## ▾ 1.1 Reading Data

```python
1  %matplotlib inline
2  import warnings
3  warnings.filterwarnings("ignore")
4
5  import sqlite3
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 from sklearn.feature_extraction.text import CountVectorizer
15 from sklearn.metrics import confusion_matrix
16 from sklearn import metrics
17 from sklearn.metrics import roc_curve, auc
18 from nltk.stem.porter import PorterStemmer
19 import re
20 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
21 import string
22 from nltk.corpus import stopwords
23 from nltk.stem import PorterStemmer
24 from nltk.stem.wordnet import WordNetLemmatizer
25 from tqdm import tqdm
26 import os
27 # from plotly import plotly
28 # import plotly.offline as offline
29 # import plotly.graph_objs as go
30 # offline.init_notebook_mode()
31 from collections import Counter
```

## ▾ importing file into googel colab

```python
1  #  resources data
2  import gdown
3
4  url = 'https://drive.google.com/uc?id=1OcMV5zjAJI7OvNxxN4Ant52BDF3jrZOZ'
5  output = 'resources.csv'
6  # https://drive.google.com/file/d/1OcMV5zjAJI7OvNxxN4Ant52BDF3jrZOZ/view?usp=sharing
7  gdown.download(url, output, quiet=False)
8
```

```
1  # test data
2  import gdown
3
4  url = 'https://drive.google.com/uc?id=1JGtsNLea4Q2HZQIgBp3pRrOfRN80qIg0'
5  # https://drive.google.com/file/d/1JGtsNLea4Q2HZQIgBp3pRrOfRN80qIg0/view?usp=sharing
6  output = 'train_data.csv'
7  gdown.download(url, output, quiet=False)
```

```
1  ls
```

```
1  project_data =pd.read_csv("train_data.csv")
2  resource_data = pd.read_csv("resources.csv")
```

```
1  print("The shape of the Train data ",project_data.shape)
2  print("-"*50)
3  print("The number of attributes in Train data","-"*5,project_data.columns.values)
```

```python
cols = ["Date" if x=="project_submitted_datetime" else x for x in list(project_data.columns)]
project_data["Date"] = pd.to_datetime(project_data["project_submitted_datetime"])
project_data.drop("project_submitted_datetime",axis=1,inplace = True )
project_data.sort_values(by=["Date"],inplace=True)

project_data = project_data[cols]
project_data.head(2)
```

```python
print("Number of data points in train data", resource_data.shape)
print(resource_data.columns.values)
resource_data.head(2)
```

⤷

```python
# we cannot remove rows where teacher prefix is not available therefore we are replacing 'nan' value with
# 'null'(string)
#https://stackoverflow.com/questions/42224700/attributeerror-float-object-has-no-attribute-split
project_data['teacher_prefix'] = project_data['teacher_prefix'].fillna('null')
```

## ▼ 1.2 preprocessing of project_subject_categories

```python
catogories = list(project_data['project_subject_categories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
cat_list = []
for i in catogories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    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 catogory based on space "Math & Science"=> "Math","&",\
13 #              "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 placeing 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('&','_') # we are replacing the & value into
18      cat_list.append(temp.strip().lower())
19
20 project_data['clean_categories'] = cat_list
21 project_data.drop(['project_subject_categories'], axis=1, inplace=True)
22
23 from collections import Counter
24 my_counter = Counter()
25 for word in project_data['clean_categories'].values:
26     my_counter.update(word.split())
27
28 cat_dict = dict(my_counter)
29 sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
30 print("The Worlds in sorted_cat_dict",sorted_cat_dict)
```

## ▾ 1.3 preprocessing of project_subject_subcategories

```
 1 catogories1 = 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 cat_list1 = []
 8 for i in catogories1:
 9     temp1 = ""
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 catogory based on space "Math & Science"=> "Math","&",\
13 #              "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 placeing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
16          temp1+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
17          temp1 = temp1.replace('&','_') # we are replacing the & value into
18      cat_list1.append(temp1.strip().lower())
19
20 project_data['clean_sub_categories'] = cat_list1
21 project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
22
23 from collections import Counter
24 my_counter1 = Counter()
25 for word in project_data['clean_sub_categories'].values:
26     my_counter1.update(word.split())
```

```
27
28  cat_dict1 = dict(my_counter1)
29  sorted_cat_dict1 = dict(sorted(cat_dict1.items(), key=lambda kv: kv[1]))
30  print("The Worlds in sorted_cat_dict1",sorted_cat_dict1)
```

## ▾ 1.4 Text preprocessing of project essay

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

```
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
```

```
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"]
```

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

⤷

```python
1  preprocessed_essays[20000]
```

⤷

## ▾ 1.5 Preprocessing of project_title

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

```
1  preprocessed_titles[1000]
```

```
1  # here we are removing the unwanted coloumns from the data that has been processed .
2  project_data["clean_titles"] = preprocessed_titles
3  project_data.drop(["project_essay_1"],axis=1,inplace=True)
4  project_data.drop(["project_essay_2"],axis=1,inplace=True)
5  project_data.drop(["project_essay_3"],axis=1,inplace=True)
6  project_data.drop(["project_essay_4"],axis=1,inplace=True)
7  project_data.drop(["project_title"],axis=1,inplace=True)
```

## ▾ 1.6 Merging the tow DataFrame (Resources.csv and Train.csv)

```
1  project_data.columns
```

```
1  resource_data.head(2)
```

```
1  price_data= resource_data.groupby("id").agg({"price" : "sum" , "quantity" : "sum"}).reset_index()
2  price_data.head(2)
```

⤷

```
1  project_data = pd.merge(project_data,price_data,on="id",how = "left")
2  project_data.head(2)
3  # project_data.drop(["Unnamed: 0"],axis=1,inplace=True)
4  # project_data.drop(["id"],axis=1,inplace=True)
5  # project_data.drop(["teacher_id"],axis=1,inplace=True)
```

⤷

## ▾ 1.7 Preprocessing of school_state

```
1  school_state = list(project_data['school_state'].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
```

```python
 8  school_state_list = []
 9  for i in school_state:
10      temp2 = ""
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 catogory 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 placeing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
16          temp2 +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
17          temp2 = temp2.replace('&','_')
18      school_state_list.append(temp2.strip().lower())
19
20  # droping the school_state column
21  project_data['School_state'] = school_state_list
22  project_data.drop(['school_state'], axis=1, inplace=True)
23
24  # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
25  my_counter3 = Counter()
26  for word in project_data['School_state'].values:
27      my_counter3.update(word.split())
28
29  school_state_dict = dict(my_counter3)
30  sorted_school_state_dict = dict(sorted(school_state_dict.items(), key=lambda kv: kv[1]))
31  print("The Values in sorted_school_state_dict : ", sorted_school_state_dict)
32
33
```

## ▾ 2.1 Splitting data

```python
1  project_data = project_data.head(50000)
2  Y = project_data["project_is_approved"].values
3  X = project_data.drop(["project_is_approved"],axis = 1)
4  print(X.columns)
```

⤷

```python
1  #Splitting the data into train and test data_set
2  from sklearn.model_selection import train_test_split
```

```
3  X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.33, stratify=Y)
4  X_train, X_cv, Y_train, Y_cv = train_test_split(X_train, Y_train, test_size=0.33, stratify=Y_train)
```

## ▾ 2.2 Resampling data using the Randomover sampler for Imblearn.over sampling

```
1  from imblearn.over_sampling import RandomOverSampler
2  from collections import Counter
3  import warnings
4  warnings.filterwarnings("ignore")
5
6  ros = RandomOverSampler(sampling_strategy='minority',random_state=42)
7  x_train, y_train = ros.fit_resample(X_train, Y_train)
8  print('Resampled dataset shape %s' % Counter(y_train))
9  print("Capitial" ,"X","represents the original train_data and lower case" ,"x", "represnts the ramdonly over-sampled data")
```

⤷

```
1  # here we have to convert x into a dataframe
2  x_train = pd.DataFrame(x_train,columns = X.columns)
3  x_train.head(1)
```

```
1  print(x_train.shape, y_train.shape)
2  print(X_cv.shape, Y_cv.shape)
3  print(X_test.shape, Y_test.shape)
4
5  print("="*100)
6
```

⤷

## ▾ 2.3 Make Data Model Ready: Vectorizing data

## ▾ Bag of Words is used for Vectorizing of Text data

```
 1  from sklearn.feature_extraction.text import CountVectorizer
 2  vectorizer = CountVectorizer(min_df=10, max_features=2000)
 3  vectorizer.fit(x_train['essay'].values) # fit has to happen only on train data
 4
 5  # we use the transform Text_data to vector , BOW CountVectorizer
 6  x_train_essay_bow = vectorizer.transform(x_train['essay'].values) # this the vectorization of the oversampled data we do it as we
 7  X_cv_essay_bow = vectorizer.transform(X_cv['essay'].values)  # here we use the X_CV for vectorization( only x_train is oversampled
 8  X_test_essay_bow = vectorizer.transform(X_test['essay'].values)
 9
10  print("="*100)
11  print("After vectorizations of the text data the shape of the data is ")
12  print(x_train_essay_bow.shape, y_train.shape)
13  print(X_cv_essay_bow.shape, Y_cv.shape)
14  print(X_test_essay_bow.shape, Y_test.shape)
15
16
17
```

⤷

```
 1  from sklearn.feature_extraction.text import CountVectorizer
 2  vectorizer = CountVectorizer(min_df=10, max_features=2000)
 3  vectorizer.fit(x_train['clean_titles'].values) # fit has to happen only on train data
 4
 5  # we use the fitted CountVectorizer to convert the text to vector
 6  x_train_titles_bow = vectorizer.transform(x_train['clean_titles'].values)
 7  X_cv_titles_bow = vectorizer.transform(X_cv['clean_titles'].values)
 8  X_test_titles_bow = vectorizer.transform(X_test['clean_titles'].values)
 9
10  print("After vectorizations of the clean_titles data the shape of the data is")
11  print(x_train_titles_bow.shape, y_train.shape)
12  print(X_cv_titles_bow.shape, Y_cv.shape)
13  print(X_test_titles_bow.shape, Y_test.shape)
14  print("*"*100)
```

⤷

```
1  from sklearn.feature_extraction.text import CountVectorizer
2  vectorizer = CountVectorizer(min_df=10, max_features=2000)
3  vectorizer.fit(x_train['project_resource_summary']) # fit has to happen only on train data
4
5  # we use the fitted CountVectorizer to convert the text to vector
6  x_train_summary_bow = vectorizer.transform(x_train['project_resource_summary'])
7  X_cv_summary_bow = vectorizer.transform(X_cv['project_resource_summary'])
8  X_test_summary_bow = vectorizer.transform(X_test['project_resource_summary'])
9
10 print("After vectorizations of the project_resource_summary data the shape of the data is")
11 print(x_train_summary_bow.shape, y_train.shape)
12 print(X_cv_summary_bow.shape, Y_cv.shape)
13 print(X_test_summary_bow.shape, Y_test.shape)
14 print("*"*100)
```

## 2.4 One-hot-encoding of the Catogorical Features

```
1  vectorizer = CountVectorizer()
2  vectorizer.fit(x_train['clean_sub_categories']) # fit has to happen only on train data
3
4
5  # we use the fitted CountVectorizer to convert the text to vector
6  x_train_clean_subcat_ohe = vectorizer.transform(x_train['clean_sub_categories'])
7  X_cv_clean_subcat_ohe = vectorizer.transform(X_cv['clean_sub_categories'])
8  X_test_clean_subcat_ohe = vectorizer.transform(X_test['clean_sub_categories'])
9
10 print("After vectorizations of the clean_sub_categories , One-hot-encoding shape of the data is")
11 print(x_train_clean_subcat_ohe.shape, y_train.shape)
12 print(X_cv_clean_subcat_ohe.shape, Y_cv.shape)
13 print(X_test_clean_subcat_ohe.shape, Y_test.shape)
14 # print(vectorizer.get_feature_names())
15 print("*"*100)
```

```
1  vectorizer = CountVectorizer()
2  vectorizer.fit(x_train['clean_categories']) # fit has to happen only on train data
3
4
5  # we use the fitted CountVectorizer to convert the text to vector
6  x_train_clean_categories_ohe = vectorizer.transform(x_train['clean_categories'])
7  X_cv_clean_categories_ohe = vectorizer.transform(X_cv['clean_categories'])
8  X_test_clean_categories_ohe = vectorizer.transform(X_test['clean_categories'])
9
10 print("After vectorizations of the clean_categories, One-hot-encoding shape of the data is")
11 print(x_train_clean_categories_ohe.shape, y_train.shape)
12 print(X_cv_clean_categories_ohe.shape, Y_cv.shape)
13 print(X_test_clean_categories_ohe.shape, Y_test.shape)
14 print(vectorizer.get_feature_names())
15 print("*"*100)
```

⊳

```
1  vectorizer = CountVectorizer()
2  vectorizer.fit(x_train['teacher_prefix'].values.astype('U')) # fit has to happen only on train data
3
4  # we use the fitted CountVectorizer to convert the text to vector
5  x_train_teacher_ohe = vectorizer.transform(x_train['teacher_prefix'].values.astype('U'))
6  X_cv_teacher_ohe = vectorizer.transform(X_cv['teacher_prefix'].values.astype('U'))
7  X_test_teacher_ohe = vectorizer.transform(X_test['teacher_prefix'].values.astype('U'))
8
9  print("After vectorizations of the teacher_prefix , One-hot-encoding shape of the data is")
10 print(x_train_teacher_ohe.shape, y_train.shape)
11 print(X_cv_teacher_ohe.shape, Y_cv.shape)
12 print(X_test_teacher_ohe.shape, Y_test.shape)
13 print(vectorizer.get_feature_names())
14 print("*"*100)
15
```

⊳

```python
1  vectorizer = CountVectorizer()
2  vectorizer.fit(x_train['School_state'].values) # fit has to happen only on train data
3
4  # we use the fitted CountVectorizer to convert the text to vector
5  x_train_state_ohe = vectorizer.transform(x_train['School_state'].values)
6  X_cv_state_ohe = vectorizer.transform(X_cv['School_state'].values)
7  X_test_state_ohe = vectorizer.transform(X_test['School_state'].values)
8
9  print("After vectorizations of the School_state , One-hot-encoding shape of the data is")
10 print(x_train_state_ohe.shape, y_train.shape)
11 print(X_cv_state_ohe.shape, Y_cv.shape)
12 print(X_test_state_ohe.shape, Y_test.shape)
13 # print(vectorizer.get_feature_names())
14 print("="*100)
```

⮕

```python
1  # #This step is to intialize a vectorizer with the vocabulary created form the project_grade_category values
2  from collections import Counter
3  my_counter5 = Counter()
4  for word in X_train['project_grade_category'].values:
5    if "Grades" in word:
6      word = word.replace("Grades","")
7    my_counter5.update(word.split())
8  # dict sort by value python: https://stackoverflow.com/a/613218/4084039
9  project_grade_category_dict = dict(my_counter5)
10 sorted_project_grade_category_dict = dict(sorted(project_grade_category_dict.items(), key=lambda kv: kv[1]))
11 print(sorted_project_grade_category_dict)
```

⮕

```python
1  vectorizer = CountVectorizer(vocabulary=list(sorted_project_grade_category_dict.keys()), lowercase=False, binary=True)
2  vectorizer.fit(x_train['project_grade_category'].values) # fit has to happen only on train data
3
4  # we use the fitted CountVectorizer to convert the text to vector
5  x_train_grade_ohe = vectorizer.transform(x_train['project_grade_category'].values)
6  X_cv_grade_ohe = vectorizer.transform(X_cv['project_grade_category'].values)
7  X_test_grade_ohe = vectorizer.transform(X_test['project_grade_category'].values)
8
```

```
 9  print("After vectorizations of the project_grade_category , One-hot-encoding shape of the data is")
10  print(x_train_grade_ohe.shape, y_train.shape)
11  print(X_cv_grade_ohe.shape, Y_cv.shape)
12  print(X_test_grade_ohe.shape,Y_test.shape)
13  print(vectorizer.get_feature_names())
14  print("="*100)
```

[→

## ▾ 2.5 Normalizing the numerical features: Price

```
 1  from sklearn.preprocessing import StandardScaler
 2  standard_vec = StandardScaler(with_mean = False)
 3  # this will rise an error Expected 2D array, got 1D array instead:
 4  # array=[105.22 215.96  96.01 ... 368.98  80.53 709.67].
 5  # Reshape your data either using
 6  # array.reshape(-1, 1) if your data has a single feature
 7  # array.reshape(1, -1)  if it contains a single sample.
 8  standard_vec.fit(x_train['price'].values.reshape(-1,1))
 9
10  x_train_price_std = standard_vec.transform(x_train['price'].values.reshape(-1,1))
11  X_cv_price_std = standard_vec.transform(X_cv['price'].values.reshape(-1,1))
12  X_test_price_std = standard_vec.transform(X_test['price'].values.reshape(-1,1))
13
14  print("After vectorizations of the price data ,  shape of the data after standazing")
15  print(x_train_price_std.shape, y_train.shape)
16  print(X_cv_price_std.shape, Y_cv.shape)
17  print(X_test_price_std.shape, Y_test.shape)
18  print("="*100)
```

[→

```
 1  from sklearn.preprocessing import StandardScaler
```

```python
 2  standard_vector = StandardScaler(with_mean = False)
 3  # this will rise an error Expected 2D array, got 1D array instead:
 4  # array=[105.22 215.96  96.01 ... 368.98  80.53 709.67].
 5  # Reshape your data either using
 6  # array.reshape(-1, 1) if your data has a single feature
 7  # array.reshape(1, -1)  if it contains a single sample.
 8  standard_vector.fit(x_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
 9
10  x_train_projects_std = standard_vector.transform(x_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
11  X_cv_projects_std = standard_vector.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
12  X_test_projects_std = standard_vector.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
13
14  print("After vectorizations of the teacher_number_of_previously_posted_projects , shape of the data after standazing")
15  print(x_train_projects_std.shape, y_train.shape)
16  print(X_cv_projects_std.shape,Y_cv.shape)
17  print(X_test_projects_std.shape, Y_test.shape)
18  print("="*100)
```

⤷

```python
 1  from sklearn.preprocessing import StandardScaler
 2  standard_vector1 = StandardScaler(with_mean = False)
 3  # this will rise an error Expected 2D array, got 1D array instead:
 4  # array=[105.22 215.96  96.01 ... 368.98  80.53 709.67].
 5  # Reshape your data either using
 6  # array.reshape(-1, 1) if your data has a single feature
 7  # array.reshape(1, -1)  if it contains a single sample.
 8  standard_vector1.fit(x_train['quantity'].values.reshape(-1,1))
 9
10  x_train_qty_std = standard_vector1.transform(x_train['quantity'].values.reshape(-1,1))
11  X_cv_qty_std = standard_vector1.transform(X_cv['quantity'].values.reshape(-1,1))
12  X_test_qty_std = standard_vector1.transform(X_test['quantity'].values.reshape(-1,1))
13
14  print("After vectorizations")
15  print(x_train_qty_std.shape, y_train.shape)
16  print(X_cv_qty_std.shape, Y_cv.shape)
17  print(X_test_qty_std.shape, Y_test.shape)
18  print("="*100)
```

⤷

## ▼ 3 Appling KNN on different kind of featurization ("BOW" and "TFIDF" )

### ▼ 3.1 **Set 1**: Categorical Vectorised data , Numerical Vectorised data, Project_title(BOW) + Preprocessed_essay (BOW)

```python
from scipy.sparse import hstack
X1_tr = hstack((x_train_clean_categories_ohe,x_train_clean_subcat_ohe,x_train_teacher_ohe,x_train_state_ohe,\
                x_train_grade_ohe,x_train_titles_bow,x_train_essay_bow,x_train_price_std,x_train_projects_std,x_train_qty_std)).to
X1_cv = hstack((X_cv_clean_categories_ohe,X_cv_clean_subcat_ohe,X_cv_teacher_ohe,X_cv_state_ohe,X_cv_grade_ohe,\
                X_cv_price_std,X_cv_projects_std,X_cv_qty_std,X_cv_essay_bow,X_cv_titles_bow)).tocsr()
X1_te =hstack((X_test_clean_categories_ohe,X_test_clean_subcat_ohe,X_test_teacher_ohe,X_test_state_ohe,\
               X_test_grade_ohe,X_test_essay_bow,X_test_titles_bow,X_test_price_std,X_test_projects_std,X_test_qty_std)).tocsr()


print("The final Data Matrix for Set:1" , " All the shapes of the data represent the merged features as mentioned in the tittle")
print("shape of X_train is : ",             X1_tr.shape)
print("shape of X_Cross validation is :" , X1_cv.shape)
print("shape of X_test is ",                X1_te.shape)
```

⮕

```python
def batch_predict(clf, data):
    y_data_pred = []
    tr_loop = data.shape[0] - data.shape[0]%1000
    # consider you X_tr shape is 49041, then your tr_loop will be 49041 - 49041%1000 = 49000
    # in this for loop we will iterate unti the last 1000 multiplier
    for i in range(0, tr_loop, 1000):
        y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
    # we will be predicting for the last data points
    if data.shape[0]%1000 !=0:
        y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])

    return y_data_pred
```

```python
import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score

```

```python
 5 train_auc = []
 6 cv_auc = []
 7 K = [3, 15, 25, 51, 101]
 8 for i in tqdm(K):
 9     neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
10     neigh.fit(X1_tr, y_train)
11
12     y_train_pred = batch_predict(neigh, X1_tr)
13     y_cv_pred = batch_predict(neigh, X1_cv)
14
15     # roc_auc_score(y_tr, y_score) the 2nd parameter should be probability estimates of the positive class
16     # not the predicted outputs
17     train_auc.append(roc_auc_score(y_train,y_train_pred))
18     cv_auc.append(roc_auc_score(Y_cv, y_cv_pred))
19
20 plt.plot(K, train_auc, label='Train AUC')
21 plt.plot(K, cv_auc, label='CV AUC')
22
23 plt.scatter(K, train_auc, label='Train AUC points')
24 plt.scatter(K, cv_auc, label='CV AUC points')
25
26 plt.legend()
27 plt.xlabel("K: hyperparameter")
28 plt.ylabel("AUC")
29 plt.title("ERROR PLOTS")
30 plt.grid(color='black', linestyle='-', linewidth=1)
31 plt.show()
```

1. The error plot above representes the Train AUC and CV AUC Curve, wherein we choose K(hyperparameter value) such that,it will be the maximum AU data and distance between the train-data line("Blue") and CV-line("orange") results be to minimum.

2. The best value of k is found to be 102

3. Here I have used the "AUCROC" curve for choosing the best hypermeter as the data is imbalanced(though i have balanced it) , the "AUCROC" curve co both the class labels equally .

4. The other hyperparameter tuning techiniques can also be choosen as "CV" and "K-fold CV"

```
1 best_k = 101
```

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2 from sklearn.metrics import roc_curve, auc
3
4
5 neigh = KNeighborsClassifier(n_neighbors=best_k, n_jobs=-1)
6 neigh.fit(X1_tr, y_train)
7 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
8 # not the predicted outputs
9
10 y_train_pred = batch_predict(neigh, X1_tr)
11 y_test_pred = batch_predict(neigh, X1_te)
12
13 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
14 test_fpr, test_tpr, te_thresholds = roc_curve(Y_test, y_test_pred)
15
16 plt.plot(train_fpr, train_tpr, label="Train AUC ="+str(auc(train_fpr, train_tpr)))
17 plt.plot(test_fpr, test_tpr, label="Test AUC ="+str(auc(test_fpr, test_tpr)))
18 plt.legend()
19 plt.xlabel("Fasle Poitive Rate")
20 plt.ylabel("True Positive Rate")
21 plt.title("ERROR PLOTS")
22 plt.grid(color='black', linestyle='-', linewidth=1)
23 plt.show()
```

**Conclusion for SET : 1**

1. The above Represents the TRP and FPR rates on the either axis and this curve is know as the AUCROC curve , it is a metrix to evaluate the performar model.

2. As the grap represents here the Train Auc = 0.67 , so from this we can conclude that the model is predecting the values with 67 % probabilty.

```
1  # we are writing our own function for predict, with defined thresould
2  # we will pick a threshold that will give the least fpr
3  def find_best_threshold(threshould, fpr, tpr):
4      t = threshould[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, threshould):
10     predictions = []
11     for i in proba:
12         if i>=threshould:
13             predictions.append(1)
14         else:
15             predictions.append(0)
16     return predictions
17
```

```
1  print("="*100)
2  from sklearn.metrics import confusion_matrix
3  best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
4  Confusion_metrix_Train_data = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
5  Confusion_metrix_Test_data = pd.DataFrame(confusion_matrix(Y_test, predict_with_best_t(y_test_pred, best_t)))
6  import seaborn as sns
7  sns.set(font_scale=1.4)#for label size
8  print("The Confusion metrix of train data")
9  sns.heatmap(Confusion_metrix_Train_data,annot=True, annot_kws={"size": 20},fmt ="g")
```

```
1 print("The Confusion metrix of test data ")
2 sns.heatmap(Confusion_metrix_Test_data,annot=True, annot_kws={"size": 20},fmt ="g")
```

⤷

1. The Ouput Stated above represents the maximun value of TRP (i.e tpr*(1-fpr) "0.3944289258349809" corrosponding to whcih the maximum thersho
.

2. The Second output represnts the Confusion metrix Based on the thershold of 0.48 , which states that the values below the thershold of 0.48 are clas
0 and the values above 0.48 are classified as 1.

## 3.2 **Set 2**: Categorical, Numerical features + Project_title(TFIDF)+ Preprocessed_essay (TFIDF)

```python
# Applying TF-IDF on Prohect title :
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(x_train["clean_titles"].values)

x_tain_project_titles_tfidf = vectorizer.transform(x_train["clean_titles"].values)
X_cv_project_titles_tfidf = vectorizer.transform(X_cv["clean_titles"].values)
X_test_project_titles_tfidf = vectorizer.transform(X_test["clean_titles"].values)


print("After TFIDF vectorizations of the clean_titles , shape of the data after standazing")
print(x_tain_project_titles_tfidf.shape, y_train.shape)
print(X_cv_project_titles_tfidf.shape,Y_cv.shape)
print(X_test_project_titles_tfidf.shape, Y_test.shape)
print("*"*100)
```

⊏→

```python
# Applying the Tfidf Vectrization on Preprocessed_essay
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(x_train["essay"])

x_tain_essay_tfidf = vectorizer.transform(x_train["essay"].values)
X_cv_essay_tfidf = vectorizer.transform(X_cv["essay"].values)
X_test_essay_tfidf = vectorizer.transform(X_test["essay"].values)


print("After TFIDF vectorizations of the essay , shape of the data after standazing")
print(x_tain_essay_tfidf.shape, y_train.shape)
print(X_cv_essay_tfidf.shape,Y_cv.shape)
```

```
14 print(X_test_essay_tfidf.shape, Y_test.shape)
15 print("*"*100)
```

⊏→

### Creating new data-set for set -2

**Catogorial data , Numerical data , Essay(TFIDF) , Project_tittle(TFIDF)**

```
 1 # the data points are merged using the Hstack which we import from sklearn
 2 from scipy.sparse import hstack
 3 X2_tr = hstack((x_train_clean_categories_ohe,x_train_clean_subcat_ohe,x_train_teacher_ohe,x_train_state_ohe,\
 4              x_train_grade_ohe,x_train_price_std,x_train_projects_std,x_train_qty_std,x_tain_project_titles_tfidf,x_tain_essay_
 5 X2_cv = hstack((X_cv_clean_categories_ohe,X_cv_clean_subcat_ohe,X_cv_teacher_ohe,X_cv_state_ohe,X_cv_grade_ohe,\
 6              X_cv_price_std,X_cv_projects_std,X_cv_qty_std,X_cv_project_titles_tfidf,X_cv_essay_tfidf)).tocsr()
 7 X2_te =hstack((X_test_clean_categories_ohe,X_test_clean_subcat_ohe,X_test_teacher_ohe,X_test_state_ohe,\
 8              X_test_grade_ohe,X_test_price_std,X_test_projects_std,X_test_qty_std,X_test_project_titles_tfidf,X_test_essay_tfid
 9
10
11 print("The final Data Matrix for Set:2" , " All the shapes of the data represent the merged features as mentioned in the tittle")
12 print("shape of X_train is : ",           X2_tr.shape)
13 print("shape of X_Cross validation is :" , X2_cv.shape)
14 print("shape of X_test is ",              X2_te.shape)
```

⊏→

```
 1 def batch_predict(clf, data):
 2   y1_data_pred = []
 3   tr_loop = data.shape[0] - data.shape[0]%1000
 4   # consider you X_tr shape is 49041, then your tr_loop will be 49041 - 49041%1000 = 49000
 5   # in this for loop we will iterate unti the last 1000 multiplier
 6   for i in range(0, tr_loop, 1000):
 7       y1_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
 8   # we will be predicting for the last data points
 9   if data.shape[0]%1000 !=0:
10       y1_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
11
```

```
12        return y1_data_pred
```

```python
1  import matplotlib.pyplot as plt
2  from sklearn.neighbors import KNeighborsClassifier
3  from sklearn.metrics import roc_auc_score
4
5  train_auc1 = []
6  cv_auc1 = []
7  K = [3, 15, 25, 51, 101]
8  for i in tqdm(K):
9      neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
10     neigh.fit(X2_tr, y_train)
11
12     y1_train_pred = batch_predict(neigh, X2_tr)
13     y1_cv_pred = batch_predict(neigh, X2_cv)
14
15     # roc_auc_score(y_tr, y_score) the 2nd parameter should be probability estimates of the positive class
16     # not the predicted outputs
17     train_auc1.append(roc_auc_score(y_train,y1_train_pred))
18     cv_auc1.append(roc_auc_score(Y_cv, y1_cv_pred))
19
20 plt.plot(K, train_auc1, label='Train AUC')
21 plt.plot(K, cv_auc1, label='CV AUC')
22
23 plt.scatter(K, train_auc1, label='Train AUC points')
24 plt.scatter(K, cv_auc1, label='CV AUC points')
25
26 plt.legend()
27 plt.xlabel("K: hyperparameter")
28 plt.ylabel("AUC")
29 plt.title("ERROR PLOTS")
30 plt.grid(color='black', linestyle='-', linewidth=1)
31 plt.show()
```

1. The error plot above representes the Train AUC and CV AUC Curve, wherein we choose K(hyperparameter value) such that,it will be the maximum AU data and distance between the train-data line("Blue") and CV-line("orange") results be to minimum.
2. The best value of k is found to be 101
3. Here I have used the "AUCROC" curve for choosing the best hypermeter as the data is imbalanced(though i have balanced it) , the "AUCROC" curve co both the class labels equally .
4. The other hyperparameter tuning techiniques can also be choosen as "CV" and "K-fold CV"

```
1  best_k1 = 101
```

```
1  # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2  from sklearn.metrics import roc_curve, auc
3
4
5  neigh = KNeighborsClassifier(n_neighbors=best_k1, n_jobs=-1)
6  neigh.fit(X2_tr, y_train)
7  # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
8  # not the predicted outputs
9
10 y2_train_pred = batch_predict(neigh, X2_tr)
11 y2_test_pred = batch_predict(neigh, X2_te)
12
13 train1_fpr, train1_tpr, tr1_thresholds = roc_curve(y_train, y2_train_pred)
14 test1_fpr, test1_tpr, te1_thresholds = roc_curve(Y_test, y2_test_pred)
15
16 plt.plot(train1_fpr, train1_tpr, label="Tain AUC ="+str(auc(train1_fpr, train1_tpr)))
17 plt.plot(test1_fpr, test1_tpr, label="Test AUC ="+str(auc(test1_fpr, test1_tpr)))
18 plt.legend()
19 plt.xlabel("Fasle Postive Rate")
20 plt.ylabel("True Pstive Rate")
21 plt.title("ROC_AUC Curve")
22 plt.grid(color='black', linestyle='-', linewidth=1)
23 plt.show()
```

↱

**Conclusion for SET : 2**

1. The above Represents the TRP and FPR rates on the either axis and this curve is know as the AUCROC curve , it is a metrix to evaluate the performar model.

2. As teh grap represents here the Train Auc = 0.67 , so from this we can conclude that the model is predecting the values with 67 % probabilty.

```
1  #this is the custom function for predecting the best thershold and sorting the values according the threshould
2  def find_best_threshold(threshould, fpr, tpr):
3      t = threshould[np.argmax(tpr*(1-fpr))]
4      # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
5      print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
6      return t
7
8
```

```
1  def predict_with_best_t(proba, threshould):
2      predictions = []
3      for i in proba:
4          if i>=threshould:
5              predictions.append(1)
6          else:
7              predictions.append(0)
8      return predictions
```

```
1  from sklearn.metrics import confusion_matrix
```

```
2  best_t = find_best_threshold(tr1_thresholds, train1_fpr, train1_tpr)
3  Confusion_metrix_Train_data = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y2_train_pred, best_t)))
4  Confusion_metrix_Test_data = pd.DataFrame(confusion_matrix(Y_test, predict_with_best_t(y2_test_pred, best_t)))
5  import seaborn as sns
6  sns.set(font_scale=1.4)#for label size
7  print("The Confusion metrix of train data")
8  sns.heatmap(Confusion_metrix_Train_data,annot=True, annot_kws={"size": 20},fmt ="g")
9
```

⤷

```
1  print("The Confusion metrix of test data ")
2  sns.heatmap(Confusion_metrix_Test_data,annot=True, annot_kws={"size": 20},fmt ="g")
```

⤷

1. The Ouput Stated above represents the maximun value of TRP (i.e tpr*(1-fpr) "0.3944" corrosponding to whcih the maximum thershold is 0.48 .

2. The Second output represnts the Confusion metrix Based on the thershold of 0.48 , which states that the values below the thershold of 0.48 are clas 0 and the values above 0.48 are classified as 1.

## ▼ 3.3 Set 3: Categorical, Numerical features + Project_title(AVG W2V)+ Preprocessed_essay (AVG W2V)

```
1  # I am using the predefined word to vector which is pre-trained , hence we use the pickel file to access the file
2  # stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
3  # please review the above link for more information
4  # code to import file in google colab form drive
5
6
7  import gdown
8
9  url = 'https://drive.google.com/uc?id=1MqUasf7jYoPbG35MJ28VQcOjjNp-ZDDp'
10 output = 'glove_vectors'
11 gdown.download(url, output, quiet=False)
```

⤷

```
1  #checking for files are present in the directory we are working or not
2  ls
```

⤷

```
1  import pickle
2  with open('glove_vectors', 'rb') as f:
3      model = pickle.load(f)
```

```
   4        glove_words =  set(model.keys())
```

```
 1  # The code below represents the Avg-word-to-vector of Project tittle .
 2  # here we are calculating the Avg-word to vec for "x_train" .
 3  x_train_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
 4  for sentence in tqdm(x_train['clean_titles'].values): # for each review/sentence
 5      vector = np.zeros(300) # as word vectors are of zero length
 6      cnt_words =0; # num of words with a valid vector in the sentence/review
 7      for word in sentence.split(): # for each word in a review/sentence
 8          if word in glove_words:
 9              vector += model[word]
10              cnt_words += 1
11      if cnt_words != 0:
12          vector /= cnt_words
13      x_train_avg_w2v_vectors.append(vector)
14
15  print(len(x_train_avg_w2v_vectors))
16  print(len(x_train_avg_w2v_vectors[0]))
17  # print(x_train_avg_w2v_vectors[0])
```

☐→

```
 1  # here we are calculating the Avg-word to vec for "X_CV"
 2  X_cv_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
 3  for sentence in tqdm(X_cv['clean_titles'].values): # for each review/sentence
 4      vector = np.zeros(300) # as word vectors are of zero length
 5      cnt_words =0; # num of words with a valid vector in the sentence/review
 6      for word in sentence.split(): # for each word in a review/sentence
 7          if word in glove_words:
 8              vector += model[word]
 9              cnt_words += 1
10      if cnt_words != 0:
11          vector /= cnt_words
12      X_cv_avg_w2v_vectors.append(vector)
13
14  print(len(X_cv_avg_w2v_vectors))
15  print(len(X_cv_avg_w2v_vectors[0]))
16  # print(X_cv_avg_w2v_vectors[0])
```

☐→

```
1  # here we are calculating the Avg-word to vec for "X_test"
2  X_test_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
3  for sentence in tqdm(X_test['clean_titles'].values): # for each review/sentence
4      vector = np.zeros(300) # as word vectors are of zero length
5      cnt_words =0; # num of words with a valid vector in the sentence/review
6      for word in sentence.split(): # for each word in a review/sentence
7          if word in glove_words:
8              vector += model[word]
9              cnt_words += 1
10     if cnt_words != 0:
11         vector /= cnt_words
12     X_test_avg_w2v_vectors.append(vector)
13
14 print(len(X_test_avg_w2v_vectors))
15 print(len(X_test_avg_w2v_vectors[0]))
16 # print(X_cv_avg_w2v_vectors[0])
```

⊏→

```
1  # The code below represents the Avg-word-to-vector of Essay.
2  # here we are calculating the Avg-word to vec for "x_train" .
3  x_train_essay_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
4  for sentence in tqdm(x_train['essay'].values): # for each review/sentence
5      vector = np.zeros(300) # as word vectors are of zero length
6      cnt_words =0; # num of words with a valid vector in the sentence/review
7      for word in sentence.split(): # for each word in a review/sentence
8          if word in glove_words:
9              vector += model[word]
10             cnt_words += 1
11     if cnt_words != 0:
12         vector /= cnt_words
13     x_train_essay_avg_w2v_vectors.append(vector)
14
15 print(len(x_train_essay_avg_w2v_vectors))
16 print(len(x_train_essay_avg_w2v_vectors[0]))
17 # print(x_train_essay_avg_w2v_vectors[0])
```

⊏→

```
1  # here we are calculating the Avg-word to vec for "X_CV"
2  X_cv_essay_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
3  for sentence in tqdm(X_cv['essay'].values): # for each review/sentence
4      vector = np.zeros(300) # as word vectors are of zero length
```

```
 5        cnt_words =0; # num of words with a valid vector in the sentence/review
 6        for word in sentence.split(): # for each word in a review/sentence
 7            if word in glove_words:
 8                vector += model[word]
 9                cnt_words += 1
10        if cnt_words != 0:
11            vector /= cnt_words
12        X_cv_essay_avg_w2v_vectors.append(vector)
13
14  print(len(X_cv_essay_avg_w2v_vectors))
15  print(len(X_cv_essay_avg_w2v_vectors[0]))
16  # print(X_cv_essay_avg_w2v_vectors[0])
```

⤷

```
 1  # here we are calculating the Avg-word to vec for "X_test"
 2  X_test_essay_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
 3  for sentence in tqdm(X_test['essay'].values): # for each review/sentence
 4      vector = np.zeros(300) # as word vectors are of zero length
 5      cnt_words =0; # num of words with a valid vector in the sentence/review
 6      for word in sentence.split(): # for each word in a review/sentence
 7          if word in glove_words:
 8              vector += model[word]
 9              cnt_words += 1
10      if cnt_words != 0:
11          vector /= cnt_words
12      X_test_essay_avg_w2v_vectors.append(vector)
13
14  print(len(X_test_essay_avg_w2v_vectors))
15  print(len(X_test_essay_avg_w2v_vectors[0]))
16  # print(X_test_essay_avg_w2v_vectors[0])
```

⤷

```
 1  from scipy.sparse import hstack
 2  X3_tr = hstack((x_train_clean_categories_ohe,x_train_clean_subcat_ohe,x_train_teacher_ohe,x_train_state_ohe,\
 3              x_train_grade_ohe,x_train_price_std,x_train_projects_std,x_train_qty_std,x_train_avg_w2v_vectors,x_train_essay_av
 4  X3_cv = hstack((X_cv_clean_categories_ohe,X_cv_clean_subcat_ohe,X_cv_teacher_ohe,X_cv_state_ohe,X_cv_grade_ohe,\
 5              X_cv_price_std,X_cv_projects_std,X_cv_qty_std,X_cv_avg_w2v_vectors,X_cv_essay_avg_w2v_vectors)).tocsr()
 6  X3_te =hstack((X_test_clean_categories_ohe,X_test_clean_subcat_ohe,X_test_teacher_ohe,X_test_state_ohe,\
 7              X_test_grade_ohe,X_test_price_std,X_test_projects_std,X_test_qty_std,X_test_avg_w2v_vectors,X_test_essay_avg_w2v_ve
 8
 9
```

```
10 print("The final Data Matrix for Set:3" , " All the shapes of the data represent the merged features as mentioned in the tittle")
11 print("shape of X_train is : ",                X3_tr.shape)
12 print("shape of X_Cross validation is :" , X3_cv.shape)
13 print("shape of X_test is ",                   X3_te.shape)
```

⟶

```
1 def batch_predict(clf, data):
2     y2_data_pred = []
3     tr_loop = data.shape[0] - data.shape[0]%1000
4     # consider you X_tr shape is 49041, then your tr_loop will be 49041 - 49041%1000 = 49000
5     # in this for loop we will iterate unti the last 1000 multiplier
6     for i in range(0, tr_loop, 1000):
7         y2_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
8     # we will be predicting for the last data points
9     if data.shape[0]%1000 !=0:
10         y2_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
11
12     return y2_data_pred
```

```
1 import matplotlib.pyplot as plt
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.metrics import roc_auc_score
4
5 train_auc2 = []
6 cv_auc2 = []
7 K = [3, 15, 25, 51, 101]
8 for i in tqdm(K):
9     neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
10     neigh.fit(X3_tr, y_train)
11
12     y2_train_pred = batch_predict(neigh, X3_tr)
13     y2_cv_pred = batch_predict(neigh, X3_cv)
14
15     # roc_auc_score(y_tr, y_score) the 2nd parameter should be probability estimates of the positive class
16     # not the predicted outputs
17     train_auc2.append(roc_auc_score(y_train,y2_train_pred))
18     cv_auc2.append(roc_auc_score(Y_cv, y2_cv_pred))
19
20 plt.plot(K, train_auc2, label='Train AUC')
21 plt.plot(K, cv_auc2, label='CV AUC')
22
23 plt.scatter(K, train_auc2, label='Train AUC points')
24 plt.scatter(K, cv_auc2, label='CV AUC points')
```

```
25
26 plt.legend()
27 plt.xlabel("K: hyperparameter")
28 plt.ylabel("AUC")
29 plt.title("ERROR PLOTS")
30 plt.grid(color='black', linestyle='-', linewidth=1)
31 plt.show()
```

↪

1. The error plot above representes the Train AUC and CV AUC Curve, wherein we choose K(hyperparameter value) such that,it will be the maximum AU data and distance between the train-data line("Blue") and CV-line("orange") results be to minimum.

2. The best value of k is found to be 101

3. Here I have used the "AUCROC" curve for choosing the best hypermeter as the data is imbalanced(though i have balanced it) , the "AUCROC" curve co both the class labels equally .

4. The other hyperparameter tuning techiniques can also be choosen as "CV" and "K-fold CV"

```
1 best_k1 = 101
```

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2 from sklearn.metrics import roc_curve, auc
3
```

```
 4
 5  neigh = KNeighborsClassifier(n_neighbors=best_k1, n_jobs=-1)
 6  neigh.fit(X3_tr, y_train)
 7  # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 8  # not the predicted outputs
 9
10  y3_train_pred = batch_predict(neigh, X3_tr)
11  y3_test_pred = batch_predict(neigh, X3_te)
12
13  train2_fpr, train2_tpr, tr2_thresholds = roc_curve(y_train, y3_train_pred)
14  test2_fpr, test2_tpr, te2_thresholds = roc_curve(Y_test, y3_test_pred)
15
16  plt.plot(train2_fpr, train2_tpr, label="Tain AUC ="+str(auc(train2_fpr, train2_tpr)))
17  plt.plot(test2_fpr, test2_tpr, label="Test AUC ="+str(auc(test2_fpr, test2_tpr)))
18  plt.legend()
19  plt.xlabel("Fasle Postive Rate")
20  plt.ylabel("True Pstive RAte")
21  plt.title("ROC_AUC Curve")
22  plt.grid(color='black', linestyle='-', linewidth=1)
23  plt.show()
```

⤷

**Conclusion for SET : 3**

1. The above Represents the TRP and FPR rates on the either axis and this curve is know as the AUCROC curve , it is a metrix to evaluate the performar model.

2. As teh grap represents here the Train Auc = 0.67 , so from this we can conclude that the model is predecting the values with 67 % probabilty.

```
1  #this is the custom function for predecting the best thershold and sorting the values according the threshould
2  def find_best_threshold(threshould, fpr, tpr):
3      t = threshould[np.argmax(tpr*(1-fpr))]
4      # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
5      print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
6      return t
```

```
1  def predict_with_best_t(proba, threshould):
2      predictions = []
3      for i in proba:
4          if i>=threshould:
5              predictions.append(1)
6          else:
7              predictions.append(0)
8      return predictions
```

```
1   print("="*100)
2   from sklearn.metrics import confusion_matrix
3   best_t = find_best_threshold(tr2_thresholds, train2_fpr, train2_tpr)
4   Confusion_metrix_Train_data = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y3_train_pred, best_t)))
5   Confusion_metrix_Test_data = pd.DataFrame(confusion_matrix(Y_test, predict_with_best_t(y3_test_pred, best_t)))
6   import seaborn as sns
7   sns.set(font_scale=1.4)#for label size
8   print("The Confusion metrix of train data")
9   sns.heatmap(Confusion_metrix_Train_data,annot=True, annot_kws={"size": 20},fmt ="g")
10
```

```
1 print("The Confusion metrix of test data ")
2 sns.heatmap(Confusion_metrix_Test_data,annot=True, annot_kws={"size": 20},fmt ="g")
```

⊏→

1. The Ouput Stated above represents the maximun value of TRP (i.e tpr*(1-fpr) "0.401" corrosponding to whcih the maximum thershold is 0.48 .

2. The Second output represnts the Confusion metrix Based on the thershold of 0.48 , which states that the values below the thershold of 0.48 are clas 0 and the values above 0.48 are classified as 1.

## ▾  3.4 Set 4: Categorical, Numerical features + Project_title(TFIDF W2V)+ Preprocessed_essay (TFIDF W2V)

```
1 # TFIDF W2V of x_train "essays"
2 x_train_tfidf_model = TfidfVectorizer()
3 x_train_tfidf_model.fit(x_train["essay"])
4 # we are converting a dictionary with word as a key, and the idf as a value
5 dictionary = dict(zip(x_train_tfidf_model.get_feature_names(), list(x_train_tfidf_model.idf_)))
6 tfidf_words = set(x_train_tfidf_model.get_feature_names())
```

```
1  x_train1_tfidf_model = []; # the avg-w2v for each sentence/review is stored in this list
2  for sentence in tqdm(x_train['essay'].values): # for each review/sentence
3      vector = np.zeros(300) # as word vectors are of zero length
4      tf_idf_weight =0; # num of words with a valid vector in the sentence/review
5      for word in sentence.split(): # for each word in a review/sentence
6          if (word in glove_words) and (word in tfidf_words):
7              vec = model[word] # getting the vector for each word
8              # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
9              tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
10             vector += (vec * tf_idf) # calculating tfidf weighted w2v
11             tf_idf_weight += tf_idf
12     if tf_idf_weight != 0:
13         vector /= tf_idf_weight
14     x_train1_tfidf_model.append(vector)
15
16 print(len(x_train1_tfidf_model))
17 print(len(x_train1_tfidf_model[0]))
```

⊡→

```
1  # TFIDF W2V of X_CV "essays"
2  # here alos we will fit in the train data only as we dont want our data to be leaked
3  X_cv_tfidf_model = []; # the avg-w2v for each sentence/review is stored in this list
4  for sentence in tqdm(X_cv['essay'].values): # for each review/sentence
5      vector = np.zeros(300) # as word vectors are of zero length
6      tf_idf_weight =0; # num of words with a valid vector in the sentence/review
7      for word in sentence.split(): # for each word in a review/sentence
8          if (word in glove_words) and (word in tfidf_words):
9              vec = model[word] # getting the vector for each word
10             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
11             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
12             vector += (vec * tf_idf) # calculating tfidf weighted w2v
13             tf_idf_weight += tf_idf
14     if tf_idf_weight != 0:
15         vector /= tf_idf_weight
16     X_cv_tfidf_model.append(vector)
17
18 print(len(X_cv_tfidf_model))
19 print(len(X_cv_tfidf_model[0]))
```

⊡→

```python
1  X_test_tfidf_model = []; # the avg-w2v for each sentence/review is stored in this list
2  for sentence in tqdm(X_test['essay'].values): # for each review/sentence
3      vector = np.zeros(300) # as word vectors are of zero length
4      tf_idf_weight =0; # num of words with a valid vector in the sentence/review
5      for word in sentence.split(): # for each word in a review/sentence
6          if (word in glove_words) and (word in tfidf_words):
7              vec = model[word] # getting the vector for each word
8              # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
9              tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
10             vector += (vec * tf_idf) # calculating tfidf weighted w2v
11             tf_idf_weight += tf_idf
12     if tf_idf_weight != 0:
13         vector /= tf_idf_weight
14     X_test_tfidf_model.append(vector)
15
16 print(len(X_test_tfidf_model))
17 print(len(X_test_tfidf_model[0]))
```

⇱

```python
1  # TFIDF W2V of x_train "essays"
2  x_train_tfidf_model1 = TfidfVectorizer()
3  x_train_tfidf_model1.fit(x_train["clean_titles"])
4  # we are converting a dictionary with word as a key, and the idf as a value
5  dictionary = dict(zip(x_train_tfidf_model1.get_feature_names(), list(x_train_tfidf_model.idf_)))
6  tfidf_words = set(x_train_tfidf_model1.get_feature_names())
```

```python
1  x_train_tfidf_clean_titles_model = []; # the avg-w2v for each sentence/review is stored in this list
2  for sentence in tqdm(x_train['clean_titles'].values): # for each review/sentence
3      vector = np.zeros(300) # as word vectors are of zero length
4      tf_idf_weight =0; # num of words with a valid vector in the sentence/review
5      for word in sentence.split(): # for each word in a review/sentence
6          if (word in glove_words) and (word in tfidf_words):
7              vec = model[word] # getting the vector for each word
8              # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
9              tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
10             vector += (vec * tf_idf) # calculating tfidf weighted w2v
11             tf_idf_weight += tf_idf
12     if tf_idf_weight != 0:
13         vector /= tf_idf_weight
14     x_train_tfidf_clean_titles_model.append(vector)
15
16 print(len(x_train_tfidf_clean_titles_model))
17 print(len(x_train_tfidf_clean_titles_model[0]))
```

```python
X_cv_tfidf_clean_titles_model = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['clean_titles'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    X_cv_tfidf_clean_titles_model.append(vector)

print(len(X_cv_tfidf_clean_titles_model))
print(len(X_cv_tfidf_clean_titles_model[0]))
```

```python
X_test_tfidf_clean_titles_model = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['clean_titles'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    X_test_tfidf_clean_titles_model.append(vector)

print(len(X_test_tfidf_clean_titles_model))
print(len(X_test_tfidf_clean_titles_model[0]))
```

```python
1  from scipy.sparse import hstack
2  X4_tr = hstack((x_train_clean_categories_ohe,x_train_clean_subcat_ohe,x_train_teacher_ohe,x_train_state_ohe,\
3                  x_train_grade_ohe,x_train_price_std,x_train_projects_std,x_train_qty_std,x_train1_tfidf_model, x_train_tfidf_clean
4  X4_cv = hstack((X_cv_clean_categories_ohe,X_cv_clean_subcat_ohe,X_cv_teacher_ohe,X_cv_state_ohe,X_cv_grade_ohe,\
5                  X_cv_price_std,X_cv_projects_std,X_cv_qty_std,X_cv_tfidf_model,X_cv_tfidf_clean_titles_model)).tocsr()
6  X4_te =hstack((X_test_clean_categories_ohe,X_test_clean_subcat_ohe,X_test_teacher_ohe,X_test_state_ohe,\
7                  X_test_grade_ohe,X_test_price_std,X_test_projects_std,X_test_qty_std,X_test_tfidf_model,X_test_tfidf_clean_titles_r
8
9
10 print("The final Data Matrix for Set:4" , " All the shapes of the data represent the merged features as mentioned in the tittle")
11 print("shape of X_train is : ",            X4_tr.shape)
12 print("shape of X_Cross validation is :" , X4_cv.shape)
13 print("shape of X_test is ",               X4_te.shape)
```

⇥

```python
1  def batch_predict(clf, data):
2      y3_data_pred = []
3      tr_loop = data.shape[0] - data.shape[0]%1000
4      # consider you X_tr shape is 49041, then your tr_loop will be 49041 - 49041%1000 = 49000
5      # in this for loop we will iterate unti the last 1000 multiplier
6      for i in range(0, tr_loop, 1000):
7          y3_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
8      # we will be predicting for the last data points
9      if data.shape[0]%1000 !=0:
10         y3_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
11
12     return y3_data_pred
```

```python
1  import matplotlib.pyplot as plt
2  from sklearn.neighbors import KNeighborsClassifier
3  from sklearn.metrics import roc_auc_score
4
5  train_auc3 = []
6  cv_auc3 = []
7  K = [3, 15, 25, 51, 101]
8  for i in tqdm(K):
9      neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
10     neigh.fit(X4_tr, y_train)
11
12     y3_train_pred = batch_predict(neigh, X4_tr)
```

```
13      y3_cv_pred = batch_predict(neigh, X4_cv)
14
15      # roc_auc_score(y_tr, y_score) the 2nd parameter should be probability estimates of the positive class
16      # not the predicted outputs
17      train_auc3.append(roc_auc_score(y_train,y3_train_pred))
18      cv_auc3.append(roc_auc_score(Y_cv, y3_cv_pred))
19
20 plt.plot(K, train_auc3, label='Train AUC')
21 plt.plot(K, cv_auc3, label='CV AUC')
22
23 plt.scatter(K, train_auc3, label='Train AUC points')
24 plt.scatter(K, cv_auc3, label='CV AUC points')
25
26 plt.legend()
27 plt.xlabel("K: hyperparameter")
28 plt.ylabel("AUC")
29 plt.title("ERROR PLOTS")
30 plt.grid(color='black', linestyle='-', linewidth=1)
31 plt.show()
```

⤷

1.  The error plot above representes the Train AUC and CV AUC Curve, wherein we choose K(hyperparameter value) such that,it will be the maximum AU
    data and distance between the train-data line("Blue") and CV-line("orange") results be to minimum.
2.  The best value of k is found to be 102

3. Here I have used the "AUCROC" curve for choosing the best hypermeter as the data is imbalanced(though i have balanced it) , the "AUCROC" curve cc
   both the class labels equally .

4. The other hyperparameter tuning techiniques can also be choosen as "CV" and "K-fold CV"

```
1 best_k2= 101
```

```
1  # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2  from sklearn.metrics import roc_curve, auc
3
4
5  neigh = KNeighborsClassifier(n_neighbors=best_k2, n_jobs=-1)
6  neigh.fit(X4_tr, y_train)
7  # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
8  # not the predicted outputs
9
10 y4_train_pred = batch_predict(neigh, X4_tr)
11 y4_test_pred = batch_predict(neigh, X4_te)
12
13 train3_fpr, train3_tpr, tr3_thresholds = roc_curve(y_train, y4_train_pred)
14 test3_fpr, test3_tpr, te3_thresholds = roc_curve(Y_test, y4_test_pred)
15
16 plt.plot(train3_fpr, train3_tpr, label="Tain AUC ="+str(auc(train3_fpr, train3_tpr)))
17 plt.plot(test3_fpr, test3_tpr, label="Test AUC ="+str(auc(test3_fpr, test3_tpr)))
18 plt.legend()
19 plt.xlabel("Fasle Postive Rate")
20 plt.ylabel("True Pstive RAte")
21 plt.title("ROC_AUC Curve")
22 plt.grid(color='black', linestyle='-', linewidth=1)
23 plt.show()
```

⤷

**Conclusion for SET : 4**

1. The above Represents the TRP and FPR rates on the either axis and this curve is know as the AUCROC curve , it is a metrix to evaluate the performar model.
2. As teh grap represents here the Train Auc = 0.67 , so from this we can conclude that the model is predecting the values with 67 % probabilty.

```
1  #this is the custom function for predecting the best thershold and sorting the values according the threshould
2  def find_best_threshold(threshould, fpr, tpr):
3      t = threshould[np.argmax(tpr*(1-fpr))]
4      # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
5      print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
6      return t
```

```
1
2  def predict_with_best_t(proba, threshould):
3      predictions = []
4      for i in proba:
5          if i>=threshould:
6              predictions.append(1)
7          else:
8              predictions.append(0)
9      return predictions
```

```
1  print("="*100)
2  from sklearn.metrics import confusion_matrix
3  best_t = find_best_threshold(tr3_thresholds, train3_fpr, train3_tpr)
4  Confusion_metrix_Train_data = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y4_train_pred, best_t)))
5  Confusion_metrix_Test_data = pd.DataFrame(confusion_matrix(Y_test, predict_with_best_t(y4_test_pred, best_t)))
6  import seaborn as sns
7  sns.set(font_scale=1.4)#for label size
8  print("The Confusion metrix of train data")
9  sns.heatmap(Confusion_metrix_Train_data,annot=True, annot_kws={"size": 20},fmt ="g")
```

⮐

```
1 print("The Confusion metrix of test data ")
2 sns.heatmap(Confusion_metrix_Test_data,annot=True, annot_kws={"size": 20},fmt ="g")
```

⤷

1. The Ouput Stated above represents the maximun value of TRP (i.e tpr*(1-fpr) "0.3939" corrosponding to whcih the maximum thershold is 0.45 .

2. The Second output represnts the Confusion metrix Based on the thershold of 0.45 , which states that the values below the thershold of 0.45 are clas
   0 and the values above 0.48 are classified as 1.

# 3.5 Applying the SelectKBest for selecting the Top best 2000 features on the set2 of the Data matrix b text

This data metrix consits of Set 2: Categorical, Numerical features + Project_title(TFIDF)+ Preprocessed_essay (TFIDF)

```
1  # Preprocessed_essay (TFIDF) and selcting the best 2000 Features using the  SelectKBest, chi2
2  %%time
3  from sklearn.feature_extraction.text import TfidfVectorizer
4  from sklearn.feature_selection import SelectKBest, chi2
5  vectorizer = TfidfVectorizer(min_df=10)
6  vectorizer.fit(x_train['essay'].values) # fit has to happen only on train data
7
8  # we use the fitted CountVectorizer to convert the text to vector
9  x_train_essay_tfidf = vectorizer.transform(x_train['essay'].values)
10 X_cv_essay_tfidf = vectorizer.transform(X_cv['essay'].values)
11 X_test_essay_tfidf = vectorizer.transform(X_test['essay'].values)
12
13
14 #Selecting top 2000 best features from the generated tfidf features
15 selector = SelectKBest(chi2, k = 2000 )
16 selector.fit(x_train_essay_tfidf,y_train)
17 x_train_essay_2000 = selector.transform(x_train_essay_tfidf)
18 X_cv_essay_2000 = selector.transform(X_cv_essay_tfidf)
19 X_test_essay_2000 = selector.transform(X_test_essay_tfidf)
20 print(x_train_essay_2000.shape)
21 print(X_cv_essay_2000.shape)
22 print(X_test_essay_2000.shape)
```

⬐

```
1  # Project_title(TFIDF) and selcting the best 2000 Features using the  SelectKBest, chi2
2  %%time
```

```
 3  from sklearn.feature_extraction.text import TfidfVectorizer
 4  from sklearn.feature_selection import SelectKBest, chi2
 5  vectorizer = TfidfVectorizer(min_df=10)
 6  vectorizer.fit(x_train['clean_titles'].values) # fit has to happen only on train data
 7
 8  # we use the fitted CountVectorizer to convert the text to vector
 9  x_train_clean_titles_tfidf = vectorizer.transform(x_train['clean_titles'].values)
10  X_cv_clean_titles_tfidf = vectorizer.transform(X_cv['clean_titles'].values)
11  X_test_clean_titles_tfidf = vectorizer.transform(X_test['clean_titles'].values)
12
13
14  #Selecting top 2000 best features from the generated tfidf features
15  selector = SelectKBest(chi2, k = 1811 )
16  selector.fit(x_train_clean_titles_tfidf,y_train)
17  x_train_clean_titles_2000 = selector.transform(x_train_clean_titles_tfidf)
18  X_cv_clean_titles_2000 = selector.transform(X_cv_clean_titles_tfidf)
19  X_test_clean_titles_2000 = selector.transform(X_test_clean_titles_tfidf)
20  print(x_train_clean_titles_2000.shape)
21  print(X_cv_clean_titles_2000.shape)
22  print(X_test_clean_titles_2000.shape)
```

⇥

```
 1  #now merging all the data matrix
 2  from scipy.sparse import hstack
 3  X5_tr = hstack((x_train_clean_categories_ohe,x_train_clean_subcat_ohe,x_train_teacher_ohe,x_train_state_ohe,\
 4             x_train_grade_ohe,x_train_price_std,x_train_projects_std,x_train_qty_std,x_train_essay_2000,x_train_clean_titles_2
 5  X5_cv = hstack((X_cv_clean_categories_ohe,X_cv_clean_subcat_ohe,X_cv_teacher_ohe,X_cv_state_ohe,\
 6             X_cv_grade_ohe,X_cv_price_std,X_cv_projects_std,X_cv_qty_std,X_cv_essay_2000,X_cv_clean_titles_2000)).tocsr()
 7  X5_te =hstack((X_test_clean_categories_ohe,X_test_clean_subcat_ohe,X_test_teacher_ohe,X_test_state_ohe,\
 8             X_test_grade_ohe,X_test_price_std,X_test_projects_std,X_test_qty_std,X_test_essay_2000,X_test_clean_titles_2000)).
 9
10
11  print("The final Data Matrix for Set:3" , " All the shapes of the data represent the merged features as mentioned in the tittle")
12  print("shape of X_train is : ",            X5_tr.shape)
13  print("shape of X_Cross validation is :" , X5_cv.shape)
14  print("shape of X_test is ",               X5_te.shape)
```

⇥

```python
1  def batch_predict(clf, data):
2      y4_data_pred = []
3      tr_loop = data.shape[0] - data.shape[0]%1000
4      # consider you X_tr shape is 49041, then your tr_loop will be 49041 - 49041%1000 = 49000
5      # in this for loop we will iterate unti the last 1000 multiplier
6      for i in range(0, tr_loop, 1000):
7          y4_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
8      # we will be predicting for the last data points
9      if data.shape[0]%1000 !=0:
10         y4_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
11
12     return y4_data_pred
```

```python
1  import matplotlib.pyplot as plt
2  from sklearn.neighbors import KNeighborsClassifier
3  from sklearn.metrics import roc_auc_score
4
5  train_auc4 = []
6  cv_auc4 = []
7  K = [3, 15, 25, 51, 101]
8  for i in tqdm(K):
9      neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
10     neigh.fit(X5_tr, y_train)
11
12     y4_train_pred = batch_predict(neigh, X5_tr)
13     y4_cv_pred = batch_predict(neigh, X5_cv)
14
15     # roc_auc_score(y_tr, y_score) the 2nd parameter should be probability estimates of the positive class
16     # not the predicted outputs
17     train_auc4.append(roc_auc_score(y_train,y4_train_pred))
18     cv_auc4.append(roc_auc_score(Y_cv, y4_cv_pred))
19
20 plt.plot(K, train_auc4, label='Train AUC')
21 plt.plot(K, cv_auc4, label='CV AUC')
22
23 plt.scatter(K, train_auc4, label='Train AUC points')
24 plt.scatter(K, cv_auc4, label='CV AUC points')
25
26 plt.legend()
27 plt.xlabel("K: hyperparameter")
28 plt.ylabel("AUC")
29 plt.title("ERROR PLOTS")
30 plt.grid(color='black', linestyle='-', linewidth=1)
31 plt.show()
```

1. The error plot above representes the Train AUC and CV AUC Curve, wherein we choose K(hyperparameter value) such that,it will be the maximum AU data and distance between the train-data line("Blue") and CV-line("orange") results be to minimum.
2. The best value of k is found to be 100
3. Here I have used the "AUCROC" curve for choosing the best hypermeter as the data is imbalanced(though i have balanced it) , the "AUCROC" curve co both the class labels equally .
4. The other hyperparameter tuning techiniques can also be choosen as "CV" and "K-fold CV"

```
1 best_k2= 101
```

```
1  # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2  from sklearn.metrics import roc_curve, auc
3
4
5  neigh = KNeighborsClassifier(n_neighbors=best_k2, n_jobs=-1)
6  neigh.fit(X5_tr, y_train)
7  # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
8  # not the predicted outputs
9
10 y5_train_pred = batch_predict(neigh, X5_tr)
11 y5_test_pred = batch_predict(neigh, X5_te)
12
13 train5_fpr, train5_tpr, tr5_thresholds = roc_curve(y_train, y5_train_pred)
14 test5_fpr, test5_tpr, te5_thresholds = roc_curve(Y_test, y5_test_pred)
15
16 plt.plot(train5_fpr, train5_tpr, label="Tain AUC ="+str(auc(train5_fpr, train5_tpr)))
```
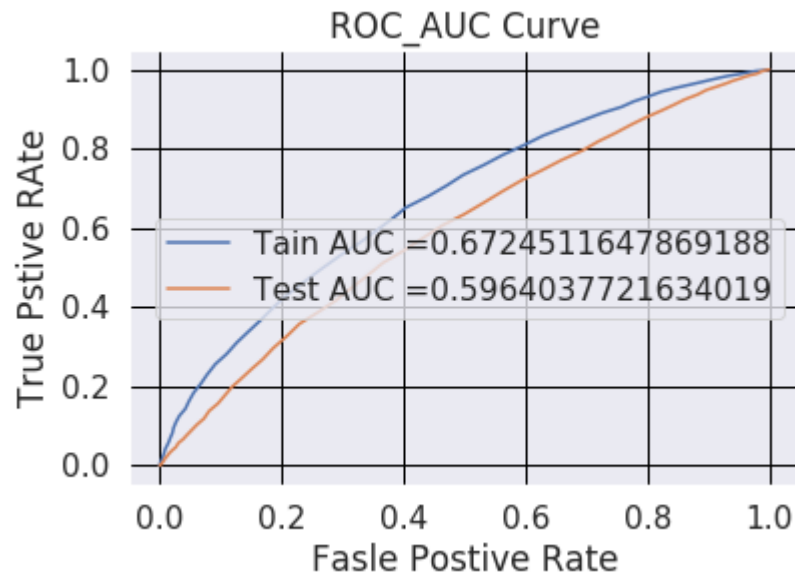
```
17 plt.plot(test5_fpr, test5_tpr, label="Test AUC ="+str(auc(test5_fpr, test5_tpr)))
18 plt.legend()
19 plt.xlabel("Fasle Postive Rate")
20 plt.ylabel("True Pstive RAte")
21 plt.title("ROC_AUC Curve")
22 plt.grid(color='black', linestyle='-', linewidth=1)
23 plt.show()
```



**Conclusion for SET : 5**

1. The above Represents the TRP and FPR rates on the either axis and this curve is know as the AUCROC curve , it is a metrix to evaluate the performar model.

2. As teh grap represents here the Train Auc = 0.67 , so from this we can conclude that the model is predecting the values with 67 % probabilty.

```
1 #this is the custom function for predecting the best thershold and sorting the values according the threshould
2 def find_best_threshold(threshould, fpr, tpr):
3     t = threshould[np.argmax(tpr*(1-fpr))]
4     # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
5     print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
6     return t
```

```
1 def predict_with_best_t(proba, threshould):
2     predictions = []
```
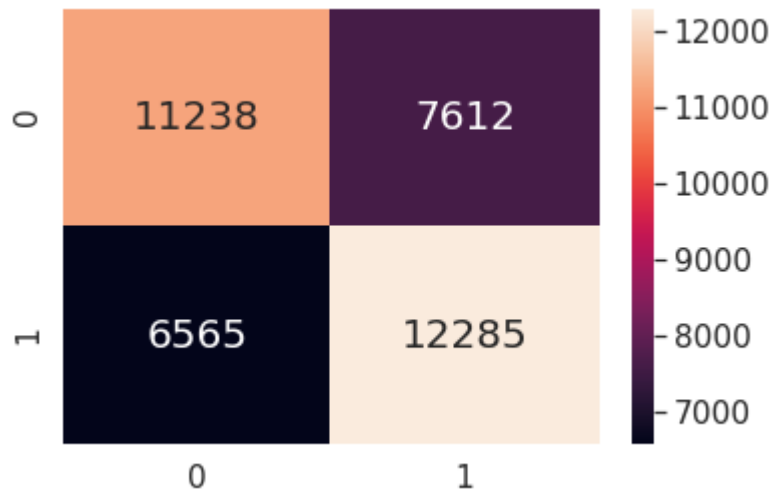
```
3      for i in proba:
4          if i>=threshould:
5              predictions.append(1)
6          else:
7              predictions.append(0)
8      return predictions
```

```
1  print("="*100)
2  from sklearn.metrics import confusion_matrix
3  best_t = find_best_threshold(tr5_thresholds, train5_fpr, train5_tpr)
4  Confusion_metrix_Train_data = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y5_train_pred, best_t)))
5  Confusion_metrix_Test_data = pd.DataFrame(confusion_matrix(Y_test, predict_with_best_t(y5_test_pred, best_t)))
6  import seaborn as sns
7  sns.set(font_scale=1.4)#for label size
8  print("The Confusion metrix of train data")
9  sns.heatmap(Confusion_metrix_Train_data,annot=True, annot_kws={"size": 20},fmt ="g")
```

```
====================================================================================================
the maximum value of tpr*(1-fpr) 0.38854513857129785 for threshold 0.475
The Confusion metrix of train data
<matplotlib.axes._subplots.AxesSubplot at 0x7fd41dba3cc0>
```
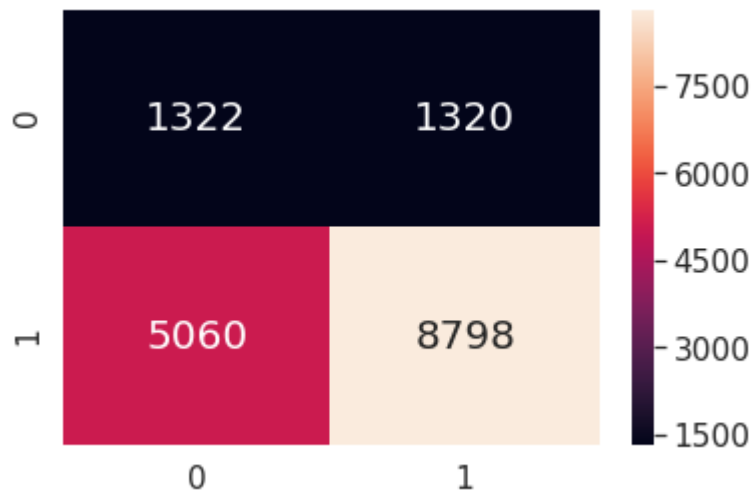


```
1  print("The Confusion metrix of test data ")
2  sns.heatmap(Confusion_metrix_Test_data,annot=True, annot_kws={"size": 20},fmt ="g")
```

```
The Confusion metrix of test data
<matplotlib.axes._subplots.AxesSubplot at 0x7fd41c69ff28>
```



1. The Ouput Stated above represents the maximun value of TRP (i.e tpr*(1-fpr) "0.3973" corrosponding to whcih the maximum thershold is 0.49 .
2. The Second output represnts the Confusion metrix Based on the thershold of 0.49 , which states that the values below the thershold of 0.49 are clas 0 and the values above 0.49 are classified as 1.

```
1  #http://zetcode.com/python/prettytable/
2  from prettytable import PrettyTable
3
4  x = PrettyTable()
5  x.field_names = ["Vectorizer", "Hyperparameter", "AUC"]
6  x.add_row(["SET 1 : project_title(BOW) + preprocessed_essay (BOW)", 101 ,  0.55])
7  x.add_row(["SET 2 : project_title(TFIDF)+ preprocessed_essay (TFIDF)", 101 , 0.59 ])
8  x.add_row(["SET 3 :  project_title(AVG W2V)+ preprocessed_essay (AVG W2V)", 101, 0.59])
9  x.add_row(["SET 4 : project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)", 101, 0.60])
10 x.add_row(["SET 5 : project_title(TFIDF)+ preprocessed_essay (TFIDF) of 2000 Words", 100, 0.59])
11
12 print(x)
```

```
+--------------------------------------------------------------+---------------+------+
|                          Vectorizer                          | Hyperparameter | AUC  |
+--------------------------------------------------------------+---------------+------+
|          SET 1 : project_title(BOW) + preprocessed_essay (BOW)          |      101      | 0.55 |
|           SET 2 : project_title(TFIDF)+ preprocessed_essay (TFIDF)       |      101      | 0.59 |
|          SET 3 :  project_title(AVG W2V)+ preprocessed_essay (AVG W2V)   |      101      | 0.59 |
|        SET 4 : project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)  |      101      | 0.6  |
| SET 5 : project_title(TFIDF)+ preprocessed_essay (TFIDF) of 2000 Words |      100      | 0.59 |
+--------------------------------------------------------------+---------------+------+
```