

```
In [0]: 1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4
5 import pandas as pd
6 import numpy as np
7 from sklearn.feature_extraction.text import TfidfVectorizer
8 from sklearn.feature_extraction.text import CountVectorizer
9 import re
10 from nltk.corpus import stopwords
11 import pickle
12 from tqdm import tqdm
13 import os
```

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

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly (http s://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly)

Enter your authorization code:

.....

Mounted at /content/drive

```
In [0]: 1 #we already have data available for all the initials we need fot project title
```

```
In [0]: 1 #Loading the data
2 project_data = pd.read_csv('train_data.csv')
3 resource_data = pd.read_csv('resources.csv')
```

```
In [0]: 1 #merging both the data to get a final dataset
2 price_data = resource_data.groupby('id').agg({'price':'sum','quantity':'sum'})
3 price_data.head(2)
```

```
Out[8]:
```

| | id | price | quantity |
|---|---------|--------|----------|
| 0 | p000001 | 459.56 | 7 |
| 1 | p000002 | 515.89 | 21 |

```
In [0]: 1 project_data = pd.merge(project_data,price_data,on = 'id',how = 'left')#getti
```

```
In [0]: 1 #print(project_data.columns)
        2 project_data.head(2)
```

```
Out[10]: Unnamed: 0      id      teacher_id  teacher_prefix  school_state  project_sul
```

| | | | | | | |
|---|--------|---------|----------------------------------|------|----|----|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | 20 |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 20 |

```
In [0]: 1 target = project_data['project_is_approved']
        2 project_data.drop(['Unnamed: 0','id','teacher_id','project_is_approved'],axis=
        3
```

```
In [0]: 1 print('shape of the final data is:',project_data.shape)
        2 print('Shape of target variable is:',target.shape)
        3 print('Features of the final dataset is:',project_data.columns)
```

shape of the final data is: (109248, 15)

Shape of target variable is: (109248,)

Features of the final dataset is: Index(['teacher_prefix', 'school_state', 'project_submitted_datetime', 'project_grade_category', 'project_subject_categories', 'project_subject_subcategories', 'project_title', 'project_essay_1', 'project_essay_2', 'project_essay_3', 'project_essay_4', 'project_resource_summary', 'teacher_number_of_previously_posted_projects', 'price', 'quantity'], dtype='object')

so we have following final features

Categorical Features:

- Teacher_prefix
- school_stae
- project_grade_category
- project_subject_categories
- project_subject_subcategories

Text Features:

- project essay 1
- project essay 2
- project essay 3
- project essay 4
- project resource summary
- project title

Numerical Feature:

- teacher number of previously posted categories
- quantity
- price

Preprocessing the features

1.Preprocessing the categorical features

project_grade_category

```
In [0]: 1 # https://stackoverflow.com/questions/36383821/pandas-dataframe-apply-function
2 project_data['project_grade_category'] = project_data['project_grade_category']
3 project_data['project_grade_category'] = project_data['project_grade_category']
4 project_data['project_grade_category'] = project_data['project_grade_category']
5 project_data['project_grade_category'].value_counts()
```

```
Out[25]: grades_prek_2    44225
grades_3_5      37137
grades_6_8      16923
grades_9_12     10963
Name: project_grade_category, dtype: int64
```

project_subject_category

```
In [0]: 1 project_data['project_subject_categories'] = project_data['project_subject_ca
2 project_data['project_subject_categories'] = project_data['project_subject_ca
3 project_data['project_subject_categories'] = project_data['project_subject_ca
4 project_data['project_subject_categories'] = project_data['project_subject_ca
5 project_data['project_subject_categories'] = project_data['project_subject_ca
6 project_data['project_subject_categories'].value_counts()
```

```
Out[26]: literacy_language      23655
math_science      17072
literacy_language_math_science  14636
health_sports      10177
music_arts         5180
specialneeds       4226
literacy_language_specialneeds  3961
appliedlearning    3771
math_science_literacy_language  2289
appliedlearning_literacy_language  2191
history_civics     1851
math_science_specialneeds      1840
literacy_language_music_arts    1757
math_science_music_arts        1642
appliedlearning_specialneeds    1467
history_civics_literacy_language  1421
health_sports_specialneeds      1391
warmth_care_hunger      1309
math_science_appliedlearning    1220
literacy_language_appliedlearning  1050
```

Teacher prefix

```
In [0]: 1 # check if we have any nan values are there
2 print(project_data['teacher_prefix'].isnull().values.any())
3 print("number of nan values",project_data['teacher_prefix'].isnull().values.s
```

```
True
number of nan values 3
```

numebr of missing values are very less in number, we can replace it with Mrs. as most of the projects are submitted by Mrs.

```
In [0]: 1 project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('Mrs.')
2 project_data['teacher_prefix'] = project_data['teacher_prefix'].str.replace('
3 project_data['teacher_prefix'] = project_data['teacher_prefix'].str.lower()
4 project_data['teacher_prefix'].value_counts()
```

```
Out[28]: mrs      57272
ms      38955
mr      10648
teacher  2360
dr       13
Name: teacher_prefix, dtype: int64
```

Project subject subcategories

```
In [0]: 1 project_data['project_subject_subcategories'] = project_data['project_subject_subcategories'].str.lower()
2 project_data['project_subject_subcategories'] = project_data['project_subject_subcategories'].str.lower()
3 project_data['project_subject_subcategories'] = project_data['project_subject_subcategories'].str.lower()
4 project_data['project_subject_subcategories'] = project_data['project_subject_subcategories'].str.lower()
5 project_data['project_subject_subcategories'] = project_data['project_subject_subcategories'].str.lower()
6 project_data['project_subject_subcategories'].value_counts()
```

```
Out[29]: literacy          9486
literacy_mathematics      8325
literature_writing_mathematics 5923
literacy_literature_writing 5571
mathematics              5379
...
gym_fitness_warmth_care_hunger 1
civics_government_foreignlanguages 1
economics_nutritioneducation 1
financialliteracy_performingarts 1
extracurricular_financialliteracy 1
Name: project_subject_subcategories, Length: 401, dtype: int64
```

School state

```
In [0]: 1 project_data['school_state'] = project_data['school_state'].str.lower()
2 project_data['school_state'].value_counts()
```

```
Out[30]: ca      15388
tx       7396
ny       7318
fl       6185
nc       5091
il       4350
ga       3963
sc       3936
mi       3161
pa       3109
in       2620
mo       2576
oh       2467
la       2394
ma       2389
wa       2334
ok       2276
nj       2237
az       2147
...       ...
```

2.Preprocessing text features

In [0]:

```

1  #writing the utility functions for text cleaning
2  from tqdm import tqdm
3  # https://gist.github.com/sebleier/554280
4  # we are removing the words from the stop words List: 'no', 'nor', 'not'
5  stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you
6              "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he
7              'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'i
8              'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', '
9              'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have',
10             'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'beca
11             'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',
12             'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on',
13             'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', '
14             'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'th
15             's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "shoul
16             've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn',
17             "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", '
18             "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shou
19             'won', "won't", 'wouldn', "wouldn't"]
20
21  #=====
22  # https://stackoverflow.com/a/47091490/4084039
23  import re
24
25  def decontracted(phrase):
26      """Function for opening thee decontracted words"""
27      # specific
28      phrase = re.sub(r"won't", "will not", phrase)#converting to will not
29      phrase = re.sub(r"can't", "can not", phrase)#converting to can not
30
31      # general
32      phrase = re.sub(r"n't", " not", phrase)#converting to not
33      phrase = re.sub(r"\re", " are", phrase)#converting to are
34      phrase = re.sub(r"\s", " is", phrase)#converting to is
35      phrase = re.sub(r"\d", " would", phrase)#converting to would
36      phrase = re.sub(r"\ll", " will", phrase)#converting to will
37      phrase = re.sub(r"\t", " not", phrase)#converting to not
38      phrase = re.sub(r"\ve", " have", phrase)#converting to have
39      phrase = re.sub(r"\m", " am", phrase)#converting to am
40      return phrase
41  #=====
42
43  # Combining all the above stundents
44  def preprocess_text(text_data):
45      """cleaning and remvng the stopwords"""
46      preprocessed_text = []
47      # tqdm is for printing the status bar
48      for sentence in tqdm(text_data):
49          sent = decontracted(sentence)
50          sent = sent.replace('\r', ' ')
51          sent = sent.replace('\n', ' ')
52          sent = sent.replace('\\"', ' ')
53          sent = re.sub('[A-Za-z0-9]+', ' ', sent)
54          # https://gist.github.com/sebleier/554280
55          sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
56          preprocessed_text.append(sent.lower().strip())

```



```
In [0]: 1 from sklearn.preprocessing import MinMaxScaler
2
3 scaler = MinMaxScaler()
4 scaler.fit(project_data['price'].values.reshape(-1, 1))
5 project_data['nrm_price'] = scaler.transform(project_data['price'].values.reshape(-1, 1))
6 project_data
7 project_data.head(2)
```

```
Out[41]:
```

| | teacher_prefix | school_state | project_submitted_datetime | project_grade_category | project_subject |
|---|----------------|--------------|----------------------------|------------------------|------------------|
| 0 | mrs | in | 2016-12-05 13:43:57 | grades_prek_2 | literac |
| 1 | mr | fl | 2016-10-25 09:22:10 | grades_6_8 | history_civics_h |

```
In [0]: 1 #Save the processed dataset into a pandas CSV file.
2 project_data.to_csv("processed_data.csv", index=False)
```

Preparing the data for modelling

- Vectorizing categorical features using one hot encoding
- Vectorizing text data using BOW or TFIDF/W2V TFIDF
- Standardizing or normalizing numerical data

```
In [0]: 1 #splitting the data
2 project_data = pd.read_csv('processed_data.csv')
3 project_data.head(2)
```

```
Out[13]:
```

| | teacher_prefix | school_state | project_submitted_datetime | project_grade_category | project_subject |
|---|----------------|--------------|----------------------------|------------------------|------------------|
| 0 | mrs | in | 2016-12-05 13:43:57 | grades_prek_2 | literac |
| 1 | mr | fl | 2016-10-25 09:22:10 | grades_6_8 | history_civics_h |


```
In [0]: 1 project_data.columns
```

```
Out[14]: Index(['teacher_prefix', 'school_state', 'project_submitted_datetime',  
               'project_grade_category', 'project_subject_categories',  
               'project_subject_subcategories',  
               'teacher_number_of_previously_posted_projects', 'price', 'quantity',  
               'cleaned_project_titles', 'cleaned_resource_summary', 'cleaned_essay',  
               'nrm_price'],  
              dtype='object')
```

```
In [0]: 1 #splitting the data in train,validation and test  
2 from sklearn.model_selection import train_test_split  
3 #80% training data and 20% test data  
4 X_train,X_test,Y_train,Y_test = train_test_split(project_data,target,test_size=0.2,  
5 X_train,X_cv,Y_train,Y_cv = train_test_split(X_train,Y_train,test_size = 0.2,  
6  
7 print('shape of training data after vectorization is:',X_train.shape)  
8 print('shape of validation data after vectorization is:',X_cv.shape)  
9 print('shape of test data after vectorization is:',X_test.shape)
```

```
In [0]: 1 print('shape of training data after vectorization is:',X_train.shape)  
2 print('shape of validation data after vectorization is:',X_cv.shape)  
3 print('shape of test data after vectorization is:',X_test.shape)
```

```
shape of training data after vectorization is: (69918, 13)  
shape of validation data after vectorization is: (17480, 13)  
shape of test data after vectorization is: (21850, 13)
```

```
In [3]: 1 #trying the embeddings
2 from keras.models import Sequential
3 from keras.preprocessing import sequence
4 from keras.initializers import he_normal
5 from keras.layers import Input
6 from keras.layers import Dense
7 from keras.layers.embeddings import Embedding
8 from keras.regularizers import L1L2
9 from keras.layers import BatchNormalization
10 from keras.layers import Dropout
11 from keras.layers import LSTM
12 from keras.layers import Flatten
13 from keras.preprocessing.text import Tokenizer
14 from keras.preprocessing.sequence import pad_sequences
15 from keras.models import Model
16 from keras.layers import concatenate
17 from sklearn.metrics import accuracy_score
18 from IPython.core.display import display,HTML
19 from numpy import zeros
20 display(HTML("<style>.container { width:100% !important; }</style>"))
```

Using TensorFlow backend.

The default version of TensorFlow in Colab will soon switch to TensorFlow 2.x.

We recommend you [upgrade](https://www.tensorflow.org/guide/migrate) (https://www.tensorflow.org/guide/migrate) now or ensure your notebook will continue to use TensorFlow 1.x via the %tensorflow_version 1.x magic: [more info](https://colab.research.google.com/notebooks/tensorflow_version.ipynb) (https://colab.research.google.com/notebooks/tensorflow_version.ipynb).

2. Vectorizing the categorical features

2.1 Project_category

```
In [0]: 1 from collections import Counter
2 def count(df):
3
4     """Returns the sorted dictionary keys"""
5     my_counter = Counter() #initiating counter to calculate occurence of words
6     for word in df.values:
7         my_counter.update(word.split())
8
9     # dict sort by value python: https://stackoverflow.com/a/613218/4084039
10    cat_dict = dict(my_counter)
11    sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
12    return sorted_cat_dict
```

```

In [0]: 1
2 # we use count vectorizer to convert the values into one hot encoded features
3 from sklearn.feature_extraction.text import CountVectorizer
4 cat_dict = count(project_data['project_subject_categories'])
5 vectorizer = CountVectorizer(vocabulary=list(cat_dict.keys()), lowercase=False)
6 vectorizer.fit(X_train['project_subject_categories'].values)
7 print(vectorizer.get_feature_names()[:10])
8
9
10 train_categories = vectorizer.transform(X_train['project_subject_categories'])
11 val_categories = vectorizer.transform(X_cv['project_subject_categories']).values
12 test_categories = vectorizer.transform(X_test['project_subject_categories']).values
13 print("Shape of train matrix after one hot encoding ", train_categories.shape)
14 print("Shape of val matrix after one hot encoding ", val_categories.shape)
15 print("Shape of test matrix after one hot encoding ", test_categories.shape)
16
17 #vocab_size_categories = train_categories.shape[1] + 1
18

```

```

['history_civics_warmth_care_hunger', 'music_arts_warmth_care_hunger', 'literacy_language_warmth_care_hunger', 'music_arts_appliedlearning', 'appliedlearning_warmth_care_hunger', 'math_science_warmth_care_hunger', 'history_civics_health_sports', 'music_arts_history_civics', 'music_arts_health_sports', 'health_sports_warmth_care_hunger']

```

Shape of train matrix after one hot encoding (69918, 51)

Shape of val matrix after one hot encoding (17480, 51)

Shape of test matrix after one hot encoding (21850, 51)

```

In [0]: 1 #getting the embedding layer
2 #Get the flattened output for clean_categories
3 input_layer_clean_categories = Input(shape=(train_categories.shape[1]), name=
4 embedding_layer_clean_categories = Embedding(input_dim=train_categories.shape
5 flatten_clean_categories = Flatten()(embedding_layer_clean_categories)

```

2.2 Project subcategories

```
In [0]: 1 # we use count vectorizer to convert the values into one hot encoded features
2 #from sklearn.feature_extraction.text import CountVectorizer
3 cat_dict = count(project_data['project_subject_subcategories'])
4 vectorizer = CountVectorizer(vocabulary=list(cat_dict.keys()), lowercase=False)
5 vectorizer.fit(X_train['project_subject_subcategories'].values)
6 print(vectorizer.get_feature_names()[:10])
7
8
9 train_subcategories = vectorizer.transform(X_train['project_subject_subcategories'].values)
10 val_subcategories = vectorizer.transform(X_cv['project_subject_subcategories'].values)
11 test_subcategories = vectorizer.transform(X_test['project_subject_subcategories'].values)
12 print("Shape of train matrix after one hot encoding ",train_subcategories.shape)
13 print("Shape of val matrix after one hot encoding ",val_subcategories.shape)
14 print("Shape of test matrix after one hot encoding ",test_subcategories.shape)
15
```

```
['gym_fitness_parentinvolvement', 'communityservice_gym_fitness', 'college_care
erprep_warmth_care_hunger', 'economics_other', 'parentinvolvement_warmth_care_h
unger', 'communityservice_music', 'gym_fitness_socialsciences', 'parentinvolvem
ent_teamsports', 'civics_government_nutritioneducation', 'financialliteracy_for
eignlanguages']
```

Shape of train matrix after one hot encoding (69918, 401)

Shape of val matrix after one hot encoding (17480, 401)

Shape of test matrix after one hot encoding (21850, 401)

```
In [0]: 1 #getting the embedding layer
2 #Get the flattened output for clean_subcategories
3 input_layer_subcategories = Input(shape=(train_subcategories.shape[1],), name='input_layer_subcategories')
4 embedding_layer_subcategories = Embedding(input_dim=train_subcategories.shape[1], output_dim=401, name='embedding_layer_subcategories')
5 flatten_subcategories = Flatten()(embedding_layer_subcategories)
```

2.3 Teacher Prefix

```
In [0]: 1 cat_dict = count(project_data['teacher_prefix'])
2 vectorizer = CountVectorizer(vocabulary=list(cat_dict.keys()), lowercase=False)
3 vectorizer.fit(X_train['teacher_prefix'].values)
4 print(vectorizer.get_feature_names()[:10])
5
6
7 train_prefix = vectorizer.transform(X_train['teacher_prefix'].values)
8 val_prefix = vectorizer.transform(X_cv['teacher_prefix'].values)
9 test_prefix= vectorizer.transform(X_test['teacher_prefix'].values)
10 print("Shape of train matrix after one hot encoding ",train_prefix.shape)
11 print("Shape of val matrix after one hot encoding ",val_prefix.shape)
12 print("Shape of test matrix after one hot encoding ",test_prefix.shape)
13
```

```
['dr', 'teacher', 'mr', 'ms', 'mrs']
```

Shape of train matrix after one hot encoding (69918, 5)

Shape of val matrix after one hot encoding (17480, 5)

Shape of test matrix after one hot encoding (21850, 5)

```
In [0]: 1 #getting the embedding layer
2 #Get the flattened output for teacher_prefix
3 input_layer_prefix = Input(shape=(train_prefix.shape[1],), name = "input_pref
4 embedding_layer_prefix = Embedding(input_dim=train_prefix.shape[1], output_di
5 flatten_prefix = Flatten()(embedding_layer_prefix)
```

2.4 School state

```
In [0]: 1 cat_dict = count(project_data['school_state'])
2 vectorizer = CountVectorizer(vocabulary=list(cat_dict.keys()), lowercase=False
3 vectorizer.fit(X_train['school_state'].values)
4 print(vectorizer.get_feature_names())
5
6
7 train_state = vectorizer.transform(X_train['school_state'].values)
8 val_state = vectorizer.transform(X_cv['school_state'].values)
9 test_state = vectorizer.transform(X_test['school_state'].values)
10 print("Shape of train matrix after one hot encodig ", train_state.shape)
11 print("Shape of val matrix after one hot encodig ", val_state.shape)
12 print("Shape of test matrix after one hot encodig ", test_state.shape)
```

```
['vt', 'wy', 'nd', 'mt', 'ri', 'sd', 'ne', 'de', 'ak', 'nh', 'wv', 'me', 'hi',
'dc', 'nm', 'ks', 'ia', 'id', 'ar', 'co', 'mn', 'or', 'ky', 'ms', 'nv', 'md',
'ct', 'tn', 'ut', 'al', 'wi', 'va', 'az', 'nj', 'ok', 'wa', 'ma', 'la', 'oh',
'mo', 'in', 'pa', 'mi', 'sc', 'ga', 'il', 'nc', 'fl', 'ny', 'tx', 'ca']
```

Shape of train matrix after one hot encodig (69918, 51)

Shape of val matrix after one hot encodig (17480, 51)

Shape of test matrix after one hot encodig (21850, 51)

```
In [0]: 1 #getting the embedding layer
2 #Get the flattened output for clean_subcategories
3 input_layer_state = Input(shape=(train_state.shape[1],), name = "input_state"
4 embedding_layer_state = Embedding(input_dim=train_state.shape[1], output_dim=
5 flatten_state = Flatten()(embedding_layer_state)
```

2.5 Project grade category

```
In [0]: 1 cat_dict = count(project_data['project_grade_category'])
2 vectorizer = CountVectorizer(vocabulary=list(cat_dict.keys()), lowercase=False)
3 vectorizer.fit(X_train['project_grade_category'].values)
4 print(vectorizer.get_feature_names()[:10])
5
6
7 train_grade_category = vectorizer.transform(X_train['project_grade_category'])
8 val_grade_category = vectorizer.transform(X_cv['project_grade_category']).values
9 test_grade_category = vectorizer.transform(X_test['project_grade_category']).values
10 print("Shape of train matrix after one hot encoding ", train_grade_category.shape)
11 print("Shape of val matrix after one hot encoding ", val_grade_category.shape)
12 print("Shape of test matrix after one hot encoding ", test_grade_category.shape)
```

```
['grades_9_12', 'grades_6_8', 'grades_3_5', 'grades_prek_2']
```

```
Shape of train matrix after one hot encoding (69918, 4)
```

```
Shape of val matrix after one hot encoding (17480, 4)
```

```
Shape of test matrix after one hot encoding (21850, 4)
```

```
In [0]: 1 #getting the embedding layer
2 #Get the flattened output for project grade category
3 input_layer_grade_category = Input(shape=(train_grade_category.shape[1],), name='input_layer_grade_category')
4 embedding_layer_grade_category = Embedding(input_dim=train_grade_category.shape[1], output_dim=EMBEDDING_DIM, name='embedding_layer_grade_category')(input_layer_grade_category)
5 flatten_grade_category = Flatten()(embedding_layer_grade_category)
```

Vectorizing the text data

as we want to input the total text that's why we must get the final text by combining all

```
In [0]: 1 #utility function
2 def final_text(df):
3     """Returns the final concatenated text"""
4     df['final_Text'] = df['cleaned_essay'].map(str) + " " + df['cleaned_resources'].map(str)
5     " " + df['cleaned_project_titles'].map(str)
6     return df
7
8 #getting the train, val and test data
9 train_text = final_text(X_train)
10 val_text = final_text(X_cv)
11 test_text = final_text(X_test)
```

```

In [0]: 1 #building the data for modelling
        2
        3
        4 #initializing the tokenizer
        5 tokens = Tokenizer()
        6 tokens.fit_on_texts(list(train_text['final_Text'].values))#will convert texts
        7
        8 #integer encoding the documents
        9 docs_train = tokens.texts_to_sequences(list(train_text['final_Text'].values))
       10 docs_val = tokens.texts_to_sequences(list(val_text['final_Text'].values))
       11 docs_test = tokens.texts_to_sequences(list(test_text['final_Text'].values))
       12
       13 vocab_size = len(tokens.word_index) + 1
       14
       15 max_length = 300#for padding the sequences
       16
       17 #finally padding the squences
       18 padded_docs_train = pad_sequences(docs_train, maxlen=max_length, padding='post')
       19 padded_docs_val = pad_sequences(docs_val, maxlen=max_length, padding='post')
       20 padded_docs_test = pad_sequences(docs_test, maxlen=max_length, padding='post')
       21
       22 #creating a matrix for each word in training data
       23 pickle_in = open("glove_vectors", "rb")
       24 glove_words = pickle.load(pickle_in)
       25 embedding_matrix = zeros((vocab_size, 300))
       26 for word, i in tokens.word_index.items(): #enumerating all unique words
       27     embedding_vector = glove_words.get(word)#Locating for loaded glove model
       28     if embedding_vector is not None:
       29         embedding_matrix[i] = embedding_vector

```

```

In [0]: 1 vocab_size

```

Out[123]: 51153

```

In [0]: 1 #as we choose the vector to be 300 dimensional
        2 #that's why the output dimension must be 300
        3 input_text_layer = Input(shape = (300,), name = 'input_total_text_sequence')
        4 embedding_layer_text = Embedding(input_dim = vocab_size, output_dim = 300, weights=
        5 #here input_dim is size of vocabulary of text data
        6 #output_dim is size of embedded vector

```

```

In [0]: 1 #lstm for this layer
        2 lstm_layer = LSTM(32, activation = 'relu', return_sequences = True)(embedding_1
        3 #and finally flattening the output
        4 flatten_text = Flatten()(lstm_layer)

```

Normalizing the numerical data

Teacher number of previously posted projects

```
In [0]: 1 from sklearn.preprocessing import Normalizer
2
3 norm = Normalizer()
4 train_ppp = norm.fit_transform(X_train['teacher_number_of_previously_posted_p
5 val_ppp = norm.transform(X_cv['teacher_number_of_previously_posted_projects'])
6 test_ppp = norm.transform(X_test['teacher_number_of_previously_posted_project
```

```
In [0]: 1 train_ppp = train_ppp.reshape(len(X_train),1)
2 val_ppp = val_ppp.reshape(len(X_cv),1)
3 test_ppp = test_ppp.reshape(len(X_test),1)
```

```
In [0]: 1 input_num_1 = Input(shape = (train_ppp.shape[1],),name = 'input_previously_po
```

Quantity

```
In [0]: 1 norm = Normalizer()
2 train_quantity = norm.fit_transform(X_train.quantity.values.reshape(-1,1))
3 val_quantity = norm.transform(X_cv.quantity.values.reshape(-1,1))
4 test_quantity = norm.transform(X_test.quantity.values.reshape(-1,1))
5
6 train_quantity = train_quantity.reshape(len(X_train),1)
7 val_quantity = val_quantity.reshape(len(X_cv),1)
8 test_quantity = test_quantity.reshape(len(X_test),1)
9
```

```
In [0]: 1 input_num_2 = Input(shape = (train_quantity.shape[1],),name = 'input_quantity
```

Price

```
In [0]: 1 norm = Normalizer()
2 train_price = norm.fit_transform(X_train.price.values.reshape(-1,1))
3 val_price = norm.fit_transform(X_cv.price.values.reshape(-1,1))
4 test_price = norm.fit_transform(X_test.price.values.reshape(-1,1))
5
6 train_price = train_price.reshape(len(X_train),1)
7 val_price = val_price.reshape(len(X_cv),1)
8 test_price = test_price.reshape(len(X_test),1)
```

```
In [0]: 1 input_num_3 = Input(shape = (train_price.shape[1],),name = 'input_price')
```

```
In [0]: 1 #concatenating all
2 num_final_concatenate = concatenate(inputs = [input_num_1,input_num_2,input_n
3 #dense layer for numerical features as described by the model architecture
4 dense_num_layer = Dense(16,activation = 'relu',kernel_initializer = 'he_norma
```

Assignment : 14

1. Preprocess all the Data we have in DonorsChoose [Dataset \(https://drive.google.com/drive/folders/1MIwK7BQMeV8f5CbDDVNLPaFGB32pFN60\)](https://drive.google.com/drive/folders/1MIwK7BQMeV8f5CbDDVNLPaFGB32pFN60) use train.csv
2. Combine 4 essay's into one column named - 'preprocessed_essays'.
3. After step 2 you have to train 3 types of models as discussed below.
4. For all the model use 'auc' (https://scikit-learn.org/stable/modules/model_evaluation.html#roc-metrics) as a metric. check [this \(https://datascience.stackexchange.com/a/20192\)](https://datascience.stackexchange.com/a/20192) for using auc as a metric
5. You are free to choose any number of layers/hiddden units but you have to use same type of architectures shown below.
6. You can use any one of the optimizers and choice of Learning rate and momentum, resources: [cs231n class notes \(http://cs231n.github.io/neural-networks-3/\)](http://cs231n.github.io/neural-networks-3/), [cs231n class video \(https://www.youtube.com/watch?v=hd_KFJ5ktUc\)](https://www.youtube.com/watch?v=hd_KFJ5ktUc).
7. For all the model's use [TensorBoard \(https://www.youtube.com/watch?v=2U6Jl7oqRkM\)](https://www.youtube.com/watch?v=2U6Jl7oqRkM) and plot the Metric value and Loss with epoch. While submitting, take a screenshot of plots and include those images in .ipynb notebook and PDF.
8. Use Categorical Cross Entropy as Loss to minimize.

Model-1

Build and Train deep neural network as shown below



ref: <https://i.imgur.com/w395Yk9.png> (<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_c** ---concatenate remaining columns and add a Dense layer after that.

- For LSTM, you can choose your sequence padding methods on your own or you can train your LSTM without padding, there is no restriction on that.

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

```
In [0]: 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/> (<https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-keras/>)

2. Please go through this link <https://keras.io/getting-started/functional-api-guide/> (<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.

```

In [0]: 1 x = concatenate(inputs = [flatten_text,flatten_prefix,flatten_clean_categorie
2         name = 'final_concatenated_layer')
3
4 x = Dense(8,activation = 'relu',kernel_initializer = 'he_normal',name = 'dens
5 x = Dropout(0.3,name = 'dropout_1')(x)
6 x = Dense(8,activation= 'relu',kernel_initializer = 'he_normal',name = 'dense
7 x = Dropout(0.3,name = 'dropout_2')(x)
8
9 output_layer = Dense(1,activation = 'sigmoid',name = 'output_layer')(x)#for b
10
11 model = Model(inputs = [input_text_layer,input_layer_prefix,input_layer_clean
12         input_num_1,input_num_2,input_num_3],outputs = [output
13 import warnings
14 warnings.filterwarnings('ignore')
15 model.summary()

```

| Layer (type) | Output Shape | Param # | Connected to |
|--|------------------|----------|---------------------------------|
| input_total_text_sequence (InputLayer) | (None, 300) | 0 | |
| embedding_19 (Embedding) | (None, 300, 300) | 15345900 | input_total_text_sequence[0][0] |
| input_prefix (InputLayer) | (None, 5) | 0 | |
| input_categories (InputLayer) | (None, 51) | 0 | |
| input_subcategories (InputLayer) | (None, 401) | 0 | |
| input_state (InputLayer) | (None, 51) | 0 | |
| input_grade_category (InputLayer) | (None, 4) | 0 | |
| input_previously_posted_project (InputLayer) | (None, 1) | 0 | |
| input_quantity (InputLayer) | (None, 1) | 0 | |
| input_price (InputLayer) | (None, 1) | 0 | |
| lstm_2 (LSTM) | (None, 300, 32) | 42624 | embedding_19[0][0] |

| | | | |
|---|----------------|------|--|
| embedding_16 (Embedding) [0][0] | (None, 5, 4) | 20 | input_prefix |
| embedding_14 (Embedding) es[0][0] | (None, 51, 4) | 204 | input_categori |
| embedding_15 (Embedding) ories[0][0] | (None, 401, 4) | 1604 | input_subcateg |
| embedding_17 (Embedding) [0] | (None, 51, 4) | 204 | input_state[0] |
| embedding_18 (Embedding) tegrity[0][0] | (None, 4, 4) | 16 | input_grade_ca |
| numerical_input_final (Concaten ly_posted_projects[[0][0] [0] | (None, 3) | 0 | input_previous input_quantity input_price[0] |
| flatten_18 (Flatten) | (None, 9600) | 0 | lstm_2[0][0] |
| flatten_15 (Flatten) [0][0] | (None, 20) | 0 | embedding_16 |
| flatten_13 (Flatten) [0][0] | (None, 204) | 0 | embedding_14 |
| flatten_14 (Flatten) [0][0] | (None, 1604) | 0 | embedding_15 |
| flatten_16 (Flatten) [0][0] | (None, 204) | 0 | embedding_17 |
| flatten_17 (Flatten) [0][0] | (None, 16) | 0 | embedding_18 |
| dense_layer_1 (Dense) t_final[0][0] | (None, 16) | 64 | numerical_inpu |
| final_concatenated_layer (Conca | (None, 11664) | 0 | flatten_18[0] |

| | | | |
|----------------------------------|-----------|-------|----------------|
| [0] | | | flatten_15[0] |
| [0] | | | flatten_13[0] |
| [0] | | | flatten_14[0] |
| [0] | | | flatten_16[0] |
| [0] | | | flatten_17[0] |
| [0] | | | dense_layer_1 |
| [0][0] | | | |
| <hr/> | | | |
| dense_layer_concatenate_1 (Dens | (None, 8) | 93320 | final_concaten |
| ated_layer[0][0] | | | |
| <hr/> | | | |
| dropout_1 (Dropout) | (None, 8) | 0 | dense_layer_co |
| ncatenate_1[0][0] | | | |
| <hr/> | | | |
| dense_layer_concatenate_2 (Dens | (None, 8) | 72 | dropout_1[0] |
| [0] | | | |
| <hr/> | | | |
| dropout_2 (Dropout) | (None, 8) | 0 | dense_layer_co |
| ncatenate_2[0][0] | | | |
| <hr/> | | | |
| output_layer (Dense) | (None, 1) | 9 | dropout_2[0] |
| [0] | | | |
| <hr/> | | | |
| ===== | | | |
| ===== | | | |
| Total params: 15,484,037 | | | |
| Trainable params: 138,137 | | | |
| Non-trainable params: 15,345,900 | | | |
| <hr/> | | | |



Defining the custom metric for AUC

```
In [0]: 1 import tensorflow as tf
2 from keras import backend as K
3 from sklearn.metrics import roc_auc_score
4
5 #https://stackoverflow.com/questions/51922500/tf-metrics-auc-yielding-very-di
6 def roc_auc(y_true, y_pred):
7     auc = tf.py_func(roc_auc_score, (y_true, y_pred), tf.double)
8     #auc = tf.metrics.auc(y_true, y_pred, num_thresholds=200)[1]
9     K.get_session().run(tf.local_variables_initializer())
10    return auc
```

Defining Callbacks

```
In [0]: 1 import datetime
        2 from tensorflow.python.keras.callbacks import TensorBoard
        3 from keras.callbacks import ModelCheckpoint
        4
        5 tensorboard = TensorBoard("logs\\fit\\" + datetime.datetime.now().strftime("%
        6 filepath="weights_best.hdf5"
        7 checkpoint = ModelCheckpoint(filepath, monitor='val_roc_auc', verbose=1, save
```

```
In [0]: 1 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=[roc_auc]
2 model.fit(x=[padded_docs_train,train_prefix,train_categories,train_subcategor
3 y=[np.array(Y_train)],
4 validation_data=([padded_docs_val,val_prefix,val_categories,val_sub
5 epochs=7,
6 batch_size=1024,
7 callbacks=[tensorboard,checkpoint])
```

Train on 69918 samples, validate on 17480 samples

Epoch 1/10

69918/69918 [=====] - ETA: 8:21 - loss: 0.6649 - roc_auc: 0.520 - ETA: 6:40 - loss: 0.6117 - roc_auc: 0.514 - ETA: 6:02 - loss: 0.5718 - roc_auc: 0.518 - ETA: 5:43 - loss: 0.5482 - roc_auc: 0.525 - ETA: 5:26 - loss: 0.5364 - roc_auc: 0.522 - ETA: 5:19 - loss: 0.5336 - roc_auc: 0.511 - ETA: 5:23 - loss: 0.5316 - roc_auc: 0.512 - ETA: 5:12 - loss: 0.5278 - roc_auc: 0.513 - ETA: 5:04 - loss: 0.5250 - roc_auc: 0.512 - ETA: 4:59 - loss: 0.5223 - roc_auc: 0.511 - ETA: 4:52 - loss: 0.5185 - roc_auc: 0.513 - ETA: 4:45 - loss: 0.5137 - roc_auc: 0.517 - ETA: 4:39 - loss: 0.5109 - roc_auc: 0.517 - ETA: 4:32 - loss: 0.5105 - roc_auc: 0.518 - ETA: 4:26 - loss: 0.5102 - roc_auc: 0.516 - ETA: 4:21 - loss: 0.5092 - roc_auc: 0.514 - ETA: 4:15 - loss: 0.5052 - roc_auc: 0.516 - ETA: 4:09 - loss: 0.5045 - roc_auc: 0.519 - ETA: 4:04 - loss: 0.5031 - roc_auc: 0.521 - ETA: 3:59 - loss: 0.5009 - roc_auc: 0.521 - ETA: 3:54 - loss: 0.5000 - roc_auc: 0.521 - ETA: 3:49 - loss: 0.4988 - roc_auc: 0.522 - ETA: 3:45 - loss: 0.4970 - roc_auc: 0.525 - ETA: 3:40 - loss: 0.4962 - roc_auc: 0.524 - ETA: 3:35 - loss: 0.4952 - roc_auc: 0.524 - ETA: 3:29 - loss: 0.4935 - roc_auc: 0.526 - ETA: 3:25 - loss: 0.4928 - roc_auc: 0.525 - ETA: 3:20 - loss: 0.4931 - roc_auc: 0.526 - ETA: 3:15 - loss: 0.4928 - roc_auc: 0.526 - ETA: 3:09 - loss: 0.4909 - roc_auc: 0.528 - ETA: 3:04 - loss: 0.4894 - roc_auc: 0.532 - ETA: 2:59 - loss: 0.4887 - roc_auc: 0.534 - ETA: 2:54 - loss: 0.4877 - roc_auc: 0.536 - ETA: 2:49 - loss: 0.4868 - roc_auc: 0.538 - ETA: 2:44 - loss: 0.4860 - roc_auc: 0.539 - ETA: 2:39 - loss: 0.4843 - roc_auc: 0.540 - ETA: 2:34 - loss: 0.4841 - roc_auc: 0.541 - ETA: 2:29 - loss: 0.4837 - roc_auc: 0.543 - ETA: 2:24 - loss: 0.4829 - roc_auc: 0.545 - ETA: 2:19 - loss: 0.4817 - roc_auc: 0.547 - ETA: 2:15 - loss: 0.4802 - roc_auc: 0.550 - ETA: 2:10 - loss: 0.4803 - roc_auc: 0.551 - ETA: 2:05 - loss: 0.4794 - roc_auc: 0.552 - ETA: 2:00 - loss: 0.4787 - roc_auc: 0.553 - ETA: 1:55 - loss: 0.4774 - roc_auc: 0.555 - ETA: 1:50 - loss: 0.4764 - roc_auc: 0.557 - ETA: 1:45 - loss: 0.4751 - roc_auc: 0.559 - ETA: 1:40 - loss: 0.4750 - roc_auc: 0.559 - ETA: 1:35 - loss: 0.4743 - roc_auc: 0.561 - ETA: 1:30 - loss: 0.4744 - roc_auc: 0.561 - ETA: 1:26 - loss: 0.4734 - roc_auc: 0.562 - ETA: 1:21 - loss: 0.4729 - roc_auc: 0.563 - ETA: 1:16 - loss: 0.4723 - roc_auc: 0.564 - ETA: 1:11 - loss: 0.4717 - roc_auc: 0.566 - ETA: 1:06 - loss: 0.4709 - roc_auc: 0.567 - ETA: 1:01 - loss: 0.4701 - roc_auc: 0.569 - ETA: 56s - loss: 0.4694 - roc_auc: 0.570 - ETA: 51s - loss: 0.4684 - roc_auc: 0.57 - ETA: 46s - loss: 0.4685 - roc_auc: 0.57 - ETA: 41s - loss: 0.4682 - roc_auc: 0.57 - ETA: 36s - loss: 0.4674 - roc_auc: 0.57 - ETA: 31s - loss: 0.4662 - roc_auc: 0.57 - ETA: 26s - loss: 0.4657 - roc_auc: 0.57 - ETA: 21s - loss: 0.4652 - roc_auc: 0.57 - ETA: 16s - loss: 0.4645 - roc_auc: 0.57 - ETA: 11s - loss: 0.4642 - roc_auc: 0.57 - ETA: 6s - loss: 0.4639 - roc_auc: 0.5798 - ETA: 1s - loss: 0.4635 - roc_auc: 0.580 - 400s 6ms/step - loss: 0.4634 - roc_auc: 0.5813 - val_loss: 0.3907 - val_roc_auc: 0.7212

Epoch 00001: val_roc_auc improved from -inf to 0.72119, saving model to weights_best.hdf5

Epoch 2/10

```
69918/69918 [=====] - ETA: 5:38 - loss: 0.4490 - roc_auc: 0.593 - ETA: 5:51 - loss: 0.4414 - roc_auc: 0.624 - ETA: 5:40 - loss: 0.4385 - roc_auc: 0.632 - ETA: 5:31 - loss: 0.4340 - roc_auc: 0.639 - ETA: 5:24 - loss: 0.4313 - roc_auc: 0.640 - ETA: 5:16 - loss: 0.4306 - roc_auc: 0.641 - ETA: 5:10 - loss: 0.4236 - roc_auc: 0.645 - ETA: 5:03 - loss: 0.4241 - roc_auc: 0.643 - ETA: 4:59 - loss: 0.4288 - roc_auc: 0.644 - ETA: 4:53 - loss: 0.4347 - roc_auc: 0.641 - ETA: 4:48 - loss: 0.4375 - roc_auc: 0.641 - ETA: 4:44 - loss: 0.4403 - roc_auc: 0.643 - ETA: 4:39 - loss: 0.4415 - roc_auc: 0.644 - ETA: 4:36 - loss: 0.4435 - roc_auc: 0.643 - ETA: 4:32 - loss: 0.4435 - roc_auc: 0.644 - ETA: 4:27 - loss: 0.4429 - roc_auc: 0.647 - ETA: 4:23 - loss: 0.4420 - roc_auc: 0.646 - ETA: 4:17 - loss: 0.4425 - roc_auc: 0.649 - ETA: 4:12 - loss: 0.4405 - roc_auc: 0.649 - ETA: 4:07 - loss: 0.4424 - roc_auc: 0.646 - ETA: 4:02 - loss: 0.4423 - roc_auc: 0.646 - ETA: 3:56 - loss: 0.4426 - roc_auc: 0.646 - ETA: 3:51 - loss: 0.4424 - roc_auc: 0.645 - ETA: 3:45 - loss: 0.4420 - roc_auc: 0.645 - ETA: 3:41 - loss: 0.4412 - roc_auc: 0.645 - ETA: 3:36 - loss: 0.4404 - roc_auc: 0.645 - ETA: 3:31 - loss: 0.4392 - roc_auc: 0.646 - ETA: 3:26 - loss: 0.4389 - roc_auc: 0.648 - ETA: 3:21 - loss: 0.4380 - roc_auc: 0.648 - ETA: 3:16 - loss: 0.4377 - roc_auc: 0.647 - ETA: 3:10 - loss: 0.4375 - roc_auc: 0.647 - ETA: 3:05 - loss: 0.4371 - roc_auc: 0.646 - ETA: 3:00 - loss: 0.4374 - roc_auc: 0.647 - ETA: 2:54 - loss: 0.4383 - roc_auc: 0.646 - ETA: 2:49 - loss: 0.4384 - roc_auc: 0.646 - ETA: 2:44 - loss: 0.4381 - roc_auc: 0.646 - ETA: 2:39 - loss: 0.4384 - roc_auc: 0.647 - ETA: 2:34 - loss: 0.4385 - roc_auc: 0.647 - ETA: 2:29 - loss: 0.4380 - roc_auc: 0.647 - ETA: 2:24 - loss: 0.4379 - roc_auc: 0.648 - ETA: 2:18 - loss: 0.4370 - roc_auc: 0.648 - ETA: 2:14 - loss: 0.4363 - roc_auc: 0.650 - ETA: 2:08 - loss: 0.4355 - roc_auc: 0.650 - ETA: 2:03 - loss: 0.4351 - roc_auc: 0.651 - ETA: 1:58 - loss: 0.4352 - roc_auc: 0.651 - ETA: 1:53 - loss: 0.4346 - roc_auc: 0.651 - ETA: 1:48 - loss: 0.4343 - roc_auc: 0.652 - ETA: 1:43 - loss: 0.4353 - roc_auc: 0.651 - ETA: 1:38 - loss: 0.4348 - roc_auc: 0.651 - ETA: 1:33 - loss: 0.4345 - roc_auc: 0.652 - ETA: 1:28 - loss: 0.4344 - roc_auc: 0.652 - ETA: 1:22 - loss: 0.4340 - roc_auc: 0.653 - ETA: 1:17 - loss: 0.4340 - roc_auc: 0.653 - ETA: 1:12 - loss: 0.4338 - roc_auc: 0.653 - ETA: 1:07 - loss: 0.4334 - roc_auc: 0.653 - ETA: 1:02 - loss: 0.4334 - roc_auc: 0.654 - ETA: 57s - loss: 0.4330 - roc_auc: 0.655 - ETA: 52s - loss: 0.4332 - roc_auc: 0.65 - ETA: 47s - loss: 0.4328 - roc_auc: 0.65 - ETA: 42s - loss: 0.4318 - roc_auc: 0.65 - ETA: 37s - loss: 0.4311 - roc_auc: 0.65 - ETA: 32s - loss: 0.4310 - roc_auc: 0.65 - ETA: 27s - loss: 0.4301 - roc_auc: 0.65 - ETA: 22s - loss: 0.4307 - roc_auc: 0.65 - ETA: 17s - loss: 0.4304 - roc_auc: 0.65 - ETA: 11s - loss: 0.4305 - roc_auc: 0.65 - ETA: 6s - loss: 0.4303 - roc_auc: 0.6588 - ETA: 1s - loss: 0.4302 - roc_auc: 0.659 - 407s 6ms/step - loss: 0.4301 - roc_auc: 0.6591 - val_loss: 0.3856 - val_roc_auc: 0.7309
```

Epoch 00002: val_roc_auc improved from 0.72119 to 0.73091, saving model to weights_best.hdf5

Epoch 3/10

```
69918/69918 [=====] - ETA: 6:04 - loss: 0.4189 - roc_auc: 0.696 - ETA: 6:11 - loss: 0.4158 - roc_auc: 0.698 - ETA: 6:14 - loss: 0.4168 - roc_auc: 0.693 - ETA: 6:14 - loss: 0.4149 - roc_auc: 0.687 - ETA: 6:04 - loss: 0.4167 - roc_auc: 0.682 - ETA: 5:54 - loss: 0.4168 - roc_auc: 0.684 - ETA: 5:48 - loss: 0.4194 - roc_auc: 0.681 - ETA: 5:44 - loss: 0.4210 - roc_auc: 0.685 - ETA: 5:39 - loss: 0.4209 - roc_auc: 0.681 - ETA: 5:34 - loss: 0.4208 - roc_auc: 0.679 - ETA: 5:28 - loss: 0.4198 - roc_auc: 0.678 - ETA: 5:24 - loss: 0.4180 - roc_auc: 0.680 - ETA: 5:18 - loss: 0.4154 - roc_auc: 0.680 - ETA: 5:14 - loss: 0.4156 - roc_auc: 0.683 - ETA: 5:09 - loss: 0.4150 - roc_auc: 0.685 - ETA: 5:04 - loss: 0.4160 - roc_auc: 0.684 - ETA: 4:57 - loss: 0.4136 - roc_auc: 0.685 - ETA: 4:53 - loss: 0.4133 - roc_auc: 0.686 - ETA: 4:46 - loss: 0.4116 - roc_auc: 0.688 - ETA: 4:41 - loss: 0.4117 - roc_auc: 0.68
```


7 - ETA: 4:35 - loss: 0.4114 - roc_auc: 0.689 - ETA: 4:31 - loss: 0.4114 - roc_auc: 0.688 - ETA: 4:25 - loss: 0.4112 - roc_auc: 0.688 - ETA: 4:22 - loss: 0.4111 - roc_auc: 0.688 - ETA: 4:16 - loss: 0.4113 - roc_auc: 0.686 - ETA: 4:09 - loss: 0.4117 - roc_auc: 0.686 - ETA: 4:03 - loss: 0.4116 - roc_auc: 0.686 - ETA: 3:57 - loss: 0.4122 - roc_auc: 0.685 - ETA: 3:51 - loss: 0.4123 - roc_auc: 0.685 - ETA: 3:45 - loss: 0.4117 - roc_auc: 0.686 - ETA: 3:41 - loss: 0.4113 - roc_auc: 0.685 - ETA: 3:35 - loss: 0.4109 - roc_auc: 0.685 - ETA: 3:30 - loss: 0.4106 - roc_auc: 0.685 - ETA: 3:24 - loss: 0.4106 - roc_auc: 0.686 - ETA: 3:18 - loss: 0.4103 - roc_auc: 0.687 - ETA: 3:12 - loss: 0.4095 - roc_auc: 0.686 - ETA: 3:07 - loss: 0.4100 - roc_auc: 0.686 - ETA: 3:01 - loss: 0.4100 - roc_auc: 0.686 - ETA: 2:55 - loss: 0.4107 - roc_auc: 0.686 - ETA: 2:49 - loss: 0.4112 - roc_auc: 0.685 - ETA: 2:43 - loss: 0.4111 - roc_auc: 0.687 - ETA: 2:37 - loss: 0.4104 - roc_auc: 0.688 - ETA: 2:31 - loss: 0.4103 - roc_auc: 0.688 - ETA: 2:26 - loss: 0.4103 - roc_auc: 0.688 - ETA: 2:20 - loss: 0.4102 - roc_auc: 0.689 - ETA: 2:13 - loss: 0.4109 - roc_auc: 0.688 - ETA: 2:07 - loss: 0.4105 - roc_auc: 0.688 - ETA: 2:01 - loss: 0.4102 - roc_auc: 0.688 - ETA: 1:55 - loss: 0.4102 - roc_auc: 0.688 - ETA: 1:50 - loss: 0.4111 - roc_auc: 0.688 - ETA: 1:44 - loss: 0.4110 - roc_auc: 0.688 - ETA: 1:38 - loss: 0.4112 - roc_auc: 0.688 - ETA: 1:32 - loss: 0.4115 - roc_auc: 0.688 - ETA: 1:27 - loss: 0.4113 - roc_auc: 0.689 - ETA: 1:20 - loss: 0.4114 - roc_auc: 0.689 - ETA: 1:14 - loss: 0.4115 - roc_auc: 0.688 - ETA: 1:09 - loss: 0.4118 - roc_auc: 0.688 - ETA: 1:02 - loss: 0.4124 - roc_auc: 0.688 - ETA: 56s - loss: 0.4120 - roc_auc: 0.688 - ETA: 50s - loss: 0.4115 - roc_auc: 0.688 - ETA: 44s - loss: 0.4122 - roc_auc: 0.688 - ETA: 38s - loss: 0.4118 - roc_auc: 0.688 - ETA: 32s - loss: 0.4116 - roc_auc: 0.688 - ETA: 26s - loss: 0.4115 - roc_auc: 0.688 - ETA: 20s - loss: 0.4115 - roc_auc: 0.688 - ETA: 14s - loss: 0.4115 - roc_auc: 0.688 - ETA: 7s - loss: 0.4115 - roc_auc: 0.688 - ETA: 1s - loss: 0.4120 - roc_auc: 0.688 - 472s 7ms/step - loss: 0.4122 - roc_auc: 0.688 - val_loss: 0.3820 - val_roc_auc: 0.7361

Epoch 00003: val_roc_auc improved from 0.73091 to 0.73612, saving model to weights_best.hdf5

Epoch 4/10

69918/69918 [=====] - ETA: 7:47 - loss: 0.4019 - roc_auc: 0.692 - ETA: 7:21 - loss: 0.4000 - roc_auc: 0.709 - ETA: 6:44 - loss: 0.4052 - roc_auc: 0.702 - ETA: 6:39 - loss: 0.4098 - roc_auc: 0.696 - ETA: 6:22 - loss: 0.4125 - roc_auc: 0.697 - ETA: 6:16 - loss: 0.4113 - roc_auc: 0.699 - ETA: 6:08 - loss: 0.4075 - roc_auc: 0.698 - ETA: 5:57 - loss: 0.4081 - roc_auc: 0.698 - ETA: 5:50 - loss: 0.4082 - roc_auc: 0.699 - ETA: 5:43 - loss: 0.4052 - roc_auc: 0.701 - ETA: 5:39 - loss: 0.4038 - roc_auc: 0.699 - ETA: 5:31 - loss: 0.4037 - roc_auc: 0.700 - ETA: 5:25 - loss: 0.4036 - roc_auc: 0.700 - ETA: 5:20 - loss: 0.4020 - roc_auc: 0.703 - ETA: 5:17 - loss: 0.4020 - roc_auc: 0.702 - ETA: 5:10 - loss: 0.4010 - roc_auc: 0.705 - ETA: 5:05 - loss: 0.4016 - roc_auc: 0.704 - ETA: 5:00 - loss: 0.4021 - roc_auc: 0.703 - ETA: 4:55 - loss: 0.4029 - roc_auc: 0.702 - ETA: 4:51 - loss: 0.4032 - roc_auc: 0.702 - ETA: 4:46 - loss: 0.4039 - roc_auc: 0.700 - ETA: 4:39 - loss: 0.4044 - roc_auc: 0.700 - ETA: 4:32 - loss: 0.4053 - roc_auc: 0.700 - ETA: 4:26 - loss: 0.4050 - roc_auc: 0.701 - ETA: 4:19 - loss: 0.4049 - roc_auc: 0.701 - ETA: 4:13 - loss: 0.4048 - roc_auc: 0.703 - ETA: 4:07 - loss: 0.4039 - roc_auc: 0.700 - ETA: 4:01 - loss: 0.4038 - roc_auc: 0.705 - ETA: 3:56 - loss: 0.4025 - roc_auc: 0.707 - ETA: 3:50 - loss: 0.4034 - roc_auc: 0.707 - ETA: 3:44 - loss: 0.4025 - roc_auc: 0.708 - ETA: 3:37 - loss: 0.4009 - roc_auc: 0.709 - ETA: 3:31 - loss: 0.4018 - roc_auc: 0.708 - ETA: 3:25 - loss: 0.4020 - roc_auc: 0.708 - ETA: 3:19 - loss: 0.4019 - roc_auc: 0.708 - ETA: 3:13 - loss: 0.4013 - roc_auc: 0.709 - ETA: 3:07 - loss: 0.4013 - roc_auc: 0.710 - ETA: 3:01 - loss: 0.4013 - roc_auc: 0.710 - ETA: 2:55 - loss: 0.4003 - roc_auc: 0.711 - ETA: 2:49 - loss: 0.3996 - roc_auc: 0.712 - ETA: 2:43 - loss: 0.3992 - roc_auc: 0.712

2 - ETA: 2:36 - loss: 0.3985 - roc_auc: 0.712 - ETA: 2:30 - loss: 0.3989 - roc_auc: 0.713 - ETA: 2:24 - loss: 0.3983 - roc_auc: 0.713 - ETA: 2:19 - loss: 0.3979 - roc_auc: 0.714 - ETA: 2:13 - loss: 0.3976 - roc_auc: 0.714 - ETA: 2:07 - loss: 0.3978 - roc_auc: 0.714 - ETA: 2:01 - loss: 0.3977 - roc_auc: 0.715 - ETA: 1:55 - loss: 0.3976 - roc_auc: 0.714 - ETA: 1:49 - loss: 0.3974 - roc_auc: 0.715 - ETA: 1:43 - loss: 0.3966 - roc_auc: 0.716 - ETA: 1:37 - loss: 0.3969 - roc_auc: 0.716 - ETA: 1:31 - loss: 0.3967 - roc_auc: 0.716 - ETA: 1:25 - loss: 0.3971 - roc_auc: 0.716 - ETA: 1:19 - loss: 0.3964 - roc_auc: 0.716 - ETA: 1:13 - loss: 0.3961 - roc_auc: 0.716 - ETA: 1:07 - loss: 0.3961 - roc_auc: 0.716 - ETA: 1:02 - loss: 0.3961 - roc_auc: 0.717 - ETA: 56s - loss: 0.3961 - roc_auc: 0.718 - ETA: 50s - loss: 0.3967 - roc_auc: 0.71 - ETA: 44s - loss: 0.3964 - roc_auc: 0.71 - ETA: 38s - loss: 0.3963 - roc_auc: 0.71 - ETA: 32s - loss: 0.3960 - roc_auc: 0.71 - ETA: 26s - loss: 0.3961 - roc_auc: 0.71 - ETA: 19s - loss: 0.3958 - roc_auc: 0.71 - ETA: 13s - loss: 0.3965 - roc_auc: 0.71 - ETA: 7s - loss: 0.3971 - roc_auc: 0.7175 - ETA: 1s - loss: 0.3968 - roc_auc: 0.717 - 468s 7ms/step - loss: 0.3968 - roc_auc: 0.7177 - val_loss: 0.3848 - val_roc_auc: 0.7336

Epoch 00004: val_roc_auc did not improve from 0.73612

Epoch 5/10

69918/69918 [=====] - ETA: 7:26 - loss: 0.3734 - roc_auc: 0.736 - ETA: 7:07 - loss: 0.3825 - roc_auc: 0.739 - ETA: 6:53 - loss: 0.3943 - roc_auc: 0.725 - ETA: 6:47 - loss: 0.3895 - roc_auc: 0.729 - ETA: 6:56 - loss: 0.3914 - roc_auc: 0.731 - ETA: 6:50 - loss: 0.3960 - roc_auc: 0.729 - ETA: 6:45 - loss: 0.3987 - roc_auc: 0.727 - ETA: 6:39 - loss: 0.3994 - roc_auc: 0.729 - ETA: 6:35 - loss: 0.4012 - roc_auc: 0.723 - ETA: 6:33 - loss: 0.3992 - roc_auc: 0.724 - ETA: 6:31 - loss: 0.3984 - roc_auc: 0.723 - ETA: 6:30 - loss: 0.3983 - roc_auc: 0.724 - ETA: 6:27 - loss: 0.3992 - roc_auc: 0.721 - ETA: 6:20 - loss: 0.3979 - roc_auc: 0.719 - ETA: 6:14 - loss: 0.3983 - roc_auc: 0.721 - ETA: 6:05 - loss: 0.3963 - roc_auc: 0.723 - ETA: 5:58 - loss: 0.3946 - roc_auc: 0.724 - ETA: 5:51 - loss: 0.3933 - roc_auc: 0.726 - ETA: 5:44 - loss: 0.3944 - roc_auc: 0.725 - ETA: 5:36 - loss: 0.3951 - roc_auc: 0.724 - ETA: 5:30 - loss: 0.3942 - roc_auc: 0.725 - ETA: 5:26 - loss: 0.3929 - roc_auc: 0.725 - ETA: 5:21 - loss: 0.3916 - roc_auc: 0.726 - ETA: 5:17 - loss: 0.3900 - roc_auc: 0.726 - ETA: 5:09 - loss: 0.3898 - roc_auc: 0.727 - ETA: 5:02 - loss: 0.3906 - roc_auc: 0.726 - ETA: 4:56 - loss: 0.3892 - roc_auc: 0.726 - ETA: 4:50 - loss: 0.3909 - roc_auc: 0.725 - ETA: 4:42 - loss: 0.3907 - roc_auc: 0.724 - ETA: 4:34 - loss: 0.3905 - roc_auc: 0.725 - ETA: 4:26 - loss: 0.3901 - roc_auc: 0.726 - ETA: 4:17 - loss: 0.3895 - roc_auc: 0.727 - ETA: 4:09 - loss: 0.3887 - roc_auc: 0.727 - ETA: 4:02 - loss: 0.3888 - roc_auc: 0.728 - ETA: 3:54 - loss: 0.3872 - roc_auc: 0.729 - ETA: 3:45 - loss: 0.3881 - roc_auc: 0.729 - ETA: 3:38 - loss: 0.3890 - roc_auc: 0.728 - ETA: 3:30 - loss: 0.3885 - roc_auc: 0.729 - ETA: 3:23 - loss: 0.3882 - roc_auc: 0.730 - ETA: 3:16 - loss: 0.3882 - roc_auc: 0.729 - ETA: 3:09 - loss: 0.3880 - roc_auc: 0.729 - ETA: 3:02 - loss: 0.3874 - roc_auc: 0.729 - ETA: 2:56 - loss: 0.3876 - roc_auc: 0.729 - ETA: 2:50 - loss: 0.3879 - roc_auc: 0.728 - ETA: 2:44 - loss: 0.3880 - roc_auc: 0.728 - ETA: 2:37 - loss: 0.3879 - roc_auc: 0.729 - ETA: 2:30 - loss: 0.3878 - roc_auc: 0.730 - ETA: 2:22 - loss: 0.3877 - roc_auc: 0.730 - ETA: 2:16 - loss: 0.3883 - roc_auc: 0.730 - ETA: 2:08 - loss: 0.3879 - roc_auc: 0.731 - ETA: 2:01 - loss: 0.3881 - roc_auc: 0.731 - ETA: 1:54 - loss: 0.3883 - roc_auc: 0.731 - ETA: 1:47 - loss: 0.3875 - roc_auc: 0.731 - ETA: 1:40 - loss: 0.3874 - roc_auc: 0.732 - ETA: 1:33 - loss: 0.3879 - roc_auc: 0.731 - ETA: 1:26 - loss: 0.3885 - roc_auc: 0.731 - ETA: 1:19 - loss: 0.3885 - roc_auc: 0.731 - ETA: 1:12 - loss: 0.3888 - roc_auc: 0.732 - ETA: 1:04 - loss: 0.3892 - roc_auc: 0.731 - ETA: 57s - loss: 0.3893 - roc_auc: 0.732 - ETA: 50s - loss: 0.3888 - roc_auc: 0.73 - ETA: 43s - loss: 0.3886 - roc_auc: 0.73 - ETA: 36s - loss: 0.3889 - roc_auc: 0.73 - ETA: 30s - loss: 0.3892 - roc_auc: 0.73

73 - ETA: 23s - loss: 0.3895 - roc_auc: 0.73 - ETA: 16s - loss: 0.3894 - roc_auc: 0.73 - ETA: 9s - loss: 0.3899 - roc_auc: 0.7327 - ETA: 1s - loss: 0.3898 - roc_auc: 0.733 - 542s 8ms/step - loss: 0.3898 - roc_auc: 0.7329 - val_loss: 0.3853 - val_roc_auc: 0.7324

Epoch 00005: val_roc_auc did not improve from 0.73612

Epoch 6/10

69918/69918 [=====] - ETA: 7:50 - loss: 0.3892 - roc_auc: 0.772 - ETA: 8:22 - loss: 0.3791 - roc_auc: 0.773 - ETA: 8:42 - loss: 0.3932 - roc_auc: 0.757 - ETA: 8:11 - loss: 0.3954 - roc_auc: 0.745 - ETA: 8:11 - loss: 0.3923 - roc_auc: 0.743 - ETA: 7:55 - loss: 0.3928 - roc_auc: 0.748 - ETA: 7:45 - loss: 0.3922 - roc_auc: 0.749 - ETA: 7:33 - loss: 0.3885 - roc_auc: 0.751 - ETA: 7:19 - loss: 0.3882 - roc_auc: 0.751 - ETA: 7:04 - loss: 0.3849 - roc_auc: 0.751 - ETA: 7:00 - loss: 0.3829 - roc_auc: 0.752 - ETA: 6:54 - loss: 0.3815 - roc_auc: 0.751 - ETA: 6:53 - loss: 0.3815 - roc_auc: 0.748 - ETA: 6:48 - loss: 0.3804 - roc_auc: 0.750 - ETA: 6:43 - loss: 0.3803 - roc_auc: 0.750 - ETA: 6:32 - loss: 0.3820 - roc_auc: 0.748 - ETA: 6:23 - loss: 0.3819 - roc_auc: 0.750 - ETA: 6:18 - loss: 0.3818 - roc_auc: 0.751 - ETA: 6:12 - loss: 0.3818 - roc_auc: 0.752 - ETA: 6:06 - loss: 0.3808 - roc_auc: 0.752 - ETA: 6:02 - loss: 0.3802 - roc_auc: 0.751 - ETA: 5:54 - loss: 0.3814 - roc_auc: 0.749 - ETA: 5:46 - loss: 0.3814 - roc_auc: 0.749 - ETA: 5:38 - loss: 0.3831 - roc_auc: 0.748 - ETA: 5:28 - loss: 0.3821 - roc_auc: 0.749 - ETA: 5:20 - loss: 0.3821 - roc_auc: 0.750 - ETA: 5:12 - loss: 0.3821 - roc_auc: 0.750 - ETA: 5:05 - loss: 0.3821 - roc_auc: 0.751 - ETA: 5:00 - loss: 0.3822 - roc_auc: 0.751 - ETA: 4:53 - loss: 0.3816 - roc_auc: 0.752 - ETA: 4:46 - loss: 0.3813 - roc_auc: 0.752 - ETA: 4:39 - loss: 0.3807 - roc_auc: 0.752 - ETA: 4:32 - loss: 0.3811 - roc_auc: 0.752 - ETA: 4:25 - loss: 0.3800 - roc_auc: 0.752 - ETA: 4:17 - loss: 0.3802 - roc_auc: 0.752 - ETA: 4:09 - loss: 0.3798 - roc_auc: 0.752 - ETA: 4:01 - loss: 0.3796 - roc_auc: 0.752 - ETA: 3:53 - loss: 0.3795 - roc_auc: 0.752 - ETA: 3:46 - loss: 0.3795 - roc_auc: 0.752 - ETA: 3:39 - loss: 0.3790 - roc_auc: 0.752 - ETA: 3:32 - loss: 0.3803 - roc_auc: 0.751 - ETA: 3:25 - loss: 0.3797 - roc_auc: 0.752 - ETA: 3:19 - loss: 0.3791 - roc_auc: 0.752 - ETA: 3:12 - loss: 0.3795 - roc_auc: 0.751 - ETA: 3:04 - loss: 0.3798 - roc_auc: 0.751 - ETA: 2:56 - loss: 0.3801 - roc_auc: 0.751 - ETA: 2:48 - loss: 0.3806 - roc_auc: 0.750 - ETA: 2:39 - loss: 0.3804 - roc_auc: 0.751 - ETA: 2:30 - loss: 0.3805 - roc_auc: 0.750 - ETA: 2:22 - loss: 0.3799 - roc_auc: 0.751 - ETA: 2:14 - loss: 0.3799 - roc_auc: 0.751 - ETA: 2:05 - loss: 0.3800 - roc_auc: 0.751 - ETA: 1:57 - loss: 0.3796 - roc_auc: 0.751 - ETA: 1:49 - loss: 0.3800 - roc_auc: 0.751 - ETA: 1:41 - loss: 0.3794 - roc_auc: 0.751 - ETA: 1:33 - loss: 0.3788 - roc_auc: 0.752 - ETA: 1:25 - loss: 0.3784 - roc_auc: 0.752 - ETA: 1:18 - loss: 0.3785 - roc_auc: 0.752 - ETA: 1:10 - loss: 0.3786 - roc_auc: 0.752 - ETA: 1:02 - loss: 0.3786 - roc_auc: 0.752 - ETA: 55s - loss: 0.3793 - roc_auc: 0.751 - ETA: 47s - loss: 0.3796 - roc_auc: 0.75 - ETA: 39s - loss: 0.3796 - roc_auc: 0.75 - ETA: 32s - loss: 0.3800 - roc_auc: 0.75 - ETA: 24s - loss: 0.3801 - roc_auc: 0.75 - ETA: 17s - loss: 0.3803 - roc_auc: 0.75 - ETA: 9s - loss: 0.3808 - roc_auc: 0.7514 - ETA: 2s - loss: 0.3808 - roc_auc: 0.751 - 554s 8ms/step - loss: 0.3806 - roc_auc: 0.7514 - val_loss: 0.3863 - val_roc_auc: 0.7327

Epoch 00006: val_roc_auc did not improve from 0.73612

Epoch 7/10

69918/69918 [=====] - ETA: 6:39 - loss: 0.4008 - roc_auc: 0.726 - ETA: 6:41 - loss: 0.3919 - roc_auc: 0.744 - ETA: 7:06 - loss: 0.3860 - roc_auc: 0.747 - ETA: 7:16 - loss: 0.3797 - roc_auc: 0.756 - ETA: 7:17 - loss: 0.3823 - roc_auc: 0.757 - ETA: 7:13 - loss: 0.3789 - roc_auc: 0.759 - ETA: 7:12 - loss: 0.3744 - roc_auc: 0.763 - ETA: 7:00 - loss: 0.3764 - roc_auc: 0.759 - ETA: 6:48 - loss: 0.3744 - roc_auc: 0.763 - ETA: 6:39 - loss: 0.3730 - roc_auc: 0.764 - ETA: 6:28 - loss: 0.3725 - roc_auc: 0.767 - ETA: 6:27 - loss: 0.37

21 - roc_auc: 0.767 - ETA: 6:22 - loss: 0.3726 - roc_auc: 0.764 - ETA: 6:11 - loss: 0.3736 - roc_auc: 0.765 - ETA: 6:03 - loss: 0.3732 - roc_auc: 0.766 - ETA: 5:55 - loss: 0.3726 - roc_auc: 0.765 - ETA: 5:47 - loss: 0.3736 - roc_auc: 0.766 - ETA: 5:38 - loss: 0.3728 - roc_auc: 0.767 - ETA: 5:34 - loss: 0.3716 - roc_auc: 0.768 - ETA: 5:26 - loss: 0.3720 - roc_auc: 0.769 - ETA: 5:19 - loss: 0.3725 - roc_auc: 0.769 - ETA: 5:11 - loss: 0.3723 - roc_auc: 0.770 - ETA: 5:03 - loss: 0.3730 - roc_auc: 0.770 - ETA: 4:55 - loss: 0.3734 - roc_auc: 0.769 - ETA: 4:48 - loss: 0.3727 - roc_auc: 0.770 - ETA: 4:40 - loss: 0.3727 - roc_auc: 0.771 - ETA: 4:33 - loss: 0.3735 - roc_auc: 0.771 - ETA: 4:25 - loss: 0.3729 - roc_auc: 0.772 - ETA: 4:18 - loss: 0.3743 - roc_auc: 0.771 - ETA: 4:11 - loss: 0.3735 - roc_auc: 0.772 - ETA: 4:05 - loss: 0.3735 - roc_auc: 0.772 - ETA: 3:58 - loss: 0.3731 - roc_auc: 0.773 - ETA: 3:51 - loss: 0.3722 - roc_auc: 0.773 - ETA: 3:44 - loss: 0.3721 - roc_auc: 0.773 - ETA: 3:37 - loss: 0.3729 - roc_auc: 0.772 - ETA: 3:30 - loss: 0.3715 - roc_auc: 0.773 - ETA: 3:24 - loss: 0.3715 - roc_auc: 0.772 - ETA: 3:17 - loss: 0.3718 - roc_auc: 0.772 - ETA: 3:10 - loss: 0.3715 - roc_auc: 0.772 - ETA: 3:04 - loss: 0.3709 - roc_auc: 0.773 - ETA: 2:57 - loss: 0.3712 - roc_auc: 0.772 - ETA: 2:50 - loss: 0.3709 - roc_auc: 0.773 - ETA: 2:44 - loss: 0.3715 - roc_auc: 0.772 - ETA: 2:38 - loss: 0.3714 - roc_auc: 0.772 - ETA: 2:32 - loss: 0.3718 - roc_auc: 0.772 - ETA: 2:27 - loss: 0.3718 - roc_auc: 0.772 - ETA: 2:20 - loss: 0.3714 - roc_auc: 0.772 - ETA: 2:15 - loss: 0.3719 - roc_auc: 0.772 - ETA: 2:08 - loss: 0.3716 - roc_auc: 0.772 - ETA: 2:02 - loss: 0.3709 - roc_auc: 0.773 - ETA: 1:55 - loss: 0.3710 - roc_auc: 0.772 - ETA: 1:48 - loss: 0.3711 - roc_auc: 0.772 - ETA: 1:42 - loss: 0.3711 - roc_auc: 0.771 - ETA: 1:35 - loss: 0.3710 - roc_auc: 0.771 - ETA: 1:28 - loss: 0.3706 - roc_auc: 0.772 - ETA: 1:21 - loss: 0.3710 - roc_auc: 0.772 - ETA: 1:15 - loss: 0.3708 - roc_auc: 0.772 - ETA: 1:08 - loss: 0.3701 - roc_auc: 0.772 - ETA: 1:01 - loss: 0.3703 - roc_auc: 0.772 - ETA: 55s - loss: 0.3702 - roc_auc: 0.772 - ETA: 48s - loss: 0.3698 - roc_auc: 0.77 - ETA: 41s - loss: 0.3699 - roc_auc: 0.77 - ETA: 35s - loss: 0.3698 - roc_auc: 0.77 - ETA: 28s - loss: 0.3700 - roc_auc: 0.77 - ETA: 21s - loss: 0.3699 - roc_auc: 0.77 - ETA: 15s - loss: 0.3700 - roc_auc: 0.77 - ETA: 8s - loss: 0.3698 - roc_auc: 0.7709 - ETA: 1s - loss: 0.3702 - roc_auc: 0.770 - 502s 7ms/step - loss: 0.3704 - roc_auc: 0.7703 - val_loss: 0.3936 - val_roc_auc: 0.7214

Epoch 00007: val_roc_auc did not improve from 0.73612

Epoch 8/10

69918/69918 [=====] - ETA: 7:01 - loss: 0.3861 - roc_auc: 0.770 - ETA: 6:51 - loss: 0.3641 - roc_auc: 0.792 - ETA: 6:52 - loss: 0.3651 - roc_auc: 0.793 - ETA: 6:46 - loss: 0.3729 - roc_auc: 0.791 - ETA: 6:39 - loss: 0.3658 - roc_auc: 0.792 - ETA: 6:46 - loss: 0.3586 - roc_auc: 0.795 - ETA: 6:47 - loss: 0.3599 - roc_auc: 0.796 - ETA: 6:43 - loss: 0.3599 - roc_auc: 0.795 - ETA: 6:46 - loss: 0.3583 - roc_auc: 0.798 - ETA: 6:36 - loss: 0.3561 - roc_auc: 0.800 - ETA: 6:27 - loss: 0.3571 - roc_auc: 0.799 - ETA: 6:20 - loss: 0.3575 - roc_auc: 0.796 - ETA: 6:11 - loss: 0.3544 - roc_auc: 0.796 - ETA: 6:03 - loss: 0.3558 - roc_auc: 0.794 - ETA: 5:55 - loss: 0.3562 - roc_auc: 0.794 - ETA: 5:47 - loss: 0.3572 - roc_auc: 0.794 - ETA: 5:41 - loss: 0.3578 - roc_auc: 0.793 - ETA: 5:35 - loss: 0.3580 - roc_auc: 0.793 - ETA: 5:28 - loss: 0.3593 - roc_auc: 0.792 - ETA: 5:20 - loss: 0.3578 - roc_auc: 0.792 - ETA: 5:12 - loss: 0.3579 - roc_auc: 0.790 - ETA: 5:05 - loss: 0.3583 - roc_auc: 0.789 - ETA: 4:57 - loss: 0.3576 - roc_auc: 0.790 - ETA: 4:51 - loss: 0.3573 - roc_auc: 0.790 - ETA: 4:44 - loss: 0.3569 - roc_auc: 0.790 - ETA: 4:37 - loss: 0.3568 - roc_auc: 0.790 - ETA: 4:30 - loss: 0.3572 - roc_auc: 0.790 - ETA: 4:23 - loss: 0.3577 - roc_auc: 0.790 - ETA: 4:16 - loss: 0.3581 - roc_auc: 0.790 - ETA: 4:09 - loss: 0.3587 - roc_auc: 0.789 - ETA: 4:02 - loss: 0.3590 - roc_auc: 0.787 - ETA: 3:56 - loss: 0.3595 - roc_auc: 0.787 - ETA: 3:49 - loss: 0.3591 - roc_auc: 0.787 - ETA: 3:44 - loss: 0.3597 - roc_auc: 0.788 - ETA: 3:37 - loss: 0.3591 - roc_auc: 0.789 - ETA: 3:30 - loss: 0.3588 - roc_auc: 0.789 - ETA: 3:23 - loss: 0.3597 - roc_auc: 0.789

auc: 0.789 - ETA: 3:18 - loss: 0.3602 - roc_auc: 0.788 - ETA: 3:12 - loss: 0.3610 - roc_auc: 0.788 - ETA: 3:05 - loss: 0.3617 - roc_auc: 0.787 - ETA: 2:58 - loss: 0.3622 - roc_auc: 0.787 - ETA: 2:52 - loss: 0.3618 - roc_auc: 0.787 - ETA: 2:45 - loss: 0.3612 - roc_auc: 0.787 - ETA: 2:38 - loss: 0.3617 - roc_auc: 0.787 - ETA: 2:32 - loss: 0.3619 - roc_auc: 0.787 - ETA: 2:26 - loss: 0.3616 - roc_auc: 0.787 - ETA: 2:20 - loss: 0.3612 - roc_auc: 0.787 - ETA: 2:13 - loss: 0.3624 - roc_auc: 0.786 - ETA: 2:07 - loss: 0.3620 - roc_auc: 0.787 - ETA: 2:00 - loss: 0.3620 - roc_auc: 0.786 - ETA: 1:54 - loss: 0.3622 - roc_auc: 0.786 - ETA: 1:47 - loss: 0.3618 - roc_auc: 0.787 - ETA: 1:41 - loss: 0.3622 - roc_auc: 0.786 - ETA: 1:34 - loss: 0.3621 - roc_auc: 0.786 - ETA: 1:28 - loss: 0.3622 - roc_auc: 0.786 - ETA: 1:21 - loss: 0.3623 - roc_auc: 0.786 - ETA: 1:14 - loss: 0.3621 - roc_auc: 0.786 - ETA: 1:08 - loss: 0.3624 - roc_auc: 0.786 - ETA: 1:01 - loss: 0.3624 - roc_auc: 0.786 - ETA: 55s - loss: 0.3622 - roc_auc: 0.787 - ETA: 48s - loss: 0.3623 - roc_auc: 0.78 - ETA: 41s - loss: 0.3617 - roc_auc: 0.78 - ETA: 35s - loss: 0.3616 - roc_auc: 0.78 - ETA: 28s - loss: 0.3615 - roc_auc: 0.78 - ETA: 21s - loss: 0.3619 - roc_auc: 0.78 - ETA: 15s - loss: 0.3623 - roc_auc: 0.78 - ETA: 8s - loss: 0.3626 - roc_auc: 0.7870 - ETA: 1s - loss: 0.3625 - roc_auc: 0.787 - 500s 7ms/step - loss: 0.3623 - roc_auc: 0.7871 - val_loss: 0.3972 - val_roc_auc: 0.7206

Epoch 00008: val_roc_auc did not improve from 0.73612

Epoch 9/10

69918/69918 [=====] - ETA: 6:58 - loss: 0.3156 - roc_auc: 0.813 - ETA: 6:47 - loss: 0.3246 - roc_auc: 0.808 - ETA: 6:41 - loss: 0.3450 - roc_auc: 0.798 - ETA: 6:35 - loss: 0.3522 - roc_auc: 0.795 - ETA: 6:29 - loss: 0.3485 - roc_auc: 0.798 - ETA: 6:30 - loss: 0.3490 - roc_auc: 0.799 - ETA: 6:24 - loss: 0.3518 - roc_auc: 0.804 - ETA: 6:22 - loss: 0.3525 - roc_auc: 0.804 - ETA: 6:20 - loss: 0.3510 - roc_auc: 0.804 - ETA: 6:14 - loss: 0.3541 - roc_auc: 0.802 - ETA: 6:12 - loss: 0.3528 - roc_auc: 0.802 - ETA: 6:14 - loss: 0.3536 - roc_auc: 0.803 - ETA: 6:08 - loss: 0.3523 - roc_auc: 0.804 - ETA: 6:03 - loss: 0.3545 - roc_auc: 0.801 - ETA: 5:59 - loss: 0.3563 - roc_auc: 0.801 - ETA: 5:59 - loss: 0.3564 - roc_auc: 0.800 - ETA: 5:51 - loss: 0.3555 - roc_auc: 0.800 - ETA: 5:43 - loss: 0.3548 - roc_auc: 0.801 - ETA: 5:38 - loss: 0.3549 - roc_auc: 0.801 - ETA: 5:30 - loss: 0.3547 - roc_auc: 0.802 - ETA: 5:24 - loss: 0.3563 - roc_auc: 0.801 - ETA: 5:16 - loss: 0.3565 - roc_auc: 0.800 - ETA: 5:13 - loss: 0.3572 - roc_auc: 0.800 - ETA: 5:07 - loss: 0.3563 - roc_auc: 0.801 - ETA: 5:00 - loss: 0.3558 - roc_auc: 0.801 - ETA: 4:52 - loss: 0.3547 - roc_auc: 0.800 - ETA: 4:44 - loss: 0.3549 - roc_auc: 0.800 - ETA: 4:36 - loss: 0.3547 - roc_auc: 0.800 - ETA: 4:28 - loss: 0.3541 - roc_auc: 0.800 - ETA: 4:20 - loss: 0.3536 - roc_auc: 0.801 - ETA: 4:13 - loss: 0.3530 - roc_auc: 0.801 - ETA: 4:06 - loss: 0.3524 - roc_auc: 0.801 - ETA: 3:58 - loss: 0.3511 - roc_auc: 0.802 - ETA: 3:51 - loss: 0.3526 - roc_auc: 0.801 - ETA: 3:44 - loss: 0.3519 - roc_auc: 0.802 - ETA: 3:37 - loss: 0.3518 - roc_auc: 0.802 - ETA: 3:30 - loss: 0.3524 - roc_auc: 0.802 - ETA: 3:22 - loss: 0.3522 - roc_auc: 0.802 - ETA: 3:15 - loss: 0.3527 - roc_auc: 0.801 - ETA: 3:08 - loss: 0.3530 - roc_auc: 0.801 - ETA: 3:01 - loss: 0.3525 - roc_auc: 0.801 - ETA: 2:55 - loss: 0.3523 - roc_auc: 0.801 - ETA: 2:48 - loss: 0.3530 - roc_auc: 0.801 - ETA: 2:41 - loss: 0.3533 - roc_auc: 0.801 - ETA: 2:34 - loss: 0.3529 - roc_auc: 0.802 - ETA: 2:29 - loss: 0.3527 - roc_auc: 0.801 - ETA: 2:22 - loss: 0.3527 - roc_auc: 0.801 - ETA: 2:15 - loss: 0.3523 - roc_auc: 0.802 - ETA: 2:09 - loss: 0.3524 - roc_auc: 0.802 - ETA: 2:02 - loss: 0.3520 - roc_auc: 0.802 - ETA: 1:56 - loss: 0.3518 - roc_auc: 0.802 - ETA: 1:49 - loss: 0.3518 - roc_auc: 0.802 - ETA: 1:42 - loss: 0.3523 - roc_auc: 0.801 - ETA: 1:35 - loss: 0.3524 - roc_auc: 0.802 - ETA: 1:29 - loss: 0.3524 - roc_auc: 0.802 - ETA: 1:22 - loss: 0.3524 - roc_auc: 0.802 - ETA: 1:15 - loss: 0.3523 - roc_auc: 0.802 - ETA: 1:08 - loss: 0.3525 - roc_auc: 0.801 - ETA: 1:02 - loss: 0.3527 - roc_auc: 0.801 - ETA: 55s - loss: 0.3530 - roc_auc: 0.802 - ETA: 49s - loss: 0.3528 - roc_auc: 0.80 - ETA: 42s - loss: 0.3528 - roc_auc: 0.80 -

ETA: 35s - loss: 0.3529 - roc_auc: 0.80 - ETA: 28s - loss: 0.3529 - roc_auc: 0.80 - ETA: 22s - loss: 0.3530 - roc_auc: 0.80 - ETA: 15s - loss: 0.3531 - roc_auc: 0.80 - ETA: 8s - loss: 0.3535 - roc_auc: 0.8024 - ETA: 1s - loss: 0.3536 - roc_auc: 0.802 - 514s 7ms/step - loss: 0.3537 - roc_auc: 0.8022 - val_loss: 0.4079 - val_roc_auc: 0.7116

Epoch 00009: val_roc_auc did not improve from 0.73612

Epoch 10/10

69918/69918 [=====] - ETA: 8:58 - loss: 0.3294 - roc_auc: 0.838 - ETA: 8:18 - loss: 0.3302 - roc_auc: 0.832 - ETA: 8:19 - loss: 0.3295 - roc_auc: 0.831 - ETA: 8:06 - loss: 0.3252 - roc_auc: 0.830 - ETA: 7:43 - loss: 0.3253 - roc_auc: 0.827 - ETA: 7:34 - loss: 0.3295 - roc_auc: 0.829 - ETA: 7:16 - loss: 0.3276 - roc_auc: 0.832 - ETA: 7:03 - loss: 0.3285 - roc_auc: 0.829 - ETA: 6:52 - loss: 0.3287 - roc_auc: 0.831 - ETA: 6:42 - loss: 0.3310 - roc_auc: 0.833 - ETA: 6:31 - loss: 0.3350 - roc_auc: 0.829 - ETA: 6:21 - loss: 0.3360 - roc_auc: 0.827 - ETA: 6:16 - loss: 0.3353 - roc_auc: 0.827 - ETA: 6:11 - loss: 0.3368 - roc_auc: 0.826 - ETA: 6:01 - loss: 0.3372 - roc_auc: 0.827 - ETA: 5:53 - loss: 0.3379 - roc_auc: 0.827 - ETA: 5:45 - loss: 0.3392 - roc_auc: 0.827 - ETA: 5:39 - loss: 0.3404 - roc_auc: 0.826 - ETA: 5:34 - loss: 0.3396 - roc_auc: 0.826 - ETA: 5:26 - loss: 0.3400 - roc_auc: 0.826 - ETA: 5:18 - loss: 0.3389 - roc_auc: 0.827 - ETA: 5:10 - loss: 0.3388 - roc_auc: 0.827 - ETA: 5:02 - loss: 0.3380 - roc_auc: 0.826 - ETA: 4:55 - loss: 0.3381 - roc_auc: 0.826 - ETA: 4:47 - loss: 0.3371 - roc_auc: 0.826 - ETA: 4:44 - loss: 0.3380 - roc_auc: 0.826 - ETA: 4:42 - loss: 0.3377 - roc_auc: 0.827 - ETA: 4:36 - loss: 0.3383 - roc_auc: 0.826 - ETA: 4:30 - loss: 0.3400 - roc_auc: 0.825 - ETA: 4:22 - loss: 0.3405 - roc_auc: 0.824 - ETA: 4:16 - loss: 0.3403 - roc_auc: 0.824 - ETA: 4:08 - loss: 0.3409 - roc_auc: 0.823 - ETA: 4:03 - loss: 0.3407 - roc_auc: 0.823 - ETA: 3:57 - loss: 0.3408 - roc_auc: 0.823 - ETA: 3:50 - loss: 0.3397 - roc_auc: 0.824 - ETA: 3:43 - loss: 0.3400 - roc_auc: 0.824 - ETA: 3:36 - loss: 0.3398 - roc_auc: 0.825 - ETA: 3:30 - loss: 0.3399 - roc_auc: 0.825 - ETA: 3:25 - loss: 0.3400 - roc_auc: 0.824 - ETA: 3:19 - loss: 0.3397 - roc_auc: 0.824 - ETA: 3:13 - loss: 0.3407 - roc_auc: 0.824 - ETA: 3:06 - loss: 0.3405 - roc_auc: 0.824 - ETA: 2:59 - loss: 0.3410 - roc_auc: 0.824 - ETA: 2:53 - loss: 0.3412 - roc_auc: 0.824 - ETA: 2:46 - loss: 0.3406 - roc_auc: 0.824 - ETA: 2:38 - loss: 0.3405 - roc_auc: 0.824 - ETA: 2:31 - loss: 0.3409 - roc_auc: 0.824 - ETA: 2:24 - loss: 0.3420 - roc_auc: 0.822 - ETA: 2:18 - loss: 0.3424 - roc_auc: 0.822 - ETA: 2:12 - loss: 0.3427 - roc_auc: 0.821 - ETA: 2:05 - loss: 0.3432 - roc_auc: 0.821 - ETA: 1:58 - loss: 0.3433 - roc_auc: 0.821 - ETA: 1:51 - loss: 0.3436 - roc_auc: 0.821 - ETA: 1:44 - loss: 0.3436 - roc_auc: 0.820 - ETA: 1:37 - loss: 0.3431 - roc_auc: 0.821 - ETA: 1:30 - loss: 0.3438 - roc_auc: 0.820 - ETA: 1:23 - loss: 0.3434 - roc_auc: 0.821 - ETA: 1:15 - loss: 0.3433 - roc_auc: 0.820 - ETA: 1:08 - loss: 0.3430 - roc_auc: 0.820 - ETA: 1:01 - loss: 0.3427 - roc_auc: 0.821 - ETA: 53s - loss: 0.3428 - roc_auc: 0.820 - ETA: 46s - loss: 0.3434 - roc_auc: 0.82 - ETA: 38s - loss: 0.3434 - roc_auc: 0.82 - ETA: 31s - loss: 0.3434 - roc_auc: 0.82 - ETA: 24s - loss: 0.3438 - roc_auc: 0.82 - ETA: 16s - loss: 0.3436 - roc_auc: 0.81 - ETA: 9s - loss: 0.3436 - roc_auc: 0.8192 - ETA: 2s - loss: 0.3439 - roc_auc: 0.819 - 553s 8ms/step - loss: 0.3437 - roc_auc: 0.8194 - val_loss: 0.4219 - val_roc_auc: 0.7169

Epoch 00010: val_roc_auc did not improve from 0.73612

Out[142]: <keras.callbacks.History at 0x28255b75e48>

```
In [0]: 1 !pip -- upgrade tensorflow
```

ERROR: unknown command "upgrade"

```
In [0]: 1 import tensorflow as tf
2 from tensorflow.python.keras.callbacks import TensorBoard
3 import datetime
4 from keras.callbacks import ModelCheckpoint
```

```
In [0]: 1 tensorboard = TensorBoard("logs\\fit\\" + datetime.datetime.now().strftime("%
```

```
In [0]: 1 % load_ext tensorboard
2 % tensorboard --logdir logs/fit
```

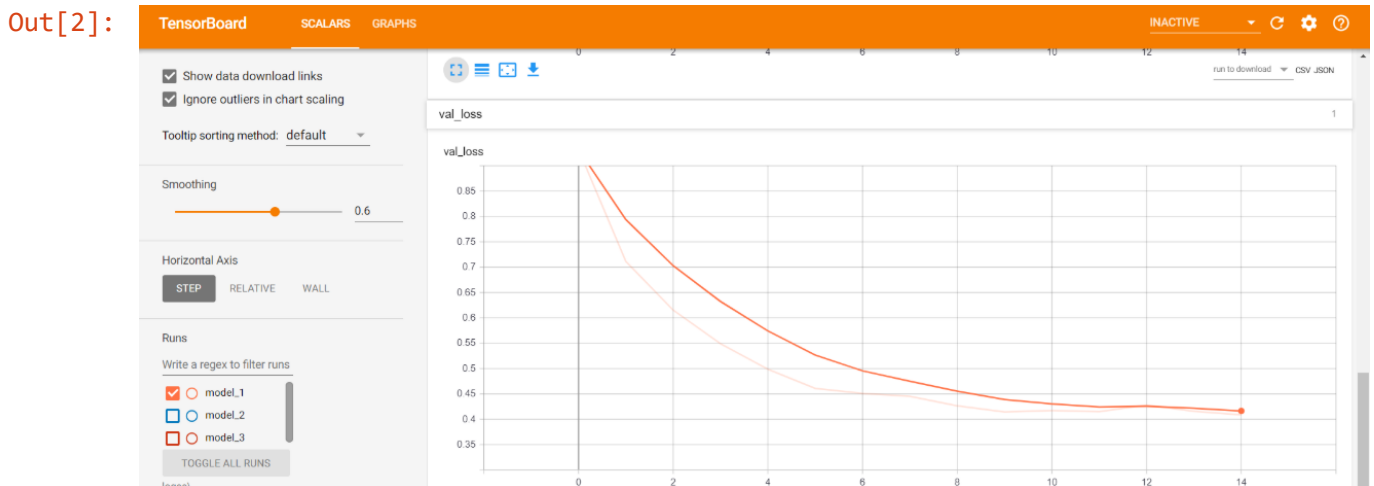
The tensorboard module is not an IPython extension.

UsageError: Line magic function `%tensorboard` not found.

```
In [0]: 1 tf.__version__
```

Out[22]: '1.11.0'

```
In [2]: 1 from IPython.display import Image
2 Image("Capture2.PNG")
```



```
In [0]: 1
```

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. Train the TF-IDF on the Train data

2. Get the idf value for each word we have in the train data.
3. Remove the low idf value and high idf value words from our data. Do some analysis on the Idf values and based on those values choose the low and high threshold value. Because very frequent words and very very rare words don't give much information. (you can plot a box plots and take only the idf scores within IQR range and corresponding words)
4. Train the LSTM after removing the Low and High idf value words. (In model-1 Train on total data but in Model-2 train on data after removing some words based on IDF values)

```
In [0]: 1 import tensorflow as tf
        2 tf.__version__
```

The default version of TensorFlow in Colab will soon switch to TensorFlow 2.x.

We recommend you [upgrade](https://www.tensorflow.org/guide/migrate) (https://www.tensorflow.org/guide/migrate) now or ensure your notebook will continue to use TensorFlow 1.x via the %tensorflow_version 1.x magic: [more info](https://colab.research.google.com/notebooks/tensorflow_version.ipynb) (https://colab.research.google.com/notebooks/tensorflow_version.ipynb).

Out[3]: '1.15.0'

```
In [0]: 1 project_data = pd.read_csv('drive/My Drive/lstm_donor/processed_data.csv')
        2 #target = project_data['']
```

```
In [0]: 1 project_data.columns
        2 target = project_data['project_is_approved']
        3 #project_data.drop(['Unnamed :0', 'id', 'teacher_id', 'project_is_approved'], axis=1)
```

```
In [0]: 1 project_data.drop(['Unnamed: 0', 'id', 'teacher_id', 'project_is_approved'], axis=1)
```

```
In [7]: 1 #splitting the data in train, validation and test
        2 from sklearn.model_selection import train_test_split
        3 #80% training data and 20% test data
        4 X_train, X_test, Y_train, Y_test = train_test_split(project_data, target, test_size=0.2)
        5 X_train, X_cv, Y_train, Y_cv = train_test_split(X_train, Y_train, test_size=0.2)
        6
        7 print('shape of training data after vectorization is:', X_train.shape)
        8 print('shape of validation data after vectorization is:', X_cv.shape)
        9 print('shape of test data after vectorization is:', X_test.shape)
```

shape of training data after vectorization is: (69918, 14)
 shape of validation data after vectorization is: (17480, 14)
 shape of test data after vectorization is: (21850, 14)

```
In [0]: 1
```


TFIDF Model

In [8]: 1 project_data.columns

Out[8]: Index(['teacher_prefix', 'school_state', 'project_submitted_datetime', 'project_grade_category', 'project_subject_categories', 'project_subject_subcategories', 'teacher_number_of_previously_posted_projects', 'price', 'quantity', 'cleaned_project_titles', 'cleaned_resource_summary', 'essay', 'cleaned_essay', 'nrm_price'], dtype='object')

```
In [0]: 1 #utility function
2 def final_text(df):
3     """Returns the final concatenated text"""
4     df['final_Text'] = df['cleaned_essay'].map(str) + " " + df['cleaned_resource_summary'].map(str) + " " + df['cleaned_project_titles'].map(str)
5     return df
6
7
8 #getting the train, val and test data
9 X_train = final_text(X_train)
10 X_cv = final_text(X_cv)
11 X_test = final_text(X_test)
```

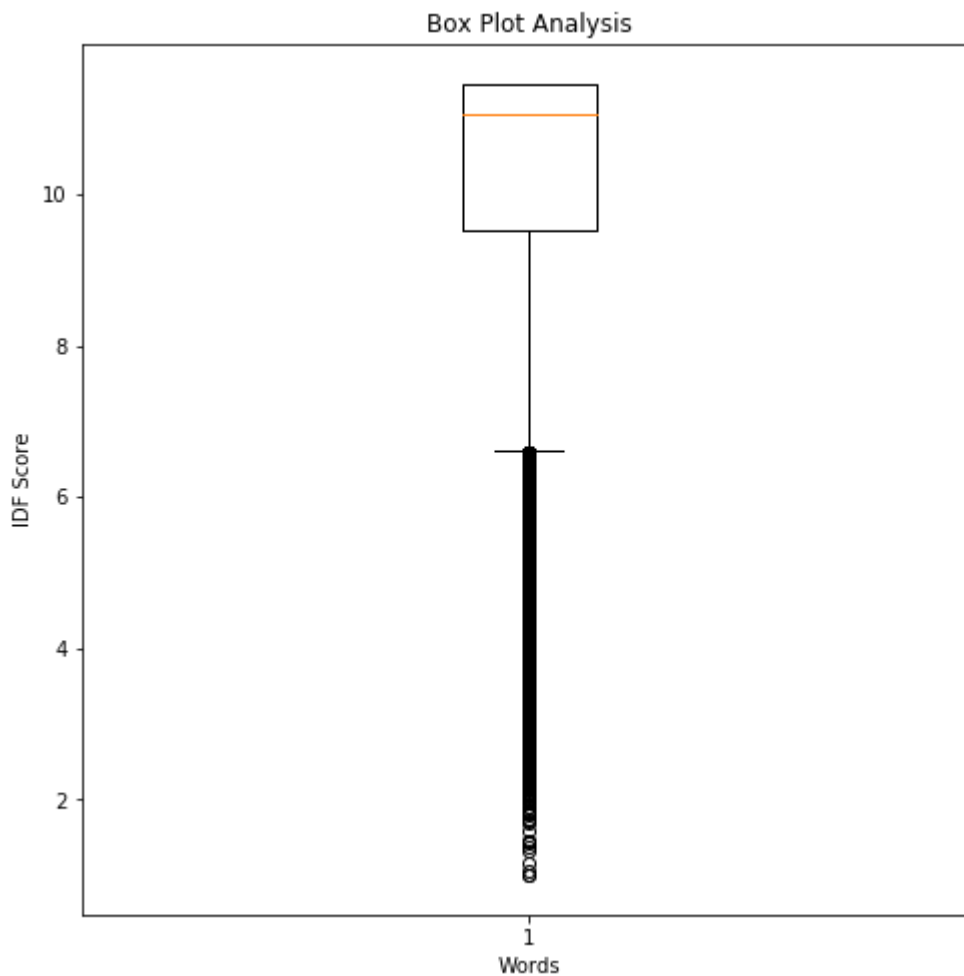
In [10]: 1 X_train.columns

Out[10]: Index(['teacher_prefix', 'school_state', 'project_submitted_datetime', 'project_grade_category', 'project_subject_categories', 'project_subject_subcategories', 'teacher_number_of_previously_posted_projects', 'price', 'quantity', 'cleaned_project_titles', 'cleaned_resource_summary', 'essay', 'cleaned_essay', 'nrm_price', 'final_Text'], dtype='object')

```
In [0]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 tfidf = TfidfVectorizer()
3
4 tfidf.fit(X_train['final_Text'])
5 corpus = tfidf.get_feature_names() #getting names of all the features
6 idf_corpus = tfidf.idf_ #getting the idf scores
7
8 dict_corpus_idf = dict(zip(corpus, idf_corpus))
9 print('Number of words in the corpus are:', len(dict_corpus_idf))
```

Number of words in the corpus are: 51124

```
In [0]: 1 p_75th = np.percentile(idf_corpus,59.81896)
2 p_75th
3 import matplotlib.pyplot as plt
4 plt.figure(figsize=(8,8))
5 plt.boxplot(idf_corpus)
6 plt.title('Box Plot Analysis')
7 plt.xlabel('Words')
8 plt.ylabel('IDF Score')
9 plt.show()
10
11 p_25th = 3
12 p_75th = np.percentile(idf_corpus,75)
13 #p_75th = 11.555
14
15 print("The lowest significant value of TF-IDF Scores: ",p_25th)
16 print("The highest significant value of TF-IDF Scores: ",p_75th)
```



The lowest significant value of TF-IDF Scores: 3

The highest significant value of TF-IDF Scores: 11.4619455276077



```
In [0]: 1 #creating the list of words to be removed from the corpus
2 removed_wordlist = []
3 for word in list(dict_corpus_idf.keys()):
4     if(dict_corpus_idf[word] < p_25th or dict_corpus_idf[word] > p_75th):
5         removed_wordlist.append(word)
6     else:
7         continue
8
9 print("Number of words to be removed: ",len(removed_wordlist))
```

Number of words to be removed: 141

```
In [0]: 1 def remove_from_text(list_of_sentences):
2     """This function will be used to remove words from text data"""
3     processed_text = []
4     for sentence in tqdm(list_of_sentences):
5         sent = ' '.join(word for word in sentence.split() if word not in removed_wordlist)
6         processed_text.append(sent)
7     return processed_text
8
9 X_train['total_text'] = remove_from_text(X_train.final_Text.values)
10 X_cv['total_text'] = remove_from_text(X_cv.final_Text.values)
11 X_test['total_text'] = remove_from_text(X_test.final_Text.values)
```

```
100%|██████████| 69918/69918 [00:15<00:00, 4432.91it/s]
100%|██████████| 17480/17480 [00:03<00:00, 4512.44it/s]
100%|██████████| 21850/21850 [00:04<00:00, 4496.13it/s]
```

```
In [0]: 1 X_train.to_csv("X_train_removed.csv", index=False)
2 X_cv.to_csv("X_val_removed.csv", index=False)
3 X_test.to_csv("X_test_removed.csv", index=False)
```

```
In [0]: 1 X_train = pd.read_csv("X_train_removed.csv")
2 X_val = pd.read_csv("X_val_removed.csv")
3 X_test = pd.read_csv("X_test_removed.csv")
```

Total Text data

```
In [0]: 1 #Get the total_text values in list
2 docs_text_train=list(X_train.total_text.values)
3 docs_text_val=list(X_val.total_text.values)
4 docs_text_test=list(X_test.total_text.values)
5 labels_train=np.array(Y_train)
6 labels_val=np.array(Y_cv)
7 labels_test=np.array(Y_test)
8
9 #Initializing the keras tokenizer and fitting it on train data
10 tokens = Tokenizer()
11 tokens.fit_on_texts(docs_text_train)
12
13 #Convert the texts to sequences using the tokenizer
14 sequences_text_train = tokens.texts_to_sequences(docs_text_train)
15 sequences_text_val = tokens.texts_to_sequences(docs_text_val)
16 sequences_text_test = tokens.texts_to_sequences(docs_text_test)
17 vocab_size_text = len(tokens.word_index) + 1
18
19 #Add padding
20 padded_text_train = pad_sequences(sequences_text_train, maxlen=300, padding='post')
21 padded_text_val = pad_sequences(sequences_text_val, maxlen=300, padding='post')
22 padded_text_test = pad_sequences(sequences_text_test, maxlen=300, padding='post')
```

```
In [0]: 1 pickle_in = open("drive/My Drive/lstm_donor/glove_vectors", "rb")
2 glove_words = pickle.load(pickle_in)
3 embedding_matrix = zeros((vocab_size_text, 300))
4 for word, i in tokens.word_index.items(): #enumerating all unique words
5     embedding_vector = glove_words.get(word)#locating for loaded glove model
6     if embedding_vector is not None:
7         embedding_matrix[i] = embedding_vector
```

```
In [0]: 1 #Get the flattened LSTM output for input text
2 import warnings
3 warnings.filterwarnings('ignore')
4 input_layer_total_text = Input(shape=(300,), name = "total_text_sequence")
5 embedding_layer_total_text = Embedding(input_dim=vocab_size_text, output_dim=
6 lstm_total_text = LSTM(16, activation="relu", return_sequences=True)(embeddi
7 flatten_total_text = Flatten()(lstm_total_text)
```

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:541: The name tf.placeholder is deprecated. Please use tf.compat.v1.placeholder instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:66: The name tf.get_default_graph is deprecated. Please use tf.compat.v1.get_default_graph instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:4432: The name tf.random_uniform is deprecated. Please use tf.random.uniform instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:190: The name tf.get_default_session is deprecated. Please use tf.compat.v1.get_default_session instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:197: The name tf.ConfigProto is deprecated. Please use tf.compat.v1.ConfigProto instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:203: The name tf.Session is deprecated. Please use tf.compat.v1.Session instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:207: The name tf.global_variables is deprecated. Please use tf.compat.v1.global_variables instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:216: The name tf.is_variable_initialized is deprecated. Please use tf.compat.v1.is_variable_initialized instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:223: The name tf.variables_initializer is deprecated. Please use tf.compat.v1.variables_initializer instead.

```

In [0]: 1 # teacher prefix
2 #Get the teacher_prefix values
3 docs_teacher_prefix_train=list(X_train.teacher_prefix.values)
4 docs_teacher_prefix_val=list(X_val.teacher_prefix.values)
5 docs_teacher_prefix_test=list(X_test.teacher_prefix.values)
6
7 #Initializing the keras tokenizer and fitting it on train data
8 tokens = Tokenizer()
9 tokens.fit_on_texts(docs_teacher_prefix_train)
10
11 #Convert the school_state to sequences using the tokenizer
12 sequences_teacher_prefix_train = np.array(tokens.texts_to_sequences(docs_tea
13 sequences_teacher_prefix_val = np.array(tokens.texts_to_sequences(docs_tea
14 sequences_teacher_prefix_test = np.array(tokens.texts_to_sequences(docs_tea
15 vocab_size_teacher_prefix = len(tokens.word_index) + 1
16 #Get the flattened output for teacher_prefix
17 input_layer_teacher_prefix = Input(shape=(1,), name = "teacher_prefix")
18 embedding_layer_teacher_prefix = Embedding(input_dim=vocab_size_teacher_prefi
19 flatten_teacher_prefix = Flatten()(embedding_layer_teacher_prefix)

```

Categorical data : school_state

```

In [0]: 1
2 #Get the school state values
3 docs_school_state_train=list(X_train.school_state.values)
4 docs_school_state_val=list(X_val.school_state.values)
5 docs_school_state_test=list(X_test.school_state.values)
6
7 #Initializing the keras tokenizer and fitting it on train data
8 tokens = Tokenizer()
9 tokens.fit_on_texts(docs_school_state_train)
10
11 #Convert the school_state to sequences using the tokenizer
12 sequences_school_train = np.array(tokens.texts_to_sequences(docs_school_state
13 sequences_school_val = np.array(tokens.texts_to_sequences(docs_school_state_v
14 sequences_school_test = np.array(tokens.texts_to_sequences(docs_school_state_
15 vocab_size_school_state = len(tokens.word_index) + 1
16 #Get the flattened output for school_state
17 input_layer_school_state = Input(shape=(1,), name = "encoded_school_state")
18 embedding_layer_school_state = Embedding(input_dim=vocab_size_school_state, c
19 flatten_school_state = Flatten()(embedding_layer_school_state)
20

```

Categorical data : project_grade_category

```

In [0]: 1 #Get the project_grade_category values
2 docs_project_grade_category_train=list(X_train.project_grade_category.values)
3 docs_project_grade_category_val=list(X_val.project_grade_category.values)
4 docs_project_grade_category_test=list(X_test.project_grade_category.values)
5
6 #Initializing the keras tokenizer and fitting it on train data
7 tokens = Tokenizer()
8 tokens.fit_on_texts(docs_project_grade_category_train)
9
10 #Convert the school_state to sequences using the tokenizer
11 sequences_project_grade_category_train = tokens.texts_to_sequences(docs_proje
12 sequences_project_grade_category_val = tokens.texts_to_sequences(docs_project
13 sequences_project_grade_category_test = tokens.texts_to_sequences(docs_projec
14 vocab_size_project_grade_category= len(tokens.word_index) + 1
15
16 #Add padding
17 padded_project_grade_category_train = pad_sequences(sequences_project_grade_c
18 padded_project_grade_category_val = pad_sequences(sequences_project_grade_cat
19 padded_project_grade_category_test = pad_sequences(sequences_project_grade_ca
20 #Get the flattened output for project_grade_category
21 input_layer_project_grade = Input(shape=(3,), name = "project_grade_category"
22 embedding_layer_project_grade = Embedding(input_dim=vocab_size_project_grade_
23 flatten_project_grade = Flatten()(embedding_layer_project_grade)

```

Categorical data : clean categories

```

In [0]: 1 X_train.columns

```

```

Out[28]: Index(['teacher_prefix', 'school_state', 'project_submitted_datetime',
               'project_grade_category', 'project_subject_categories',
               'project_subject_subcategories',
               'teacher_number_of_previously_posted_projects', 'price', 'quantity',
               'cleaned_project_titles', 'cleaned_resource_summary', 'essay',
               'cleaned_essay', 'nrm_price', 'final_Text', 'total_text'],
              dtype='object')

```

```

In [0]: 1 #Get the clean_categories values
2 docs_clean_categories_train=list(X_train.project_subject_categories.values)
3 docs_clean_categories_val=list(X_val.project_subject_categories.values)
4 docs_clean_categories_test=list(X_test.project_subject_categories.values)
5
6 #Initializing the keras tokenizer and fitting it on train data
7 tokens = Tokenizer()
8 tokens.fit_on_texts(docs_clean_categories_train)
9
10 #Convert the school_state to sequences using the tokenizer
11 sequences_clean_categories_train = tokens.texts_to_sequences(docs_clean_cate
12 sequences_clean_categories_val = tokens.texts_to_sequences(docs_clean_categor
13 sequences_clean_categories_test = tokens.texts_to_sequences(docs_clean_catego
14 vocab_size_clean_categories = len(tokens.word_index) + 1
15
16 #Add padding
17 padded_clean_categories_train = pad_sequences(sequences_clean_categories_train,
18 padded_clean_categories_val = pad_sequences(sequences_clean_categories_val, m
19 padded_clean_categories_test = pad_sequences(sequences_clean_categories_test,
20 #Get the flattened output for clean_categories
21 input_layer_clean_categories = Input(shape=(3,), name = "clean_categories")
22 embedding_layer_clean_categories = Embedding(input_dim=vocab_size_clean_cate
23 flatten_clean_categories = Flatten()(embedding_layer_clean_categories)

```

Categorical data : project subcategories

```

In [0]: 1 #Get the clean_subcategories values
2 docs_clean_subcategories_train=list(X_train.project_subject_subcategories.val
3 docs_clean_subcategories_val=list(X_val.project_subject_subcategories.values)
4 docs_clean_subcategories_test=list(X_test.project_subject_subcategories.value
5
6 #Initializing the keras tokenizer and fitting it on train data
7 tokens = Tokenizer()
8 tokens.fit_on_texts(docs_clean_subcategories_train)
9
10 #Convert the school_state to sequences using the tokenizer
11 sequences_clean_subcategories_train = tokens.texts_to_sequences(docs_clean_su
12 sequences_clean_subcategories_val = tokens.texts_to_sequences(docs_clean_subc
13 sequences_clean_subcategories_test = tokens.texts_to_sequences(docs_clean_sub
14 vocab_size_clean_subcategories = len(tokens.word_index) + 1
15
16 padded_clean_subcategories_train = pad_sequences(sequences_clean_subcategorie
17 padded_clean_subcategories_val = pad_sequences(sequences_clean_subcategories_
18 padded_clean_subcategories_test = pad_sequences(sequences_clean_subcategories
19 #Get the flattened output for clean_subcategories
20 input_layer_clean_subcategories = Input(shape=(3,), name = "clean_subcategori
21 embedding_layer_clean_subcategories = Embedding(input_dim=vocab_size_clean_su
22 flatten_clean_subcategories = Flatten()(embedding_layer_clean_subcategories)

```



```
In [0]: 1 from sklearn.preprocessing import Normalizer
2
3 def normalize_vars(data):
4     """This function is used to normalize all the input datas between 0 and 1
5     normalizer = Normalizer()
6     data_normalized = normalizer.fit_transform(data.reshape(1, -1))
7     return data_normalized, normalizer
```

```
In [0]: 1 previous_projects_train = X_train.teacher_number_of_previously_posted_project
2 previous_projects_val = X_val.teacher_number_of_previously_posted_projects.va
3 previous_projects_test = X_test.teacher_number_of_previously_posted_projects.
4
5 norm_previous_projects_train, normalizer = normalize_vars(previous_projects_t
6 norm_previous_projects_val = normalizer.transform(previous_projects_val.resha
7 norm_previous_projects_test = normalizer.transform(previous_projects_test.res
8
9 norm_previous_projects_train = norm_previous_projects_train.reshape(len(X_tra
10 norm_previous_projects_val = norm_previous_projects_val.reshape(len(X_val),1)
11 norm_previous_projects_test = norm_previous_projects_test.reshape(len(X_test)
12 #Input Layer for teacher_number_of_previously_posted_projects
13 input_layer_previous_projects = Input(shape=(1,), name = "previous_projects")
14 #price
15 #Building train, test and validation data
16 price_train = X_train.price.values
17 price_val = X_val.price.values
18 price_test = X_test.price.values
19
20 norm_price_train, normalizer = normalize_vars(price_train.reshape(1,-1))
21 norm_price_val = normalizer.transform(price_val.reshape(1,-1))
22 norm_price_test = normalizer.transform(price_test.reshape(1,-1))
23
24 norm_price_train = norm_price_train.reshape(len(X_train),1)
25 norm_price_val = norm_price_val.reshape(len(X_val),1)
26 norm_price_test = norm_price_test.reshape(len(X_test),1)
27 #Input Layer for price
28 input_layer_price = Input(shape=(1,), name = "price")
29 #quantity
30 #Building train and validation data
31 quantity_train = X_train.quantity.values
32 quantity_val = X_val.quantity.values
33 quantity_test = X_test.quantity.values
34
35 norm_quantity_train, normalizer = normalize_vars(quantity_train.reshape(1,-1)
36 norm_quantity_val = normalizer.transform(quantity_val.reshape(1,-1))
37 norm_quantity_test = normalizer.transform(quantity_test.reshape(1,-1))
38
39 norm_quantity_train = norm_quantity_train.reshape(len(X_train),1)
40 norm_quantity_val = norm_quantity_val.reshape(len(X_val),1)
41 norm_quantity_test = norm_quantity_test.reshape(len(X_test),1)
42 #Input Layer for quantity
43 input_layer_quantity = Input(shape=(1,), name = "quantity")
```

```
In [0]: 1 # concatenating all
        2 numerical_features_layers_concat = concatenate([input_layer_previous_projects
        3 dense_layer_numerical = Dense(6, activation='relu',kernel_initializer='he_nor
```

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:4479: The name tf.truncated_normal is deprecated. Please use tf.random.truncated_normal instead.

```
In [0]: 1 # Merge all the layers according to the architecture diagram
2 x = concatenate([flatten_total_text, flatten_teacher_prefix, flatten_school_s
3 x = Dense(10, activation='relu', kernel_initializer='he_normal', name='dense_la
4 x = Dropout(0.3, name='dropout_1')(x)
5 x = Dense(10, activation='relu', kernel_initializer='he_normal', name='dense_la
6 x = Dropout(0.3, name='dropout_2')(x)
7 output_layer = Dense(1, activation='sigmoid', name='output_layer')(x)
8
9 # Final model
10 model = Model(inputs=[input_layer_total_text, input_layer_teacher_prefix, input
11                     input_layer_clean_subcategories, input_layer_previous_pr
12 model.summary()
```

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:148: The name tf.placeholder_with_default is deprecated. Please use tf.compat.v1.placeholder_with_default instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:3733: calling dropout (from tensorflow.python.ops.nn_ops) with keep_prob is deprecated and will be removed in a future version.

Instructions for updating:

Please use `rate` instead of `keep_prob`. Rate should be set to `rate = 1 - keep_prob`.

Model: "model_1"

| Layer (type) | Output Shape | Param # | Connected to |
|--|------------------|----------|---------------------------|
| ===== | | | |
| total_text_sequence (InputLayer (None, 300)) | | 0 | |
| embedding_1 (Embedding) | (None, 300, 300) | 15303600 | total_text_sequence[0][0] |
| teacher_prefix (InputLayer) | (None, 1) | 0 | |
| encoded_school_state (InputLayer) | (None, 1) | 0 | |
| project_grade_category (InputLayer) | (None, 3) | 0 | |
| clean_categories (InputLayer) | (None, 3) | 0 | |
| clean_subcategories (InputLayer) | (None, 3) | 0 | |
| previous_projects (InputLayer) | (None, 1) | 0 | |
| price (InputLayer) | (None, 1) | 0 | |

| | | | |
|---|-----------------|-------|---|
| quantity (InputLayer) | (None, 1) | 0 | |
| lstm_1 (LSTM) [0] | (None, 300, 16) | 20288 | embedding_1[0] |
| embedding_2 (Embedding) [0][0] | (None, 1, 4) | 24 | teacher_prefix |
| embedding_3 (Embedding) _state[0][0] | (None, 1, 4) | 208 | encoded_school |
| embedding_4 (Embedding) category[0][0] | (None, 3, 4) | 40 | project_grade_ |
| embedding_5 (Embedding) es[0][0] | (None, 3, 4) | 64 | clean_categori |
| embedding_6 (Embedding) ories[0][0] | (None, 3, 4) | 152 | clean_subcateg |
| concatenate_1 (Concatenate) cts[0][0] | (None, 3) | 0 | previous_proje price[0][0] quantity[0][0] |
| flatten_1 (Flatten) | (None, 4800) | 0 | lstm_1[0][0] |
| flatten_2 (Flatten) [0] | (None, 4) | 0 | embedding_2[0] |
| flatten_3 (Flatten) [0] | (None, 4) | 0 | embedding_3[0] |
| flatten_4 (Flatten) [0] | (None, 12) | 0 | embedding_4[0] |
| flatten_5 (Flatten) [0] | (None, 12) | 0 | embedding_5[0] |
| flatten_6 (Flatten) [0] | (None, 12) | 0 | embedding_6[0] |

| | | | |
|------------------------------------|--------------|-------|---------------|
| dense_1 (Dense) [0][0] | (None, 6) | 24 | concatenate_1 |
| concatenate_2 (Concatenate) [0] | (None, 4850) | 0 | flatten_1[0] |
| [0] | | | flatten_2[0] |
| [0] | | | flatten_3[0] |
| [0] | | | flatten_4[0] |
| [0] | | | flatten_5[0] |
| [0] | | | flatten_6[0] |
| [0] | | | dense_1[0][0] |
| dense_layer_1 (Dense) [0][0] | (None, 10) | 48510 | concatenate_2 |
| dropout_1 (Dropout) [0][0] | (None, 10) | 0 | dense_layer_1 |
| dense_layer_2 (Dense) [0] | (None, 10) | 110 | dropout_1[0] |
| dropout_2 (Dropout) [0][0] | (None, 10) | 0 | dense_layer_2 |
| output_layer (Dense) [0] | (None, 1) | 11 | dropout_2[0] |
| ===== | | | |
| ===== | | | |
| Total params: 15,373,031 | | | |
| Trainable params: 69,431 | | | |
| Non-trainable params: 15,303,600 | | | |
| ===== | | | |
| ◀ | | | |
| ▶ | | | |

```
In [0]: 1 from keras import backend as K
2 from sklearn.metrics import roc_auc_score
3 #https://stackoverflow.com/questions/51922500/tf-metrics-auc-yielding-very-di
4 def roc_auc(y_true, y_pred):
5     auc = tf.py_func(roc_auc_score, (y_true, y_pred), tf.double)
6     #auc = tf.metrics.auc(y_true, y_pred, num_thresholds=200)[1]
7     K.get_session().run(tf.local_variables_initializer())
8     return auc
9 #Defining callbacks
10 from time import time
11 #from tensorflow.python.keras.callbacks import TensorBoard
12 from keras.callbacks import ModelCheckpoint
13 from keras.callbacks import TensorBoard
14
15 tensorboard = TensorBoard(log_dir="logs".format(time))
16 filepath="weights_best.hdf5"
17 checkpoint = ModelCheckpoint(filepath, monitor='val_roc_auc', verbose=1, save
18 import gc
19 gc.collect()
```

Out[35]: 5

```
In [0]: 1 import tensorflow as tf
```

```
In [0]: 1 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=[roc_auc]
2 model.fit(x=[padded_text_train, sequences_teacher_prefix_train, sequences_schoc
3 norm_previous_projects_train, norm_price_train, norm_quantity_train],
4 y=[labels_train],
5 validation_data=([padded_text_val, sequences_teacher_prefix_val, sequ
6 norm_previous_projects_val, norm_price_val, norm_qu
7 epochs=7,
8 batch_size=1024,
9 callbacks=[tensorboard]))
```

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/optimizer.py:793: The name tf.train.Optimizer is deprecated. Please use tf.compat.v1.train.Optimizer instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:3657: The name tf.log is deprecated. Please use tf.math.log instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/tensorflow_core/python/ops/nn_impl.py:183: where (from tensorflow.python.ops.array_ops) is deprecated and will be removed in a future version.

Instructions for updating:

Use tf.where in 2.0, which has the same broadcast rule as np.where

WARNING:tensorflow:From <ipython-input-35-b70f758d7380>:5: 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

o

- 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.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:1033: The name tf.assign_add is deprecated. Please use tf.compat.v1.assign_add instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:1020: The name tf.assign is deprecated. Please use tf.compat.v1.assign instead.

Train on 69918 samples, validate on 17480 samples

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/callbacks.py:1122: The name tf.summary.merge_all is deprecated. Please use tf.compat.v1.summary.merge_all instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/callbacks.py:1125: The name tf.summary.FileWriter is deprecated. Please use tf.compat.v1.summary.FileWriter instead.

```

Epoch 1/7
69918/69918 [=====] - 39s 557us/step - loss: 0.4618
- roc_auc: 0.5689 - val_loss: 0.4135 - val_roc_auc: 0.7006
WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/callback
s.py:1265: The name tf.Summary is deprecated. Please use tf.compat.v1.Summary
instead.

Epoch 2/7
69918/69918 [=====] - 36s 516us/step - loss: 0.4259
- roc_auc: 0.6409 - val_loss: 0.3959 - val_roc_auc: 0.7145
Epoch 3/7
69918/69918 [=====] - 36s 518us/step - loss: 0.4117
- roc_auc: 0.6753 - val_loss: 0.3884 - val_roc_auc: 0.7197
Epoch 4/7
69918/69918 [=====] - 37s 524us/step - loss: 0.4012
- roc_auc: 0.6970 - val_loss: 0.3836 - val_roc_auc: 0.7267
Epoch 5/7
69918/69918 [=====] - 36s 513us/step - loss: 0.3962
- roc_auc: 0.7081 - val_loss: 0.3843 - val_roc_auc: 0.7276
Epoch 6/7
69918/69918 [=====] - 36s 516us/step - loss: 0.3916
- roc_auc: 0.7194 - val_loss: 0.3862 - val_roc_auc: 0.7255
Epoch 7/7
69918/69918 [=====] - 36s 508us/step - loss: 0.3871
- roc_auc: 0.7287 - val_loss: 0.3866 - val_roc_auc: 0.7262

```

Out[37]: <keras.callbacks.History at 0x7f6f48f5f1d0>

```
In [0]: 1 !tensorboard --logdir=logs/ --host=127.0.0.1
```

TensorBoard 1.15.0 at <http://127.0.0.1:6006/> (<http://127.0.0.1:6006/>) (Press CT
RL+C to quit)

Auc on test data

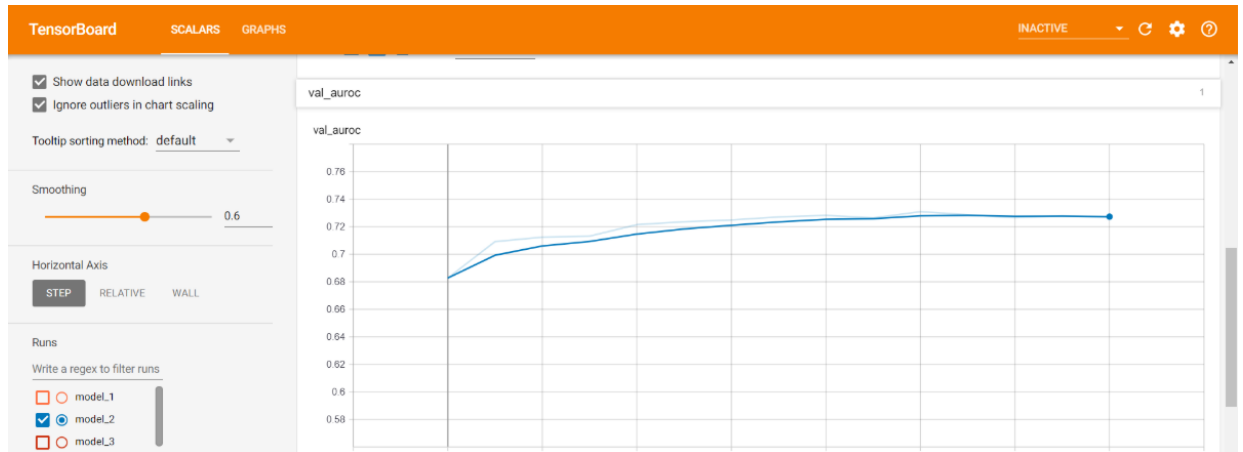
```
In [0]: 1 test_data=[padded_text_test,sequences_teacher_prefix_test,sequences_school_te
2           padded_clean_categories_test,padded_clean_subcategories_test]
```

```
In [0]: 1 #Test AUC
2 y_pred= model.predict(test_data)
3 print("AUC on unseen test data: ",roc_auc_score(Y_test,y_pred))
4
```

AUC on unseen test data: 0.7275554325048124

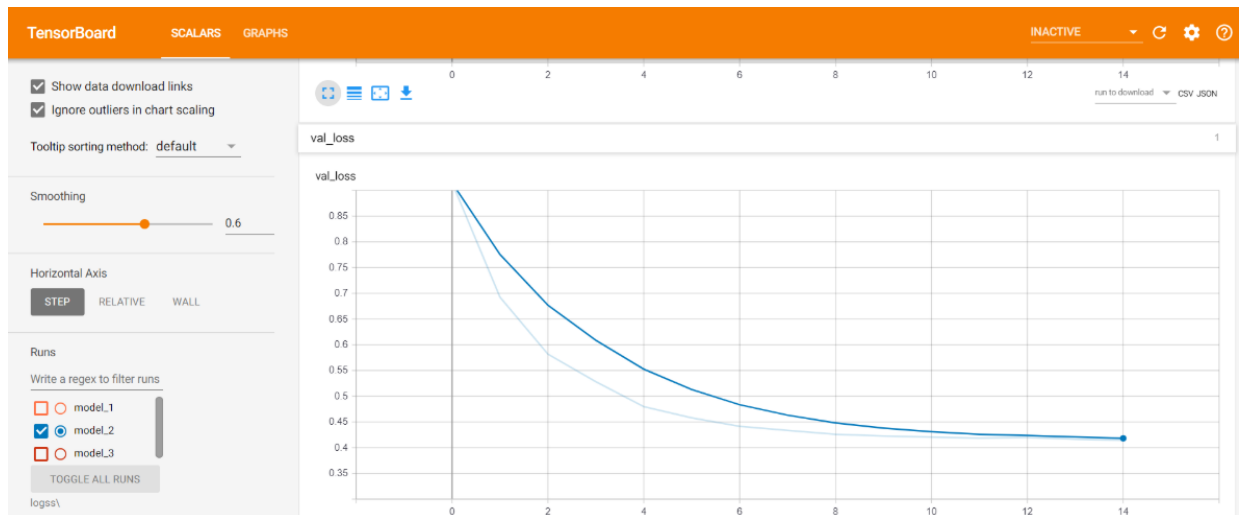

```
In [3]: 1 from IPython.display import Image
        2 Image("Capture3.PNG")
```

Out[3]:

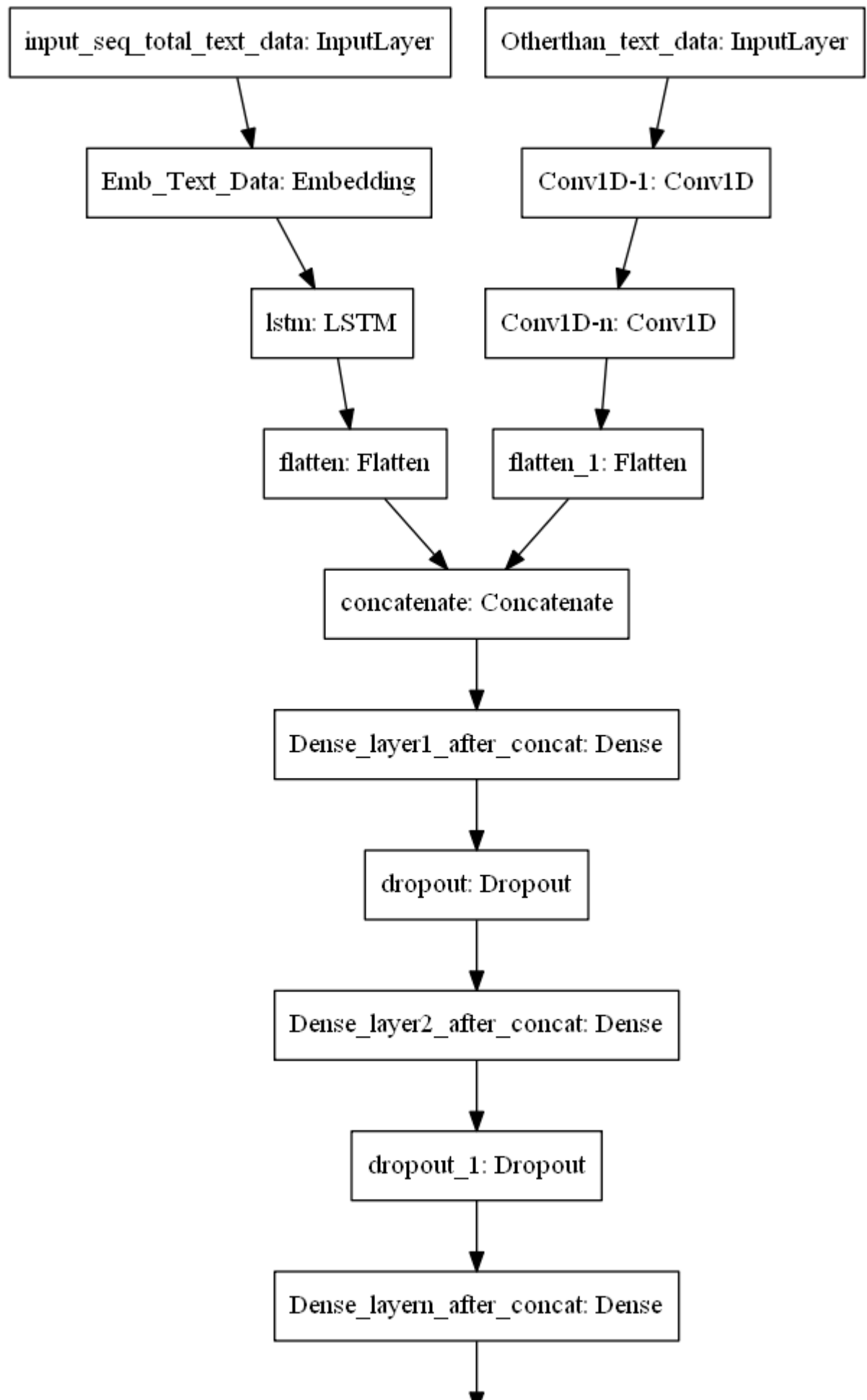


```
In [4]: 1 from IPython.display import Image
        2 Image("Capture4.PNG")
```

Out[4]:



Model-3



output_layer_to_classify_with_soft_max: Dense

ref: <https://i.imgur.com/fkQ8nGo.png> (<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 \(https://keras.io/getting-started/sequential-model-guide/#sequence-classification-with-1d-convolution\)](https://keras.io/getting-started/sequential-model-guide/#sequence-classification-with-1d-convolution) as shown in above figure.
- . You are free to choose all CNN parameters like kernel sizes, stride.

In [0]: 1 X_train.columns

Out[11]: Index(['teacher_prefix', 'school_state', 'project_submitted_datetime', 'project_grade_category', 'project_subject_categories', 'project_subject_subcategories', 'teacher_number_of_previously_posted_projects', 'price', 'quantity', 'cleaned_project_titles', 'cleaned_resource_summary', 'essay', 'cleaned_essay', 'nrm_price', 'final_Text'], dtype='object')

```
In [0]: 1 tokenizer = Tokenizer()
2         tokenizer.fit_on_texts(fil_words)
3         seq_train = tokenizer.texts_to_sequences(X_train["final_Text"])
4         seq_cv = tokenizer.texts_to_sequences(X_cv['final_Text'])
5         seq_test = tokenizer.texts_to_sequences(X_test["final_Text"])
6
7         vocab_size_text = len(tokenizer.word_index)+1 #vocablury size of the data
8
9         #Add padding
10        padded_text_train = pad_sequences(seq_train, maxlen=300, padding='post')
11        padded_text_val = pad_sequences(seq_cv, maxlen=300, padding='post')
12        padded_text_test = pad_sequences(seq_test, maxlen=300, padding='post')
```

```
In [15]: 1 pickle_in = open("drive/My Drive/lstm_donor/glove_vectors", "rb")
2         glove_words = pickle.load(pickle_in)
3         embedding_matrix = zeros((vocab_size_text, 300))
4         for word, i in tokenizer.word_index.items(): #enumerating all unique words
5             embedding_vector = glove_words.get(word) #locating for loaded glove model
6             if embedding_vector is not None:
7                 embedding_matrix[i] = embedding_vector
8
9         print('the shappe of embedding matrix is:', embedding_matrix.shape)
```

the shappe of embedding matrix is: (24730, 300)

```
In [0]: 1 # input_text = Input(shape=(800,), name="input_text")
2         import warnings
3         warnings.filterwarnings('ignore')
4         from keras import regularizers
5         input_text = Input(shape=(300,), name="input_text")
6         embedding_layer = Embedding(vocab_size_text, output_dim = 300, weights=[embeddi
7         x = embedding_layer(input_text)
8         x = LSTM(256, recurrent_dropout=0.5, kernel_regularizer=regularizers.l2(0.001),
9         flat_1 = Flatten()(x)
```

```
In [0]: 1 X_train.columns
```

```
Out[19]: Index(['teacher_prefix', 'school_state', 'project_submitted_datetime',
               'project_grade_category', 'project_subject_categories',
               'project_subject_subcategories',
               'teacher_number_of_previously_posted_projects', 'price', 'quantity',
               'cleaned_project_titles', 'cleaned_resource_summary', 'essay',
               'cleaned_essay', 'nrm_price', 'final_Text'],
              dtype='object')
```

```

In [0]: 1 vect = CountVectorizer(binary=True)
2 vect.fit(X_train["teacher_prefix"])
3 train_prefix = vect.transform(X_train["teacher_prefix"])
4 val_prefix = vect.transform(X_cv["teacher_prefix"])
5 test_prefix = vect.transform(X_test["teacher_prefix"])
6
7 vect = CountVectorizer(binary=True)
8 vect.fit(X_train["school_state"])
9 train_state = vect.transform(X_train["school_state"])
10 val_state = vect.transform(X_cv["school_state"])
11 test_state = vect.transform(X_test["school_state"])
12
13 vect = CountVectorizer(binary=True)
14 vect.fit(X_train["school_state"])
15 train_grade = vect.transform(X_train["project_grade_category"])
16 val_grade = vect.transform(X_cv["project_grade_category"])
17 test_grade = vect.transform(X_test["project_grade_category"])
18
19 vect = CountVectorizer(binary=True)
20 vect.fit(X_train["project_subject_categories"])
21 train_subcat = vect.transform(X_train["project_subject_categories"])
22 val_subcat = vect.transform(X_cv["project_subject_categories"])
23 test_subcat = vect.transform(X_test["project_subject_categories"])
24
25 vect = CountVectorizer(binary=True)
26 vect.fit(X_train["project_subject_subcategories"])
27 train_subcat_1 = vect.transform(X_train["project_subject_subcategories"])
28 val_subcat_1 = vect.transform(X_cv["project_subject_subcategories"])
29 test_subcat_1 = vect.transform(X_test["project_subject_subcategories"])

```

```

In [0]: 1 #for all the numerical features
2
3 #numerical_train_a=X_train['digits_in_summary'].values.reshape(-1, 1)
4 num_train_1=X_train['price'].values.reshape(-1, 1)
5 num_train_2=X_train['quantity'].values.reshape(-1, 1)
6 num_train_3=X_train['teacher_number_of_previously_posted_projects'].values.re
7
8
9 num_val_1=X_cv['price'].values.reshape(-1, 1)
10 num_val_2=X_cv['quantity'].values.reshape(-1, 1)
11 num_val_3=X_cv['teacher_number_of_previously_posted_projects'].values.reshape
12
13 #numerical_test_a=X_test['digits_in_summary'].values.reshape(-1, 1)
14 num_test_1=X_test['price'].values.reshape(-1, 1)
15 num_test_2=X_test['quantity'].values.reshape(-1, 1)
16 num_test_3=X_test['teacher_number_of_previously_posted_projects'].values.res

```

```
In [ ]: 1 numerical_train=np.concatenate((num_train_1,num_train_2,num_train_3),axis=1)
2 numerical_val=np.concatenate((num_val_1,num_val_2,num_val_3),axis=1)
3 numerical_test=np.concatenate((num_test_1,num_test_2,num_test_3),axis=1)
4
5 from sklearn.preprocessing import StandardScaler
6 normal=StandardScaler()
7 normal_train=normal.fit_transform(numerical_train)
8 normal_val=normal.fit_transform(numerical_val)
9 normal_test=normal.transform(numerical_test)
```

```
In [0]: 1 from scipy.sparse import hstack
2 other_train = hstack([train_prefix,train_state,train_grade,train_subcat,train
3 other_val = hstack([val_prefix,val_state,val_grade,val_subcat,val_subcat_1]).
4 other_test = hstack([test_prefix,test_state,test_grade,test_subcat,test_subca
5 other_train = np.hstack((ot_train,normal_train))
6 other_val = np.hstack((ot_val,normal_val))
7 other_test = np.hstack((ot_test,normal_test))
8 other_train_all = np.expand_dims(ot_all_train,2)
9 other_val_all = np.expand_dims(ot_all_val,2)
10 other_test_all = np.expand_dims(ot_all_test,2)
```

```
In [22]: 1 from keras.layers import Input, Embedding, LSTM, Dropout, BatchNormalization,
2 inp_conv=Input(shape=(other_all_train.shape[1], 1))
3 x1 = Conv1D(filters=128, kernel_size=3, activation='relu',kernel_initializer=
4 x1 = Conv1D(filters=128, kernel_size=3, activation='relu',kernel_initializer=
5 x1 = Flatten()(x1)
```

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:4479: The name tf.truncated_normal is deprecated. Please use tf.random.truncated_normal instead.

```
In [0]: 1
2 x_concatenate = concatenate([flat_1,x1])
3 x = Dense(128,activation="relu",kernel_initializer="he_normal",kernel_regularizer
4 x=Dropout(0.5)(x)
5 x = Dense(64,activation="relu",kernel_initializer="he_normal",kernel_regularizer
6 x=Dropout(0.3)(x)
7 x = Dense(32,activation="relu",kernel_initializer="he_normal",kernel_regularizer
8 output = Dense(2, activation='softmax', name='output')(x)
9 model_1 = Model(inputs=[input_text,inp_conv],outputs=[output])
```

```
In [0]: 1 train_data_3 = [padded_text_train,ot_train_all]
2 val_data_3 = [padded_text_val,ot_val_all]
3 test_data_3 = [padded_text_test,ot_test_all]
4 from keras.utils import np_utils
5 Y_train = np_utils.to_categorical(Y_train, 2)
6 Y_val = np_utils.to_categorical(Y_cv, 2)
7 Y_test = np_utils.to_categorical(Y_test, 2)
```

```
In [0]: 1 from sklearn.metrics import roc_auc_score
2 def auc1(y_true, y_pred):
3     if len(np.unique(y_true[:,1])) == 1:
4         return 0.5
5     else:
6         return roc_auc_score(y_true, y_pred)
7 def auROC(y_true, y_pred):
8     return tf.py_func(auc1, (y_true, y_pred), tf.double)
```

```
In [0]: 1 from keras.optimizers import Adam
2 import tensorflow as tf
```

```
In [27]: 1 model_1.compile(optimizer=Adam(lr=0.001, beta_1=0.9, beta_2=0.999, epsilon=No
```

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/optimizer_s.py:793: The name tf.train.Optimizer is deprecated. Please use tf.compat.v1.train.Optimizer instead.

WARNING:tensorflow:From /usr/local/lib/python3.6/dist-packages/keras/backend/tensorflow_backend.py:3576: The name tf.log is deprecated. Please use tf.math.log instead.

WARNING:tensorflow:From <ipython-input-25-b3d027b8e5dd>:8: 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 [0]: 1
2 checkpoint_3 = ModelCheckpoint("model_3.h5",monitor="val_auroc",mode="max",sa
3 NAME = 'model_3'
4 tensorboard_2 = TensorBoard(log_dir='logs')
5 callbacks_2 = [tensorboard_2,checkpoint_3]
```

```
In [37]: 1 history_1 = model_1.fit(train_data_3,Y_train,batch_size=512,epochs=15,validat
```

Train on 69918 samples, validate on 17480 samples

Epoch 1/15

69918/69918 [=====] - 76s 1ms/step - loss: 0.8016 - au
roc: 0.6520 - val_loss: 0.6379 - val_auroc: 0.7352

Epoch 2/15

69918/69918 [=====] - 75s 1ms/step - loss: 0.5536 - au
roc: 0.7216 - val_loss: 0.5570 - val_auroc: 0.7476

Epoch 3/15

69918/69918 [=====] - 74s 1ms/step - loss: 0.4930 - au
roc: 0.7324 - val_loss: 0.5011 - val_auroc: 0.7417

Epoch 4/15

69918/69918 [=====] - 75s 1ms/step - loss: 0.4623 - au
roc: 0.7346 - val_loss: 0.4821 - val_auroc: 0.7493

Epoch 5/15

69918/69918 [=====] - 75s 1ms/step - loss: 0.4403 - au
roc: 0.7460 - val_loss: 0.4691 - val_auroc: 0.7615

Epoch 6/15

69918/69918 [=====] - 75s 1ms/step - loss: 0.4241 - au
roc: 0.7544 - val_loss: 0.4552 - val_auroc: 0.7610

Epoch 7/15

69918/69918 [=====] - 75s 1ms/step - loss: 0.4128 - au
roc: 0.7587 - val_loss: 0.4182 - val_auroc: 0.7620

Epoch 8/15

69918/69918 [=====] - 75s 1ms/step - loss: 0.4047 - au
roc: 0.7627 - val_loss: 0.4197 - val_auroc: 0.7631

Epoch 9/15

69918/69918 [=====] - 76s 1ms/step - loss: 0.4015 - au
roc: 0.7617 - val_loss: 0.4119 - val_auroc: 0.7602

Epoch 10/15

69918/69918 [=====] - 76s 1ms/step - loss: 0.3976 - au
roc: 0.7625 - val_loss: 0.4145 - val_auroc: 0.7628

Epoch 11/15

69918/69918 [=====] - 76s 1ms/step - loss: 0.3933 - au
roc: 0.7635 - val_loss: 0.4053 - val_auroc: 0.7654

Epoch 12/15

69918/69918 [=====] - 78s 1ms/step - loss: 0.3916 - au
roc: 0.7661 - val_loss: 0.4001 - val_auroc: 0.7633

Epoch 13/15

69918/69918 [=====] - 78s 1ms/step - loss: 0.3892 - au
roc: 0.7655 - val_loss: 0.4177 - val_auroc: 0.7619

Epoch 14/15

69918/69918 [=====] - 77s 1ms/step - loss: 0.3866 - au
roc: 0.7691 - val_loss: 0.3923 - val_auroc: 0.7681

Epoch 15/15

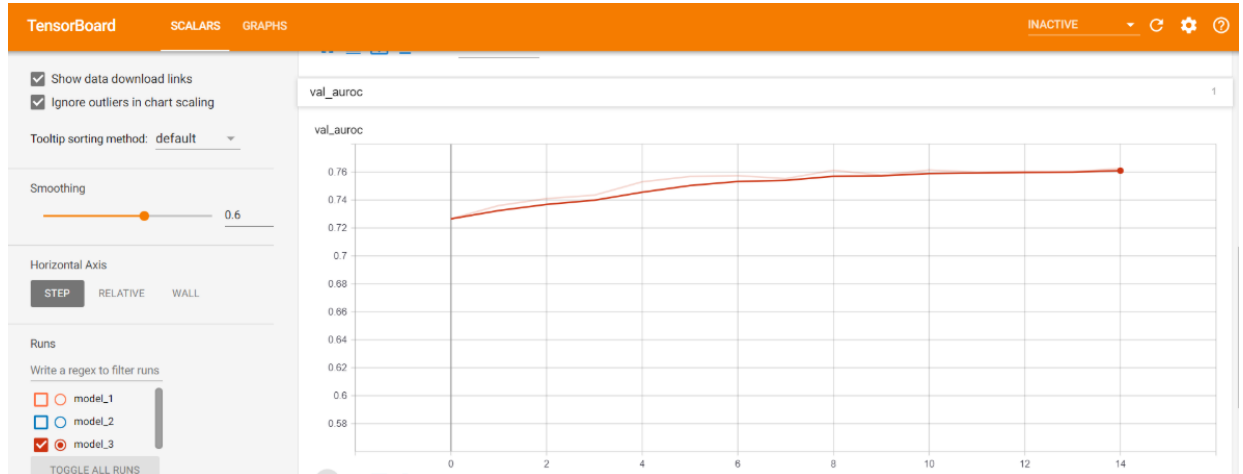
69918/69918 [=====] - 78s 1ms/step - loss: 0.3855 - au
roc: 0.7678 - val_loss: 0.3983 - val_auroc: 0.7677

```
In [40]: 1 #Test AUC
2 y_pred= model_1.predict(test_data_3)
3 print("AUC on unseen test data: ",roc_auc_score(Y_test,y_pred))
```

AUC on unseen test data: 0.7633947879300654

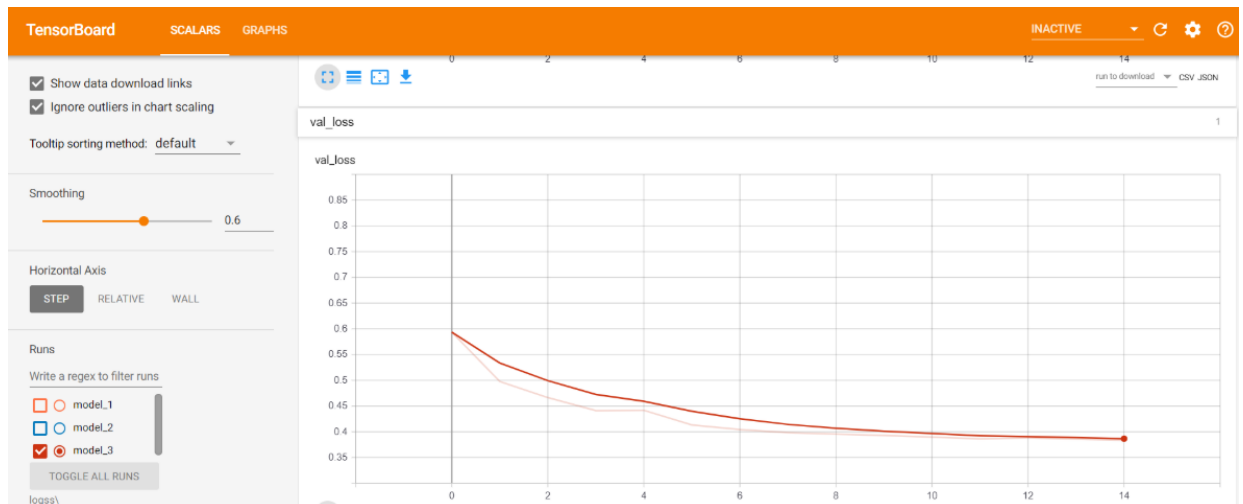

```
In [5]: 1 from IPython.display import Image
        2 Image("Capture5.PNG")
```

Out[5]:



```
In [6]: 1 from IPython.display import Image
        2 Image("Capture6.PNG")
```

Out[6]:



```
In [12]: 1 from prettytable import PrettyTable
2 table = PrettyTable()
3 table.field_names = ["Model No", "Trained for Epochs", "Train ROC-AUC", "Validation ROC-AUC", "Test ROC-AUC"]
4 table.add_row(["Model 1", 10, 0.7988, 0.7169, 0.7148])
5 table.add_row(["Model 2", 7, 0.7287, 0.7262, 0.7275])
6 table.add_row(["Model 3", 15, 0.7678, 0.7677, 0.7633])
7 print(table)
```

```
+-----+-----+-----+-----+-----+
+-----+
| Model No | Trained for Epochs | Train ROC-AUC | Validation ROC-AUC | Test ROC-AUC |
+-----+-----+-----+-----+-----+
+-----+
| Model 1 | 10 | 0.7988 | 0.7169 | 0.7148 |
8 |
| Model 2 | 7 | 0.7287 | 0.7262 | 0.7275 |
5 |
| Model 3 | 15 | 0.7678 | 0.7677 | 0.7633 |
3 |
+-----+-----+-----+-----+-----+
+-----+
```

What was the case study all about and what we learnt

So this case study was an application of machine learning for social good where we were able to automate the process of approving the projects being submitted for evaluation based on lots of different factors.

We tokenized the text data and one hot encoded the categorical variables and also performed normalization on the numerical features primarily, along with that we performed the feature engineering with digits present in the text or not.

We also used pre embedded glove vectors as the main embedding layer for the text data.

Amongst all the three models, the 3rd model with 1D CNN layers seems to perform the best as we have got the maximum value of ROC-AUC for it.

Using RSM Prop with proper weight initialization was also resulting in exploding gradients for Model 1. Changing the optimizer to adam has changed this problem. Got the loss curves and score curve using Tensorboard. Used a custom metric function for training the model with a custom roc-auc score. For TFIDF analysis, I first tried with the 25 percentile threshold for tfidf scores and the model performed very poorly despite trying my best to optimize it. The lowest threshold I have considered for this assignment is 6.5 approximately. Model 3 has given us the best value of ROC-AUC - just under 0.75.

For TFIDF analysis, I had previously tried with IDF score between 7 and max value, but it resulted in huge loss of data. Hence I took IDF values which were greater than 3 and the model has improved a lot.

