

In [0]:

```
from google.colab import drive
drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

In [0]:

```
%cd drive/My Drive/Assignments_DonorsChoose_2018
```

/content/drive/My Drive/Assignments_DonorsChoose_2018

In [0]:

```
%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
import sklearn
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
```

In [0]:

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

In [0]:

```
# 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['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)
```

In [0]:

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

In [0]:

```
project_data = pd.merge(project_data, price_data, on='id', how='left')
```

1.2 preprocessing of project_subject_categories

In [0]:

```
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 & H
unger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Scienc
e"=> "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]))
```

1.3 preprocessing of project_subject_subcategories

In [0]:

```
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 replacing 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('&', '_')
    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

In [0]:

```
# 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 [0]:

```
# 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 [0]:

```
# 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',
'throughout', 'till', 'towards', 'against', 'across', 'among', 'beside', 'by', 'for', 'from', 'in', 'of', 'on', 'off', 'over', 'per', 'through', 'till', 'towards', 'under', 'until', 'up', 'via', 'within']
```

```
'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', 'each', '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', "doesn'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 [0]:

```
def find_num(text):
    if re.findall(r'\d+', text):
        return 1
    return 0
```

```
project_data['numerical_digits'] = project_data['project_resource_summary'].apply(lambda x: find_num(x))
```

In [0]:

```
project_data['project_grade_category']=project_data['project_grade_category'].str.replace(' ','_')
project_data['project_grade_category']=project_data['project_grade_category'].str.replace('-', 'to')
set(project_data['project_grade_category'])
```

Out[0]:

```
{'Grades_3to5', 'Grades_6to8', 'Grades_9to12', 'Grades_PreKto2'}
```

In [0]:

```
project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('Mrs.')
```

In [0]:

```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import accuracy_score
from collections import Counter
import matplotlib.pyplot as plt
import numpy as np
from scipy.sparse import csr_matrix
import time

project_data_features = project_data.copy()
project_data_features.drop('project_is_approved', axis=1, inplace=True)
y=list(project_data['project_is_approved'])
X_train, X_test, y_train, y_test = train_test_split(project_data_features, y, stratify=y, test_size=0.20)
X_train, X_val, y_train, y_val = train_test_split(X_train, y_train, stratify=y_train, test_size=0.30)
```

In [0]:

```
print("Train:",X_train.shape,len(y_train))
print("CV:",X_val.shape,len(y_val))
print("Test:",X_test.shape,len(y_test))
```

```
Train: (61178, 20) 61178
CV: (26220, 20) 26220
Test: (21850, 20) 21850
```

Function for standardizing/normalizing the data

In [0]:

```
from sklearn.preprocessing import StandardScaler
def standardize_data(df_tr,df_cv,df_te,column_name):
    standardized_vec = StandardScaler(with_mean=False)
    # here it will learn mu and sigma
    standardized_vec.fit(df_tr[column_name].values.reshape(-1,1))

    # with the learned mu and sigma it will do std on train data
    standardized_data_train = standardized_vec.transform(df_tr[column_name].values.reshape(-1,1))
    print(standardized_data_train.shape)

    # with the same learned mu and sigma it will do std on cv data
    standardized_data_traincv = standardized_vec.transform(df_cv[column_name].values.reshape(-1,1))
    print(standardized_data_traincv.shape)

    # with the same learned mu and sigma it will do std on test data
    standardized_data_test =standardized_vec.transform(df_te[column_name].values.reshape(-1,1))
    print(standardized_data_test.shape)

    return standardized_data_train, standardized_data_traincv, standardized_data_test
```

In [0]:

```
from sklearn.preprocessing import Normalizer
def normalize_data(df,column_data):
    normalizer_vec = Normalizer()

    normalizer_vec.fit(df_tr[column_name].values.reshape(-1,1))
    normalizer_data_train = normalizer_vec.transform(df_tr[column_name].values.reshape(-1,1))
    print(normalizer_data_train.shape)

    normalizer_data_traincv = normalizer_vec.transform(df_cv[column_name].values.reshape(-1,1))
    print(normalizer_data_traincv.shape)

    normalizer_data_test = normalizer_vec.transform(df_te[column_name].values.reshape(-1,1))
    print(normalizer_data_test.shape)

    return normalizer_data_train, normalizer_data_traincv, normalizer_data_test
```

Function for vectorizing the text data

In [0]:

```
from sklearn.feature_extraction.text import CountVectorizer
def vectorized_data(df_train,df_cv,df_test,column_name,vocab=False):
    if(vocab):
        vectorizer = CountVectorizer(vocabulary=list(vocab.keys()), lowercase=False, binary=True)
    else:
        vectorizer = CountVectorizer(lowercase=False, binary=True)
    categories_one_hot_tr = vectorizer.fit_transform(df_train[column_name].values)
    print(vectorizer.get_feature_names())
    print("Shape of matrix after one hot encoding ",categories_one_hot_tr.shape)

    categories_one_hot_cv = vectorizer.transform(df_cv[column_name].values)
    print(vectorizer.get_feature_names())
    print("Shape of matrix after one hot encoding ",categories_one_hot_cv.shape)

    categories_one_hot_te = vectorizer.transform(df_test[column_name].values)
    print(vectorizer.get_feature_names())
    print("Shape of matrix after one hot encoding ",categories_one_hot_te.shape)
    return categories_one_hot_tr,categories_one_hot_cv, categories_one_hot_te
```

In [0]:

```
def create_dict(df,column_name):
    my_counter = Counter()
    for word in df[column_name].values:
```

```

my_counter.update(word.split())

my_dict = dict(my_counter)
sorted_dict = dict(sorted(my_dict.items(), key=lambda kv: kv[1]))
return sorted_dict

```

In [0]:

```

def num_hot_encode(df,column_name):
    one_hot_num_dig = pd.get_dummies(df[column_name].values)
    print("Shape of matrix after one hot encoding ",one_hot_num_dig.shape)
    return one_hot_num_dig

```

In [0]:

```

from tqdm import tqdm
def textpreprocessed(df,column_name):
    # Combining all the above students
    preprocessed_list = []
    # tqdm is for printing the status bar
    for sentence in tqdm(df[column_name].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_list.append(sent.lower().strip())
    return preprocessed_list

```

In [0]:

```

price_standardized_tr, price_standardized_val, price_standardized_te =
standardize_data(X_train,X_val, X_test,'price')
print()
project_standardized_tr, project_standardized_val, project_standardized_te =
standardize_data(X_train,X_val, X_test,'teacher_number_of_previously_posted_projects')

```

```

(61178, 1)
(26220, 1)
(21850, 1)

```

```

(61178, 1)
(26220, 1)
(21850, 1)

```

In [0]:

```

cat_one_hot_tr,cat_one_hot_val,cat_one_hot_te = vectorized_data(X_train,X_val,X_test,'clean_categories',vocab=sorted_cat_dict)
cat_sub_one_hot_tr,cat_sub_one_hot_val,cat_sub_one_hot_te =
vectorized_data(X_train,X_val,X_test,'clean_subcategories',vocab=sorted_cat_dict)

```

```

['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding (61178, 9)
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding (26220, 9)
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding (21850, 9)
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding (61178, 9)
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding (26220, 9)
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding (21850, 9)

```

Shape of matrix after one hot encoding (61178, 5)

In [0]:

```
school_state_dict = create_dict(X_train,'school_state')
teacher_prefix_dict = create_dict(X_train,'teacher_prefix')
```

```
state_one_hot_tr,state_one_hot_val,state_one_hot_te =
vectorized_data(X_train,X_val,X_test,'school_state',vocab=school_state_dict)
teacher_one_hot_tr,teacher_one_hot_val,teacher_one_hot_te =
vectorized_data(X_train,X_val,X_test,'teacher_prefix',vocab=teacher_prefix_dict)
```

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

Shape of matrix after one hot encoding (61178, 51)

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

Shape of matrix after one hot encoding (26220, 51)

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

Shape of matrix after one hot encoding (21850, 51)

```
['Dr.', 'Teacher', 'Mr.', 'Ms.', 'Mrs.']
```

Shape of matrix after one hot encoding (61178, 5)

```
['Dr.', 'Teacher', 'Mr.', 'Ms.', 'Mrs.']
```

Shape of matrix after one hot encoding (26220, 5)

```
['Dr.', 'Teacher', 'Mr.', 'Ms.', 'Mrs.']
```

Shape of matrix after one hot encoding (21850, 5)

In [0]:

```
grade_dict = create_dict(X_train,'project_grade_category')
grade_one_hot_tr,grade_one_hot_val,grade_one_hot_te =
vectorized_data(X_train,X_val,X_test,'project_grade_category',vocab=grade_dict)
```

```
['Grades_9to12', 'Grades_6to8', 'Grades_3to5', 'Grades_PreKto2']
```

Shape of matrix after one hot encoding (61178, 4)

```
['Grades_9to12', 'Grades_6to8', 'Grades_3to5', 'Grades_PreKto2']
```

Shape of matrix after one hot encoding (26220, 4)

```
['Grades_9to12', 'Grades_6to8', 'Grades_3to5', 'Grades_PreKto2']
```

Shape of matrix after one hot encoding (21850, 4)

In [0]:

```
one_hot_num_dig_tr = num_hot_encode(X_train,'numerical_digits')
one_hot_num_dig_te = num_hot_encode(X_test,'numerical_digits')
one_hot_num_dig_val = num_hot_encode(X_val,'numerical_digits')
```

Shape of matrix after one hot encoding (61178, 2)

Shape of matrix after one hot encoding (21850, 2)

Shape of matrix after one hot encoding (26220, 2)

In [0]:

```
one_hot_num_dig_te.shape, grade_one_hot_te.shape \
,state_one_hot_te.shape,cat_one_hot_te.shape,\
cat_sub_one_hot_te.shape,teacher_one_hot_te.shape,\
price_standardized_te.shape,project_standardized_te.shape
```

Out[0]:

```
((21850, 2),
 (21850, 4),
 (21850, 51),
 (21850, 9),
 (21850, 9),
 (21850, 5),
```

```
(21850, 1),
(21850, 1))
```

In [0]:

```
from scipy.sparse import hstack

f_tr =
hstack((one_hot_num_dig_tr,grade_one_hot_tr,state_one_hot_tr,cat_one_hot_tr,cat_sub_one_hot_tr,teacher_one_hot_tr,price_standardized_tr,project_standardized_tr))
f_cr =
hstack((one_hot_num_dig_val,grade_one_hot_val,state_one_hot_val,cat_one_hot_val,cat_sub_one_hot_val,teacher_one_hot_val,price_standardized_val,project_standardized_val))
f_te =
hstack((one_hot_num_dig_te,grade_one_hot_te,state_one_hot_te,cat_one_hot_te,cat_sub_one_hot_te,teacher_one_hot_te,price_standardized_te,project_standardized_te))

def hstack_data(f1_tr, f1_cr, f1_te, f2_tr, f2_cr, f2_te, f3_tr, f3_cr, f3_te):
    X_tr = hstack((f1_tr, f1_cr, f2_tr, f3_tr)).tocsr()
    X_cr = hstack((f1_cr, f1_cr, f2_cr, f3_cr)).tocsr()
    X_te = hstack((f1_te, f1_te, f2_te, f3_te)).tocsr()
    return X_tr, X_cr, X_te
```

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

In [0]:

```
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]%5000
    # 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, 5000):
        y_data_pred.extend(clf.predict_proba(data[i:i+5000])[:,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 [0]:

```
##cd -
```

In [0]:

```
import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score
import time
#f = open("myfile.txt", "a")
def optimised_k(X_tr, y_train, X_cr, y_val):
    """
    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 non-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.

    """
    start = time.time()

    train_auc = []
    cv_auc = []
    X = [f1, f2, f3, f4, f5, f6, f7, f8, f9, f10, f11, f12, f13, f14, f15, f16, f17, f18, f19, f20, f21, f22, f23, f24, f25, f26, f27, f28, f29, f30, f31, f32, f33, f34, f35, f36, f37, f38, f39, f40, f41, f42, f43, f44, f45, f46, f47, f48, f49, f50, f51, f52, f53, f54, f55, f56, f57, f58, f59, f60, f61, f62, f63, f64, f65, f66, f67, f68, f69, f70, f71, f72, f73, f74, f75, f76, f77, f78, f79, f80, f81, f82, f83, f84, f85, f86, f87, f88, f89, f90, f91, f92, f93, f94, f95, f96, f97, f98, f99, f100, f101, f102, f103, f104, f105, f106, f107, f108, f109, f110, f111, f112, f113, f114, f115, f116, f117, f118, f119, f120, f121, f122, f123, f124, f125, f126, f127, f128, f129, f130, f131, f132, f133, f134, f135, f136, f137, f138, f139, f140, f141, f142, f143, f144, f145, f146, f147, f148, f149, f150, f151, f152, f153, f154, f155, f156, f157, f158, f159, f160, f161, f162, f163, f164, f165, f166, f167, f168, f169, f170, f171, f172, f173, f174, f175, f176, f177, f178, f179, f180, f181, f182, f183, f184, f185, f186, f187, f188, f189, f190, f191, f192, f193, f194, f195, f196, f197, f198, f199, f200, f201, f202, f203, f204, f205, f206, f207, f208, f209, f210, f211, f212, f213, f214, f215, f216, f217, f218, f219, f220, f221, f222, f223, f224, f225, f226, f227, f228, f229, f230, f231, f232, f233, f234, f235, f236, f237, f238, f239, f240, f241, f242, f243, f244, f245, f246, f247, f248, f249, f250, f251, f252, f253, f254, f255, f256, f257, f258, f259, f260, f261, f262, f263, f264, f265, f266, f267, f268, f269, f270, f271, f272, f273, f274, f275, f276, f277, f278, f279, f280, f281, f282, f283, f284, f285, f286, f287, f288, f289, f290, f291, f292, f293, f294, f295, f296, f297, f298, f299, f300, f301, f302, f303, f304, f305, f306, f307, f308, f309, f310, f311, f312, f313, f314, f315, f316, f317, f318, f319, f320, f321, f322, f323, f324, f325, f326, f327, f328, f329, f330, f331, f332, f333, f334, f335, f336, f337, f338, f339, f340, f341, f342, f343, f344, f345, f346, f347, f348, f349, f350, f351, f352, f353, f354, f355, f356, f357, f358, f359, f360, f361, f362, f363, f364, f365, f366, f367, f368, f369, f370, f371, f372, f373, f374, f375, f376, f377, f378, f379, f380, f381, f382, f383, f384, f385, f386, f387, f388, f389, f390, f391, f392, f393, f394, f395, f396, f397, f398, f399, f400, f401, f402, f403, f404, f405, f406, f407, f408, f409, f410, f411, f412, f413, f414, f415, f416, f417, f418, f419, f420, f421, f422, f423, f424, f425, f426, f427, f428, f429, f430, f431, f432, f433, f434, f435, f436, f437, f438, f439, f440, f441, f442, f443, f444, f445, f446, f447, f448, f449, f450, f451, f452, f453, f454, f455, f456, f457, f458, f459, f460, f461, f462, f463, f464, f465, f466, f467, f468, f469, f470, f471, f472, f473, f474, f475, f476, f477, f478, f479, f480, f481, f482, f483, f484, f485, f486, f487, f488, f489, f490, f491, f492, f493, f494, f495, f496, f497, f498, f499, f500, f501, f502, f503, f504, f505, f506, f507, f508, f509, f510, f511, f512, f513, f514, f515, f516, f517, f518, f519, f520, f521, f522, f523, f524, f525, f526, f527, f528, f529, f530, f531, f532, f533, f534, f535, f536, f537, f538, f539, f540, f541, f542, f543, f544, f545, f546, f547, f548, f549, f550, f551, f552, f553, f554, f555, f556, f557, f558, f559, f560, f561, f562, f563, f564, f565, f566, f567, f568, f569, f570, f571, f572, f573, f574, f575, f576, f577, f578, f579, f580, f581, f582, f583, f584, f585, f586, f587, f588, f589, f590, f591, f592, f593, f594, f595, f596, f597, f598, f599, f600, f601, f602, f603, f604, f605, f606, f607, f608, f609, f610, f611, f612, f613, f614, f615, f616, f617, f618, f619, f620, f621, f622, f623, f624, f625, f626, f627, f628, f629, f630, f631, f632, f633, f634, f635, f636, f637, f638, f639, f640, f641, f642, f643, f644, f645, f646, f647, f648, f649, f650, f651, f652, f653, f654, f655, f656, f657, f658, f659, f660, f661, f662, f663, f664, f665, f666, f667, f668, f669, f670, f671, f672, f673, f674, f675, f676, f677, f678, f679, f680, f681, f682, f683, f684, f685, f686, f687, f688, f689, f690, f691, f692, f693, f694, f695, f696, f697, f698, f699, f700, f701, f702, f703, f704, f705, f706, f707, f708, f709, f710, f711, f712, f713, f714, f715, f716, f717, f718, f719, f720, f721, f722, f723, f724, f725, f726, f727, f728, f729, f730, f731, f732, f733, f734, f735, f736, f737, f738, f739, f740, f741, f742, f743, f744, f745, f746, f747, f748, f749, f750, f751, f752, f753, f754, f755, f756, f757, f758, f759, f760, f761, f762, f763, f764, f765, f766, f767, f768, f769, f770, f771, f772, f773, f774, f775, f776, f777, f778, f779, f780, f781, f782, f783, f784, f785, f786, f787, f788, f789, f790, f791, f792, f793, f794, f795, f796, f797, f798, f799, f800, f801, f802, f803, f804, f805, f806, f807, f808, f809, f810, f811, f812, f813, f814, f815, f816, f817, f818, f819, f820, f821, f822, f823, f824, f825, f826, f827, f828, f829, f830, f831, f832, f833, f834, f835, f836, f837, f838, f839, f840, f841, f842, f843, f844, f845, f846, f847, f848, f849, f850, f851, f852, f853, f854, f855, f856, f857, f858, f859, f860, f861, f862, f863, f864, f865, f866, f867, f868, f869, f870, f871, f872, f873, f874, f875, f876, f877, f878, f879, f880, f881, f882, f883, f884, f885, f886, f887, f888, f889, f890, f891, f892, f893, f894, f895, f896, f897, f898, f899, f900, f901, f902, f903, f904, f905, f906, f907, f908, f909, f910, f911, f912, f913, f914, f915, f916, f917, f918, f919, f920, f921, f922, f923, f924, f925, f926, f927, f928, f929, f930, f931, f932, f933, f934, f935, f936, f937, f938, f939, f940, f941, f942, f943, f944, f945, f946, f947, f948, f949, f950, f951, f952, f953, f954, f955, f956, f957, f958, f959, f960, f961, f962, f963, f964, f965, f966, f967, f968, f969, f970, f971, f972, f973, f974, f975, f976, f977, f978, f979, f980, f981, f982, f983, f984, f985, f986, f987, f988, f989, f990, f991, f992, f993, f994, f995, f996, f997, f998, f999, f1000, f1001, f1002, f1003, f1004, f1005, f1006, f1007, f1008, f1009, f1010, f1011, f1012, f1013, f1014, f1015, f1016, f1017, f1018, f1019, f1020, f1021, f1022, f1023, f1024, f1025, f1026, f1027, f1028, f1029, f1030, f1031, f1032, f1033, f1034, f1035, f1036, f1037, f1038, f1039, f1040, f1041, f1042, f1043, f1044, f1045, f1046, f1047, f1048, f1049, f1050, f1051, f1052, f1053, f1054, f1055, f1056, f1057, f1058, f1059, f1060, f1061, f1062, f1063, f1064, f1065, f1066, f1067, f1068, f1069, f1070, f1071, f1072, f1073, f1074, f1075, f1076, f1077, f1078, f1079, f1080, f1081, f1082, f1083, f1084, f1085, f1086, f1087, f1088, f1089, f1090, f1091, f1092, f1093, f1094, f1095, f1096, f1097, f1098, f1099, f1100, f1101, f1102, f1103, f1104, f1105, f1106, f1107, f1108, f1109, f1110, f1111, f1112, f1113, f1114, f1115, f1116, f1117, f1118, f1119, f1120, f1121, f1122, f1123, f1124, f1125, f1126, f1127, f1128, f1129, f1130, f1131, f1132, f1133, f1134, f1135, f1136, f1137, f1138, f1139, f1140, f1141, f1142, f1143, f1144, f1145, f1146, f1147, f1148, f1149, f1150, f1151, f1152, f1153, f1154, f1155, f1156, f1157, f1158, f1159, f1160, f1161, f1162, f1163, f1164, f1165, f1166, f1167, f1168, f1169, f1170, f1171, f1172, f1173, f1174, f1175, f1176, f1177, f1178, f1179, f1180, f1181, f1182, f1183, f1184, f1185, f1186, f1187, f1188, f1189, f1190, f1191, f1192, f1193, f1194, f1195, f1196, f1197, f1198, f1199, f1200, f1201, f1202, f1203, f1204, f1205, f1206, f1207, f1208, f1209, f1210, f1211, f1212, f1213, f1214, f1215, f1216, f1217, f1218, f1219, f1220, f1221, f1222, f1223, f1224, f1225, f1226, f1227, f1228, f1229, f1230, f1231, f1232, f1233, f1234, f1235, f1236, f1237, f1238, f1239, f1240, f1241, f1242, f1243, f1244, f1245, f1246, f1247, f1248, f1249, f1250, f1251, f1252, f1253, f1254, f1255, f1256, f1257, f1258, f1259, f1260, f1261, f1262, f1263, f1264, f1265, f1266, f1267, f1268, f1269, f1270, f1271, f1272, f1273, f1274, f1275, f1276, f1277, f1278, f1279, f1280, f1281, f1282, f1283, f1284, f1285, f1286, f1287, f1288, f1289, f1290, f1291, f1292, f1293, f1294, f1295, f1296, f1297, f1298, f1299, f1300, f1301, f1302, f1303, f1304, f1305, f1306, f1307, f1308, f1309, f1310, f1311, f1312, f1313, f1314, f1315, f1316, f1317, f1318, f1319, f1320, f1321, f1322, f1323, f1324, f1325, f1326, f1327, f1328, f1329, f1330, f1331, f1332, f1333, f1334, f1335, f1336, f1337, f1338, f1339, f1340, f1341, f1342, f1343, f1344, f1345, f1346, f1347, f1348, f1349, f1350, f1351, f1352, f1353, f1354, f1355, f1356, f1357, f1358, f1359, f1360, f1361, f1362, f1363, f1364, f1365, f1366, f1367, f1368, f1369, f1370, f1371, f1372, f1373, f1374, f1375, f1376, f1377, f1378, f1379, f1380, f1381, f1382, f1383, f1384, f1385, f1386, f1387, f1388, f1389, f1390, f1391, f1392, f1393, f1394, f1395, f1396, f1397, f1398, f1399, f1400, f1401, f1402, f1403, f1404, f1405, f1406, f1407, f1408, f1409, f1410, f1411, f1412, f1413, f1414, f1415, f1416, f1417, f1418, f1419, f1420, f1421, f1422, f1423, f1424, f1425, f1426, f1427, f1428, f1429, f1430, f1431, f1432, f1433, f1434, f1435, f1436, f1437, f1438, f1439, f1440, f1441, f1442, f1443, f1444, f1445, f1446, f1447, f1448, f1449, f1450, f1451, f1452, f1453, f1454, f1455, f1456, f1457, f1458, f1459, f1460, f1461, f1462, f1463, f1464, f1465, f1466, f1467, f1468, f1469, f1470, f1471, f1472, f1473, f1474, f1475, f1476, f1477, f1478, f1479, f1480, f1481, f1482, f1483, f1484, f1485, f1486, f1487, f1488, f1489, f1490, f1491, f1492, f1493, f1494, f1495, f1496, f1497, f1498, f1499, f1500, f1501, f1502, f1503, f1504, f1505, f1506, f1507, f1508, f1509, f1510, f1511, f1512, f1513, f1514, f1515, f1516, f1517, f1518, f1519, f1520, f1521, f1522, f1523, f1524, f1525, f1526, f1527, f1528, f1529, f1530, f1531, f1532, f1533, f1534, f1535, f1536, f1537, f1538, f1539, f1540, f1541, f1542, f1543, f1544, f1545, f1546, f1547, f1548, f1549, f1550, f1551, f1552, f1553, f1554, f1555, f1556, f1557, f1558, f1559, f1560, f1561, f1562, f1563, f1564, f1565, f1566, f1567, f1568, f1569, f1570, f1571, f1572, f1573, f1574, f1575, f1576, f1577, f1578, f1579, f1580, f1581, f1582, f1583, f1584, f1585, f1586, f1587, f1588, f1589, f1590, f1591, f1592, f1593, f1594, f1595, f1596, f1597, f1598, f1599, f1600, f1601, f1602, f1603, f1604, f1605, f1606, f1607, f1608, f1609, f1610, f1611, f1612, f1613, f1614, f1615, f1616, f1617, f1618, f1619, f1620, f1621, f1622, f1623, f1624, f1625, f1626, f1627, f1628, f1629, f1630, f1631, f1632, f1633, f1634, f1635, f1636, f1637, f1638, f1639, f1640, f1641, f1642, f1643, f1644, f1645, f1646, f1647, f1648, f1649, f1650, f1651, f1652, f1653, f1654, f1655, f1656, f1657, f1658, f1659, f1660, f1661, f1662, f1663, f1664, f1665, f1666, f1667, f1668, f1669, f1670, f1671, f1672, f1673, f1674, f1675, f1676, f1677, f1678, f1679, f1680, f1681, f1682, f1683, f1684, f1685, f1686, f1687, f1688, f1689, f1690, f1691, f1692, f1693, f1694, f1695, f1696, f1697, f1698, f1699, f1700, f1701, f1702, f1703, f1704, f1705, f1706, f1707, f1708, f1709, f1710, f1711, f1712, f1713, f1714, f1715, f1716, f1717, f1718, f1719, f1720, f1721, f1722, f1723, f1724, f1725, f1726, f1727, f1728, f1729, f1730, f1731, f1732, f1733, f1734, f1735, f1736, f1737, f1738, f1739, f1740, f1741, f1742, f1743, f1744, f1745, f1746, f1747, f1748, f1749, f1750, f1751, f1752, f1753, f1754, f1755, f1756, f1757, f1758, f1759, f1760, f1761, f1762, f1763, f1764, f1765, f1766, f1767, f1768, f1769, f1770, f1771, f1772, f1773, f1774, f1775, f1776, f1777, f1778, f1779, f1780, f1781, f1782, f1783, f1784, f1785, f1786, f1787, f1788, f1789, f1790, f1791, f1792, f1793, f1794, f1795, f1796, f1797, f1798, f1799, f1800, f1801, f1802, f1803, f1804, f1805, f1806, f1807, f1808, f1809, f1810, f1811, f1812, f1813, f1814, f1815, f1816, f1817, f1818, f1819, f1820, f1821, f1822, f1823, f1824, f1825, f1826, f1827, f1828, f1829, f1830, f1831, f1832, f1833, f1834, f1835, f1836, f1837, f1838, f1839, f1840, f1841, f1842, f1843, f1844, f1845, f1846, f1847, f1848, f1849, f1850, f1851, f1852, f1853, f1854, f1855, f1856, f1857, f1858, f1859, f1860, f1861, f1862, f1863, f1864, f1865, f1866, f1867, f1868, f1869, f1870, f1871, f1872, f1873, f1874, f1875, f1876, f1877, f1878, f1879, f1880, f1881, f1882, f1883, f1884, f1885, f1886, f1887, f1888, f1889, f1890, f1891, f1892, f1893, f1894, f1895, f1896, f1897, f1898, f1899, f1900, f1901, f1902, f1903, f1904, f1905, f1906, f1907, f1908, f1909, f1910, f1911, f1912, f1913, f1914, f1915, f1916, f1917, f1918, f1919, f1920, f1921, f1922, f1923, f1924, f1925, f1926, f1927, f1928, f1929, f1930, f1931, f1932, f1933, f1934, f1935, f1936, f1937, f1938, f1939, f1940, f1941, f1942, f1943, f1944, f1945, f1946, f1947, f1948, f1949, f1950, f1951, f1952, f1953, f1954, f1955, f1956, f1957, f1958, f1959, f1960, f1961, f1962, f1963, f1964, f1965, f1966, f1967, f1968, f1969, f1970, f1971, f1972, f1973, f1974, f1975, f1976, f1977, f1978, f1979, f1980, f1981, f1982, f1983, f1984, f1985, f1986, f1987, f1988, f1989, f1990, f1991, f1992, f1993, f1994, f1995, f1996, f1997, f1998, f1999, f2000, f2001, f2002, f2003, f2004, f2005, f2006, f2007, f2008, f2009, f2010, f2011, f2012, f2013, f2014, f2015, f2016, f2017, f2018, f2019, f2020, f2021, f2022, f2023, f2024, f2025, f2026, f2027, f2028, f2029, f2030, f2031, f2032, f2033, f2034, f2035, f2036, f2037, f2038, f2039, f2040, f2041, f2042, f2043, f2044, f2045, f2046, f2047, f2048, f2049, f2050, f2051, f2052, f2053, f2054, f2055, f2056, f2057, f2058, f2059, f2060, f2061, f2062, f2063, f2064, f2065, f2066, f2067, f2068, f2069, f2070, f2071, f2
```

```

K = [1, 15, 51, 101, 151, 201, 251, 301, 401]
for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i, algorithm = "brute")
    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_val, y_cv_pred))
    #f.write('> %d\n' % (i,))
    #f.write('> %s\n' % (roc_auc_score(y_train, y_train_pred),))
    #f.write('%s\n' % (roc_auc_score(y_val, 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()
end = time.time()
minutes = int((end - start)/60)
print("execution time in minutes:", minutes)
print("execution time in hours:", float(minutes/60))
print(max(cv_auc))
optimalK = K[np.argmax(cv_auc)]
print("optimal K is: ", optimalK)
return optimalK

```

In [0]:

```

# 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 [0]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.metrics import confusion_matrix
import seaborn as sns
def error_plot(X_train, X_te, y_train, y_test):
    neigh = KNeighborsClassifier(n_neighbors=best_k, algorithm = "brute")
    neigh.fit(X_train, 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_train)
    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()
print("="*100)

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)))

class_label = ["0:rejected", "1:approved"]
cm = confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_tpr))
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()

class_label = ["0:rejected", "1:approved"]
cm = confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_tpr))
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()

```

Set 3: categorical, numerical features + project_title(AVG W2V)+ preprocessed_essay (AVG W2V)

In [0]:

```

#%cd Assignments_DonorsChoose_2018
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
# make sure you have the glove_vectors file
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())

```

In [0]:

```

#%cd -

```

In [0]:

```

def avg_w2v(preprocessed_list):
    # average Word2Vec
    preprocessed_list = preprocessed_list[:]
    # compute average word2vec for each review.
    avg_w2v_vectors_list = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in tqdm(preprocessed_list): # for each review/sentence
        vector = np.zeros(30) # 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][:30]
                cnt_words += 1
        if cnt_words != 0:
            vector /= cnt_words
        avg_w2v_vectors_list.append(vector)

    print(len(avg_w2v_vectors_list))
    print(len(avg_w2v_vectors_list[0]))
    return avg_w2v_vectors_list

```

In [0]:

```
avgw2v_vec_essay_tr = avg_w2v(textpreprocessed(X_train,'essay'))
avgw2v_vec_essay_te = avg_w2v(textpreprocessed(X_test,'essay'))
avgw2v_vec_essay_val = avg_w2v(textpreprocessed(X_val,'essay'))
```

```
100%|██████████| 61178/61178 [00:38<00:00, 1588.16it/s]
100%|██████████| 61178/61178 [00:16<00:00, 3786.16it/s]
 1%|█          | 158/21850 [00:00<00:13, 1573.61it/s]
```

61178
30

```
100%|██████████| 21850/21850 [00:13<00:00, 1588.33it/s]
100%|██████████| 21850/21850 [00:05<00:00, 3737.31it/s]
 1%|█          | 164/26220 [00:00<00:15, 1637.83it/s]
```

21850
30

```
100%|██████████| 26220/26220 [00:16<00:00, 1588.60it/s]
100%|██████████| 26220/26220 [00:07<00:00, 3307.26it/s]
```

26220
30

In [0]:

```
avgw2v_vec_titles_tr = avg_w2v(textpreprocessed(X_train,'project_title'))
avgw2v_vec_titles_te = avg_w2v(textpreprocessed(X_test,'project_title'))
avgw2v_vec_titles_val = avg_w2v(textpreprocessed(X_val,'project_title'))
```

```
100%|██████████| 61178/61178 [00:01<00:00, 33165.44it/s]
100%|██████████| 61178/61178 [00:00<00:00, 86803.29it/s]
16%|███        | 3394/21850 [00:00<00:00, 33936.63it/s]
```

61178
30

```
100%|██████████| 21850/21850 [00:00<00:00, 34205.62it/s]
100%|██████████| 21850/21850 [00:00<00:00, 84418.17it/s]
13%|███        | 3388/26220 [00:00<00:00, 33873.73it/s]
```

21850
30

```
100%|██████████| 26220/26220 [00:00<00:00, 33631.82it/s]
100%|██████████| 26220/26220 [00:00<00:00, 85093.35it/s]
```

26220
30

In [0]:

```
avgw2v_vec_resource_tr = avg_w2v(textpreprocessed(X_train,'project_resource_summary'))
avgw2v_vec_resource_te = avg_w2v(textpreprocessed(X_test,'project_resource_summary'))
avgw2v_vec_resource_val = avg_w2v(textpreprocessed(X_val,'project_resource_summary'))
```

```
100%|██████████| 61178/61178 [00:04<00:00, 15195.72it/s]
100%|██████████| 61178/61178 [00:01<00:00, 35479.11it/s]
 7%|██         | 1500/21850 [00:00<00:01, 14986.15it/s]
```

61178
30

```
100%|██████████| 21850/21850 [00:01<00:00, 15182.10it/s]
100%|██████████| 21850/21850 [00:00<00:00, 36305.57it/s]
 6%|███████| 1504/26220 [00:00<00:01, 15039.05it/s]
```

21850
30

```
100%|██████████| 26220/26220 [00:01<00:00, 15343.34it/s]
100%|██████████| 26220/26220 [00:00<00:00, 35143.17it/s]
```

26220
30

In [0]:

```
X_tr, X_cr, X_te = hstack_data(avgw2v_vec_titles_tr, avgw2v_vec_titles_val, avgw2v_vec_titles_te, \
                               avgw2v_vec_essay_tr, avgw2v_vec_essay_val, avgw2v_vec_essay_te, \
                               avgw2v_vec_resource_tr, avgw2v_vec_resource_val, avgw2v_vec_resource_te)
```

In [0]:

```
X_tr.shape,X_cr.shape,X_te.shape
```

Out[0]:

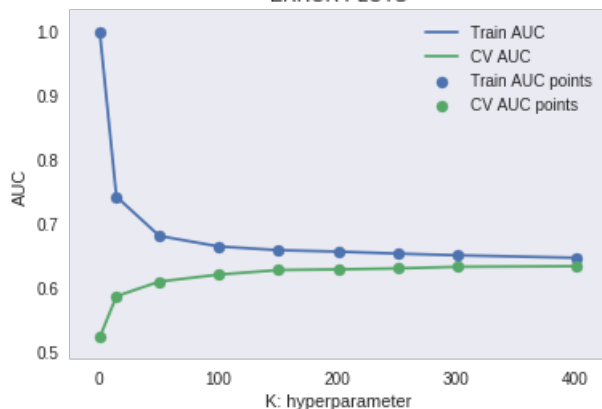
```
((61178, 172), (26220, 172), (21850, 172))
```

In [0]:

```
best_k = optimised_k(X_tr,y_train,X_cr,y_val)
```

```
 0%|          | 0/9 [00:00<?, ?it/s]
11%|███       | 1/9 [45:40<6:05:21, 2740.17s/it]
22%|████      | 2/9 [1:32:06<5:21:17, 2753.89s/it]
33%|██████    | 3/9 [2:18:21<4:36:01, 2760.22s/it]
44%|████████  | 4/9 [3:04:41<3:50:31, 2766.21s/it]
56%|██████████| 5/9 [3:51:02<3:04:43, 2770.81s/it]
67%|███████████| 6/9 [4:37:27<2:18:44, 2774.93s/it]
78%|███████████| 7/9 [5:23:51<1:32:35, 2777.79s/it]
89%|███████████| 8/9 [6:10:10<46:18, 2778.13s/it]
100%|███████████| 9/9 [6:56:38<00:00, 2781.07s/it]
```

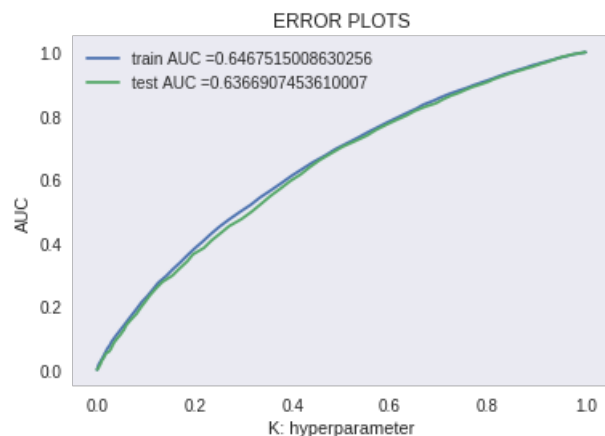
ERROR PLOTS



execution time in minutes: 416
execution time in hours: 6.933333333333333
0.6335998301870773
optimal K is: 401

In [0]:

```
error_plot(x_tr,x_te,y_train,y_test)
```



Train confusion matrix

the maximum value of $tpr \cdot (1-fpr)$ 0.24986156181299424 for threshold 0.835

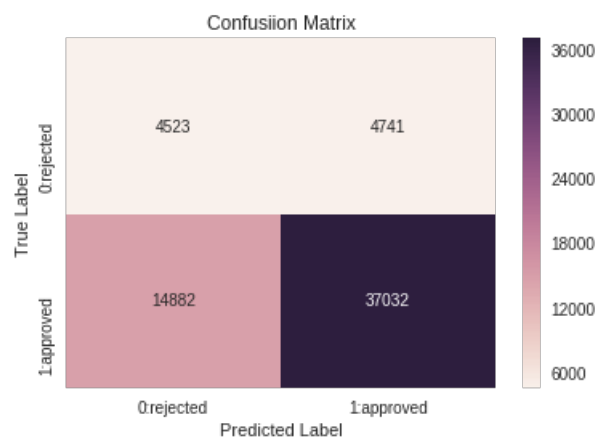
```
[[ 4523  4741]
 [14882 37032]]
```

Test confusion matrix

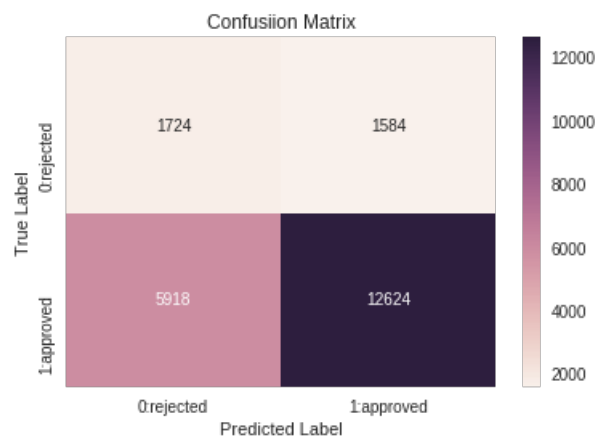
the maximum value of $tpr \cdot (1-fpr)$ 0.24999415143969622 for threshold 0.84

```
[[ 1724  1584]
 [ 5918 12624]]
```

the maximum value of $tpr \cdot (1-fpr)$ 0.24986156181299424 for threshold 0.835



the maximum value of $tpr \cdot (1-fpr)$ 0.24999415143969622 for threshold 0.84



Set 2: categorical, numerical features + project_title(TFIDF)+ preprocessed_essay (TFIDF)

1.5.1 Vectorizing Categorical data

In [0]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
def tfidf_vec(preprocessed_data_tr,preprocessed_data_val,preprocessed_data_te):
    vectorizer = TfidfVectorizer(min_df=10)
    text_tfidf_tr = vectorizer.fit_transform(preprocessed_data_tr)
    print("Shape of matrix after one hot encodig ",text_tfidf_tr.shape)

    text_tfidf_val = vectorizer.transform(preprocessed_data_val)
    print("Shape of matrix after one hot encodig ",text_tfidf_val.shape)

    text_tfidf_te = vectorizer.transform(preprocessed_data_te)
    print("Shape of matrix after one hot encodig ",text_tfidf_te.shape)
    return text_tfidf_tr,text_tfidf_val, text_tfidf_te
```

In [0]:

```
tfidf_vec_essay_tr,tfidf_vec_essay_val,tfidf_vec_essay_te =
tfidf_vec(textpreprocessed(X_train,'essay'),textpreprocessed(X_val,'essay'), textpreprocessed(X_test,'essay'))
tfidf_vec_titles_tr,tfidf_vec_titles_val,tfidf_vec_titles_te = tfidf_vec(textpreprocessed(X_train,
'project_title'),textpreprocessed(X_val,'project_title'), textpreprocessed(X_test,'project_title')
)
tfidf_vec_resource_tr,tfidf_vec_resource_val,tfidf_vec_resource_te = tfidf_vec(textpreprocessed(X_train,
'project_resource_summary'),textpreprocessed(X_val,'project_resource_summary'),
textpreprocessed(X_test,'project_resource_summary'))
```

```
100%|██████████| 61178/61178 [00:37<00:00, 1646.00it/s]
100%|██████████| 26220/26220 [00:15<00:00, 1644.27it/s]
100%|██████████| 21850/21850 [00:13<00:00, 1642.86it/s]
```

```
Shape of matrix after one hot encodig (61178, 13150)
Shape of matrix after one hot encodig (26220, 13150)
```

```
6%|███████| 3512/61178 [00:00<00:01, 35119.11it/s]
```

```
Shape of matrix after one hot encodig (21850, 13150)
```

```
100%|██████████| 61178/61178 [00:01<00:00, 35300.24it/s]
100%|██████████| 26220/26220 [00:00<00:00, 35023.29it/s]
100%|██████████| 21850/21850 [00:00<00:00, 34172.18it/s]
```

```
Shape of matrix after one hot encodig (61178, 2274)
Shape of matrix after one hot encodig (26220, 2274)
```

```
3%|███████| 1550/61178 [00:00<00:03, 15495.10it/s]
```

```
Shape of matrix after one hot encodig (21850, 2274)
```

```
100%|██████████| 61178/61178 [00:03<00:00, 15655.91it/s]
100%|██████████| 26220/26220 [00:01<00:00, 15722.07it/s]
100%|██████████| 21850/21850 [00:01<00:00, 15475.88it/s]
```

```
Shape of matrix after one hot encodig (61178, 4348)
Shape of matrix after one hot encodig (26220, 4348)
Shape of matrix after one hot encodig (21850, 4348)
```

In [0]:

```
print(one_hot_num_dig_tr.shape)
print(grade_one_hot_tr.shape)
print(state_one_hot_tr.shape)
print(cat_one_hot_tr.shape)
print(cat_sub_one_hot_tr.shape)
print(teacher_one_hot_tr.shape)
print(price_standardized_tr.shape)
print(project_standardized_tr.shape)
print(tfidf_vec_titles_tr.shape)
```

```
print(tfidf_vec_essay_tr.shape)
```

```
(61178, 2)
(61178, 4)
(61178, 51)
(61178, 9)
(61178, 9)
(61178, 5)
(61178, 1)
(61178, 1)
(61178, 2274)
(61178, 13150)
```

In [0]:

```
X_tr, X_cr, X_te = hstack_data(tfidf_vec_titles_tr, tfidf_vec_titles_val, tfidf_vec_titles_te, \
                                tfidf_vec_essay_tr, tfidf_vec_essay_val, tfidf_vec_essay_te, \
                                tfidf_vec_resource_tr, tfidf_vec_resource_val, tfidf_vec_resource_te)
```

In [0]:

```
X_tr.shape,X_cr.shape,X_te.shape
```

Out[0]:

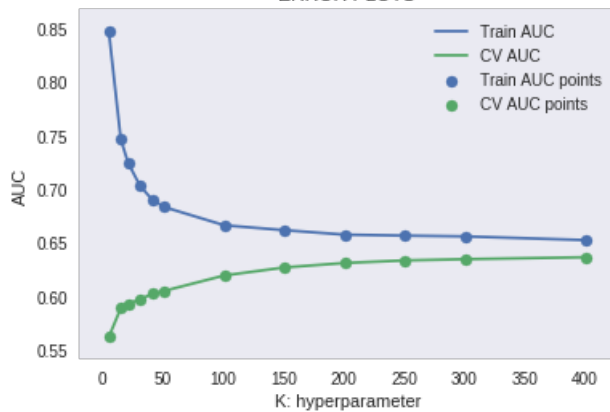
```
((61178, 19854), (26220, 19854), (21850, 19854))
```

In [0]:

```
best_k = optimised_k(X_tr,y_train,X_cr,y_val)
```

100%|██████████| 12/12 [2:39:39<00:00, 803.62s/it]

ERROR PLOTS



execution time in minutes: 159

execution time in hours: 2

0.6359927037047518

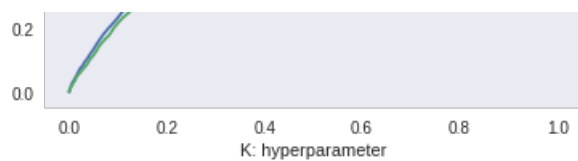
optimal K is: 401

In [0]:

```
error_plot(X_tr,X_te,y_train,y_test)
```

ERROR PLOTS

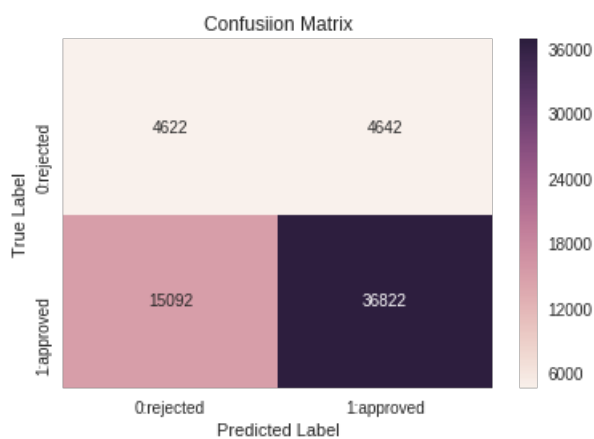




```

=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999883479347693 for threshold 0.83
[[ 4622  4642]
 [15092 36822]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.24998208878406972 for threshold 0.835
[[ 1741  1567]
 [ 6053 12489]]
the maximum value of tpr*(1-fpr) 0.24999883479347693 for threshold 0.83

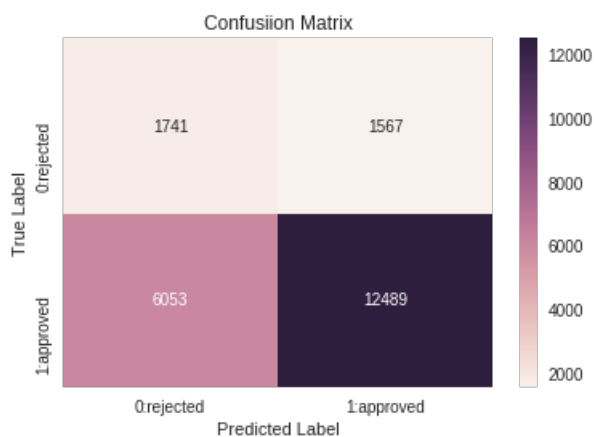
```



```

the maximum value of tpr*(1-fpr) 0.24998208878406972 for threshold 0.835

```



[Task-2]

In [0]:

```

from sklearn.feature_selection import SelectKBest, chi2
X, y = X_tr, y_train
print(X.shape)
X_new = SelectKBest(chi2, k=2000).fit_transform(X, y)
print(X_new.shape)
selector = SelectKBest(chi2, k=2000)
selector.fit(X, y)

X_new = selector.transform(X)
X_new.shape
print(selector.get_support(indices=True))

'''indexes selected = selector.get_support(indices=True) # extracting features indexes

```

```
selected_features = [] # it will contain features names
for i in indexes_selected:
    selected_features.append(features[i])'''
```

```
(61178, 19854)
(61178, 2000)
[ 0 1 2 ... 19835 19838 19841]
```

Out[0]:

```
'indexes_selected = selector.get_support(indices=True) # extracting features
indexes\nselected_features = [] # it will contain features names\nfor i in indexes_selected:\n se
lected_features.append(features[i])'
```

In [0]:

```
indexes_selected = selector.get_support(indices=True) # extracting features indexes
selected_features = [] # it will contain features names
for i in indexes_selected:
    selected_features.append(i)
```

In [0]:

```
df = X_te.toarray()
df1 = pd.DataFrame(df)
X_te_new = csr_matrix(df1[list(selected_features)].values)
X_te.shape,X_te_new.shape
```

Out[0]:

```
((21850, 19854), (21850, 2000))
```

In [0]:

```
df = X_cr.toarray()
df1 = pd.DataFrame(df)
X_cr_new = csr_matrix(df1[list(selected_features)].values)
X_cr.shape,X_cr_new.shape
```

Out[0]:

```
((26220, 19854), (26220, 2000))
```

In [0]:

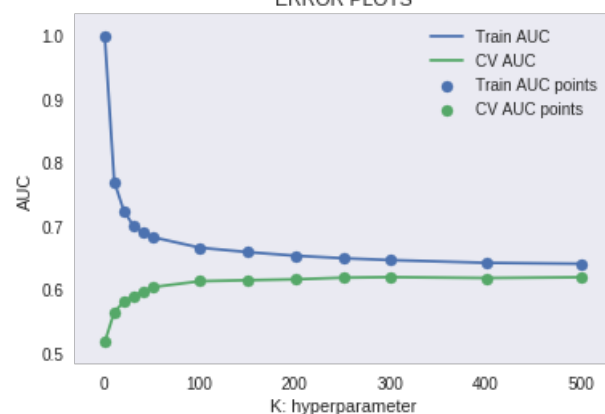
```
X_tr_new = X_new
```

In [0]:

```
best_k = optimised_k(X_tr_new,y_train,X_cr_new,y_val)
```

```
100%|██████████| 13/13 [1:36:22<00:00, 452.57s/it]
```

ERROR PLOTS



execution time in minutes: 96
execution time in hours: 1.6
0.6186178529986132
optimal K is: 301

In [0]:

```
error_plot(X_tr_new,X_te_new,y_train,y_test)
```



Train confusion matrix

the maximum value of $\text{tpr} \cdot (1 - \text{fpr})$ 0.2499731536417085 for threshold 0.837

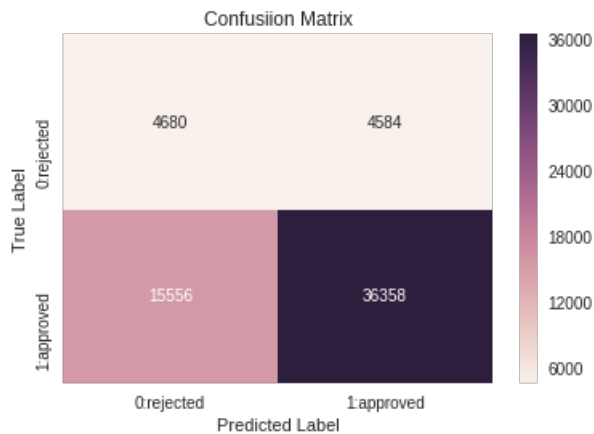
```
[[ 4680  4584]  
 [15556 36358]]
```

Test confusion matrix

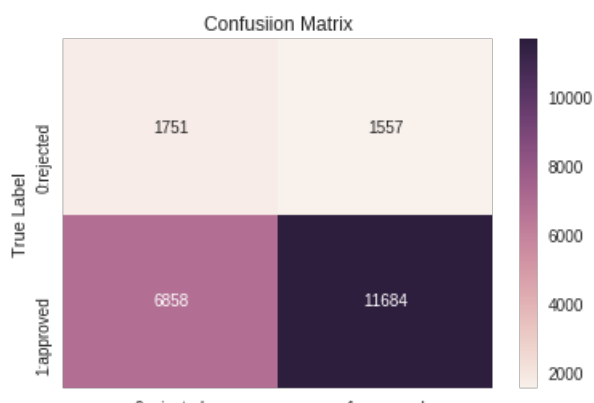
the maximum value of $\text{tpr} \cdot (1 - \text{fpr})$ 0.24998894256567566 for threshold 0.844

```
[[ 1751  1557]  
 [ 6858 11684]]
```

the maximum value of $\text{tpr} \cdot (1 - \text{fpr})$ 0.2499731536417085 for threshold 0.837



the maximum value of $\text{tpr} \cdot (1 - \text{fpr})$ 0.24998894256567566 for threshold 0.844



Set 1: categorical, numerical features + project_title(BOW) + preprocessed_essay (BOW)

In [0]:

```
def bow_vec(preprocessed_data_tr,preprocessed_data_val,preprocessed_data_te):
    vectorizer = CountVectorizer(min_df=10)
    text_bow_tr = vectorizer.fit_transform(preprocessed_data_tr)
    print("Shape of matrix after one hot encodig ",text_bow_tr.shape)

    text_bow_val = vectorizer.transform(preprocessed_data_val)
    print("Shape of matrix after one hot encodig ",text_bow_val.shape)

    text_bow_te = vectorizer.transform(preprocessed_data_te)
    print("Shape of matrix after one hot encodig ",text_bow_te.shape)
    return text_bow_tr,text_bow_val, text_bow_te
```

In [0]:

```
bow_vec_essay_tr,bow_vec_essay_val,bow_vec_essay_te = bow_vec(textpreprocessed(X_train,'essay'),textpreprocessed(X_val,'essay'), textpreprocessed(X_test,'essay'))
bow_vec_titles_tr,bow_vec_titles_val,bow_vec_titles_te =
bow_vec(textpreprocessed(X_train,'project_title'),textpreprocessed(X_val,'project_title'),
textpreprocessed(X_test,'project_title'))
bow_vec_resource_tr,bow_vec_resource_val,bow_vec_resource_te =
bow_vec(textpreprocessed(X_train,'project_resource_summary'),textpreprocessed(X_val,'project_resource_summary'), textpreprocessed(X_test,'project_resource_summary'))
```

```
100%|██████████| 61178/61178 [00:38<00:00, 1603.54it/s]
100%|██████████| 26220/26220 [00:15<00:00, 1640.37it/s]
100%|██████████| 21850/21850 [00:13<00:00, 1636.87it/s]
```

```
Shape of matrix after one hot encodig (61178, 13145)
Shape of matrix after one hot encodig (26220, 13145)
```

```
6%|███████| 3556/61178 [00:00<00:01, 35553.42it/s]
```

```
Shape of matrix after one hot encodig (21850, 13145)
```

```
100%|██████████| 61178/61178 [00:01<00:00, 35204.18it/s]
100%|██████████| 26220/26220 [00:00<00:00, 34550.18it/s]
100%|██████████| 21850/21850 [00:00<00:00, 34856.56it/s]
```

```
Shape of matrix after one hot encodig (61178, 2295)
Shape of matrix after one hot encodig (26220, 2295)
```

```
3%|███████| 1606/61178 [00:00<00:03, 16053.74it/s]
```

```
Shape of matrix after one hot encodig (21850, 2295)
```

```
100%|██████████| 61178/61178 [00:03<00:00, 15566.92it/s]
100%|██████████| 26220/26220 [00:01<00:00, 15555.07it/s]
100%|██████████| 21850/21850 [00:01<00:00, 15482.56it/s]
```

```
Shape of matrix after one hot encodig (61178, 4353)
Shape of matrix after one hot encodig (26220, 4353)
Shape of matrix after one hot encodig (21850, 4353)
```

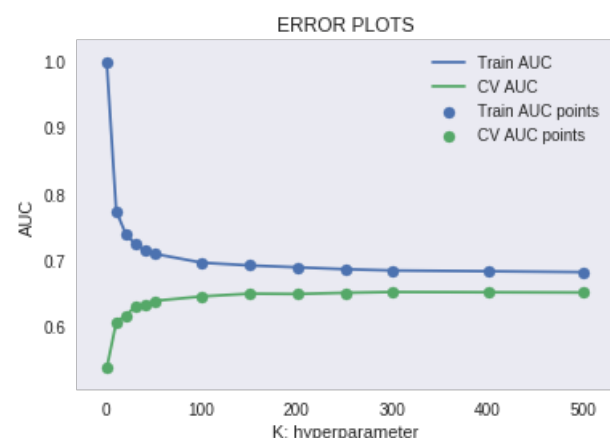
In [0]:

```
X_tr, X_cr, X_te = hstack_data(bow_vec_titles_tr, bow_vec_titles_val, bow_vec_titles_te, \
                               bow_vec_essay_tr, bow_vec_essay_val, bow_vec_essay_te, \
                               bow_vec_resource_tr, bow_vec_resource_val, bow_vec_resource_te)
```

In [0]:

```
best_k = optimised_k(X_tr,y_train,X_cr,y_val)
```

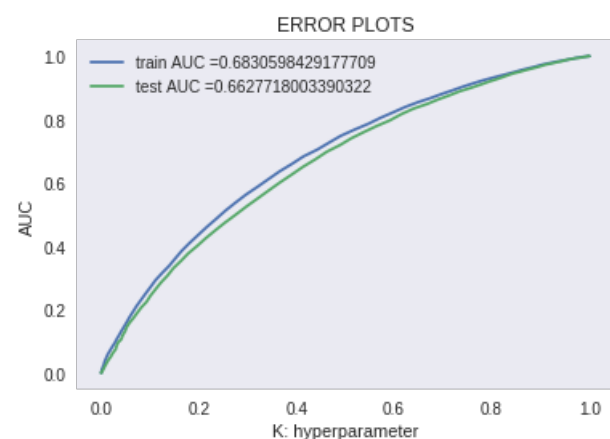
100%|██████████| 13/13 [2:54:35<00:00, 812.13s/it]



execution time in minutes: 174
execution time in hours: 2.9
0.6504490815951093
optimal K is: 301

In [0]:

```
error_plot(X_tr,X_te,y_train,y_test)
```



Train confusion matrix

the maximum value of $tpr \cdot (1 - fpr)$ 0.2499056182716315 for threshold 0.781

```
[[ 4722  4542]  
 [13089 38825]]
```

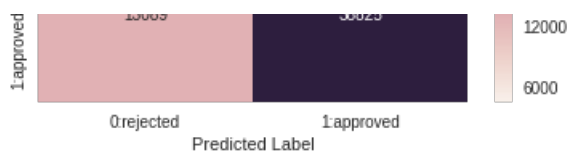
Test confusion matrix

the maximum value of $tpr \cdot (1 - fpr)$ 0.2499004830910811 for threshold 0.794

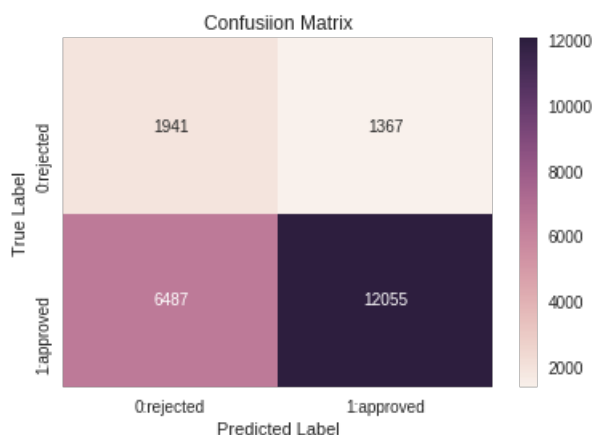
```
[[ 1941  1367]  
 [ 6487 12055]]
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.2499056182716315 for threshold 0.781





the maximum value of $tpr \cdot (1 - fpr)$ 0.2499004830910811 for threshold 0.794



Set 4: categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

In [0]:

```
# average Word2Vec
# compute average word2vec for each review.
def tfidf_w2v(preprocessed_list):
    # average Word2Vec
    preprocessed_list = preprocessed_list[:]
    tfidf_model = TfidfVectorizer()
    tfidf_model.fit(preprocessed_list)
    # 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())
    tfidf_w2v_vectors_list = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in tqdm(preprocessed_list): # for each review/sentence
        vector = np.zeros(30) # 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][:30] # getting the vector for each word
                # here we are multiplying idf value(dictionary[word]) and the tf
                value((sentence.count(word)/len(sentence.split())))
                tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the
                tfidf value for each word
                vector += (vec * tf_idf) # calculating tfidf weighted w2v
                tf_idf_weight += tf_idf
        if tf_idf_weight != 0:
            vector /= tf_idf_weight
        tfidf_w2v_vectors_list.append(vector)

    print(len(tfidf_w2v_vectors_list))
    print(len(tfidf_w2v_vectors_list[0]))
    return tfidf_w2v_vectors_list
```

In [0]:

```
tfidf_w2v_vec_essay_tr = tfidf_w2v(textpreprocessed(X_train, 'essay'))
tfidf_w2v_vec_essay_te = tfidf_w2v(textpreprocessed(X_test, 'essay'))
tfidf_w2v_vec_essay_val = tfidf_w2v(textpreprocessed(X_val, 'essay'))
```

```
100%|██████████| 61178/61178 [00:37<00:00, 1651.96it/s]
100%|██████████| 61178/61178 [01:43<00:00, 590.46it/s]
1%|          | 162/21850 [00:00<00:13, 1616.00it/s]
```

61178
30

100%|██████████| 21850/21850 [00:13<00:00, 1649.94it/s]
100%|██████████| 21850/21850 [00:38<00:00, 572.61it/s]
1%|███████| 166/26220 [00:00<00:15, 1654.00it/s]

21850
30

100%|██████████| 26220/26220 [00:15<00:00, 1644.64it/s]
100%|██████████| 26220/26220 [00:45<00:00, 578.58it/s]

26220
30

In [0]:

```
tfidf_w2v_vec_titles_tr = tfidf_w2v(textpreprocessed(X_train,'project_title'))  
tfidf_w2v_vec_titles_te = tfidf_w2v(textpreprocessed(X_test,'project_title'))  
tfidf_w2v_vec_titles_val = tfidf_w2v(textpreprocessed(X_val,'project_title'))
```

100%|██████████| 61178/61178 [00:01<00:00, 35317.46it/s]
100%|██████████| 61178/61178 [00:01<00:00, 45150.33it/s]
16%|███████| 3547/21850 [00:00<00:00, 35463.69it/s]

61178
30

100%|██████████| 21850/21850 [00:00<00:00, 35215.37it/s]
100%|██████████| 21850/21850 [00:00<00:00, 44991.44it/s]
14%|███████| 3555/26220 [00:00<00:00, 35541.56it/s]

21850
30

100%|██████████| 26220/26220 [00:00<00:00, 35337.65it/s]
100%|██████████| 26220/26220 [00:00<00:00, 45371.03it/s]

26220
30

In [0]:

```
tfidf_w2v_vec_resource_tr = tfidf_w2v(textpreprocessed(X_train,'project_resource_summary'))  
tfidf_w2v_vec_resource_te = tfidf_w2v(textpreprocessed(X_test,'project_resource_summary'))  
tfidf_w2v_vec_resource_val = tfidf_w2v(textpreprocessed(X_val,'project_resource_summary'))
```

100%|██████████| 61178/61178 [00:03<00:00, 15673.21it/s]
100%|██████████| 61178/61178 [00:04<00:00, 14329.76it/s]
7%|███████| 1528/21850 [00:00<00:01, 15277.57it/s]

61178
30

100%|██████████| 21850/21850 [00:01<00:00, 15461.22it/s]
100%|██████████| 21850/21850 [00:01<00:00, 14284.73it/s]
6%|███████| 1568/26220 [00:00<00:01, 15675.01it/s]

21850
30

100%|██████████| 26220/26220 [00:01<00:00, 15586.98it/s]
100%|██████████| 26220/26220 [00:01<00:00, 14346.76it/s]

26220
30

In [0]:

```
X_tr, X_cr, X_te = hstack_data(tfidf2v_vec_titles_tr, tfidf2v_vec_titles_val,
                               tfidf2v_vec_titles_te, \
                               tfidf2v_vec_essay_tr, tfidf2v_vec_essay_val, tfidf2v_vec_essay_te, \
                               tfidf2v_vec_resource_tr, tfidf2v_vec_resource_val, tfidf2v_vec_resource_te)
```

In [0]:

```
X_tr.shape,X_cr.shape,X_te.shape
```

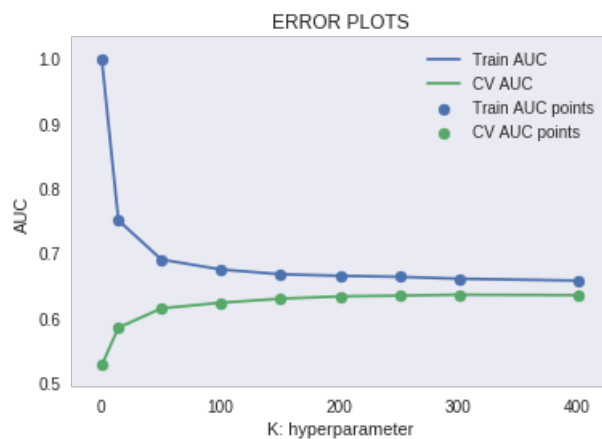
Out[0]:

```
((61178, 172), (26220, 172), (21850, 172))
```

In [0]:

```
best_k = optimised_k(X_tr,y_train,X_cr,y_val)
```

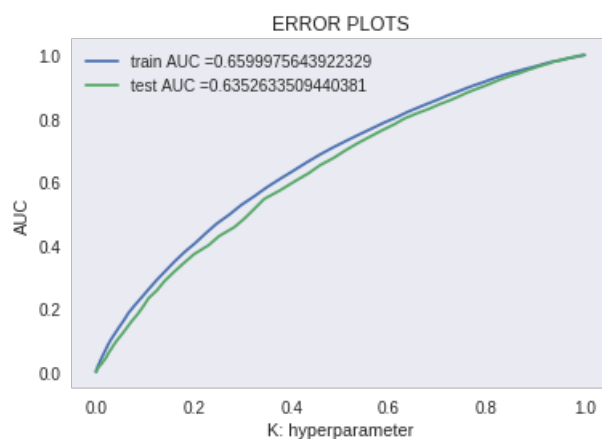
```
100%|██████████| 9/9 [6:49:47<00:00, 2730.41s/it]
```



execution time in minutes: 409
execution time in hours: 6.816666666666666
0.6353322729459712
optimal K is: 301

In [0]:

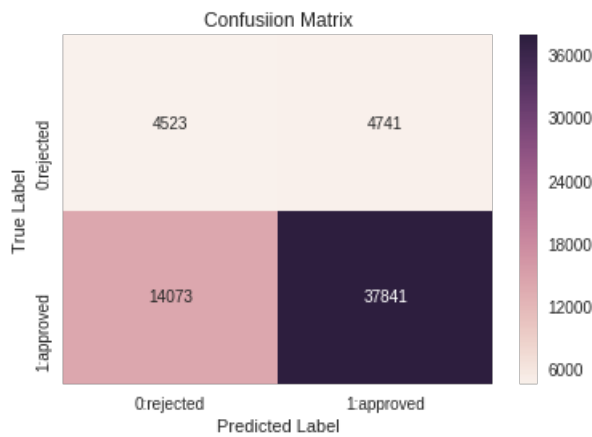
```
error_plot(X_tr,X_te,y_train,y_test)
```



```

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24986156181299424 for threshold 0.834
[[ 4523  4741]
 [14073 37841]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.24987489563975207 for threshold 0.841
[[ 1697  1611]
 [ 6012 12530]]
the maximum value of tpr*(1-fpr) 0.24986156181299424 for threshold 0.834

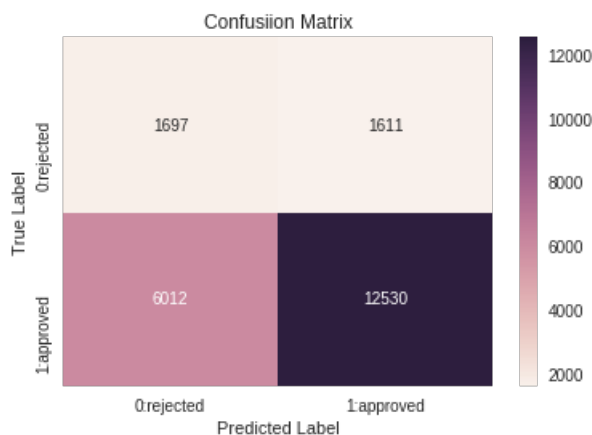
```



```

the maximum value of tpr*(1-fpr) 0.24987489563975207 for threshold 0.841

```



Summary

```
In [0]:
```

```

from prettytable import PrettyTable
import sys
sys.stdout.write("\033[1;30m")

x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Hyper parameter", "AUC"]

x.add_row(["BOW", "Brute", 301, 0.6627])
x.add_row(["TFIDF", "Brute", 401, 0.6378])
x.add_row(["W2V", "Brute", 401, 0.6366])
x.add_row(["TFIDFW2V", "Brute", 301, 0.6352])

print(x)

```

Vectorizer	Model	Hyper parameter	AUC
BOW	Brute	301	0.6627
TFIDF	Brute	401	0.6378
W2V	Brute	401	0.6366
TFIDFW2V	Brute	301	0.6352

