

```
In [1]: 1  ## Importing libraries
2  %matplotlib inline
3  import warnings
4  warnings.filterwarnings("ignore")
5
6  import pandas as pd
7  import numpy as np
8  import nltk
9  import string
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 import tensorflow as tf
13 from sklearn.feature_extraction.text import TfidfTransformer
14 from sklearn.feature_extraction.text import TfidfVectorizer
15
16 from sklearn.feature_extraction.text import CountVectorizer
17 from sklearn.metrics import confusion_matrix
18 from sklearn import metrics
19 from sklearn.metrics import roc_curve, auc
20 from nltk.stem.porter import PorterStemmer
21
22 import re
23 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
24 import string
25 from nltk.corpus import stopwords
26 from nltk.stem import PorterStemmer
27 from nltk.stem.wordnet import WordNetLemmatizer
28
29 from gensim.models import Word2Vec
30 from gensim.models import KeyedVectors
31 import pickle
32
33 from tqdm import tqdm
34 import os
35
36 # from chart_studio.plotly import plotly
37 # import plotly.offline as offline
38 # import plotly.graph_objs as go
39 #offline.init_notebook_mode()
40 from collections import Counter
```

```
In [3]: 1  from google.colab import drive
2  drive.mount('/content/drive')
```

Mounted at /content/drive

```
In [2]: 1  # !pip install tensorflow==1.15.0
2  # !pip install keras==2.3.1
3  # # import tensorflow
4  # tensorflow.__version__
```

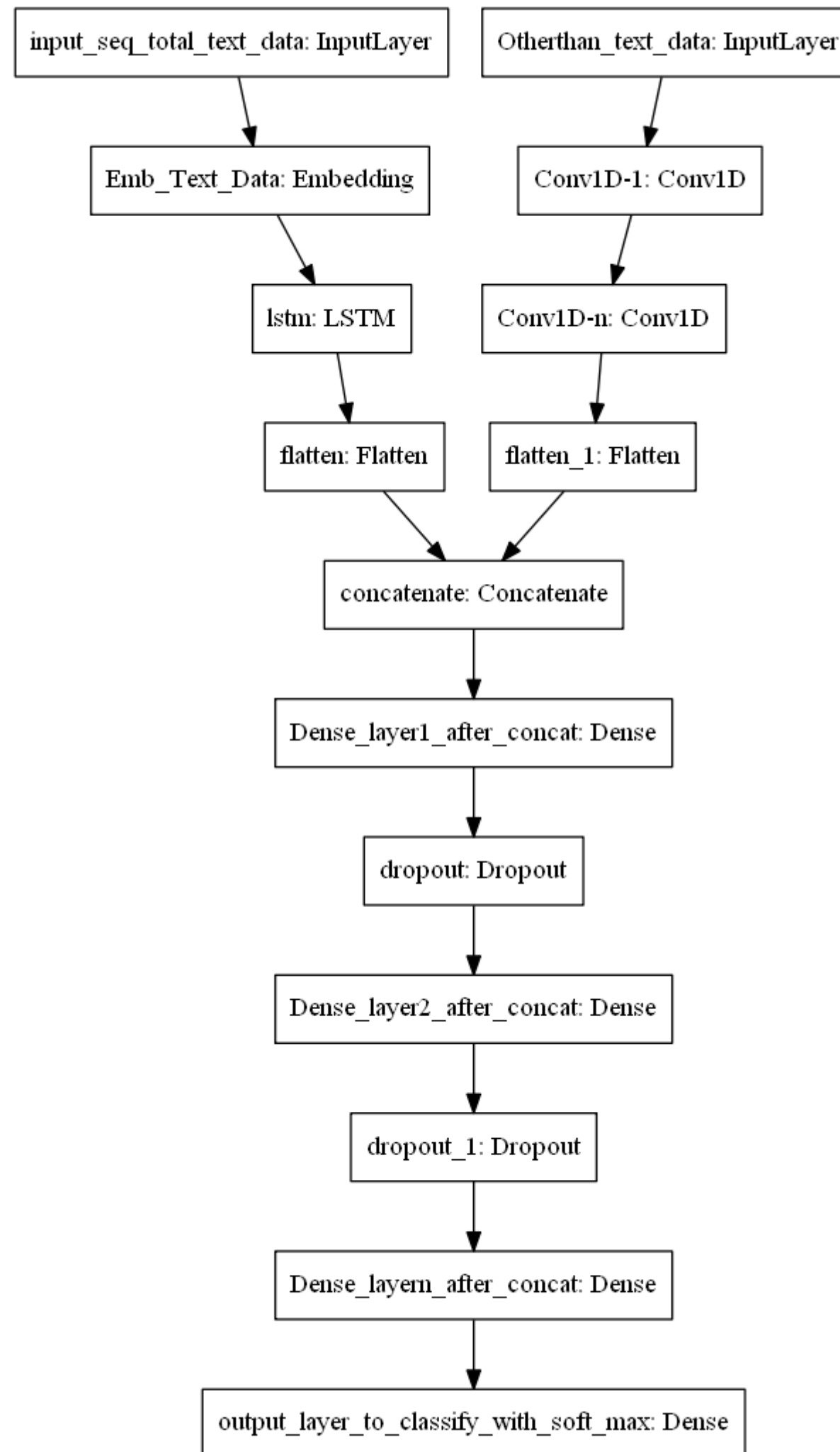
```
In [5]: 1  # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model1/processed_data_split.h2' '/content/'
2  # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model1/model_inputs.pkl' '/content/'
3  # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model1/model_input_cat_labels.pkl' '/content/'
4  # !cp -r '/content/drive/My Drive/_datasets/glove_vectors' '/content/'
```

```
In [6]: ► 1  ## Load data after split
2  x_train = pd.read_hdf('processed_data_split.h2', 'x_train',mode='r')
3  x_test = pd.read_hdf('processed_data_split.h2', 'x_test',mode='r')
4  x_cv = pd.read_hdf('processed_data_split.h2', 'x_cv',mode='r')
5  y_train =pd.read_hdf('processed_data_split.h2', 'y_train',mode='r')
6  y_test =pd.read_hdf('processed_data_split.h2', 'y_test',mode='r')
7  y_cv =pd.read_hdf('processed_data_split.h2', 'y_cv',mode='r')
8  print('*'*50)
9  print(' Successfully loaded processed split data')
10 emd_i,embedding_matrix,seq_x_train,seq_x_test,seq_x_cv,padseq_x_train,sklstate_train,proj_grade_train,train_categories,train_subcategories,teacher_prefix_train,numerical_train
11 print('*'*50)
12 print(' Successfully loaded model input variables')
13 y_train_cat,y_test_cat,y_cv_cat = pickle.load(open('model_input_cat_lables.pkl', 'rb'))
14 print('*'*50)
15 print('Successfully loaded split y labels')
```

```
*****
Successfully loaded processed split data
*****
Successfully loaded model input variables
*****
Successfully loaded split y labels
```

```
In [12]: ► 1 vocab_size = embedding_matrix.shape[0]
```

Model-3



ref: <https://i.imgur.com/fkQ8nGo.png>

- **input_seq_total_text_data:**

- . Use text column('essay'), and use the Embedding layer to get word vectors.
- . Use given predefined glove word vectors, don't train any word vectors.
- . Use LSTM that is given above, get the LSTM output and Flatten that output.
- . You are free to preprocess the input text as you needed.

- **Other_than_text_data:**

- . Convert all your Categorical values to onehot coded and then concatenate all these onehot vectors
- . Neumerical values and use [CNN1D \(https://keras.io/getting-started/sequential-model-guide/#sequence-classification-with-1d-convolutions\)](https://keras.io/getting-started/sequential-model-guide/#sequence-classification-with-1d-convolutions) as shown in above figure.
- . You are free to choose all CNN parameters like kernel sizes, stride.

3.1 Hyperparameter Tuning

```
In [8]: ▶ 1 from scipy.sparse import hstack
2 other_than_text_train = hstack((sklstate_train,proj_grade_train,train_categories,train_subcategories,teacher_prefix_train,numerical_train)).todense()
3 other_than_text_test = hstack((sklstate_test,proj_grade_test,test_categories,test_subcategories,teacher_prefix_test,numerical_test)).todense()
4 other_than_text_cv = hstack((sklstate_cv,proj_grade_cv,cv_categories,cv_subcategories,teacher_prefix_cv,numerical_cv)).todense()
```

```
In [9]: ▶ 1 other_than_text_train = np.expand_dims(other_than_text_train,2)
2 other_than_text_test = np.expand_dims(other_than_text_test,2)
3 other_than_text_cv = np.expand_dims(other_than_text_cv,2)
```

```
In [10]: ▶ 1 print(other_than_text_train.shape)
2 print(other_than_text_test.shape)
3 print(other_than_text_cv.shape)
```

```
(69918, 102, 1)
(21850, 102, 1)
(17480, 102, 1)
```

```
In [11]: ▶ 1 import tensorflow as tf
2 from keras.callbacks import TensorBoard, ModelCheckpoint
3 import keras
4 import datetime, os
5 import keras.backend as K
6 from keras.regularizers import l2
7 from tensorflow import set_random_seed
8 from sklearn.metrics import roc_auc_score
9 from keras.layers import Dropout, Input, Activation, Dense, Embedding, concatenate, LSTM, Flatten, BatchNormalization
10 from keras.layers import Conv1D, MaxPool1D
11 from keras.models import Model
12 # Load the TensorBoard notebook extension
13 %load_ext tensorboard
14
15 def aucroc(y_true, y_pred):
16     try:
17         return tf.py_func(roc_auc_score, (y_true, y_pred), tf.double)
18     except ValueError:
19         pass
```

Using TensorFlow backend.

```

In [9]: 1  ## clear the graph of the tensorflow
2  K.clear_session()
3  ### defining all the Input Layer
4  input_seq_total_text_data = Input(shape=padseq_x_train[0].shape,name='text_Input')
5  input_other_than_text = Input(shape=(102,1),name='other_than_text')
6
7  set_random_seed(5)
8  auc_scores_model3=[]
9  if not os.path.isfile('tuning_output_model3.pkl'):
10     for k_s in [3,5,10]:
11         for stride in [1,2,3]:
12             for mx_p in [3,5]:
13                 #definig embedding layer
14                 embed_text_data = Embedding(vocab_size,300,weights = [embedding_matrix],trainable=False)(input_seq_total_text_data)
15                 #embed_other_than_text = Embedding(other_than_text_train.shape[1],other_than_text_train.shape[1]//e)(input_other_than_text)
16                 ##defining LSTM layer
17                 lstm_layer = LSTM(128,return_sequences=True)(embed_text_data)
18                 ##convolutional layer
19                 for c_i in [64,128,256]:
20                     conv_1 = Conv1D(c_i,kernel_size=k_s,strides=stride,padding='same',activation='relu',kernel_initializer='he_normal')(input_other_than_text)
21                     conv_2 = Conv1D(c_i,kernel_size=k_s,strides=stride,padding='same',activation='relu',kernel_initializer='he_normal')(conv_1)
22                     max_pool = MaxPool1D(pool_size=3,padding='same')(conv_2)
23                     conv_3 = Conv1D(c_i,kernel_size=k_s,strides=stride,padding='same',activation='relu',kernel_initializer='he_normal')(max_pool)
24                     conv_4 = Conv1D(c_i,kernel_size=k_s,strides=stride,padding='same',activation='relu',kernel_initializer='he_normal')(conv_3)
25                     max_pool = MaxPool1D(pool_size=3,padding='same')(conv_4)
26                 #max_pool = MaxPool1D(pool_size=(mx_p),padding='same')(conv_2)
27                 #flatten layer
28                 flatten_1 = Flatten()(lstm_layer)
29                 flatten_2 = Flatten()(max_pool)
30                 # concat layer
31                 concat_layer = concatenate([flatten_1,flatten_2])
32                 #dense layers
33                 dense_layer_1 = Dense(512,activation='relu',kernel_initializer='he_normal',kernel_regularizer=l2(0.0001))(concat_layer)
34                 regularization = BatchNormalization()(dense_layer_1)
35                 regularization = Dropout(0.25)(regularization)
36                 dense_layer_2 = Dense(256,activation='relu',kernel_initializer='he_normal',kernel_regularizer=l2(0.0001))(regularization)
37                 regularization = BatchNormalization()(dense_layer_2)
38                 dense_layer_3 = Dense(128,activation='relu',kernel_initializer='he_normal',kernel_regularizer=l2(0.0001))(regularization)
39                 regularization = BatchNormalization()(dense_layer_3)
40                 regularization = Dropout(0.25)(dense_layer_3)
41                 out_layer = Dense(2,activation='softmax',kernel_initializer='glorot_normal',kernel_regularizer=l2(0.0001))(dense_layer_3)
42
43                 model3 = Model(inputs=[input_seq_total_text_data,input_other_than_text],outputs=[out_layer])
44
45                 # https://stackoverflow.com/questions/48285129/saving-best-model-in-keras
46                 # Whether to restore model weights from the epoch with the best value of the monitored quantity. If False, the model weights obtained at the last step of training
47                 model3.compile(optimizer=keras.optimizers.Adam(),loss='categorical_crossentropy',metrics=['accuracy',aucroc])
48
49                 callback = tf.keras.callbacks.EarlyStopping(monitor='val_aucroc',verbose=1, patience=2,restore_best_weights=True,mode='max')
50
51                 history = model3.fit([padseq_x_train,other_than_text_train],y_train_cat,epochs=10,batch_size=1000,verbose=1,
52                                     validation_data=[padseq_x_cv,other_than_text_cv],y_cv_cat,
53                                     callbacks=[callback])
54                 max_ = np.argmax(history.history['val_aucroc'])
55                 print('Validation loss for LSTM_untis={0}, kernel_size={1}, stride_size={2}, max_pool={3}, Conv1_units={4}'
56                       ''.format(128,k_s,stride,mx_p,c_i), ' is : ',history.history['val_loss'][max_])
57
58                 auc_scores_model3.append((k_s,stride,mx_p,c_i,history.history['accuracy'][max_],
59                                           ,history.history['loss'][max_],history.history['aucroc'][max_],
60                                           history.history['val_accuracy'][max_],history.history['val_loss'][max_],history.history['val_aucroc'][max_]))

```

```

61 df = pd.DataFrame(data=auc_scores_model3,columns=['Kernel_size','stride_size','pool_size','CNN_units',
62                                                    'Train Accuracy','Train Loss','Train auc','Test Accuracy','Test Loss','Test auc'])
63 best_param = df[df['Test auc'] == df['Test auc'].max()]
64 with open('tuning_output_model3.pkl', 'wb') as f:
65     pickle.dump([df, auc_scores_model3, best_param], f)
66 else:
67     df, auc_scores_model3, best_param = pickle.load(open('tuning_output_model3.pkl', 'rb'))
68     print('----Tuning output loaded -----')
69
70

```

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/tensorflow_core/python/ops/resource_variable_ops.py:1630: calling BaseResourceVariable.__init__ (from tensorflow.python.ops.resource_variable_ops) with constraint is deprecated and will be removed in a future version.

Instructions for updating:

If using Keras pass *_constraint arguments to layers.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:4070: The name tf.nn.max_pool is deprecated. Please use tf.nn.max_pool2d instead.

WARNING:tensorflow:From <ipython-input-8-ed55f280d3db>:11: py_func (from tensorflow.python.ops.script_ops) is deprecated and will be removed in a future version.

Instructions for updating:

tf.py_func is deprecated in TF V2. Instead, there are two

options available in V2.

- tf.py_function takes a python function which manipulates tf eager tensors instead of numpy arrays. It's easy to convert a tf eager tensor to an ndarray (just call tensor.numpy()) but having access to eager tensors means `tf.py_function`s can use accelerators such as GPUs as well as being differentiable using a gradient tape.

- tf.numpy_function maintains the semantics of the deprecated tf.py_func (it is not differentiable, and manipulates numpy arrays). It drops the stateful argument making all functions stateful.

```

In [14]: 1 # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model3/tuning_output_model3.pkl' '/content/'
        2 # df, auc_scores_model3, best_param = pickle.load(open('tuning_output_model3.pkl', 'rb'))

```

```

In [14]: 1 best_param

```

```

Out[14]:
   Kernel_size  stride_size  pool_size  CNN_units  Train Accuracy  Train Loss  Train auc  Test Accuracy  Test Loss  Test auc
13           10           1           5         256         0.854358     0.516614     0.756903         0.837185     0.597083     0.752001

```

3.1 Model Trials after hyperparameter tuning

Trail :

- * Our primary task is to improve Test AUC.
- * After hyperparameter tuning we found best kernel size ,stride size ,poolsize etc. Now lets tune our dense layers without using for loops.
- * lets increase layer1 units to 854 from 512

```

In [75]: 1 #!cp -r '/content/model3_weights.best.hdf5' '/content/drive/My Drive/LSTM_preprocessed/model3/'

```

In [21]:

```

1  # Clear any logs from previous runs
2  !rm -rf ./logs/
3
4  ## clear the graph of the tensorflow
5  K.clear_session()
6  ### defining all the Input Layer
7  input_seq_total_text_data = Input(shape=padseq_x_train[0].shape,name='text_Input')
8  input_other_than_text = Input(shape=(102,1),name='other_than_text')
9
10
11 set_random_seed(5)
12 #definig embedding Layer
13 embed_text_data = Embedding(vocab_size,300,weights = [embedding_matrix],trainable=False)(input_seq_total_text_data)
14 #embed_other_than_text = Embedding(other_than_text_train.shape[1],other_than_text_train.shape[1]//e)(input_other_than_text)
15 ##defining LSTM Layer
16 lstm_layer = LSTM(128,return_sequences=True)(embed_text_data)
17 ##convolutional Layer
18
19 conv_1 = Conv1D(int(best_param['CNN_units']),
20                 kernel_size=int(best_param['Kernel_size']),
21                 strides=int(best_param['stride_size']),
22                 padding='same',activation='relu',kernel_initializer='he_normal')(input_other_than_text)
23 conv_2 = Conv1D(int(best_param['CNN_units']),
24                 kernel_size=int(best_param['Kernel_size']),
25                 strides=int(best_param['stride_size']),
26                 padding='same',activation='relu',kernel_initializer='he_normal')(conv_1)
27 max_pool = MaxPool1D(pool_size=(int(best_param['pool_size'])),padding='same')(conv_2)
28 conv_3 = Conv1D(int(best_param['CNN_units']),
29                 kernel_size=int(best_param['Kernel_size']),
30                 strides=int(best_param['stride_size']),
31                 padding='same',activation='relu',kernel_initializer='he_normal')(max_pool)
32 conv_4 = Conv1D(int(best_param['CNN_units']),
33                 kernel_size=int(best_param['Kernel_size']),
34                 strides=int(best_param['stride_size']),
35                 padding='same',activation='relu',kernel_initializer='he_normal')(conv_3)
36 max_pool = MaxPool1D(pool_size=((int(best_param['pool_size']))),padding='same')(conv_4)
37 #max_pool = MaxPool1D(pool_size=(mx_p),padding='same')(conv_2)
38 #flatten Layer
39 flatten_1 = Flatten()(lstm_layer)
40 flatten_2 = Flatten()(max_pool)
41 # concat Layer
42 concat_layer = concatenate([flatten_1,flatten_2])
43 #dense Layers
44 dense_layer_1 = Dense(854,activation='relu',kernel_initializer='he_normal',kernel_regularizer=l2(0.0001))(concat_layer)
45 regularization = BatchNormalization()(dense_layer_1)
46 regularization = Dropout(0.45)(regularization)
47 dense_layer_2 = Dense(512,activation='relu',kernel_initializer='he_normal',kernel_regularizer=l2(0.0001))(regularization)
48 regularization = BatchNormalization()(dense_layer_2)
49 dense_layer_3 = Dense(256,activation='relu',kernel_initializer='he_normal',kernel_regularizer=l2(0.0001))(regularization)
50 regularization = BatchNormalization()(dense_layer_3)
51 regularization = Dropout(0.25)(regularization)
52 out_layer = Dense(2,activation='softmax',kernel_initializer='glorot_normal',kernel_regularizer=l2(0.0001))(regularization)
53
54 model3 = Model(inputs=[input_seq_total_text_data,input_other_than_text],outputs=[out_layer])
55
56 # https://stackoverflow.com/questions/48285129/saving-best-model-in-keras
57 # Whether to restore model weights from the epoch with the best value of the monitored quantity. If False, the model weights obtained at the last step of training are used.
58 model3.compile(optimizer=keras.optimizers.Adam(),loss='categorical_crossentropy',metrics=['accuracy',aucroc])
59
60 callback = tf.keras.callbacks.EarlyStopping(monitor='val_aucroc',verbose=1, patience=2,restore_best_weights=True,mode='max')

```



```

61 log_dir="logs/" + datetime.datetime.now().strftime("%Y%m%d-%H%M%S")
62 tensorboard_callback = tf.keras.callbacks.TensorBoard(log_dir=log_dir,histogram_freq=0, write_graph=True,write_grads=True)
63
64 #https://machinelearningmastery.com/check-point-deep-learning-models-keras/
65 filepath="model3_weights.best.hdf5"
66 checkpoint = ModelCheckpoint(filepath, monitor='val_aucroc', verbose=1, save_best_only=True, mode='max')
67
68 callback_list = [callback,tensorboard_callback,checkpoint]
69 history = model3.fit([padseq_x_train,other_than_text_train],y_train_cat,epochs=10,batch_size=1000,verbose=1,
70     validation_data=([padseq_x_cv,other_than_text_cv],y_cv_cat),callbacks=callback_list)

```

Train on 69918 samples, validate on 17480 samples

Epoch 1/10

69918/69918 [=====] - 38s 550us/step - loss: 1.0177 - accuracy: 0.7194 - aucroc: 0.5127 - val_loss: 0.9752 - val_accuracy: 0.8486 - val_aucroc: 0.5377

Epoch 00001: val_aucroc improved from -inf to 0.53766, saving model to model3_weights.best.hdf5

Epoch 2/10

69918/69918 [=====] - 37s 534us/step - loss: 0.7582 - accuracy: 0.8357 - aucroc: 0.6052 - val_loss: 0.7427 - val_accuracy: 0.8371 - val_aucroc: 0.6949

Epoch 00002: val_aucroc improved from 0.53766 to 0.69492, saving model to model3_weights.best.hdf5

Epoch 3/10

69918/69918 [=====] - 37s 523us/step - loss: 0.6614 - accuracy: 0.8470 - aucroc: 0.7099 - val_loss: 0.7183 - val_accuracy: 0.8486 - val_aucroc: 0.7344

Epoch 00003: val_aucroc improved from 0.69492 to 0.73442, saving model to model3_weights.best.hdf5

Epoch 4/10

69918/69918 [=====] - 37s 535us/step - loss: 0.6001 - accuracy: 0.8513 - aucroc: 0.7438 - val_loss: 0.5829 - val_accuracy: 0.8490 - val_aucroc: 0.7577

Epoch 00004: val_aucroc improved from 0.73442 to 0.75773, saving model to model3_weights.best.hdf5

Epoch 5/10

69918/69918 [=====] - 37s 530us/step - loss: 0.5574 - accuracy: 0.8559 - aucroc: 0.7629 - val_loss: 0.6824 - val_accuracy: 0.8486 - val_aucroc: 0.7412

Epoch 00005: val_aucroc did not improve from 0.75773

Epoch 6/10

69918/69918 [=====] - 37s 527us/step - loss: 0.5247 - accuracy: 0.8584 - aucroc: 0.7779 - val_loss: 0.8514 - val_accuracy: 0.8486 - val_aucroc: 0.7352

Restoring model weights from the end of the best epoch.

Epoch 00006: val_aucroc did not improve from 0.75773

Epoch 00006: early stopping

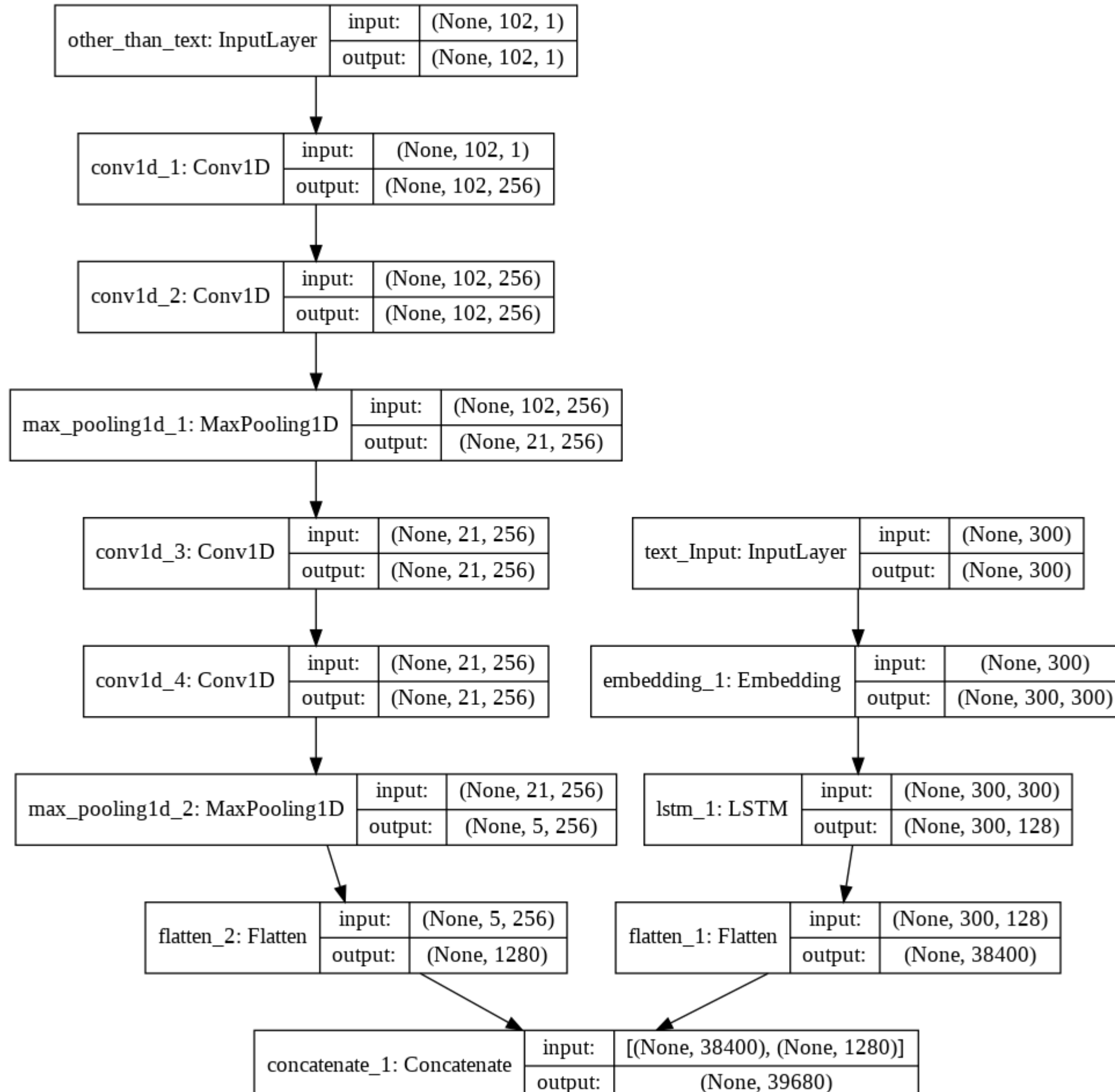
In [23]: 1 model3.summary()

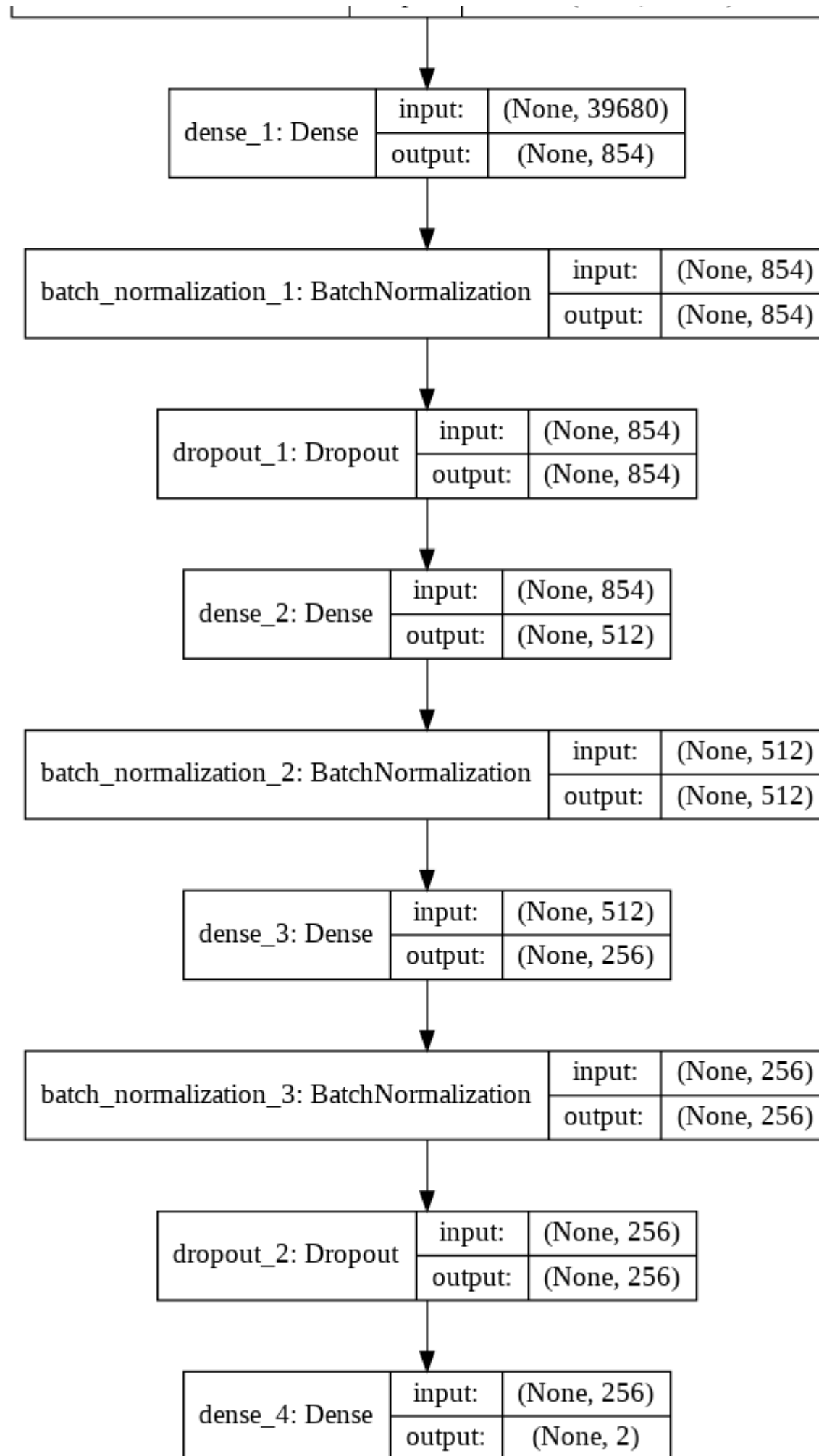
Model: "model_1"

Layer (type)	Output Shape	Param #	Connected to
=====			
other_than_text (InputLayer)	(None, 102, 1)	0	
conv1d_1 (Conv1D)	(None, 102, 256)	2816	other_than_text[0][0]
conv1d_2 (Conv1D)	(None, 102, 256)	655616	conv1d_1[0][0]
max_pooling1d_1 (MaxPooling1D)	(None, 21, 256)	0	conv1d_2[0][0]
text_input (InputLayer)	(None, 300)	0	
conv1d_3 (Conv1D)	(None, 21, 256)	655616	max_pooling1d_1[0][0]
embedding_1 (Embedding)	(None, 300, 300)	14160600	text_input[0][0]
conv1d_4 (Conv1D)	(None, 21, 256)	655616	conv1d_3[0][0]
lstm_1 (LSTM)	(None, 300, 128)	219648	embedding_1[0][0]
max_pooling1d_2 (MaxPooling1D)	(None, 5, 256)	0	conv1d_4[0][0]
flatten_1 (Flatten)	(None, 38400)	0	lstm_1[0][0]
flatten_2 (Flatten)	(None, 1280)	0	max_pooling1d_2[0][0]
concatenate_1 (Concatenate)	(None, 39680)	0	flatten_1[0][0] flatten_2[0][0]
dense_1 (Dense)	(None, 854)	33887574	concatenate_1[0][0]
batch_normalization_1 (BatchNor	(None, 854)	3416	dense_1[0][0]
dropout_1 (Dropout)	(None, 854)	0	batch_normalization_1[0][0]
dense_2 (Dense)	(None, 512)	437760	dropout_1[0][0]
batch_normalization_2 (BatchNor	(None, 512)	2048	dense_2[0][0]
dense_3 (Dense)	(None, 256)	131328	batch_normalization_2[0][0]
batch_normalization_3 (BatchNor	(None, 256)	1024	dense_3[0][0]
dropout_2 (Dropout)	(None, 256)	0	batch_normalization_3[0][0]
dense_4 (Dense)	(None, 2)	514	dropout_2[0][0]
=====			
Total params: 50,813,576			
Trainable params: 36,649,732			
Non-trainable params: 14,163,844			

```
In [24]: 1 #https://machinelearningmastery.com/visualize-deep-learning-neural-network-model-keras/
2 from keras.utils.vis_utils import plot_model
3 plot_model(model3, to_file='/content/drive/My Drive/LSTM_preprocessed/model3/model3_t2.png', show_shapes=True, show_layer_names=True)
```

Out[24]:



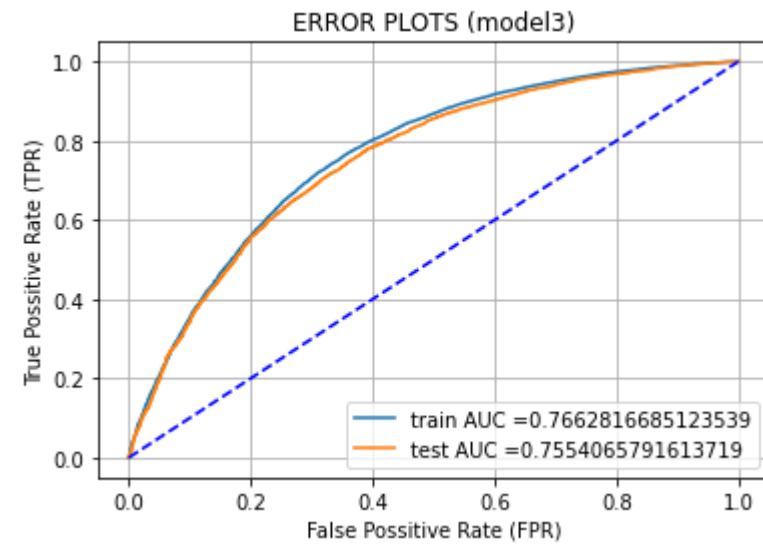


In [26]:

```

1  2  ## Predict the test and train
3  from sklearn.metrics import confusion_matrix
4  ## Finding best threshold for predictions
5  def best_threshold(thresholds,fpr,tpr):
6      t=thresholds[np.argmax(tpr*(1-fpr))]
7      # (tpr*(1-fpr)) will be maximum if your fpr is very Low and tpr is very high
8      print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
9      return t
10
11 def predict_with_best_t(proba, threshold):
12     predictions = []
13     for i in proba:
14         if i>=threshold:
15             predictions.append(1)
16         else:
17             predictions.append(0)
18     return predictions
19
20
21 y_test_predict = model3.predict([padseq_x_test,other_than_text_test],use_multiprocessing=True)[: ,1]
22 y_train_predict = model3.predict([padseq_x_train,other_than_text_train],use_multiprocessing=True)[: ,1]
23
24 # if os.path.isfile('model_predictions.pkl'):
25 #     os.remove('model_predictions.pkl')
26 #     print("File model_predictions Removed!")
27 #     with open('model_predictions.pkl','wb') as f:
28 #         pickle.dump([y_train_predict,y_test_predict],f)
29
30 ## Store fpr and tpr rates
31
32 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_predict)
33 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_predict)
34
35
36 #plot
37 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
38 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
39 plt.legend()
40 plt.xlabel("False Possitive Rate (FPR)")
41 plt.ylabel("True Possitive Rate (TPR)")
42 plt.title("ERROR PLOTS (model3)")
43 plt.plot([0, 1], [0, 1], 'b--')
44 plt.grid()
45 plt.show()
46
47 print("="*100)
48
49 best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
50

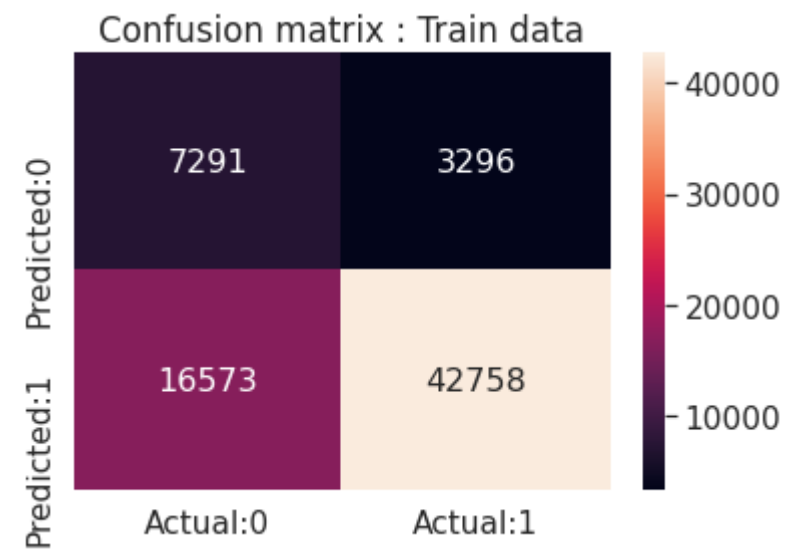
```



=====

the maximum value of $tpr \cdot (1 - fpr)$ 0.4963064277331075 for threshold 0.811

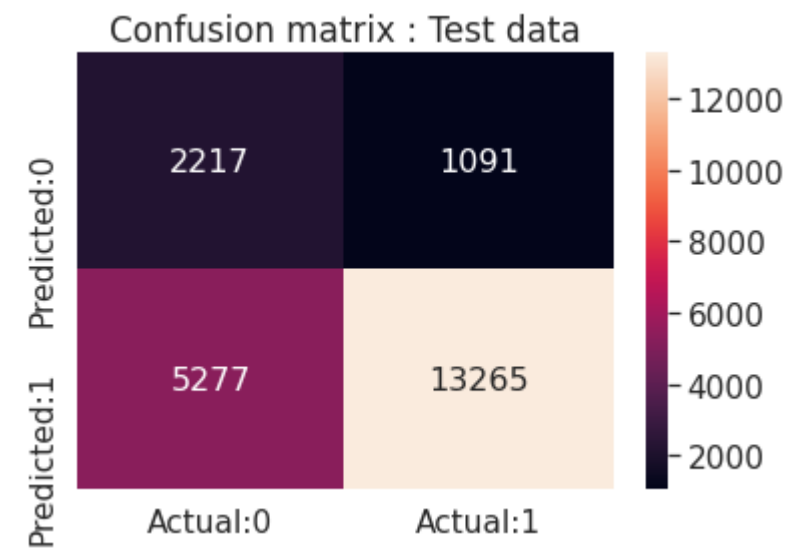
```
In [27]: ▶ 1 ### PLOT the matrix for Train
2 #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
3 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
4 df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_predict, best_t))
5                       , range(2), range(2))
6 # plt.figure(figsize=(10,7))
7 sns.set(font_scale=1.4) # for label size
8 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16}, fmt='g',
9             xticklabels=['Actual:0', 'Actual:1']
10             , yticklabels=['Predicted:0', 'Predicted:1']) # font size
11 plt.title('Confusion matrix : Train data')
12 plt.show()
```



```

In [28]: 1  ### PLOT the matrix for Test
          2  #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
          3  # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
          4  df_cm = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(y_test_predict, best_t))
          5                               , range(2), range(2))
          6  # plt.figure(figsize=(10,7))
          7  sns.set(font_scale=1.4) # for label size
          8  sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g',
          9                      xticklabels=['Actual:0','Actual:1']
         10                      ,yticklabels=['Predicted:0','Predicted:1']) # font size
         11  plt.title('Confusion matrix : Test data')
         12  plt.show()

```



```

1  <li><b>Validation Loss,Validation aucroc,Train Loss,Train aucroc</b>
2  <img src='model3.jpg' width="800" height="800">
3
4

```

```
In [1]: 1 from prettytable import PrettyTable
2 pt=PrettyTable()
3 pt.field_names=["model","train_auc","test_auc"]
4 pt.add_row(["model_1","0.74","0.71"])
5 pt.add_row(["model_2","0.78","0.72"])
6 pt.add_row(["model_3","0.766","0.7554"])
7
8 print(pt)
```

model	train_auc	test_auc
model_1	0.74	0.71
model_2	0.78	0.72
model_3	0.766	0.7554

Observations

- * Model3 performs best among all on test data with 0.76 auc with CNN and max pooling layers as features to our dense layers.
- * As the number of units in dense layer increases the model tends to overfit hence we regularize using dropouts in between layers and l2 regularizer in all the dense layers.
- * using maxpool represents data in convoltional layer 1 and layer 2 in a abstracted way and avoids overfitting and reduces computational time while evaluation.
- * Using filtered words in Model2 in our dataset reduced train time and produced good results.
- * Model3 can be considered best fit among all as train and test AUC's are more close than Model1 and Model2 and also resulted with highest AUC.
- * using early stopping we are fitting data in epoch which produced best validation auc and stop the training when there is no improvement.