

صورت سوال: یک شبکه عصبی بازگشتی LSTM با استفاده از Keras طراحی نمایید و از آن برای تحلیل نظرات دیتابیس IMdB استفاده نمایید. شبکه را با استفاده از این دیتابیس آموزش داده و در نهایت برنامه باید جملات را از کاربر دریافت نموده و با امتیاز دهی در صورت کسب امتیاز کمتر مساوی ۰/۳ آن را منفی و در صورتیکه بالاتر از ۰/۳ باشد آن را مثبت قلمداد کند.

این دیتابیس در آدرس زیر موجود است:

[/http://ai.stanford.edu/~amaas/data/sentiment](http://ai.stanford.edu/~amaas/data/sentiment)

پاسخ: کد ها رو اینجا میذارم. و همچنین هر قسمت را همراه با تئوری ای که دارد کنار آن توضیح می دهم.

```
1 import numpy as np
2 import tensorflow as tf
3 import numpy as np
4 from tensorflow import keras
5 from keras.preprocessing import sequence
6 from keras import layers
7 from keras.layers import Embedding
8 import os
9 import time
10 from keras.preprocessing.sequence import pad_sequences
11 import nltk
12 from nltk.corpus import stopwords
13 stop = set (stopwords.words("english"))
14 import pandas as pd
15 import matplotlib.pyplot as plt
```

```
1 from keras.preprocessing.text import Tokenizer
2
3 def remove_stopwords (text):
4     filtered_words= [word.lower() for word in text.split() \
5                     if word.lower() not in stop]
6     return " ".join(filtered_words)
```

ابتدا کتابخانه هایی که لازم داریم را وارد کرده ایم. بیشتر کتابخانه ها معروف هستند و نیازی به توضیح ندارند اما چهار کتابخانه ی `nltk` و `pad_sequence` و `embedding` و `Tokenizer` مخصوص این سوال هستند که توضیح می دهم.

**کتابخانه nltk:** برای این استفاده کرده ام که کلمات رایج و تکراری (stopwords) را از نظرات حذف کنم. الان میتوانیم چندتا از این کلمات را ببینیم که کمکی به ما نمی کنند چون تاثیر زیادی در نظرات ندارند.

```
1 stop

{'a',
 'about',
 'above',
 'after',
 'again',
 'against',
 'ain',
 'all',
 'am',
 'an',
 'and',
```

تابعی که صفحه قبل بود این کلمات را حذف می کند. مثال:

```
1 remove_stopwords("Arash is a student.")

'arash student.'
```

**کتابخانه Tokenizer:** این کتابخانه نیاز هست چون ما نمیتونیم برداری از کلمات به شبکه بازگشتی LSTM ورودی بدیم. و باید برداری از اعداد (یا بردارها) به شبکه بدیم. این کتاب خانه میاد و به هر کلمه یک عدد تخصیص میده و هر نظر کاربر میشه یک دنباله ای از اعداد. مثال:

```
1 #EXAMPLE
2 t = Tokenizer()
3
4 test_text = ['Deep Learning',
5             'Machine Learning',
6             'Data Mining ! $ % *']
7
8 t.fit_on_texts(test_text)
9
10 sequences = t.texts_to_sequences(test_text)
11
12 print("The sequences generated from text are : ",sequences)
13 print('word_index ', t.word_index)
```

The sequences generated from text are : `[[2, 1], [3, 1], [4, 5]]`  
word\_index `{'learning': 1, 'deep': 2, 'machine': 3, 'data': 4, 'mining': 5}`

همان طور که می بینیم به هر کلمه یک integer تخصیص داده و `[2, 1]` => "Deep learning" هم چنین خوبی این تابع این هست که کاراکترهای اضافی را در نظر نمی گیرد مثلا `[4, 5]` => "Data Mining ! \$ % \*" و به خوبی این تبدیل را برای ما انجام می دهد.

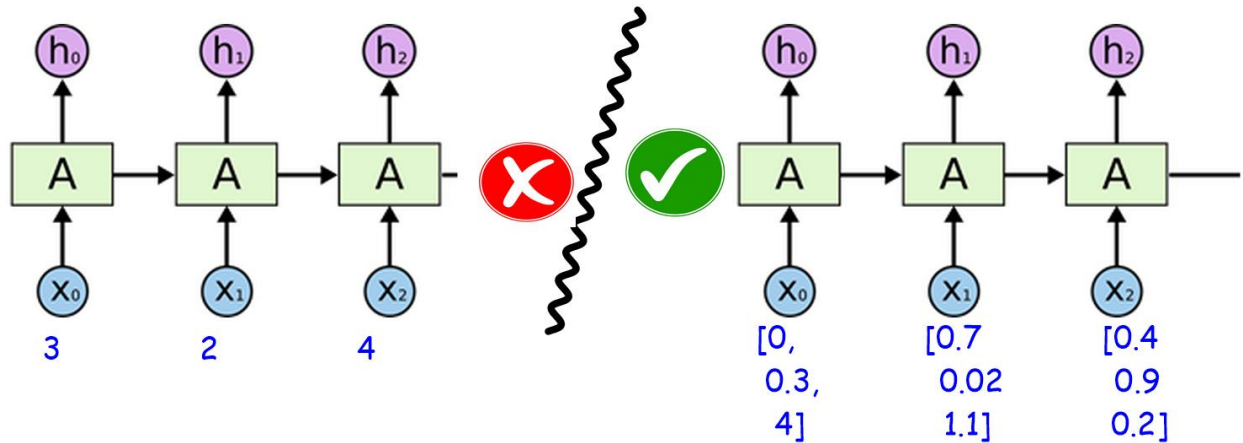
کتابخانه ی **pad\_sequence**: این کتابخانه هم به ما کمک می کند که طول ورودی ها را یکسان کنیم. چون ما باید به شبکه عصبی بازگشتی LSTM ورودی ها با طول یکسان بدیم از نظر کامپیوتری. مثال:

```
1 KNTU = [np.array([1,5,3]),np.array([2]),np.array([4,1,3,5,8,7,6])]
2 KNTU_padded=pad_sequences(KNTU, maxlen=5, \
3 padding="post", truncating="post")
```

```
1 print(KNTU[0], " "*9, KNTU_padded[0])
2 print(KNTU[1], " "*13, KNTU_padded[1])
3 print(KNTU[2], " "*1, KNTU_padded[2])
```

```
[1 5 3]          [1 5 3 0 0]
[2]              [2 0 0 0 0]
[4 1 3 5 8 7 6]  [4 1 3 5 8]
```

همان طور که می بینید وقتی pad می شوند (در اینجا با حداکثر طول=5) همه طول ها یکسان 5 می شوند آنهایی که کمتر دارند صفر می گیرند و آنهایی که بیشتر دارند از کاراکتر 5 ام به بعد بریده می شوند. حال می توانیم این ورودی را به شبکه عصبی بدهیم. البته الان یک مشکلی هست ما نمیخواهیم یک سری integer به واحد های ورودی شبکه عصبی بدهیم. می خواهیم یکسری بردار به واحد های شبکه عصبی بدهیم.



کتابخانه ی **embedding**: همانطور که در بالا گفتیم باید کلمات یا همون integer ها رو embed کنیم به یک بردار و بردار ها را به شبکه بازگشتی LSTM بدهیم که بتواند آن ها را در ماتریس های weights ضرب کند و واحد های hidden states را حساب کند. برای این تبدیل (کلمات-integer) به بردار ها از embedding استفاده می کنیم. مثال صفحه ی بعد:

```
1 embedder=tf.keras.layers.Embedding(4, 3, input_length=3)
```

چون ورودی بعد 3 داره

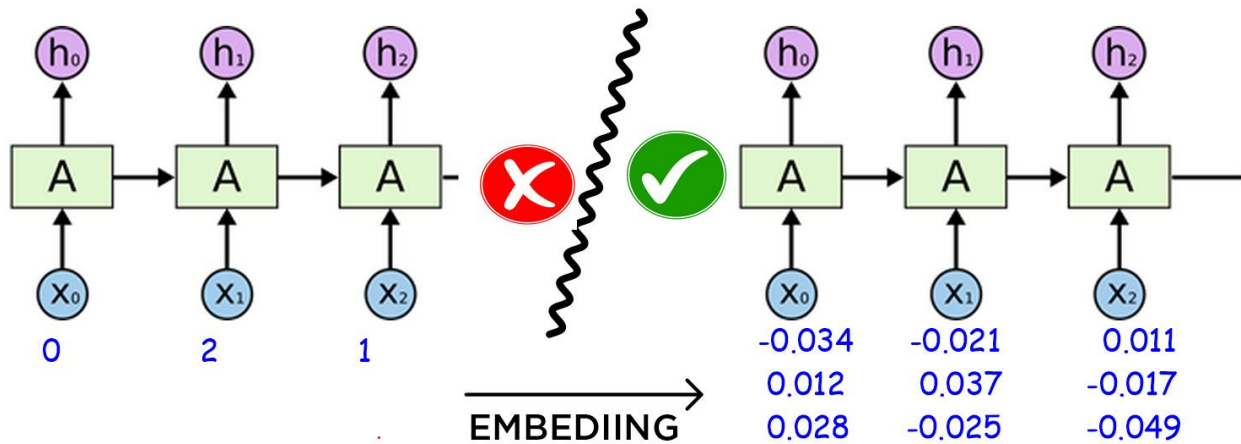
```
1 data= np.array ([[0,2,1]])
```

```
2 embedder(data)
```

چون تعداد کلمات اینجا کمتر از 4 هست کل کلمات=0 1 2 (تا 3)

چون می خواهیم هر عدد را به یک بردار 3 تایی مپ کنه \*\*\*

```
<tf.Tensor: shape=(1, 3, 3), dtype=float32, numpy=
array([[[[-0.0347145,  0.01235352,  0.02839499],
          [-0.02170154,  0.03740683, -0.02561101],
          [ 0.01157119, -0.01725961, -0.04921926]]]], dtype=float32)>
```



خوب توضیحات کتابخانه ها و روش استفاده از آن ها تمام شد. حالا وارد کد می شویم. از اینجا به بعد کار ساده هست چون فقط می خواهیم از این تکنیک ها استفاده کنیم تا داده ها آماده شوند و فیت کنیم به مدل LSTM ای که با keras به راحتی می توان طراحی کرد.

## READ TRAIN

```
1 path_train_positive = \
2 r"D:\ALL ABOUT UNIVERSITY\KNTUmaster's\Term2\DL\LMDB\aclImdb\train\pos"
3 path_train_negative = \
4 r"D:\ALL ABOUT UNIVERSITY\KNTUmaster's\Term2\DL\LMDB\aclImdb\train\neg"
5
6
7 x_train_positive=[]
8 x_train_negative=[]
9
10 start = time.time()
11
12 os.chdir(path_train_positive)
13 for file in os.listdir():
14     if file.endswith(".txt"):
15         file = open(f"{path_train_positive}\{file}", 'r', encoding='utf8')
16         x_train_positive.append(remove_stopwords(file.read()))
```

```

17
18 os.chdir(path_train_negative)
19 for file in os.listdir():
20     if file.endswith(".txt"):
21         file = open(f"{path_train_negative}\\{file}", 'r', encoding='utf8')
22         x_train_negative.append(remove_stopwords(file.read()))
23
24 end = time.time()
25 print('It took', round((end-start),2), \
26       'secs to read positive and negative train')

```

It took 135.74 secs to read positive and negative train

```

1 x_train = x_train_positive + x_train_negative
2 y_train = [1]*12500 + [0]*12500

```

در این قسمت تمام ورودی های train هم مثبت pos و هم منفی neg را خوانده ایم و همچنین ۱۲۵۰۰ تای اول را کامنت مثبت = ۱ و ۱۲۵۰۰ تای دوم را کامنت منفی = ۰ برچسب داده ایم. فقط دقت کنید که کلمات رایج (stopwords) قبل از وارد شدن دیتا به لیست حذف شده اند.

```
append(remove_stopwords(file.read()))
```

## READ TEST

```

1 path_test_positive = \
2 r"D:\ALL ABOUT UNIVERSITY\KNTUmaster's\Term2\DL\LMDB\aclImdb\test\pos"
3 path_test_negative = \
4 r"D:\ALL ABOUT UNIVERSITY\KNTUmaster's\Term2\DL\LMDB\aclImdb\test\neg"
5
6
7 x_test_positive=[]
8 x_test_negative=[]
9
10 start = time.time()
11
12 os.chdir(path_test_positive)
13 for file in os.listdir():
14     if file.endswith(".txt"):
15         file = open(f"{path_test_positive}\\{file}", 'r', encoding='utf8')
16         x_test_positive.append(remove_stopwords(file.read()))

```

```

17
18 os.chdir(path_test_negative)
19 for file in os.listdir():
20     if file.endswith(".txt"):
21         file = open(f"{path_test_negative}\\{file}", 'r', encoding='utf8')
22         x_test_negative.append(remove_stopwords(file.read()))
23
24 end = time.time()
25 print('It took', round((end-start),2), \
26       'secs to read positive and negative test')

```

It took 133.94 secs to read positive and negative test

مشابهاً برای داده ی آزمون هم این کد تکرار شده است.

## review => bag of words => id's of word

```

1 x_full = x_train + x_test
2
3 t1 = Tokenizer()
4 t1.fit_on_texts(x_full)
5 num_of_unique_words= len(t1.word_counts)
6 print("Number of unique words in All Data: ", num_of_unique_words)
7 #124,252 unique words =stopwords=> 123,277

```

Number of unique words in All Data: 123277

در این قسمت همه ی دیتا را نگاه کرده ایم و با یک tokenizer تعداد کل کلمات متمایز را حساب کرده ایم  
(۱۲۳۲۷۷).

```

1 tokenizer = Tokenizer (num_of_unique_words)
2 tokenizer.fit_on_texts(x_train)

```

```

1 x_train_sequences = tokenizer.texts_to_sequences(x_train)
2 x_test_sequences = tokenizer.texts_to_sequences(x_test)

```

سپس توسط یک tokenizer دیگر تمامی کامنت ها به اعداد تبدیل شده اند. یک مثال میبینیم:

```
1 x_train[13568]
```

'movie terrible good effects.'

```
1 np.array(x_train_sequences[13568])
```

array([ 2, 306, 6, 213])

```

1 for i in range (0,len(x_train_sequences)):    #change to numpy format
2     x_train_sequences[i] = np.array( x_train_sequences[i])
3
4 for i in range (0,len(x_test_sequences)):
5     x_test_sequences[i] = np.array( x_test_sequences[i])

```

سپس همه چیز را بفرمت numpy تبدیل می کنیم تا کار کردن شبکه عصبی با آن راحت شود.

```

1 max, sum, min, ind=0, 0, 2000, -1
2
3 for j in range (0, len(x_train)):
4
5     sum = sum + len(x_train_sequences[j])
6     if len(x_train_sequences[j])>max:
7         max=len(x_train_sequences[j])
8     if len(x_train_sequences[j])<min:
9         min=len(x_train_sequences[j])
10    ind=j
11
12 print("max length of a message: ", max)
13 print("min length of a message: ", min)
14 print("avg length of a message: ", sum/25000)

```

```

max length of a message:  1495
min length of a message:  4
avg length of a message:  132.69968

```

این قسمت من کل train رو گشتم ببینم بیشترین طول نظر و کمترین طول نظر و طول متوسط نظرات چقدر هست که یک عدد مناسب برای padding پیدا کنم.

```

1 max_length= 140
2 x_train_padded = pad_sequences(x_train_sequences, \
3                               maxlen=max_length, padding="post", truncating="post")
4 x_test_padded = pad_sequences(x_test_sequences, \
5                               maxlen=max_length, padding="post", truncating="post")
6 x_train_padded.shape, x_test_padded.shape

```

```
((25000, 140), (25000, 140))
```

من دیدم بطور متوسط طول کامنت ها 132 هست پس بیشترین طول کامنت رو 140 گذاشتم. یعنی اگر کمتر بود padding صفر اضافه کند. اگر بیشتر بود که از 140 اُمین کلمه به بعد cut کن. البته اعداد دیگری هم امتحان کردم این ظاهراً از همه بهتر بود. اعداد خیلی زیاد هم واقعاً شبکه را کند می کرد.

## LSTM MODEL

```

1 model = keras.models.Sequential()
2 model.add(layers.Embedding(num_of_unique_words, 32, \
3                             input_length=max_length))
4 model.add(layers.LSTM(80, dropout=0.15, activation="tanh"))
5 model.add(layers.Dense(1, activation="sigmoid"))
6
7 model.summary()

```

خوب یک مدل می سازیم که این کلمات را ابتدا خودش embed کنه و سپس وارد lstm شود. هر کلمه (integer) متناظر با اون کلمه) به یک بردار ۳۲ بُعدی در embedding تبدیل می شود. سپس به LSTM(80, ...) داده می شود. 80 یعنی اون خروجی hidden state یک بردار 80 تایی خواهد بود. می توان اعداد دیگر هم امتحان کرد. Dropout=0.15 یعنی با احتمال ۱۵ درصد یک نرون ممکن است حذف شود. (جلوگیری از overfitting)

Model: "sequential\_4"

Layer (type)	Output Shape	Param #
embedding_35 (Embedding)	(None, 140, 32)	3944864
lstm_4 (LSTM)	(None, 80)	36160
dense_4 (Dense)	(None, 1)	81
Total params: 3,981,105		
Trainable params: 3,981,105		
Non-trainable params: 0		

تعداد پارامترهای لایه ی embedder رو نمی تونیم بشماریم. Embedder خودش یک مقاله هست با نام [Efficient Estimation of Word Representations in Vector Space](#) که خودش یک شبکه عصبی جدا هست. ولی تعداد پارامترهای LSTM را با درسی که داشتیم میتونیم راحت بشماریم.



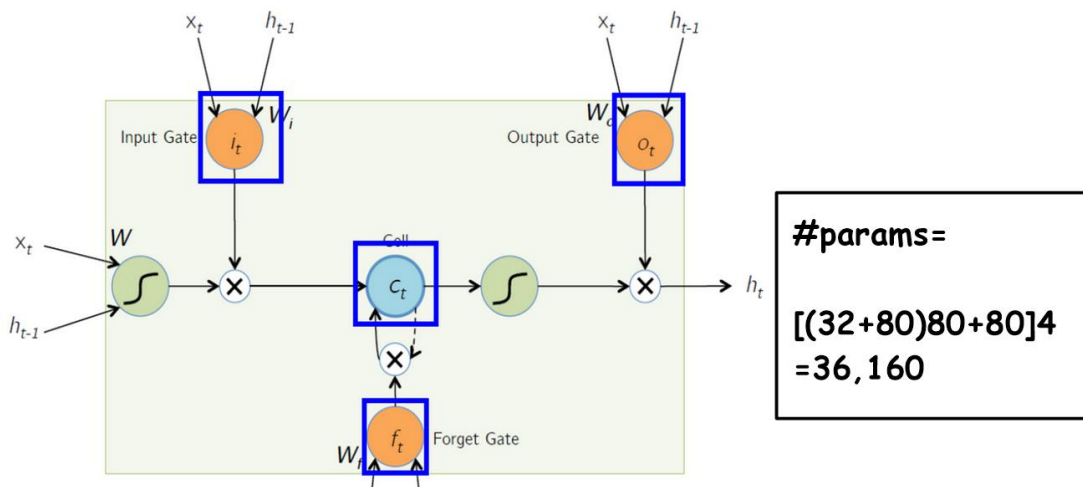
الان ورودی ها بعد از embed شدن 32 تایی هستند. و چون LSTM(80,...) هست یعنی این بردار ۳۲ تایی در کنار یک بردار hidden state قبلی که ۸۰ تایی هست قرار می گیرد و سپس باید ضرب در یک ماتریسی بشه که این بردار طولانی را دوباره به ۸۰ تایی تبدیل کنه (hidden state جدید). پس این ماتریس ابعاد

$$(80+32)*80$$

دارد همچنین 80 تا پارامتر هم برای bias ها هست. پس مجموعاً

$$(80+32)*80+80=9040$$

پارامتر داریم. اما پارامتر ها گفته 36,160 این به این دلیل است که این محاسبات برای یک RNN معمولی simpleRNN بود اینجا LSTM هست. هر LSTM چهار برابر یک simpleRNN پارامتر داره بخاطر cell, forget, input, output هر کدام به همین تعداد پارامتر دارند.



```
1 loss = keras.losses.BinaryCrossentropy(from_logits=False)
2 #optim = keras.optimizers.Adam(learning_rate=0.001)
3 optim = keras.optimizers.Adamax(learning_rate=0.001)
4 metrics = ["accuracy"]
5
6 model.compile(loss=loss, optimizer=optim, metrics=metrics)
```

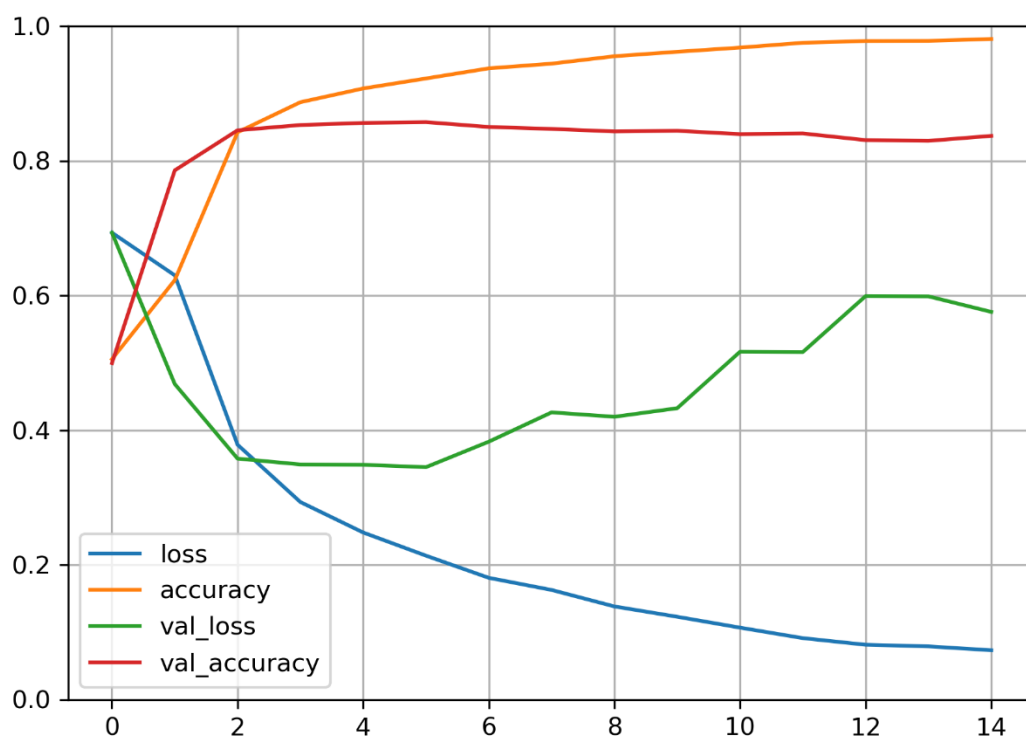
```
1 kee = model.fit(x_train_padded, np.array(y_train), \
2                 epochs=15, batch_size=150, \
3                 validation_data=\
4                 (x_test_padded, np.array(y_test)), verbose=2 )
```

خوب تابع loss و optimizer تعریف می شوند. ابتدا Adam استفاده کردم خیلی دیر همگرا می شود (بعد از ۱۰ مرحله تقریباً) سپس تو سایت خوندم نوشته بود Adamax برای مدل هایی که embedder دارند سریع تر و بهتر کار

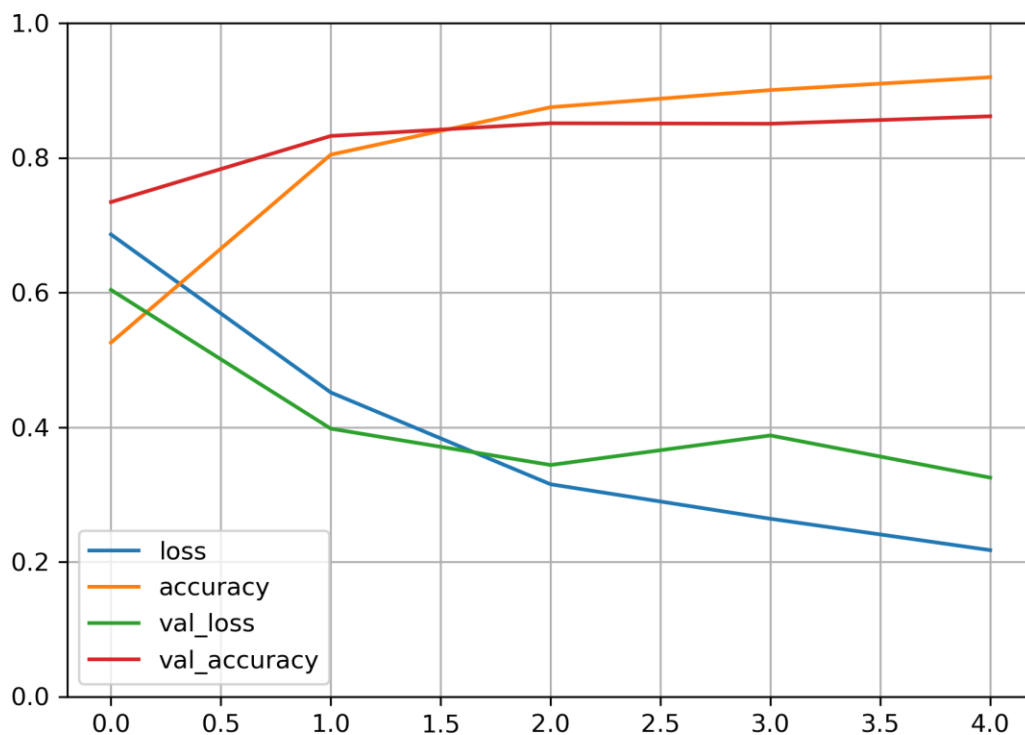
می‌کنه. پس از optimizer Adamax استفاده کردم. سپس دیتا را فیت می‌کنیم. و بعد از هر ۱۵۰ تا نمونه وزن‌ها را آپدیت می‌کنیم. نتیجه:

```
Epoch 1/15
loss: 0.6931 - accuracy: 0.5054 - val_loss: 0.6931 - val_accuracy: 0.4995
Epoch 2/15
loss: 0.6299 - accuracy: 0.6227 - val_loss: 0.4684 - val_accuracy: 0.7857
Epoch 3/15
loss: 0.3783 - accuracy: 0.8423 - val_loss: 0.3577 - val_accuracy: 0.8452
Epoch 4/15
loss: 0.2932 - accuracy: 0.8869 - val_loss: 0.3490 - val_accuracy: 0.8530
Epoch 5/15
loss: 0.2480 - accuracy: 0.9073 - val_loss: 0.3485 - val_accuracy: 0.8561
Epoch 6/15
loss: 0.2136 - accuracy: 0.9222 - val_loss: 0.3450 - val_accuracy: 0.8573
Epoch 7/15
loss: 0.1807 - accuracy: 0.9372 - val_loss: 0.3828 - val_accuracy: 0.8502
Epoch 8/15
loss: 0.1627 - accuracy: 0.9442 - val_loss: 0.4263 - val_accuracy: 0.8472
```

همان‌طور که می‌بینیم بعد از ۴ یا پنج مرحله دیگه روی test بهبودی نداریم.



پس من یکبار دیگه فقط تا epoch ۵ آموزش میدم که overfit نشده باشه. نتیجه میشه:



و نتیجه ی آخرین epoch:

Epoch 5/5

loss: 0.2173 - accuracy: 0.9195 - val\_loss: 0.3250 - val\_accuracy: 0.8615

پس به دقت ۹۱ درصد روی داده ی آموزش و ۸۶ درصد روی داده ی آزمون رسیده ایم.

حالا چند مورد هم خودم میخوام بنویسم ببینم این شبکه چقدر میتونه خوب عمل کنه و کجاها دیگه نمیتونه تحلیل کنه نظر من رو. من یک نظر می نویسم و کارهای tokenizing و padding و remove\_stopwords و embedding همه رو انجام میدم و به شبکه عصبی می دهم. انتظار دارم اگر کامنتم بد باشه کمتر از 0.3 بده اگر خوب باشه بیشتر از 0.3 خروجی دهد. تابع فعالساز تک واحد خروجی sigmoid هست پس همیشه بین 0 و 1 می دهد. این کار را با کد زیر انجام می دهم.

## TEST ON MY COMMENTS

```
1 tt = [" It was the best movie I have seen in my life"]
2 tt = [remove_stopwords(tt[0])]
3 tkn=tokenizer.texts_to_sequences(tt)
4 tkn_np = np.array (tkn)
5 arash_padded = pad_sequences(tkn_np, maxlen=max_length, \
6                             padding="post", truncating="post")
```

```
1 model.predict(arash_padded)
```

array([[0.8846246]], dtype=float32)

## TEST ON MY COMMENTS

```
1 tt = [" It was the worst movie I have seen in my life"]
2 tt = [remove_stopwords(tt[0])]
3 tkn=tokenizer.texts_to_sequences(tt)
4 tkn_np = np.array (tkn)
5 arash_padded = pad_sequences(tkn_np, maxlen=max_length, \
6                             padding="post", truncating="post")
```

```
1 model.predict(arash_padded)
```

array([[0.05919594]], dtype=float32)

خوب بنظر میرسه که داره خوب کار می کنه. یه چندتا کامنت دیگه هم بدم:

Comment	RNN OUTPUT	PERFORMANCE
"Although the movie seemed to be bad at first I really loved it by the end of it. The ending was nice."	0.592	ACCEPTABLE
"My friends said it was not good at all, but what I thought was quite different I think it was good but it could've been even better"	0.827	ACCEPTABLE
"Everyone should see this movie"	0.709	ACCEPTABLE
"Everyone should see this movie to see how bad it is. Then they can appreciate other good movies we have."	0.749*	NOT ACCEPTABLE
"To be honest, I didn't like this movie. The director didn't do well neither did the actors. It was boring and long. They could've conveyed the message in much less time."	0.067	ACCEPTABLE

نسبتاً خوب کار می کند. فقط یک نظر که با ستاره \* مشخص کردم. منفی هست ولی مثبت تشخیص داده.

**"Everyone should see this movie to see how bad it is. Then they can appreciate other good movies we have."**

البته تشخیص همچنین نظری منطقی هست که برای ماشین باید سخت باشد.

(کد ها را جدا ارسال می کنم.)

پایان