

## ▼ DonorsChoose

DonorsChoose.org receives hundreds of thousands of project proposals each year for classroom projects. A large 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, there are several challenges to solve:

- How to scale current manual processes and resources to screen 500,000 projects so that they can be screened as fast as possible
- How to increase the consistency of project vetting across different volunteers to improve the experience
- 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 will be approved. The text of project descriptions as well as additional metadata about the project, teacher, and school is available 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. <b>Example:</b> p036502
<code>project_title</code>	Title of the project. <b>Examples:</b> <ul style="list-style-type: none"><li>• Art Will Make You Happy!</li><li>• First Grade Fun</li></ul>
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following categories: <ul style="list-style-type: none"><li>• Grades PreK-2</li><li>• Grades 3-5</li><li>• Grades 6-8</li><li>• Grades 9-12</li></ul>

Feature	Description
<code>project_subject_categories</code>	<p>One or more (comma-separated) subject categories for the project for</p> <ul style="list-style-type: none"> <li>• Applied Learning</li> <li>• Care &amp; Hunger</li> <li>• Health &amp; Sports</li> <li>• History &amp; Civics</li> <li>• Literacy &amp; Language</li> <li>• Math &amp; Science</li> <li>• Music &amp; The Arts</li> <li>• Special Needs</li> <li>• Warmth</li> </ul> <p><b>Examples:</b></p> <ul style="list-style-type: none"> <li>• Music &amp; The Arts</li> <li>• Literacy &amp; Language, Math &amp; Science</li> </ul>
<code>school_state</code>	State where school is located ( <a href="#">Two-letter U.S. postal code</a> ). <b>Example:</b>
<code>project_subject_subcategories</code>	<p>One or more (comma-separated) subject subcategories for the project</p> <ul style="list-style-type: none"> <li>• Literacy</li> <li>• Literature &amp; Writing, Social Sciences</li> </ul>
<code>project_resource_summary</code>	<p>An explanation of the resources needed for the project. <b>Example:</b></p> <ul style="list-style-type: none"> <li>• My students need hands on literacy materi</li> </ul>
<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. <b>Example:</b> 2016-04
<code>teacher_id</code>	A unique identifier for the teacher of the proposed project. <b>Example:</b>
<code>teacher_prefix</code>	<p>Teacher's title. One of the following enumerated values:</p> <ul style="list-style-type: none"> <li>• nan</li> <li>• Dr.</li> <li>• Mr.</li> <li>• Mrs.</li> <li>• Ms.</li> <li>• Teacher.</li> </ul>
<code>teacher_number_of_previously_posted_projects</code>	Number of project applications previously submitted by the same tea

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

Additionally, the `resources.csv` data set provides more data about the resources required for each project resource required by a project:

Feature	Description
<code>id</code>	A <code>project_id</code> value from the <code>train.csv</code> file. <b>Example:</b> p036502

## ▼ Notes on the Essay Data

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

- \_\_project\_essay\_1:\_\_ "Introduce us to your classroom"
- \_\_project\_essay\_2:\_\_ "Tell us more about your students"
- \_\_project\_essay\_3:\_\_ "Describe how your students will use the materials you're requesting"
- \_\_project\_essay\_3:\_\_ "Close by sharing why your project will make a difference"

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

- \_\_project\_essay\_1:\_\_ "Describe your students: What makes your students special? Specific details about your neighborhood, and your school are all helpful."
- \_\_project\_essay\_2:\_\_ "About your project: How will these materials make a difference in your students' lives?"

For all projects with project\_submitted\_datetime of 2016-05-17 and later, the values of project\_essay\_1 and project\_essay\_2 are:

```
%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
```

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

## ▼ 1.1 Reading Data

```
project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')
```

```
print("Number of data points in train data", project_data.shape)
print('-'*50)
print("The attributes of data :", project_data.columns.values)
```



Number of data points in train data (109248, 17)

```
-----
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
'project_submitted_datetime' 'project_grade_category'
'project_subject_categories' 'project_subject_subcategories'
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

```
# how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
cols = ['Date' if x=='project_submitted_datetime' else x for x in list(project_data.columns)]
```

```
#sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
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/13148611/4084039
project_data = project_data[cols]
```

```
project_data.head(2)
```



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

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



Number of data points in train data (1541272, 4)  
['id' 'description' 'quantity' 'price']

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

## ▼ 1.2 preprocessing of project\_subject\_categories

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

# 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", "Ca
        if 'The' in j.split(): # this will split each of the catogory based on space "Math &
            j=j.replace('The','') # if we have the words "The" we are going to replace it wit
            j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
            temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
            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)
```

```

from collections import Counter
my_counter = Counter()
for word in project_data['clean_categories'].values:
    my_counter.update(word.split())

cat_dict = dict(my_counter)
sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))

```

## ▼ 1.3 preprocessing of project\_subject\_subcategories

```

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

# 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

sub_cat_list = []
for i in sub_categories:
    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", "Ca
        if 'The' in j.split(): # this will split each of the category based on space "Math &
            j=j.replace('The','') # if we have the words "The" we are going to replace it wit
            j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math &
            temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
            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)

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in project_data['clean_subcategories'].values:
    my_counter.update(word.split())

sub_cat_dict = dict(my_counter)
sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))

```

## ▼ 1.3 Text preprocessing

```

# merge two column text dataframe:
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)
```

```
project_data.head(2)
```



Unnamed:  
0

id

teacher\_id

teacher\_prefix

school\_state

55660

8393

p205479

2bf07ba08945e5d8b2a3f269b2b3cfe5

Mrs.

CA

76127

37728

p043609

3f60494c61921b3b43ab61bdde2904df

Ms.

UT

#### 1.4.2.3 Using Pretrained Models: TFIDF weighted W2V

```
# printing some random reviews
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
print("="*50)
print(project_data['essay'].values[99999])
print("="*50)
```




I have been fortunate enough to use the Fairy Tale STEM kits in my classroom as well as  
 =====  
 I teach high school English to students with learning and behavioral disabilities. My st  
 =====  
 \"Life moves pretty fast. If you don't stop and look around once in awhile, you could mi  
 =====  
 \"A person's a person, no matter how small.\" (Dr.Seuss) I teach the smallest students w  
 =====  
 My classroom consists of twenty-two amazing sixth graders from different cultures and ba  
 =====

```
# 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


sent = decontracted(project_data['essay'].values[20000])
print(sent)
print("="*50)
```

 \ "A person is a person, no matter how small.\" (Dr.Seuss) I teach the smallest students  
=====

```
# \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
sent = sent.replace('\r', ' ')
sent = sent.replace('\n', ' ')
sent = sent.replace('\t', ' ')
print(sent)
```

 A person is a person, no matter how small. (Dr.Seuss) I teach the smallest students wi

```
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

 A person is a person no matter how small Dr Seuss I teach the smallest students with th

```
# 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', "you're", "yo
    'you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',
    'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they',
    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll"
    'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'h
    'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'unt
    'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'dur
    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', '

```



```
'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'bo
'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'ver
's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd
've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'does
'hadn't', 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "
'mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn',
'won', "won't", 'wouldn', "wouldn't"]
```

```
# Combining all the above stundents
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('\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
    preprocessed_essays.append(sent.lower().strip())
```



100% | 109248/10

```
# after preprocessing
preprocessed_essays[20000]
```



'a person person no matter small dr seuss i teach smallest students biggest enthusiasm l

## 1.4 Preprocessing of `project\_title`

```
# similarly you can preprocess the titles also
```

## ▼ 1.5 Preparing data for models

```
project_data.columns
```



```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
'Date', 'project_grade_category', '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',
'clean_categories', 'clean_subcategories', 'essay'],
dtype='object')
```

we are going to consider

- school\_state : categorical data
- clean\_categories : categorical data
- clean\_subcategories : categorical data
- project\_grade\_category : categorical data
- teacher\_prefix : categorical data
  
- project\_title : text data
- text : text data
- project\_resource\_summary: text data (optinal)
  
- quantity : numerical (optinal)
- teacher\_number\_of\_previously\_posted\_projects : numerical
- price : numerical

### ► 1.5.1 Vectorizing Categorical data

↳ 4 cells hidden

### ► 1.5.2 Vectorizing Text data

↳ 13 cells hidden

### ▼ 1.5.3 Vectorizing Numerical features

```
price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
project_data = pd.merge(project_data, price_data, on='id', how='left')
```

```
# check this one: https://www.youtube.com/watch?v=0H0q0cIn3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing
from sklearn.preprocessing import StandardScaler
```

```
# price_standardized = standardScaler.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399.]
# Reshape your data either using array.reshape(-1, 1)
```

```
price_scalar = StandardScaler()
price_scalar.fit(project_data['price'].values.reshape(-1,1)) # finding the mean and standard
print(f"Mean : {price_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])}")
```

```
# Now standardize the data with above maen and variance.
price_standardized = price_scalar.transform(project_data['price'].values.reshape(-1, 1))
```

price\_standardized

```
array([[4.63560392e-03, 1.36200635e-03, 2.10346002e-03, ...,
        2.55100471e-03, 1.83960046e-03, 3.51642253e-05]])
```

## ▼ 1.5.4 Merging all the above features

- we need to merge all the numerical vectors i.e catogorical, text, numerical vectors

```
print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(text_bow.shape)
print(price_standardized.shape)
```

```
(109248, 9)
(109248, 30)
(109248, 16623)
(109248, 1)
```

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X = hstack((categories_one_hot, sub_categories_one_hot, text_bow, price_standardized))
X.shape
```

```
(109248, 16663)
```

## ▼ Assignment 3: Apply KNN



### 1. [Task-1] Apply KNN(brute force version) on these feature sets

- **Set 1:** categorical, numerical features + project\_title(BOW) + preprocessed\_essay (BOW)
- **Set 2:** categorical, numerical features + project\_title(TFIDF)+ preprocessed\_essay (TFIDF)
- **Set 3:** categorical, numerical features + project\_title(AVG W2V)+ preprocessed\_essay (AVG W2V)
- **Set 4:** categorical, numerical features + project\_title(TFIDF W2V)+ preprocessed\_essay (TFIDF W2V)

### 2. Hyper paramter tuning to find best K

- Find the best hyper parameter which results in the maximum [AUC](#) value
- Find the best hyper paramter using k-fold cross validation (or) simple cross validation data
- Use gridsearch-cv or randomsearch-cv or write your own for loops to do this task

### 3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for figure 
- Once you find the best hyper parameter, you need to train your model-M using the best hyper parameter and plot the ROC curve on both train and test using model-M. 
- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and actual

#### 4. [Task-2]

- Select top 2000 features from feature **Set 2** using `SelectKBest` and then apply KNN on to


```

from sklearn.datasets import load_digits
from sklearn.feature_selection import SelectKBest, chi2
X, y = load_digits(return_X_y=True)
X.shape
X_new = SelectKBest(chi2, k=20).fit_transform(X, y)
X_new.shape
=====
output:
(1797, 64)
(1797, 20)

```

- Repeat the steps 2 and 3 on the data matrix after feature selection

#### 5. Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table for this prettytable library [link](#) 

### Note: Data Leakage

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
2. To avoid the issue of data-leakage, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on your train data, and apply the method `transform()` on your test data.
4. For more details please go through this [link](#).

## 2. K Nearest Neighbor

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

```
%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_notebook as tqdm1
from tqdm import tqdm
import time
import os

from plotly import 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 train_test_split
```



```
C:\Users\LENOVO\Anaconda3\lib\site-packages\smart_open\ssh.py:34: UserWarning: paramiko
warnings.warn('paramiko missing, opening SSH/SCP/SFTP paths will be disabled. `pip in
C:\Users\LENOVO\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected
warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

```
project_data = pd.read_csv('../train_data.csv', nrows=50000)
resource_data = pd.read_csv('../resources.csv')

print("Number of data points in train data", project_data.shape)
print('- '*50)
print("The attributes of data :", project_data.columns.values)
```



Number of data points in train data (50000, 17)

```
-----
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
'project_submitted_datetime' 'project_grade_category'
'project_subject_categories' 'project_subject_subcategories'
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

## ▼ Text preprocessing(1)

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

# 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", "Ca
        if 'The' in j.split(): # this will split each of the catogory based on space "Math &
            j=j.replace('The','') # if we have the words "The" we are going to replace it wit
            j = j.replace(' ', '') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
            temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
            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)
project_data.head(5)
```



	Unnamed: 0		id	teacher_id	teacher_prefix	school_state	pr
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc		Mrs.	IN	
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a		Mr.	FL	
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0		Ms.	AZ	
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60		Mrs.	KY	
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec		Mrs.	TX	

# count of all the words in corpus python: <https://stackoverflow.com/a/22898595/4084039>

```
from collections import Counter
my_counter = Counter()
for word in project_data['clean_categories'].values:
    my_counter.update(word.split())
my_counter
```

 Counter({'Literacy\_Language': 23998, 'History\_Civics': 2689, 'Health\_Sports': 6538, 'Math\_Science': 18874, 'SpecialNeeds': 6233, 'AppliedLearning': 5569, 'Music\_Arts': 4699, 'Warmth': 643, 'Care\_Hunger': 643})

# dict sort by value python: <https://stackoverflow.com/a/613218/4084039>

```
cat_dict = dict(my_counter)
sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
```

```

# ind = np.arange(len(sorted_cat_dict))
# plt.figure(figsize=(20,5))
# p1 = plt.bar(ind, list(sorted_cat_dict.values()))

# plt.ylabel('Projects')
# plt.title('% of projects aproved category wise')
# plt.xticks(ind, list(sorted_cat_dict.keys()))
# plt.show()
# print(sorted_cat_dict)

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

# 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

sub_cat_list = []
for i in sub_categories:
    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", "Ca
        if 'The' in j.split(): # this will split each of the catogory based on space "Math &
            j=j.replace('The','') # if we have the words "The" we are going to replace it wit
            j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
            temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
            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)
project_data.head(2)

```



Unnamed:  
0

id

teacher\_id teacher\_prefix school\_state pr

0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL

```

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()

```



```

for word in project_data['clean_subcategories'].values:
    my_counter.update(word.split())

# dict sort by value python: https://stackoverflow.com/a/613218/4084039
sub_cat_dict = dict(my_counter)
sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))

# ind = np.arange(len(sorted_sub_cat_dict))
# plt.figure(figsize=(20,5))
# p1 = plt.bar(ind, list(sorted_sub_cat_dict.values()))

# plt.ylabel('Projects')
# plt.title('% of projects aproved state wise')
# plt.xticks(ind, list(sorted_sub_cat_dict.keys()))
# plt.show()

# merge two column text dataframe:
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)

# https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indexes-for-all-grou
price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
price_data.head(2)

```

	id	price	quantity
0	p000001	459.56	7
1	p000002	515.89	21

```

# join two dataframes in python:
project_data = pd.merge(project_data, price_data, on='id', how='left')

#presence of the numerical digits in a strings with numeric : https://stackoverflow.com/a/198
def hasNumbers(inputString):
    return any(i.isdigit() for i in inputString)
p1 = project_data[['id','project_resource_summary']]
p1 = pd.DataFrame(data=p1)
p1.columns = ['id','digits_in_summary']
p1['digits_in_summary'] = p1['digits_in_summary'].map(hasNumbers)
# https://stackoverflow.com/a/17383325/8089731
p1['digits_in_summary'] = p1['digits_in_summary'].astype(int)
project_data = pd.merge(project_data, p1, on='id', how='left')
project_data.head(5)

```

	Unnamed: 0		id	teacher_id	teacher_prefix	school_state	pr
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc		Mrs.	IN	
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a		Mr.	FL	
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0		Ms.	AZ	
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60		Mrs.	KY	
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec		Mrs.	TX	

5 rows × 21 columns

## ▼ Text preprocessing(2)

# <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
```





	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	pr
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL	

2 rows × 21 columns

```
project_data.drop(['project_essay_1', 'project_essay_2', 'project_essay_3', 'project_essay_4'],
project_data.head(2)
```



	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	pr
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL	

```
#Replacing Nan's with maximum occurred value: https://stackoverflow.com/a/51053916/8089731
project_data['teacher_prefix'].value_counts().argmax()
project_data.fillna(value=project_data['teacher_prefix'].value_counts().argmax(),axis=1,inpla
```

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

```
project_data.columns
```

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'project_submitted_datetime', 'project_title',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'essay', 'price', 'quantity',
      'digits_in_summary', 'clean_project_grade_category',
      'preprocessed_essays', 'preprocessed_titles'],
      dtype='object')
```

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label
```

## 2.2 Make Data Model Ready: encoding numerical, categorical features

```
X_train, X_test, y_train, y_test = train_test_split(project_data,project_data['project_is_app
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_
```

```
X_train.drop(['project_is_approved'], axis=1, inplace=True)
X_test.drop(['project_is_approved'], axis=1, inplace=True)
X_cv.drop(['project_is_approved'], axis=1, inplace=True)
```

### ▼ 1.4.1 Vectorizing Categorical data


```
# we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary
vectorizer.fit(X_train['clean_categories'].values)
print(vectorizer.get_feature_names())
```

```
categories_one_hot_train = vectorizer.transform(X_train['clean_categories'].values)
categories_one_hot_cv = vectorizer.transform(X_cv['clean_categories'].values)
categories_one_hot_test = vectorizer.transform(X_test['clean_categories'].values)
print("Shape of matrix after one hot encodig_train ",categories_one_hot_train.shape)
print("Shape of matrix after one hot encodig_cv ",categories_one_hot_cv.shape)
print("Shape of matrix after one hot encodig_test ",categories_one_hot_test.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNee
Shape of matrix after one hot encodig_train (22445, 9)
Shape of matrix after one hot encodig_cv (11055, 9)
Shape of matrix after one hot encodig_test (16500, 9)
```


```
# we use count vectorizer to convert the values into one hot encoded features
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, bi
vectorizer.fit(X_train['clean_subcategories'].values)
print(vectorizer.get_feature_names())
```

```
sub_categories_one_hot_train = vectorizer.transform(X_train['clean_subcategories'].values)
sub_categories_one_hot_cv = vectorizer.transform(X_cv['clean_subcategories'].values)
sub_categories_one_hot_test = vectorizer.transform(X_test['clean_subcategories'].values)
print("Shape of matrix after one hot encodig_train ",sub_categories_one_hot_train.shape)
print("Shape of matrix after one hot encodig_cv ",sub_categories_one_hot_cv.shape)
print("Shape of matrix after one hot encodig_test ",sub_categories_one_hot_test.shape)
```

```
 ['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurric
Shape of matrix after one hot encodig_train (22445, 30)
Shape of matrix after one hot encodig_cv (11055, 30)
Shape of matrix after one hot encodig_test (16500, 30)
```

```
# we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer( lowercase=False, binary=True)
vectorizer.fit(X_train['school_state'].values)
print(vectorizer.get_feature_names())
```

```
school_state_one_hot_train = vectorizer.transform(X_train['school_state'].values)
school_state_one_hot_cv = vectorizer.transform(X_cv['school_state'].values)
school_state_one_hot_test = vectorizer.transform(X_test['school_state'].values)
print("Shape of matrix after one hot encodig_train ",school_state_one_hot_train.shape)
print("Shape of matrix after one hot encodig_cv ",school_state_one_hot_cv.shape)
print("Shape of matrix after one hot encodig_test ",school_state_one_hot_test.shape)
```

```
 ['AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DC', 'DE', 'FL', 'GA', 'HI', 'IA', 'ID', 'IL
Shape of matrix after one hot encodig_train (22445, 51)
Shape of matrix after one hot encodig_cv (11055, 51)
Shape of matrix after one hot encodig_test (16500, 51)
```

```
# #Replacing Nan's with maximum occured value: https://stackoverflow.com/a/51053916/8089731
# project_data['teacher_prefix'].value_counts().argmax()
# project_data.fillna(value=project_data['teacher_prefix'].value_counts().argmax(),axis=1,inp
```

```
# we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
```

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer( lowercase=False, binary=True)
vectorizer.fit(X_train['teacher_prefix'].values.astype('U'))
print(vectorizer.get_feature_names())
```

#<https://stackoverflow.com/a/39308809/8089731>

```
teacher_prefix_one_hot_train = vectorizer.transform(X_train['teacher_prefix'].values.astype('U'))
teacher_prefix_one_hot_cv = vectorizer.transform(X_cv['teacher_prefix'].values.astype('U'))
teacher_prefix_one_hot_test = vectorizer.transform(X_test['teacher_prefix'].values.astype('U'))
print("Shape of matrix after one hot encodig_train ",teacher_prefix_one_hot_train.shape)
print("Shape of matrix after one hot encodig_cv ",teacher_prefix_one_hot_cv.shape)
print("Shape of matrix after one hot encodig_test ",teacher_prefix_one_hot_test.shape)
```

```

[ 'Dr', 'Mr', 'Mrs', 'Ms', 'Teacher' ]
Shape of matrix after one hot encodig_train (22445, 5)
Shape of matrix after one hot encodig_cv (11055, 5)
Shape of matrix after one hot encodig_test (16500, 5)
```

```
print(project_data['clean_project_grade_category'].unique())
```

```

[ 'GradesPreK-2' 'Grades6-8' 'Grades3-5' 'Grades9-12' ]
```

# we use count vectorizer to convert the values into one hot encoded features

```
from sklearn.feature_extraction.text import CountVectorizer
```

# <https://stackoverflow.com/a/38161028/8089731>

```
pattern = "(?u)\\b[\\w-]+\\b"
vectorizer = CountVectorizer(token_pattern=pattern, lowercase=False, binary=True)
vectorizer.fit(X_train['clean_project_grade_category'].values)
print(vectorizer.get_feature_names())
```

#<https://stackoverflow.com/a/39308809/8089731>

```
project_grade_category_one_hot_train = vectorizer.transform(X_train['clean_project_grade_category'])
project_grade_category_one_hot_cv = vectorizer.transform(X_cv['clean_project_grade_category'])
project_grade_category_one_hot_test = vectorizer.transform(X_test['clean_project_grade_category'])
print("Shape of matrix after one hot encodig_train ",project_grade_category_one_hot_train.shape)
print("Shape of matrix after one hot encodig_cv ",project_grade_category_one_hot_cv.shape)
print("Shape of matrix after one hot encodig_test ",project_grade_category_one_hot_test[:5,:])
```

```

[ 'Grades3-5', 'Grades6-8', 'Grades9-12', 'GradesPreK-2' ]
Shape of matrix after one hot encodig_train (22445, 4)
Shape of matrix after one hot encodig_cv (11055, 4)
Shape of matrix after one hot encodig_test (0, 0) 1
(1, 3) 1
(2, 1) 1
(3, 3) 1
(4, 1) 1
```

## ▼ Vectorizing Numerical features

# check this one: <https://www.youtube.com/watch?v=0H0q0c1n3Z4&t=530s>

# standardization sklearn: <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html>



```

from sklearn.preprocessing import StandardScaler

# price_standardized = standardScalar.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399.
# Reshape your data either using array.reshape(-1, 1)

price_scalar = StandardScaler()
price_scalar.fit(X_train['price'].values.reshape(-1,1)) # finding the mean and standard devia
# print(f"Mean : {price_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])}

# Now standardize the data with above maen and variance.
price_standardized_train = price_scalar.transform(X_train['price'].values.reshape(-1, 1))
price_standardized_cv = price_scalar.transform(X_cv['price'].values.reshape(-1, 1))
price_standardized_test = price_scalar.transform(X_test['price'].values.reshape(-1, 1))
print(price_standardized_train.shape)
print(price_standardized_cv.shape)
print(price_standardized_test.shape)

(22445, 1)
(11055, 1)
(16500, 1)

```

```

# check this one: https://www.youtube.com/watch?v=0H0qOcln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler
from sklearn.preprocessing import StandardScaler

```

```

# price_standardized = standardScalar.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399.
# Reshape your data either using array.reshape(-1, 1)

quantity_scalar = StandardScaler()
quantity_scalar.fit(X_train['quantity'].values.reshape(-1,1)) # finding the mean and standard
# print(f"Mean : {quantity_scalar.mean_[0]}, Standard deviation : {np.sqrt(quantity_scalar.va

# Now standardize the data with above maen and variance.
quantity_standardized_train = quantity_scalar.transform(X_train['quantity'].values.reshape(-1
quantity_standardized_cv = quantity_scalar.transform(X_cv['quantity'].values.reshape(-1, 1))
quantity_standardized_test = quantity_scalar.transform(X_test['quantity'].values.reshape(-1,
print(quantity_standardized_train.shape)
print(quantity_standardized_cv.shape)
print(quantity_standardized_test.shape)

```

```
C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
(22445, 1)
(11055, 1)
(16500, 1)
```

```
# check this one: https://www.youtube.com/watch?v=0H0qOcln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler
from sklearn.preprocessing import StandardScaler
```

```
# price_standardized = standardScaler.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399.]
# Reshape your data either using array.reshape(-1, 1)
```

```
teacher_number_of_previously_posted_projects_scalar = StandardScaler()
teacher_number_of_previously_posted_projects_scalar.fit(X_train['teacher_number_of_previously_posted_projects'])
# print(f"Mean : {teacher_number_of_previously_posted_projects_scalar.mean_[0]}, Standard dev
```

```
# Now standardize the data with above mean and variance.
teacher_number_of_previously_posted_projects_standardized_train = teacher_number_of_previously_posted_projects_scalar.transform(X_train['teacher_number_of_previously_posted_projects'])
teacher_number_of_previously_posted_projects_standardized_cv = teacher_number_of_previously_posted_projects_scalar.transform(X_cv['teacher_number_of_previously_posted_projects'])
teacher_number_of_previously_posted_projects_standardized_test = teacher_number_of_previously_posted_projects_scalar.transform(X_test['teacher_number_of_previously_posted_projects'])
print(teacher_number_of_previously_posted_projects_standardized_train.shape)
print(teacher_number_of_previously_posted_projects_standardized_cv.shape)
print(teacher_number_of_previously_posted_projects_standardized_test.shape)
```



```
C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.

(22445, 1)
(11055, 1)
(16500, 1)
```

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly
```

```
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

## 2.3 Make Data Model Ready: encoding eassay, and project\_title

```
X_train.head(2)
```



## ▼ Bag of Words(BOW) on project\_TEXT/ESSAYS (Train,Cv,Test)

```
# We are considering only the words which appeared in at least 10 documents(rows or projects)
vectorizer = CountVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_essays'])
```

```
text_bow_train = vectorizer.transform(X_train['preprocessed_essays'])
text_bow_cv = vectorizer.transform(X_cv['preprocessed_essays'])
text_bow_test = vectorizer.transform(X_test['preprocessed_essays'])
print("Shape of matrix after BOW_text_train ",text_bow_train.shape)
print("Shape of matrix after BOW_text_cv ",text_bow_cv.shape)
print("Shape of matrix after BOW_text_test ",text_bow_test.shape)
```



```
Shape of matrix after BOW_text_train (22445, 8894)
Shape of matrix after BOW_text_cv (11055, 8894)
Shape of matrix after BOW_text_test (16500, 8894)
```

## ▼ Bag of Words(BOW) on project\_title (Train,Cv,Test)

```
# We are considering only the words which appeared in at least 10 documents(rows or projects)
vectorizer = CountVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_titles'])
```

```
title_bow_train = vectorizer.transform(X_train['preprocessed_titles'])
title_bow_cv = vectorizer.transform(X_cv['preprocessed_titles'])
title_bow_test = vectorizer.transform(X_test['preprocessed_titles'])
print("Shape of matrix after BOW_title_train ",title_bow_train.shape)
print("Shape of matrix after BOW_title_cv ",title_bow_cv.shape)
print("Shape of matrix after BOW_title_test ",title_bow_test.shape)
```



```
Shape of matrix after BOW_title_train (22445, 1249)
Shape of matrix after BOW_title_cv (11055, 1249)
Shape of matrix after BOW_title_test (16500, 1249)
```

## ▼ TFIDF Vectorizer on project\_TEXT/ESSAYS (Train,Cv,Test)

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_essays'])
```

```
text_tfidf_train = vectorizer.transform(X_train['preprocessed_essays'])
text_tfidf_cv = vectorizer.transform(X_cv['preprocessed_essays'])
text_tfidf_test = vectorizer.transform(X_test['preprocessed_essays'])
print("Shape of matrix after tfidf_text_train ",text_tfidf_train.shape)
print("Shape of matrix after tfidf_text_cv ",text_tfidf_cv.shape)
```

```
print("Shape of matrix after tfidf_text_test ",text_tfidf_test.shape)
```



```
Shape of matrix after tfidf_text_train (22445, 8894)
Shape of matrix after tfidf_text_cv (11055, 8894)
Shape of matrix after tfidf_text_test (16500, 8894)
```

## ▼ TFIDF Vectorizer on project\_title(Train,Cv,Test)

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_titles'])
```

```
title_tfidf_train = vectorizer.transform(X_train['preprocessed_titles'])
title_tfidf_cv = vectorizer.transform(X_cv['preprocessed_titles'])
title_tfidf_test = vectorizer.transform(X_test['preprocessed_titles'])
print("Shape of matrix after tfidf_title_train ",title_tfidf_train.shape)
print("Shape of matrix after tfidf_title_cv ",title_tfidf_cv.shape)
print("Shape of matrix after tfidf_title_test ",title_tfidf_test.shape)
```



```
Shape of matrix after tfidf_title_train (22445, 1249)
Shape of matrix after tfidf_title_cv (11055, 1249)
Shape of matrix after tfidf_title_test (16500, 1249)
```

```
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-
# make sure you have the glove_vectors file
with open('../glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())
```

## ▼ Avg W2V on TEXT/ESSAYS(Train,cv,test)

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_essays_vectors_train = []; # the avg-w2v for each sentence/review is stored in this l
for sentence in tqdm1(X_train['preprocessed_essays']): # for each review/sentence
    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_essays_vectors_train.append(vector)
```

```
avg_w2v_essays_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm1(X_cv['preprocessed_essays']): # for each review/sentence
    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_essays_vectors_cv.append(vector)

```

```

avg_w2v_essays_vectors_test = []; # the avg-w2v for each sentence/review is stored in this li
for sentence in tqdm1(X_test['preprocessed_essays']): # for each review/sentence
    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_essays_vectors_test.append(vector)

```

```

print(len(avg_w2v_essays_vectors_train))
print(len(avg_w2v_essays_vectors_cv))
print(len(avg_w2v_essays_vectors_test))
print(len(avg_w2v_essays_vectors_train[0]))
print(len(avg_w2v_essays_vectors_cv[0]))
print(len(avg_w2v_essays_vectors_test[0]))

```



```
HBox(children=(IntProgress(value=0, max=22445), HTML(value='')))
```

```
HBox(children=(IntProgress(value=0, max=11055), HTML(value='')))
```

```
HBox(children=(IntProgress(value=0, max=16500), HTML(value='')))
```

```
22445
```

```
11055
```

```
16500
```

```
300
```

```
300
```

```
300
```

## ▼ Avg W2V on TITLES(Train,cv,test)

```

# average Word2Vec
# compute average word2vec for each review.
avg_w2v_titles_vectors_train = []; # the avg-w2v for each sentence/review is stored in this l
for sentence in tqdm1(X_train['preprocessed_titles']): # for each review/sentence
    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

```

```

    1+ word in glove_words:
        vector += model[word]
        cnt_words += 1
if cnt_words != 0:
    vector /= cnt_words
avg_w2v_titles_vectors_train.append(vector)

avg_w2v_titles_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm1(X_cv['preprocessed_titles']): # for each review/sentence
    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_titles_vectors_cv.append(vector)

avg_w2v_titles_vectors_test = []; # the avg-w2v for each sentence/review is stored in this li
for sentence in tqdm1(X_test['preprocessed_titles']): # for each review/sentence
    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_titles_vectors_test.append(vector)

print(len(avg_w2v_titles_vectors_train))
print(len(avg_w2v_titles_vectors_cv))
print(len(avg_w2v_titles_vectors_test))
print(len(avg_w2v_titles_vectors_train[0]))
print(len(avg_w2v_titles_vectors_cv[0]))
print(len(avg_w2v_titles_vectors_test[0]))

```



HBox(children=(IntProgress(value=0, max=22445), HTML(value='')))

HBox(children=(IntProgress(value=0, max=11055), HTML(value='')))

HBox(children=(IntProgress(value=0, max=16500), HTML(value='')))

22445

11055

16500

300

300

300

## ▼ TFIDF weighted W2V on TEXT/ESSAYS(Train,cv,test)

```
# 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_essays_vectors_train = []; # the avg-w2v for each sentence/review is stored in this
for sentence in tqdm1(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((sentence.
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_essays_vectors_train.append(vector)

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_essays_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this li
for sentence in tqdm1(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((sentence.
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_essays_vectors_cv.append(vector)

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_essays_vectors_test = []; # the avg-w2v for each sentence/review is stored in this
for sentence in tqdm1(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
```



```

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.
        tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting
        vector += (vec * tf_idf) # calculating tfidf weighted w2v
        tf_idf_weight += tf_idf
if tf_idf_weight != 0:
    vector /= tf_idf_weight
tfidf_w2v_essays_vectors_test.append(vector)

print(len(tfidf_w2v_essays_vectors_train))
print(len(tfidf_w2v_essays_vectors_cv))
print(len(tfidf_w2v_essays_vectors_test))
print(len(tfidf_w2v_essays_vectors_train[0]))
print(len(tfidf_w2v_essays_vectors_cv[0]))
print(len(tfidf_w2v_essays_vectors_test[0]))

```



HBox(children=(IntProgress(value=0, max=22445), HTML(value='')))

HBox(children=(IntProgress(value=0, max=11055), HTML(value='')))

HBox(children=(IntProgress(value=0, max=16500), HTML(value='')))

22445

11055

16500

300

300

300

## ▼ TFIDF weighted W2V on TEXT/ESSAYS(Train,cv,test)

```

# 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_essays_vectors_train = []; # the avg-w2v for each sentence/review is stored in th
# for sentence in tqdm1(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((sentenc
#             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # gettin

```

```

#         vector += (vec * tt_idf) # calculating ttidf weighted w2v
#         tf_idf_weight += tf_idf
#     if tf_idf_weight != 0:
#         vector /= tf_idf_weight
#     tfidf_w2v_essays_vectors_train.append(vector)

# # average Word2Vec
# # compute average word2vec for each review.
# tfidf_w2v_essays_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this
# for sentence in tqdm1(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((sentence
#             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # gettin
#             vector += (vec * tf_idf) # calculating tfidf weighted w2v
#             tf_idf_weight += tf_idf
#     if tf_idf_weight != 0:
#         vector /= tf_idf_weight
#     tfidf_w2v_essays_vectors_cv.append(vector)

# # average Word2Vec
# # compute average word2vec for each review.
# tfidf_w2v_essays_vectors_test = []; # the avg-w2v for each sentence/review is stored in thi
# for sentence in tqdm1(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((sentenc
#             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # gettin
#             vector += (vec * tf_idf) # calculating tfidf weighted w2v
#             tf_idf_weight += tf_idf
#     if tf_idf_weight != 0:
#         vector /= tf_idf_weight
#     tfidf_w2v_essays_vectors_test.append(vector)

# print(len(tfidf_w2v_essays_vectors_train))
# print(len(tfidf_w2v_essays_vectors_cv))
# print(len(tfidf_w2v_essays_vectors_test))
# print(len(tfidf_w2v_essays_vectors_train[0]))
# print(len(tfidf_w2v_essays_vectors_cv[0]))
# print(len(tfidf_w2v_essays_vectors_test[0]))

```

## ▼ TFIDF weighted W2V on TITLES(Train,cv,test)

```

# 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_titles_vectors_train = []; # the avg-w2v for each sentence/review is stored in this
for sentence in tqdm1(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((sentence.
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_titles_vectors_train.append(vector)

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_titles_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this li
for sentence in tqdm1(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((sentence.
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_titles_vectors_cv.append(vector)

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_titles_vectors_test = []; # the avg-w2v for each sentence/review is stored in this
for sentence in tqdm1(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((sentence.
    tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting
    vector += (vec * tf_idf) # calculating tfidf weighted w2v
    tf_idf_weight += tf_idf
if tf_idf_weight != 0:
    vector /= tf_idf_weight
tfidf_w2v_titles_vectors_test.append(vector)

```

```

print(len(tfidf_w2v_titles_vectors_train))
print(len(tfidf_w2v_titles_vectors_cv))
print(len(tfidf_w2v_titles_vectors_test))
print(len(tfidf_w2v_titles_vectors_train[0]))
print(len(tfidf_w2v_titles_vectors_cv[0]))
print(len(tfidf_w2v_titles_vectors_test[0]))

```



```
HBox(children=(IntProgress(value=0, max=22445), HTML(value='')))
```

```
HBox(children=(IntProgress(value=0, max=11055), HTML(value='')))
```

```
HBox(children=(IntProgress(value=0, max=16500), HTML(value='')))
```

```

22445
11055
16500
300
300
300

```

```

# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly

```

```

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label

```

## 2.4 Applying KNN on different kind of featurization as mentione

Apply KNN on different kind of featurization as mentioned in the instructions

For Every model that you work on make sure you do the step 2 and step 3 of instructions

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

## ► 2.4.1 Applying KNN brute force on BOW, SET 1

↳ 2 cells hidden

### ▼ 1.1 Method 1: Simple for loop (if you are having memory limitations us

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # 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 cr_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
    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
    return y_data_pred

import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score
"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence values,
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

train_auc = []
cv_auc = []
```

```

K = [1, 5, 10, 15, 21, 31, 41, 51, 101]
for i in (K):
    neigh = KNeighborsClassifier(n_neighbors=i)
    neigh.fit(X_tr[:, :], y_train[:])

    y_train_pred = batch_predict(neigh, X_tr[:, :])
    y_cv_pred = batch_predict(neigh, X_cr[:, :])

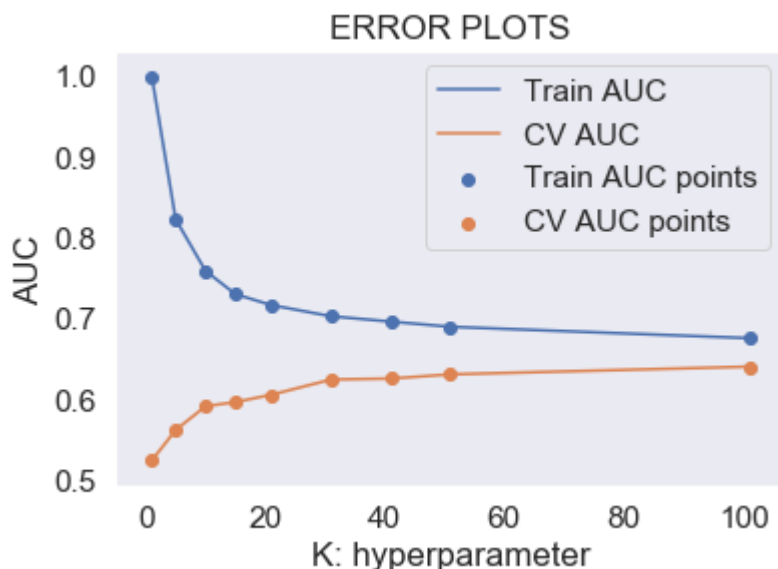
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train[:, :], y_train_pred))
    cv_auc.append(roc_auc_score(y_cv[:, :], y_cv_pred))

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()

```



# [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\\_curve.html#sklearn.me](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve)  
 from sklearn.metrics import roc\_curve, auc

```

neigh = KNeighborsClassifier(n_neighbors=51)
neigh.fit(X_tr[:, :], y_train[:])
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the pos
# not the predicted outputs

```

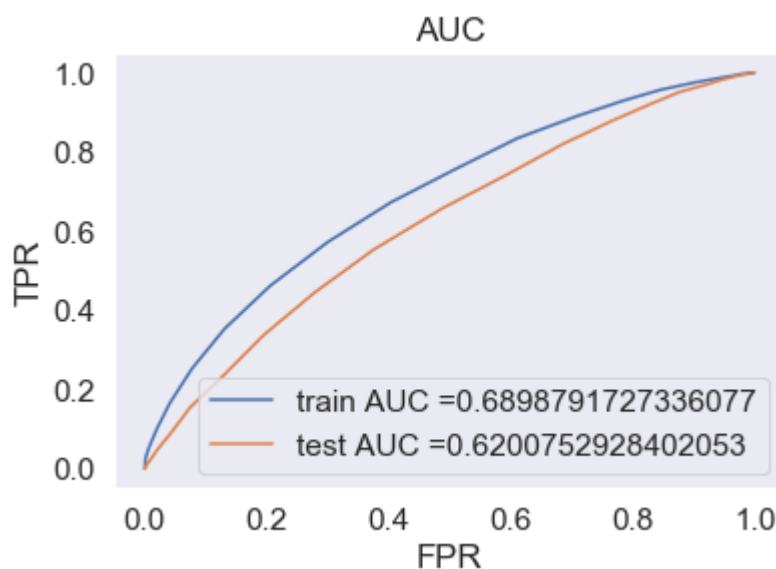
```

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_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("FPR")
plt.ylabel("TPR")
plt.title("AUC")
plt.grid()
plt.show()

```



```

# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, 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", np.round(t,3))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions

print("="*100)
from sklearn.metrics import confusion_matrix

```

```
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr))
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

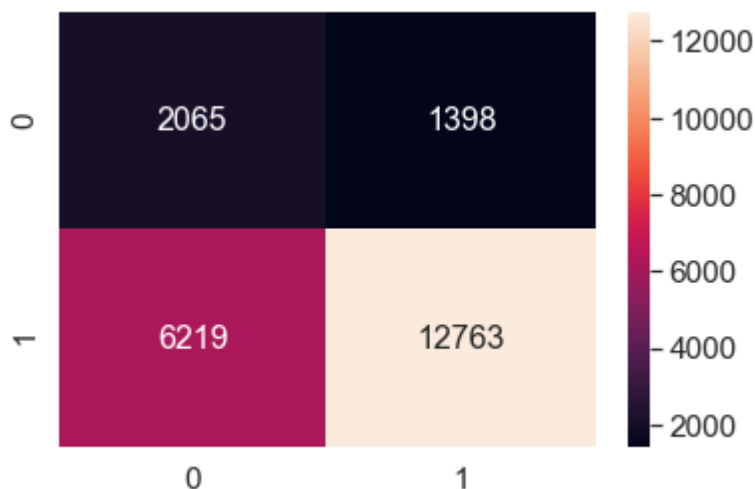


```
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.4009390570266228 for threshold 0.784
[[ 2065  1398]
 [ 6219 12763]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.3453216161597511 for threshold 0.804
[[1591  955]
 [6243 7711]]
```

```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thre
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```



```
the maximum value of tpr*(1-fpr) 0.4009390570266228 for threshold 0.784
<matplotlib.axes._subplots.AxesSubplot at 0x19785933240>
```

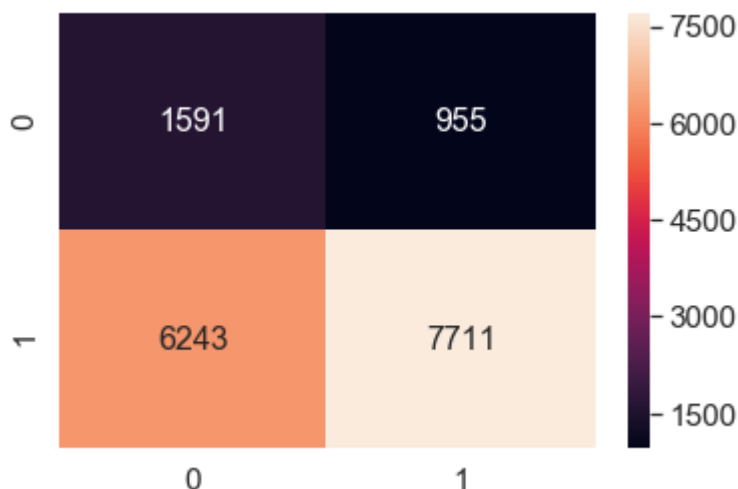


```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresho
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```





the maximum value of  $\text{tpr} \cdot (1 - \text{fpr})$  0.3453216161597511 for threshold 0.804  
 <matplotlib.axes.\_subplots.AxesSubplot at 0x197853de940>



```
print(train_fpr.shape)
print(train_tpr.shape)
print(len(y_train_pred))
```

(29,)  
 (29,)  
 22445

## ► 2.4.2 Applying KNN brute force on TFIDF, SET 2

↳ 2 cells hidden

## ▼ 1.1 Method 1: Simple for loop (if you are having memory limitations us

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # 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 cr_loop will be 49041 - 49041%1000 = 49000
    # in this for loop we will iterate until 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
    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
    return y_data_pred
```

```
import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
```

```
from sklearn.metrics import roc_auc_score
"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence values,
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

train_auc = []
cv_auc = []
K = [1, 5, 10, 15, 21, 31, 41, 51, 101]
for i in (K):
    neigh = KNeighborsClassifier(n_neighbors=i)
    neigh.fit(X_tr[:, :], y_train[:])

    y_train_pred = batch_predict(neigh, X_tr[:, :])
    y_cv_pred = batch_predict(neigh, X_cr[:, :])

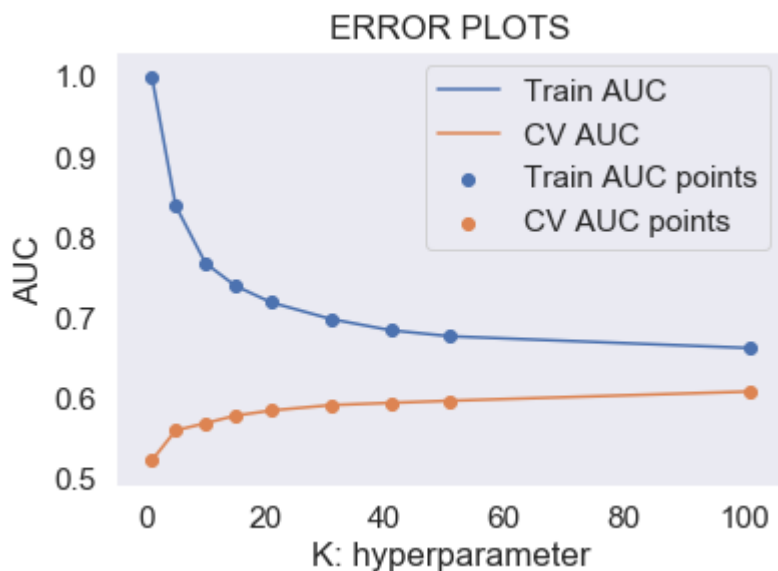
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train[:], y_train_pred))
    cv_auc.append(roc_auc_score(y_cv[:], y_cv_pred))

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()
```





# [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\\_curve.html#sklearn.me](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.me)

```
from sklearn.metrics import roc_curve, auc
from sklearn.neighbors import KNeighborsClassifier
```

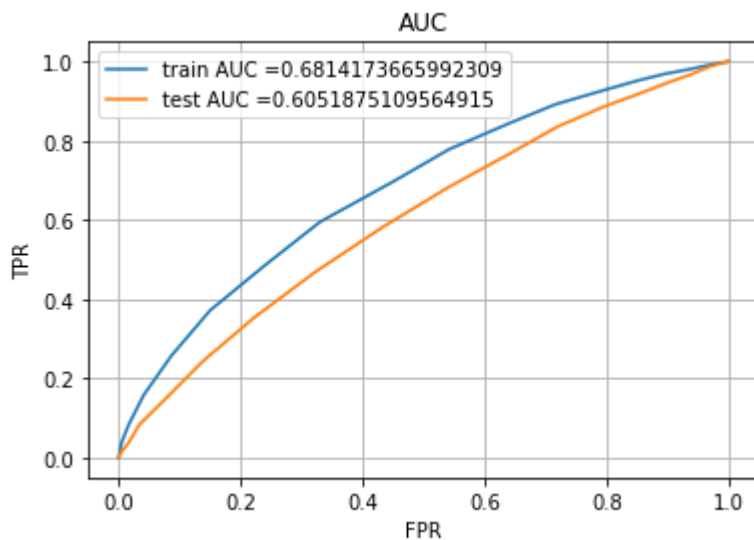
```
neigh = KNeighborsClassifier(n_neighbors=55)
neigh.fit(X_tr[:, :], y_train[:])
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the pos
# not the predicted outputs
```

```
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_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("FPR")
plt.ylabel("TPR")
plt.title("AUC")
plt.grid()
plt.show()
```





```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, 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", np.round(t,3))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions

print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

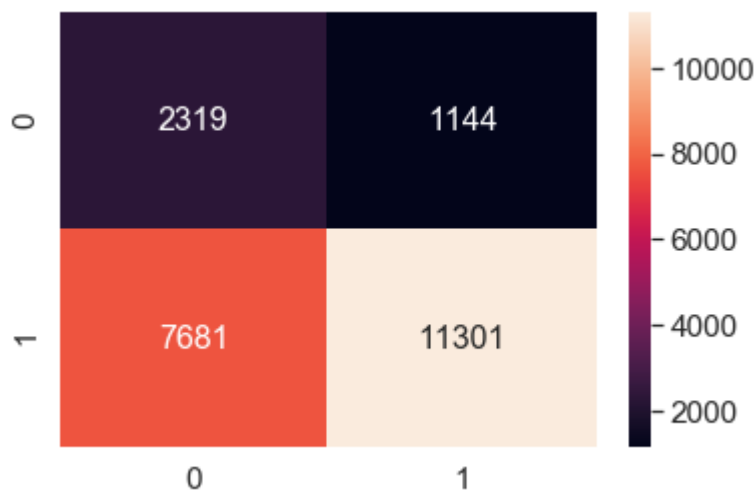


```
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.3986788188746559 for threshold 0.855
[[ 2319  1144]
 [ 7681 11301]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.33006483202973835 for threshold 0.855
[[1445  1101]
 [5839  8115]]
```

```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thre
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```



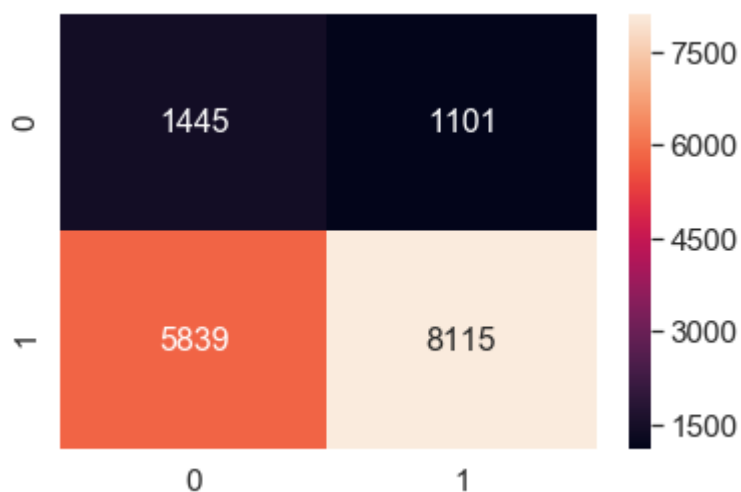
the maximum value of  $tpr*(1-fpr)$  0.3986788188746559 for threshold 0.855  
 <matplotlib.axes.\_subplots.AxesSubplot at 0x1d709a211d0>



```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresho
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```



the maximum value of  $tpr*(1-fpr)$  0.33006483202973835 for threshold 0.855  
 <matplotlib.axes.\_subplots.AxesSubplot at 0x1d71169c780>



### ► 2.4.3 Applying KNN brute force on AVG W2V, SET 3

↳ 2 cells hidden

### ▼ 1.1 Method 1: Simple for loop (if you are having memory limitations us

```
def batch_predict(clf, data):
```

```
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
# 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 cr_loop will be 49041 - 49041%1000 = 49000
# in this for loop we will iterate until 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
y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
return y_data_pred
```

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

```
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.
```

```
y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence values,
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.
```

```
"""
```

```
train_auc = []
cv_auc = []
K = [1, 5, 10, 15, 21, 31, 41, 51, 101]
start_time = time.time()
for i in (K):
    neigh = KNeighborsClassifier(n_neighbors=i)
    neigh.fit(X_tr[:5747,:],y_train[:5747])

    y_train_pred = batch_predict(neigh, X_tr[:5747,:])
    y_cv_pred = batch_predict(neigh, X_cr[:5747,:])

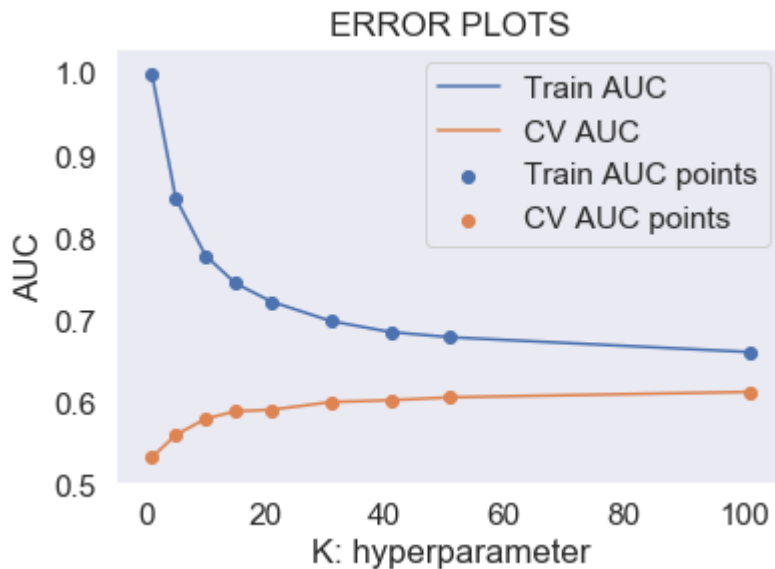
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train[:5747],y_train_pred))
    cv_auc.append(roc_auc_score(y_cv[:5747], y_cv_pred))

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()
print("Execution time: " + str((time.time() - start_time)) + ' ms')
```



Execution time: 981.1687150001526 ms

# [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\\_curve.html#sklearn.me](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.me)

```
from sklearn.metrics import roc_curve, auc
from sklearn.neighbors import KNeighborsClassifier
```

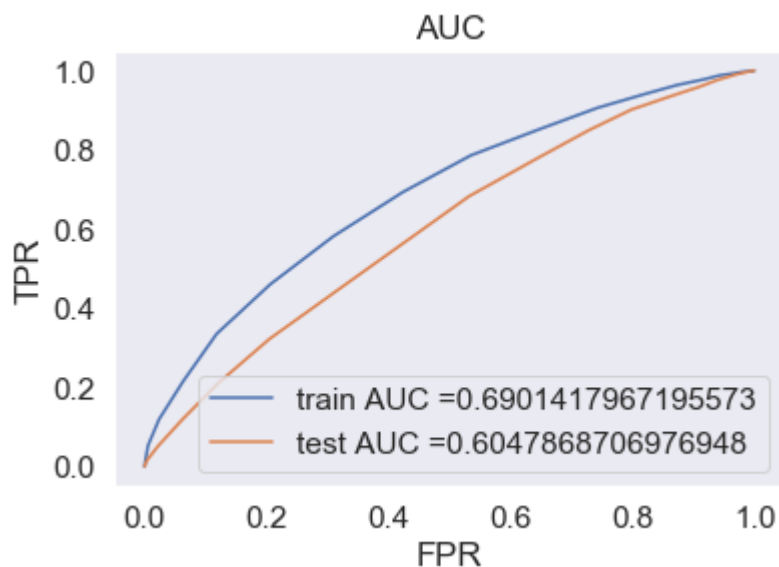
```
start_time = time.time()
neigh = KNeighborsClassifier(n_neighbors=53)
neigh.fit(X_tr[:, :], y_train[:])
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the pos
# not the predicted outputs
```

```
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_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("FPR")
plt.ylabel("TPR")
plt.title("AUC")
plt.grid()
plt.show()
print("Execution time: " + str((time.time() - start_time)) + ' ms')
```





Execution time: 1464.733214378357 ms

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, 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", np.round(t,3)
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions

print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```





```
=====
```

Train confusion matrix

the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.40157172168487176 for threshold 0.868

```
[[ 2393  1070]
 [ 7951 11031]]
```

Test confusion matrix

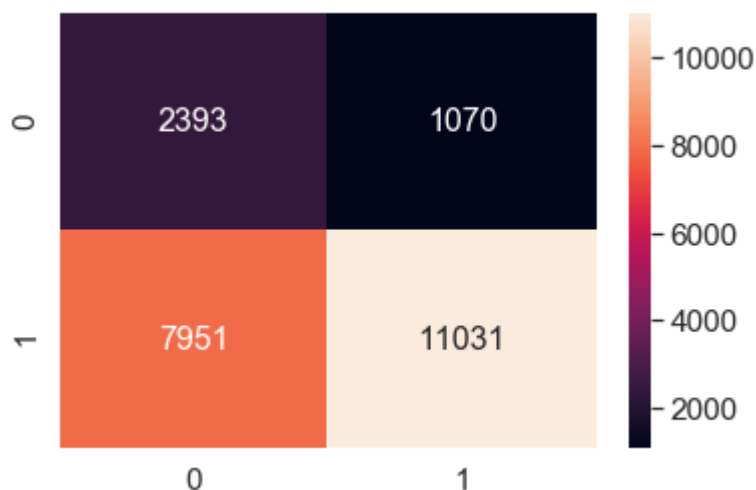
the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.32514391073531806 for threshold 0.868

```
[[1453 1093]
 [6004 7950]]
```

```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thre
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```



the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.40157172168487176 for threshold 0.868  
<matplotlib.axes.\_subplots.AxesSubplot at 0x1d712735048>



```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresho
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```



## ▶ 2.4.4 Applying KNN brute force on TFIDF W2V, SET 4

↳ 2 cells hidden

## ▼ 1.1 Method 1: Simple for loop (if you are having memory limitations us

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # 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 cr_loop will be 49041 - 49041%1000 = 49000
    # in this for loop we will iterate until 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
    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
    return y_data_pred

import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score
"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence values,
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

train_auc = []
cv_auc = []
K = [1, 5, 10, 15, 21, 31, 41, 51, 101]
start_time = time.time()
for i in (K):
    neigh = KNeighborsClassifier(n_neighbors=i)
    neigh.fit(X_tr[:5747,:], y_train[:5747])

    y_train_pred = batch_predict(neigh, X_tr[:5747,:])
    y_cv_pred = batch_predict(neigh, X_cr[:5747,:])

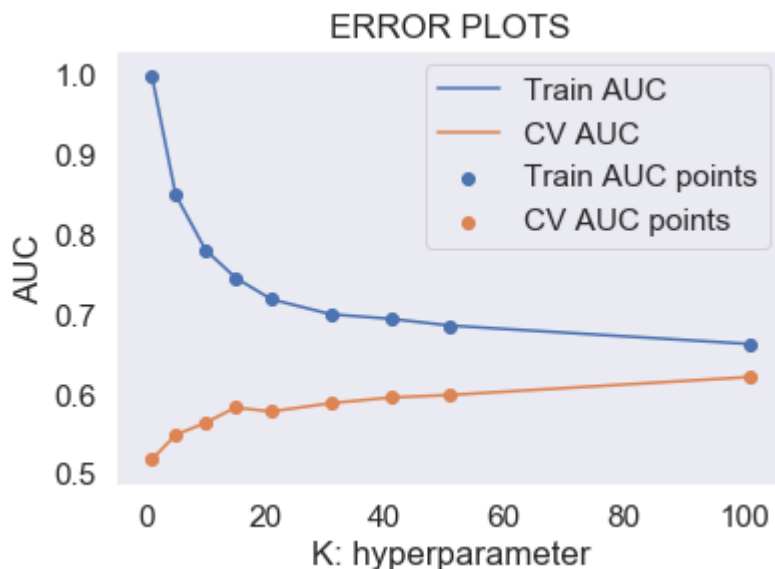
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # not the predicted outputs
```

```
train_auc.append(roc_auc_score(y_train[:5747],y_train_pred))
cv_auc.append(roc_auc_score(y_cv[:5747], 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()
print("Execution time: " + str((time.time() - start_time)) + ' ms')
```



Execution time: 996.8478593826294 ms

# [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\\_curve.html#sklearn.me](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.me)  
 from sklearn.metrics import roc\_curve, auc

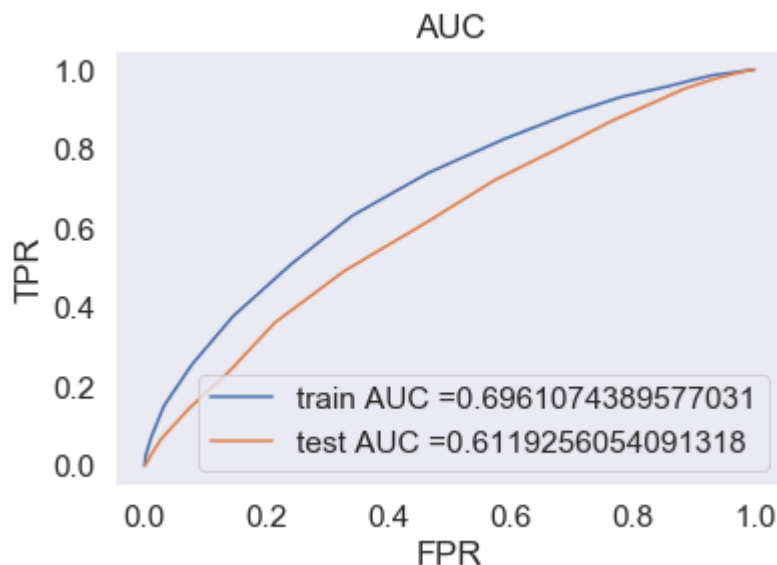
```
start_time = time.time()
neigh = KNeighborsClassifier(n_neighbors=50)
neigh.fit(X_tr[:, :], y_train[:])
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the pos
# not the predicted outputs
```

```
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_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
```

```
plt.legend()
plt.xlabel("FPR")
plt.ylabel("TPR")
plt.title("AUC")
plt.grid()
plt.show()
print("Execution time: " + str((time.time() - start_time)) + ' ms')
```



Execution time: 1695.6296381950378 ms

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, 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", np.round(t,3)
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions

print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```



```
=====
```

Train confusion matrix

the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.41619309361060725 for threshold 0.86

```
[[ 2281  1182]
 [ 6988 11994]]
```

Test confusion matrix

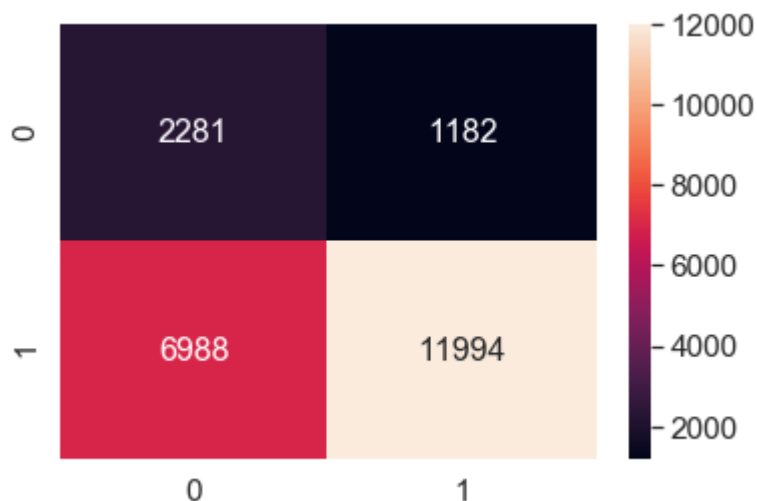
the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.3307163104988324 for threshold 0.86

```
[[1380 1166]
 [5440 8514]]
```

```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thre
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```



the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.41619309361060725 for threshold 0.86  
<matplotlib.axes.\_subplots.AxesSubplot at 0x1d711d98c18>



```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresho
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```



## 2.5 Feature selection with `SelectKBest`

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train,
              ,project_grade_category_one_hot_train,price_standardized_train,quantity_standardized_train,
              ,teacher_number_of_previously_posted_projects_standardized_train,text_tfidf_train))
X_cr = hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv,teacher_number_of_previously_posted_projects_standardized_cv,
              ,project_grade_category_one_hot_cv,price_standardized_cv,quantity_standardized_cv,
              ,teacher_number_of_previously_posted_projects_standardized_cv,text_tfidf_cv,text_tfidf_test))
X_te = hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test,
              ,project_grade_category_one_hot_test,price_standardized_test,quantity_standardized_test,
              ,teacher_number_of_previously_posted_projects_standardized_test,text_tfidf_test))

#####
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
X_tr = scaler.fit_transform(X_tr,y_train)
X_cr = scaler.transform(X_cr)
X_te = scaler.transform(X_te)
#####
from sklearn.feature_selection import SelectKBest, chi2
t = SelectKBest(chi2,k=2000).fit(X_tr, y_train)
X_tr = t.transform(X_tr)
X_te = t.transform(X_te)
X_cr = t.transform(X_cr)
#####
print("Final Data matrix on TFIDF")
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 on TFIDF

(22445, 2000) (22445,)

(11055, 2000) (11055,)

(16500, 2000) (16500,)

=====

```
# #####
# from sklearn.preprocessing import MinMaxScaler
# scaler = MinMaxScaler()
# X_tr = scaler.fit_transform(X_tr)
# X_cr = scaler.fit_transform(X_cr)
# X_te = scaler.fit_transform(X_te)
# #####
# from sklearn.feature_selection import SelectKBest, chi2
# t = SelectKBest(chi2,k=2000)
# X_tr = t.fit_transform(X_tr,y_train)
# X_te = t.fit_transform(X_te,y_test)
# X_cr = t.fit_transform(X_cr,y_cv)
# #####
# print("Final Data matrix on TFIDF")
# print(X_tr.shape, y_train.shape)
# print(X_cr.shape, y_cv.shape)
# print(X_te.shape, y_test.shape)
# print("="*100)
```

## ▼ 1.1 Method 1: Simple for loop (if you are having memory limitations us

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # 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 cr_loop will be 49041 - 49041%1000 = 49000
    # in this for loop we will iterate until 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
    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
    return y_data_pred

import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score
"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.
```

```
y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence values,
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

train_auc = []
cv_auc = []
K = [1, 5, 10, 15, 21, 31, 41, 51, 101]
start_time = time.time()
for i in (K):
    neigh = KNeighborsClassifier(n_neighbors=i)
    neigh.fit(X_tr[:5747,:],y_train[:5747])

    y_train_pred = batch_predict(neigh, X_tr[:5747,:])
    y_cv_pred = batch_predict(neigh, X_cr[:5747,:])

    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train[:5747],y_train_pred))
    cv_auc.append(roc_auc_score(y_cv[:5747], y_cv_pred))

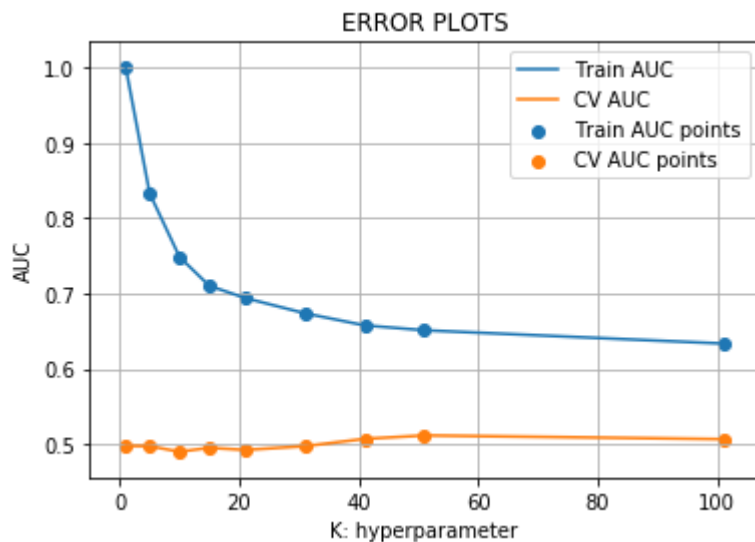
plt.plot(K, train_auc, label='Train AUC')
plt.plot(K, cv_auc, label='CV AUC')

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()
print("Execution time: " + str((time.time() - start_time)) + ' ms')
```







Execution time: 2072.0459604263306 ms

```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.me
from sklearn.metrics import roc_curve, auc
from sklearn.neighbors import KNeighborsClassifier

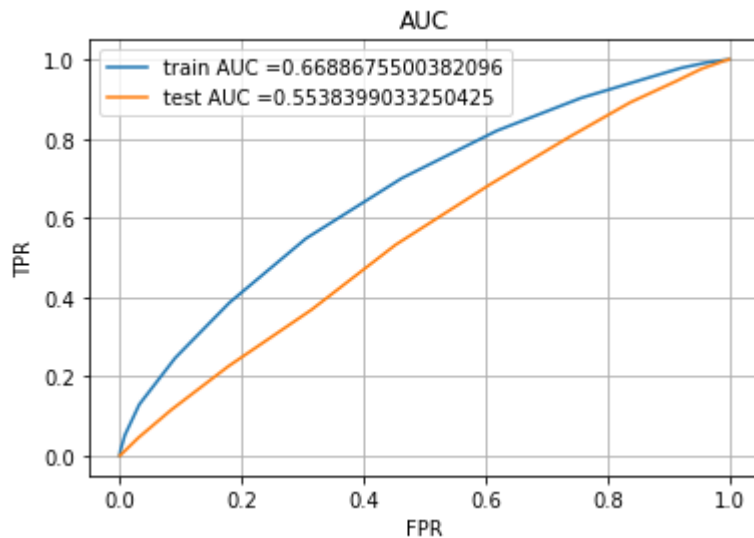
start_time = time.time()
neigh = KNeighborsClassifier(n_neighbors=49)
neigh.fit(X_tr[:, :], y_train[:])
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the pos
# not the predicted outputs

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_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("FPR")
plt.ylabel("TPR")
plt.title("AUC")
plt.grid()
plt.show()
print("Execution time: " + str((time.time() - start_time)) + ' ms')
```





Execution time: 3397.581614255905 ms

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, 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", np.round(t,3)
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions

print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

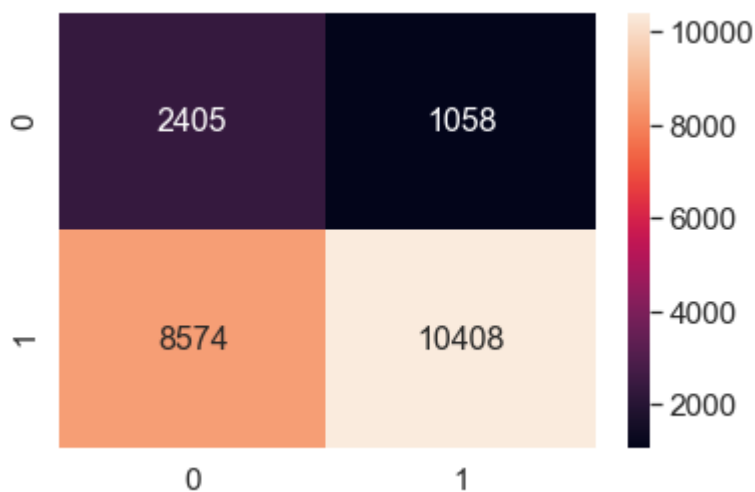


```
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.38079207704501006 for threshold 0.878
[[ 2405  1058]
 [ 8574 10408]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.2914590539378573 for threshold 0.878
[[1394 1152]
 [6526 7428]]
```

```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thre
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```



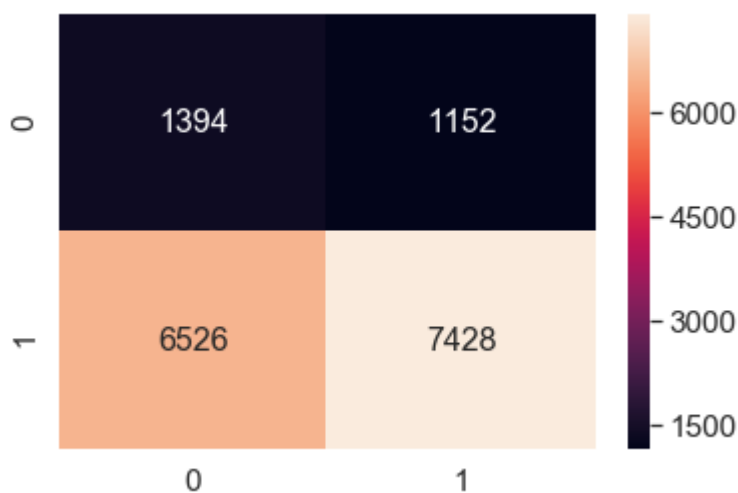
the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.38079207704501006 for threshold 0.878  
 <matplotlib.axes.\_subplots.AxesSubplot at 0x1f589e16a20>



```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresho
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```



the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.2914590539378573 for threshold 0.878  
 <matplotlib.axes.\_subplots.AxesSubplot at 0x1f58c005518>



### 3. Conclusions

# Please compare all your models using Prettytable library

```
from prettytable import PrettyTable
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable
```

```
x = PrettyTable()  
x.field_names = ["Vectorizer", "Model", "Hyper Parameter", "AUC"]  
x.add_row(["BOW", "Brute", 51, 0.62])  
x.add_row(["TFIDF", "Brute", 55, 0.60])  
x.add_row(["AVG W2V", "Brute", 53, 0.60])  
x.add_row(["TFIDF W2V", "Brute", 50, 0.61])  
x.add_row(["TFIDF", "Top 2000", 49, 0.55])  
print(x)
```



Vectorizer	Model	Hyper Parameter	AUC
BOW	Brute	51	0.62
TFIDF	Brute	55	0.6
AVG W2V	Brute	53	0.6
TFIDF W2V	Brute	50	0.61
TFIDF	Top 2000	49	0.55