

## Assignment : 14

1. You can work with `preprocessed_data.csv` for the assignment. You can get the data from the link.
2. Load the data in your notebook.
3. After step 2 you have to train 3 types of models as discussed below.
4. For all the model use `'auc'` as a metric. check [this](#) and [this](#) for using auc as a metric.
5. You are free to choose any number of layers/hiddenn units but you have to use `'relu'` as an activation function.
6. You can use any one of the optimizers and choice of Learning rate and momentum.
7. For all the model's use [TensorBoard](#) and plot the Metric value and Loss with epoch.
8. Make sure that you are using GPU to train the given models.

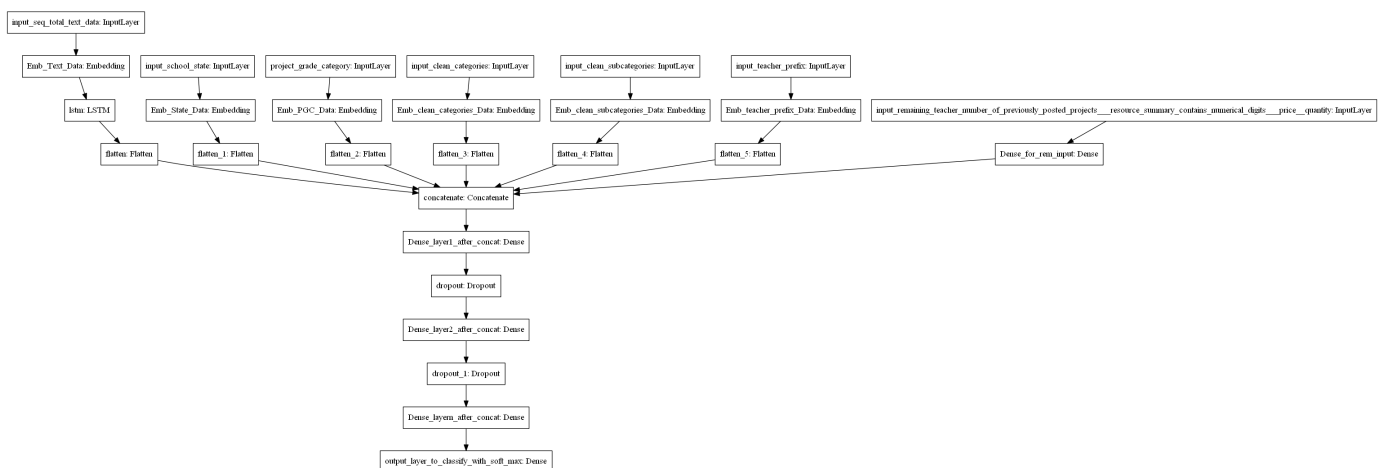
```

1 #you can use gdown modules to import dataset for the assignment
2 #for importing any file from drive to Colab you can write the syntax as !gdown
3 #you can run the below cell to import the required preprocessed data.csv file

```

## Model-1

Build and Train deep neural network as shown below



ref: <https://i.imgur.com/w395Yk9.png>

- **Input\_seq\_total\_text\_data** --- You have to give Total text data columns. After this use the Embedding layer to get word vectors. Use given predefined glove word vectors, don't train any word vectors. After this use LSTM and get the LSTM output and Flatten that output.
- **Input\_school\_state** --- Give 'school\_state' column as input to embedding layer and Train the Keras Embedding layer.

- **Project\_grade\_category** --- Give 'project\_grade\_category' column as input to embedding layer and Train the Keras Embedding layer.
- **Input\_clean\_categories** --- Give 'input\_clean\_categories' column as input to embedding layer and Train the Keras Embedding layer.
- **Input\_clean\_subcategories** --- Give 'input\_clean\_subcategories' column as input to embedding layer and Train the Keras Embedding layer.
- **Input\_clean\_subcategories** --- Give 'input\_teacher\_prefix' column as input to embedding layer and Train the Keras Embedding layer.
- **Input\_remaining\_teacher\_number\_of\_previously\_posted\_projects.\_resource\_summary\_contains\_numerical\_digits.\_price.\_quantity** ---concatenate remaining columns and add a Dense layer after that.

Below is an example of embedding layer for a categorical columns. In below code all are dummy values, we gave only for reference.

```
1 # # https://stats.stackexchange.com/questions/270546/how-does-keras-embedding-l
2 # input_layer = Input(shape=(n,))
3 # embedding = Embedding(no_1, no_2, input_length=n)(input_layer)
4 # flatten = Flatten()(embedding)
```

1. Go through this blog, if you have any doubt on using predefined Embedding values in Embedding layer - <https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-keras/>
2. Please go through this link <https://keras.io/getting-started/functional-api-guide/> and check the 'Multi-input and multi-output models' then you will get to know how to give multiple inputs.

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call

## ▼ Model-1

```
1 # import all the libraries
2 #make sure that you import your libraries from tf.keras and not just keras
3 import tensorflow as tf
```

```

1
2 # path = "/Users/yamasanimanoj-kumarreddy/Documents/AAIC/Named_Assignments/LSTM
3 path = "/content/drive/MyDrive/Named_Assignments/LSTM_DonorsChoose"
4 # path = "/content/drive/MyDrive/DeepLearn/LSTM_DonorsChoose"

```

```

1 #read the csv file
2 import pandas as pd
3
4 data = pd.read_csv(path+'/preprocessed_data.csv')
5 data.head(1)

```

school_state	teacher_prefix	project_grade_category	teacher_number_of_pre
--------------	----------------	------------------------	-----------------------

0	ca	mrs	grades_prek_2
---	----	-----	---------------



```

1 Y=data['project_is_approved']
2 X=data.drop('project_is_approved',axis=1)

```

```
1 print(Y.value_counts())
```

```

1    92706
0    16542
Name: project_is_approved, dtype: int64

```

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 pos_class = Y.value_counts()[1]/sum(Y.value_counts())
5 neg_class = Y.value_counts()[0]/sum(Y.value_counts())
6 y = np.array([neg_class, pos_class])
7 print(y)
8 mylabels = ["Negative","Positive"]
9
10 plt.pie(y, labels = mylabels)
11 plt.legend(title = "Distribution of classes")
12 plt.show()

```

```
[0.15141696 0.84858304]
```



From the above plot we can conclude that 85% of the times a project is approved.

Positive

From the above pie chart, we can observe that given data is highly imbalanced data set with 85% points of positive class and remaining 15% of points as negative class.

```
1 # Scaling by total/2 helps keep the loss to a similar magnitude.
2 # The sum of the weights of all examples stays the same.
3
4 total = sum(Y.value_counts())
5 weight_for_0 = (1 / Y.value_counts()[0]) * (total / 2.0)
6 weight_for_1 = (1 / Y.value_counts()[1]) * (total / 2.0)
7
8 class_weight = {0: weight_for_0, 1: weight_for_1}
9
10 print('Weight for class 0: {:.2f}'.format(weight_for_0))
11 print('Weight for class 1: {:.2f}'.format(weight_for_1))
```

```
Weight for class 0: 3.30
Weight for class 1: 0.59
```

## ▼ perform stratified train test split on the dataset

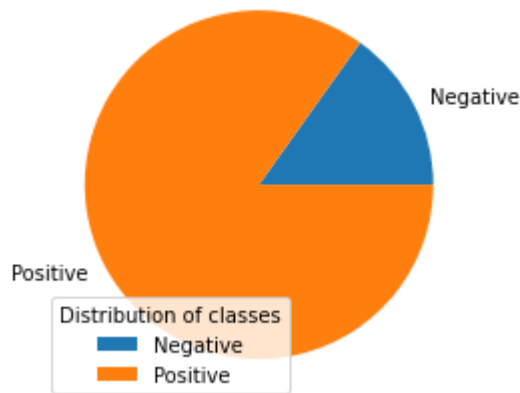
```
1 from sklearn.model_selection import train_test_split
2
3 X_train,X_test,Y_train,Y_test = train_test_split(X,Y,test_size=0.33, stratify=Y)
```

```
1 print(X_train.shape,Y_train.shape,X_test.shape,Y_test.shape,end=" ")

(73196, 8) (73196,) (36052, 8) (36052,)
```

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 pos_class = Y_train.value_counts()[1]/sum(Y_train.value_counts())
5 neg_class = Y_train.value_counts()[0]/sum(Y_train.value_counts())
6 y = np.array([neg_class, pos_class])
7 print(y)
8 mylabels = ["Negative","Positive"]
9
10 plt.pie(y, labels = mylabels)
11 plt.legend(title = "Distribution of classes")
12 plt.show()
```

```
[0.15141538 0.84858462]
```



We got similar distribution of points in class due to stratified sampling.

## ▼ 1.1 Text Vectorization

```
1 #since the data is already preprocessed, we can directly move to vectorization
2 #first we will vectorize the text data
3 #for vectorization of text data in deep learning we use tokenizer, you can go
4 # https://www.kdnuggets.com/2020/03/tensorflow-keras-tokenization-text-data-pre
5 # https://stackoverflow.com/questions/51956000/what-does-keras-tokenizer-method-
6 # after text vectorization you should get train_padded_docs and test_padded_do
```

```
1 from tensorflow.keras.preprocessing.text import Tokenizer
2 from tensorflow.keras.preprocessing.sequence import pad_sequences
3
4 tokenizer = Tokenizer()
5 tokenizer.fit_on_texts(X_train['essay'])
6 vocab_size = len(tokenizer.word_index) + 1
7 print(vocab_size)
8 # integer encode the documents
9 train_encoded_docs = tokenizer.texts_to_sequences(X_train['essay'])
10 test_encoded_docs = tokenizer.texts_to_sequences(X_test['essay'])
```

```
48232
```

```
1 max_length = max([len(seq) for seq in train_encoded_docs])
2 print(max_length)
3 train_padded_docs = pad_sequences(train_encoded_docs, maxlen=max_length, padding=
4 test_padded_docs = pad_sequences(test_encoded_docs, maxlen=max_length, padding=
5 print(train_padded_docs.shape)
6 print(test_padded_docs.shape)
7 print(train_padded_docs[0][-10:])
8 print(test_padded_docs[0][-10:])
```

```
339
(73196, 339)
(36052, 339)
```

```
[0 0 0 0 0 0 0 0 0 0]
[0 0 0 0 0 0 0 0 0 0]
```

```
1 #after getting the padded_docs you have to use predefined glove vectors to get
2 # we will be storing this data in form of an embedding matrix and will use it v
3 # Please go through following blog's 'Example of Using Pre-Trained GloVe Embed
4 # https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-ke
5 import pickle
6 import numpy as np
7
8 with open(path+'/glove_vectors', 'rb') as f:
9     model = pickle.load(f)
10     glove_words = set(model.keys())
```

```
1 embedding_matrix = np.zeros((vocab_size, 300))
2 for word, i in tokenizer.word_index.items():
3     embedding_vector = model.get(word)
4     if embedding_vector is not None:
5         embedding_matrix[i] = embedding_vector
6 print(embedding_matrix[1][:5])
7 print(embedding_matrix.shape)
8 print("The number of tokens present in this vocabulary are", embedding_matrix.s
```

```
[ 0.15243 -0.16945 -0.022748 -0.25051 -0.15213 ]
(48232, 300)
```

The number of tokens present in this vocabulary are 48232

## ▼ 1.2 Categorical feature Vectorization

```
1 # for model 1 and model 2, we have to assign a unique number to each feature in
2 # you can either use tokenizer,label encoder or ordinal encoder to perform the
3 # label encoder gives an error for 'unseen values' (values present in test but
4 # handle unseen values with label encoder - https://stackoverflow.com/a/5687635
5 # ordinal encoder also gives error with unseen values but you can use modify ha
6 # documentation of ordianl encoder https://scikit-learn.org/stable/modules/gene
7 # after categorical feature vectorization you will have column_train_data and c
8
```

### ▼ School\_State

```
1 from sklearn.preprocessing import OrdinalEncoder
2 encoder = OrdinalEncoder(handle_unknown = 'use_encoded_value',unknown_value=-1)
3 train_school_state = encoder.fit_transform(X_train['school_state'].to_numpy()).r
4 test_school_state = encoder.transform(X_test['school_state'].to_numpy()).resha
5 no_of_school_state = len(encoder.categories_[0])
6 print("No of categories present in school state are ",no_of_school_state)
7 print(train_school_state.shape)
8 print(test_school_state.shape)
```

```
No of categories present in school state are 51
(73196, 1)
(36052, 1)
```

### ▼ project\_grade\_category

```
1 from sklearn.preprocessing import OrdinalEncoder
2 encoder = OrdinalEncoder(handle_unknown = 'use_encoded_value',unknown_value=-1)
3 train_project_grade = encoder.fit_transform(X_train['project_grade_category'].to_numpy())
4 test_project_grade = encoder.transform(X_test['project_grade_category'].to_numpy())
5 no_of_project_grade = len(encoder.categories_[0])
6 print("No of categories present in project grade are ",no_of_project_grade)
7 print(train_project_grade.shape)
8 print(test_project_grade.shape)
```

```
No of categories present in project grade are 4
(73196, 1)
(36052, 1)
```

### ▼ clean\_categories

```
1 from sklearn.preprocessing import OrdinalEncoder
2 encoder = OrdinalEncoder(handle_unknown = 'use_encoded_value',unknown_value=-1)
3 train_clean_categories = encoder.fit_transform(X_train['clean_categories'].to_numpy())
4 test_clean_categories = encoder.transform(X_test['clean_categories'].to_numpy())
5 no_of_clean_categories = len(encoder.categories_[0])
6 print("No of categories present in clean categories are ",no_of_clean_categories)
7 print(train_clean_categories.shape)
8 print(test_clean_categories.shape)
```

```
No of categories present in clean categories are 51
(73196, 1)
(36052, 1)
```

### ▼ clean\_subcategories

```
1 from sklearn.preprocessing import OrdinalEncoder
2 encoder = OrdinalEncoder(handle_unknown = 'use_encoded_value',unknown_value=-1)
3 train_clean_subcategories = encoder.fit_transform(X_train['clean_subcategories'].to_numpy())
4 test_clean_subcategories = encoder.transform(X_test['clean_subcategories'].to_numpy())
5 no_of_clean_subcategories = len(encoder.categories_[0])
6 print("No of categories present in clean sub categories are ",no_of_clean_subcategories)
7 print(train_clean_subcategories.shape)
8 print(test_clean_subcategories.shape)
```

```
No of categories present in clean sub categories are 390
(73196, 1)
(36052, 1)
```

## ▼ teacher\_prefix

```
1 from sklearn.preprocessing import OrdinalEncoder
2 encoder = OrdinalEncoder(handle_unknown = 'use_encoded_value',unknown_value=-1)
3 train_teacher_prefix = encoder.fit_transform(X_train['teacher_prefix'].to_numpy
4 test_teacher_prefix = encoder.transform(X_test['teacher_prefix'].to_numpy().re
5 no_of_teacher_prefix = len(encoder.categories_[0])
6 print("No of categories present in teacher prefix are ",no_of_teacher_prefix)
7 print(train_teacher_prefix.shape)
8 print(test_teacher_prefix.shape)
```

```
No of categories present in teacher prefix are 5
(73196, 1)
(36052, 1)
```

## ▼ 1.3 Numerical feature Vectorization

```
1 # you have to standardise the numerical columns
2 # stack both the numerical features
3 #after numerical feature vectorization you will have numerical_data_train and r
```

## ▼ teacher\_number\_of\_previously\_posted\_projects

### ▶ Min-Max scaling of features

```
[ ] ↪ 1 cell hidden
```

### ▶ Standard scaling of features

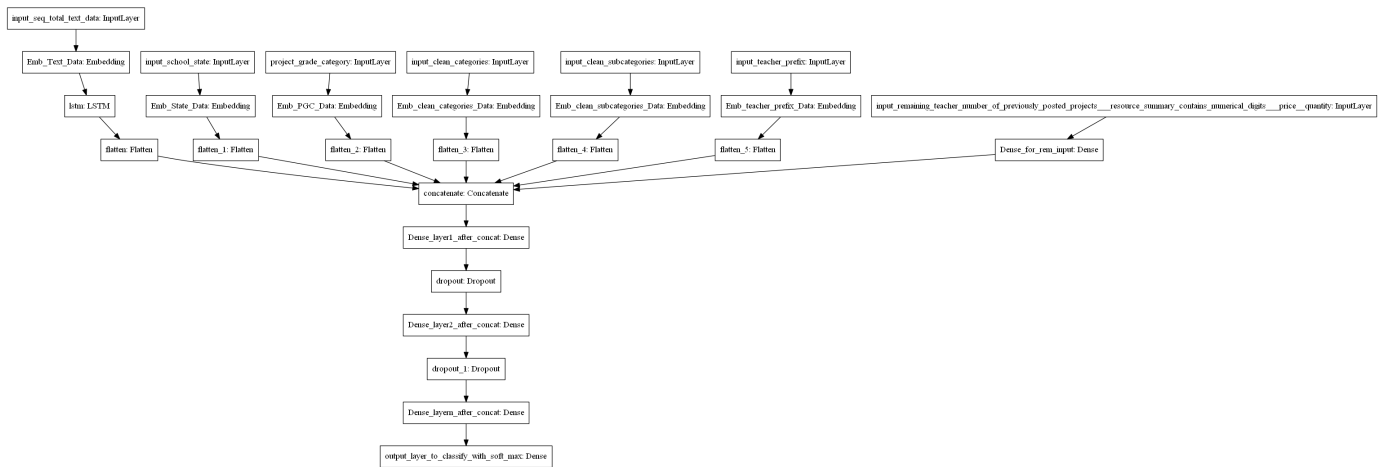
```
[ ] ↪ 2 cells hidden
```

### ▶ price

```
[ ] ↪ 6 cells hidden
```

## ▼ 1.4 Defining the model





```

1 # as of now we have vectorized all our features now we will define our model.
2 # as it is clear from above image that the given model has multiple input layer
3 # Please go through - https://keras.io/guides/functional\_api/
4 # it is a good programming practise to define your complete model i.e all input
5 # while defining your model make sure that you use variable names while defining
6 # for ex.- you should write the code as 'input_text = Input(shape=(pad_length,))'
7 # the embedding layer for text data should be non trainable
8 # the embedding layer for categorical data should be trainable
9 # https://stats.stackexchange.com/questions/270546/how-does-keras-embedding-layer
10 # https://towardsdatascience.com/deep-embeddings-for-categorical-variables-categorical-embeddings
11 # print model.summary() after you have defined the model
12 # plot the model using utils.plot_model module and make sure that it is similar

```

```

1 from tensorflow.keras import layers, Input, Model
2
3 # Hyper Parameters
4 Embedding_vector_size_text = 300
5 categorical_ftr_embedding_size = 20
6 max_seq_length = train_padded_docs.shape[1]
7 no_of_lstm_units_on_text = 150
8 no_of_Dense_units_on_numerical_ftr = 150
9 no_of_units_on_dense_layer1 = 128
10 no_of_units_on_dense_layer2 = 64
11 no_of_units_on_dense_layer3 = 32
12 no_of_units_on_output_dense_layer = 1
13
14 vocab_size_school_state = no_of_school_state
15 vocab_size_project_grade = no_of_project_grade
16 vocab_size_input_clean = no_of_clean_categories
17 vocab_size_input_clean_sub = no_of_clean_subcategories
18 vocab_size_teacher_prefix = no_of_teacher_prefix
19
20 def create_model(max_seq_length, vocab_size_text, embedding_matrix, categories_
21                  categorical_ftr_embedding_size, initializer, numerical_ftr_dim,
22                  tf.keras.backend.clear_session())

```

```

23 # Text Feature block
24 input1 = Input(shape = (max_seq_length,), name = "Text_Input")
25 x = layers.Embedding(input_dim = vocab_size_text, output_dim=embedding_matrix
26                       input_length = max_seq_length, weights=[embedding_matrix
27                       name = "Text_Embedding_layer", trainable = False)(input1
28 x = layers.LSTM(no_of_lstm_units_on_text, return_sequences = True, name = "Text
29 text_lstm_output = layers.Flatten(name = "Output_from_text_lstm")(x)
30
31 # Categorical feature school_state block
32 input2 = Input(shape = (categories_dim,), name = "School_state_Input")
33 x = layers.Embedding(input_dim = vocab_size_school_state,
34                       output_dim= categorical_ftr_embedding_size,
35                       embeddings_initializer = initializer,
36                       input_length = categories_dim,
37                       name = "School_state_Embedding_layer")(input2)
38 school_state_output = layers.Flatten(name = "school_state_output")(x)
39
40 # Categorical feature project_grade block
41 input3 = Input(shape = (categories_dim,), name = "project_grade_Input")
42 x = layers.Embedding(input_dim = vocab_size_project_grade,
43                       output_dim = categorical_ftr_embedding_size,
44                       embeddings_initializer = initializer,
45                       input_length = categories_dim,
46                       name = "Project_grade_Embedding_layer")(input3)
47 project_grade_output = layers.Flatten(name = "project_grade_output")(x)
48
49 # Categorical feature input_clean_categories block
50 input4 = Input(shape = (categories_dim,), name = "input_clean_categories_Input")
51 x = layers.Embedding(input_dim = vocab_size_input_clean,
52                       output_dim = categorical_ftr_embedding_size,
53                       embeddings_initializer = initializer,
54                       input_length = categories_dim,
55                       name = "Clean_categories_Embedding_layer")(input4)
56 clean_cat_output = layers.Flatten(name = "clean_cat_output")(x)
57
58 # Categorical feature input_clean_subcategories block
59 input5 = Input(shape = (categories_dim,), name = "input_clean_subcategories_Input")
60 x = layers.Embedding(input_dim = vocab_size_input_clean_sub,
61                       output_dim = categorical_ftr_embedding_size,
62                       embeddings_initializer = initializer,
63                       input_length = categories_dim,
64                       name = "Clean_subcategories_Embedding_layer")(input5)
65 clean_subcat_output = layers.Flatten(name = "clean_subcat_output")(x)
66
67 # Categorical feature teacher_prefix block
68 input6 = Input(shape = (categories_dim,), name = "teacher_prefix_Input")
69 x = layers.Embedding(input_dim = vocab_size_teacher_prefix,
70                       output_dim = categorical_ftr_embedding_size,
71                       embeddings_initializer = initializer,
72                       input_length = categories_dim,
73                       name = "Teacher_prefix_Embedding_layer")(input6)
74 teacher_prefix_output = layers.Flatten(name = "teacher_prefix_output")(x)
75
76 # Numerical features block
77 input7 = Input(shape = (numerical_ftr_dim,), name = "numerical_ftrs_Input")

```

```

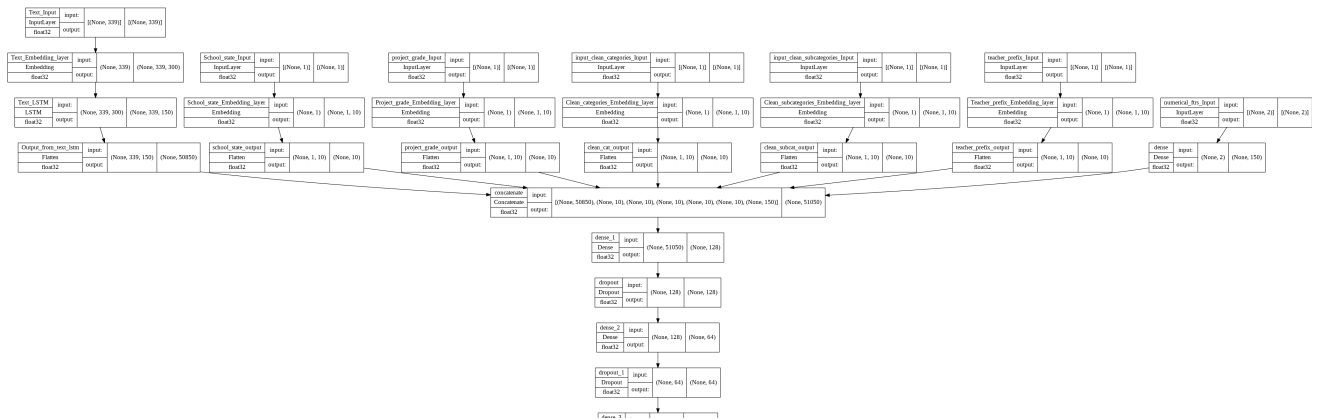
78 numerical_ftr_output = layers.Dense(no_of_Dense_units_on_numerical_ftr,
79                                     activation= activation_function,
80                                     kernel_initializer = initializer)(input7)
81
82 concat_output = layers.Concatenate(axis=1)([text_lstm_output, school_state_o
83                                             clean_cat_output, clean_subcat_o
84                                             numerical_ftr_output])
85
86 dense_layer1_after_concat = layers.Dense(no_of_units_on_dense_layer1,
87                                           activation = activation_function,
88                                           kernel_initializer=initializer)(conca
89 dropout_layer1 = layers.Dropout(dropout_rate)(dense_layer1_after_concat)
90
91 dense_layer2_after_concat = layers.Dense(no_of_units_on_dense_layer2,
92                                           activation = activation_function,
93                                           kernel_initializer=initializer)(dropo
94 dropout_layer2 = layers.Dropout(dropout_rate)(dense_layer2_after_concat)
95
96 dense_layer3_after_concat = layers.Dense(no_of_units_on_dense_layer3,
97                                           activation = activation_function,
98                                           kernel_initializer=initializer)(dropo
99 output = layers.Dense(no_of_units_on_output_dense_layer,
100                       activation = 'sigmoid',
101                       kernel_initializer=initializer)(dense_layer3_after_conca
102
103 return Model(inputs = [input1 input2 input3 input4 input5 input6 input7] outp

```

```

1 from tensorflow.keras import utils,initializers
2
3 max_seq_length = train_padded_docs.shape[1]
4 vocab_size_text = vocab_size
5 categories_dim = 1
6 categorical_ftr_embedding_size = 10
7 initializer = initializers.HeUniform()
8 activation_function = 'relu'
9 dropout_rate = 0.3
10 numerical_ftr_dim = 2
11
12 model = create_model(max_seq_length, vocab_size_text, embedding_matrix, categor
13                     categorical_ftr_embedding_size, initializer,numerical_ftr_dim,
14
15 utils.plot_model(model, "Model_1.png", show_shapes=True,show_dtype = True)

```



## 1.5 Compiling and fitting your model

```

1 #define custom auc as metric , do not use tf.keras.metrics
2 # https://stackoverflow.com/a/46844409 - custom AUC reference 1
3 # https://www.kaggle.com/c/santander-customer-transaction-prediction/discussion
4 # compile and fit your model

```

```

1 from sklearn.metrics import roc_auc_score,roc_curve,auc
2 from keras.callbacks import Callback
3 class RocCallback(Callback):
4     def __init__(self,training_data,validation_data):
5         self.x = training_data[0]
6         self.y = training_data[1]
7         self.x_val = validation_data[0]
8         self.y_val = validation_data[1]
9
10
11     def on_train_begin(self, logs={}):
12         self.model.auc_train = []
13         self.model.auc_val = []
14         return
15
16     def on_train_end(self, logs={}):
17         return
18
19     def on_epoch_begin(self, epoch, logs={}):
20         return
21
22     def on_epoch_end(self, epoch, logs={}):
23         y_pred_train = self.model.predict(self.x)
24         auc_train = roc_auc_score(self.y, y_pred_train)
25         y_pred_val = self.model.predict(self.x_val)
26         auc_val = roc_auc_score(self.y_val, y_pred_val)
27         self.model.auc_train.append(auc_train)
28         self.model.auc_val.append(auc_val)
29         print('\rauc_train: %s - auc_val: %s' % (str(round(auc_train,4)),str(round(auc_val,4))))
30         return
31
32     def on_batch_begin(self, batch, logs={}):
33         return

```

```

34
35     def on_batch_end(self, batch, logs={}):
36         return

1 from tensorflow.keras import callbacks
2 import os
3
4 # path = "/Users/yamasanimanoj-kumarreddy/Documents/AAIC/Named_Assignments/LSTM
5 # path = "/content/drive/MyDrive/Named_Assignments/LSTM_DonorsChoose"
6
7 graph_saving_dir = os.path.join(path, 'model_1', "Tensorboard_graphs")
8
9 roc_callback = RocCallback(training_data=( [train_padded_docs, train_school_state
10                                     train_project_grade, train_clean_categories,
11                                     train_clean_subcategories, train_teacher_prefi
12                                     numerical_data_train ],
13                                     Y_train),
14                               validation_data=( [test_padded_docs, test_school_state,
15                                     test_project_grade, test_clean_categories,
16                                     test_clean_subcategories, test_teacher_prefix,
17                                     numerical_data_test ],
18                                     Y_test))
19
20
21 tensorboard_callback = callbacks.TensorBoard(log_dir=graph_saving_dir, histogram
22
23 callbacksList = [roc_callback, tensorboard_callback]

```

```

1  from tensorflow.keras import initializers, optimizers, metrics, losses
2  from sklearn.utils.class_weight import compute_class_weight
3
4  # Types of intializers to try
5  # initializers.HeUniform()
6  # initializers.HeNormal()
7  # initializers.GlorotNormal()
8  # initializers.GlorotUniform()
9  # initializers.LecunNormal()
10 # initializers.LecunUniform()
11 # initializers.RandomNormal(mean=0., stddev=1.)
12 # initializers.RandomUniform(minval=0., maxval=1.)
13
14 # Types of Optimizers to try
15 # optimizers.Adadelata(learning_rate=0.001)
16 # optimizers.Adagrad(learning_rate=0.001)
17 # optimizers.Adamax(learning_rate=0.001)
18 # optimizers.RMSprop(learning_rate=0.001)
19 # tf.keras.optimizers.SGD(learning_rate=0.01, momentum=0.0)
20
21 opt = optimizers.Adam()
22 max_seq_length = train_padded_docs.shape[1]
23 vocab_size_text = vocab_size
24 categories_dim = 1
25 categorical_ftr_embedding_size = 20
26 initializer = initializers.HeNormal()
27 activation_function = 'relu'

```

```

27 activation_function = 'relu'
28 dropout_rate = 0.7
29 numerical_ftr_dim = 2
30 class_weights = compute_class_weight(class_weight = "balanced", classes = np.un
31 class_weights = dict(zip(np.unique(Y_train), class_weights))
32 print(class_weights)
33
34 model = create_model(max_seq_length, vocab_size_text, embedding_matrix, catego
35 categorical_ftr_embedding_size, initializer, numerical_ftr_dim
36
37 batch_size = 1024
38 epochs = 10
39
40 model.compile(optimizer=opt,
41               metrics = [metrics.BinaryAccuracy()],
42               loss=losses.BinaryCrossentropy())
43
44 # fit the model
45 model.fit([train_padded_docs, train_school_state,
46           train_project_grade, train_clean_categories,
47           train_clean_subcategories, train_teacher_prefix,
48           numerical_data_train ],
49         Y_train, batch_size = batch_size,
50         epochs=epochs, callbacks=callbacksList,
51         class_weight = class_weights,
52         validation_data = ([ test_padded_docs, test_school_state,
53                             test_project_grade, test_clean_categories,
54                             test_clean_subcategories, test_teacher_prefix,
55                             numerical_data_test ],
56                             Y_test))
57 # evaluate the model
58 loss, accuracy = model.evaluate([ test_padded_docs, test_school_state,
59                                   test_project_grade, test_clean_categories,
60                                   test_clean_subcategories, test_teacher_prefix,
61                                   numerical_data_test ],
62                                   Y_test, batch_size = batch_size)
63 print('Test Accuracy: %f' % (accuracy*100))

```

```
{0: 3.3021745014887665, 1: 0.5892164281229372}
```

```
Epoch 1/10
```

```
6/72 [=>.....] - ETA: 16s - loss: 0.9706 - binary_accu
auc_train: 0.5047 - auc_val: 0.5051
```

```
72/72 [=====] - 59s 791ms/step - loss: 0.7168 - binar
```

```
Epoch 2/10
```

```
auc_train: 0.509 - auc_val: 0.5069
```

```
72/72 [=====] - 56s 789ms/step - loss: 0.6936 - binar
```

```
Epoch 3/10
```

```
auc_train: 0.5068 - auc_val: 0.5072
```

```
72/72 [=====] - 54s 758ms/step - loss: 0.6935 - binar
```

```
Epoch 4/10
```

```
auc_train: 0.5181 - auc_val: 0.5192
```

```
72/72 [=====] - 55s 764ms/step - loss: 0.6931 - binar
```

```
Epoch 5/10
```

```
auc_train: 0.5629 - auc_val: 0.5615
```

```
72/72 [=====] - 54s 758ms/step - loss: 0.6931 - binar
```

```
Epoch 6/10
```

```
auc_train: 0.7013 - auc_val: 0.6898
```

```

72/72 [=====] - 49s 688ms/step - loss: 0.6895 - binar
Epoch 7/10
auc_train: 0.718 - auc_val: 0.7005
72/72 [=====] - 56s 783ms/step - loss: 0.6850 - binar
Epoch 8/10
auc_train: 0.7322 - auc_val: 0.7158
72/72 [=====] - 55s 772ms/step - loss: 0.6695 - binar
Epoch 9/10
auc_train: 0.7667 - auc_val: 0.7397
72/72 [=====] - 51s 715ms/step - loss: 0.6516 - binar
Epoch 10/10
auc_train: 0.777 - auc_val: 0.7488
72/72 [=====] - 56s 781ms/step - loss: 0.6367 - binar
36/36 [=====] - 3s 83ms/step - loss: 0.6533 - binary_
Test Accuracy: 68.492734

```



```

1 #After fit method accessing the auc scores of train and validation
2 auc_train = model.auc_train
3 auc_val = model.auc_val
4
5 writer=tf.summary.create_file_writer(graph_saving_dir)
6 for idx in range(len(auc_train)):
7     with writer.as_default(step=idx+1):
8         tf.summary.scalar('auc_train', auc_train[idx])
9         tf.summary.scalar('auc_val', auc_val[idx])
10 writer.flush ()

```

```
1 %load_ext tensorboard
```

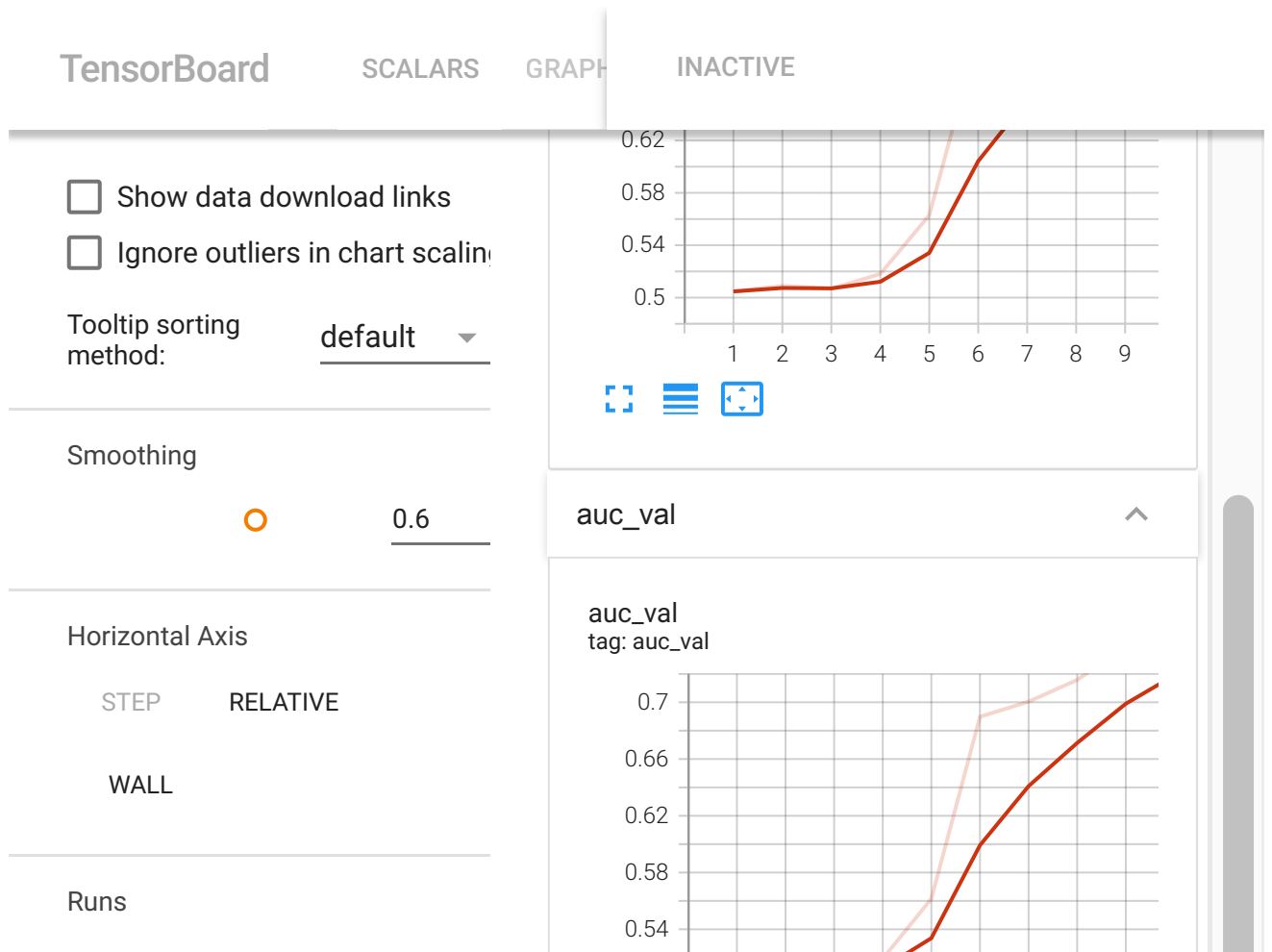
The tensorboard extension is already loaded. To reload it, use:

```
%reload_ext tensorboard
```

```

1 %tensorboard --logdir /content/drive/MyDrive/Named_Assignments/LSTM_DonorsChoose
2 # %tensorboard --logdir /content/drive/MyDrive/DeepLearn/LSTM_DonorsChoose/mode
3 # %tensorboard --logdir /Users/yamasanimanoj-kumarreddy/Documents/AAIC/Named_As

```



```

1 !tensorboard dev upload --logdir /content/drive/MyDrive/Named Assignments/LSTM
2
3 # !tensorboard dev upload --logdir %tensorboard --logdir /content/drive/MyDrive
4 # !tensorboard dev upload --logdir /Users/yamasanimanoj-kumarreddy/Documents/A

```

Upload started and will continue reading any new data as it's added to the log

To stop uploading, press Ctrl-C.

New experiment created. View your TensorBoard at: <https://tensorboard.dev/expe>

[2022-07-26T21:30:19] Started scanning logdir.

[2022-07-26T21:30:23] Total uploaded: 80 scalars, 190 tensors (136.4 kB), 1 bi

Interrupted. View your TensorBoard at <https://tensorboard.dev/experiment/heE6t>

Traceback (most recent call last):

```

File "/usr/local/bin/tensorboard", line 8, in <module>
    sys.exit(run_main())
File "/usr/local/lib/python3.7/dist-packages/tensorboard/main.py", line 46,
    app.run(tensorboard.main, flags_parser=tensorboard.configure)
File "/usr/local/lib/python3.7/dist-packages/absl/app.py", line 308, in run
    _run_main(main, args)
File "/usr/local/lib/python3.7/dist-packages/absl/app.py", line 254, in _run
    sys.exit(main(argv))
File "/usr/local/lib/python3.7/dist-packages/tensorboard/program.py", line 2
    return runner(self.flags) or 0
File "/usr/local/lib/python3.7/dist-packages/tensorboard/uploader/uploader_s
    return _run(flags, self._experiment_url_callback)
File "/usr/local/lib/python3.7/dist-packages/tensorboard/uploader/uploader_s

```



```
intent.execute(server_info, channel)
KeyboardInterrupt
```

The above Tensorboard logs can be observed at this url

<https://tensorboard.dev/experiment/heE6bKqjS4yFjs10JVqL0Q/>

Observations from the above Training of model

1. By using different initializers for kernel in Dense layers we could train different patterns of input data which might increase auc and accuracy.
2. At first we have tried without including `return_sequences = True` parameter in LSTM of text data. Due to this we get only 1 vector as output from LSTM cell and as a single vector cannot capture the entire essence of text sequence we got very less accuracy and auc values. But when we include the parameter `return_sequences = True` we get the output from every timestamp which increases the pattern information gained and thereby we got significant improvement in Test Accuracy and Auc.
3. From the training results we can observe that model is not overfitted and achieved **validation auc = 0.75**

## ▼ Model-2

Use the same model as above but for 'input\_seq\_total\_text\_data' give only some words in the sentence not all the words. Filter the words as below.

1. Fit TF-IDF vectorizer on the Train data
2. Get the idf value for each word we have in the train data. Please go through [t](#)
3. Do some analysis on the Idf values and based on those values choose the low and frequent words and very very rare words don't give much information.  
Hint - A preferable IDF range is 2-11 for model 2.
4. Remove the low idf value and high idf value words from the train and test data.
5. Perform tokenization on the modified text data same as you have done for previous
6. Create embedding matrix for model 2 and then use the rest of the features similar
7. Define the model, compile and fit the model.

```

1 from sklearn.feature_extraction.text import TfidfVectorizer
2
3 vectorizer = TfidfVectorizer()
4 X = vectorizer.fit(X_train['essay'])
5 idf = vectorizer.idf_
6 word_idf = dict(zip(vectorizer.get_feature_names(), idf))

```

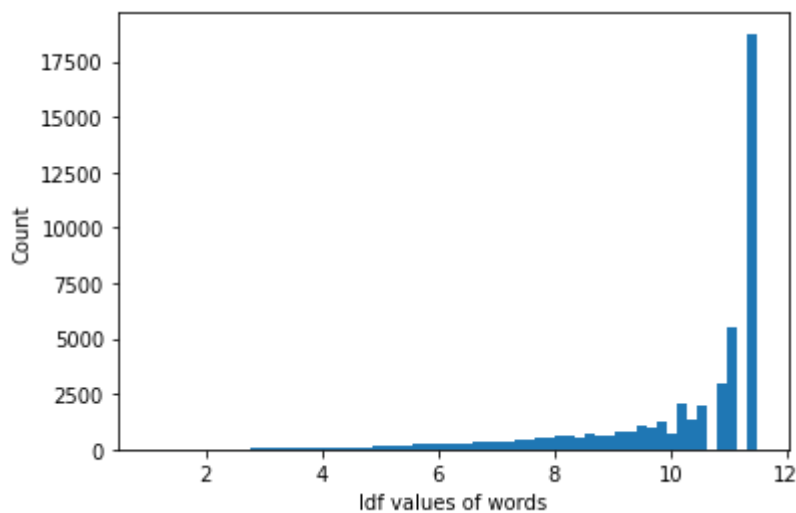
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: warnings.warn(msg, category=FutureWarning)

## ▼ plot histogram of idf values.

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 %matplotlib inline
4 import seaborn as sns
5
6 np.random.seed(42)
7
8 plt.hist(idf, density=False, bins=60) # density=False would make counts
9 plt.ylabel('Count')
10 plt.xlabel('Idf values of words');
11

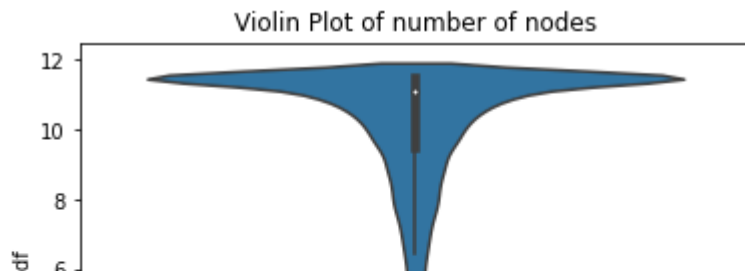
```



```

1 sns.violinplot(y='idf', data=pd.DataFrame({'idf':idf}))
2 plt.title("Violin Plot of number of nodes")
3 plt.show()

```



From the above cell output we can observe that most of the idf values are present in between 2 and 12.

```
1 remove_words = []
2 low_idf_threshold = 4
3 high_idf_threshold = 11.5
4 print(len(word_idf))
5 for k,v in word_idf.items():
6     if (v>high_idf_threshold or v<low_idf_threshold):
7         remove_words.append(k)
8 remove_words = set(remove_words)
9 print("No of words to be removed from training and test data are :",len(remove
10 # print(remove_words)
```

48195

No of words to be removed from training and test data are : 19213

```
1 if 'classroom' in remove_words:
2     print("word 'classroom' is present in remove_words set")
```

word 'classroom' is present in remove\_words set

+ Code

+ Text

```
1 check_word = "classroom"
2 sample = X_train['essay'].iloc[0]
3 print("Index of word '{0}' is {1}".format(check_word,sample.find(check_word)))
```

Index of word 'classroom' is 153

```
1 func = lambda x: ' '.join([item for item in x.split() if item.lower() not in re
2 X_train['modified_essay'] = X_train['essay'].apply(func)
3 X_test['modified_essay'] = X_test['essay'].apply(func)
```

```
1 check_word = "classroom"
2 sample = X_train['modified_essay'].iloc[0]
3 if (sample.find(check_word) == -1):
4     print("word '{0}' is not found after modifying".format(check_word))
```

word 'classroom' is not found after modifying

```
1 print(X_train['essay'].iloc[1])
2 print(X_train['modified_essay'].iloc[1])
3 # checking whether words like 'we','service' exists in remove_words
4 print('we' in remove_words)
```

```
we service highly impoverished area 83 families participating free lunch too m
service highly impoverished 83 participating too deal homelessness hunger torr
True
```

## ▼ Text Vectorization

```
1 from tensorflow.keras.preprocessing.text import Tokenizer
2 from tensorflow.keras.preprocessing.sequence import pad_sequences
3
4 tokenizer = Tokenizer()
5 tokenizer.fit_on_texts(X_train['modified_essay'])
6 vocab_size = len(tokenizer.word_index) + 1
7 # integer encode the documents
8 train_encoded_docs = tokenizer.texts_to_sequences(X_train['modified_essay'])
9 test_encoded_docs = tokenizer.texts_to_sequences(X_test['modified_essay'])
```

```
1 max_length = max([len(seq) for seq in train_encoded_docs])
2 print(max_length)
3 train_padded_docs = pad_sequences(train_encoded_docs, maxlen=max_length, padding=
4 test_padded_docs = pad_sequences(test_encoded_docs, maxlen=max_length, padding=
5 print(train_padded_docs.shape)
6 print(test_padded_docs.shape)
7 print(train_padded_docs[0][-10:])
8 print(test_padded_docs[0][-10:])
```

```
192
(73196, 192)
(36052, 192)
[0 0 0 0 0 0 0 0 0 0]
[0 0 0 0 0 0 0 0 0 0]
```

```
1 #after getting the padded_docs you have to use predefined glove vectors to get
2 # we will be storing this data in form of an embedding matrix and will use it v
3 # Please go through following blog's 'Example of Using Pre-Trained GloVe Embed
4 # https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-k
5 import pickle
6 import numpy as np
7
8 # path = "/Users/yamasanimanoj-kumarreddy/Documents/AAIC/Named_Assignments/LSTM
9 # path = "/content/drive/MyDrive/Named_Assignments/LSTM_DonorsChoose/"
10
11 with open(path+'/glove_vectors', 'rb') as f:
12     model = pickle.load(f)
13     glove_words = set(model.keys())
```

```
1 embedding_matrix_model2 = np.zeros((vocab_size, 300))
2 for word, i in tokenizer.word_index.items():
3     embedding_vector = model.get(word)
4     if embedding_vector is not None:
```

```

5     embedding_matrix_model2[i] = embedding_vector
6 print(embedding_matrix_model2[1][:5])
7 print(embedding_matrix_model2.shape)

```

```

[-0.043504 -0.18484 -0.14613 -0.21751  0.2025 ]
(29019, 300)

```

## ▼ Defining model and compiling the model

```

1 from tensorflow.keras import callbacks
2 import os
3
4 # path = "/Users/yamasanimanoj-kumarreddy/Documents/AAIC/Named_Assignments/LSTM
5 # path = "/content/drive/MyDrive/Named_Assignments/LSTM_DonorsChoose"
6
7 graph_saving_dir = os.path.join(path, 'model_2', "Tensorboard_graphs")
8
9 roc_callback = RocCallback(training_data=( [train_padded_docs, train_school_state,
10                                           train_project_grade, train_clean_categories,
11                                           train_clean_subcategories, train_teacher_prefix,
12                                           numerical_data_train ],
13                                           Y_train),
14                               validation_data=( [test_padded_docs, test_school_state,
15                                                  test_project_grade, test_clean_categories,
16                                                  test_clean_subcategories, test_teacher_prefix,
17                                                  numerical_data_test ],
18                                                  Y_test))
19
20
21 tensorboard_callback = callbacks.TensorBoard(log_dir=graph_saving_dir, histogram
22
23 callbacksList = [roc_callback, tensorboard_callback]

```

```

1 from tensorflow.keras import initializers, optimizers, metrics, losses
2 from sklearn.utils.class_weight import compute_class_weight
3
4 # Types of initializers to try
5 # initializers.HeUniform()
6 # initializers.HeNormal()
7 # initializers.GlorotNormal()
8 # initializers.GlorotUniform()
9 # initializers.LecunNormal()
10 # initializers.LecunUniform()
11 # initializers.RandomNormal(mean=0., stddev=1.)
12 # initializers.RandomUniform(minval=0., maxval=1.)
13
14 # Types of Optimizers to try
15 # optimizers.Adadelta(learning_rate=0.001)
16 # optimizers.Adagrad(learning_rate=0.001)
17 # optimizers.Adamax(learning_rate=0.001)
18 # optimizers.RMSprop(learning_rate=0.001)
19 # tf.keras.optimizers.SGD(learning_rate=0.01, momentum=0.0)

```

```

20
21 opt = optimizers.Adam()
22 max_seq_length = train_padded_docs.shape[1]
23 vocab_size_text = vocab_size
24 categories_dim = 1
25 categorical_ftr_embedding_size = 20
26 initializer = initializers.HeNormal()
27 activation_function = 'relu'
28 dropout_rate = 0.7
29 numerical_ftr_dim = 2
30 class_weights = compute_class_weight(class_weight = "balanced", classes = np.uni
31 class_weights = dict(zip(np.unique(Y_train), class_weights))
32 print(class_weights)
33
34 model = create_model(max_seq_length, vocab_size_text, embedding_matrix_model2,
35                       categorical_ftr_embedding_size, initializer, numerical_ftr_dim,
36
37 batch_size = 1024
38 epochs = 10
39
40 model.compile(optimizer=opt,
41               metrics = [metrics.BinaryAccuracy()],
42               loss=losses.BinaryCrossentropy())
43
44 # fit the model
45 model.fit([train_padded_docs, train_school_state,
46           train_project_grade, train_clean_categories,
47           train_clean_subcategories, train_teacher_prefix,
48           numerical_data_train ],
49         Y_train, batch_size = batch_size,
50         epochs=epochs, callbacks=callbacksList,
51         class_weight = class_weights,
52         validation_data = ([ test_padded_docs, test_school_state,
53                             test_project_grade, test_clean_categories,
54                             test_clean_subcategories, test_teacher_prefix,
55                             numerical_data_test ],
56                             Y_test))
57 # evaluate the model
58 loss, accuracy = model.evaluate([ test_padded_docs, test_school_state,
59                                   test_project_grade, test_clean_categories,
60                                   test_clean_subcategories, test_teacher_prefix,
61                                   numerical_data_test ],
62                                   Y_test, batch_size = batch_size)
63 print('Test Accuracy: %f' % (accuracy*100))

{0: 3.3021745014887665, 1: 0.5892164281229372}
Epoch 1/10
 6/72 [=>.....] - ETA: 9s - loss: 0.8705 - binary_accu
auc_train: 0.6831 - auc_val: 0.6839
72/72 [=====] - 48s 646ms/step - loss: 0.7052 - binar
Epoch 2/10
auc_train: 0.7105 - auc_val: 0.7017
72/72 [=====] - 45s 626ms/step - loss: 0.6725 - binar
Epoch 3/10
auc_train: 0.7329 - auc_val: 0.7105
72/72 [=====] - 33s 457ms/step - loss: 0.6542 - binar

```

```

Epoch 4/10
auc_train: 0.7489 - auc_val: 0.7207
72/72 [=====] - 32s 444ms/step - loss: 0.6397 - binar
Epoch 5/10
auc_train: 0.762 - auc_val: 0.7219
72/72 [=====] - 32s 451ms/step - loss: 0.6267 - binar
Epoch 6/10
auc_train: 0.781 - auc_val: 0.7173
72/72 [=====] - 46s 640ms/step - loss: 0.6141 - binar
Epoch 7/10
auc_train: 0.7985 - auc_val: 0.7213
72/72 [=====] - 33s 455ms/step - loss: 0.6038 - binar
Epoch 8/10
auc_train: 0.8183 - auc_val: 0.7175
72/72 [=====] - 46s 642ms/step - loss: 0.5886 - binar
Epoch 9/10
auc_train: 0.8383 - auc_val: 0.706
72/72 [=====] - 46s 639ms/step - loss: 0.5734 - binar
Epoch 10/10
auc_train: 0.8545 - auc_val: 0.7088
72/72 [=====] - 37s 525ms/step - loss: 0.5573 - binar
36/36 [=====] - 2s 47ms/step - loss: 0.5535 - binary_
Test Accuracy: 71.729726

```

```

1 #After fit method accessing the auc scores of train and validation
2 auc_train = model.auc_train
3 auc_val = model.auc_val
4
5 writer=tf.summary.create_file_writer(graph_saving_dir)
6 for idx in range(len(auc_train)):
7     with writer.as_default(step=idx+1):
8         tf.summary.scalar('auc_train', auc_train[idx])
9         tf.summary.scalar('auc_val', auc_val[idx])
10 writer.flush()

```

```
1 %load_ext tensorboard
```

The tensorboard extension is already loaded. To reload it, use:

```
%reload_ext tensorboard
```

```

1 %tensorboard --logdir /content/drive/MyDrive/Named_Assignments/LSTM_DonorsChoose
2 # %tensorboard --logdir /Users/yamasanimanoj-kumarreddy/Documents/AAIC/Named_Assignments

```

Reusing TensorBoard on port 6006 (pid 810), started 0:17:29 ago. (Use '!kill 810' to kill it.)

TensorBoard

SCALARS

GRAPH

INACTIVE

- ☐ Show data download links
- ☐ Ignore outliers in chart scaling

Tooltip sorting method: default

Smoothing



0.6

Horizontal Axis

STEP

RELATIVE

WALL

Runs

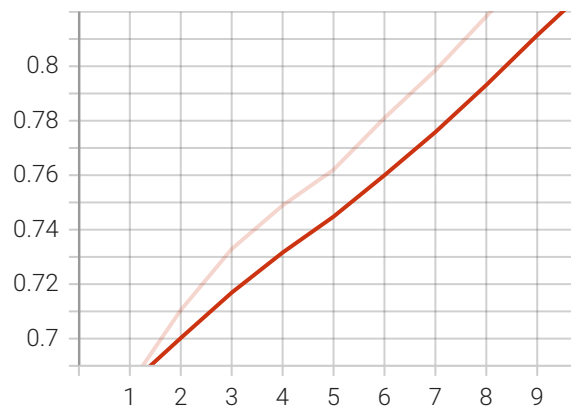
Write a regex to filter runs

☐ ☐ train☐ ☐ validation☐ ☐

Filter tags (regular expressions supported)

auc\_train

auc\_train  
tag: auc\_train



auc\_val

auc\_val  
tag: auc\_val

```
1 !tensorboard dev upload --logdir /content/drive/MyDrive/Named Assignments/LSTM
```

Upload started and will continue reading any new data as it's added to the log

To stop uploading, press Ctrl-C.

New experiment created. View your TensorBoard at: <https://tensorboard.dev/expe>

[2022-07-26T21:29:01] Started scanning logdir.

[2022-07-26T21:29:04] Total uploaded: 80 scalars, 190 tensors (136.4 kB), 1 bi

Interrupted. View your TensorBoard at <https://tensorboard.dev/experiment/4A5yC>

Traceback (most recent call last):

File "/usr/local/bin/tensorboard", line 8, in <module>

sys.exit(run\_main())

File "/usr/local/lib/python3.7/dist-packages/tensorboard/main.py", line 46,

app.run(tensorboard.main, flags\_parser=tensorboard.configure)

File "/usr/local/lib/python3.7/dist-packages/absl/app.py", line 308, in run

\_run\_main(main, args)

File "/usr/local/lib/python3.7/dist-packages/absl/app.py", line 254, in \_run

sys.exit(main(argv))



```
File "/usr/local/lib/python3.7/dist-packages/tensorboard/program.py", line 2
    return runner(self.flags) or 0
File "/usr/local/lib/python3.7/dist-packages/tensorboard/uploader/uploader_s
    return _run(flags, self._experiment_url_callback)
File "/usr/local/lib/python3.7/dist-packages/tensorboard/uploader/uploader_s
    intent.execute(server_info, channel)
KeyboardInterrupt
```

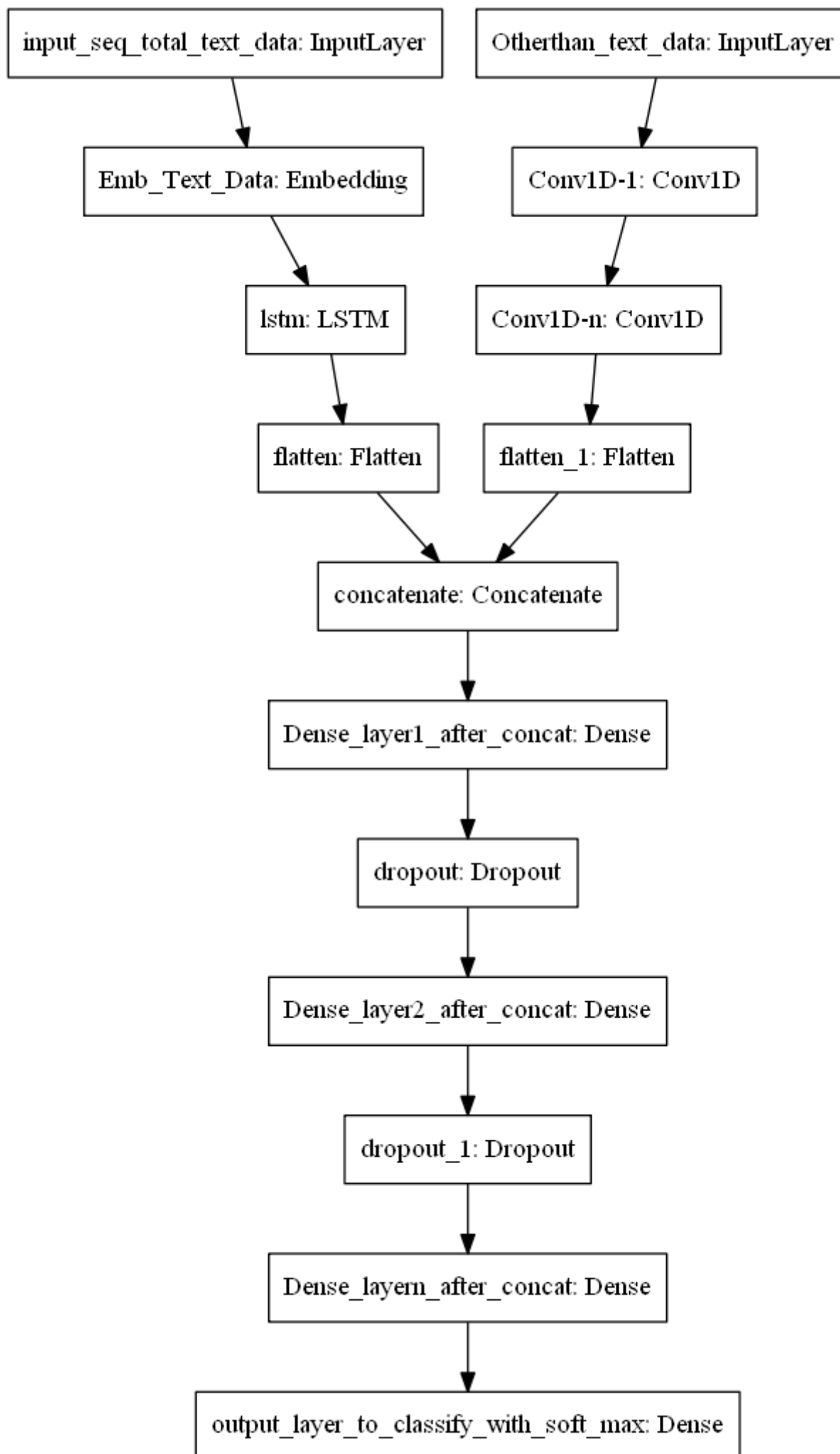
The above Tensorboard logs can be observed at this url

<https://tensorboard.dev/experiment/4A5yGvWDTLOqikhI3XZoHA/>

Observations from the above Training of model

1. In this model due to the removal of stop words, unnecessary information is removed from the text data and thereby our model gets trained faster and we achieve result at a less epoch number.
2. By using different initializers for kernel in Dense layers we could train different patterns of input data which might increase auc and accuracy.
3. At first we have tried without including `return_sequences = True` parameter in LSTM of text data. Due to this we get only 1 vector as output from LSTM cell and as a single vector cannot capture the entire essence of text sequence we got very less accuracy and auc values. But when we include the parameter `return_sequences = True` we get the output from every timestamp which increases the pattern information gained and thereby we got significant improvement in Test Accuracy and Auc.
4. From the training results we can observe that model gets overfitted in first five epochs only and thereby we achieved **validation auc = 0.70**

## ▼ Model-3



ref:

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

```
1 #in this model you can use the text vectorized data from model1
2 #for other than text data consider the following steps
3 # you have to perform one hot encoding of categorical features. You can use one
4 # Stack up standardised numerical features and all the one hot encoded categori
5 #the input to conv1d layer is 3d, you can convert your 2d data to 3d using np.r
6 # Note - deep learning models won't work with sparse features, you have to conv
```

## ▼ 1.1 Text Vectorization

```
1 from tensorflow.keras.preprocessing.text import Tokenizer
2 from tensorflow.keras.preprocessing.sequence import pad_sequences
3
4 tokenizer = Tokenizer()
5 tokenizer.fit_on_texts(X_train['essay'])
6 vocab_size = len(tokenizer.word_index) + 1
7 print(vocab_size)
8 # integer encode the documents
9 train_encoded_docs = tokenizer.texts_to_sequences(X_train['essay'])
10 test_encoded_docs = tokenizer.texts_to_sequences(X_test['essay'])
```

48232

```
1 max_length = max([len(seq) for seq in train_encoded_docs])
2 print(max_length)
3 train_padded_docs = pad_sequences(train_encoded_docs, maxlen=max_length, padding=
4 test_padded_docs = pad_sequences(test_encoded_docs, maxlen=max_length, padding=
5 print(train_padded_docs.shape)
6 print(test_padded_docs.shape)
7 print(train_padded_docs[0][-10:])
8 print(test_padded_docs[0][-10:])
```

```
339
(73196, 339)
(36052, 339)
[0 0 0 0 0 0 0 0 0 0]
[0 0 0 0 0 0 0 0 0 0]
```

```
1 #after getting the padded_docs you have to use predefined glove vectors to get
2 # we will be storing this data in form of an embedding matrix and will use it v
3 # Please go through following blog's 'Example of Using Pre-Trained GloVe Embed
4 # https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-ke
5 import pickle
```

```

6 import numpy as np
7
8 with open(path+'/glove_vectors', 'rb') as f:
9     model = pickle.load(f)
10    glove_words = set(model.keys())

```

```

1 embedding_matrix_model3 = np.zeros((vocab_size, 300))
2 for word, i in tokenizer.word_index.items():
3     embedding_vector = model.get(word)
4     if embedding_vector is not None:
5         embedding_matrix_model3[i] = embedding_vector
6 print(embedding_matrix_model3[1][:5])
7 print(embedding_matrix_model3.shape)
8 print("The number of tokens present in this vocabulary are", embedding_matrix_r

```

```

[ 0.15243 -0.16945 -0.022748 -0.25051 -0.15213 ]
(48232, 300)
The number of tokens present in this vocabulary are 48232

```

## ► 1.2 Categorical feature Vectorization

```
[ ] ↪ 11 cells hidden
```

## ► 1.3 Numerical feature Vectorization

```
[ ] ↪ 13 cells hidden
```

## ▼ Non-Text Data Concatenation

```

1 from sklearn.decomposition import PCA
2
3 nontext_data_train = np.concatenate((train_school_state, train_project_grade,
4                                     train_clean_categories, train_clean_subc
5                                     train_tchr_prj, train_price), axis=1)
6 nontext_data_test = np.concatenate((test_school_state, test_project_grade,
7                                     test_clean_categories, test_clean_subcate
8                                     test_tchr_prj, test_price), axis=1)
9
10

```

```

1 # Due to onehot encoding, the features are sparse. To avoid this we use pca to
2 # We are taking 30 components as they combinedly retain more than 80% of the in
3 reduced_dim = 30
4 pca = PCA(n_components = reduced_dim)
5 pca.fit(nontext_data_train)
6

```

```
7 nontext_data_train = pca.transform(nontext_data_train)
```

```
1 print(sum(pca.explained_variance_ratio_))
```

```
0.8052539834276029
```

```
1 nontext_data_train = np.expand_dims(nontext_data_train, axis=2)
```

```
2 nontext_data_test = np.expand_dims(nontext_data_test, axis=2)
```

```
3
```

```
4 print(nontext_data_train.shape)
```

```
5 print(nontext_data_test.shape)
```

```
(73196, 30, 1)
```

```
(36052, 30, 1)
```

## ▼ 1.4 Defining the model

```
1 from tensorflow.keras import layers, Input, Model
2
3 # Hyper Parameters
4 Embedding_vector_size_text = 300
5 categorical_ftr_embedding_size = 20
6 max_seq_length = train_padded_docs.shape[1]
7 no_of_lstm_units_on_text = 150
8 no_of_Dense_units_on_numerical_ftr = 150
9 no_of_units_on_dense_layer1 = 128
10 no_of_units_on_dense_layer2 = 64
11 no_of_units_on_dense_layer3 = 32
12 no_of_units_on_output_dense_layer = 1
13
14 def create_model(max_seq_length, vocab_size_text, embedding_matrix, nontextdata_dim,
15                  no_of_filters, kernel_size,
16                  initializer, activation_function, dropout_rate):
17     tf.keras.backend.clear_session()
18     # Text Feature block
19     input1 = Input(shape = (max_seq_length,), name = "Text_Input")
20     x = layers.Embedding(input_dim = vocab_size_text, output_dim=embedding_matrix.shape[1],
21                          input_length = max_seq_length, weights=[embedding_matrix],
22                          name = "Text_Embedding_layer", trainable = False)(input1)
23     x = layers.LSTM(no_of_lstm_units_on_text, return_sequences=True, name = "Text_LSTM")(x)
24     text_lstm_output = layers.Flatten(name = "Output_from_text_lstm")(x)
25
26     # Nontext data block
27     input2 = Input(shape = (nontextdata_dim,1), name = "nontextdata_Input")
28     x = layers.Conv1D(no_of_filters, kernel_size, activation=activation_function,
29                      name = "nontextdata_Convolution_layer1")(input2)
30     x = layers.Conv1D(no_of_filters, kernel_size, activation=activation_function,
31                      name = "nontextdata_Convolution_layer2")(x)
32     nontextdata_output = layers.Flatten(name = "nontextdata_output")(x)
33
34     # Concatenate two flatten inputs from text and nontext data
35     concat_output = layers.Concatenate(axis=1)([text_lstm_output, nontextdata_output])
```

```

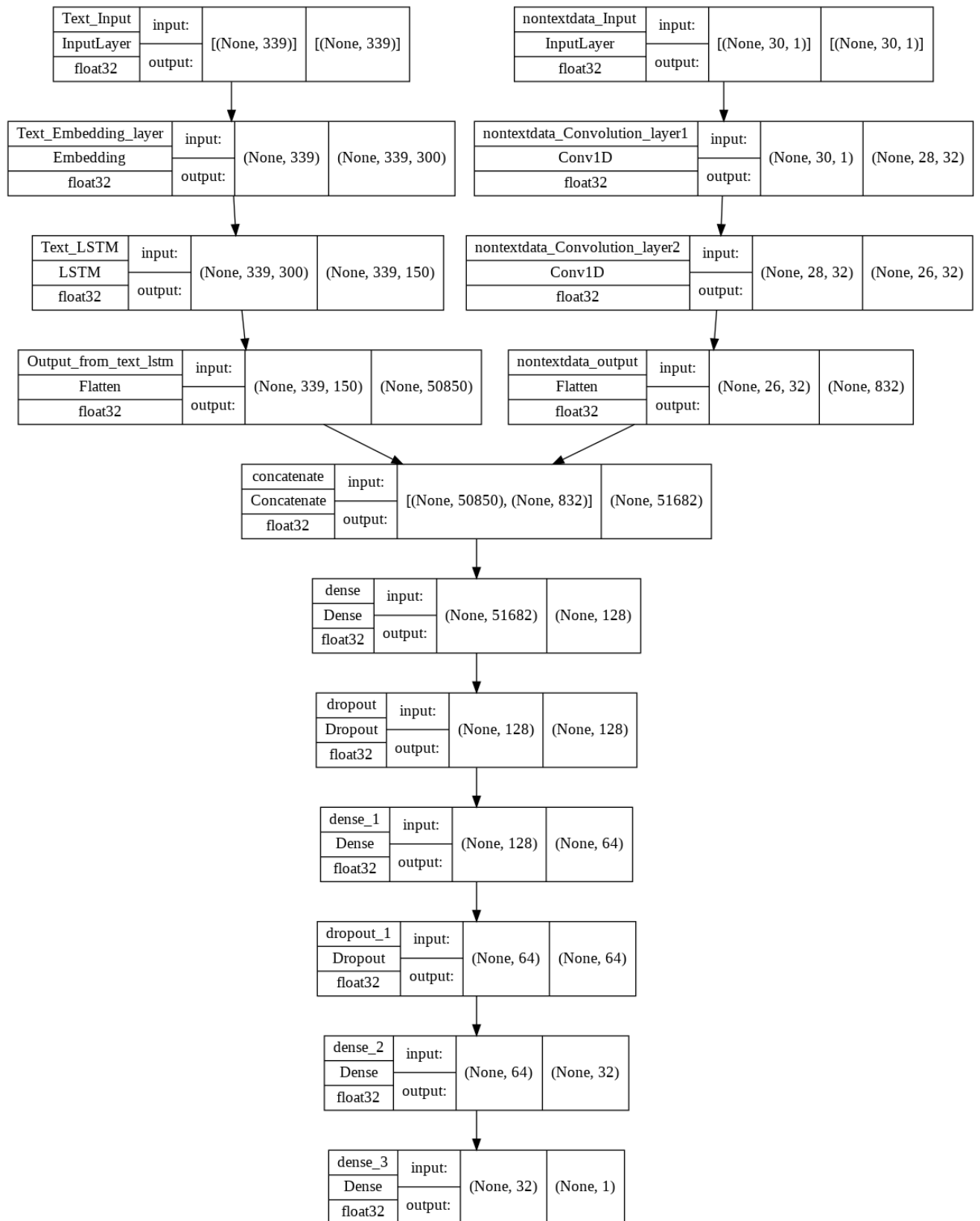
36
37 dense_layer1_after_concat = layers.Dense(no_of_units_on_dense_layer1,
38                                           activation = activation_function,
39                                           kernel_initializer=initializer)(concat)
40 dropout_layer1 = layers.Dropout(dropout_rate)(dense_layer1_after_concat)
41
42 dense_layer2_after_concat = layers.Dense(no_of_units_on_dense_layer2,
43                                           activation = activation_function,
44                                           kernel_initializer=initializer)(dropout_layer1)
45 dropout_layer2 = layers.Dropout(dropout_rate)(dense_layer2_after_concat)
46
47 dense_layer3_after_concat = layers.Dense(no_of_units_on_dense_layer3,
48                                           activation = activation_function,
49                                           kernel_initializer=initializer)(dropout_layer2)
50 output = layers.Dense(no_of_units_on_output_dense_layer,
51                       activation = 'sigmoid',
52                       kernel_initializer=initializer)(dense_layer3_after_concat)
53
54 return Model(inputs = [input1,input2],outputs = output, name = "Model_3")

```

```

1 from tensorflow.keras import utils,initializers
2
3 max_seq_length = train_padded_docs.shape[1]
4 vocab_size_text = vocab_size
5 nontextdata_dim = nontext_data_train.shape[1]
6 initializer = initializers.HeUniform()
7 activation_function = 'relu'
8 dropout_rate = 0.3
9 no_of_filters = 32
10 kernel_size = 3
11
12 model = create_model(max_seq_length, vocab_size_text, embedding_matrix_model3,
13                     no_of_filters, kernel_size,
14                     initializer, activation_function, dropout_rate)
15
16 utils.plot_model(model, "Model_3.png", show_shapes=True,show_dtype = True)

```



## ▼ 1.5 Compiling and fitting your model

```

1 from sklearn.metrics import roc_auc_score, roc_curve, auc
2 from keras.callbacks import Callback
3 class RocCallback(Callback):
4     def __init__(self, training_data, validation_data):
5         self.x = training_data[0]
6         self.y = training_data[1]
7         self.x_val = validation_data[0]
8         self.y_val = validation_data[1]
9
10
11     def on_train_begin(self, logs={}):
12         self.model.auc_train = []
13         self.model.auc_val = []
14         return
15
16     def on_train_end(self, logs={}):
17         return
18
19     def on_epoch_begin(self, epoch, logs={}):
20         return
21
22     def on_epoch_end(self, epoch, logs={}):
23         y_pred_train = self.model.predict(self.x)
24         auc_train = roc_auc_score(self.y, y_pred_train)
25         y_pred_val = self.model.predict(self.x_val)
26         auc_val = roc_auc_score(self.y_val, y_pred_val)
27         self.model.auc_train.append(auc_train)
28         self.model.auc_val.append(auc_val)
29         print('\rauc_train: %s - auc_val: %s' % (str(round(auc_train, 4)), str(round(auc_val, 4))))
30         return
31
32     def on_batch_begin(self, batch, logs={}):
33         return
34
35     def on_batch_end(self, batch, logs={}):
36         return

```

```

1 from tensorflow.keras import callbacks
2 import os
3
4 # path = "/Users/yamasanimanoj-kumarreddy/Documents/AAIC/Named_Assignments/LSTM
5 # path = "/content/drive/MyDrive/Named_Assignments/LSTM_DonorsChoose"
6
7 graph_saving_dir = os.path.join(path, 'model_3', "Tensorboard_graphs")
8
9 roc_callback = RocCallback(training_data=(train_padded_docs, nontext_data_train

```



```

10                                     Y_train),
11                                 validation_data=([test_padded_docs,nontext_data_test
12                                     Y_test]))
13
14 tensorboard_callback = callbacks.TensorBoard(log_dir=graph_saving_dir,histogram
15
16 callbacksList = [roc_callback, tensorboard_callback]

```

```

1 from tensorflow.keras import initializers,optimizers, metrics,losses
2 from sklearn.utils.class_weight import compute_class_weight
3
4 # Types of intializers to try
5 # initializers.HeUniform()
6 # initializers.HeNormal()
7 # initializers.GlorotNormal()
8 # initializers.GlorotUniform()
9 # initializers.LecunNormal()
10 # initializers.LecunUniform()
11 # initializers.RandomNormal(mean=0., stddev=1.)
12 # initializers.RandomUniform(minval=0., maxval=1.)
13
14 # Types of Optimizers to try
15 # optimizers.Adadelta(learning_rate=0.001)
16 # optimizers.Adagrad(learning_rate=0.001)
17 # optimizers.Adamax(learning_rate=0.001)
18 # optimizers.RMSprop(learning_rate=0.001)
19 # tf.keras.optimizers.SGD(learning_rate=0.01,momentum=0.0)
20
21 opt = optimizers.Adam()
22 max_seq_length = train_padded_docs.shape[1]
23 vocab_size_text = vocab_size
24 nontextdata_dim = nontext_data_train.shape[1]
25 initializer = initializers.HeUniform()
26 activation_function = 'relu'
27 dropout_rate = 0.7
28 no_of_filters = 64
29 kernel_size = 10
30 class_weights = compute_class_weight(class_weight = "balanced",classes = np.uni
31 class_weights = dict(zip(np.unique(Y_train), class_weights))
32 print(class_weights)
33
34 model = create_model(max_seq_length, vocab_size_text, embedding_matrix_model3,
35                       no_of_filters, kernel_size,
36                       initializer, activation_function, dropout_rate)
37
38 batch_size = 1024
39 epochs = 10
40
41 model.compile(optimizer=opt,
42               metrics = [metrics.BinaryAccuracy()],
43               loss=losses.BinaryCrossentropy())
44
45 # fit the model
46 model.fit([train_padded_docs,nontext_data_train],

```

```

47     Y_train, batch_size = batch_size,
48     epochs=epochs, callbacks=callbacksList,
49     class_weight = class_weights,
50     validation_data = ([test_padded_docs, nontext_data_test],
51                        Y_test))
52 # evaluate the model
53 loss, accuracy = model.evaluate([test_padded_docs, nontext_data_test],
54                                Y_test, batch_size = batch_size)
55 print('Test Accuracy: %f' % (accuracy*100))

{0: 3.3021745014887665, 1: 0.5892164281229372}
Epoch 1/10
 6/72 [=>.....] - ETA: 16s - loss: 1.0236 - binary_accu
auc_train: 0.6123 - auc_val: 0.6068
72/72 [=====] - 67s 826ms/step - loss: 0.7196 - binar
Epoch 2/10
auc_train: 0.6685 - auc_val: 0.6481
72/72 [=====] - 52s 733ms/step - loss: 0.6808 - binar
Epoch 3/10
auc_train: 0.7172 - auc_val: 0.6991
72/72 [=====] - 55s 772ms/step - loss: 0.6578 - binar
Epoch 4/10
auc_train: 0.7509 - auc_val: 0.7289
72/72 [=====] - 55s 772ms/step - loss: 0.6335 - binar
Epoch 5/10
auc_train: 0.772 - auc_val: 0.7333
72/72 [=====] - 57s 803ms/step - loss: 0.6181 - binar
Epoch 6/10
auc_train: 0.7927 - auc_val: 0.7365
72/72 [=====] - 53s 740ms/step - loss: 0.6020 - binar
Epoch 7/10
auc_train: 0.8153 - auc_val: 0.7303
72/72 [=====] - 57s 803ms/step - loss: 0.5803 - binar
Epoch 8/10
auc_train: 0.8389 - auc_val: 0.7334
72/72 [=====] - 55s 771ms/step - loss: 0.5608 - binar
Epoch 9/10
auc_train: 0.8456 - auc_val: 0.727
72/72 [=====] - 55s 768ms/step - loss: 0.5488 - binar
Epoch 10/10
auc_train: 0.8648 - auc_val: 0.7191
72/72 [=====] - 57s 799ms/step - loss: 0.5326 - binar
36/36 [=====] - 3s 82ms/step - loss: 0.4896 - binary_
Test Accuracy: 83.731830

```

```

1 #After fit method accessing the auc scores of train and validation
2 auc_train = model.auc_train
3 auc_val = model.auc_val
4
5 writer=tf.summary.create_file_writer(graph_saving_dir)
6 for idx in range(len(auc_train)):
7     with writer.as_default(step=idx+1):
8         tf.summary.scalar('auc_train', auc_train[idx])
9         tf.summary.scalar('auc_val', auc_val[idx])
10 writer.flush ()

```

```
1 %load_ext tensorboard
```

The tensorboard extension is already loaded. To reload it, use:

```
%reload_ext tensorboard
```

```
1 %tensorboard --logdir /content/drive/MyDrive/Named_Assignments/LSTM_DonorsChoose/
2 # %tensorboard --logdir /content/drive/MyDrive/DeepLearn/LSTM_DonorsChoose/model_3/
3 # %tensorboard --logdir /Users/yamasanimanoj-kumarreddy/Documents/AAIC/Named_Assignments/LSTM_DonorsChoose/
```

## TensorBoard

SCALARS

GRAPHS

INACTIVE

- ☐ Show data download links
- ☐ Ignore outliers in chart scaling

Tooltip sorting method: default

Smoothing



0.6

Horizontal Axis

STEP

RELATIVE

WALL

Runs

Write a regex to filter runs

☐ ○ train

☐ ○ validation

☐ ○ .

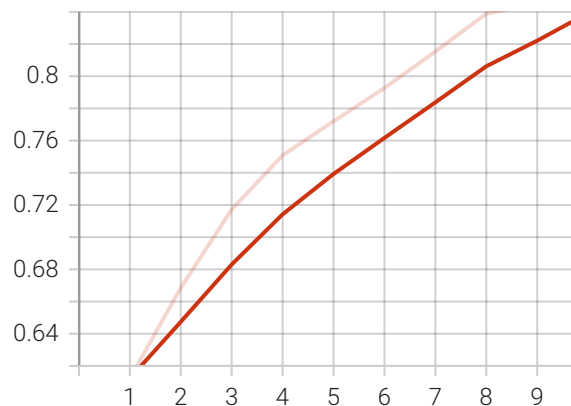
TOGGLE ALL RUNS

/content/drive/MyDrive/Named\_Assignments/LSTM\_DonorsChoose/  
model\_3/Tensorboard\_graphs

Filter tags (regular expressions supported)

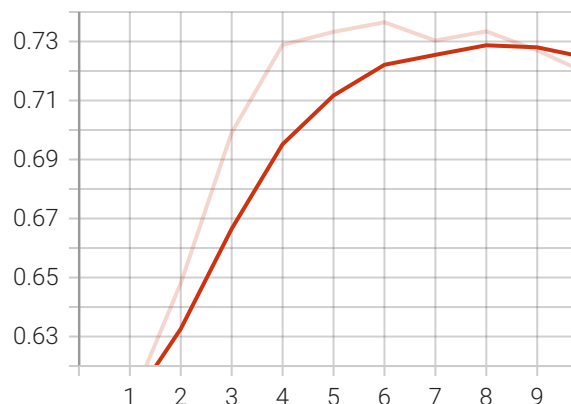
auc\_train

auc\_train  
tag: auc\_train



auc\_val

auc\_val  
tag: auc\_val



```
1 !tensorboard dev upload --logdir /content/drive/MyDrive/Named_Assignments/LSTM_
```

```
2
3 # !tensorboard dev upload --logdir %tensorboard --logdir /content/drive/MyDrive
4 # !tensorboard dev upload --logdir /Users/yamasanimanoj-kumarreddy/Documents/AI
```

The above Tensorboard logs can be observed at this url

<https://tensorboard.dev/experiment/g1T08oQxRESPPyEW0qvBnw/>

Observations from the above Training of model

1. If we increase the no of components in PCA there is more information included in the components and thus we get more accuracy and auc.
2. If we increase the no of filters and kernel size the information gained from the words increases upto a certain point and later decreases due to addition of unnecessary data.
3. By using different intializers for kernel in Convolution layer we could train different patterns of input data which might increase auc and accuracy.
4. At first i have tried without including `return_sequences = True` paramter in LSTM of text data. Due to this we get only 1 vector as output from LSTM cell and as a single vector cannot capture the entire esence of text sequence we got very less accuracy and auc values. But when we include the parameter `return_sequences = True` we get the output from every timestamp which increases the pattern information gained and thereby we got significant improvement in Test Accuracy and Auc.
5. From the training results we can observe that model is not overfitted and achieved

**validation auc = 0.7323**