

# DonorsChoose

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

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

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

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

## About the DonorsChoose Data Set

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

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. <b>Example:</b> p036502
<code>project_title</code>	<ul style="list-style-type: none"><li>•</li><li>•</li></ul>	Title of the project. <b>Examples:</b> <code>Art Will Make You Happy!</code> <code>First Grade Fun</code>
<code>project_grade_category</code>	<ul style="list-style-type: none"><li>•</li><li>•</li><li>•</li><li>•</li></ul>	Grade level of students for which the project is targeted. One of the following enumerated values: <code>Grades PreK-2</code> <code>Grades 3-5</code> <code>Grades 6-8</code> <code>Grades 9-12</code>
<code>project_subject_categories</code>	<ul style="list-style-type: none"><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li></ul>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <code>Applied Learning</code> <code>Care &amp; Hunger</code> <code>Health &amp; Sports</code> <code>History &amp; Civics</code> <code>Literacy &amp; Language</code> <code>Math &amp; Science</code> <code>Music &amp; The Arts</code> <code>Special Needs</code> <code>Warmth</code>  <b>Examples:</b> <ul style="list-style-type: none"><li>• <code>Music &amp; The Arts</code></li><li>• <code>Literacy &amp; Language, Math &amp; Science</code></li></ul>
<code>school_state</code>		State where school is located ( <a href="#">Two-letter U.S. postal code</a> ). <b>Example:</b> WY
<code>project_subject_subcategories</code>	<ul style="list-style-type: none"><li>•</li><li>•</li></ul>	One or more (comma-separated) subject subcategories for the project. <b>Examples:</b> <code>Literacy</code> <code>Literature &amp; Writing, Social Sciences</code>
<code>project_resource_summary</code>	<ul style="list-style-type: none"><li>•</li></ul>	An explanation of the resources needed for the project. <b>Example:</b> <code>My students need hands on literacy materials to manage sensory needs!</code>
<code>project_essay_1</code>		First application essay*
<code>project_essay_2</code>		Second application essay*
<code>project_essay_3</code>		Third application essay*

Feature	Description
project_essay_4	Fourth application essay
project_submitted_datetime	Datetime when project application was submitted. <b>Example:</b> 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <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>
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. <b>Example:</b> 2

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

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

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

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

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

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

## Notes on the Essay Data

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

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

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

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

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

In [1]:

```
%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

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

In [2]:

```

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

```

In [3]:

```
project_data.shape
```

Out[3]:

```
(109248, 17)
```

In [4]:

```

project_data=project_data.sample(n=10000)
project_data.shape

```

Out[4]:

```
(10000, 17)
```

In [5]:

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

Out[5]:

```

1      8471
0      1529
Name: project_is_approved, dtype: int64

```

In [6]:

```
resource_data.shape
```

Out[6]:

```
(1541272, 4)
```

In [7]:

```
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 (10000, 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']

In [8]:

```
# 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)
```

Out[8]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
64637	146737	p224791	ff5d658932d9ad0d9ebedabea582648e	Mrs.	MI	2016-04-27 08:51:57	Grades PreK-2	Applied
40180	86786	p028642	6b518fbb85fcb1587c8ee36e31b990e	Ms.	OK	2016-04-27 09:36:31	Grades PreK-2	

In [9]:

```
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']

Out[9]:

	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

In [10]:

```
print(project_data['project_subject_categories'].head(5))
```

```

64637      Applied Learning, Health & Sports
40180      Music & The Arts
92524      Literacy & Language, Math & Science
5364      Applied Learning, Music & The Arts
96705      Literacy & Language, Special Needs
Name: project_subject_categories, dtype: object

```

In [11]:

```

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

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
cat_list = []
for i in 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", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science" => "Math", "&", "Science"
            j = j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
            j = j.replace(' ', '') # we are placing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            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]))

```

In [12]:

```
print(project_data['clean_categories'].head(5))
```

```

64637      AppliedLearning Health_Sports
40180      Music_Arts
92524      Literacy_Language Math_Science
5364      AppliedLearning Music_Arts
96705      Literacy_Language SpecialNeeds
Name: clean_categories, dtype: object

```

## 1.3 preprocessing of project\_subject\_subcategories

In [13]:

```
print(project_data['project_subject_subcategories'].head(5))
```

```

64637      Early Development, Gym & Fitness
40180      Music
92524      ESL, Environmental Science
5364      College & Career Prep, Visual Arts
96705      Literature & Writing, Special Needs
Name: project_subject_subcategories, dtype: object

```

In [14]:

```

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

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

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", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science" => "Math", "&", "Science"
            j = j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
            j = j.replace(' ', '') # we are placing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " #"
        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]))

```

In [15]:

```
print(project_data['clean_subcategories'].head(5))
```

```

64637      EarlyDevelopment Gym_Fitness
40180                                     Music
92524      ESL EnvironmentalScience
5364      College_CareerPrep VisualArts
96705      Literature_Writing SpecialNeeds
Name: clean_subcategories, dtype: object

```

## 1.4 preprocessing of school\_state

In [16]:

```

my_counter = Counter()
for word in project_data['school_state'].values:
    my_counter.update(word.split())

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

```

## 1.5 preprocessing of project\_grade\_category

In [17]:

```

preproc = []
# tqdm is for printing the status bar
for sent in project_data['project_grade_category']:
    sent = sent.replace('Grades ', '')
    sent = sent.replace('PreK-2', 'PreKto2')
    sent = sent.replace('3-5', '3to5')
    sent = sent.replace('6-8', '6to8')
    sent = sent.replace('9-12', '9to12')

```

```
preproc.append(sent)
project_data['project_grade_category']=preproc
```

In [18]:

```
my_counter = Counter()
for word in project_data['project_grade_category'].values:
    my_counter.update(word.split())

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

## 1.6 preprocessing of teacher\_prefix

In [19]:

```
project_data['teacher_prefix'] = project_data['teacher_prefix'].astype(str)
preproc = []
# tqdm is for printing the status bar
for sent in project_data['teacher_prefix']:
    sent = sent.replace('Mr.', 'Mr')
    sent = sent.replace('Mrs.', 'Mrs')
    sent = sent.replace('Dr.', 'Dr')
    sent = sent.replace('Ms.', 'Ms')
    sent = sent.replace('nan', '')
    preproc.append(sent)
project_data['teacher_prefix']=preproc
```

In [20]:

```
#[ 'Teacher', 'Mrs.', 'Dr.', 'Mr.', 'Ms.' ]
project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('')
my_counter = Counter()
for word in project_data['teacher_prefix'].values:
    my_counter.update(word.split())

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

## 1.3 Text preprocessing

In [21]:

```
# 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)
```

In [22]:

```
# 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
```

In [23]:

```
# 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", "you've",
\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',
'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them',
'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll",
'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having',
'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', '
while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during',
'before', 'after',\
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under'
, 'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'e
ach', 'few', 'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll'
, 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "dc
esn't", 'hadn',\
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn',
"mightn't", 'mustn',\
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn',
"wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [24]:

```
# Combining all the above students
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%|██████████| 10000/10000 [00:06<00:00, 1612.27it/s]

In [25]:

```
# after preprocesing
preprocessed_essays[2000]
```

Out[25]:

'high school longest running public charter high school city detroit michigan city declaring bankr  
uptcy regular public city needing emergency financial manager city failing students commute across  
city chance quality education 75 students school receive free reduced lunch regardless may hear me  
dia students amazing day students striving not knock walls change negative stereotypes associated  
students live making difference outside classroom student internships volunteer community give bac  
k neighborhoods daily clothing drives collect books read children collect bottle tops ronald  
mcdonald house name things students bright futures many school place receive encouragement fulfill  
dreams k 6 education calculators not required strong foundation skills embedded education higher l  
evel mathematics students need calculators explore higher level used compare graphs calculate larg  
e values much teaching students math related real world ultimate goal calculators applications stu  
dents learn stocks interest personal finance velocities even self guided explorations inner city s  
tudents deserve chance compete nationally internationally lack resources not affect futures  
standardized tests make break many students statistics show students know navigate calculators per  
form better standardized exams class set calculators diminished years mitch matched calculators co



roded calculators stopped working reasons unknown want student opportunities students first year teaching advanced placement calculus give students opportunity receive college credit classes high school passing exam previous year students not perform well exam unfortunately not working class s et calculators attended professional development summer strongly encouraged ask department get wor king calculators class unfortunately told not budget strongly believe allowing students proper sup plies class allow compete nationally increase test scores need help nannan'

In [26]:

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

## 1.4 Preprocessing of `project\_title`

In [27]:

```
# similarly you can preprocess the titles also
preprocessed_titles = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_title'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\\"', ' ')
    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 not in stopwords)
    preprocessed_titles.append(sent.lower().strip())
```

100%|██████████| 10000/10000 [00:00<00:00, 38351.43it/s]

In [28]:

```
project_data['project_title']=preprocessed_titles
```

## 1.5 Preparing data for models

In [29]:

```
project_data.columns
```

Out[29]:

```
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 (optional)
- quantity : numerical (optional)
- teacher\_number\_of\_previously\_posted\_projects : numerical
- price : numerical

## 1.5.1 Vectorizing Categorical data

- <https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>

## 1.5.2 Vectorizing Text data

In [30]:

```
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/
# make sure you have the glove_vectors file
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f,encoding = "ISO-8859-1")
    glove_words = set(model.keys())
```

# 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 each hyper parameter, as shown in the figure
- Once you find the best hyper parameter, you need to train your model-M using the best hyper-param. Now, find the AUC on test data 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 original labels of test data points

### 4. [Task-2]

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

```
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 format. To print out a table please refer to 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 you train data, and apply the method `transform()` on cv/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 Sampling

In [31]:

```
y = project_data['project_is_approved']
print(y.shape)
```

(10000,)

In [32]:

```
project_data.drop(['project_is_approved'],axis=1,inplace=True)
```

In [33]:

```
X=project_data
print(X.shape)
```

(10000, 17)

In [34]:

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

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

In [35]:

```
print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)
print("="*100)
```

(4489, 17) (4489,)  
(2211, 17) (2211,)  
(3300, 17) (3300,)



#### Encoding of Text Data

In [36]:

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

**BOW of Essay**

In [37]:

```
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
```

In [38]:

```
vectorizer.fit(X_train['essay'].values) # fit has to happen only on train data
```

Out[38]:

```
CountVectorizer(analyzer='word', binary=False, decode_error='strict',
dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
lowercase=True, max_df=1.0, max_features=5000, min_df=10,
ngram_range=(1, 4), preprocessor=None, stop_words=None,
strip_accents=None, token_pattern='(?u)\\b\\w\\w+\\b',
tokenizer=None, vocabulary=None)
```

In [39]:

```
# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vectorizer.transform(X_train['essay'].values)
```

In [40]:

```
X_cv_essay_bow = vectorizer.transform(X_cv['essay'].values)
```

In [41]:

```
X_test_essay_bow = vectorizer.transform(X_test['essay'].values)
```

In [42]:

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

```
After vectorizations
(4489, 5000) (4489,)
(2211, 5000) (2211,)
(3300, 5000) (3300,)
```



## BOW of Title

In [43]:

```
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
```

In [44]:

```
vectorizer.fit(X_train['project_title'].values) # fit has to happen only on train data
```

Out[44]:

```
CountVectorizer(analyzer='word', binary=False, decode_error='strict',
dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
lowercase=True, max_df=1.0, max_features=5000, min_df=10,
ngram_range=(1, 4), preprocessor=None, stop_words=None,
strip_accents=None, token_pattern='(?u)\\b\\w\\w+\\b',
tokenizer=None, vocabulary=None)
```

In [45]:

```
# we use the fitted CountVectorizer to convert the text to vector
```

```
# we use the fitted CountVectorizer to convert the text to vector
X_train_title_bow = vectorizer.transform(X_train['project_title'].values)
```

In [46]:

```
X_cv_title_bow = vectorizer.transform(X_cv['project_title'].values)
```

In [47]:

```
X_test_title_bow = vectorizer.transform(X_test['project_title'].values)
```

In [48]:

```
print("After vectorizations")
print(X_train_title_bow.shape, y_train.shape)
print(X_cv_title_bow.shape, y_cv.shape)
print(X_test_title_bow.shape, y_test.shape)
print("=="*100)
```

```
After vectorizations
(4489, 430) (4489,)
(2211, 430) (2211,)
(3300, 430) (3300,)
=====
```



## TFIDF of Essay

In [49]:

```
vectorizer = TfidfVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
```

In [50]:

```
vectorizer.fit(X_train['essay'].values) # fit has to happen only on train data
```

Out[50]:

```
TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
               dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
               lowercase=True, max_df=1.0, max_features=5000, min_df=10,
               ngram_range=(1, 4), norm='l2', preprocessor=None, smooth_idf=True,
               stop_words=None, strip_accents=None, sublinear_tf=False,
               token_pattern='(?u)\\b\\w\\w+\\b', tokenizer=None, use_idf=True,
               vocabulary=None)
```

In [51]:

```
# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tfidf = vectorizer.transform(X_train['essay'].values)
```

In [52]:

```
X_cv_essay_tfidf = vectorizer.transform(X_cv['essay'].values)
```

In [53]:

```
X_test_essay_tfidf = vectorizer.transform(X_test['essay'].values)
```

In [54]:

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

```
After vectorizations
(4489, 5000) (4489,)
(2211, 5000) (2211,)
(3300, 5000) (3300,)
=====
```

## TFIDF of Title

In [55]:

```
vectorizer = TfidfVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
```

In [56]:

```
vectorizer.fit(X_train['project_title'].values) # fit has to happen only on train data
```

Out[56]:

```
TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
lowercase=True, max_df=1.0, max_features=5000, min_df=10,
ngram_range=(1, 4), norm='l2', preprocessor=None, smooth_idf=True,
stop_words=None, strip_accents=None, sublinear_tf=False,
token_pattern='(?u)\\b\\w\\w+\\b', tokenizer=None, use_idf=True,
vocabulary=None)
```

In [57]:

```
# we use the fitted CountVectorizer to convert the text to vector
X_train_title_tfidf = vectorizer.transform(X_train['project_title'].values)
```

In [58]:

```
X_cv_title_tfidf = vectorizer.transform(X_cv['project_title'].values)
```

In [59]:

```
X_test_title_tfidf = vectorizer.transform(X_test['project_title'].values)
```

In [60]:

```
print("After vectorizations")
print(X_train_title_tfidf.shape, y_train.shape)
print(X_cv_title_tfidf.shape, y_cv.shape)
print(X_test_title_tfidf.shape, y_test.shape)
print("=*100)
```

```
After vectorizations
(4489, 430) (4489,)
(2211, 430) (2211,)
(3300, 430) (3300,)
=====
```

## Avg W2V of Essay

In [61]:

```
# average Word2Vec
# compute average word2vec for each essay.
avg_w2v_essay_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['essay'].values): # 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
```

```

    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_essay_train.append(vector)

print(len(avg_w2v_essay_train))
print(len(avg_w2v_essay_train[0]))
print(type(avg_w2v_essay_train))

```

100%|██████████| 4489/4489 [00:01<00:00, 3507.86it/s]

```

4489
300
<class 'list'>

```

In [62]:

```

# average Word2Vec
# compute average word2vec for each essay.
avg_w2v_essay_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['essay'].values): # 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_essay_test.append(vector)

print(len(avg_w2v_essay_test))
print(len(avg_w2v_essay_test[0]))
print(type(avg_w2v_essay_test))

```

100%|██████████| 3300/3300 [00:00<00:00, 3419.41it/s]

```

3300
300
<class 'list'>

```

In [63]:

```

# average Word2Vec
# compute average word2vec for each essay.
avg_w2v_essay_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['essay'].values): # 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_essay_cv.append(vector)

print(len(avg_w2v_essay_cv))
print(len(avg_w2v_essay_cv[0]))
print(type(avg_w2v_essay_cv))

```

100%|██████████| 2211/2211 [00:00<00:00, 3423.25it/s]

```

2211
300

```

```
<class 'list'>
```

## Avg W2V of Title

In [64]:

```
# average Word2Vec
# compute average word2vec for each essay.
avg_w2v_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['project_title'].values): # 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_title_train.append(vector)

print(len(avg_w2v_title_train))
print(len(avg_w2v_title_train[0]))
print(type(avg_w2v_title_train))
```

```
100%|██████████| 4489/4489 [00:00<00:00, 68459.32it/s]
```

```
4489
```

```
300
```

```
<class 'list'>
```

In [65]:

```
# average Word2Vec
# compute average word2vec for each essay.
avg_w2v_title_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['project_title'].values): # 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_title_test.append(vector)

print(len(avg_w2v_title_test))
print(len(avg_w2v_title_test[0]))
print(type(avg_w2v_title_test))
```

```
100%|██████████| 3300/3300 [00:00<00:00, 69169.01it/s]
```

```
3300
```

```
300
```

```
<class 'list'>
```

In [66]:

```
# average Word2Vec
# compute average word2vec for each essay.
avg_w2v_title_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['project_title'].values): # 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:
```



```

    if word in glove_words:
        vector += model[word]
        cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_title_cv.append(vector)

print(len(avg_w2v_title_cv))
print(len(avg_w2v_title_cv[0]))
print(type(avg_w2v_title_cv))

```

100%|██████████| 2211/2211 [00:00<00:00, 65659.89it/s]

```

2211
300
<class 'list'>

```

## TFIDF-W2V of Essay

In [67]:

```

tfidf_model = TfidfVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
tfidf_model.fit(X_train['essay'].values)
# 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())

```

In [68]:

```

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_train_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_train_essay.append(vector)

print(len(tfidf_w2v_train_essay))
print(len(tfidf_w2v_train_essay[0]))

```

100%|██████████| 4489/4489 [00:07<00:00, 614.97it/s]

```

4489
300

```

In [69]:

```

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_test_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf

```

```

        # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split()))))
        tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
        vector += (vec * tf_idf) # calculating tfidf weighted w2v
        tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_test_essay.append(vector)

print(len(tfidf_w2v_test_essay))
print(len(tfidf_w2v_test_essay[0]))

```

100%|██████████| 3300/3300 [00:05<00:00, 617.26it/s]

3300

300

In [70]:

```

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_cv_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split()))))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_cv_essay.append(vector)

print(len(tfidf_w2v_cv_essay))
print(len(tfidf_w2v_cv_essay[0]))

```

100%|██████████| 2211/2211 [00:03<00:00, 646.46it/s]

2211

300

## TFIDF-W2V of Title

In [71]:

```

tfidf_model = TfidfVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
tfidf_model.fit(X_train['project_title'].values)
# 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())

```

In [72]:

```

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_train_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):

```

```

        vec = model[word] # getting the vector for each word
        # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
        tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
        vector += (vec * tf_idf) # calculating tfidf weighted w2v
        tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_train_title.append(vector)

print(len(tfidf_w2v_train_title))
print(len(tfidf_w2v_train_title[0]))

```

100%|██████████| 4489/4489 [00:00<00:00, 45792.73it/s]

4489  
300

In [73]:

```

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_test_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_test_title.append(vector)

print(len(tfidf_w2v_test_title))
print(len(tfidf_w2v_test_title[0]))

```

100%|██████████| 3300/3300 [00:00<00:00, 47855.68it/s]

3300  
300

In [74]:

```

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_cv_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_cv_title.append(vector)

```

```
print(len(tfidf_w2v_cv_title))
print(len(tfidf_w2v_cv_title[0]))
```

```
100%|██████████| 2211/2211 [00:00<00:00, 42944.10it/s]
```

```
2211
300
```

## 2.3 Make Data Model Ready: encoding numerical and categorical features

### Vectorizing Numerical features

In [75]:

```
price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
```

In [76]:

```
price_data.head(5)
```

Out[76]:

	id	quantity	price
0	p000001	7	459.56
1	p000002	21	515.89
2	p000003	4	298.97
3	p000004	98	1113.69
4	p000005	8	485.99

In [77]:

```
X_train=pd.merge(X_train,price_data,on='id',how='left')
X_test=pd.merge(X_test,price_data,on='id',how='left')
X_cv=pd.merge(X_cv,price_data,on='id',how='left')
```

In [78]:

```
X_train=X_train.fillna(0)
X_cv=X_cv.fillna(0)
X_test=X_test.fillna(0)
```

### Normalizing the numerical features: Price

In [79]:

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

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

print("After vectorizations")
print(X_train_price_norm.shape, y_train.shape)
print(X_cv_price_norm.shape, y_cv.shape)
```

```
print(X_test_price_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(4489, 1) (4489,)
(2211, 1) (2211,)
(3300, 1) (3300,)
```

=====

## Normalizing the numerical features: Number of previously posted projects

In [80]:

```
normalizer = Normalizer()
normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
X_train_project_norm = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
X_cv_project_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
X_test_project_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

print("After vectorizations")
print(X_train_project_norm.shape, y_train.shape)
print(X_cv_project_norm.shape, y_cv.shape)
print(X_test_project_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(4489, 1) (4489,)
(2211, 1) (2211,)
(3300, 1) (3300,)
```

=====

## Vectorizing Categorical features

- school\_state : categorical data
- clean\_categories : categorical data
- clean\_subcategories : categorical data
- project\_grade\_category : categorical data
- teacher\_prefix : categorical data

### Vectorizing Categorical features: project grade category

In [81]:

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

In [82]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_grade_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['project_grade_category'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_grade_oh = vectorizer.transform(X_train['project_grade_category'].values)
X_cv_grade_oh = vectorizer.transform(X_cv['project_grade_category'].values)
X_test_grade_oh = vectorizer.transform(X_test['project_grade_category'].values)

print("After vectorizations")
print(X_train_grade_oh.shape, y_train.shape)
print(X_cv_grade_oh.shape, y_cv.shape)
print(X_test_grade_oh.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

```
After vectorizations
(4489, 4) (4489,)
(2211, 4) (2211,)
(3300, 4) (3300,)
['9to12', '3to5', '6to8', 'PreKto2']
=====
```

### Vectorizing Categorical features: teacher prefix

In [83]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_teacher_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['teacher_prefix'].values) # fit has to happen only on train data

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

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

```
After vectorizations
(4489, 5) (4489,)
(2211, 5) (2211,)
(3300, 5) (3300,)
['Ms', 'Dr', 'Mr', 'Mrs', 'Teacher']
=====
```

### Vectorizing Categorical features: school state

In [84]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_state_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['school_state'].values) # fit has to happen only on train data

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

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

```
After vectorizations
(4489, 51) (4489,)
(2211, 51) (2211,)
(3300, 51) (3300,)
['GA', 'MT', 'AZ', 'OR', 'OK', 'MS', 'MI', 'MA', 'FL', 'NV', 'DE', 'NJ', 'WI', 'IA', 'ID', 'AK', 'NH', 'CO', 'NM', 'NE', 'MD', 'WV', 'DC', 'NC', 'KY', 'TN', 'KS', 'AR', 'ME', 'UT', 'VA', 'RI', 'MN', 'WA', 'AL', 'VT', 'HI', 'OH', 'MO', 'IN', 'ND', 'SC', 'CT', 'WY', 'NY', 'CA', 'PA', 'SD', 'LA', 'TX', 'IL']
=====
```

## Vectorizing Categorical features: clean categories

In [85]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['clean_categories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_cat_ohe = vectorizer.transform(X_train['clean_categories'].values)
X_cv_cat_ohe = vectorizer.transform(X_cv['clean_categories'].values)
X_test_cat_ohe = vectorizer.transform(X_test['clean_categories'].values)

print("After vectorizations")
print(X_train_cat_ohe.shape, y_train.shape)
print(X_cv_cat_ohe.shape, y_cv.shape)
print(X_test_cat_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations

```
(4489, 9) (4489,)
(2211, 9) (2211,)
(3300, 9) (3300,)
['Care_Hunger', 'SpecialNeeds', 'Music_Arts', 'Literacy_Language', 'History_Civics',
'AppliedLearning', 'Warmth', 'Math_Science', 'Health_Sports']
=====
```

## Vectorizing Categorical features: clean subcategories

In [86]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['school_state'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_sub_ohe = vectorizer.transform(X_train['clean_subcategories'].values)
X_cv_sub_ohe = vectorizer.transform(X_cv['clean_subcategories'].values)
X_test_sub_ohe = vectorizer.transform(X_test['clean_subcategories'].values)

print("After vectorizations")
print(X_train_sub_ohe.shape, y_train.shape)
print(X_cv_sub_ohe.shape, y_cv.shape)
print(X_test_sub_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations

```
(4489, 30) (4489,)
(2211, 30) (2211,)
(3300, 30) (3300,)
['History_Geography', 'EarlyDevelopment', 'Extracurricular', 'EnvironmentalScience',
'Literature_Writing', 'SpecialNeeds', 'Economics', 'Civics_Government', 'Other',
'CharacterEducation', 'NutritionEducation', 'VisualArts', 'ParentInvolvement',
'College_CareerPrep', 'Care_Hunger', 'FinancialLiteracy', 'AppliedSciences', 'Warmth',
'SocialSciences', 'Literacy', 'TeamSports', 'PerformingArts', 'CommunityService', 'Gym_Fitness',
'Mathematics', 'Health_LifeScience', 'Music', 'ESL', 'ForeignLanguages', 'Health_Wellness']
=====
```

## Concatinating all the features

## 2.4 Appling KNN on different kind of featurization as mentioned in the instructions

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

## Applying KNN brute force on BOW, SET 1

In [87]:

```
# Please write all the code with proper documentation

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((X_train_essay_bow,X_train_title_bow, X_train_state_ohc, X_train_teacher_ohc,
X_train_grade_ohc,X_train_cat_ohc,X_train_sub_ohc, X_train_price_norm,X_train_project_norm)).tocsr()
X_cr = hstack((X_cv_essay_bow,X_cv_title_bow, X_cv_state_ohc, X_cv_teacher_ohc, X_cv_grade_ohc,X_cv
_cat_ohc,X_cv_sub_ohc, X_cv_price_norm,X_cv_project_norm)).tocsr()
X_te = hstack((X_test_essay_bow,X_test_title_bow, X_test_state_ohc, X_test_teacher_ohc, X_test_grad
e_ohc,X_test_cat_ohc,X_test_sub_ohc, X_test_price_norm,X_test_project_norm)).tocsr()

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

```
Final Data matrix
(4489, 5531) (4489,)
(2211, 5531) (2211,)
(3300, 5531) (3300,)
```

=====



### Hyperparameter Tuning: Simple for loop (if you are having memory limitations use this)

In [88]:

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi
    tive class
    # 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
```

In [89]:

```
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, or no
n-thresholded measure of
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, 9, 15, 51, 71, 101, 151, 251, 451, 651]

for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i,p=2)
    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 posi
    tive class
    # 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))

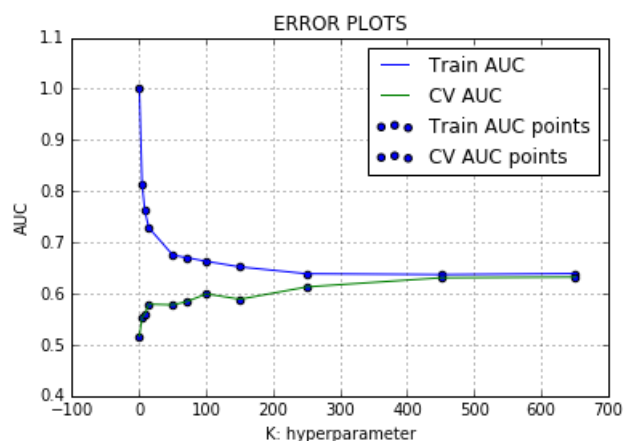
plt.plot(K, train_auc, label='Train AUC')
plt.plot(K, cv_auc, label='CV AUC')

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()

```

100%|██████████| 11/11 [00:38<00:00, 3.69s/it]



In [96]:

```
best_k=775
```

## Train The Model

In [97]:

```

from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k,p=2)
neigh.fit(X_tr, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

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

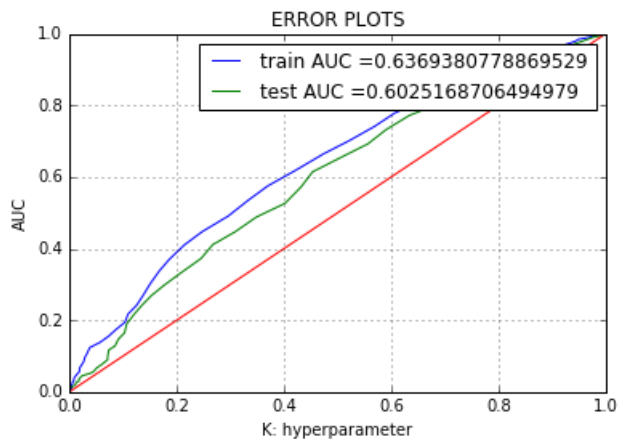
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

x=[0.0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0]

```

In [98]:

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



## Confusion Matrix

In [99]:

```
# 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(fpr*(1-tpr))]

    # (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
```

In [100]:

```
print("Train confusion matrix")
conf_matr_df_train_2=pd.DataFrame(confusion_matrix(y_train,predict(y_train_pred,tr_thresholds,train_fpr,train_fpr)),range(2),range(2))
sns.set(font_scale=1)#for label size
sns.heatmap(conf_matr_df_train_2,annot=True,annot_kws={"size":30},fmt='g')
```

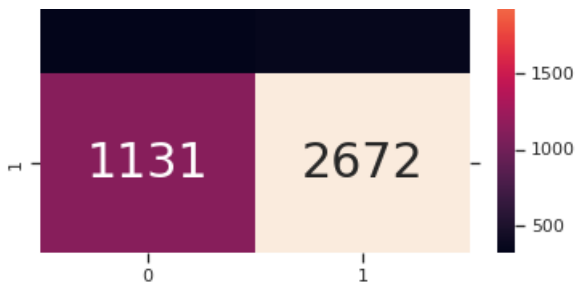
Train confusion matrix

the maximum value of  $tpr \cdot (1-fpr)$  0.2494560089758519 for threshold 0.814

Out[100]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb459b0c4e0>





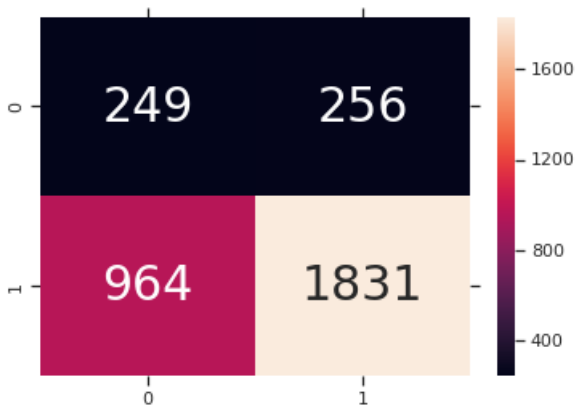
In [101]:

```
print("Test confusion matrix")
conf_matr_df_train_2=pd.DataFrame(confusion_matrix(y_test,predict(y_test_pred,tr_thresholds,test_fpr,
test_fpr)),range(2),range(2))
sns.set(font_scale=1)#for label size
sns.heatmap(conf_matr_df_train_2,annot=True,annot_kws={"size":30},fmt='g')
```

Test confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24995196549357906 for threshold 0.815

Out[101]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb459b64f28>



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

In [102]:

```
# Please write all the code with proper documentation

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((X_train_essay_tfidf,X_train_title_tfidf, X_train_state_ohe, X_train_teacher_ohe, X_train_grade_ohe,X_train_cat_ohe,X_train_sub_ohe, X_train_price_norm,X_train_project_norm)).tocsr()
X_cr = hstack((X_cv_essay_tfidf,X_cv_title_tfidf, X_cv_state_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_cat_ohe,X_cv_sub_ohe, X_cv_price_norm,X_cv_project_norm)).tocsr()
X_te = hstack((X_test_essay_tfidf,X_test_title_tfidf, X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe,X_test_cat_ohe,X_test_sub_ohe, X_test_price_norm,X_test_project_norm)).tocsr()

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

Final Data matrix  
(4489, 5531) (4489,)  
(2211, 5531) (2211,)  
(3300, 5531) (3300,)

## Hyperparameter Tuning: Simple for loop (if you are having memory limitations use this)

In [103]:

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # 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
```

In [125]:

```
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, or no
n-thresholded measure of
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, 9, 15, 51, 71, 101, 151, 251, 451, 651,851,1051]
for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i,p=2)
    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 positive class
    # 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))
```

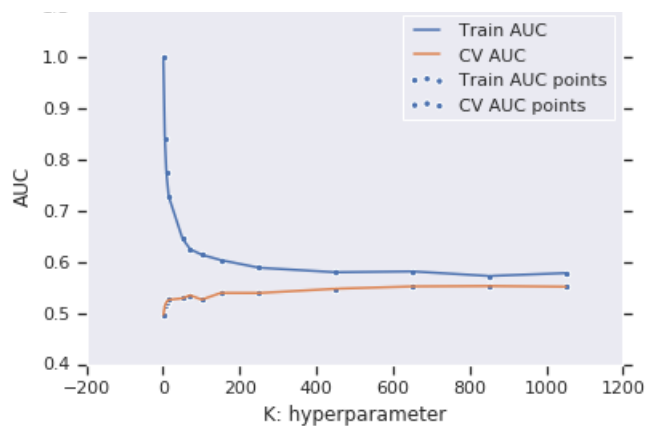
100%|██████████| 13/13 [00:50<00:00, 4.32s/it]

In [126]:

```
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()
```



In [130]:

```
best_k=851
```

## Train The Model

In [131]:

```
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k,p=2)
neigh.fit(X_tr, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# 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)
```

In [132]:

```
x=[0.0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0]
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.plot(x,x)
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



## Confusion Matrix

In [133]:

```
# 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(fpr*(1-tpr))]

    # (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
```

In [134]:

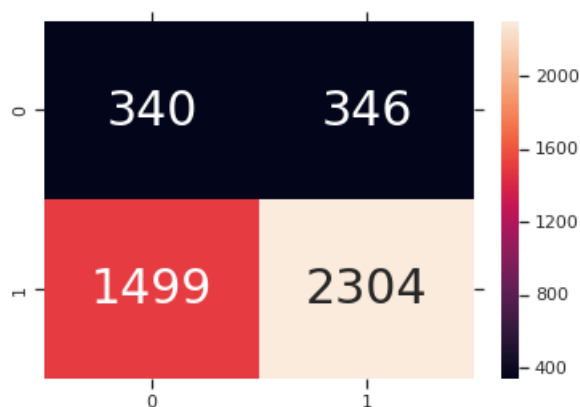
```
print("Train confusion matrix")
conf_matr_df_train_2=pd.DataFrame(confusion_matrix(y_train,predict(y_train_pred,tr_thresholds,train_fpr,train_fpr)),range(2),range(2))
sns.set(font_scale=1)#for label size
sns.heatmap(conf_matr_df_train_2,annot=True,annot_kws={"size":30},fmt='g')
```

Train confusion matrix

the maximum value of tpr\*(1-fpr) 0.2499808753155573 for threshold 0.854

Out[134]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb45ab122b0>



In [135]:

```
print("Test confusion matrix")
conf_matr_df_train_2=pd.DataFrame(confusion_matrix(y_test,predict(y_test_pred,tr_thresholds,test_fpr,test_fpr)),range(2),range(2))
sns.set(font_scale=1)#for label size
sns.heatmap(conf_matr_df_train_2,annot=True,annot_kws={"size":30},fmt='g')
```

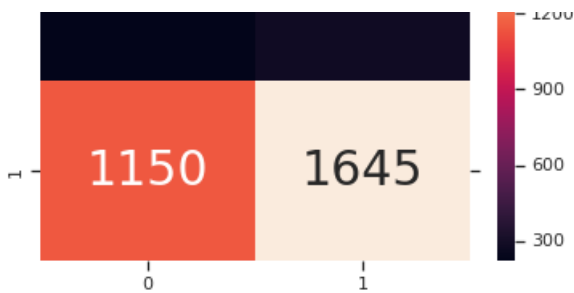
Test confusion matrix

the maximum value of tpr\*(1-fpr) 0.24997549259876484 for threshold 0.854

Out[135]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb45aad7710>





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

In [136]:

```
# Please write all the code with proper documentation

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((avg_w2v_essay_train, avg_w2v_title_train, X_train_state_ohe, X_train_teacher_ohe, X_train_grade_ohe, X_train_cat_ohe, X_train_sub_ohe, X_train_price_norm, X_train_project_norm)).tocsr()
X_cr = hstack((avg_w2v_essay_cv, avg_w2v_title_cv, X_cv_state_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_cat_ohe, X_cv_sub_ohe, X_cv_price_norm, X_cv_project_norm)).tocsr()
X_te = hstack((avg_w2v_essay_test, avg_w2v_title_test, X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe, X_test_cat_ohe, X_test_sub_ohe, X_test_price_norm, X_test_project_norm)).tocsr()

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

```
Final Data matrix
(4489, 701) (4489,)
(2211, 701) (2211,)
(3300, 701) (3300,)
```

### Hyperparameter Tuning: Simple for loop (if you are having memory limitations use this)

In [137]:

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # 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
```

In [138]:

```
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, or no
```

*n*-thresholded measure of 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, 9, 15, 51, 71, 101, 151, 251, 451, 651]
for i in tqdm(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)

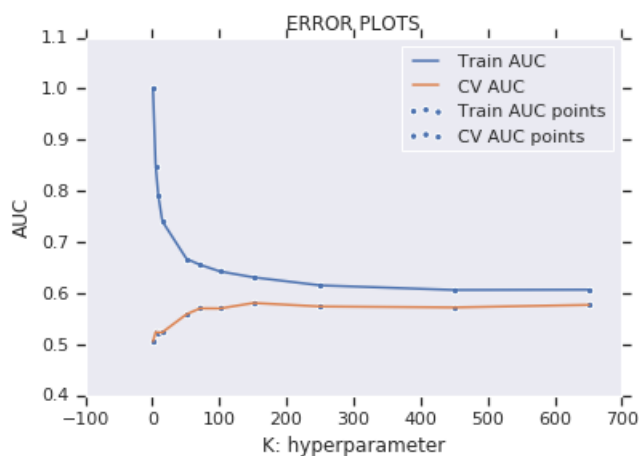
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

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()
```

100%|██████████| 11/11 [10:39<00:00, 58.93s/it]



In [142]:

```
best_k=351
```

## Train The Model

In [143]:

```
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k,p=2)
neigh.fit(X_tr, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# 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)

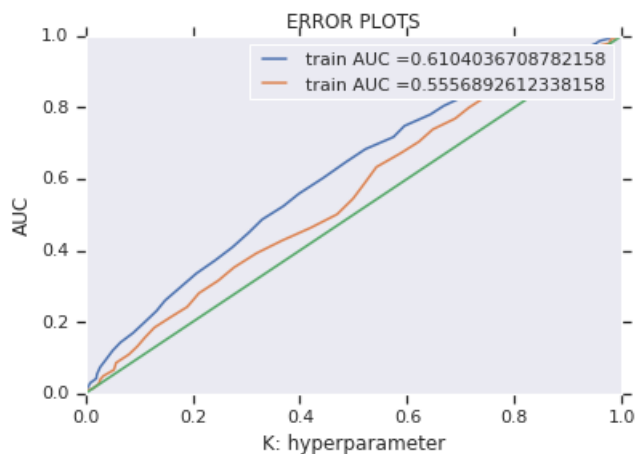
```

In [144]:

```

x=[0.0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0]
plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="train AUC =" +str(auc(test_fpr, test_tpr)))
plt.plot(x,x)
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



## Confusion Matrix

In [145]:

```

# 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(fpr*(1-tpr))]

    # (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

```

In [146]:

```

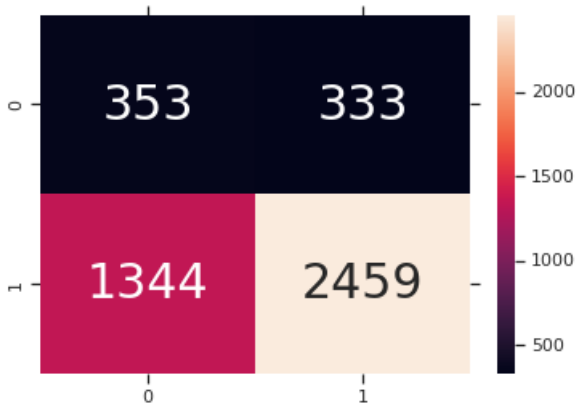
print("Train confusion matrix")
conf_matr_df_train_2=pd.DataFrame(confusion_matrix(y_train,predict(y_train_pred,tr_thresholds,train_fpr,train_tpr)),range(2),range(2))
sns.set(font_scale=1)#for label size
sns.heatmap(conf_matr_df_train_2,annot=True,annot_kws={"size":30},fmt='g')

```

Train confusion matrix  
the maximum value of tpr\*(1-fpr) 0.24978750350619217 for threshold 0.855

Out[146]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb459bbca90>



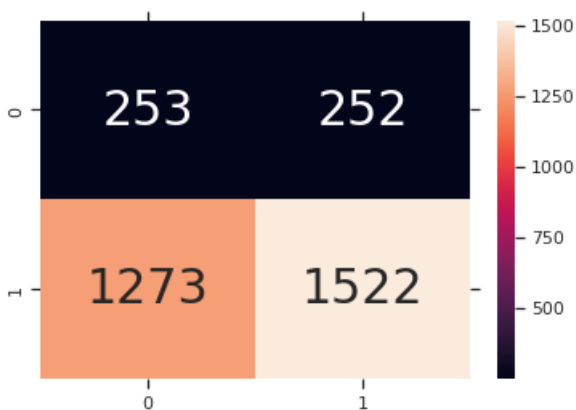
In [147]:

```
print("Test confusion matrix")
conf_matr_df_train_2=pd.DataFrame(confusion_matrix(y_test,predict(y_test_pred,tr_thresholds,test_fpr,test_fpr)),range(2),range(2))
sns.set(font_scale=1)#for label size
sns.heatmap(conf_matr_df_train_2,annot=True,annot_kws={"size":30},fmt='g')
```

Test confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24999901970395058 for threshold 0.86

Out[147]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb45ab71b00>



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

In [148]:

```
# Please write all the code with proper documentation

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((tfidf_w2v_train_essay,tfidf_w2v_train_title, X_train_state_ohe, X_train_teacher_ohe,
X_train_grade_ohe,X_train_cat_ohe,X_train_sub_ohe, X_train_price_norm,X_train_project_norm)).tocsr()
X_cr = hstack((tfidf_w2v_cv_essay,tfidf_w2v_cv_title, X_cv_state_ohe, X_cv_teacher_ohe, X_cv_grade_ohe,
X_cv_cat_ohe,X_cv_sub_ohe, X_cv_price_norm,X_cv_project_norm)).tocsr()
X_te = hstack((tfidf_w2v_test_essay,tfidf_w2v_test_title, X_test_state_ohe, X_test_teacher_ohe,
X_test_grade_ohe,X_test_cat_ohe,X_test_sub_ohe, X_test_price_norm,X_test_project_norm)).tocsr()

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

Final Data matrix

```
=====  
(4489, 701) (4489,)  
(2211, 701) (2211,)  
(3300, 701) (3300,)  
=====
```

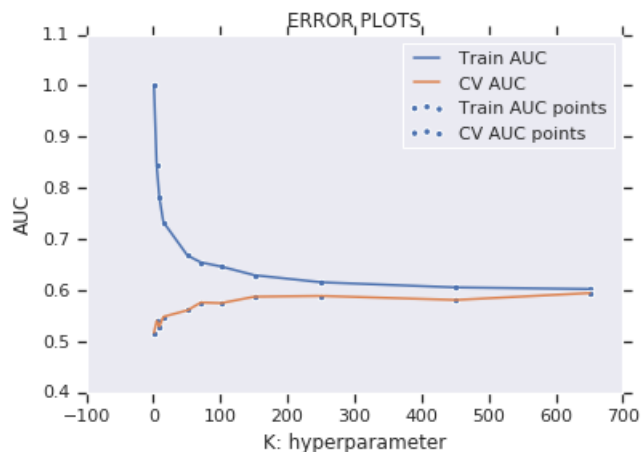
## Hyperparameter Tuning: Simple for loop (if you are having memory limitations use this)

In [149]:

```
def batch_predict(clf, data):  
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi  
    # tive class  
    # 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
```

In [150]:

```
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, or no  
n-thresholded measure of  
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, 9, 15, 51, 71, 101, 151, 251, 451, 651]  
for i in tqdm(K):  
    neigh = KNeighborsClassifier(n_neighbors=i,p=2)  
    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 posi  
    # tive class  
    # 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()
```



In [161]:

```
best_k=101
```

## Train The Model

In [162]:

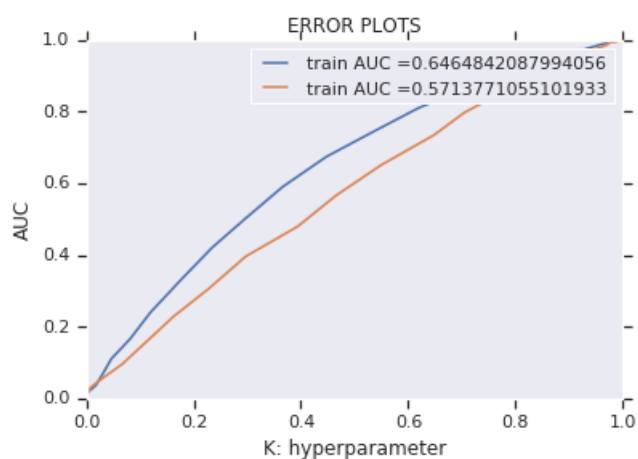
```
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k)
neigh.fit(X_tr, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# 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="train AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



## Confusion Matrix

In [163]:

```
# 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(fpr*(1-tpr))]

    # (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
```

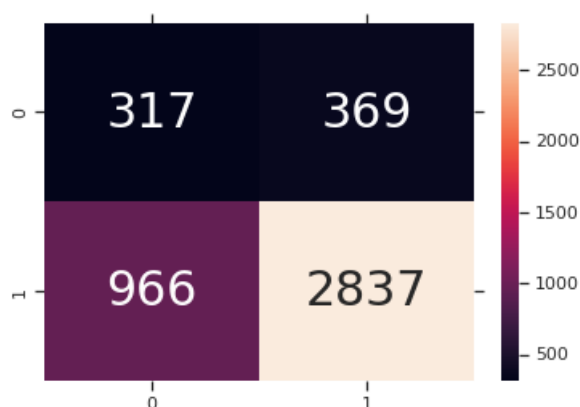
In [164]:

```
print("Train confusion matrix")
conf_matr_df_train_2=pd.DataFrame(confusion_matrix(y_train,predict(y_train_pred,tr_thresholds,train_fpr,train_fpr)),range(2),range(2))
sns.set(font_scale=1)#for label size
sns.heatmap(conf_matr_df_train_2,annot=True,annot_kws={"size":30},fmt='g')
```

Train confusion matrix  
the maximum value of tpr\*(1-fpr) 0.2485635237018589 for threshold 0.832

Out[164]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb45ab9ba58>



In [165]:

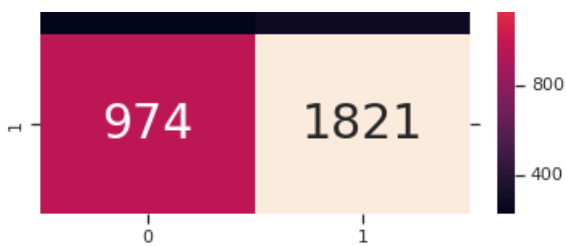
```
print("Test confusion matrix")
conf_matr_df_train_2=pd.DataFrame(confusion_matrix(y_test,predict(y_test_pred,tr_thresholds,test_fpr,test_fpr)),range(2),range(2))
sns.set(font_scale=1)#for label size
sns.heatmap(conf_matr_df_train_2,annot=True,annot_kws={"size":30},fmt='g')
```

Test confusion matrix  
the maximum value of tpr\*(1-fpr) 0.24879913733947656 for threshold 0.842

Out[165]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb45ab96ac8>





## 2.5 Feature selection with `SelectKBest`

In [166]:

```
# 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
```

In [167]:

```
X_tr = hstack((X_train_essay_tfidf,X_train_title_tfidf, X_train_state_ohe, X_train_teacher_ohe, X_train_grade_ohe,X_train_cat_ohe,X_train_sub_ohe, X_train_price_norm,X_train_project_norm)).tocsr()
X_cr = hstack((X_cv_essay_tfidf,X_cv_title_tfidf, X_cv_state_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_cat_ohe,X_cv_sub_ohe, X_cv_price_norm,X_cv_project_norm)).tocsr()
X_te = hstack((X_test_essay_tfidf,X_test_title_tfidf, X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe,X_test_cat_ohe,X_test_sub_ohe, X_test_price_norm,X_test_project_norm)).tocsr()

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

```
Final Data matrix
(4489, 5531) (4489,)
(2211, 5531) (2211,)
(3300, 5531) (3300,)
=====
```



In [170]:

```
from sklearn.feature_selection import SelectKBest,chi2,f_classif
best_feature=SelectKBest(score_func=f_classif,k=2000)
```

In [172]:

```
best_feature.fit(X_tr,y_train)
```

Out[172]:

```
SelectKBest(k=2000, score_func=<function f_classif at 0x7fb464f96a60>)
```

In [176]:

```
X_tr=best_feature.transform(X_tr)
X_te=best_feature.transform(X_te)
X_cr=best_feature.transform(X_cr)
```

In [178]:

```
print("Final Data matrix")
```

```

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

```

```

Final Data matrix
(4489, 2000) (4489,)
(2211, 2000) (2211,)
(3300, 2000) (3300,)
=====

```

## Hyperparameter Tuning: Simple for loop (if you are having memory limitations use this)

In [179]:

```

def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi
    # 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

```

In [180]:

```

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, or no
n-thresholded measure of
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 = [1,5,10,15,21,31,41,51,71,101,121,151]
for i in tqdm(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 posi
    # 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')

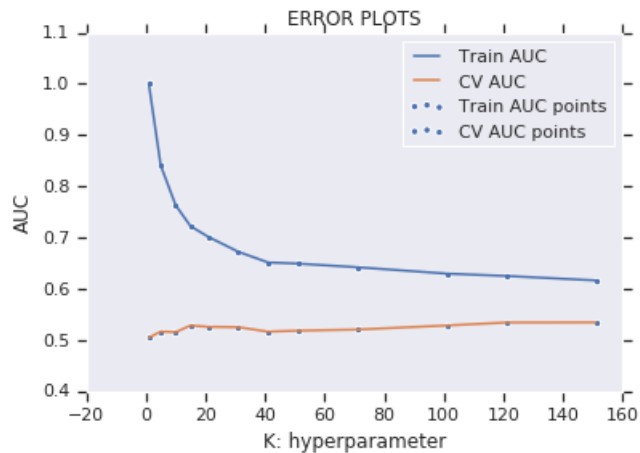
plt.scatter(K, train_auc, label='Train AUC points')
plt.scatter(K, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("K: hyperparameter")

```

```
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

100%|██████████| 12/12 [00:29<00:00, 2.50s/it]



In [185]:

```
best_k=151
```

## Train The Model

In [186]:

```
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k)
neigh.fit(X_tr, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# 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("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```





0.0 0.2 0.4 0.6 0.8 1.0  
K: hyperparameter

## Confusion Matrix

In [187]:

```
# 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(fpr*(1-tpr))]

    # (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
```

In [188]:

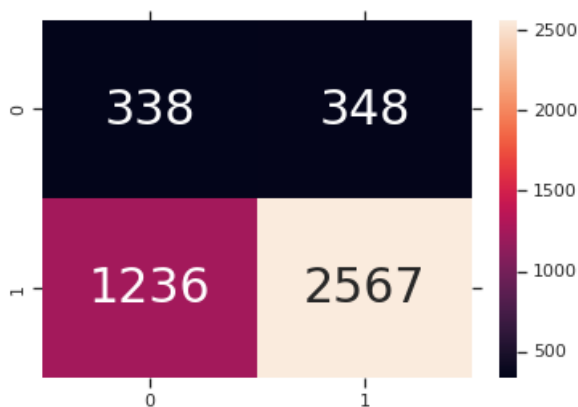
```
print("Train confusion matrix")
conf_matr_df_train_2=pd.DataFrame(confusion_matrix(y_train,predict(y_train_pred,tr_thresholds,train_fpr,train_fpr)),range(2),range(2))
sns.set(font_scale=1)#for label size
sns.heatmap(conf_matr_df_train_2,annot=True,annot_kws={"size":30},fmt='g')
```

Train confusion matrix

the maximum value of tpr\*(1-fpr) 0.24994687587654804 for threshold 0.854

Out[188]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb45aafe390>



In [189]:

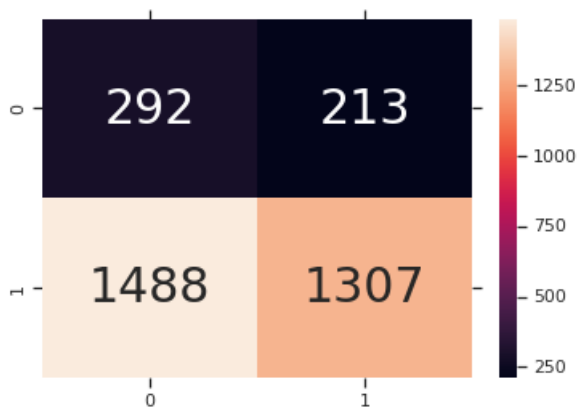
```
print("Test confusion matrix")
conf_matr_df_train_2=pd.DataFrame(confusion_matrix(y_test,predict(y_test_pred,tr_thresholds,test_fpr,test_fpr)),range(2),range(2))
sns.set(font_scale=1)#for label size
sns.heatmap(conf_matr_df_train_2,annot=True,annot_kws={"size":30},fmt='g')
```

Test confusion matrix

the maximum value of tpr\*(1-fpr) 0.24995196549357906 for threshold 0.874

Out[189]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb459bddc18>



### 3. Conclusions

In [190]:

```
# http://zetcode.com/python/prettitable/

from prettytable import PrettyTable
#If you get a ModuleNotFoundError error , install prettifytable using: pip3 install prettifytable
x=PrettyTable()
x.field_names=["Vectorizer","Model","Hyper Parameter","AUC"]
x.add_row(["BOW","Brute",120,0.60])
x.add_row(["TFIDF","Brute",130,0.53])
x.add_row(["AVG W2V","Brute",120,0.56])
x.add_row(["TFIDF W2V","Brute",110,0.57])
x.add_row(["TFIDF","Top 2000",100,0.53])
print(x)
```

Vectorizer	Model	Hyper Parameter	AUC
BOW	Brute	120	0.6
TFIDF	Brute	130	0.53
AVG W2V	Brute	120	0.56
TFIDF W2V	Brute	110	0.57
TFIDF	Top 2000	100	0.53

In [ ]: