

▼ DonorsChoose

DonorsChoose.org receives hundreds of thousands of project proposals each year for classroom number of volunteers is needed to manually screen each submission before it's approved to be posted. Next year, DonorsChoose.org expects to receive close to 500,000 project proposals. As a result, they need to solve:

- How to scale current manual processes and resources to screen 500,000 projects so that they can be as possible
- How to increase the consistency of project vetting across different volunteers to improve the quality of the projects
- 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. The text of project descriptions as well as additional metadata about the project, teacher, and school 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>	Title of the project. Examples: <ul style="list-style-type: none"> • Art Will Make You Happy! • First Grade Fun
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following: <ul style="list-style-type: none"> • Grades PreK-2 • Grades 3-5 • Grades 6-8 • Grades 9-12
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project. <ul style="list-style-type: none"> • Applied Learning • Care & Hunger • Health & Sports • History & Civics • Literacy & Language • Math & Science • Music & The Arts • Special Needs • Warmth Examples: <ul style="list-style-type: none"> • Music & The Arts • Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (Two-letter U.S. postal code). Example: CA

Feature	Description
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. <ul style="list-style-type: none"> • Literacy • Literature & Writing, Social Science
<code>project_resource_summary</code>	An explanation of the resources needed for the project. Example: <ul style="list-style-type: none"> • My students need hands on literacy r
<code>project_essay_1</code>	First application essay [*]
<code>project_essay_2</code>	Second application essay [*]
<code>project_essay_3</code>	Third application essay [*]
<code>project_essay_4</code>	Fourth application essay [*]
<code>project_submitted_datetime</code>	Datetime when project application was submitted. Example: 2016
<code>teacher_id</code>	A unique identifier for the teacher of the proposed project. Examp
<code>teacher_prefix</code>	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> • nan • Dr. • Mr. • Mrs. • Ms. • Teacher.
<code>teacher_number_of_previously_posted_projects</code>	Number of project applications previously submitted by the same

^{*} 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 resource required by a project:

Feature	Description
<code>id</code>	A <code>project_id</code> value from the <code>train.csv</code> file. Example: p0365
<code>description</code>	Description of the resource. Example: Tenor Saxophone Reeds,
<code>quantity</code>	Quantity of the resource required. Example: 3
<code>price</code>	Price of the resource required. Example: 9.95

Note: Many projects require multiple resources. The `id` value corresponds to a `project_id` in the `resources` needed for a project:

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

Label	Description
<code>project_is_approved</code>	A binary flag indicating whether DonorsChoose approved the project. A value of 0 indicates the project

▼ 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 f following:

- __project_essay_1:__ "Describe your students: What makes your students special? Specific c neighborhood, and your school are all helpful."
- __project_essay_2:__ "About your project: How will these materials make a difference in your lives?"

For all projects with project_submitted_datetime of 2016-05-17 and later, the values of proje

```
# importing required libraries

%matplotlib inline
import warnings
warnings.filterwarnings("ignore")

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os

import chart_studio.plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
```

```
from sklearn.model_selection import GridSearchCV
```



▼ 1.1 Reading Data

```
from google.colab import drive
```

```
# This will prompt for authorization.
drive.mount('/content/drive', force_remount=True)
```

```
↳ Mounted at /content/drive
```

```
# Reading data from project and resources data file
```

```
project_data = pd.read_csv('/content/drive/My Drive/Colab Notebooks/Dataset/Assign
resource_data = pd.read_csv('/content/drive/My Drive/Colab Notebooks/Dataset/Assig
```

```
# Getting basic information about the data
```

```
print("Number of data points in Project_train data", project_data.shape)
print('- '*100)
print("The attributes of Project_train data :", project_data.columns.values)
print('='*100)
print("Number of data points in Resource_train data", resource_data.shape)
print('- '*100)
print("The attributes of Resource_train data :", resource_data.columns.values)
```

```
↳ Number of data points in Project_train data (109248, 17)
```

```
-----
The attributes of Project_train data : ['Unnamed: 0' 'id' 'teacher_id' 'teach
'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']
=====
```

```
Number of data points in Resource_train data (1541272, 4)
```

```
-----
The attributes of Resource_train data : ['id' 'description' 'quantity' 'price
```

▼ 1.2 Data Pre-Processing

```
# Merge two column text dataframe:
```

```
# Merge 4 essays into one:
```

```
project_data["essay"] = project_data["project_essay_1"].map(str) + \
                        project_data["project_essay_2"].map(str) + \
                        project_data["project_essay_3"].map(str) + \
                        project_data["project_essay_4"].map(str)
```

```
# Merge Price information from resource data to project data
price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).re
project_data = pd.merge(project_data, price_data, on='id', how='left')

# find how many digits are present in each project_resource_summary coloumn
summary = list(project_data['project_resource_summary'].values)
presence_of_numeric_data=[]
for i in summary:
    count = 0
    for j in i.split(' '):
        if j.isdigit():
            count+=1
    presence_of_numeric_data.append(count)

# Replace Text summary coloumn with new numerical coloumn presence_of_numeric_data
project_data['numerical_data_in_resource_summary'] = presence_of_numeric_data
project_data.drop(['project_resource_summary'], axis=1, inplace=True)

# how to replace elements in list python: https://stackoverflow.com/a/2582163/4084
cols = ['Date' if x=='project_submitted_datetime' else x for x in list(project_dat

#sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/
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)

# how to reorder columns pandas python: https://stackoverflow.com/a/13148611/40840
project_data = project_data[cols]

# https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.drop
# Here we drop 3 rows where teacher_prefix is having np.nan value
project_data.dropna(axis=0, subset=['teacher_prefix'], inplace=True)

project_data.head(2)
```

```
↳
```

	Unnamed: 0	id	teacher_id	teacher_prefix	school
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	

▼ 1.2.1 Pre-Processing Essay Text

```
# printing some random essays.
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
```

```
☞ I have been fortunate enough to use the Fairy Tale STEM kits in my classroom
=====
I teach high school English to students with learning and behavioral disabili
=====
```

```
# https://stackoverflow.com/a/47091490/4084039
import re
```

```
def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"\ 're", " are", phrase)
    phrase = re.sub(r"\ 's", " is", phrase)
    phrase = re.sub(r"\ 'd", " would", phrase)
    phrase = re.sub(r"\ 'll", " will", phrase)
    phrase = re.sub(r"\ 't", " not", phrase)
    phrase = re.sub(r"\ 've", " have", phrase)
    phrase = re.sub(r"\ 'm", " am", phrase)
    return phrase
```

```
# https://gist.github.com/sebleier/554280
```

```
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "y
    'you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'h
    'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself
    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that'
    'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has'
    'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because',
    'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'thr
    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off'
    'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all',
    'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than',
    's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've"
    've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "did
    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma',
    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't
    'won', "won't", 'wouldn', "wouldn't"]
```

```
from tqdm import tqdm
preprocessed_essays = []
# tqdm is for printing the status bar
```

```

for sentence in tqdm(project_data['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\r', ' ')
    sent = sent.replace('\n', ' ')
    sent = sent.replace('\t', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    preprocessed_essays.append(sent.lower().strip())

```

Adding preprocessed_essays column to our data matrix

```
project_data['preprocessed_essays']=preprocessed_essays
```

```

↳ 100%|██████████| 109245/109245 [01:03<00:00, 1714.06it/s]

```

```

# after preproceasing
preprocessed_essays[100]

```

```

↳ 'a typical day campus exciting my students love learning always put smile fac

```

▼ 1.2.2 Pre-Processing Project Title Text

```

from tqdm import tqdm
preprocessed_titles = []
# tqdm is for printing the status bar
for title in tqdm(project_data['project_title'].values):
    title = decontracted(title)
    title = title.replace('\r', ' ')
    title = title.replace('\n', ' ')
    title = title.replace('\t', ' ')
    title = re.sub('[^A-Za-z0-9]+', ' ', title)
    # https://gist.github.com/sebleier/554280
    title = ' '.join(e for e in title.split() if e not in stopwords)
    preprocessed_titles.append(title.lower().strip())

```

Adding preprocessed_titles column to our data matrix

```
project_data['preprocessed_titles']=preprocessed_titles
preprocessed_titles[1000]
```

```

↳ 100%|██████████| 109245/109245 [00:02<00:00, 41023.70it/s]
'empowering students through art learning about then now'

```

▼ 1.2.3 Pre-Processing Project Grades

```

# Remove special characters from grades
from tqdm import tqdm
preprocessed_grade_categories = []
# tadm is for printing the status bar

```

```

for categories in tqdm(project_data['project_grade_category'].values):
    categories = decontracted(categories)
    # https://gist.github.com/sebleier/554280
    categories = '_'.join(e for e in categories.split(' ') if e not in stopwords)
    categories = '_'.join(e for e in categories.split('-') if e not in stopwords)
    preprocessed_grade_categories.append(categories.lower().strip())

```

```
# Adding preprocessed_titles coloumn to our data matrix
```

```
project_data['preprocessed_grade_category']=preprocessed_grade_categories
```

```
project_data.head(5)
```

```

[ ] 100%|██████████| 109245/109245 [00:02<00:00, 52758.59it/s]

```

	Unnamed: 0	id	teacher_id	teacher_prefix	school
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5		Mrs.
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df		Ms.
51140	74477	p189804	4a97f3a390bfe21b99cf5e2b81981c73		Mrs.
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3		Mrs.
41558	33679	p137682	06f6e62e17de34fcf81020c77549e1d5		Mrs.

▼ 1.2.4 preprocessing of project_subject_categories

```
catogories = list(project_data['project_subject_categories'].values)
```

```
# remove special characters from list of strings python: https://stackoverflow.com
```



```
# remove special characters from list of strings python: https://stackoverflow.com

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in
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", "W
        if 'The' in j.split(): # this will split each of the catogory based on spa
            j=j.replace('The','') # if we have the words "The" we are going to rep
            j = j.replace(' ', '') # we are placeing all the ' '(space) with ''(empty)
            temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailin
            temp = temp.replace('&','_') # we are replacing the & value into
    cat_list.append(temp.strip())

project_data['clean_categories'] = cat_list
project_data.drop(['project_subject_categories'], axis=1, inplace=True)
```

▼ 1.2.5 preprocessing of project_subject_subcategories

```
sub_catogories = list(project_data['project_subject_subcategories'].values)
# remove special characters from list of strings python: https://stackoverflow.com

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in

sub_cat_list = []
for i in sub_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", "W
        if 'The' in j.split(): # this will split each of the catogory based on spa
            j=j.replace('The','') # if we have the words "The" we are going to rep
            j = j.replace(' ', '') # we are placeing all the ' '(space) with ''(empty)
            temp +=j.strip()+" " # " abc ".strip() will return "abc", remove the trailin
            temp = temp.replace('&','_')
    sub_cat_list.append(temp.strip())

project_data['clean_subcategories'] = sub_cat_list
project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)

# Drop all unnecessary featurus like project_grade_category, project_essay_1, etc.
project_data.drop(['project_grade_category'], axis=1, inplace=True)
project_data.drop(['project_essay_1'], axis=1, inplace=True)
project_data.drop(['project_essay_2'], axis=1, inplace=True)
project_data.drop(['project_essay_3'], axis=1, inplace=True)
project_data.drop(['project_essay_4'], axis=1, inplace=True)
project_data.drop(['essay'], axis=1, inplace=True)
```

```
project_data.head(5)
```

	Unnamed: 0	id	teacher_id	teacher_prefix	school
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5		Mrs.
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df		Ms.
51140	74477	p189804	4a97f3a390bfe21b99cf5e2b81981c73		Mrs.
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3		Mrs.
41558	33679	p137682	06f6e62e17de34fcf81020c77549e1d5		Mrs.

▼ 1.3 Sampling data for KNN Assignment

```
project_data['project_is_approved'].value_counts()
```

```
1    92703
0    16542
Name: project_is_approved, dtype: int64
```

```
# take first 50k points from project_data
data = project_data.head(50000)
data['project_is_approved'].value_counts()
```

```
1    41992
0    8008
Name: project_is_approved, dtype: int64
```

```
data.head(5)
```

	Unnamed: 0		id	teacher_id	teacher_prefix	school
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5		Mrs.	
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df		Ms.	
51140	74477	p189804	4a97f3a390bfe21b99cf5e2b81981c73		Mrs.	
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3		Mrs.	
41558	33679	p137682	06f6e62e17de34fcf81020c77549e1d5		Mrs.	

```
# Split the class label from data
y = data['project_is_approved'].values
X = data.drop(['project_is_approved'], axis=1)
X.head(1)
```



	Unnamed: 0		id	teacher_id	teacher_prefix	school
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5		Mrs.	

2.1 Splitting data into Train and cross validation(or test): Stratified Sampling

```
# train test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify=y)
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train)
```

```
# Simple Upsampling for negative class data points in training dataset
# https://www.kaggle.com/rafjaa/resampling-strategies-for-imbalanced-datasets
```

```

from sklearn.utils import resample
#df3 = pd.DataFrame(y_train,columns=['project_is_approved'],dtype = int)

#X = pd.concat([X_train,df3],axis = 1)
X_train['project_is_approved']=y_train
Accepted, Rejected = X_train.project_is_approved.value_counts()

# Divide by class
df_class_0 = X_train[X_train['project_is_approved'] == 0]
df_class_1 = X_train[X_train['project_is_approved'] == 1]

upsampled_data = df_class_0.sample(Accepted, replace=True,)
X_train = pd.concat([df_class_1, upsampled_data], axis=0)
print(X_train.project_is_approved.value_counts())

↳ 1      18850
   0      18850
   Name: project_is_approved, dtype: int64

y_train = X_train.project_is_approved
X_train = X_train.drop('project_is_approved', axis=1)
X_train.shape

↳ (37700, 16)

```

▼ 2.2 Make Data Model Ready:

▼ 2.2.1 Encoding numerical, categorical features

▼ 2.2.1.1 Encoding School State

```

# Encoding School State

vectorizer = CountVectorizer()
vectorizer.fit(X_train['school_state'].values) # fit has to happen only on train d

# we use the fitted CountVectorizer to convert the text to vector
X_train_state_ohe = vectorizer.transform(X_train['school_state'].values)
X_cv_state_ohe = vectorizer.transform(X_cv['school_state'].values)
X_test_state_ohe = vectorizer.transform(X_test['school_state'].values)

print("After vectorizations")
print(X_train_state_ohe.shape, y_train.shape)
print(X_cv_state_ohe.shape, y_cv.shape)
print(X_test_state_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

```

```

↳ After vectorizations
(37700, 51) (37700,)
(11055, 51) (11055,)
(16500, 51) (16500,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia']
=====

```

▼ 2.2.1.2 Encoding Teacher Prefix

```

vectorizer = CountVectorizer()
vectorizer.fit(X_train['teacher_prefix'].values) # fit has to happen only on train

# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_ohe = vectorizer.transform(X_train['teacher_prefix'].values)
X_cv_teacher_ohe = vectorizer.transform(X_cv['teacher_prefix'].values)
X_test_teacher_ohe = vectorizer.transform(X_test['teacher_prefix'].values)

print("After vectorizations")
print(X_train_teacher_ohe.shape, y_train.shape)
print(X_cv_teacher_ohe.shape, y_cv.shape)
print(X_test_teacher_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

```

```

↳ After vectorizations
(37700, 5) (37700,)
(11055, 5) (11055,)
(16500, 5) (16500,)
['dr', 'mr', 'mrs', 'ms', 'teacher']
=====

```

▼ 2.2.1.3 Encoding preprocessed_grade_category

```

vectorizer = CountVectorizer()
vectorizer.fit(X_train['preprocessed_grade_category'].values) # fit has to happen

# we use the fitted CountVectorizer to convert the text to vector
X_train_grade_ohe = vectorizer.transform(X_train['preprocessed_grade_category'].va
X_cv_grade_ohe = vectorizer.transform(X_cv['preprocessed_grade_category'].values)
X_test_grade_ohe = vectorizer.transform(X_test['preprocessed_grade_category'].valu

print("After vectorizations")
print(X_train_grade_ohe.shape, y_train.shape)
print(X_cv_grade_ohe.shape, y_cv.shape)
print(X_test_grade_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

```

```

↳

```

```

After vectorizations
(37700, 4) (37700,)
(11055, 4) (11055,)
(16500, 4) (16500,)
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
=====

```

▼ 2.2.1.4 Encoding numerical feature Price

```

from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['price'].values.reshape(1,-1))

X_train_price_norm = normalizer.transform(X_train['price'].values.reshape(1,-1))
X_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(1,-1))
X_test_price_norm = normalizer.transform(X_test['price'].values.reshape(1,-1))

X_train_price_norm = X_train_price_norm.reshape(-1,1)
X_test_price_norm = X_test_price_norm.reshape(-1,1)
X_cv_price_norm = X_cv_price_norm.reshape(-1,1)

print("After vectorizations")
print(X_train_price_norm.shape, y_train.shape)
print(X_train_price_norm)
print(X_cv_price_norm.shape, y_cv.shape)
print(X_test_price_norm.shape, y_test.shape)
print("="*100)

```

```

☞ After vectorizations
(37700, 1) (37700,)
[[0.00472546]
 [0.00013186]
 [0.00149878]
 ...
 [0.00131354]
 [0.00208145]
 [0.00672036]]
(11055, 1) (11055,)
(16500, 1) (16500,)
=====

```

▼ 2.2.1.5 Encoding numeric feature Quantity

```

from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)

```

```
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['quantity'].values.reshape(1,-1))

X_train_quantity_norm = normalizer.transform(X_train['quantity'].values.reshape(1,
X_train_quantity_norm = X_train_quantity_norm.reshape(-1,1)
X_cv_quantity_norm = normalizer.transform(X_cv['quantity'].values.reshape(1,-1))
X_cv_quantity_norm = X_cv_quantity_norm.reshape(-1,1)
X_test_quantity_norm = normalizer.transform(X_test['quantity'].values.reshape(1,-1)
X_test_quantity_norm = X_test_quantity_norm.reshape(-1,1)
print(X_train_quantity_norm)
print("After vectorizations")
print(X_train_quantity_norm.shape, y_train.shape)
print(X_cv_quantity_norm.shape, y_cv.shape)
print(X_test_quantity_norm.shape, y_test.shape)
print("=="*100)
```

```
↳ [[0.00170828]
    [0.0042707 ]
    [0.00170828]
    ...
    [0.00113885]
    [0.00085414]
    [0.00370127]]
After vectorizations
(37700, 1) (37700,)
(11055, 1) (11055,)
(16500, 1) (16500,)
=====
```

▼ 2.2.1.6 Encoding numeric feature teacher_number_of_previously_posted_projects

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
#List_of_imp_features.append('teacher_number_of_previously_posted_projects')
X_train_teacher_number_of_previously_posted_projects_norm = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
X_cv_teacher_number_of_previously_posted_projects_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
X_test_teacher_number_of_previously_posted_projects_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
X_train_teacher_number_of_previously_posted_projects_norm = X_train_teacher_number_of_previously_posted_projects_norm.reshape(-1,1)
X_test_teacher_number_of_previously_posted_projects_norm = X_test_teacher_number_of_previously_posted_projects_norm.reshape(-1,1)
X_cv_teacher_number_of_previously_posted_projects_norm=X_cv_teacher_number_of_previously_posted_projects_norm.reshape(-1,1)

print(X_test_teacher_number_of_previously_posted_projects_norm)
print("After vectorizations")
```

```
print(X_train_teacher_number_of_previously_posted_projects_norm.shape, y_train.shape)
```

```
print(X_train_teacher_number_of_previously_posted_projects_norm.shape, y_train.shape)
print(X_cv_teacher_number_of_previously_posted_projects_norm.shape, y_cv.shape)
print(X_test_teacher_number_of_previously_posted_projects_norm.shape, y_test.shape)
print("="*100)
```

```
↳ [[0.00091876]
    [0.00245002]
    [0.00091876]
    ...
    [0.         ]
    [0.         ]
    [0.         ]]
After vectorizations
(37700, 1) (37700,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

▼ 2.2.1.7 Encoding numeric feature numerical_data_in_resource_summary

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['numerical_data_in_resource_summary'].values.reshape(1, -1))
X_train_numerical_data_in_resource_summary_norm = normalizer.transform(X_train['nu
X_cv_numerical_data_in_resource_summary_norm = normalizer.transform(X_cv['numera
X_test_numerical_data_in_resource_summary_norm = normalizer.transform(X_test['nume
X_cv_numerical_data_in_resource_summary_norm = X_cv_numerical_data_in_resource_sum
X_train_numerical_data_in_resource_summary_norm = X_train_numerical_data_in_resour
X_test_numerical_data_in_resource_summary_norm = X_test_numerical_data_in_resource

print(X_test_numerical_data_in_resource_summary_norm)
print("After vectorizations")
print(X_train_numerical_data_in_resource_summary_norm.shape, y_train.shape)
#print(X_cv_numerical_data_in_resource_summary_norm.shape, y_cv.shape)
print(X_test_numerical_data_in_resource_summary_norm.shape, y_test.shape)
print("="*100)
```

```
↳ [[0.]
    [0.]
    [0.]
    ...
    [0.]
    [0.]
    [0.]]
After vectorizations
(37700, 1) (37700,)
(16500, 1) (16500,)
```


▼ 2.3 Applying KNN on different kind of featurization as mentio

```
# Define Functions for Train KNN model, Test KNN Model and Plot the graphs for dif
from sklearn.model_selection import GridSearchCV
import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score
from sklearn.model_selection import RandomizedSearchCV
from scipy.stats import randint as sp_randint

# Function for batchwise prediction

def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estim
    # not the predicted outputs

    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%10
    # 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

# Train KNN model using training data, and calculate AUC scores for different valu

def train_KNN(X_tr,y_train,X_cr,y_cv):
    train_auc = []
    cv_auc = []
    K= [7,11,17,21,27,31,37,41,47,51]
    for i in tqdm(K):
        neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
        neigh.fit(X_tr, y_train)

        y_train_pred = batch_predict(neigh, X_tr)
        y_cv_pred = batch_predict(neigh, X_cr)

        # roc_auc_score(y_true, y_score) the 2nd parameter should be probability e
        # not the predicted outputs
        train_auc.append(roc_auc_score(y_train,y_train_pred))
        cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

    plt.plot(K, train_auc, label='Train AUC')
    plt.plot(K, cv_auc, label='CV AUC')

    plt.scatter(K, train_auc, label='Train AUC points')
    plt.scatter(K, cv_auc, label='CV AUC points')
```

```
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

Test the model with optimal K found out using training data. Plot FPR vs TPR(ROC

```
def Testing_KNN(X_tr,X_te,best_k):
    from sklearn.metrics import roc_curve, auc

    #here we are choosing the best_k based on forloop results
    neigh = KNeighborsClassifier(n_neighbors=best_k, n_jobs=-1)
    neigh.fit(X_tr, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estim
    # not the predicted outputs

    y_train_pred = batch_predict(neigh, X_tr)
    y_test_pred = batch_predict(neigh, X_te)

    train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
    test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

    plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tp
    plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
    plt.legend()
    plt.xlabel("tpr")
    plt.ylabel("fpr")
    plt.title("ERROR PLOTS")
    plt.grid()
    plt.show()
    return train_fpr,train_tpr,tr_thresholds,y_train_pred,y_test_pred

# we are writing our own function for predict, with defined thresould
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", n
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

▼ 2.3.1 Applying KNN brute force on BOW encoding essay, and project

▼ 2.3.1.1 Encoding preprocessed_essays BoW

```
print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=10000)
vectorizer.fit(X_train['preprocessed_essays'].values) # fit has to happen only on

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vectorizer.transform(X_train['preprocessed_essays'].values)
X_cv_essay_bow = vectorizer.transform(X_cv['preprocessed_essays'].values)
X_test_essay_bow = vectorizer.transform(X_test['preprocessed_essays'].values)

print("After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
print(X_cv_essay_bow.shape, y_cv.shape)
print(X_test_essay_bow.shape, y_test.shape)
print("="*100)

# print("YOU SHOULD NOT DO SOMETHING LIKE THIS")
# vectorizer = CountVectorizer()
# x_train_bow = vectorizer.fit_transform(X_train['essay'].values)
# x_cv_bow = vectorizer.fit_transform(X_cv['essay'].values)
# x_test_bow = vectorizer.fit_transform(X_test['essay'].values)

# print(x_train_bow.shape, y_train.shape)
# print(x_cv_bow.shape, y_cv.shape)
# print(x_test_bow.shape, y_test.shape)

print("NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME")
```

```
↳ (37700, 16) (37700,)
   (11055, 16) (11055,)
   (16500, 16) (16500,)
=====
After vectorizations
(37700, 10000) (37700,)
(11055, 10000) (11055,)
(16500, 10000) (16500,)
=====
NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME
```

▼ 2.3.1.2 Encoding preprocessed_titles BoW

```

vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=10000)
vectorizer.fit(X_train['preprocessed_titles'].values) # fit has to happen only on

# we use the fitted CountVectorizer to convert the text to vector
X_train_titles_bow = vectorizer.transform(X_train['preprocessed_titles'].values)
X_cv_titles_bow = vectorizer.transform(X_cv['preprocessed_titles'].values)
X_test_titles_bow = vectorizer.transform(X_test['preprocessed_titles'].values)

print("After vectorizations")
print(X_train_titles_bow.shape, y_train.shape)
print(X_cv_titles_bow.shape, y_cv.shape)
print(X_test_titles_bow.shape, y_test.shape)
print("="*100)

# print("YOU SHOULD NOT DO SOMETHING LIKE THIS")
# vectorizer = CountVectorizer()
# x_train_bow = vectorizer.fit_transform(X_train['essay'].values)
# x_cv_bow = vectorizer.fit_transform(X_cv['essay'].values)
# x_test_bow = vectorizer.fit_transform(X_test['essay'].values)

# print(x_train_bow.shape, y_train.shape)
# print(x_cv_bow.shape, y_cv.shape)
# print(x_test_bow.shape, y_test.shape)

print("NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME")

↳ After vectorizations
(37700, 4616) (37700,)
(11055, 4616) (11055,)
(16500, 4616) (16500,)
=====
NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME

```

▼ 2.3.1.3 Merge all the features and obtain final data matrix

```

# Merge all the features:
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((X_train_titles_bow,X_train_essay_bow, X_train_state_ohe, X_train_te
X_cr = hstack((X_cv_titles_bow,X_cv_essay_bow, X_cv_state_ohe, X_cv_teacher_ohe, X
X_te = hstack((X_test_titles_bow,X_test_essay_bow, X_test_state_ohe, X_test_teach

print("Final Data matrix")
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
print("="*100)

```

↳

```
Final Data matrix
(37700, 14680) (37700,)
(11055, 14680) (11055,)
(16500, 14680) (16500,)
```

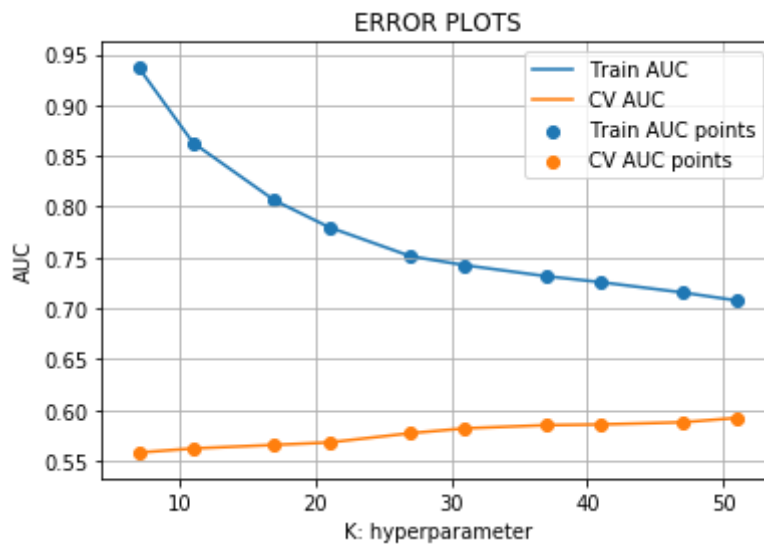
```
=====
```

▼ 2.3.1.4 Training the data model and find best hyperparameter using ROC-AUC

```
# Call train_KNN function on above data
```

```
train_KNN(X_tr,y_train,X_cr,y_cv)
```

```
100%|██████████| 10/10 [29:53<00:00, 176.74s/it]
```

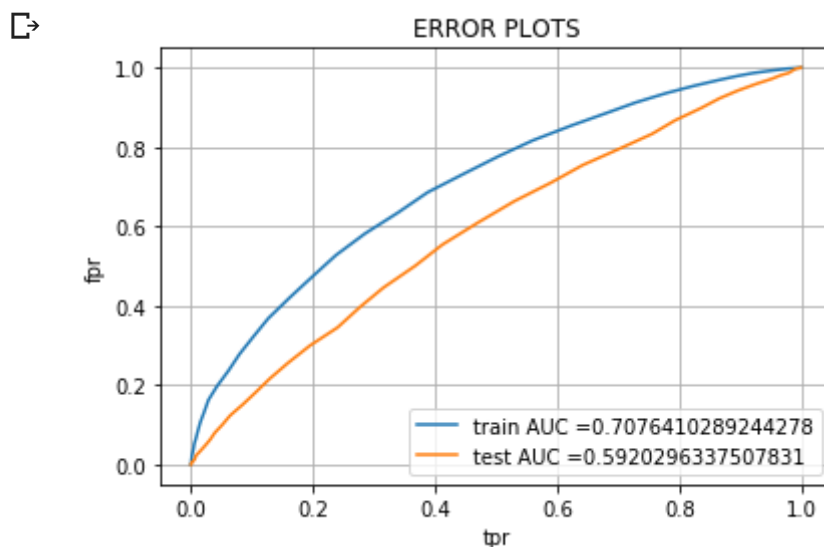


▼ 2.3.1.5 Testing the performance of the model on test data, plotting ROC Curve

```
# Call test_KNN for K obtained by training the data
```

```
best_k=51
```

```
train_fpr,train_tpr,tr_thresholds,y_train_pred,y_test_pred=Testing_KNN(X_tr,X_te,b
```



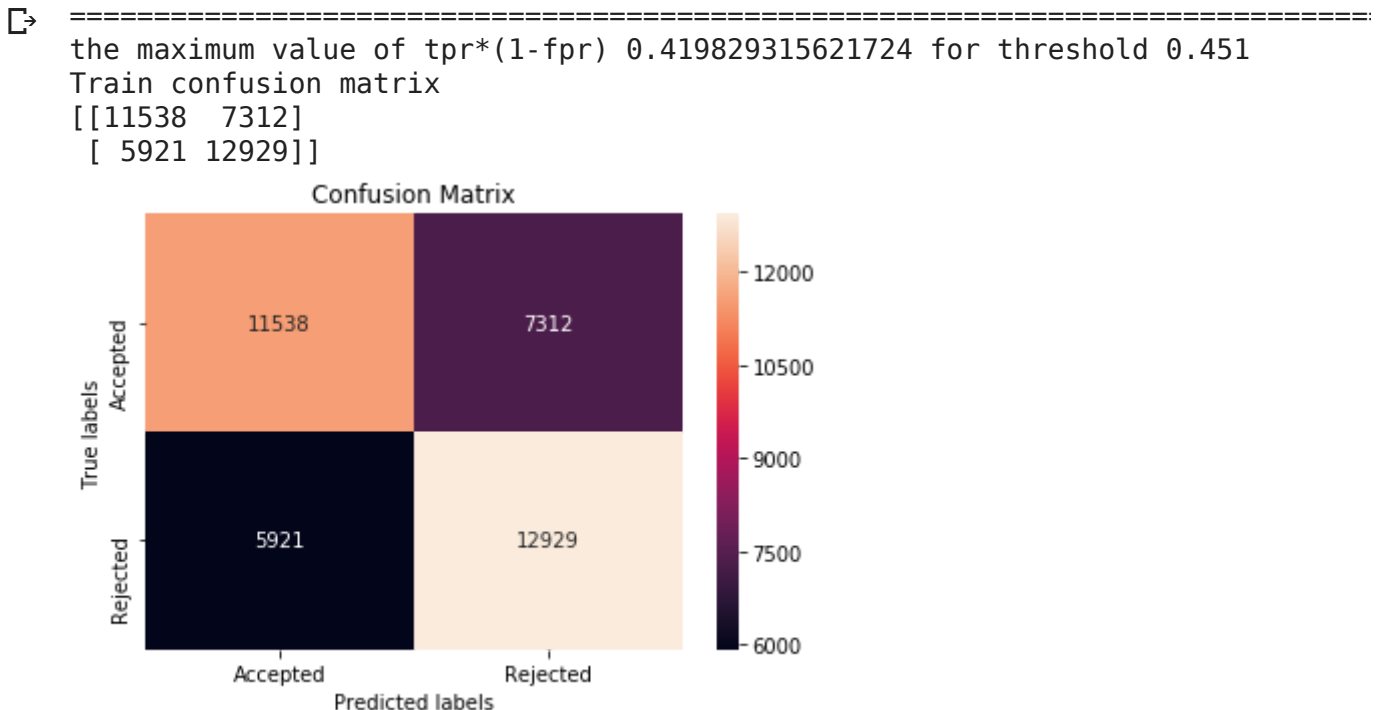
```

print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")

ax= plt.subplot()
cm=confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
print(cm)
sns.heatmap(cm, annot=True, ax = ax,fmt='d'); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
ax.xaxis.set_ticklabels(['Accepted', 'Rejected']); ax.yaxis.set_ticklabels(['Accep

```



```

print("Test confusion matrix")

cm_test = confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
print(cm_test)
ax= plt.subplot()
sns.heatmap(cm_test, annot=True, ax = ax,fmt='d'); #annot=True to annotate cells

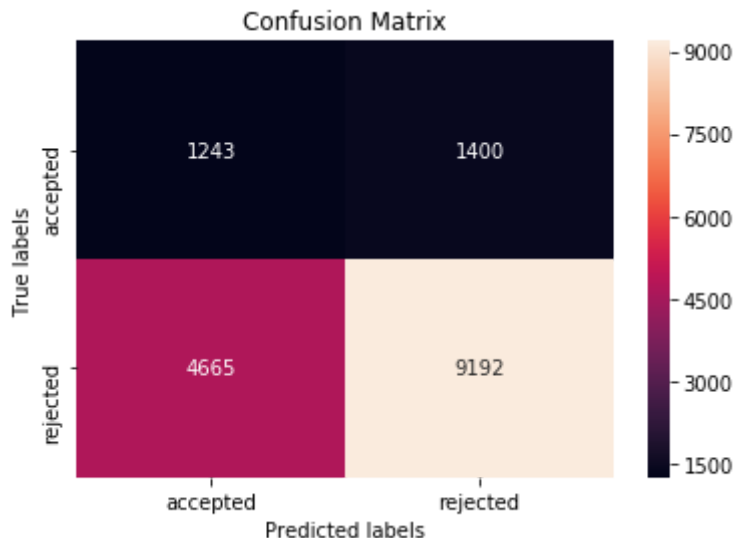
# labels, title and ticks
ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
ax.xaxis.set_ticklabels(['accepted', 'rejected']); ax.yaxis.set_ticklabels(['accep

```

```

↳ Test confusion matrix
[[1243 1400]
 [4665 9192]]

```



▼ 2.3.2 Applying KNN brute force on TFIDF encoding eassay, and proje

▼ 2.3.2.1 Encoding preprocessed_titles TFIDF

```

from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)

#vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=10000)
vectorizer.fit(X_train['preprocessed_titles'].values) # fit has to happen only on

# we use the fitted CountVectorizer to convert the text to vector
X_train_titles_tfidf = vectorizer.transform(X_train['preprocessed_titles'].values)
X_cv_titles_tfidf = vectorizer.transform(X_cv['preprocessed_titles'].values)
X_test_titles_tfidf = vectorizer.transform(X_test['preprocessed_titles'].values)

print("After vectorizations")
print(X_train_titles_tfidf.shape, y_train.shape)
print(X_cv_titles_tfidf.shape, y_cv.shape)
print(X_test_titles_tfidf.shape, y_test.shape)
print("="*100)

```

```

↳ After vectorizations
(37700, 1818) (37700,)
(11055, 1818) (11055,)
(16500, 1818) (16500,)
=====

```

▼ 2.3.2.2 Encoding preprocessed_essays TFIDF

```
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_essays'].values) # fit has to happen only on

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tfidf = vectorizer.transform(X_train['preprocessed_essays'].values)
X_cv_essay_tfidf = vectorizer.transform(X_cv['preprocessed_essays'].values)
X_test_essay_tfidf = vectorizer.transform(X_test['preprocessed_essays'].values)

print("After vectorizations")
print(X_train_essay_tfidf.shape, y_train.shape)
print(X_cv_essay_tfidf.shape, y_cv.shape)
print(X_test_essay_tfidf.shape, y_test.shape)
print("="*100)
```

```
☐➤ After vectorizations
(37700, 11008) (37700,)
(11055, 11008) (11055,)
(16500, 11008) (16500,)
=====
```

▼ 2.3.2.3 Merge all the features and obtain final data matrix

```
from scipy.sparse import hstack
X_tr = hstack((X_train_titles_tfidf, X_train_essay_tfidf, X_train_state_ohe, X_train_teacher_oh))
X_cr = hstack((X_cv_titles_tfidf, X_cv_essay_tfidf, X_cv_state_ohe, X_cv_teacher_oh))
X_te = hstack((X_test_titles_tfidf, X_test_essay_tfidf, X_test_state_ohe, X_test_teacher_oh))

print("Final Data matrix")
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
print("="*100)
```

```
☐➤ Final Data matrix
(37700, 12890) (37700,)
(11055, 12890) (11055,)
(16500, 12890) (16500,)
=====
```

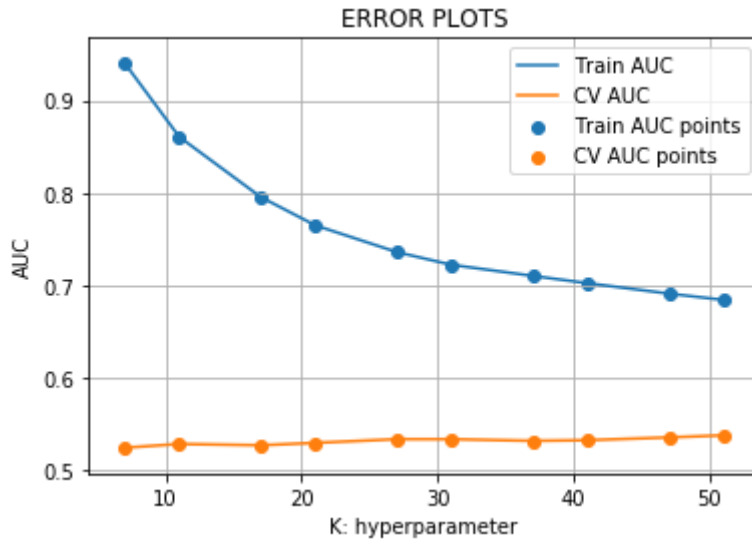
▼ 2.3.2.4 Training the data model and find best hyperparameter using ROC-AUC

```
# Call train_KNN function on above data

train_KNN(X_tr, y_train, X_cr, y_cv)
```

```
☐➤
```


100% | 10/10 [27:00<00:00, 161.70s/it]

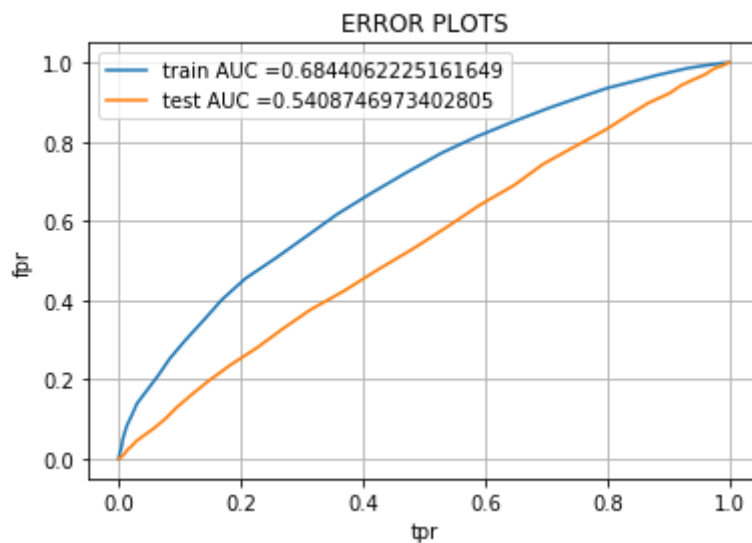


▼ 2.3.2.5 Testing the performance of the model on test data, plotting ROC Curve

Call test_KNN for K obtained by training the data

best_k=51

train_fpr,train_tpr,tr_thresholds,y_train_pred,y_test_pred=Testing_KNN(X_tr,X_te,b



```
print("="*100)
```

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

```
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
```

```
print("Train confusion matrix")
```

```
ax= plt.subplot()
```

```
cm=confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
```

```
print(cm)
```

```
sns.heatmap(cm, annot=True, ax = ax,fmt='d'); #annot=True to annotate cells
```

```
# labels, title and ticks
```

```
ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
```

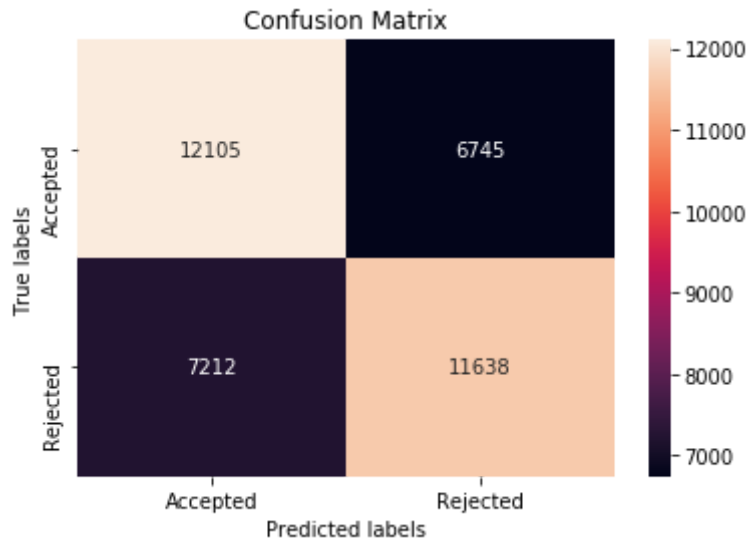
```
ax.set_title('Confusion Matrix');
```

```
ax.set_title('CONFUSION MATRIX');
ax.xaxis.set_ticklabels(['Accepted', 'Rejected']); ax.yaxis.set_ticklabels(['Accep
```

```

↳ =====
the maximum value of tpr*(1-fpr) 0.39647922661807233 for threshold 0.51
Train confusion matrix
[[12105  6745]
 [ 7212 11638]]

```



```
print("Test confusion matrix")
```

```

cm_test = confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
print(cm_test)
ax= plt.subplot()
sns.heatmap(cm_test, annot=True, ax = ax,fmt='d'); #annot=True to annotate cells

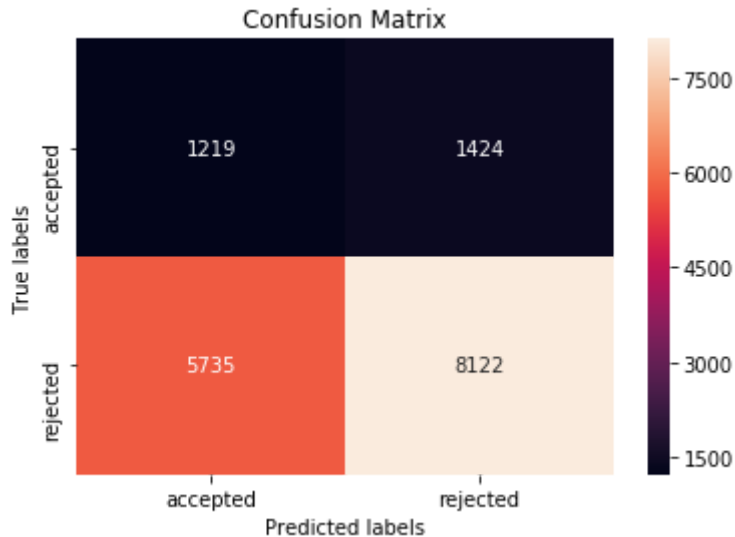
# labels, title and ticks
ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
ax.xaxis.set_ticklabels(['accepted', 'rejected']); ax.yaxis.set_ticklabels(['accep

```

```
↳
```

Test confusion matrix

```
[[1219 1424]
 [5735 8122]]
```



▼ 2.3.3 Applying KNN brute force on AVG W2V

▼ 2.3.3.1 Encoding preprocessed_essays AVG W2V

```
with open('/content/drive/My Drive/Colab Notebooks/Dataset/Assignments_DonorsChoos
model = pickle.load(f)
glove_words = set(model.keys())
```

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors_essays_train = []; # the avg-w2v for each sentence/review is store
for sentence in tqdm(X_train['preprocessed_essays'].values): # for each review/sen
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 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:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_essays_train.append(vector)

print(len(avg_w2v_vectors_essays_train))
print(len(avg_w2v_vectors_essays_train[0]))
```

100% |██████████| 37700/37700 [00:09<00:00, 4009.22it/s]37700
300

```

avg_w2v_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this
for sentence in tqdm(X_cv['preprocessed_essays'].values): # for each review/senten
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 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:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_cv.append(vector)
print(len(avg_w2v_vectors_cv))

```

```

[> 100%|██████████| 11055/11055 [00:02<00:00, 4069.15it/s]11055

```

```

avg_w2v_vectors_essays_test = []; # the avg-w2v for each sentence/review is stored
for sentence in tqdm(X_test['preprocessed_essays'].values): # for each review/sent
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 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:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_essays_test.append(vector)
print(len(avg_w2v_vectors_essays_test))

```

```

[> 100%|██████████| 16500/16500 [00:04<00:00, 4023.13it/s]16500

```

▼ 2.3.3.2 Encoding preprocessed_titles AVG W2V

```

avg_w2v_vectors_titles_train = []; # the avg-w2v for each sentence/review is store
for sentence in tqdm(X_train['preprocessed_titles'].values): # for each review/sen
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 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:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_titles_train.append(vector)
print(len(avg_w2v_vectors_titles_train))

```

```

[> 100%|██████████| 37700/37700 [00:00<00:00, 74438.51it/s]37700

```

```

avg_w2v_vectors_titles_cv = []; # the avg-w2v for each sentence/review is stored i
for sentence in tqdm(X_cv['preprocessed_titles'].values): # for each review/senten
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 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:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_titles_cv.append(vector)

print(len(avg_w2v_vectors_titles_cv))

```

```

[> 100%|██████████| 11055/11055 [00:00<00:00, 56772.64it/s]11055

```

```

avg_w2v_vectors_titles_test = []; # the avg-w2v for each sentence/review is stored
for sentence in tqdm(X_test['preprocessed_titles'].values): # for each review/sent
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 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:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_titles_test.append(vector)

```

```

[> 100%|██████████| 16500/16500 [00:00<00:00, 71154.22it/s]

```

▼ 2.3.3.3 Merge all the features and obtain final data matrix

```

from scipy.sparse import hstack
X_tr = hstack((avg_w2v_vectors_titles_train, avg_w2v_vectors_essays_train, X_train_
X_cr = hstack((avg_w2v_vectors_titles_cv, avg_w2v_vectors_cv, X_cv_state_ohe, X_cv_
X_te = hstack((avg_w2v_vectors_titles_test, avg_w2v_vectors_essays_test, X_test_sta

print("Final Data matrix")
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
print("=="*100)

```

```

[> Final Data matrix
(37700, 664) (37700,)
(11055, 664) (11055,)
(16500, 664) (16500,)
=====

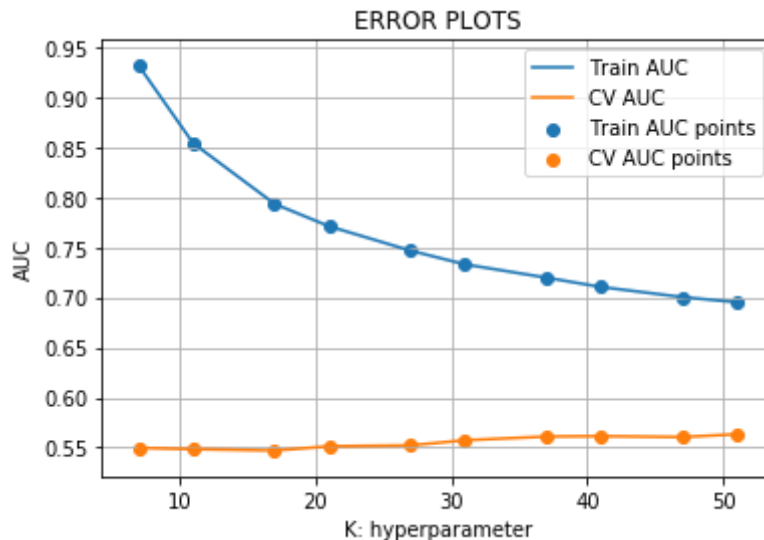
```

▼ 2.3.3.4 Training the data model and find best hyperparameter using ROC-AUC

```
# Call train_KNN function on above data
```

```
train_KNN(X_tr,y_train,X_cr,y_cv)
```

```
100%|██████████| 10/10 [6:08:29<00:00, 2224.91s/it]
```

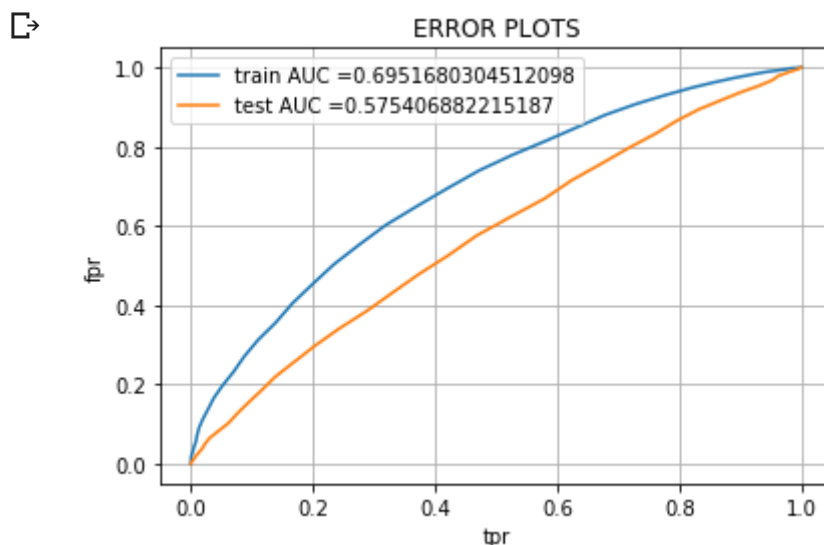


▼ 2.3.3.5 Testing the performance of the model on test data, plotting ROC Curve

```
# Call test_KNN for K obtained by training the data
```

```
best_k=51
```

```
train_fpr,train_tpr,tr_thresholds,y_train_pred,y_test_pred=Testing_KNN(X_tr,X_te,b
```



```
print("="*100)
```

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

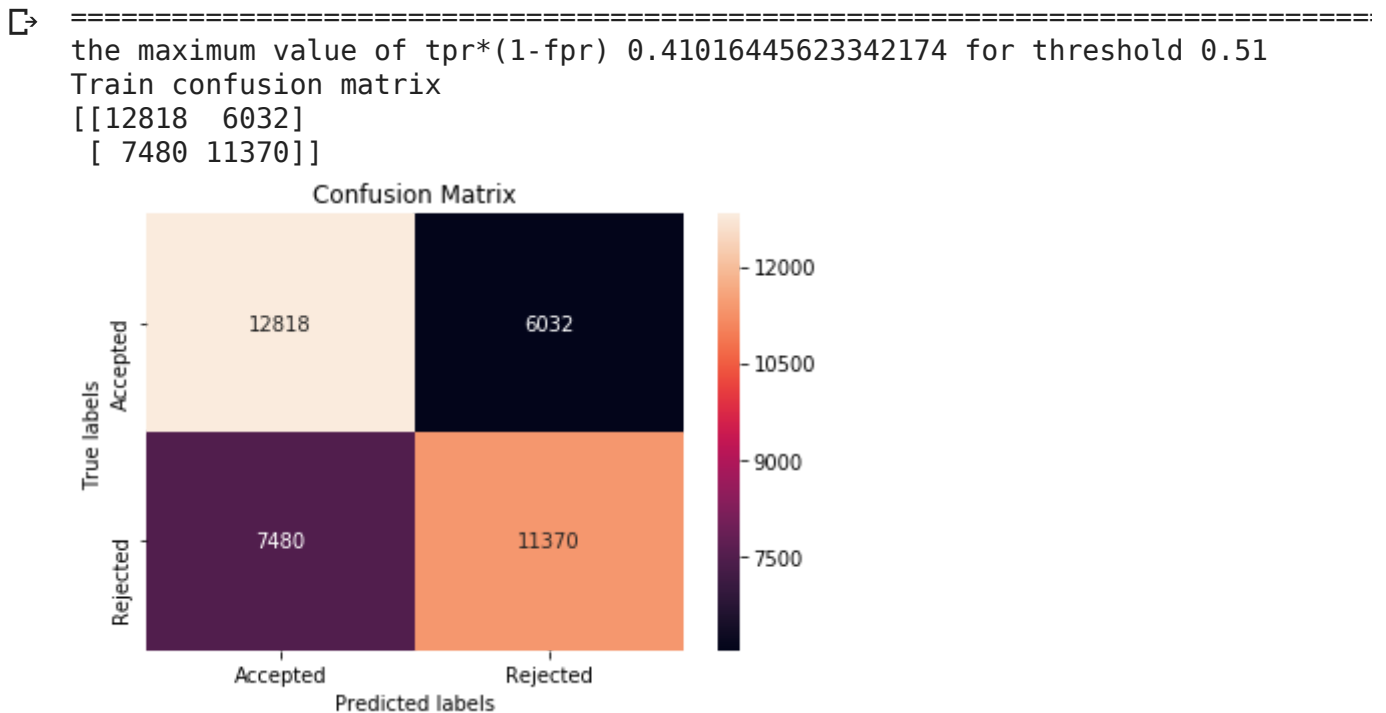
```
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
```

```
print("Train confusion matrix")
```

```
ax = plt.subplot(1,1)
```

```
ax= plt.subplot()
cm=confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
print(cm)
sns.heatmap(cm, annot=True, ax = ax,fmt='d'); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
ax.xaxis.set_ticklabels(['Accepted', 'Rejected']); ax.yaxis.set_ticklabels(['Accep
```



```
print("Test confusion matrix")
```

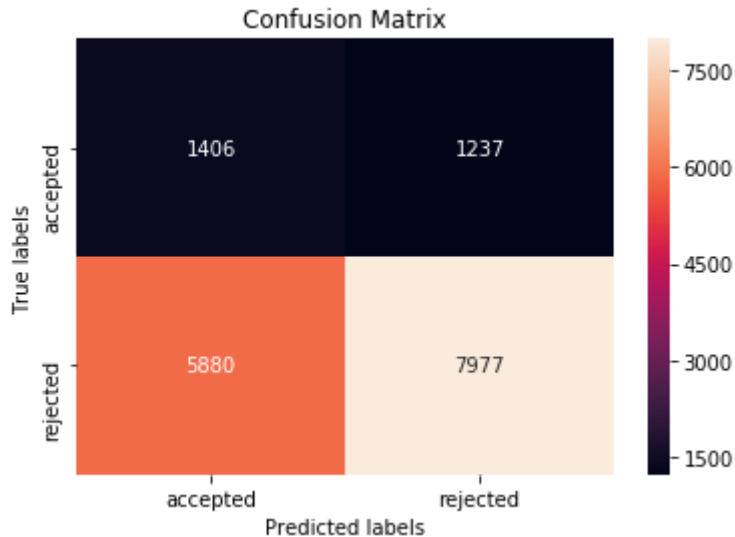
```
cm_test = confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
print(cm_test)
ax= plt.subplot()
sns.heatmap(cm_test, annot=True, ax = ax,fmt='d'); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
ax.xaxis.set_ticklabels(['accepted', 'rejected']); ax.yaxis.set_ticklabels(['accep
```



Test confusion matrix

```
[[1406 1237]
 [5880 7977]]
```



▼ 2.3.4 Applying KNN brute force on TFIDF W2V

▼ 2.3.4.1 Encoding preprocessed_titles tfidf W2V

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['preprocessed_titles'])
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_titles_train = []; # the avg-w2v for each sentence/review is stored
for sentence in tqdm(X_train['preprocessed_titles']): # 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
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_titles_train.append(vector)

print(len(tfidf_w2v_vectors_titles_train))
print(len(tfidf_w2v_vectors_titles_train[0]))
```



```

↳ 100%|██████████| 37700/37700 [00:01<00:00, 26805.96it/s]37700
300

```

```

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_titles_cv = []; # the avg-w2v for each sentence/review is stored
for sentence in tqdm(X_cv['preprocessed_titles']): # 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
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_titles_cv.append(vector)

print(len(tfidf_w2v_vectors_titles_cv))
print(len(tfidf_w2v_vectors_titles_cv[0]))

```

```

↳ 100%|██████████| 11055/11055 [00:00<00:00, 27105.85it/s]11055
300

```

```

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_titles_test = []; # the avg-w2v for each sentence/review is stor
for sentence in tqdm(X_test['preprocessed_titles']): # 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
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_titles_test.append(vector)

print(len(tfidf_w2v_vectors_titles_test))
print(len(tfidf_w2v_vectors_titles_test[0]))

```

```

↳ 100%|██████████| 16500/16500 [00:00<00:00, 25467.36it/s]16500
300

```

▼ 2.3.4.2 Encoding preprocessed_essays tfidf W2V

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['preprocessed_essays'])
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_essays_train = []; # the avg-w2v for each sentence/review is stored
for sentence in tqdm(X_train['preprocessed_essays']): # 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
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_essays_train.append(vector)

print(len(tfidf_w2v_vectors_essays_train))
print(len(tfidf_w2v_vectors_essays_train[0]))
```

```
100%|██████████| 37700/37700 [01:19<00:00, 475.91it/s]37700
300
```

```
tfidf_w2v_vectors_essays_cv = []; # the avg-w2v for each sentence/review is stored
for sentence in tqdm(X_cv['preprocessed_essays']): # 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
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_essays_cv.append(vector)

print(len(tfidf_w2v_vectors_essays_cv))
print(len(tfidf_w2v_vectors_essays_cv[0]))
```

```
100%|██████████| 11055/11055 [00:23<00:00, 461.55it/s]11055
300
```

```
tfidf_w2v_vectors_essays_test = []; # the avg-w2v for each sentence/review is stor
for sentence in tqdm(X_test['preprocessed_essays']): # 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
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_essays_test.append(vector)

print(len(tfidf_w2v_vectors_essays_test))
print(len(tfidf_w2v_vectors_essays_test[0]))
```

```
100%|██████████| 16500/16500 [00:35<00:00, 467.80it/s]16500
300
```

▼ 2.3.4.3 Merge all the features and obtain final data matrix

```
from scipy.sparse import hstack
X_tr = hstack((tfidf_w2v_vectors_titles_train,tfidf_w2v_vectors_essays_train, X_tr
X_cr = hstack((tfidf_w2v_vectors_titles_cv,tfidf_w2v_vectors_essays_cv, X_cv_state
X_te = hstack((tfidf_w2v_vectors_titles_test,tfidf_w2v_vectors_essays_test, X_test

print("Final Data matrix")
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(37700, 664) (37700,)
(11055, 664) (11055,)
(16500, 664) (16500,)
=====
```

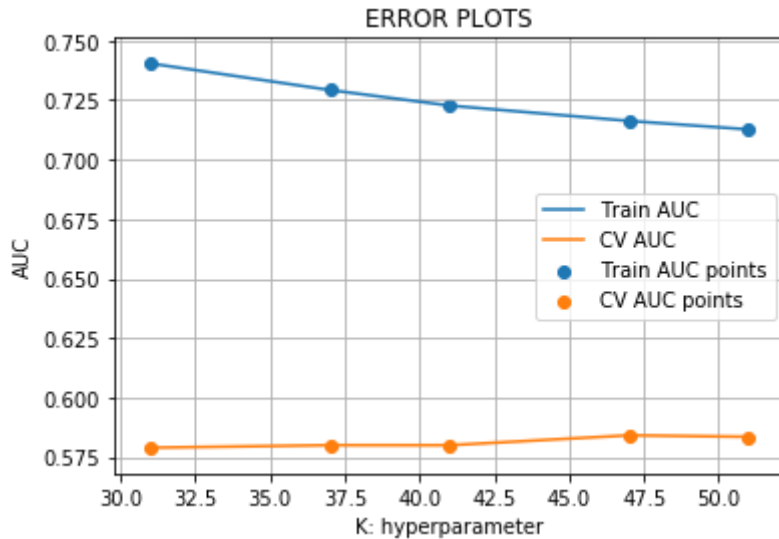
▼ 2.3.4.4 Training the data model and find best hyperparameter using ROC-AUC

```
# Call train_KNN function on above data
```

```
train_KNN(X_tr,y_train,X_cr,y_cv)
```

```
↳
```

100% | 5/5 [3:27:12<00:00, 2487.32s/it]

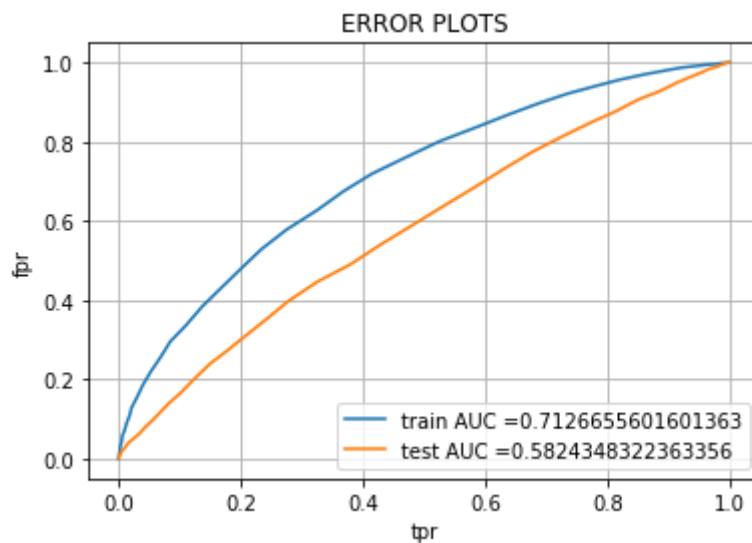


2.3.4.5 Testing the performance of the model on test data, plotting ROC Curve

Call test_KNN for K obtained by training the data

best_k=51

train_fpr,train_tpr,tr_thresholds,y_train_pred,y_test_pred=Testing_KNN(X_tr,X_te,b



```
print("="*100)
```

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

```
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
```

```
print("Train confusion matrix")
```

```
ax= plt.subplot()
```

```
cm=confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
```

```
print(cm)
```

```
sns.heatmap(cm, annot=True, ax = ax,fmt='d'); #annot=True to annotate cells
```

```
# labels, title and ticks
```

```
ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
```

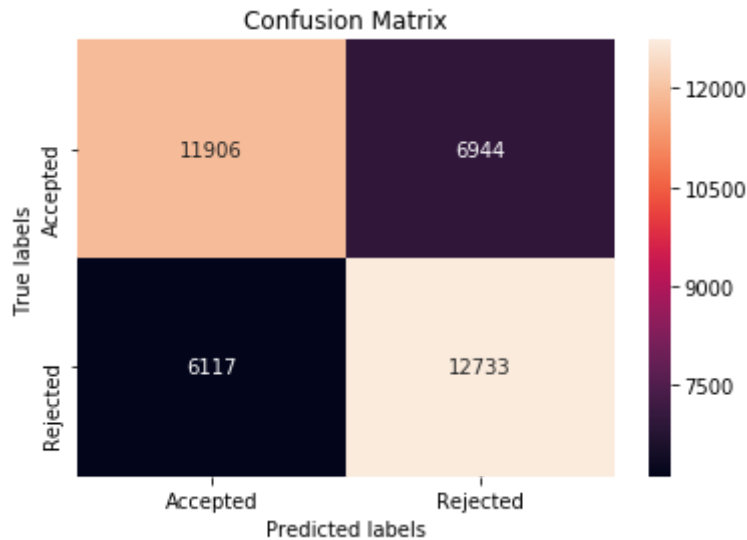
```
ax.set_title('Confusion Matrix');
```

```
ax.set_title('CONFUSION MATRIX');
ax.xaxis.set_ticklabels(['Accepted', 'Rejected']); ax.yaxis.set_ticklabels(['Accep
```

```

↳ =====
the maximum value of tpr*(1-fpr) 0.4266521202569497 for threshold 0.471
Train confusion matrix
[[11906  6944]
 [ 6117 12733]]

```



```
print("Test confusion matrix")
```

```

cm_test = confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
print(cm_test)
ax= plt.subplot()
sns.heatmap(cm_test, annot=True, ax = ax,fmt='d'); #annot=True to annotate cells

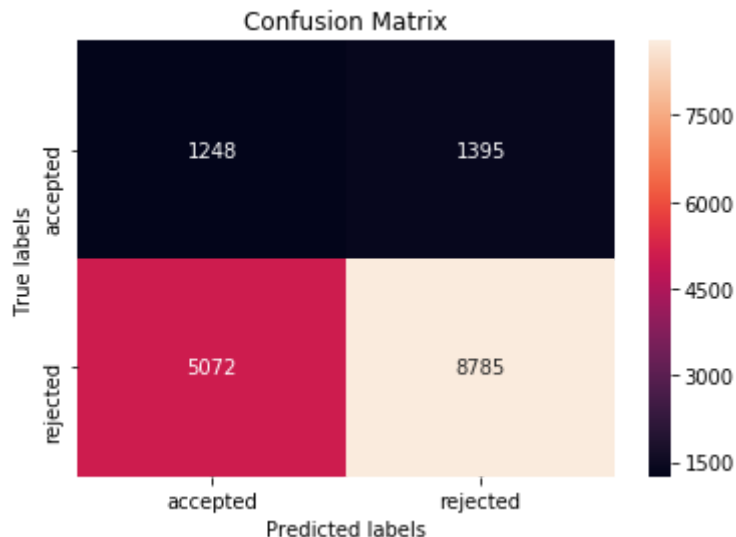
# labels, title and ticks
ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
ax.xaxis.set_ticklabels(['accepted', 'rejected']); ax.yaxis.set_ticklabels(['accep

```

```
↳
```

Test confusion matrix

```
[[1248 1395]
 [5072 8785]]
```



▼ 2.3.5 Applying KNN brute force on TFIDF with feature selection

▼ 2.3.5.1 Merge all the features and obtain final data matrix and perform feature

```
from sklearn.feature_selection import SelectKBest, chi2
from scipy.sparse import hstack
X_tr = hstack((X_train_titles_tfidf, X_train_essay_tfidf, X_train_state_ohe, X_train_teacher_ohe))
X_cr = hstack((X_cv_titles_tfidf, X_cv_essay_tfidf, X_cv_state_ohe, X_cv_teacher_ohe))
X_te = hstack((X_test_titles_tfidf, X_test_essay_tfidf, X_test_state_ohe, X_test_teacher_ohe))

print("Final Data matrix")
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
print("="*100)

# Select best K features
# https://www.w3cschool.cn/doc\_scikit\_learn/scikit\_learn-modules-generated-sklearn
# https://www.programcreek.com/python/example/93974/sklearn.feature\_selection.SelectKBest

X_new = SelectKBest(chi2, k=2000).fit(X_tr, y_train)
X_tr = X_new.transform(X_tr)
X_cr = X_new.transform(X_cr)
X_te = X_new.transform(X_te)

print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
```



Final Data matrix

```
(37700, 12890) (37700,)
(11055, 12890) (11055,)
(16500, 12890) (16500,)
```

```
=====
(37700, 2000) (37700,)
(11055, 2000) (11055,)
(16500, 2000) (16500,)
```

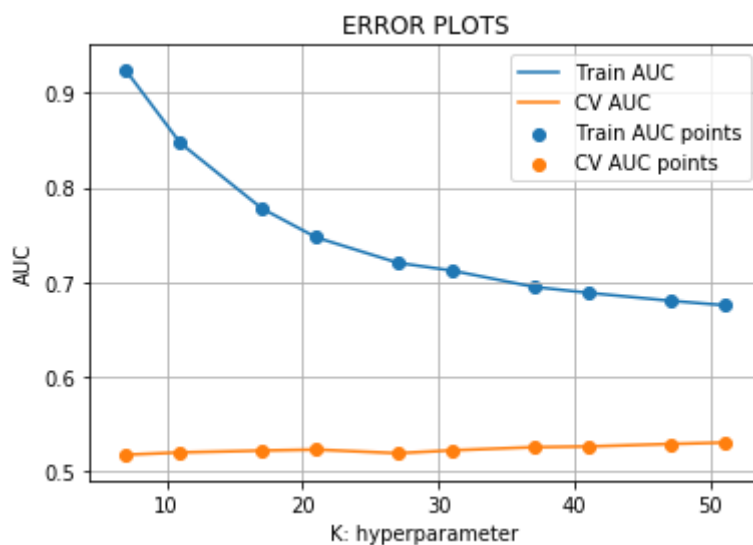
▼ 2.3.5.2 Training the data model and find best hyperparameter using ROC-AUC

Call train_KNN function on above data

```
train_KNN(X_tr,y_train,X_cr,y_cv)
```



0%		0/10 [00:00<?, ?it/s]
10%		1/10 [02:04<18:43, 124.82s/it]
20%		2/10 [04:09<16:37, 124.64s/it]
30%		3/10 [06:14<14:33, 124.85s/it]
40%		4/10 [08:19<12:29, 124.86s/it]
50%		5/10 [10:24<10:25, 125.04s/it]
60%		6/10 [12:30<08:20, 125.14s/it]
70%		7/10 [14:34<06:15, 125.06s/it]
80%		8/10 [16:40<04:10, 125.25s/it]
90%		9/10 [18:45<02:05, 125.24s/it]
100%		10/10 [20:50<00:00, 125.19s/it]



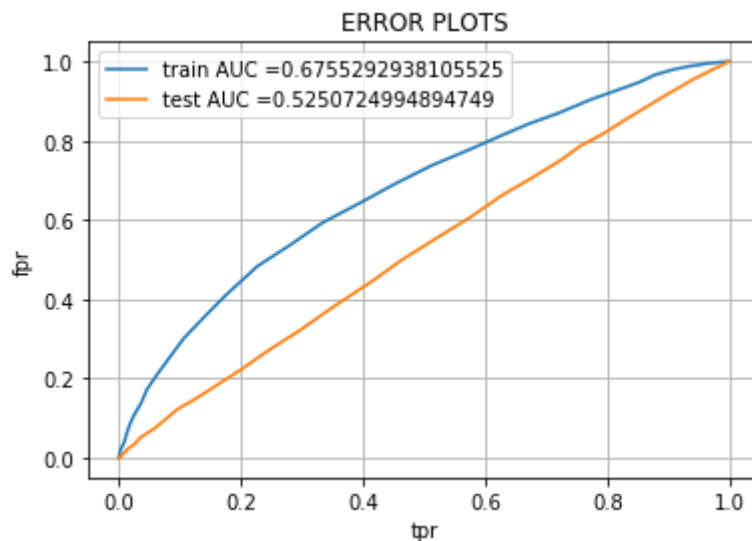
▼ 2.3.5.3 Testing the performance of the model on test data, plotting ROC Curve

Call test_KNN for K obtained by training the data

```
best_k=51
```

```
train_fpr,train_tpr,tr_thresholds,y_train_pred,y_test_pred=Testing_KNN(X_tr,X_te,b
```





```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")

ax= plt.subplot()
cm=confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
print(cm)
sns.heatmap(cm, annot=True, ax = ax,fmt='d'); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
ax.xaxis.set_ticklabels(['Accepted', 'Rejected']); ax.yaxis.set_ticklabels(['Accep
```

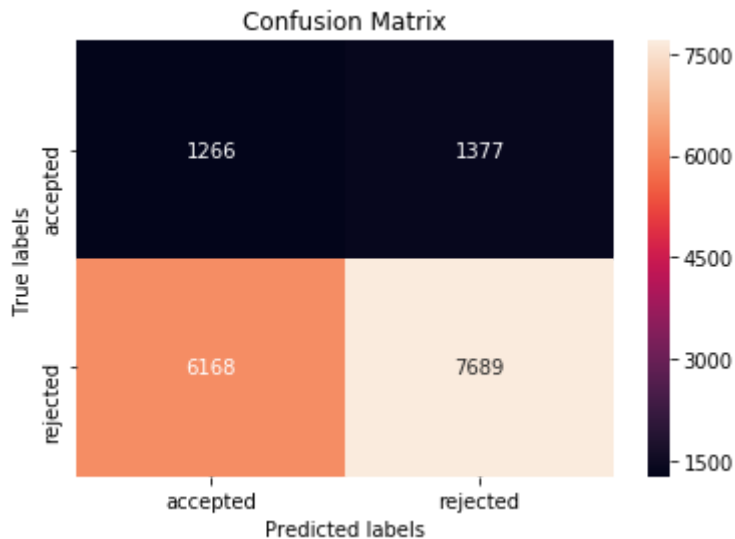



```
print("Test confusion matrix")
```

```
cm_test = confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
print(cm_test)
ax= plt.subplot()
sns.heatmap(cm_test, annot=True, ax = ax,fmt='d'); #annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
ax.xaxis.set_ticklabels(['accepted', 'rejected']); ax.yaxis.set_ticklabels(['accep
```

```
☞ Test confusion matrix
[[1266 1377]
 [6168 7689]]
```



▼ 3.1 Summary

```
# To summarize the results:
# summary table in jupyter notebook
# http://zetcode.com/python/prettytable/
# https://stackoverflow.com/questions/35160256/how-do-i-output-lists-as-a-table-in
```

```
from prettytable import PrettyTable
```

```
x = PrettyTable()
```

```
x.field_names = ["Vectorizer", "Model", "Hyperparameter", "Train_AUC", "Test_AUC"]
```

```
x.add_row(["Bag of Words", "Brute Force", 51, 0.70, 0.59])
```

```
x.add_row(["TF-IDF", "KNN", 51, 0.68, 0.54])
```

```

x.add_row(["Avg W2V", "KNN",51 , 0.69,0.57])
x.add_row(["TF-IDF W2V", "KNN",51 , 0.71,0.58])
x.add_row(["TF-IDF with Feature Selection", "KNN",51 , 0.67,0.52])

print(x)

```

↗

Vectorizer	Model	Hyperparameter	Train_AUC
Bag of Words	Brute Force	51	0.7
TF-IDF	KNN	51	0.68
Avg W2V	KNN	51	0.69
TF-IDF W2V	KNN	51	0.71
TF-IDF with Feature Selection	KNN	51	0.67