

DonorsChoose

DonorsChoose.org receives hundreds of thousands of project proposals each year for classroom projects in need of funding. Right now, a large number of volunteers is needed to manually screen each submission before it's approved to be posted on the DonorsChoose.org website.

Next year, DonorsChoose.org expects to receive close to 500,000 project proposals. As a result, there are three main problems they need to solve:

- How to scale current manual processes and resources to screen 500,000 projects so that they can be posted as quickly and as efficiently as possible
- How to increase the consistency of project vetting across different volunteers to improve the experience for teachers
- How to focus volunteer time on the applications that need the most assistance

The goal of the competition is to predict whether or not a DonorsChoose.org project proposal submitted by a teacher will be approved, using the text of project descriptions as well as additional metadata about the project, teacher, and school. DonorsChoose.org can then use this information to identify projects most likely to need further review before approval.

About the DonorsChoose Data Set

The `train.csv` data set provided by DonorsChoose contains the following features:

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	<ul style="list-style-type: none">••	Title of the project. Examples: Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: Grades PreK-2 Grades 3-5 Grades 6-8 Grades 9-12
<code>project_subject_categories</code>	<ul style="list-style-type: none">•••••••••	One or more (comma-separated) subject categories for the project from the following enumerated list of values: Applied Learning Care & Hunger Health & Sports History & Civics Literacy & Language Math & Science Music & The Arts Special Needs Warmth
	<ul style="list-style-type: none">••	Examples: Music & The Arts Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (Two-letter U.S. postal code (https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes)). Example: WY	
<code>project_subject_subcategories</code>	<ul style="list-style-type: none">••	One or more (comma-separated) subject subcategories for the project. Examples: Literacy Literature & Writing, Social Sciences
<code>project_resource_summary</code>	<ul style="list-style-type: none">•	An explanation of the resources needed for the project. Example: My students need hands on literacy materials to manage sensory needs!

Feature	Description
project_essay_1	First application essay*
project_essay_2	Second application essay*
project_essay_3	Third application essay*
project_essay_4	Fourth application essay*
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <div><ul style="list-style-type: none">nanDr.Mr.Mrs.Ms.Teacher.</div>
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. Example: 2

* See the section **Notes on the Essay Data** for more details about these features.

Additionally, the `resources.csv` data set provides more data about the resources required for each project. Each line in this file represents a resource required by a project:

Feature	Description
id	A <code>project_id</code> value from the <code>train.csv</code> file. Example: p036502
description	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
quantity	Quantity of the resource required. Example: 3
price	Price of the resource required. Example: 9.95

Note: Many projects require multiple resources. The `id` value corresponds to a `project_id` in `train.csv`, so you use it as a key to retrieve all resources needed for a project:

The data set contains the following label (the value you will attempt to predict):

Label	Description
project_is_approved	A binary flag indicating whether DonorsChoose approved the project. A value of <code>0</code> indicates the project was not approved, and a value of <code>1</code> indicates the project was approved.

Notes on the Essay Data

Prior to May 17, 2016, the prompts for the essays were as follows:

- `__project_essay_1:__` "Introduce us to your classroom"
- `__project_essay_2:__` "Tell us more about your students"
- `__project_essay_3:__` "Describe how your students will use the materials you're requesting"
- `__project_essay_3:__` "Close by sharing why your project will make a difference"

Starting on May 17, 2016, the number of essays was reduced from 4 to 2, and the prompts for the first 2 essays were changed to the following:

- `__project_essay_1:__` "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- `__project_essay_2:__` "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

For all projects with `project_submitted_datetime` of 2016-05-17 and later, the values of `project_essay_3` and `project_essay_4` will be NaN.

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

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

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

Mounted at /content/drive

```
In [ ]: 1 !cp -r '/content/drive/My Drive/_datasets/glove_vectors' '/content/'
2 !cp -r '/content/drive/My Drive/_datasets/resources.csv' '/content/'
3 !cp -r '/content/drive/My Drive/_datasets/train_data.csv' '/content/'
```

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

```
In [ ]: 1 ## drop unnecesary columns
2 project_data.drop(columns=['Unnamed: 0'],inplace=True)
```

```
In [ ]: 1 resource_data.head(2)
```

Out[12]:

	id	description	quantity	price
0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.00
1	p069063	Bouncy Bands for Desks (Blue support pipes)	3	14.95

```
In [ ]: 1 print('Number of train points in train data',project_data.shape)
2 print('***50')
3 print('Number of train points in resource data',project_data.shape)
4 print('***50')
5 print('Columns in resource data',resource_data.columns.values)
6 print('***50')
7 print('Columns in project data',project_data.columns.values)
```

Number of train points in train data (109248, 16)

Number of train points in resource data (109248, 16)

Columns in resource data ['id' 'description' 'quantity' 'price']

Columns in project data ['id' 'teacher_id' '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' 'project_is_approved']

2. Preprocessing Categorical Features: project_grade_category

```
In [ ]: 1 project_data['project_grade_category'].value_counts()
```

Out[14]:

Grades PreK-2	44225
Grades 3-5	37137
Grades 6-8	16923
Grades 9-12	10963

Name: project_grade_category, dtype: int64

```
In [ ]: 1 # https://stackoverflow.com/questions/36383821/pandas-dataframe-apply-function-to-column-strings-based-on-other-column-value
2 project_data['project_grade_category'] = project_data['project_grade_category'].str.replace(' ', '_')
3 project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('-', '_')
4 project_data['project_grade_category'] = project_data['project_grade_category'].str.lower()
5 project_data['project_grade_category'].value_counts()
```

```
Out[15]: grades_prek_2    44225
grades_3_5      37137
grades_6_8      16923
grades_9_12     10963
Name: project_grade_category, dtype: int64
```

3. Preprocessing Categorical Features: project_subject_categories

```
In [ ]: 1 project_data['project_subject_categories'].value_counts(dropna=False)
```

```
Out[16]: Literacy & Language                23655
Math & Science                17072
Literacy & Language, Math & Science    14636
Health & Sports                10177
Music & The Arts                5180
Special Needs                 4226
Literacy & Language, Special Needs    3961
Applied Learning              3771
Math & Science, Literacy & Language    2289
Applied Learning, Literacy & Language  2191
History & Civics               1851
Math & Science, Special Needs        1840
Literacy & Language, Music & The Arts  1757
Math & Science, Music & The Arts      1642
Applied Learning, Special Needs      1467
History & Civics, Literacy & Language  1421
Health & Sports, Special Needs       1391
Warmth, Care & Hunger              1309
Math & Science, Applied Learning     1220
Applied Learning, Math & Science     1052
Literacy & Language, History & Civics   809
Health & Sports, Literacy & Language   803
Applied Learning, Music & The Arts     758
Math & Science, History & Civics       652
Literacy & Language, Applied Learning  636
Applied Learning, Health & Sports      608
Math & Science, Health & Sports        414
History & Civics, Math & Science       322
History & Civics, Music & The Arts     312
Special Needs, Music & The Arts       302
Health & Sports, Math & Science       271
History & Civics, Special Needs       252
Health & Sports, Applied Learning     192
Applied Learning, History & Civics     178
Health & Sports, Music & The Arts     155
Music & The Arts, Special Needs       138
Literacy & Language, Health & Sports   72
Health & Sports, History & Civics      43
History & Civics, Applied Learning    42
Special Needs, Health & Sports        42
Special Needs, Warmth, Care & Hunger  23
Health & Sports, Warmth, Care & Hunger 23
Music & The Arts, Health & Sports     19
Music & The Arts, History & Civics    18
History & Civics, Health & Sports     13
Math & Science, Warmth, Care & Hunger  11
Applied Learning, Warmth, Care & Hunger 10
Music & The Arts, Applied Learning    10
Literacy & Language, Warmth, Care & Hunger 9
Music & The Arts, Warmth, Care & Hunger 2
History & Civics, Warmth, Care & Hunger 1
Name: project_subject_categories, dtype: int64
```

```

In [ ]: ▶ 1 categories = list(project_data['project_subject_categories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
3
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7 cat_list = []
8 for i in categories:
9     temp = ""
10    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
11    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
12        if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
13            j=j.replace('The', '') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
14            j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
15            temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing spaces
16            temp = temp.replace('&','_') # we are replacing the & value into
17    cat_list.append(temp.strip().lower())
18
19 project_data['clean_categories'] = cat_list
20 project_data.drop(['project_subject_categories'], axis=1, inplace=True)
21
22 from collections import Counter
23 my_counter = Counter()
24 for word in project_data['clean_categories'].values:
25     my_counter.update(word.split())
26
27 cat_dict = dict(my_counter)
28 sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))

```

4. Preprocessing Categorical Features: teacher_prefix

```

In [ ]: ▶ 1 project_data['teacher_prefix'].value_counts(dropna=False)

```

```

Out[18]: Mrs.      57269
Ms.       38955
Mr.       10648
Teacher   2360
Dr.        13
NaN         3
Name: teacher_prefix, dtype: int64

```

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 [ ]: ▶ 1 project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('Mrs.')

```

```
In [ ]: 1 project_data['teacher_prefix'].value_counts(dropna=False)
```

Out[20]: Mrs. 57272
Ms. 38955
Mr. 10648
Teacher 2360
Dr. 13
Name: teacher_prefix, dtype: int64

Remove '.'
convert all the chars to small

```
In [ ]: 1 project_data['teacher_prefix']=project_data['teacher_prefix'].str.replace('.', '')  
2 project_data['teacher_prefix']=project_data['teacher_prefix'].str.lower()  
3 project_data['teacher_prefix'].value_counts()
```

Out[21]: mrs 57272
ms 38955
mr 10648
teacher 2360
dr 13
Name: teacher_prefix, dtype: int64

5. Preprocessing Categorical Features: project_subject_subcategories

```
In [ ]: 1 project_data['project_subject_subcategories'].value_counts()
```

Out[22]: Literacy 9486
Literacy, Mathematics 8325
Literature & Writing, Mathematics 5923
Literacy, Literature & Writing 5571
Mathematics 5379
...
Economics, Nutrition Education 1
Community Service, Music 1
Gym & Fitness, Social Sciences 1
Parent Involvement, Team Sports 1
Extracurricular, Financial Literacy 1
Name: project_subject_subcategories, Length: 401, dtype: int64


```

In [ ]: ▶ 1 sub_categories = list(project_data['project_subject_subcategories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
3
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7
8 sub_cat_list = []
9 for i in sub_categories:
10     temp = ""
11     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
12     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
13         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
14             j=j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
15             j = j.replace(' ', '') # we are placing all the ' ' (space) with '' (empty) ex: "Math & Science"=>"Math&Science"
16             temp +=j.strip()+" #" "abc ".strip() will return "abc", remove the trailing spaces
17             temp = temp.replace('&', '_')
18     sub_cat_list.append(temp.strip().lower())
19
20 project_data['clean_subcategories'] = sub_cat_list
21 project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
22
23 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
24 my_counter = Counter()
25 for word in project_data['clean_subcategories'].values:
26     my_counter.update(word.split())
27
28 sub_cat_dict = dict(my_counter)
29 sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))

```

6. Preprocessing Categorical Features: school_state

In []:

1 project_data['school_state'].value_counts()

Out[24]:

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
VA	2045
WI	1827
AL	1762
UT	1731
TN	1688
CT	1663
MD	1514
NV	1367
MS	1323
KY	1304
OR	1242
MN	1208
CO	1111
AR	1049
ID	693
IA	666
KS	634
NM	557
DC	516
HI	507
ME	505
WV	503
NH	348
AK	345
DE	343
NE	309
SD	300
RI	285
MT	245
ND	143
WY	98
VT	80

Name: school_state, dtype: int64

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

Out[25]:

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
va	2045
wi	1827
al	1762
ut	1731
tn	1688
ct	1663
md	1514
nv	1367
ms	1323
ky	1304
or	1242
mn	1208
co	1111
ar	1049
id	693
ia	666
ks	634
nm	557
dc	516
hi	507
me	505
wv	503
nh	348
ak	345
de	343
ne	309
sd	300
ri	285
mt	245
nd	143
wy	98
vt	80

Name: school_state, dtype: int64

7. Preprocessing Categorical Features: project_title

```
In [ ]: ▶ 1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\'re", " are", phrase)
12    phrase = re.sub(r"\'s", " is", phrase)
13    phrase = re.sub(r"\'d", " would", phrase)
14    phrase = re.sub(r"\'ll", " will", phrase)
15    phrase = re.sub(r"\'t", " not", phrase)
16    phrase = re.sub(r"\'ve", " have", phrase)
17    phrase = re.sub(r"\'m", " am", phrase)
18    return phrase
```

```
In [ ]: ▶ 1 # https://gist.github.com/sebleier/554280
2 # we are removing the words from the stop words list: 'no', 'nor', 'not'
3 stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
4             "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
5             'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', \
6             'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
7             'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
8             'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
9             'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
10            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
11            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \
12            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
13            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
14            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', \
15            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', \
16            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
17            'won', "won't", 'wouldn', "wouldn't"]
```

```
In [ ]: ▶ 1 print("printing some random reviews")
2 print(9, project_data['project_title'].values[9])
3 print(34, project_data['project_title'].values[34])
4 print(147, project_data['project_title'].values[147])
```

```
printing some random reviews
9 Just For the Love of Reading--\r\nPure Pleasure
34 \"Have A Ball!!!\"
147 Who needs a Chromebook?\r\nWE DO!!
```

```
In [ ]: 1 # Combining all the above stundents
2 from tqdm import tqdm
3 def preprocess_text(text_data):
4     preprocessed_text = []
5     # tqdm is for printing the status bar
6     for sentence in tqdm(text_data):
7         sent = decontracted(sentence)
8         sent = sent.replace('\r', ' ')
9         sent = sent.replace('\n', ' ')
10        sent = sent.replace('\\"', ' ')
11        sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12        # https://gist.github.com/sebleier/554280
13        sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
14        preprocessed_text.append(sent.lower().strip())
15    return preprocessed_text
```

```
In [ ]: 1 preprocessed_titles = preprocess_text(project_data['project_title'])
2
```

100%|██████████| 109248/109248 [00:02<00:00, 40305.44it/s]

```
In [ ]: 1 ## combine all the essays into 1
2 # merge two column text dataframe:
3 project_data["essay"] = project_data["project_essay_1"].map(str) + \
4                          project_data["project_essay_2"].map(str) + \
5                          project_data["project_essay_3"].map(str) + \
6                          project_data["project_essay_4"].map(str)
```

```
In [ ]: 1 preprocessed_essays = preprocess_text(project_data['essay'])
```

100%|██████████| 109248/109248 [01:00<00:00, 1803.75it/s]

```
In [ ]: 1 ## drop unnecessary columns
2 project_data.drop(columns=['teacher_id', "project_essay_1", "project_essay_2",
3                          "project_essay_3", "project_essay_4"], inplace=True)
4 project_data.head(2)
```

Out[33]:

	id	teacher_prefix	school_state	project_submitted_datetime	project_grade_category	project_title	project_resource_summary	teacher_number_of_previously_posted_projects	project_is_approved	clean_categories
0	p253737	mrs	in	2016-12-05 13:43:57	grades_prek_2	Educational Support for English Learners at Home	My students need opportunities to practice beg...	0	0	literacy_language
1	p258326	mr	fl	2016-10-25 09:22:10	grades_6_8	Wanted: Projector for Hungry Learners	My students need a projector to help with view...	7	1	history_civics health_sports

```
In [ ]: 1 project_data['preprocessed_essay'] = preprocessed_essays
```

8. Preprocessing Numerical Values: price

```
In [ ]: 1 # https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indexes-for-all-groups-in-one-step
2 price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
3 price_data.head(2)
```

```
Out[35]:
```

	id	price	quantity
0	p000001	459.56	7
1	p000002	515.89	21

```
In [ ]: 1 # join two dataframes in python:
2 project_data = pd.merge(project_data, price_data, on='id', how='left')
```

```
In [ ]: 1 project_data['price'].head()
```

```
Out[37]: 0    154.60
1    299.00
2    516.85
3    232.90
4     67.98
Name: price, dtype: float64
```

```
In [ ]: 1 if (not os.path.isfile('processed_data_before_split.h1')):
2     hdf = pd.HDFStore('processed_data_before_split.h1')
3     hdf.put('project_data',project_data, format='table', data_columns=True)
4     hdf.close()
```

9. Split the data into Train,Test and CV

```
In [ ]: 1
2 if (not os.path.isfile('processed_data_split.h2')):
3     x_train,x_test,y_train,y_test = train_test_split(project_data,project_data['project_is_approved'],test_size=0.2,stratify=project_data['project_is_approved'])
4     x_train,x_cv,y_train,y_cv = train_test_split(x_train,y_train,test_size=0.2,stratify=y_train)
5     hdf = pd.HDFStore('processed_data_split.h2')
6     hdf.put('x_train',x_train, format='table', data_columns=True)
7     hdf.put('x_test',x_test, format='table', data_columns=True)
8     hdf.put('x_cv',x_cv, format='table', data_columns=True)
9     hdf.put('y_train',y_train, format='table', data_columns=True)
10    hdf.put('y_test',y_test, format='table', data_columns=True)
11    hdf.put('y_cv',y_cv, format='table', data_columns=True)
12    hdf.close()
13 else:
14     x_train = pd.read_hdf('processed_data_split.h2', 'x_train',mode='r')
15     x_test = pd.read_hdf('processed_data_split.h2', 'x_test',mode='r')
16     x_cv = pd.read_hdf('processed_data_split.h2', 'x_cv',mode='r')
17     y_train =pd.read_hdf('processed_data_split.h2', 'y_train',mode='r')
18     y_test =pd.read_hdf('processed_data_split.h2', 'y_test',mode='r')
19     y_cv =pd.read_hdf('processed_data_split.h2', 'y_cv',mode='r')
20     print(' Successfully loaded processed data split')
```

```
In [ ]: 1 x_train.head(2)
```

Out[40]:

	id	teacher_prefix	school_state	project_submitted_datetime	project_grade_category	project_title	project_resource_summary	teacher_number_of_previously_posted_projects	project_is_approved	clean_catego	
	54916	p243348	mrs	mo	2016-09-09 12:14:53	grades_3_5	Help Us Wiggle While We Work! We Need *Flexibl...	My students need flexible seating options such...	0	1	literacy_langu math_scie
	33715	p173833	mrs	ca	2016-09-21 17:35:27	grades_prek_2	Using iPads as a Superpower!	My students need iPads as a teaching tool, as ...	1	1	health_sp specialne

10. Vectorization on categorical(label encoding) and numerical features

```
In [ ]: 1 ## clean categories
2 from sklearn.preprocessing import LabelEncoder
3 t = LabelEncoder()
4 t.fit(x_train['clean_categories'].values)
5
6 train_categories = t.transform(x_train['clean_categories'].values)
7 test_categories = t.transform(x_test['clean_categories'].values)
8 cv_categories = t.transform(x_cv['clean_categories'].values)
9
10 print("Shape of matrix of Train data after one hot encoding ",train_categories.shape)
11 print("Shape of matrix of Test data after one hot encoding ",test_categories.shape)
12 print("Shape of matrix of CV data after one hot encoding ",cv_categories.shape)
```

Shape of matrix of Train data after one hot encoding (69918,)
Shape of matrix of Test data after one hot encoding (21850,)
Shape of matrix of CV data after one hot encoding (17480,)

```

In [ ]: ▶ 1 ## clean subcategories
2
3 t = LabelEncoder()
4 t.fit(x_train['clean_subcategories'].values)
5
6 # https://www.thetopsites.net/article/51321922.shtml
7 # few categories which are not in train but are in test are labeled as unknown
8 x_test['clean_subcategories'] = x_test['clean_subcategories'].map(lambda s: '<unknown>' if s not in t.classes_ else s)
9 x_cv['clean_subcategories'] = x_cv['clean_subcategories'].map(lambda s: '<unknown>' if s not in t.classes_ else s)
10 t.classes_ = np.append(t.classes_, '<unknown>')
11
12 train_subcategories = t.transform(x_train['clean_subcategories'].values)
13 test_subcategories = t.transform(x_test['clean_subcategories'].values)
14 cv_subcategories = t.transform(x_cv['clean_subcategories'].values)
15
16 print("Shape of matrix of Train data after one hot encoding ",train_subcategories.shape)
17 print("Shape of matrix of Test data after one hot encoding ",test_subcategories.shape)
18 print("Shape of matrix of CV data after one hot encoding ",cv_categories.shape)

```

Shape of matrix of Train data after one hot encoding (69918,)
 Shape of matrix of Test data after one hot encoding (21850,)
 Shape of matrix of CV data after one hot encoding (17480,)

```

In [ ]: ▶ 1 # we use count vectorizer to convert the values into one hot vectors
2 ## school state
3
4 t = LabelEncoder()
5 t.fit(x_train['school_state'].values)
6
7 # https://www.thetopsites.net/article/51321922.shtml
8 # few categories which are not in train but are in test are labeled as unknown
9 x_test['school_state'] = x_test['school_state'].map(lambda s: '<unknown>' if s not in t.classes_ else s)
10 x_cv['school_state'] = x_cv['school_state'].map(lambda s: '<unknown>' if s not in t.classes_ else s)
11 t.classes_ = np.append(t.classes_, '<unknown>')
12
13 sklstate_train = t.transform(x_train['school_state'].values)
14 sklstate_test = t.transform(x_test['school_state'].values)
15 sklstate_cv = t.transform(x_cv['school_state'].values)
16
17 print("Shape of matrix of Train data after one hot encoding ",sklstate_train.shape)
18 print("Shape of matrix of Test data after one hot encoding ",sklstate_test.shape)
19 print("Shape of matrix of CV data after one hot encoding ",sklstate_cv.shape)

```

Shape of matrix of Train data after one hot encoding (69918,)
 Shape of matrix of Test data after one hot encoding (21850,)
 Shape of matrix of CV data after one hot encoding (17480,)


```

In [ ]: ▶ 1 # we use count vectorizer to convert the values into one hot vectors
          2 ## teacher_prefix
          3
          4
          5 t = LabelEncoder()
          6 t.fit(x_train['teacher_prefix'].values)
          7
          8 # https://www.thetopsites.net/article/51321922.shtml
          9 # few categories which are not in train but are in test are labeled as unknown
         10
         11 x_test['teacher_prefix'] = x_test['teacher_prefix'].map(lambda s: '<unknown>' if s not in t.classes_ else s)
         12 x_cv['teacher_prefix'] = x_cv['teacher_prefix'].map(lambda s: '<unknown>' if s not in t.classes_ else s)
         13 t.classes_ = np.append(t.classes_, '<unknown>')
         14
         15 teacher_prefix_train = t.transform(x_train['teacher_prefix'].values)
         16 teacher_prefix_test = t.transform(x_test['teacher_prefix'].values)
         17 teacher_prefix_cv = t.transform(x_cv['teacher_prefix'].values)
         18
         19 print("Shape of matrix of Train data after one hot encoding ",teacher_prefix_train.shape)
         20 print("Shape of matrix of Test data after one hot encoding ",teacher_prefix_test.shape)
         21 print("Shape of matrix of CV data after one hot encoding ",teacher_prefix_cv.shape)

```

Shape of matrix of Train data after one hot encoding (69918,)

Shape of matrix of Test data after one hot encoding (21850,)

Shape of matrix of CV data after one hot encoding (17480,)

```

In [ ]: ▶ 1 # we use count vectorizer to convert the values into one hot vectors
          2 ## project_grade
          3
          4
          5 t = LabelEncoder()
          6 t.fit(x_train['project_grade_category'].values)
          7
          8 # https://www.thetopsites.net/article/51321922.shtml
          9 # few categories which are not in train but are in test are labeled as unknown
         10
         11 x_test['project_grade_category'] = x_test['project_grade_category'].map(lambda s: '<unknown>' if s not in t.classes_ else s)
         12 x_cv['project_grade_category'] = x_cv['project_grade_category'].map(lambda s: '<unknown>' if s not in t.classes_ else s)
         13 t.classes_ = np.append(t.classes_, '<unknown>')
         14
         15 proj_grade_train = t.transform(x_train['project_grade_category'].values)
         16 proj_grade_test = t.transform(x_test['project_grade_category'].values)
         17 proj_grade_cv = t.transform(x_cv['project_grade_category'].values)
         18
         19 print("Shape of matrix of Train data after one hot encoding ",proj_grade_train.shape)
         20 print("Shape of matrix of Test data after one hot encoding ",proj_grade_test.shape)
         21 print("Shape of matrix of CV data after one hot encoding ",proj_grade_cv.shape)

```

Shape of matrix of Train data after one hot encoding (69918,)

Shape of matrix of Test data after one hot encoding (21850,)

Shape of matrix of CV data after one hot encoding (17480,)

Vectorize Numerical features

```
In [ ]: 1 ### price
2 from sklearn.preprocessing import StandardScaler
3 price_vectorize = StandardScaler()
4
5 price_vectorize.fit(x_train['price'].values.reshape(-1,1))
6 proj_price_train = price_vectorize.transform(x_train['price'].values.reshape(-1,1))
7 proj_price_test = price_vectorize.transform(x_test['price'].values.reshape(-1,1))
8 proj_price_cv = price_vectorize.transform(x_cv['price'].values.reshape(-1,1))
9
10 #print(price_vectorize.get_feature_names())
11 print("Shape of matrix of Train data after one hot encoding ",proj_price_train.shape)
12 print("Shape of matrix of Test data after one hot encoding ",proj_price_test.shape)
13 print("Shape of matrix of CV data after one hot encoding ",proj_price_cv.shape)
```

Shape of matrix of Train data after one hot encoding (69918, 1)
Shape of matrix of Test data after one hot encoding (21850, 1)
Shape of matrix of CV data after one hot encoding (17480, 1)

```
In [ ]: 1 print(proj_price_train[0], '\n', proj_price_train[1])
```

[0.4741993]
[-0.07774343]

```
In [ ]: 1 ## quantity
2 quantity_vectorize = StandardScaler()
3
4 quantity_vectorize.fit(x_train['quantity'].values.reshape(-1,1))
5 proj_quantity_train = quantity_vectorize.transform(x_train['quantity'].values.reshape(-1,1))
6 proj_quantity_test = quantity_vectorize.transform(x_test['quantity'].values.reshape(-1,1))
7 proj_quantity_cv = quantity_vectorize.transform(x_cv['quantity'].values.reshape(-1,1))
8
9 #print(quantity_vectorize.get_feature_names())
10 print("Shape of matrix of Train data after one hot encoding ",proj_quantity_train.shape)
11 print("Shape of matrix of Test data after one hot encoding ",proj_quantity_test.shape)
12 print("Shape of matrix of CV data after one hot encoding ",proj_quantity_cv.shape)
```

Shape of matrix of Train data after one hot encoding (69918, 1)
Shape of matrix of Test data after one hot encoding (21850, 1)
Shape of matrix of CV data after one hot encoding (17480, 1)

```
In [ ]: 1 print(proj_quantity_train[0], '\n', proj_quantity_train[1])
```

[0.04025279]
[-0.57287171]

```
In [ ]: 1 ## teacher_number_of_previously_posted_projects
2 prev_proj_vectorize = StandardScaler()
3
4 prev_proj_vectorize.fit(x_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
5 proj_prev_proj_train = prev_proj_vectorize.transform(x_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
6 proj_prev_proj_test = prev_proj_vectorize.transform(x_test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
7 proj_prev_proj_cv = prev_proj_vectorize.transform(x_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
8
9 #print(prev_proj_vectorize.get_feature_names())
10 print("Shape of matrix of Train data after one hot encoding ",proj_prev_proj_train.shape)
11 print("Shape of matrix of Test data after one hot encoding ",proj_prev_proj_test.shape)
12 print("Shape of matrix of CV data after one hot encoding ",proj_prev_proj_cv.shape)
```

Shape of matrix of Train data after one hot encoding (69918, 1)
Shape of matrix of Test data after one hot encoding (21850, 1)
Shape of matrix of CV data after one hot encoding (17480, 1)

```
In [ ]: 1
2 ### concatenating numerical features for building model
3
4 numerical_train = np.concatenate((proj_price_train,proj_quantity_train,proj_prev_proj_train),axis=1)
5 numerical_test = np.concatenate((proj_price_test,proj_quantity_test,proj_prev_proj_test),axis=1)
6 numerical_cv = np.concatenate((proj_price_cv,proj_quantity_cv,proj_prev_proj_cv),axis=1)
7
8 print("Shape of matrix of Train data after one hot encoding and concatination ",numerical_train.shape)
9 print("Shape of matrix of Test data after one hot encoding and concatination ",numerical_test.shape)
10 print("Shape of matrix of CV data after one hot encoding and concatination ",numerical_cv.shape)
```

Shape of matrix of Train data after one hot encoding and concatination (69918, 3)
Shape of matrix of Test data after one hot encoding and concatination (21850, 3)
Shape of matrix of CV data after one hot encoding and concatination (17480, 3)

Model-1

Build and Train deep neural network as shown below

Type *Markdown* and LaTeX: α^2



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

- **Input_seq_total_text_data** --- You have to give Total text data columns. After this use the Embedding layer to get word vectors. Use given predefined glove word vectors, don't train any word vectors. After this use LSTM and get the LSTM output and Flatten that output.
- **Input_school_state** --- Give 'school_state' column as input to embedding layer and Train the Keras Embedding layer.
- **Project_grade_category** --- Give 'project_grade_category' column as input to embedding layer and Train the Keras Embedding layer.
- **Input_clean_categories** --- Give 'input_clean_categories' column as input to embedding layer and Train the Keras Embedding layer.
- **Input_clean_subcategories** --- Give 'input_clean_subcategories' column as input to embedding layer and Train the Keras Embedding layer.
- **Input_clean_subcategories** --- Give 'input_teacher_prefix' column as input to embedding layer and Train the Keras Embedding layer.
- **Input_remaining_teacher_number_of_previously_posted_projects__resource_summary_contains_numerical_digits__price__quantity** ---concatenate remaining columns and add a Dense layer after that.

```
In [ ]: 1 ## convert essay to sequences
2 '''The essay is in textual format ,we need to convert to sequences of index and pad them'''
3 from keras.preprocessing.text import Tokenizer
4 tok = Tokenizer()
5 tok.fit_on_texts(x_train['preprocessed_essay'].values)
6 seq_x_train = tok.texts_to_sequences(x_train['preprocessed_essay'].values)
7 seq_x_test = tok.texts_to_sequences(x_test['preprocessed_essay'].values)
8 seq_x_cv = tok.texts_to_sequences(x_cv['preprocessed_essay'].values)
9 vocab_size = len(tok.word_index) + 1
```

Using TensorFlow backend.

```
In [ ]: 1 print(vocab_size)
```

47255

```
In [ ]: 1 print('Train data after padding and sequencing')
        2 print('*'*50)
        3 print(padseq_x_train[1],len(padseq_x_train[1]))
```

```
In [ ]: 1 '''using glove vectors lets create a embedding matrix such that for every word
2         in vocabulary we store its corresponding glove vector in matrix form'''
3 with open('glove_vectors', 'rb') as f:
4     model = pickle.load(f)
```

```

In [ ]: ▶ 1 import numpy as np
2 #embedd_matrix = np.zeros((vocab_len,max_review_Length))
3 glove_words = model.keys()
4 emd_i =dict()
5
6 ## Lets create a dictionary that stores the 300 dim glove vector as value and the word's index as key
7 for i,w in tok.index_word.items():
8     #if w in glove_words:
9     emd_i[i] = model.get(w)
10    #else: emd_i[i] = np.zeros((1,max_review_Length))
11
12 ## emd_matrix stores all the 300 dimensional glove vectors of words based on their rank from the tokenizer.
13 ## the most frequent word is given the highest rank
14
15 # emd_matrix = np.zeros((vocab_len,max_review_Length))
16 # print(emd_matrix.shape)
17 # for i in range(1,vocab_len+1):
18 #     emd_matrix[i-1] = emd_i[i]
19
20 # create a weight matrix for words in training docs
21 print('Loaded %s word vectors.' % len(emd_i))
22 # create a weight matrix for words in training docs
23 embedding_matrix = np.zeros((vocab_size, 300))
24 for word, i in tok.word_index.items():
25     embedding_vector = model.get(word)
26     if embedding_vector is not None:
27         embedding_matrix[i] = embedding_vector

```

Loaded 47254 word vectors.

```

In [ ]: ▶ 1 print('shape of embedding matrix = ',embedding_matrix.shape)

```

shape of embedding matrix = (47255, 300)

```

In [ ]: ▶ 1 if os.path.isfile('model_inputs_labelencode.pkl') :
2     os.remove("model_inputs_labelencode.pkl")
3     print("File model_inputs Removed!")
4 with open('model_inputs_labelencode.pkl', 'wb') as f:
5     pickle.dump([emd_i,embedding_matrix,seq_x_train,seq_x_test,seq_x_cv,
6                 padseq_x_train,sklstate_train,proj_grade_train,train_categories,train_subcategories,
7                 teacher_prefix_train,numerical_train,
8                 padseq_x_test,sklstate_test,proj_grade_test,test_categories,test_subcategories,
9                 teacher_prefix_test,numerical_test,
10                padseq_x_cv,sklstate_cv,proj_grade_cv,cv_categories,cv_subcategories,
11                teacher_prefix_cv,numerical_cv],f)
12

```

Model 1

```

In [ ]: ▶ 1 import tensorflow as tf
          2 from keras.callbacks import TensorBoard
          3 import keras
          4 import keras.backend as k
          5 from sklearn.metrics import roc_auc_score
          6 from keras.layers import Dropout, Input, Activation, Dense, Embedding, concatenate, LSTM, Flatten, BatchNormalization
          7 from keras.models import Model
          8 def aucroc(y_true, y_pred):
          9     try:
         10         return tf.py_func(roc_auc_score, (y_true, y_pred), tf.double)
         11     except ValueError:
         12         pass

```

Using TensorFlow backend.

Hyperparameter tuning

- * hyperparameters:
- * embedding output dimension
- * number of lstm layers
- * numerical layer for numerical input

```

In [ ]: ▶ 1 ### convert y_train, y_test and y_cv to categorical
          2 ## one hot encode the target labels
          3 from keras.utils import np_utils
          4 from keras.regularizers import l2
          5 classes=2
          6
          7 if not os.path.isfile('model_input_cat_labels.pkl') :
          8     y_train_cat = np_utils.to_categorical(y_train, classes)
          9     y_test_cat = np_utils.to_categorical(y_test, classes)
         10     y_cv_cat = np_utils.to_categorical(y_cv, classes)
         11     with open('model_input_cat_labels.pkl', 'wb') as f:
         12         pickle.dump([y_train_cat, y_test_cat, y_cv_cat], f)
         13 else:
         14     y_train_cat, y_test_cat, y_cv_cat = pickle.load(open('model_input_cat_labels.pkl', 'rb'))
         15     print('-----Successfully loaded one hot y labels-----')

```


In []: ▶

```

1
2
3 ## clear the graph of the tensorflow
4 k.clear_session()
5 ### defining all the Input Layer
6 set_random_seed(2)
7 input_seq_total_text_data = Input(shape=padseq_x_train[0].shape,name='text_Input')
8 input_school_state = Input(shape=(1,),name='school_state_Input')
9 input_project_grade_category = Input(shape=(1,),name='project_grade_category_Input')
10 input_clean_categories = Input(shape=(1,),name='input_clean_categories_Input')
11 input_clean_subcategories = Input(shape=(1,),name='input_clean_subcategories_Input')
12 input_teacher_prefix = Input(shape=(1,),name='input_teacher_prefix')
13 input_numerical = Input(shape=(numerical_train.shape[1],),name='input_numerical')
14
15 auc_scores_model1 = []
16 if (not os.path.isfile('tuning_output.pkl')):
17     for embedding_index in [128,64,32]:
18         for lstm_index in [128,64,32]:
19             for num_dense_index in [128,64,32]:
20                 ## Define embedding layers for all inputs
21                 embedding_layer_text = Embedding(input_dim=vocab_size,output_dim=300,weights = [embedding_matrix]
22                                                    ,trainable=False)(input_seq_total_text_data)
23                 embedding_layer_school_state = Embedding(input_dim=1,output_dim=embedding_index,
24                                                         input_length=1)(input_school_state)
25                 embedding_layer_project_grade_category = Embedding(input_dim=1,output_dim=embedding_index,
26                                                                    input_length=1)(input_project_grade_category)
27                 embedding_layer_clean_categories = Embedding(input_dim=1,output_dim=embedding_index,
28                                                             input_length=1)(input_clean_categories)
29                 embedding_layer_clean_subcategories = Embedding(input_dim=1,output_dim=embedding_index,
30                                                                input_length=1)(input_clean_subcategories)
31                 embedding_layer_teacher_prefix = Embedding(input_dim=1,output_dim=embedding_index,
32                                                           input_length=1)(input_teacher_prefix)
33
34                 ### Define LSTM for the text
35                 '''Return sequences = True ensure output from all theLSTM is returned not just the final output from last LSTM'''
36                 lstm_layer_text = LSTM(lstm_index,return_sequences=True)(embedding_layer_text)
37
38                 ### Define flatten layer and Dense layer for numerical input
39                 flatten_text = Flatten()(lstm_layer_text)
40                 flatten_school_state = Flatten()(embedding_layer_school_state)
41                 flatten_project_grade_category = Flatten()(embedding_layer_project_grade_category)
42                 flatten_clean_categories = Flatten()(embedding_layer_clean_categories)
43                 flatten_clean_subcategories = Flatten()(embedding_layer_clean_subcategories)
44                 flatten_teacher_prefix = Flatten()(embedding_layer_teacher_prefix)
45                 rem_input_dense = Dense(num_dense_index,activation='relu',kernel_initializer='he_normal')(input_numerical)
46
47                 ##Concatenate all the layers
48                 concat_layer = concatenate([flatten_text,flatten_school_state,flatten_project_grade_category,flatten_clean_categories,
49                                           flatten_clean_subcategories,flatten_teacher_prefix,rem_input_dense])
50
51                 ##define three dense layers with dropout
52                 dense1_layer = Dense(256,activation='relu',kernel_initializer='he_normal')(concat_layer)
53                 regularization_layer1 = BatchNormalization()(dense1_layer)
54                 regularization_layer1 = Dropout(0.35)(regularization_layer1)
55                 dense2_layer = Dense(128,activation='relu',kernel_initializer='he_normal')(regularization_layer1)
56                 regularization_layer2 = BatchNormalization()(dense2_layer)
57                 regularization_layer2 = Dropout(0.35)(regularization_layer2)
58                 dense3_layer = Dense(64,activation='relu',kernel_initializer='he_normal')(regularization_layer2)
59                 regularization_layer2 = BatchNormalization()(dense3_layer)
60                 #regularization_layer2 = Dropout(0.25)(regularization_layer2)

```



```

61 output_layer = Dense(2,activation='sigmoid',kernel_initializer='glorot_normal',activity_regularizer=l2(0.0001))(regularization_layer2)
62
63 model1 = Model(inputs=[input_seq_total_text_data,
64                       input_school_state,input_project_grade_category,
65                       input_clean_categories,input_clean_subcategories,
66                       input_teacher_prefix,input_numerical],outputs=output_layer)
67
68 ## Compile the model1 with default Learning rate
69 model1.compile(optimizer=keras.optimizers.Adam(),loss='categorical_crossentropy',metrics=['accuracy',aucroc])
70 callback = tf.keras.callbacks.EarlyStopping(monitor='val_aucroc',verbose=1, patience=3,restore_best_weights=True,mode='max')
71 history = model1.fit([padseq_x_train,sklstate_train,proj_grade_train,train_categories,train_subcategories,
72                     teacher_prefix_train,numerical_train],y_train_cat,epochs=10,batch_size=1000,verbose=1,
73                     validation_data=[padseq_x_cv,sklstate_cv,proj_grade_cv,cv_categories,
74                                     cv_subcategories,teacher_prefix_cv,numerical_cv],y_cv_cat],
75                     callbacks=[callback])
76 max_ = np.argmax(history.history['val_aucroc'])
77 print('Validation loss for embedding units={0}, lstm layer={1} ,numerical dense units={2} '.format(embedding_index,lstm_index,num_dense_index),
78       ' is :',history.history['val_loss'][max_])
79
80 auc_scores_model1.append((embedding_index,lstm_index,num_dense_index,history.history['accuracy'][max_]
81                          ,history.history['loss'][max_],history.history['aucroc'][max_],
82                          history.history['val_accuracy'][max_],history.history['val_loss'][max_],history.history['val_aucroc'][max_]))
83
84 df = pd.DataFrame(data=auc_scores_model1,columns=['Embedding units','LSTM units','Dense numerical units',
85                                                  'Train Accuracy','Train Loss','Train auc','Test Accuracy','Test Loss','Test auc'])
86 best_param = df[df['Test auc'] == df['Test auc'].max()]
87 with open('tuning_output.pkl','wb') as f:
88     pickle.dump([df,auc_scores_model1,best_param] , f)
89 else:
90     df, auc_scores_model1, best_param = pickle.load(open('tuning_output.pkl','rb'))
91     print('----Tuning output loaded -----')

```

Train on 69918 samples, validate on 17480 samples

Epoch 1/10

69918/69918 [=====] - 30s 430us/step - loss: 0.5324 - accuracy: 0.7301 - aucroc: 0.5473 - val_loss: 0.5243 - val_accuracy: 0.7248 - val_aucroc: 0.5070

Epoch 2/10

69918/69918 [=====] - 28s 407us/step - loss: 0.4343 - accuracy: 0.7942 - aucroc: 0.5990 - val_loss: 0.4122 - val_accuracy: 0.8062 - val_aucroc: 0.5379

Epoch 3/10

69918/69918 [=====] - 29s 410us/step - loss: 0.4042 - accuracy: 0.8100 - aucroc: 0.6251 - val_loss: 0.4498 - val_accuracy: 0.8096 - val_aucroc: 0.5930

Epoch 4/10

69918/69918 [=====] - 28s 407us/step - loss: 0.3785 - accuracy: 0.8239 - aucroc: 0.6499 - val_loss: 0.4575 - val_accuracy: 0.8146 - val_aucroc: 0.6483

Epoch 5/10

69918/69918 [=====] - 28s 406us/step - loss: 0.3474 - accuracy: 0.8390 - aucroc: 0.7008 - val_loss: 0.4596 - val_accuracy: 0.8190 - val_aucroc: 0.6524

Epoch 6/10

69918/69918 [=====] - 29s 408us/step - loss: 0.3052 - accuracy: 0.8578 - aucroc: 0.7496 - val_loss: 0.4588 - val_accuracy: 0.8188 - val_aucroc: 0.6898

Epoch 7/10

In []:

```
1 df.head(10)
```

Out[101]:

	Embedding units	LSTM units	Dense numerical units	Train Accuracy	Train Loss	Train auc	Test Accuracy	Test Loss	Test auc
0	128	128	128	0.857833	0.305186	0.749588	0.818821	0.458777	0.689762
1	128	128	64	0.830687	0.344531	0.707659	0.806121	0.471380	0.671242
2	128	128	32	0.814297	0.408000	0.599179	0.814130	0.421263	0.657687
3	128	64	128	0.889242	0.236805	0.849639	0.782666	0.531789	0.683967
4	128	64	64	0.861180	0.348306	0.663887	0.841133	0.392900	0.687134
5	128	64	32	0.868060	0.297007	0.751041	0.829519	0.451504	0.686675
6	128	32	128	0.822392	0.389784	0.644577	0.827689	0.386795	0.706951
7	128	32	64	0.847121	0.389075	0.562900	0.854348	0.392449	0.613339
8	128	32	32	0.818473	0.391271	0.634026	0.817334	0.407131	0.683609
9	64	128	128	0.822864	0.362238	0.692673	0.810412	0.417715	0.686193

Best Hyperparameter values

In []:

```
1 df[df['Test auc'] == df['Test auc'].max()]
```

Out[25]:

	Embedding units	LSTM units	Dense numerical units	Train Accuracy	Train Loss	Train auc	Test Accuracy	Test Loss	Test auc
24	32	32	128	0.859121	0.327327	0.725428	0.834153	0.409079	0.714769

In []:

```
1 # Load the TensorBoard notebook extension
2 %load_ext tensorboard
3 import datetime, os
4 from tensorflow import set_random_seed
5
```

Train the Model with best hyperparameters

In []: ▶

```

1 from keras.callbacks import ModelCheckpoint, ReduceLROnPlateau
2 # Clear any logs from previous runs
3 !rm -rf ./logs/
4 #set_random_seed(65)
5 ## clear the graph of the tensorflow
6 tf.keras.backend.clear_session()
7 ### defining all the Input Layer
8 input_seq_total_text_data = Input(shape=padseq_x_train[0].shape, name='text_Input')
9 input_school_state = Input(shape=(1,), name='school_state_Input')
10 input_project_grade_category = Input(shape=(1,), name='project_grade_category_Input')
11 input_clean_categories = Input(shape=(1,), name='input_clean_categories_Input')
12 input_clean_subcategories = Input(shape=(1,), name='input_clean_subcategories_Input')
13 input_teacher_prefix = Input(shape=(1,), name='input_teacher_prefix')
14 input_numerical = Input(shape=(numerical_train.shape[1],), name='input_numerical')
15
16 ## Define embedding layers for all inputs
17 embedding_layer_text = Embedding(input_dim=vocab_size, output_dim=300, weights = [embedding_matrix], trainable=False)(input_seq_total_text_data)
18 embedding_layer_school_state = Embedding(input_dim=1, output_dim=int(best_param['Embedding units']))(input_school_state)
19 embedding_layer_project_grade_category = Embedding(input_dim=1, output_dim=int(best_param['Embedding units']))(input_project_grade_category)
20 embedding_layer_clean_categories = Embedding(input_dim=1, output_dim=int(best_param['Embedding units']))(input_clean_categories)
21 embedding_layer_clean_subcategories = Embedding(input_dim=1, output_dim=int(best_param['Embedding units']))(input_clean_subcategories)
22 embedding_layer_teacher_prefix = Embedding(input_dim=1, output_dim=int(best_param['Embedding units']))(input_teacher_prefix)
23
24 ### Define LSTM for the text
25 '''Return sequences = True ensure output from all the LSTM is returned not just the final output from last LSTM'''
26 lstm_layer_text = LSTM(int(best_param['LSTM units']), return_sequences=True)(embedding_layer_text)
27
28 ### Define flatten layer and Dense layer for numerical input
29 flatten_text = Flatten()(lstm_layer_text)
30 flatten_school_state = Flatten()(embedding_layer_school_state)
31 flatten_project_grade_category = Flatten()(embedding_layer_project_grade_category)
32 flatten_clean_categories = Flatten()(embedding_layer_clean_categories)
33 flatten_clean_subcategories = Flatten()(embedding_layer_clean_subcategories)
34 flatten_teacher_prefix = Flatten()(embedding_layer_teacher_prefix)
35 rem_input_dense = Dense(int(best_param['Dense numerical units']), activation='relu', kernel_initializer='he_normal')(input_numerical)
36
37 ##Concatenate all the layers
38 concat_layer = concatenate([flatten_text, flatten_school_state, flatten_project_grade_category, flatten_clean_categories,
39                             flatten_clean_subcategories, flatten_teacher_prefix, rem_input_dense])
40
41 ##define three dense layers with dropout
42 dense1_layer = Dense(512, kernel_initializer='he_normal')(concat_layer)
43 activation = LeakyReLU(0.3)(dense1_layer)
44 regularization_layer1 = BatchNormalization()(activation)
45 regularization_layer1 = Dropout(0.15)(regularization_layer1)
46 dense2_layer = Dense(256, kernel_initializer='he_normal')(regularization_layer1)
47 activation = LeakyReLU(0.3)(dense2_layer)
48 regularization_layer2 = BatchNormalization()(activation)
49 regularization_layer2 = Dropout(0.25)(regularization_layer2)
50 dense3_layer = Dense(128, kernel_initializer='he_normal')(regularization_layer2)
51 activation = LeakyReLU(0.3)(dense3_layer)
52 regularization_layer2 = BatchNormalization()(activation)
53 regularization_layer2 = Dropout(0.15)(regularization_layer2)
54 dense4_layer = Dense(64, kernel_initializer='he_normal')(regularization_layer2)
55 activation = LeakyReLU(0.3)(dense4_layer)
56 regularization_layer2 = BatchNormalization()(activation)
57 regularization_layer2 = Dropout(0.25)(regularization_layer2)
58 dense5_layer = Dense(32, kernel_initializer='he_normal')(regularization_layer2)
59 activation = LeakyReLU(0.3)(dense5_layer)
60 regularization_layer2 = BatchNormalization()(activation)

```

```

61 regularization_layer2 = Dropout(0.25)(regularization_layer2)
62 output_layer = Dense(2,activation='sigmoid',kernel_initializer='glorot_normal',activity_regularizer=l2(0.002))(regularization_layer2)
63
64 if not os.path.isfile('best_model_output.pkl'):
65     model1 = Model(inputs=[input_seq_total_text_data,
66                           input_school_state,input_project_grade_category,
67                           input_clean_categories,input_clean_subcategories,
68                           input_teacher_prefix,input_numerical],outputs=output_layer)
69
70     model1.compile(optimizer=keras.optimizers.Adam(),loss='categorical_crossentropy',metrics=['accuracy',aucroc])
71
72     log_dir="logs/" + datetime.datetime.now().strftime("%Y%m%d-%H%M%S")
73     tensorboard_callback = tf.keras.callbacks.TensorBoard(log_dir=log_dir,histogram_freq=0, write_graph=True,write_grads=True)
74
75     ## early stopping
76     #https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/EarlyStopping
77     #https://stackoverflow.com/questions/48285129/saving-best-model-in-keras
78     # https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/ReduceLRonPlateau
79     callback = tf.keras.callbacks.EarlyStopping(monitor='val_aucroc',verbose=1, patience=4,restore_best_weights=True,mode='max')
80     mcp_save = ModelCheckpoint('mdl_wts.hdf5', save_best_only=True, monitor='val_aucroc', mode='max')
81     reduce_lr_2 = ReduceLRonPlateau(monitor='val_aucroc', factor=0.2,patience=2, min_lr=0.001,verbose = 1,mode='max')
82     history = model1.fit([padseq_x_train,sklstate_train,proj_grade_train,train_categories,train_subcategories,
83                         teacher_prefix_train,numerical_train],y_train_cat,epochs=25,batch_size=1000,verbose=1,
84                         validation_data=[[padseq_x_cv,sklstate_cv,proj_grade_cv,cv_categories,
85                                         cv_subcategories,teacher_prefix_cv,numerical_cv],y_cv_cat],
86                         callbacks=[tensorboard_callback,callback,mcp_save,reduce_lr_2])
87
88     hist = history.history
89     with open('best_model_output.pkl','wb') as f:
90         pickle.dump(hist, f)
91
92 else:
93     hist = pickle.load(open('best_model_output.pkl', 'rb'))
94     print('----Model output loaded after tuning-----')

```

Train on 69918 samples, validate on 17480 samples

Epoch 1/25

69918/69918 [=====] - 29s 419us/step - loss: 1.3231 - accuracy: 0.7031 - aucroc: 0.5269 - val_loss: 1.0148 - val_accuracy: 0.7858 - val_aucroc: 0.5286

Epoch 2/25

69918/69918 [=====] - 28s 399us/step - loss: 0.8121 - accuracy: 0.8006 - aucroc: 0.5615 - val_loss: 0.6074 - val_accuracy: 0.8150 - val_aucroc: 0.5537

Epoch 3/25

69918/69918 [=====] - 28s 398us/step - loss: 0.6442 - accuracy: 0.8168 - aucroc: 0.5961 - val_loss: 0.5800 - val_accuracy: 0.8226 - val_aucroc: 0.6332

Epoch 4/25

69918/69918 [=====] - 28s 394us/step - loss: 0.5620 - accuracy: 0.8264 - aucroc: 0.6276 - val_loss: 0.5295 - val_accuracy: 0.8288 - val_aucroc: 0.6885

Epoch 5/25

69918/69918 [=====] - 28s 394us/step - loss: 0.5042 - accuracy: 0.8357 - aucroc: 0.6500 - val_loss: 0.5388 - val_accuracy: 0.8020 - val_aucroc: 0.7339

Epoch 6/25

69918/69918 [=====] - 28s 399us/step - loss: 0.4630 - accuracy: 0.8429 - aucroc: 0.6716 - val_loss: 0.4603 - val_accuracy: 0.8395 - val_aucroc: 0.7073

Epoch 7/25

69918/69918 [=====] - 28s 398us/step - loss: 0.4340 - accuracy: 0.8486 - aucroc: 0.6872 - val_loss: 0.4614 - val_accuracy: 0.8377 - val_aucroc: 0.6988

Epoch 00007: ReduceLRonPlateau reducing learning rate to 0.001.

Epoch 8/25

69918/69918 [=====] - 28s 399us/step - loss: 0.4078 - accuracy: 0.8546 - aucroc: 0.7137 - val_loss: 0.4485 - val_accuracy: 0.8372 - val_aucroc: 0.6979

Epoch 9/25

69918/69918 [=====] - 28s 402us/step - loss: 0.3803 - accuracy: 0.8627 - aucroc: 0.7334 - val_loss: 1.4165 - val_accuracy: 0.5404 - val_aucroc: 0.5253

Restoring model weights from the end of the best epoch.

Epoch 00009: ReduceLRonPlateau reducing learning rate to 0.001.

Epoch 00009: early stopping

In []: ▶

1

%tensorboard --logdir logs/

2

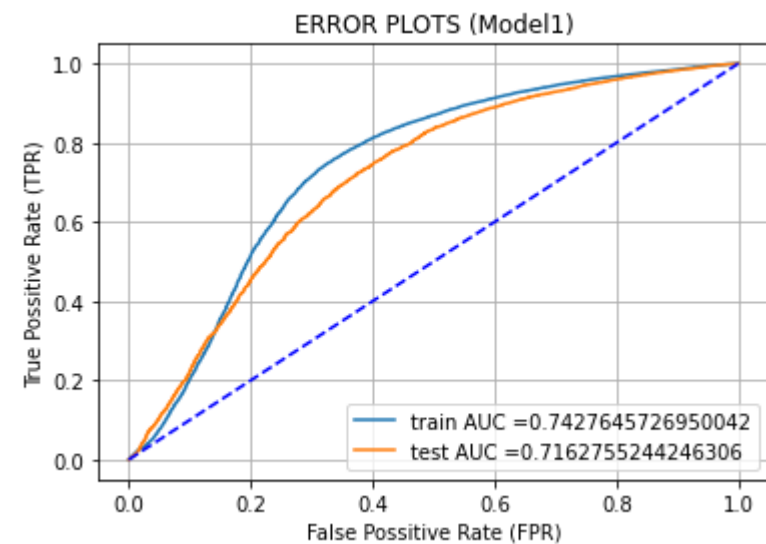
Confusion Matrix and ROC_AUC curve

In []: ▶

```

1 from sklearn.metrics import confusion_matrix
2 ## Finding best threshold for predictions
3 def best_threshold(thresholds,fpr,tpr):
4     t=thresholds[np.argmax(tpr*(1-fpr))]
5     # (tpr*(1-fpr)) will be maximum if your fpr is very Low and tpr is very high
6     print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
7     return t
8
9 def predict_with_best_t(proba, threshold):
10    predictions = []
11    for i in proba:
12        if i>=threshold:
13            predictions.append(1)
14        else:
15            predictions.append(0)
16    return predictions
17
18
19 model1.load_weights('mdl_wts.hdf5')
20 ## Predict the test and train
21 # if not os.path.file('model_predictions.pkl'):
22 y_test_predict = model1.predict([padseq_x_test,sklstate_test,proj_grade_test,test_categories,test_subcategories,
23                                 teacher_prefix_test,numerical_test],use_multiprocessing=True)[:,1]
24 y_train_predict = model1.predict([padseq_x_train,sklstate_train,proj_grade_train,train_categories,train_subcategories,
25                                 teacher_prefix_train,numerical_train],use_multiprocessing=True)[:,1]
26     #
27 # else:
28 #     y_train_predict,y_test_predict = pickle.load(open('model_predictions.pkl','rb'))
29 #     print('-----Predicts Loaded-----')
30 ## Store fpr and tpr rates
31
32 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_predict)
33 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_predict)
34
35
36 #plot
37 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
38 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
39 plt.legend()
40 plt.xlabel("False Possitive Rate (FPR)")
41 plt.ylabel("True Possitive Rate (TPR)")
42 plt.title("ERROR PLOTS (Model1)")
43 plt.plot([0, 1], [0, 1], 'b--')
44 plt.grid()
45 plt.show()
46
47 print("=*100)
48
49 best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
50

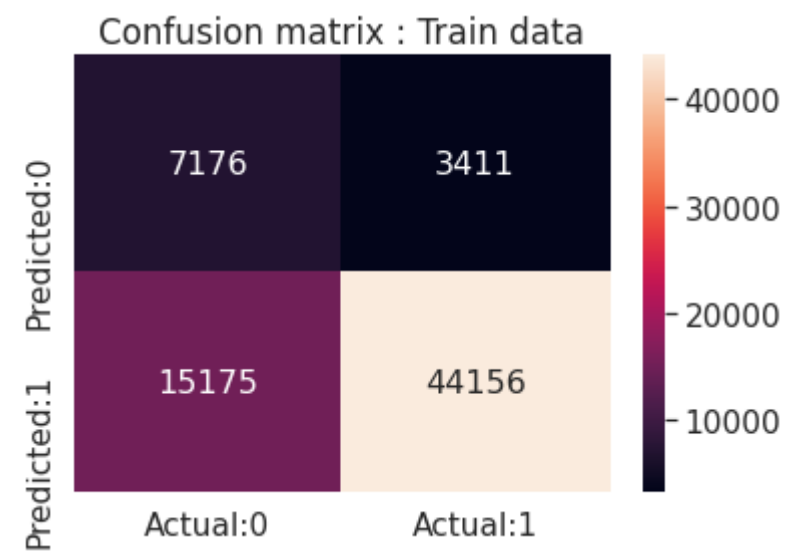
```



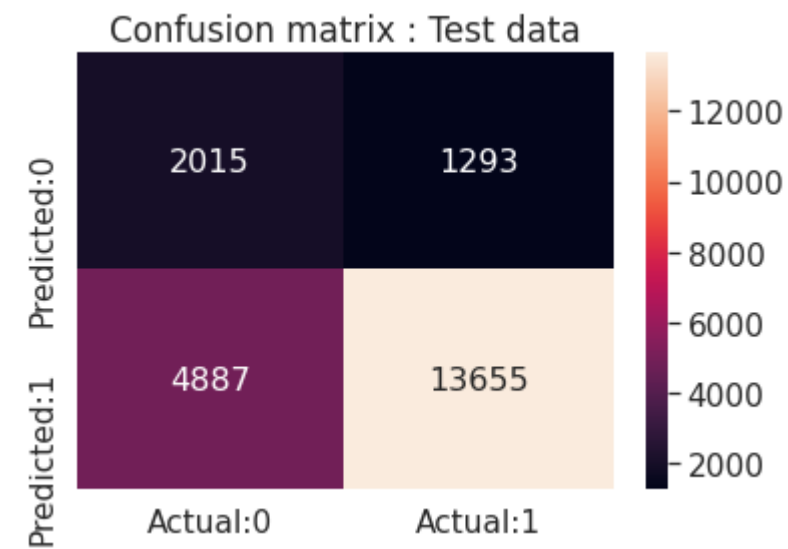
=====

the maximum value of $tpr \cdot (1 - fpr)$ 0.5044493576696496 for threshold 0.08

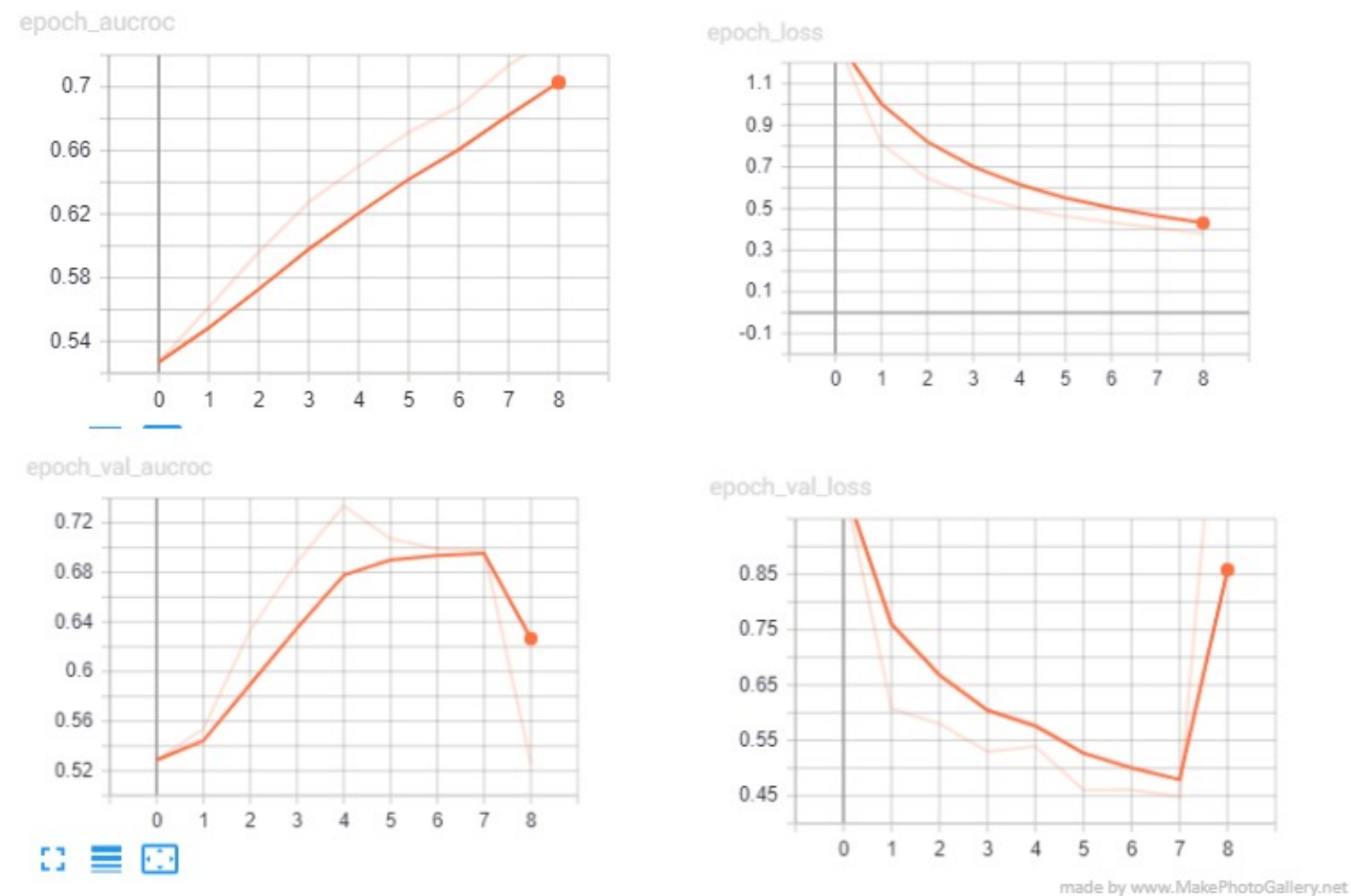
```
In [ ]: 1 ### PLOT the matrix for Train
2 #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
3 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
4 df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_predict, best_t))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



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



- Validation Loss, Validation aucroc, Train Loss, Train aucroc



Model1 Summary :

- * Preprocessed and vectorized the input variables along with label encoding categorical data.
- * We Tuned our model with various embedding input dimensions,dense numerical units and LSTM units with early stopping techniques and found the best combination with highest Test auc values .
- * After finding the best combination we tuned our best model with this combination increasing the number of layers, dropout and batch normalization layers,including Leaky relu activation and activity regularizers to improve the model performance
- * Used keras callback methods to perform early stopping,reduce learning rate,model checkpoints to monitor our model to attain maximum validation auc.
- * We see that as the number of epochs goes above 7 our Validation auc reduces and loss increases.Hence we save the model at the best epoch.
- * When we plot FPR against TPR our model gives 0.74 as train auc and 0.71 as test auc.

```
In [ ]: 1 ### load evrything to drive
2 # !cp -r '/content/tuning_output.pkl' '/content/drive/My Drive/LSTM_preprocessed/model1/'
3 # !cp -r '/content/processed_data_split.h2' '/content/drive/My Drive/LSTM_preprocessed/model1/'
4 # #!cp -r '/content/processed_data_before_split.h1' '/content/drive/My Drive/LSTM_preprocessed/model1/'
5 # #!cp -r '/content/model_predictions.pkl' '/content/drive/My Drive/LSTM_preprocessed/model1/'
6 # !cp -r '/content/model_inputs_labelencode.pkl' '/content/drive/My Drive/LSTM_preprocessed/model1/'
7 # !cp -r '/content/model_input_cat_labels.pkl' '/content/drive/My Drive/LSTM_preprocessed/model1/'
8 # !cp -r '/content/best_model_output.pkl' '/content/drive/My Drive/LSTM_preprocessed/model1/'u
9 # !cp -r '/content/mdl_wts.hdf5' '/content/drive/My Drive/LSTM_preprocessed/model1/'
```

```
In [ ]: 1 ### Load evrything from drive
2 # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model1/tuning_output.pkl' '/content/'
3 # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model1/processed_data_split.h2' '/content/'
4 # # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model1/processed_data_before_split.h1' '/content/'
5 # # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model1/model_predictions.pkl' '/content/'
6 # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model1/model_inputs_labelencode.pkl' '/content/'
7 # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model1/model_input_cat_labels.pkl' '/content/'
8 # !cp -r '/content/drive/My Drive/LSTM_preprocessed/model1/best_model_output.pkl' '/content/'
9
```