

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. Example: p036502
<code>project_title</code>	Title of the project. Examples: Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following enumerated values: Grades PreK-2 Grades 3-5 Grades 6-8 Grades 9-12
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: Applied Learning Care & Hunger Health & Sports History & Civics Literacy & Language Math & Science Music & The Arts Special Needs Warmth Examples: Music & The Arts Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. Examples: Literacy Literature & Writing, Social Sciences
<code>project_resource_summary</code>	An explanation of the resources needed for the project. Example: My students need hands on literacy materials to manage sensory needs!
<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. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> nan Dr. Mr. Mrs. Ms. Teacher.
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. Example: 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. Example: p036502
description	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
quantity	Quantity of the resource required. Example: 3
price	Price of the resource required. Example: 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 [69]:

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
```

```

import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

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

```

1.1 Reading Data

In [93]:

```

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

```

In [94]:

```

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

```

```

Number of data points in train data (109248, 17)
Number of data points in resource data (1541272, 4)
-----
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 [95]:

```

# 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 [95]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_:
55660	8393 p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2	
76127	37728 p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5	

In [96]:

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

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

Out [96]:

	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

Handling Missing Values

In [97]:

```
#https://stackoverflow.com/questions/29530232/how-to-check-if-any-value-is-nan-in-a-pandas-dataframe
me

project_data[project_data['teacher_prefix'].isnull()]
#Handle null values in pandas https://www.geeksforgeeks.org/python-pandas-dataframe-fillna-to-replace-null-values-in-dataframe/

project_data['teacher_prefix'].fillna( method = 'ffill', inplace = True)
```

1.2 preprocessing of project_subject_categories

In [98]:

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

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
cat_list = []
for i in catogories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the catogory 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 placeing 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 [99]:

```

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 [100]:

```

grades = list(project_data['project_grade_category'].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
grade_list = []
for i in grades:
    if 'Grades' in i.split(): # this will split each of the category based on space "Math & Science" => "Math", "&", "Science"
        i = i.replace('Grades', '') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
        i = i.replace(' ', '') # we are placing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
        grade_list.append(i.strip())

project_data['project_grade_category'] = grade_list

```

```

from collections import Counter
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]))

```

1.3 Text preprocessing

In [101]:

```

# 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 [102]:

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

Out[102]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_
55660	8393 p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	PreK-2	Enginee STEAM the Prin Classr
76127	37728 p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	3-5	Sen Tools Fo

In [103]:

```

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

```

I have been fortunate enough to use the Fairy Tale STEM kits in my classroom as well as the STEM journals, which my students really enjoyed. I would love to implement more of the Lakeshore STEM kits in my classroom for the next school year as they provide excellent and engaging STEM lessons. My students come from a variety of backgrounds, including language and socioeconomic status. Many of them don't have a lot of experience in science and engineering and these kits give me the materials to provide these exciting opportunities for my students. Each month I try to do several science or STEM/STEAM projects. I would use the kits and robot to help guide my science instruction in engaging and meaningful ways. I can adapt the kits to my current language arts pacing guide where we already teach some of the material in the kits like tall tales (Paul Bunyan) or Johnny Appleseed. The following units will be taught in the next school year where I will implement these kits: magnets, motion, sink vs. float, robots. I often get to these units and don't know if I am teaching the right way or using the right materials. The kits will give me additional ideas, strategies, and lessons to prepare my students in science. It is challenging to develop high quality science activities. These kits give me the materials I need to provide my students with science activities that will go along with the curriculum in my classroom. Although I have some things (like magnets) in my classroom, I don't know how to use them effectively. The kits will provide me with the right amount of materials and show me how to use them in an

kits will provide me with the right amount of materials and show me how to use them in an appropriate way.

I teach high school English to students with learning and behavioral disabilities. My students all vary in their ability level. However, the ultimate goal is to increase all students literacy level s. This includes their reading, writing, and communication levels.I teach a really dynamic group of students. However, my students face a lot of challenges. My students all live in poverty and in a dangerous neighborhood. Despite these challenges, I have students who have the the desire to defeat these challenges. My students all have learning disabilities and currently all are performing below grade level. My students are visual learners and will benefit from a classroom that fulfills their preferred learning style.The materials I am requesting will allow my students to be prepared for the classroom with the necessary supplies. Too often I am challenged with students who come to school unprepared for class due to economic challenges. I want my students to be able to focus on learning and not how they will be able to get school supplies. The supplies will last all year . Students will be able to complete written assignments and maintain a classroom journal. The chart paper will be used to make learning more visual in class and to create posters to aid students in their learning. The students have access to a classroom printer. The toner will be used to print student work that is completed on the classroom Chromebooks.I want to try and remove all barriers for the students learning and create opportunities for learning. One of the biggest barriers is the students not having the resources to get pens, paper, and folders. My students will be able to increase their literacy skills because of this project.

"Life moves pretty fast. If you don't stop and look around once in awhile, you could miss it.\" from the movie, Ferris Bueller's Day Off. Think back...what do you remember about your grandparents? How amazing would it be to be able to flip through a book to see a day in their lives?My second graders are voracious readers! They love to read both fiction and nonfiction books . Their favorite characters include Pete the Cat, Fly Guy, Piggie and Elephant, and Mercy Watson. They also love to read about insects, space and plants. My students are hungry bookworms! My students are eager to learn and read about the world around them. My kids love to be at school and are like little sponges absorbing everything around them. Their parents work long hours and usually do not see their children. My students are usually cared for by their grandparents or a family friend. Most of my students do not have someone who speaks English at home. Thus it is difficult for my students to acquire language.Now think forward... wouldn't it mean a lot to your kids, nieces or nephews or grandchildren, to be able to see a day in your life today 30 years from now? Memories are so precious to us and being able to share these memories with future generations will be a rewarding experience. As part of our social studies curriculum, students will be learning about changes over time. Students will be studying photos to learn about how their community has changed over time. In particular, we will look at photos to study how the land, buildings, clothing, and schools have changed over time. As a culminating activity, my students will capture a slice of their history and preserve it through scrap booking. Key important events in their young lives will be documented with the date, location, and names. Students will be using photos from home and from school to create their second grade memories. Their scrap books will preserve their unique stories for future generations to enjoy.Your donation to this project will provide my second graders with an opportunity to learn about social studies in a fun and creative manner. Through their scrapbooks, children will share their story with others and have a historical document for the rest of their lives.

"A person's a person, no matter how small.\" (Dr.Seuss) I teach the smallest students with the biggest enthusiasm for learning. My students learn in many different ways using all of our senses and multiple intelligences. I use a wide range of techniques to help all my students succeed. \r\nStudents in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures, including Native Americans.\r\nOur school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom. Kindergarteners in my class love to work with hands-on materials and have many different opportunities to practice a skill before it is mastered. Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum.Montana is the perfect place to learn about agriculture and nutrition. My students love to role play in our pretend kitchen in the early childhood classroom. I have had several kids ask me, \"Can we try cooking with REAL food?\" I will take their idea and create \"Common Core Cooking Lessons\" where we learn important math and writing concepts while cooking delicious healthy food for snack time. My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it's healthy for their bodies. This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce, make our own bread, and mix up healthy plants from our classroom garden in the spring. We will also create our own cookbooks to be printed and shared with families. \r\nStudents will gain math and literature skills as well as a life long enjoyment for healthy cooking.nannan

My classroom consists of twenty-two amazing sixth graders from different cultures and backgrounds. They are a social bunch who enjoy working in partners and working with groups. They are hard-working and eager to head to middle school next year. My job is to get them ready to make this transition and make it as smooth as possible. In order to do this, my students need to come to school every day and feel safe and ready to learn. Because they are getting ready to head to middle school, I give them lots of choice- choice on where to sit and work, the order to complete assignments, choice of projects, etc. Part of the students feeling safe is the ability for them to come into a welcoming, encouraging environment. My room is colorful and the atmosphere is casual. I want them to take ownership of the classroom because we ALL share it together. Because my time with them is limited I want to ensure they get the most of this time and enjoy it to the best of t

their abilities. Currently, we have twenty-two desks of differing sizes, yet the desks are similar to the ones the students will use in middle school. We also have a kidney table with crates for seating. I allow my students to choose their own spots while they are working independently or in groups. More often than not, most of them move out of their desks and onto the crates. Believe it or not, this has proven to be more successful than making them stay at their desks! It is because of this that I am looking toward the "Flexible Seating" option for my classroom. The students look forward to their work time so they can move around the room. I would like to get rid of the constricting desks and move toward more "fun" seating options. I am requesting various seating so my students have more options to sit. Currently, I have a stool and a papasan chair I inherited from the previous sixth-grade teacher as well as five milk crate seats I made, but I would like to give them more options and reduce the competition for the "good seats". I am also requesting two rugs as not only more seating options but to make the classroom more welcoming and appealing. In order for my students to be able to write and complete work without desks, I am requesting a class set of clipboards. Finally, due to curriculum that requires groups to work together, I am requesting tables that we can fold up when we are not using them to leave more room for our flexible seating options. I know that with more seating options, they will be that much more excited about coming to school! Thank you for your support in making my classroom one students will remember forever!nannan

In [104]:

```
# 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 [105]:

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

"A person is a person, no matter how small." (Dr.Seuss) I teach the smallest students with the biggest enthusiasm for learning. My students learn in many different ways using all of our senses and multiple intelligences. I use a wide range of techniques to help all my students succeed. Students in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures, including Native Americans. Our school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom. Kindergarteners in my class love to work with hands-on materials and have many different opportunities to practice a skill before it is mastered. Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum. Montana is the perfect place to learn about agriculture and nutrition. My students love to role play in our pretend kitchen in the early childhood classroom. I have had several kids ask me, "Can we try cooking with REAL food?" I will take their idea and create "Common Core Cooking Lessons" where we learn important math and writing concepts while cooking delicious healthy food for snack time. My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it is healthy for their bodies. This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce, make our own bread, and mix up healthy plants from our classroom garden in the spring. We will also create our own cookbooks to be printed and shared with families. Students will gain math and literature skills as well as a life long enjoyment for healthy cooking.nannan

In [106]:

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



```
sent = sent.replace("\\\\", ' ')
sent = sent.replace("\\n", ' ')
print(sent)
```

A person is a person, no matter how small. (Dr.Seuss) I teach the smallest students with the biggest enthusiasm for learning. My students learn in many different ways using all of our senses and multiple intelligences. I use a wide range of techniques to help all my students succeed. Students in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures, including Native Americans. Our school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom. Kindergarteners in my class love to work with hands-on materials and have many different opportunities to practice a skill before it is mastered. Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum. Montana is the perfect place to learn about agriculture and nutrition. My students love to role play in our pretend kitchen in the early childhood classroom. I have had several kids ask me, Can we try cooking with REAL food? I will take their idea and create Common Core Cooking Lessons where we learn important math and writing concepts while cooking delicious healthy food for snack time. My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it is healthy for their bodies. This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce, make our own bread, and mix up healthy plants from our classroom garden in the spring. We will also create our own cookbooks to be printed and shared with families. Students will gain math and literature skills as well as a life long enjoyment for healthy cooking.annan

In [107]:

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

A person is a person no matter how small Dr Seuss I teach the smallest students with the biggest enthusiasm for learning My students learn in many different ways using all of our senses and multiple intelligences I use a wide range of techniques to help all my students succeed Students in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures including Native Americans Our school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom Kindergarteners in my class love to work with hands on materials and have many different opportunities to practice a skill before it is mastered Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum Montana is the perfect place to learn about agriculture and nutrition My students love to role play in our pretend kitchen in the early childhood classroom I have had several kids ask me Can we try cooking with REAL food I will take their idea and create Common Core Cooking Lessons where we learn important math and writing concepts while cooking delicious healthy food for snack time My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it is healthy for their bodies This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce make our own bread and mix up healthy plants from our classroom garden in the spring We will also create our own cookbooks to be printed and shared with families Students will gain math and literature skills as well as a life long enjoyment for healthy cooking.annan

In [108]:

```
# 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
```

◀ ▶

1.5 Preparing data for models

In [114]:

```
project_data.columns
```

Out[114]:

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

we are going to consider

- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data
- project_title : text data
- text : text data
- project_resource_summary: text data (optinal)
- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical

In [115]:

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

Out[115]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',  
      'Date', 'project_grade_category', 'project_title',  
      'project_resource_summary',  
      'teacher_number_of_previously_posted_projects', 'project_is_approved',  
      'clean_categories', 'clean_subcategories', 'essay', 'price',  
      'quantity'],  
      dtype='object')
```

In [116]:

```
# move columns in pandas dataframe https://stackoverflow.com/questions/35321812/move-column-in-pandas-dataframe/35321983  
project_data = project_data[[c for c in project_data if c not in ['project_is_approved']]  
                             + ['project_is_approved']]  
project_data.columns
```

Out[116]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',  
      'Date', 'project_grade_category', 'project_title',  
      'project_resource_summary',  
      'teacher_number_of_previously_posted_projects', 'clean_categories',  
      'clean_subcategories', 'essay', 'price', 'quantity',  
      'project_is_approved'],  
      dtype='object')
```

Selecting 50000 data points

In [117]:

```
project_data = project_data.iloc[0:50000 , :]
```

Due to memory issues I selected first 50K points

In [118]:

```
project_data.shape
```

Out[118]:

```
(50000, 16)
```

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 your train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](#).

2. K Nearest Neighbour

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

In [119]:

```
#importing necessary modules

from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_val_score
from collections import Counter
from sklearn.metrics import accuracy_score

# Splitting the data into X , Y labels

# create design matrix X and target vector y
X = np.array(project_data.iloc[:, :-1]) # end index is exclusive
y = np.array(project_data['project_is_approved']) # showing you two ways of indexing a pandas df
# split the data set into train and test
X_1, X_test, y_1, y_test = train_test_split(X, y, test_size=0.3, random_state=0, stratify = y)

# split the train data set into cross validation train and cross validation test
X_tr, X_cv, y_tr, y_cv = train_test_split(X_1, y_1, test_size=0.3, stratify = y_1)
```

In [120]:

```
print(len(X_tr))
print(len(X_cv))
print(len(X_test))
```

```
24500
10500
15000
```

In [121]:

```
X_tr = pd.DataFrame(data=X_tr[0:,0:], columns=project_data.columns[0:-1])
X_cv = pd.DataFrame(data=X_cv[0:,0:], columns=project_data.columns[0:-1])
X_test = pd.DataFrame(data=X_test[0:,0:], columns=project_data.columns[0:-1])
```

2.2 Make Data Model Ready: encoding numerical, categorical features

In [122]:

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

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if neededtHuWDX6yizwIhai
# c. X-axis label
# d. Y-axis label
print("="*25+"encoding categorical features"+"="*25)
#Vectorizing categorical data :
# 1 Clean_Categories

# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary=list(sorted(set(dict_keys(A))), lowercase=False, binary=True))
```

```

vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_tr['clean_categories'].values)

categories_one_hot_train = vectorizer.transform(X_tr['clean_categories'].values)
categories_one_hot_test = vectorizer.transform(X_test['clean_categories'].values)
categories_one_hot_cv = vectorizer.transform(X_cv['clean_categories'].values)

print(vectorizer.get_feature_names())
print("Shape of train matrix after one hot encoding ",categories_one_hot_train.shape)
print("Shape of test matrix after one hot encoding ",categories_one_hot_test.shape)
print("Shape of cv matrix after one hot encoding ",categories_one_hot_cv.shape)
print("="*100)

```

```

=====encoding categorical features=====
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of train matrix after one hot encoding (24500, 9)
Shape of test matrix after one hot encoding (15000, 9)
Shape of cv matrix after one hot encoding (10500, 9)
=====

```

In [123]:

```

# 2 clean_subcategories

vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_tr['clean_subcategories'].values)
print(vectorizer.get_feature_names())

sub_categories_one_hot_train = vectorizer.transform(X_tr['clean_subcategories'].values)
sub_categories_one_hot_test = vectorizer.transform(X_test['clean_subcategories'].values)
sub_categories_one_hot_cv = vectorizer.transform(X_cv['clean_subcategories'].values)

print("Shape of train matrix after one hot encoding ",sub_categories_one_hot_train.shape)
print("Shape of test matrix after one hot encoding ",sub_categories_one_hot_test.shape)
print("Shape of cv matrix after one hot encoding ",sub_categories_one_hot_cv.shape)
print("="*100)

```

```

['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular',
'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger',
'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other',
'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL',
'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences',
'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
Shape of train matrix after one hot encoding (24500, 30)
Shape of test matrix after one hot encoding (15000, 30)
Shape of cv matrix after one hot encoding (10500, 30)
=====

```

In [124]:

```

my_counter = Counter()
for state in project_data['school_state'].values:
    my_counter.update(state.split())
school_state_cat_dict = dict(my_counter)
sorted_school_state_cat_dict = dict(sorted(school_state_cat_dict.items(), key=lambda kv: kv[1]))

```

In [125]:

```

# 3 school_state

vectorizer = CountVectorizer(vocabulary=list(sorted_school_state_cat_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_tr['school_state'].values)
print(vectorizer.get_feature_names())
school_state_one_hot_train = vectorizer.transform(X_tr['school_state'].values)
school_state_one_hot_test = vectorizer.transform(X_test['school_state'].values)
school_state_one_hot_cv = vectorizer.transform(X_cv['school_state'].values)

```

```
print("Shape of train matrix after one hot encoding ",school_state_one_hot_train.shape)
print("Shape of test matrix after one hot encoding ",school_state_one_hot_test.shape)
print("Shape of cv matrix after one hot encoding ",school_state_one_hot_cv.shape)
print("=="*100)
```

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

```
Shape of train matrix after one hot encoding (24500, 51)
```

```
Shape of test matrix after one hot encoding (15000, 51)
```

```
Shape of cv matrix after one hot encoding (10500, 51)
```

```
=====
```

In [126]:

```
my_counter = Counter()
for teacher in project_data['teacher_prefix'].values:
    my_counter.update(teacher.split())
teacher_prefix_cat_dict = dict(my_counter)
sorted_teacher_prefix_cat_dict = dict(sorted(teacher_prefix_cat_dict.items(), key=lambda kv: kv[1])
)
```

In [127]:

```
# 4 teacher_prefix

#one hot encoding for teacher_prefix feature

vectorizer = CountVectorizer(vocabulary=list(sorted_teacher_prefix_cat_dict.keys()),lowercase=False, binary=True)
vectorizer.fit(X_tr['teacher_prefix'].values)
print(vectorizer.get_feature_names())
teacher_prefix_one_hot_train = vectorizer.transform(X_tr['teacher_prefix'].values)
teacher_prefix_one_hot_test = vectorizer.transform(X_test['teacher_prefix'].values)
teacher_prefix_one_hot_cv = vectorizer.transform(X_cv['teacher_prefix'].values)

print("Shape of train matrix after one hot encoding ",teacher_prefix_one_hot_train.shape)
print("Shape of test matrix after one hot encoding ",teacher_prefix_one_hot_test.shape)
print("Shape of cv matrix after one hot encoding ",teacher_prefix_one_hot_cv.shape)
print("=="*100)
```

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

```
Shape of train matrix after one hot encoding (24500, 5)
```

```
Shape of test matrix after one hot encoding (15000, 5)
```

```
Shape of cv matrix after one hot encoding (10500, 5)
```

```
=====
```

In [128]:

```
print(teacher_prefix_one_hot_train.toarray()[0:5,:])
```

```
[[0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]]
```

In [129]:

```
# 5 project_grade_category

vectorizer = CountVectorizer(vocabulary=list(sorted_grade_dict.keys()),lowercase=False, binary=True)
vectorizer.fit(X_tr['project_grade_category'].values)
print(vectorizer.get_feature_names())
project_grade_one_hot_train = vectorizer.transform(X_tr['project_grade_category'].values)
project_grade_one_hot_test = vectorizer.transform(X_test['project_grade_category'].values)
project_grade_one_hot_cv = vectorizer.transform(X_cv['project_grade_category'].values)
```

```

print("Shape of train matrix after one hot encodig ",project_grade_one_hot_train.shape)
print("Shape of test matrix after one hot encodig ",project_grade_one_hot_test.shape)
print("Shape of cv matrix after one hot encodig ",project_grade_one_hot_cv.shape)
print("="*100)

```

```

['9-12', '6-8', '3-5', 'PreK-2']
Shape of train matrix after one hot encodig (24500, 4)
Shape of test matrix after one hot encodig (15000, 4)
Shape of cv matrix after one hot encodig (10500, 4)
=====

```

In [131]:

```

print("_"*25+"encoding numerical features"+"_"*25)
# check this one: https://www.youtube.com/watch?v=0HOqOcln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html
from sklearn.preprocessing import StandardScaler

# price_standardized = standardScaler.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.
73 5.5 ].
# Reshape your data either using array.reshape(-1, 1)

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

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

print("Shape of train matrix after one hot encodig ",price_standardized_train.shape)
print("Shape of test matrix after one hot encodig ",price_standardized_test.shape)
print("Shape of cv matrix after one hot encodig ",price_standardized_cv.shape)
print("="*100)
#2 teacher_number_of_previously_posted_projects

previous_project_scalar = StandardScaler()
previous_project_scalar.fit(X_tr['teacher_number_of_previously_posted_projects'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
print(f"Mean : {previous_project_scalar.mean_[0]}, Standard deviation :
{np.sqrt(previous_project_scalar.var_[0])}")

# Now standardize the data with above maen and variance.
previous_project_standardized_train =
previous_project_scalar.transform(X_tr['teacher_number_of_previously_posted_projects'].values.res
hape(-1, 1))
previous_project_standardized_test =
previous_project_scalar.transform(X_test['teacher_number_of_previously_posted_projects'].values.re
shape(-1, 1))
previous_project_standardized_cv =
previous_project_scalar.transform(X_cv['teacher_number_of_previously_posted_projects'].values.res
hape(-1, 1))

print("Shape of train matrix after one hot encodig ",previous_project_standardized_train.shape)
print("Shape of test matrix after one hot encodig ",previous_project_standardized_test.shape)
print("Shape of cv matrix after one hot encodig ",previous_project_standardized_cv.shape)

```

```

encoding numerical features
Mean : 312.27854857142853, Standard deviation : 371.80026210042183
Shape of train matrix after one hot encodig (24500, 1)
Shape of test matrix after one hot encodig (15000, 1)
Shape of cv matrix after one hot encodig (10500, 1)
=====

Mean : 9.473632653061225, Standard deviation : 23.973543212274038
Shape of train matrix after one hot encodig (24500, 1)

```



```
Shape of test matrix after one hot encoding (15000, 1)
Shape of cv matrix after one hot encoding (10500, 1)
```

2.3 Make Data Model Ready: encoding essay, and project_title

Encoding Essay and Project_title columns using Bag Of Words

In [132]:

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

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
print("_"*25+"Essay BOW"+"_"*25)
vectorizer = CountVectorizer(min_df=10)
vectorizer.fit(X_tr['essay'])
essay_bow_train = vectorizer.transform(X_tr['essay'])
essay_bow_test = vectorizer.transform(X_test['essay'])
essay_bow_cv = vectorizer.transform(X_cv['essay'])
print("Shape of train matrix after one hot encoding ",essay_bow_train.shape)
print("Shape of test matrix after one hot encoding ",essay_bow_test.shape)
print("Shape of cv matrix after one hot encoding ",essay_bow_cv.shape)
print("="*100)
print("_"*25+"Project_Title BOW"+"_"*25)
vectorizer = CountVectorizer(min_df=10)
vectorizer.fit(X_tr['project_title'])
title_bow_train = vectorizer.transform(X_tr['project_title'])
title_bow_test = vectorizer.transform(X_test['project_title'])
title_bow_cv = vectorizer.transform(X_cv['project_title'])
print("Shape of train matrix after one hot encoding ",title_bow_train.shape)
print("Shape of test matrix after one hot encoding ",title_bow_test.shape)
print("Shape of cv matrix after one hot encoding ",title_bow_cv.shape)
```

Essay BOW

```
Shape of train matrix after one hot encoding (24500, 9188)
Shape of test matrix after one hot encoding (15000, 9188)
Shape of cv matrix after one hot encoding (10500, 9188)
=====
```

Project_Title BOW

```
Shape of train matrix after one hot encoding (24500, 1258)
Shape of test matrix after one hot encoding (15000, 1258)
Shape of cv matrix after one hot encoding (10500, 1258)
```

Encoding Essay and Project_title columns using TFIDF

In [133]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
print("_"*25+"Essay TFIDF"+"_"*25)
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(X_tr['essay'])
essay_tfidf_train = vectorizer.transform(X_tr['essay'])
essay_tfidf_test = vectorizer.transform(X_test['essay'])
essay_tfidf_cv = vectorizer.transform(X_cv['essay'])
print("Shape of train matrix after one hot encoding ",essay_tfidf_train.shape)
print("Shape of test matrix after one hot encoding ",essay_tfidf_test.shape)
print("Shape of cv matrix after one hot encoding ",essay_tfidf_cv.shape)
print("="*100)
print("_"*25+"Project_Title TFIDF"+"_"*25)
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit_transform(X_tr['project_title'])
```

```

      Essay_TFIDF
Shape of train matrix after one hot encodig (24500, 9188)
Shape of test matrix after one hot encodig (15000, 9188)
Shape of cv matrix after one hot encodig (10500, 9188)
=====
      Project_Title_TFIDF
Shape of train matrix after one hot encodig (24500, 1258)
Shape of test matrix after one hot encodig (15000, 1258)
Shape of cv matrix after one hot encodig (10500, 1258)

```

```
# 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())
```

```

words = []
for i in X_tr['essay']:
    words.extend(i.split(' '))

for i in X_tr['project_title']:
    words.extend(i.split(' '))
print("all the words in the coupus", len(words))
words = set(words)
print("the unique words in the coupus", len(words))

inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove vectors and our coupus", \
      len(inter_words), "(" , np.round(len(inter_words)/len(words)*100,3), "%) ")

words_courpus = {}
words_glove = set(model.keys())
for i in words:
    if i in words_glove:
        words_courpus[i] = model[i]
print("word 2 vec length", len(words_courpus))

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/

#import pickle
#with open('./My Drive/glove_vectors', 'wb') as f:
#    pickle.dump(words_courpus, f)

```

Encoding Essay column using AVG WORD2VEC

```
# Similarly you can vectorize for title also
print("_"*25+"Essay AVG_W2V"+"_"*25)
essay_avg_w2v_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_tr['essay']): # for each review/sentence
```

Essay AVG W2V

24500
300

Essay AVG W2V

15000
300

Essay AVG W2V

[illegible]

10500
300

In [139]:

Title AVG W2V

24500
300

In [140]:

Title AVG W2V

15000
300

In [141]:

```
print("_"*25+"Title AVG_W2V"+"_"*25)
title_avg_w2v_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['project_title']): # 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:
```

_____Title AVG_W2V_____

10500
300

In [142]:

24500
300

In [143]:

```
essay_tfidf_w2v_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['essay']): # 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
    essay_tfidf_w2v_test.append(vector)
```

```
print(len(essay_tfidf_w2v_test))
print(len(essay_tfidf_w2v_test[0]))
```

```
100%|███████████████████████████████████████████████████████████| 15000/15000 [00:  
58<00:00, 254.91it/s]
```

15000
300

In [144]:

```
essay_tfidf_w2v_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['essay']): # 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
    essay_tfidf_w2v_cv.append(vector)

print(len(essay_tfidf_w2v_cv))
print(len(essay_tfidf_w2v_cv[0]))
```

[illegible]

10500
300

Encoding Project title column using TFIDF WORD2VEC

In [145]:

```
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_tr['project_title'])
# 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())
title_tfidf_w2v_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_tr['project_title']): # 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
    title_tfidf_w2v_train.append(vector)

print(len(title_tfidf_w2v_train))
print(len(title_tfidf_w2v_train[0]))
```

```
100%|██████████████████████████████████████████████████████████████████████████| 24500/24500  
[00:01<00:00, 19782.46it/s]
```

24500
300

In [146]:

```
title_tfidf_w2v_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['project_title']): # 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
    title_tfidf_w2v_test.append(vector)

print(len(title_tfidf_w2v_test))
print(len(title_tfidf_w2v_test[0]))
```

100% |██| 15000/15000
[00:00<00:00, 17260.61it/s]

15000
300

In [147]:

```
title_tfidf_w2v_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['project_title']): # 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
    title_tfidf_w2v_cv.append(vector)

print(len(title_tfidf_w2v_cv))
print(len(title_tfidf_w2v_cv[0]))
```

100% |██| 10500/10500
[00:00<00:00, 18085.94it/s]

10500
300

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

In [148]:

```
# please write all the code with proper documentation, and proper titles for each subsection
```

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

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

2.4.1 Applying KNN brute force on BOW, SET 1

SET 1 : Combining all the features

In [149]:

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

from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_tr = hstack((school_state_one_hot_train, categories_one_hot_train, sub_categories_one_hot_train,
teacher_prefix_one_hot_train, project_grade_one_hot_train, essay_bow_train, title_bow_train,
price_standardized_train, previous_project_standardized_train))
X_cv = hstack((school_state_one_hot_cv, categories_one_hot_cv, sub_categories_one_hot_cv,
teacher_prefix_one_hot_cv, project_grade_one_hot_cv, essay_bow_cv, title_bow_cv,
price_standardized_cv, previous_project_standardized_cv))
X_test = hstack((school_state_one_hot_test, categories_one_hot_test, sub_categories_one_hot_test, t
eacher_prefix_one_hot_test, project_grade_one_hot_test, essay_bow_test, title_bow_test,
price_standardized_test, previous_project_standardized_test))
X_tr = X_tr.tocsr()
X_cv = X_cv.tocsr()
X_test = X_test.tocsr()
```

In [150]:

```
print(X_tr.shape , y_tr.shape)
print(X_cv.shape , y_cv.shape)
print(X_test.shape , y_test.shape)
```

```
(24500, 10547) (24500,)
(10500, 10547) (10500,)
(15000, 10547) (15000,)
```

In [151]:

```
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 tr_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
    if data.shape[0]%1000 != 0:
        y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])

    return y_data_pred
```

Applying KNN on train data and finding best hyperparameter with CV data

In [152]:

```
import matplotlib.pyplot as plt
```



```

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 = [3, 15, 25, 51, 101, 130]
for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
    neigh.fit(X_tr, y_tr)

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

    # 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_tr, y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

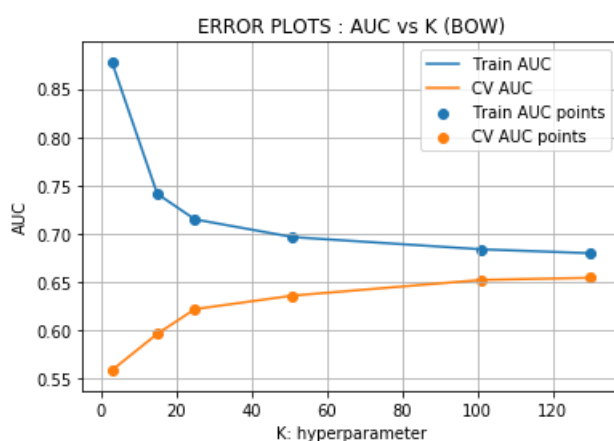
plt.plot(K, train_auc, label='Train AUC')
plt.plot(K, cv_auc, label='CV AUC')

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 : AUC vs K (BOW)")
plt.grid()
plt.show()

```

100% | 6/6 [09:40<00:00, 96.17s/it]



Best Hyperparameter K

In [153]:

```
best_k1 = 130
```

Training model with best K and finding the train and test AUC

In [154]:

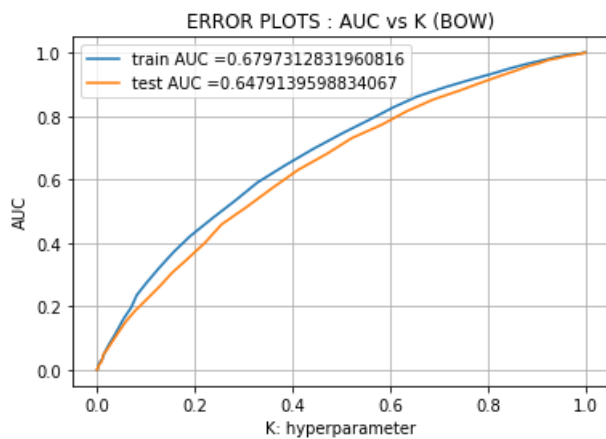
```
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k1, n_jobs=-1)
neigh.fit(X_tr, y_tr)
# 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_test)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_tr, 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 : AUC vs K (BOW)")
plt.grid()
plt.show()
```



In [155]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    return t

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

Confusion Matrix

In [156]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t1 = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_tr, predict_with_best_t(y_train_pred, best_t1)))
```

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t1)))
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.3967667015127465 for threshold 0.777

Train confusion matrix

```
[[ 2635  1288]
 [ 8422 12155]]
```

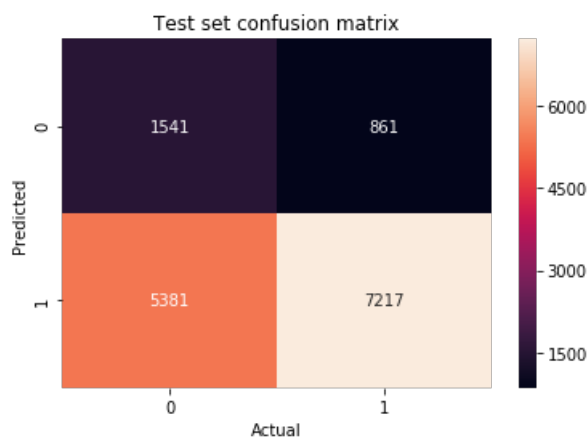
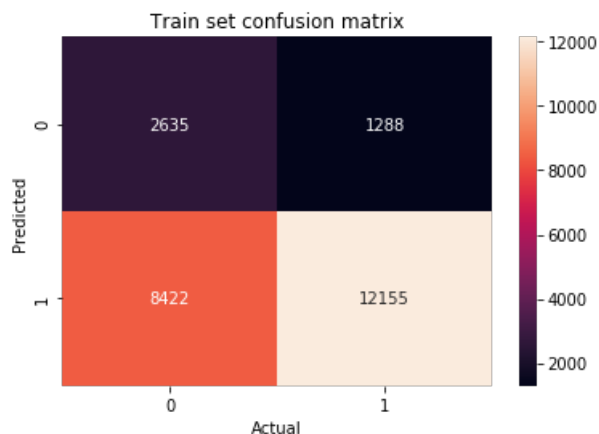
Test confusion matrix

```
[[1541   861]
 [5381  7217]]
```

In [159]:

```
# Heatmap for train set confusion matrix(Select K best)
heatmap_train = sns.heatmap(confusion_matrix(y_tr, predict_with_best_t(y_train_pred,best_t1)),
annot=True, fmt="d")
plt.title("Train set confusion matrix")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.show()

# Heatmap for test set confusion matrix(Select K best)
heatmap_train = sns.heatmap(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t1)), an
not=True, fmt="d")
plt.title("Test set confusion matrix")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.show()
```



2.4.2 Applying KNN brute force on TFIDF, SET 2

SET 2 : Combining all the features

In [160]:

In [160]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_tr = hstack((school_state_one_hot_train, categories_one_hot_train, sub_categories_one_hot_train,
teacher_prefix_one_hot_train, project_grade_one_hot_train, essay_tfidf_train, title_tfidf_train, price_standardized_train, previous_project_standardized_train))
X_cv = hstack((school_state_one_hot_cv, categories_one_hot_cv, sub_categories_one_hot_cv,
teacher_prefix_one_hot_cv, project_grade_one_hot_cv, essay_tfidf_cv, title_tfidf_cv, price_standardized_cv, previous_project_standardized_cv))
X_test = hstack((school_state_one_hot_test, categories_one_hot_test, sub_categories_one_hot_test, teacher_prefix_one_hot_test, project_grade_one_hot_test, essay_tfidf_test, title_tfidf_test, price_standardized_test, previous_project_standardized_test))
X_tr = X_tr.tocsr()
X_cv = X_cv.tocsr()
X_test = X_test.tocsr()
```

In [161]:

```
print(X_tr.shape , y_tr.shape)
print(X_cv.shape , y_cv.shape)
print(X_test.shape , y_test.shape)
```

```
(24500, 10547) (24500,)
(10500, 10547) (10500,)
(15000, 10547) (15000,)
```

Applying KNN on train data and finding best hyperparameter with CV data

In [162]:

```
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 = [3, 15, 25, 51, 101, 130]
for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
    neigh.fit(X_tr, y_tr)

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

    # 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_tr, 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 : AUC vs K (TFIDF)")
plt.grid()
plt.show()
```

[illegible]

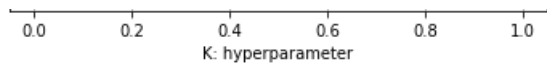
```
best_k2 = 130
```

In [164]:

ERROR PLOTS : AUC vs K (TFIDF)

train AUC = 0.6434399191780057
test AUC = 0.6130994782751686

K	train AUC	test AUC
1	0.05	0.05
2	0.15	0.15
3	0.25	0.25
4	0.35	0.35
5	0.45	0.45
6	0.55	0.55
7	0.65	0.65
8	0.75	0.75
9	0.85	0.85
10	0.95	0.95



Confusion Matrix

In [165]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t2 = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_tr, predict_with_best_t(y_train_pred, best_t2)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t2)))
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.36372127293526196 for threshold 0.846

Train confusion matrix

```
[[ 2484  1439]
 [ 8757 11820]]
```

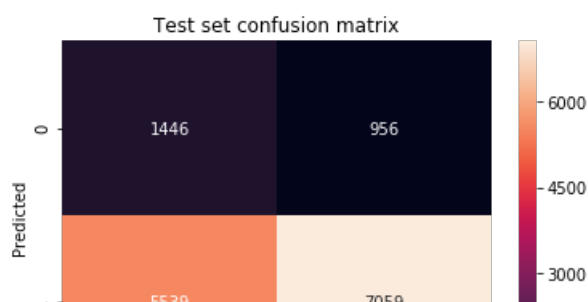
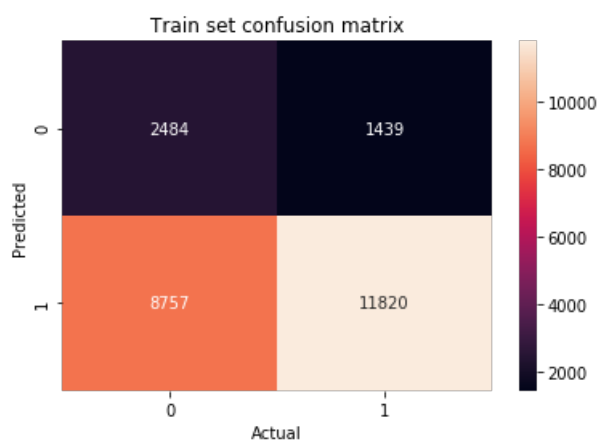
Test confusion matrix

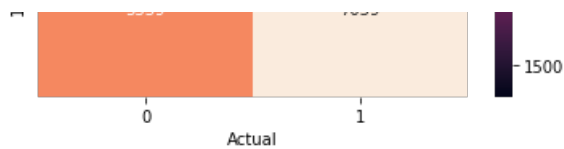
```
[[1446  956]
 [5539 7059]]
```

In [166]:

```
# Heatmap for train set confusion matrix(Select K best)
heatmap_train = sns.heatmap(confusion_matrix(y_tr, predict_with_best_t(y_train_pred,best_t2)),
annot=True, fmt="d")
plt.title("Train set confusion matrix")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.show()

# Heatmap for test set confusion matrix(Select K best)
heatmap_train = sns.heatmap(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t2)), an
not=True, fmt="d")
plt.title("Test set confusion matrix")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.show()
```





2.4.3 Applying KNN brute force on AVG W2V, SET 3

SET 3 : Combining all the features

In [167]:

```
# Set 3: categorical, numerical features + project_title(AVG W2V)+ preprocessed_essay (AVG W2V)
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_tr = hstack((school_state_one_hot_train, categories_one_hot_train, sub_categories_one_hot_train,
teacher_prefix_one_hot_train, project_grade_one_hot_train, essay_avg_w2v_train, title_avg_w2v_train,
price_standardized_train, previous_project_standardized_train))
X_cv = hstack((school_state_one_hot_cv, categories_one_hot_cv, sub_categories_one_hot_cv,
teacher_prefix_one_hot_cv, project_grade_one_hot_cv, essay_avg_w2v_cv, title_avg_w2v_cv,
price_standardized_cv, previous_project_standardized_cv))
X_test = hstack((school_state_one_hot_test, categories_one_hot_test, sub_categories_one_hot_test, t
eacher_prefix_one_hot_test, project_grade_one_hot_test, essay_avg_w2v_test, title_avg_w2v_test,
price_standardized_test, previous_project_standardized_test))
X_tr = X_tr.tocsr()
X_cv = X_cv.tocsr()
X_test = X_test.tocsr()
```

In [168]:

```
print(X_tr.shape , y_tr.shape)
print(X_cv.shape , y_cv.shape)
print(X_test.shape , y_test.shape)
```

```
(24500, 701) (24500,)
(10500, 701) (10500,)
(15000, 701) (15000,)
```

Applying KNN on train data and finding best hyperparameter with CV data

In [169]:

```
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 = [3, 15, 25, 51, 101, 130]
for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
    neigh.fit(X_tr, y_tr)

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

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

```

five class
# not the predicted outputs
train_auc.append(roc_auc_score(y_tr,y_train_pred))
cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

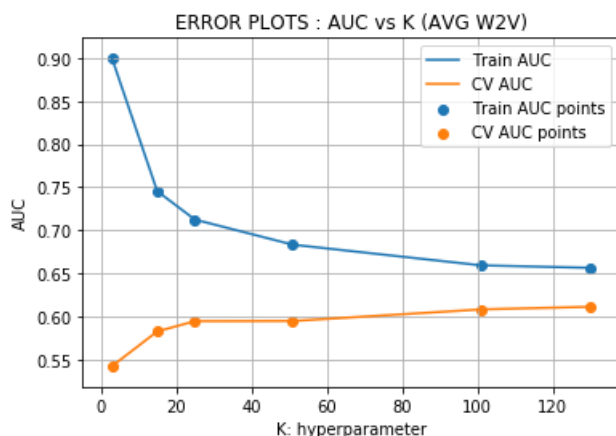
plt.plot(K, train_auc, label='Train AUC')
plt.plot(K, cv_auc, label='CV AUC')

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 : AUC vs K (AVG W2V)")
plt.grid()
plt.show()

```

100% | 6/6 [59:59<00:00, 594.49s/it]



Best Hyperparameter K

In [170]:

```
best_k3 = 130
```

Training model with best K and finding the train and test AUC

In [171]:

```

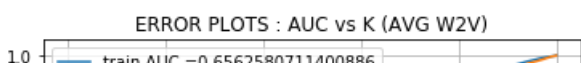
from sklearn.metrics import roc_curve, auc

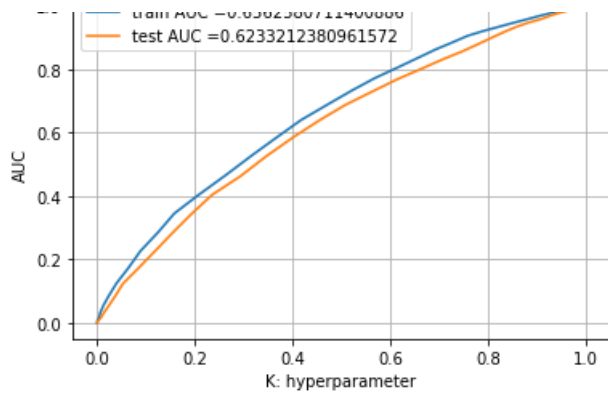
neigh = KNeighborsClassifier(n_neighbors=best_k3, n_jobs=-1)
neigh.fit(X_tr, y_tr)
y_train_pred = batch_predict(neigh, X_tr)
y_test_pred = batch_predict(neigh, X_test)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_tr, 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 : AUC vs K (AVG W2V)")
plt.grid()
plt.show()

```





Confusion Matrix

In [172]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t3 = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_tr, predict_with_best_t(y_train_pred, best_t3)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t3)))
```

=====

the maximum value of $tpr \cdot (1 - fpr)$ 0.3723568918922083 for threshold 0.846

Train confusion matrix

```
[[ 2283  1640]
 [ 7411 13166]]
```

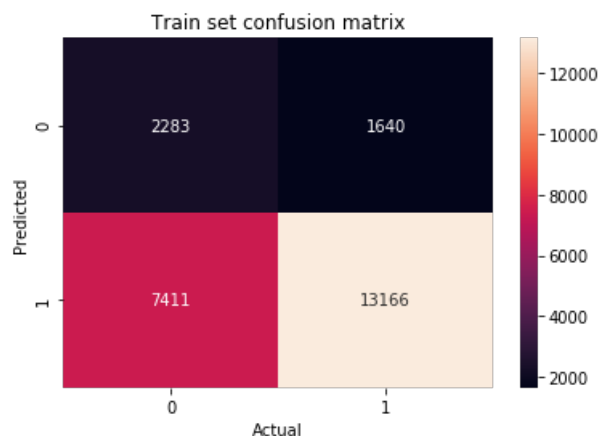
Test confusion matrix

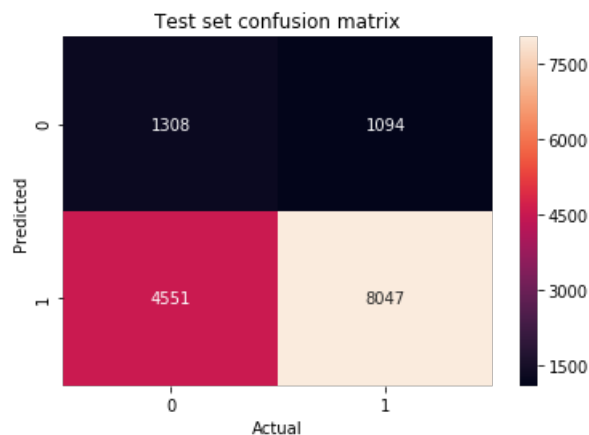
```
[[1308 1094]
 [4551 8047]]
```

In [173]:

```
# Heatmap for train set confusion matrix(Select K best)
heatmap_train = sns.heatmap(confusion_matrix(y_tr, predict_with_best_t(y_train_pred, best_t3)),
annot=True, fmt="d")
plt.title("Train set confusion matrix")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.show()

# Heatmap for test set confusion matrix(Select K best)
heatmap_test = sns.heatmap(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t3)),
annot=True, fmt="d")
plt.title("Test set confusion matrix")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.show()
```





2.4.4 Applying KNN brute force on TFIDF W2V, SET 4

SET 4 : Combining all the features

In [176]:

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

from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_tr = hstack((school_state_one_hot_train, categories_one_hot_train, sub_categories_one_hot_train,
teacher_prefix_one_hot_train,
project_grade_one_hot_train, essay_tfidf_w2v_train, title_tfidf_w2v_train, price_standardized_train,
previous_project_standardized_train))
X_cv = hstack((school_state_one_hot_cv, categories_one_hot_cv, sub_categories_one_hot_cv,
teacher_prefix_one_hot_cv, project_grade_one_hot_cv, essay_tfidf_w2v_cv, title_tfidf_w2v_cv,
price_standardized_cv, previous_project_standardized_cv))
X_test = hstack((school_state_one_hot_test, categories_one_hot_test, sub_categories_one_hot_test, teacher_prefix_one_hot_test,
project_grade_one_hot_test, essay_tfidf_w2v_test, title_tfidf_w2v_test, price_standardized_test, previous_project_standardized_test))
X_tr = X_tr.tocsr()
X_cv = X_cv.tocsr()
X_test = X_test.tocsr()
print(X_tr.shape, y_tr.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)
```

```
(24500, 701) (24500,)
(10500, 701) (10500,)
(15000, 701) (15000,)
```

Applying KNN on train data and finding best hyperparameter with CV data

In [177]:

```
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 = []
```

```

train_auc = []
cv_auc = []
K = [3, 15, 25, 51, 101, 130]
for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
    neigh.fit(X_tr, y_tr)

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

    # 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_tr, 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')

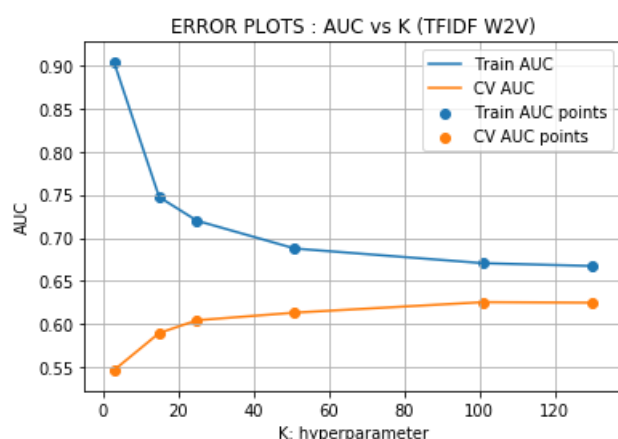
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS : AUC vs K (TFIDF W2V)")
plt.grid()
plt.show()

```

```

0%|
[00:00<?, ?it/s]
17%|
[11:28<57:24, 688.85s/it]
33%|
52:11, 782.92s/it]
50%|
<37:51, 757.25s/it]
67%|
[51:00<24:23, 731.77s/it]
83%|
6<11:45, 705.95s/it]
100%|
[1:42:45<00:00, 1231.89s/it]

```



Best Hyperparameter K

In [178]:

```
best_k4 = 130
```

Training model with best K and finding the train and test AUC

In [179]:

```

from sklearn.metrics import roc_curve, auc

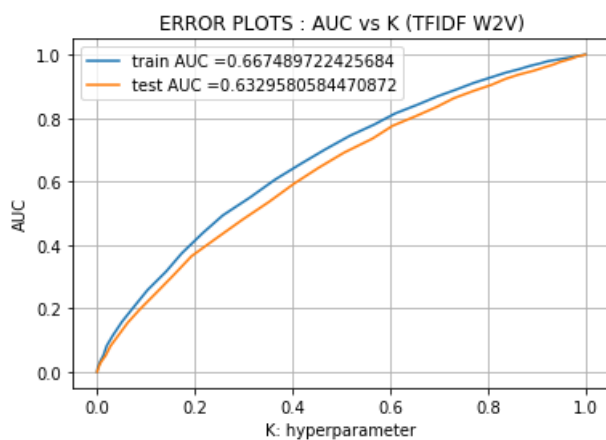
neigh = KNeighborsClassifier(n_neighbors=best_k4, n_jobs=-1)
neigh.fit(X_tr, y_tr)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

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

train_fpr, train_tpr, tr_thresholds = roc_curve(y_tr, 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 : AUC vs K (TFIDF W2V)")
plt.grid()
plt.show()

```



Confusion Matrix

In [180]:

```

print("="*100)
from sklearn.metrics import confusion_matrix
best_t4 = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_tr, predict_with_best_t(y_train_pred, best_t4)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t4)))

```

=====

the maximum value of tpr*(1-fpr) 0.3852985393820102 for threshold 0.846

Train confusion matrix

```
[[ 2494  1429]
 [ 8106 12471]]
```

Test confusion matrix

```
[[1432  970]
 [5115 7483]]
```

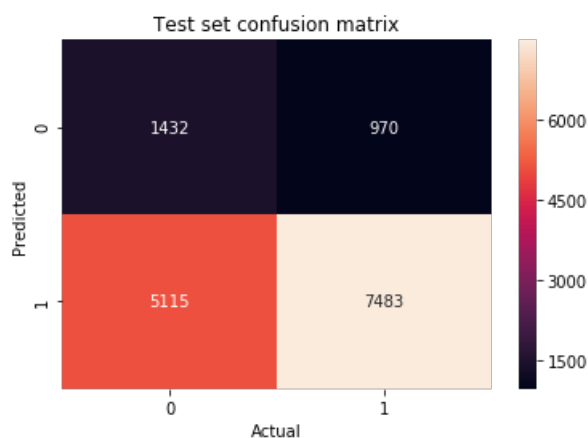
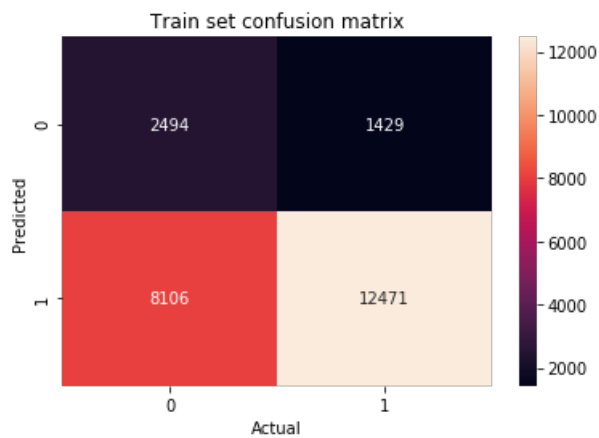
In [181]:

```

# Heatmap for train set confusion matrix(Select K best)
heatmap_train = sns.heatmap(confusion_matrix(y_tr, predict_with_best_t(y_train_pred, best_t4)),
annot=True, fmt="d")
plt.title("Train set confusion matrix")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.show()

```

```
# Heatmap for test set confusion matrix(Select K best)
heatmap_train = sns.heatmap(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t4)),
                             not=True, fmt="d")
plt.title("Test set confusion matrix")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.show()
```



2.5 Feature selection with `SelectKBest`

In [196]:

```
# 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
print("="*40 + "Before Feature Selection" + "="*40)
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_tr = hstack((school_state_one_hot_train, categories_one_hot_train, sub_categories_one_hot_train,
               teacher_prefix_one_hot_train, project_grade_one_hot_train, essay_tfidf_train, title_tfidf_train, price_standardized_train, previous_project_standardized_train))
X_cv = hstack((school_state_one_hot_cv, categories_one_hot_cv, sub_categories_one_hot_cv,
               teacher_prefix_one_hot_cv, project_grade_one_hot_cv, essay_tfidf_cv, title_tfidf_cv, price_standardized_cv, previous_project_standardized_cv))
X_test = hstack((school_state_one_hot_test, categories_one_hot_test, sub_categories_one_hot_test, teacher_prefix_one_hot_test, project_grade_one_hot_test, essay_tfidf_test, title_tfidf_test, price_standardized_test, previous_project_standardized_test))
X_tr = X_tr.tocsr()
X_cv = X_cv.tocsr()
X_test = X_test.tocsr()
print(X_tr.shape, X_cv.shape, X_test.shape)
```

```
print(X_tr.shape , y_tr.shape)
print(X_cv.shape , y_cv.shape)
print(X_test.shape , y_test.shape)
```

```
=====Before Feature
Selection=====
(24500, 10547) (24500,)
(10500, 10547) (10500,)
(15000, 10547) (15000,)
```

Selecting top 2000 features using 'SelectKBest'

In [197]:

```
#Feature Selection https://scikit-learn.org/stable/modules/generated/sklearn.feature\_selection.SelectKBest.html
print("="*40 + "After Feature Selection" + "="*40)

from sklearn.feature_selection import SelectKBest, f_classif

selector = SelectKBest(f_classif, k=2000)
selector.fit(X_tr, y_tr)
X_tr_new = selector.transform(X_tr)
X_cv_new = selector.transform(X_cv)
X_test_new = selector.transform(X_test)
print(X_tr_new.shape , y_tr.shape)
print(X_cv_new.shape , y_cv.shape)
print(X_test_new.shape , y_test.shape)
```

```
=====After Feature
Selection=====
(24500, 2000) (24500,)
(10500, 2000) (10500,)
(15000, 2000) (15000,)
```

Applying KNN on train data after feature selection and finding best hyperparameter with CV data

In [198]:

```
train_auc = []
cv_auc = []
K = [3, 15, 25, 51, 101, 130]
for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
    neigh.fit(X_tr_new, y_tr)

    y_train_pred = batch_predict(neigh, X_tr_new)
    y_cv_pred = batch_predict(neigh, X_cv_new)

    # 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_tr, 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')

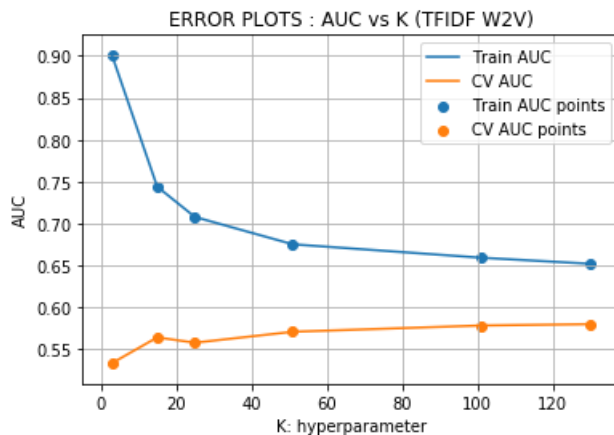
plt.scatter(K, train_auc, label='Train AUC points')
plt.scatter(K, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS : AUC vs K (TFIDF W2V)")
plt.grid()
plt.show()
```

```

0%|
[00:00<?, ?it/s]
17%|
[01:17<06:25, 77.11s/it]
33%|
<05:15, 78.83s/it]
50%|
2<03:59, 79.85s/it]
67%|
[05:26<02:42, 81.11s/it]
83%|
49<01:21, 81.79s/it]
100%|
:13<00:00, 82.28s/it]

```



Best Hyperparameter K

In [199]:

```
best_k5 = 130
```

Training model with best K and finding the train and test AUC

In [200]:

```

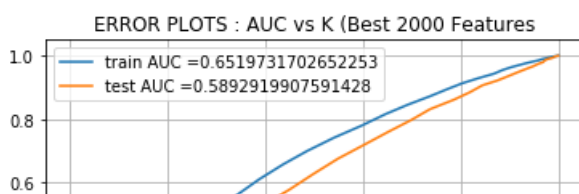
neigh = KNeighborsClassifier(n_neighbors=best_k5, n_jobs=-1)
neigh.fit(X_tr_new, y_tr)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

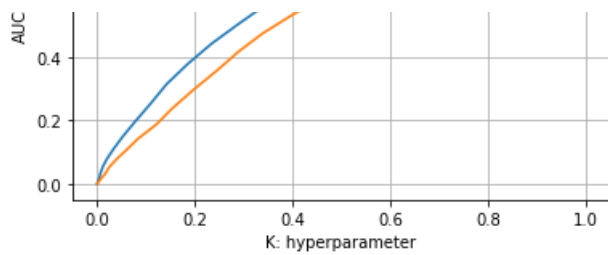
y_train_pred = batch_predict(neigh, X_tr_new)
y_test_pred = batch_predict(neigh, X_test_new)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_tr, 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 : AUC vs K (Best 2000 Features)")
plt.grid()
plt.show()

```





Confusion Matrix

In [201]:

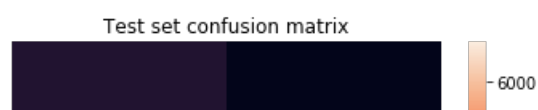
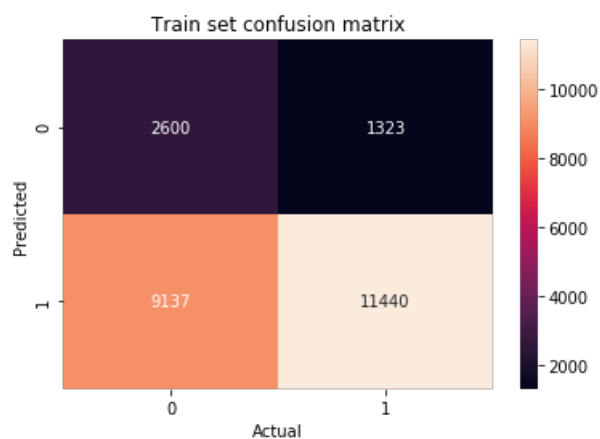
```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t5 = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_tr, predict_with_best_t(y_train_pred, best_t5)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t5)))
```

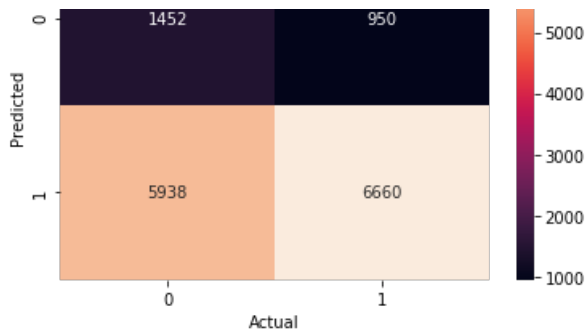
```
=====
the maximum value of tpr*(1-fpr) 0.37399194096604077 for threshold 0.838
Train confusion matrix
[[ 2405  1518]
 [ 8024 12553]]
Test confusion matrix
[[1322 1080]
 [5313 7285]]
```

In [202]:

```
# Heatmap for train set confusion matrix(Select K best)
heatmap_train = sns.heatmap(confusion_matrix(y_tr, predict_with_best_t(y_train_pred,best_t4)),
annot=True, fmt="d")
plt.title("Train set confusion matrix")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.show()

# Heatmap for test set confusion matrix(Select K best)
heatmap_train = sns.heatmap(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t4)), an
not=True, fmt="d")
plt.title("Test set confusion matrix")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.show()
```





3. Conclusions

In [203]:

```
# Please compare all your models using Prettytable library

from prettytable import PrettyTable

model_compare = PrettyTable()
model_compare.field_names = ["Feature_sets", "Best_k_value", "Best_threshold"]
model_compare.add_row(["Bag of words", best_k1, np.round(best_t1,3)])
model_compare.add_row(["TF-IDF", best_k2, np.round(best_t2,3)])
model_compare.add_row(["Average word2vector", best_k3, np.round(best_t3,3)])
model_compare.add_row(["TF-IDF Average word2vector", best_k4, np.round(best_t4,3)])
model_compare.add_row(["Select k best", best_k5, np.round(best_t5,3)])

print(model_compare)
```

Feature_sets	Best_k_value	Best_threshold
Bag of words	130	0.777
TF-IDF	130	0.846
Average word2vector	130	0.846
TF-IDF Average word2vector	130	0.846
Select k best	130	0.838

Summary

1) The Best Hyperparameter K is found to be 130 in all the cases. 2) The Best threshold value is found to be 0.846 in majority cases.