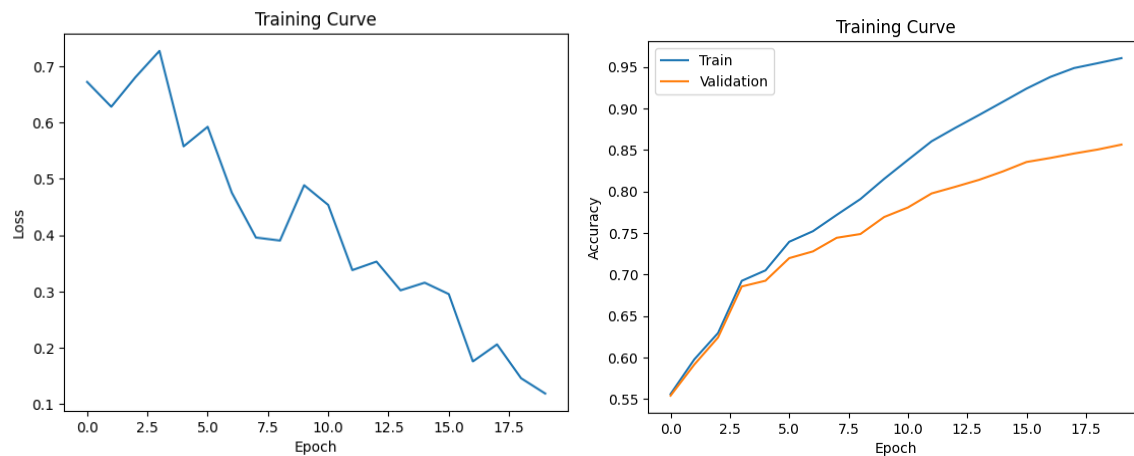
همچنین ماتریس در هم‌ریختگی برای این مدل به شکل زیر خواهد شد:



شکل ۹ ماتریس در هم‌ریختگی برای مدل LSTM و one hot

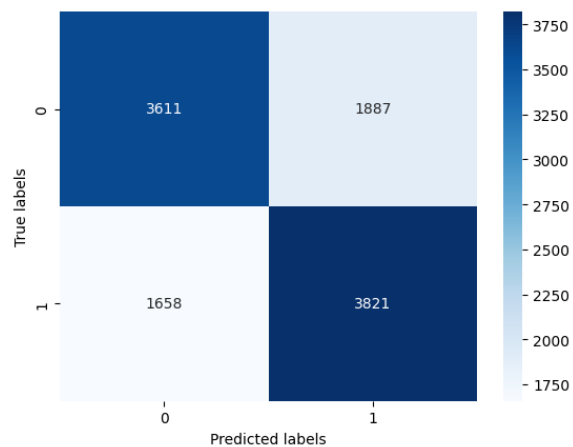
Glove:

نتایج مدل



شکل ۱۰ عملکرد *lstm* برای *glove*

همچنین ماتریس درهم‌ریختگی این مدل به شکل زیر می‌باشد:



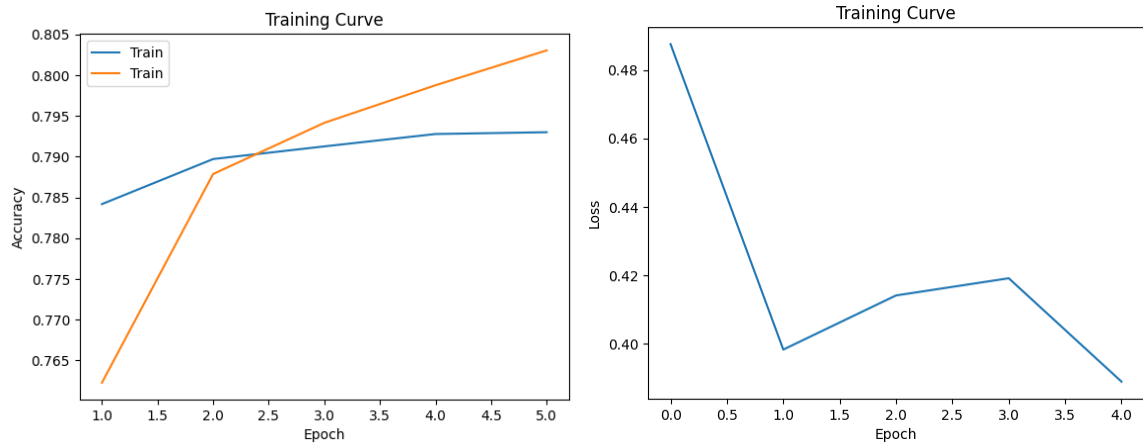
شکل ۱۱ ماتریس درهم‌ریختگی *LSMT* با *glove*

این مدل به شدت بهتر از مدل *RNN* عمل می‌کند که با توجه به پیچیده‌تر بودن آن همین نیز انتظار می‌رفت.

**GRU:**

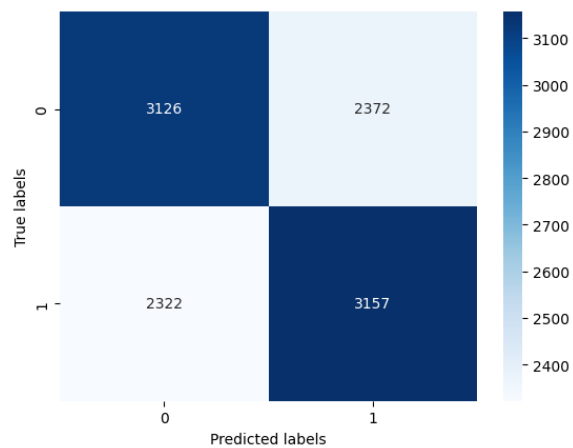
One Hot encoding:

عملکرد مدل به شکل زیر است:



شکل ۱۲ عملکرد مدل gru برای one hot

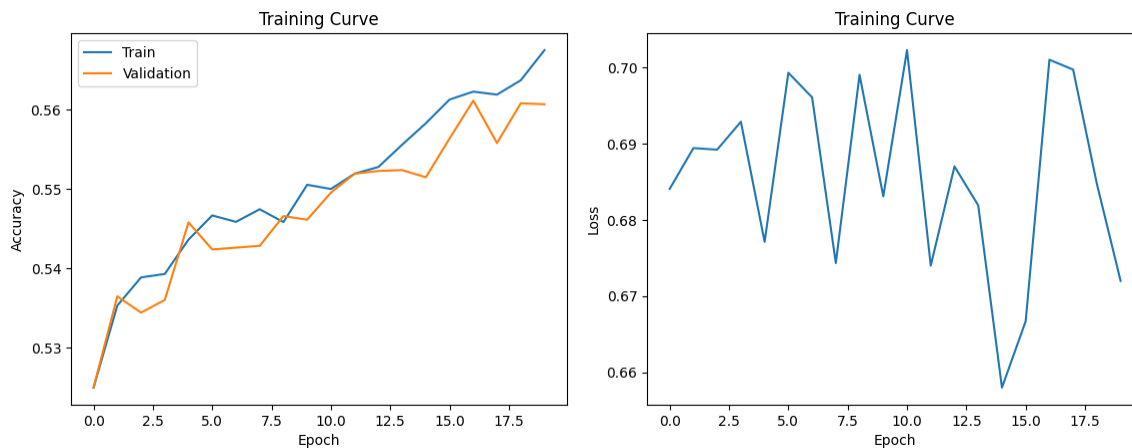
ماتریس درهم‌ریختگی به شکل زیر می‌باشد:



شکل ۱۳ ماتریس درهم‌ریختگی مدل gru و one hot

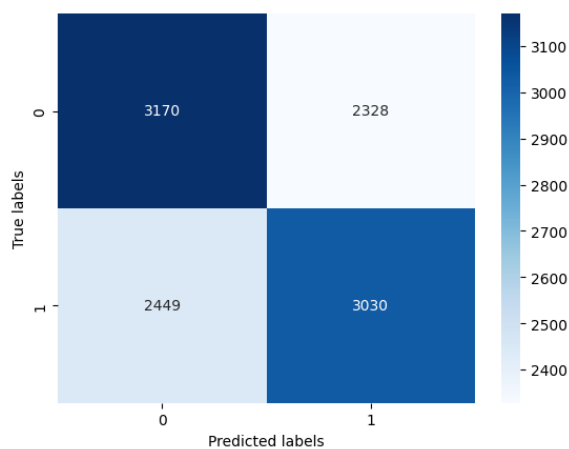
Glove:

نتیجه حاصل از این مدل به شکل زیر می‌باشد:



شکل ۴ عملکرد gru برای glove

همچنین ماتریس در هم‌ریختگی این مدل به شکل زیر می‌باشد:



شکل ۵ ماتریس در هم‌ریختگی gru با glove

## مقایسه دو مدل:

همانطور که پیش‌بینی می‌شد مدل GRU به نسبت سریع‌تر از مدل LSTM می‌باشد. اما این سریع‌تر بودن باعث کمی ضعیف‌تر عمل کردن این مدل نسبت به مدل LSTM نیز می‌شود که مورد انتظار بود. به طور کلی مدل lstm بهتر عمل کرد اما کمی کندتر نیز بود.