

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 | Descri |
|--|--|
| <code>project_id</code> | A unique identifier for the proposed project. Example: p03 |
| <code>project_title</code> | Title of the project. Exam Art Will Make You Ha First Grade |
| <code>project_grade_category</code> | Grade level of students for which the project is targeted. One of the foll enumerated va Grades Pr Grades Grades Grades |
| <code>project_subject_categories</code> | One or more (comma-separated) subject categories for the project fro following enumerated list of va Applied Lear Care & Hu Health & Sp History & Ci Literacy & Lang Math & Sci Music & The Special N Wa |
| <code>school_state</code> | State where school is located (Two-letter U.S. postal (https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_co) Example Music & The Literacy & Language, Math & Sci |
| <code>project_subject_subcategories</code> | One or more (comma-separated) subject subcategories for the pr Exam Lite Literature & Writing, Social Scie |
| <code>project_resource_summary</code> | An explanation of the resources needed for the project. Exam My students need hands on literacy materials to mar sensory ne |
| <code>project_essay_1</code> | First application e |
| <code>project_essay_2</code> | Second application e |
| <code>project_essay_3</code> | Third application e |
| <code>project_essay_4</code> | Fourth application e |

| Feature | Description |
|---|---|
| project_submitted_datetime | Datetime when project application was submitted. Example: 2016-04-12:43:56 |
| teacher_id | A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c1 |
| teacher_prefix | Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> • • • • • • |
| teacher_number_of_previously_posted_projects | Number of project applications previously submitted by the same teacher. Example: 1 |

* 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 <code>0</code> indicates the project was not approved, and a value of <code>1</code> 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 [0]: %matplotlib inline
import warnings
warnings.filterwarnings("ignore")

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

```
D:\installed\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

1.1 Reading Data

```
In [0]: project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')
```

```
In [0]: print("Number of data points in train data", project_data.shape)
print('-'*50)
print("The attributes of data :", project_data.columns.values)
```

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

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

```
In [0]: # how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
cols = ['Date' if x=='project_submitted_datetime' else x for x in list(project_data.columns)]
```

```
#sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
```

```
project_data['Date'] = pd.to_datetime(project_data['project_submitted_datetime'])
```

```
project_data.drop('project_submitted_datetime', axis=1, inplace=True)
```

```
project_data.sort_values(by=['Date'], inplace=True)
```

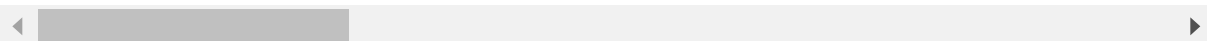
```
# how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
```

```
project_data = project_data[cols]
```

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

Out[0]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | |
|-------|------------|---------|----------------------------------|----------------|--------------|----------------|
| 55660 | 8393 | p205479 | 2bf07ba08945e5d8b2a3f269b2b3cfe5 | Mrs. | CA | 2 0 00:2 |
| 76127 | 37728 | p043609 | 3f60494c61921b3b43ab61bdde2904df | Ms. | UT | 2 0 00:3 |



```
In [0]: print("Number of data points in train data", resource_data.shape)
print(resource_data.columns.values)
resource_data.head(2)
```

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

Out[0]:

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

1.2 preprocessing of project_subject_categories

```
In [0]: categories = list(project_data['project_subject_categories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039

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

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

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

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

1.3 preprocessing of project_subject_subcategories

```
In [0]: sub_categories = list(project_data['project_subject_subcategories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039

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

sub_cat_list = []
for i in sub_categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
            j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
            temp +=j.strip()+" #" "abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&','_')
            sub_cat_list.append(temp.strip())

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

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

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

1.3 Text preprocessing

```
In [0]: # merge two column text dataframe:
project_data["essay"] = project_data["project_essay_1"].map(str) + \
    project_data["project_essay_2"].map(str) + \
    project_data["project_essay_3"].map(str) + \
    project_data["project_essay_4"].map(str)
```



```
In [0]: project_data.head(2)
```

Out[0]:

| Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | |
|------------|---------------|----------------------------------|----------------|--------------|----------------|
| 55660 | 8393 p205479 | 2bf07ba08945e5d8b2a3f269b2b3cfe5 | Mrs. | CA | 2 0 00:2 |
| 76127 | 37728 p043609 | 3f60494c61921b3b43ab61bdde2904df | Ms. | UT | 2 0 00:3 |

```
In [0]: ##### 1.4.2.3 Using Pretrained Models: TFIDF weighted W2V
```

```
In [0]: # 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 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 levels. 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 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 scrapbooking. 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. Kindergarten 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 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 [0]: `# https://stackoverflow.com/a/47091490/4084039`
`import re`

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

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

```
In [0]: sent = decontracted(project_data['essay'].values[20000])
print(sent)
print("="*50)
```

\n"A person is a person, no matter how small.\n" (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. \n\nStudents in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures, including Native Americans.\n\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, \n"Can we try cooking with REAL food?\n" I will take their idea and create \n"Common Core Cooking Lessons\n" 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. \n\nStudents will gain math and literature skills as well as a life long enjoyment for healthy cooking.nannan

=====

```
In [0]: # \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
sent = sent.replace('\r', ' ')
sent = sent.replace('\n', ' ')
sent = sent.replace('\t', ' ')
print(sent)
```

A person is a person, no matter how small. (Dr.Seuss) I teach the smallest students 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 [0]: #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 nannan


```
In [0]: # https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you',
, "you're", "you've", \
, "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he',
, 'him', 'his', 'himself', \
, 'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'it
self', 'they', 'them', 'their', \
, 'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 't
hat', "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', 'becau
se', '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', 'a
ll', 'any', 'both', 'each', 'few', 'more', \
, 'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'tha
n', 'too', 'very', \
, 's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "shoul
d've", 'now', 'd', 'll', 'm', 'o', 're', \
, 've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn',
"didn't", 'doesn', "doesn't", 'hadn', \
, 'hadn't', 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'm
a', 'mightn', "mightn't", 'mustn', \
, 'mustn't', 'needn', "needn't", 'shan', "shan't", 'shouldn', "shoul
dn't", 'wasn', "wasn't", 'weren', "weren't", \
, 'won', "won't", 'wouldn', "wouldn't"]
```

```
In [0]: # Combining all the above students
from tqdm import tqdm
preprocessed_essays = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
    preprocessed_essays.append(sent.lower().strip())
```

```
100%|███████████████████████████████████████████████████  
| 109248/109248 [01:48<00:00, 1005.43it/s]
```

```
In [0]: # after preprocessing
preprocessed_essays[20000]
```

```
Out[0]: 'a person person no matter small dr seuss i teach smallest students biggest e
nthusiasm learning my students learn many different ways using senses multipl
e intelligences i use wide range techniques help students succeed students cl
ass come variety different backgrounds makes wonderful sharing experiences cu
ltures including native americans our school caring community successful lear
ners seen collaborative student project based learning classroom kindergarten
ers class love work hands materials many different opportunities practice ski
ll mastered having social skills work cooperatively friends crucial aspect ki
ndergarten curriculum montana perfect place learn agriculture nutrition my st
udents love role play pretend kitchen early childhood classroom i several kid
s ask can try cooking real food i take idea create common core cooking lesson
s learn important math writing concepts cooking delicious healthy food snack
time my students grounded appreciation work went making food knowledge ingred
ients came well healthy bodies this project would expand learning nutrition a
gricultural cooking recipes us peel apples make homemade applesauce make brea
d mix healthy plants classroom garden spring we also create cookbooks printed
shared families students gain math literature skills well life long enjoyment
healthy cooking nannan'
```

1.4 Preprocessing of `project_title`

```
In [0]: # similarly you can preprocess the titles also
```

1.5 Preparing data for models

```
In [0]: project_data.columns
```

```
Out[0]: Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
              'Date', 'project_grade_category', 'project_title', 'project_essay_1',
              'project_essay_2', 'project_essay_3', 'project_essay_4',
              'project_resource_summary',
              'teacher_number_of_previously_posted_projects', 'project_is_approved',
              'clean_categories', 'clean_subcategories', 'essay'],
              dtype='object')
```

we are going to consider

- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data

- project_title : text data
- text : text data
- project_resource_summary: text data (optinal)

- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical

1.5.1 Vectorizing Categorical data

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

```
In [0]: # we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
categories_one_hot = vectorizer.fit_transform(project_data['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encoding ", categories_one_hot.shape)

['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning',
 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding (109248, 9)
```

```
In [0]: # we use count vectorizer to convert the values into one
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowe
rcase=False, binary=True)
sub_categories_one_hot = vectorizer.fit_transform(project_data['clean_subcateg
ories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",sub_categories_one_hot.shape)
```

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement',
'Extracurricular', 'Civics_Government', 'ForeignLanguages', 'NutritionEducati
on', 'Warmth', 'Care_Hunger', 'SocialSciences', 'PerformingArts', 'CharacterE
ducation', 'TeamSports', 'Other', 'College_CareerPrep', 'Music', 'History_Geo
graphy', 'Health_LifeScience', 'EarlyDevelopment', 'ESL', 'Gym_Fitness', 'Env
ironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences', 'Spec
ialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
Shape of matrix after one hot encodig (109248, 30)
```

```
In [0]: # you can do the similar thing with state, teacher_prefix and project_grade_ca
tegory also
```

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

```
In [0]: # We are considering only the words which appeared in at least 10 documents(ro
ws or projects).
vectorizer = CountVectorizer(min_df=10)
text_bow = vectorizer.fit_transform(preprocessed_essays)
print("Shape of matrix after one hot encodig ",text_bow.shape)
```

```
Shape of matrix after one hot encodig (109248, 16623)
```

```
In [0]: # you can vectorize the title also
# before you vectorize the title make sure you preprocess it
```

1.5.2.2 TFIDF vectorizer

```
In [0]: from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
text_tfidf = vectorizer.fit_transform(preprocessed_essays)
print("Shape of matrix after one hot encodig ",text_tfidf.shape)
```

```
Shape of matrix after one hot encodig (109248, 16623)
```

1.5.2.3 Using Pretrained Models: Avg W2V

```

In [0]: '''
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
def loadGloveModel(gloveFile):
    print ("Loading Glove Model")
    f = open(gloveFile, 'r', encoding="utf8")
    model = {}
    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding
    print ("Done.", len(model), " words loaded!")
    return model
model = loadGloveModel('glove.42B.300d.txt')

# =====
Output:

Loading Glove Model
1917495it [06:32, 4879.69it/s]
Done. 1917495 words loaded!

# =====

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

for i in preproced_titles:
    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 coup
us", \
      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('glove_vectors', 'wb') as f:
    pickle.dump(words_courpus, f)

```

```

Out[0]: "\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/40
84039\ndef loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n
f = open(gloveFile,\r', encoding="utf8")\n    model = {}\n    for line in t
qdm(f):\n        splitLine = line.split()\n        word = splitLine[0]\n
embedding = np.array([float(val) for val in splitLine[1:]])\n        model[word]
= embedding\n    print ("Done.",len(model)," words loaded!")\n    return
model\nmodel = loadGloveModel('glove.42B.300d.txt')\n\n# =====
=====
\nOutput:\n    \nLoading Glove Model\n1917495it [06:32, 4879.69it/s]
\nDone. 1917495 words loaded!\n\n# =====
=====
\n\nwords =
[]\nfor i in preproced_texts:\n    words.extend(i.split(' '))\n\nfor i in p
reproced_titles:\n    words.extend(i.split(' '))\n\nprint("all the words in t
he coupus", len(words))\nwords = set(words)\nprint("the unique words in the c
oupus", len(words))\n\ninter_words = set(model.keys()).intersection(words)\npr
int("The number of words that are present in both glove vectors and our coup
us", len(inter_words), "(", np.round(len(inter_words)/len(words)*100,
3), "%")\n\nwords_courpus = {}\nwords_glove = set(model.keys())\nfor i in wor
ds:\n    if i in words_glove:\n        words_courpus[i] = model[i]\n\nprint("wo
rd 2 vec length", len(words_courpus))\n\n\n# stronging variables into pickle
files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-v
ariables-in-python/\n\nimport pickle\nwith open('glove_vectors', 'wb') as
f:\n    pickle.dump(words_courpus, f)\n\n\n"

```

```

In [0]: # stronging variables into pickle files python: http://www.jessicayung.com/how
-to-use-pickle-to-save-and-load-variables-in-python/
# make sure you have the glove_vectors file
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())

```

```

In [0]: # average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this
list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors.append(vector)

print(len(avg_w2v_vectors))
print(len(avg_w2v_vectors[0]))

```

```

100%|████████████████████████████████████████████████████████████████████████████████
| 109248/109248 [00:59<00:00, 1830.39it/s]

```

```

109248
300

```



```
In [0]: # check this one: https://www.youtube.com/watch?v=0H0q0cLn3Z4&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 32
9. ... 399. 287.73 5.5 ].
# Reshape your data either using array.reshape(-1, 1)

price_scalar = StandardScaler()
price_scalar.fit(project_data['price'].values.reshape(-1,1)) # finding the mean and standard 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 mean and variance.
price_standardized = price_scalar.transform(project_data['price'].values.reshape(-1, 1))
```

```
In [0]: price_standardized
```

```
Out[0]: array([[4.63560392e-03, 1.36200635e-03, 2.10346002e-03, ...,
2.55100471e-03, 1.83960046e-03, 3.51642253e-05]])
```

1.5.4 Merging all the above features

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

```
In [0]: print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(text_bow.shape)
print(price_standardized.shape)
```

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

```
In [0]: # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X = hstack((categories_one_hot, sub_categories_one_hot, text_bow, price_standardized))
X.shape
```

```
Out[0]: (109248, 16663)
```


Assignment 3: Apply KNN


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

- **Set 1**: categorical, numerical features + project_title(BOW) + preprocessed_essay (BOW)
- **Set 2**: categorical, numerical features + project_title(TFIDF)+ preprocessed_essay (TFIDF)
- **Set 3**: categorical, numerical features + project_title(AVG W2V)+ preprocessed_essay (AVG W2V)
- **Set 4**: categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

2. Hyper paramter tuning to find best K

- Find the best hyper parameter which results in the maximum [AUC](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/>) 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](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/>) with predicted and original labels of test data points

4. [Task-2]

- Select top 2000 features from feature **Set 2** using `'SelectKBest'` (https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.SelectKBest.html) 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](http://zetcode.com/python/prettytable/) (<http://zetcode.com/python/prettytable/>)

**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-leakag, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on you train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link. \(https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf\)](https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf)

## 2. K Nearest Neighbor

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

```

In [1]: %matplotlib inline
import warnings
warnings.filterwarnings("ignore")

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

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

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

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

from tqdm import tqdm_notebook as tqdm1
from tqdm import tqdm
import time
import os

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

from sklearn.model_selection import train_test_split

```

```

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

```

```
In [2]: project_data = pd.read_csv('../train_data.csv', nrows=50000)
resource_data = pd.read_csv('../resources.csv')
```

```
In [3]: print("Number of data points in train data", project_data.shape)
print('-'*50)
print("The attributes of data :", project_data.columns.values)
```

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

```

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

```
In []:
```

## Text preprocessing(1)

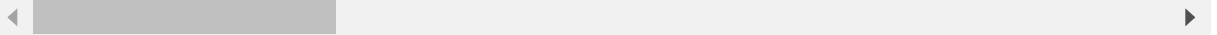
```
In [4]: categories = list(project_data['project_subject_categories'].values)
remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039

https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
cat_list = []
for i in categories:
 temp = ""
 # consider we have text like this "Math & Science, Warmth, Care & Hunger"
 for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
 if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
 j=j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
 j = j.replace(' ', '') # we are placing all the ' ' (space) with '' (empty) ex: "Math & Science"=> "Math&Science"
 temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing spaces
 temp = temp.replace('&', '_') # we are replacing the & value into
 cat_list.append(temp.strip())
```

```
In [5]: project_data['clean_categories'] = cat_list
project_data.drop(['project_subject_categories'], axis=1, inplace=True)
project_data.head(5)
```

Out[5]:

|   | Unnamed: 0 | id      | teacher_id                       | teacher_prefix | school_state | project_: |
|---|------------|---------|----------------------------------|----------------|--------------|-----------|
| 0 | 160221     | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs.           | IN           |           |
| 1 | 140945     | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr.            | FL           |           |
| 2 | 21895      | p182444 | 3465aaf82da834c0582ebd0ef8040ca0 | Ms.            | AZ           |           |
| 3 | 45         | p246581 | f3cb9bffbba169bef1a77b243e620b60 | Mrs.           | KY           |           |
| 4 | 172407     | p104768 | be1f7507a41f8479dc06f047086a39ec | Mrs.           | TX           |           |



```
In [6]: # count of all the words in corpus python: https://stackoverflow.com/a/2289859
5/4084039
from collections import Counter
my_counter = Counter()
for word in project_data['clean_categories'].values:
 my_counter.update(word.split())
my_counter
```

```
Out[6]: Counter({'Literacy_Language': 23998,
 'History_Civics': 2689,
 'Health_Sports': 6538,
 'Math_Science': 18874,
 'SpecialNeeds': 6233,
 'AppliedLearning': 5569,
 'Music_Arts': 4699,
 'Warmth': 643,
 'Care_Hunger': 643})
```

```
In [7]: # dict sort by value python: https://stackoverflow.com/a/613218/4084039
cat_dict = dict(my_counter)
sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))

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

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

```
In [8]: 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-fr
om-a-string
https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-strin
g-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 & Scienc
e", "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 ''(emp
ty) ex:"Math & Science"=>"Math&Science"
 temp +=j.strip()+" #" abc ".strip() will return "abc", remove the tra
iling spaces
 temp = temp.replace('&','_')
 sub_cat_list.append(temp.strip())
```

```
In [9]: project_data['clean_subcategories'] = sub_cat_list
project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
project_data.head(2)
```

Out[9]:

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

```
In [10]: # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()
for word in project_data['clean_subcategories'].values:
 my_counter.update(word.split())
```

```
In [11]: # dict sort by value python: https://stackoverflow.com/a/613218/4084039
sub_cat_dict = dict(my_counter)
sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))

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

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

```
In [12]: # 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 [13]: # https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indices-for-all-groups-in-one-step
price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'})
.reset_index()
price_data.head(2)
```

Out[13]:

|   | id      | price  | quantity |
|---|---------|--------|----------|
| 0 | p000001 | 459.56 | 7        |
| 1 | p000002 | 515.89 | 21       |

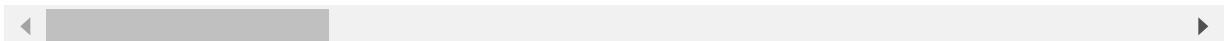
```
In [14]: # join two dataframes in python:
project_data = pd.merge(project_data, price_data, on='id', how='left')
```

```
In [15]: #presence of the numerical digits in a strings with numeric : https://stackove
rflow.com/a/19859308/8089731
def hasNumbers(inputString):
 return any(i.isdigit() for i in inputString)
p1 = project_data[['id', 'project_resource_summary']]
p1 = pd.DataFrame(data=p1)
p1.columns = ['id', 'digits_in_summary']
p1['digits_in_summary'] = p1['digits_in_summary'].map(hasNumbers)
https://stackoverflow.com/a/17383325/8089731
p1['digits_in_summary'] = p1['digits_in_summary'].astype(int)
project_data = pd.merge(project_data, p1, on='id', how='left')
project_data.head(5)
```

Out[15]:

|   | Unnamed:<br>0 |         | id                               | teacher_id | teacher_prefix | school_state | project_: |
|---|---------------|---------|----------------------------------|------------|----------------|--------------|-----------|
| 0 | 160221        | p253737 | c90749f5d961ff158d4b4d1e7dc665fc |            | Mrs.           | IN           |           |
| 1 | 140945        | p258326 | 897464ce9ddc600bced1151f324dd63a |            | Mr.            | FL           |           |
| 2 | 21895         | p182444 | 3465aaf82da834c0582ebd0ef8040ca0 |            | Ms.            | AZ           |           |
| 3 | 45            | p246581 | f3cb9bffbba169bef1a77b243e620b60 |            | Mrs.           | KY           |           |
| 4 | 172407        | p104768 | be1f7507a41f8479dc06f047086a39ec |            | Mrs.           | TX           |           |

5 rows × 21 columns



## Text preprocessing(2)

In [16]: [# https://stackoverflow.com/a/47091490/4084039](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 [17]: [# https://gist.github.com/sebleier/554280](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', 'it
self', 'they', 'them', 'their', \
 'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 't
hat', "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', 'becau
se', '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', 'a
ll', 'any', 'both', 'each', 'few', 'more', \
 'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'tha
n', 'too', 'very', \
 's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "shoul
d've", 'now', 'd', 'll', 'm', 'o', 're', \
 've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn',
"didn't", 'doesn', "doesn't", 'hadn', \
 "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'm
a', 'mightn', "mightn't", 'mustn', \
 "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shoul
dn't", 'wasn', "wasn't", 'weren', "weren't", \
 'won', "won't", 'wouldn', "wouldn't"]
```

```
In [18]: # Combining all the above statements
from tqdm import tqdm
preprocessed_essays = []
tqdm is for printing the status bar
for sentence in tqdm(project_data['essay'].values):
 sent = decontracted(sentence)
 sent = sent.replace('\\r', ' ')
 sent = sent.replace('\\\"', ' ')
 sent = sent.replace('\\n', ' ')
 sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
 # https://gist.github.com/sebleier/554280
 sent = ' '.join(e for e in sent.split() if e not in stopwords)
 preprocessed_essays.append(sent.lower().strip())
```

```
100%|███████████ 50000/50000 [00:28<00:00, 1738.97it/s]
```

```
In [19]: from tqdm import tqdm
preprocessed_titles = []
tqdm is for printing the status bar
for title in tqdm(project_data['project_title'].values):
 _title = decontracted(title)
 _title = _title.replace('\r', ' ')
 _title = _title.replace('\n', ' ')
 _title = _title.replace('\n', ' ')
 _title = re.sub('[^A-Za-z0-9]+', ' ', _title)
 # https://gist.github.com/sebleier/554280
 _title = ' '.join(e for e in _title.split() if e not in stopwords)
 preprocessed_titles.append(_title.lower().strip())
```

```
100%|██████████| 50000/50000 [00:01<00:00, 36002.42it/s]
```

```
In [20]: preprocessed_titles[1000]
```

```
Out[20]: 'sailing into super 4th grade year'
```

```
In [21]: project_grade_catogories = 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

project_grade_cat_list = []
for i in tqdm1(project_grade_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('&','_')
 project_grade_cat_list.append(temp.strip())
```

```
In [22]: project_data['clean_project_grade_category'] = project_grade_cat_list
project_data.drop(['project_grade_category'], axis=1, inplace=True)
project_data.head(2)
```

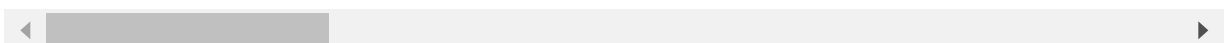
Out[22]:

|  | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_ |
|--|------------|----|------------|----------------|--------------|----------|
|--|------------|----|------------|----------------|--------------|----------|

|   |        |         |                                  |      |    |  |
|---|--------|---------|----------------------------------|------|----|--|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN |  |
|---|--------|---------|----------------------------------|------|----|--|

|   |        |         |                                 |     |    |  |
|---|--------|---------|---------------------------------|-----|----|--|
| 1 | 140945 | p258326 | 897464ce9ddc600bcd1151f324dd63a | Mr. | FL |  |
|---|--------|---------|---------------------------------|-----|----|--|

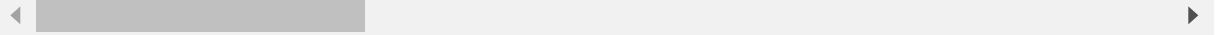
2 rows × 21 columns



```
In [23]: project_data.drop(['project_essay_1', 'project_essay_2', 'project_essay_3', 'project_essay_4'], axis=1, inplace=True)
project_data.head(2)
```

Out[23]:

|   | Unnamed: 0 | id      | teacher_id                       | teacher_prefix | school_state | project_ |
|---|------------|---------|----------------------------------|----------------|--------------|----------|
| 0 | 160221     | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs.           | IN           |          |
| 1 | 140945     | p258326 | 897464ce9ddc600bcd1151f324dd63a  | Mr.            | FL           |          |



```
In [24]: #Replacing Nan's with maximum occured value: https://stackoverflow.com/a/51053916/8089731
project_data['teacher_prefix'].value_counts().argmax()
project_data.fillna(value=project_data['teacher_prefix'].value_counts().argmax(), axis=1, inplace=True)
```

```
In [25]: project_data['preprocessed_essays'] = preprocessed_essays
project_data['preprocessed_titles'] = preprocessed_titles
```

```
In [26]: project_data.columns
```

```
Out[26]: Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
 'project_submitted_datetime', 'project_title',
 'project_resource_summary',
 'teacher_number_of_previously_posted_projects', 'project_is_approved',
 'clean_categories', 'clean_subcategories', 'essay', 'price', 'quantity',
 'digits_in_summary', 'clean_project_grade_category',
 'preprocessed_essays', 'preprocessed_titles'],
 dtype='object')
```

```
In [27]: # 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 helpful in debugging your code
when you plot any graph make sure you use
a. Title, that describes your plot, this will be very helpful to the reader
b. Legends if needed
c. X-axis Label
d. Y-axis Label
```

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

```
In [28]: X_train, X_test, y_train, y_test = train_test_split(project_data, project_data[
'project_is_approved'], test_size=0.33, stratify = project_data['project_is_approved'])
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train)

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

### 1.4.1 Vectorizing Categorical data

```
In [29]: # we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['clean_categories'].values)
print(vectorizer.get_feature_names())

categories_one_hot_train = vectorizer.transform(X_train['clean_categories'].values)
categories_one_hot_cv = vectorizer.transform(X_cv['clean_categories'].values)
categories_one_hot_test = vectorizer.transform(X_test['clean_categories'].values)
print("Shape of matrix after one hot encoding_train ", categories_one_hot_train.shape)
print("Shape of matrix after one hot encoding_cv ", categories_one_hot_cv.shape)
print("Shape of matrix after one hot encoding_test ", categories_one_hot_test.shape)

['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding_train (22445, 9)
Shape of matrix after one hot encoding_cv (11055, 9)
Shape of matrix after one hot encoding_test (16500, 9)
```

```
In [30]: # we use count vectorizer to convert the values into one hot encoded features
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['clean_subcategories'].values)
print(vectorizer.get_feature_names())

sub_categories_one_hot_train = vectorizer.transform(X_train['clean_subcategories'].values)
sub_categories_one_hot_cv = vectorizer.transform(X_cv['clean_subcategories'].values)
sub_categories_one_hot_test = vectorizer.transform(X_test['clean_subcategories'].values)
print("Shape of matrix after one hot encoding_train ", sub_categories_one_hot_train.shape)
print("Shape of matrix after one hot encoding_cv ", sub_categories_one_hot_cv.shape)
print("Shape of matrix after one hot encoding_test ", sub_categories_one_hot_test.shape)

['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 matrix after one hot encoding_train (22445, 30)
Shape of matrix after one hot encoding_cv (11055, 30)
Shape of matrix after one hot encoding_test (16500, 30)
```



```
In [31]: # we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(lowercase=False, binary=True)
vectorizer.fit(X_train['school_state'].values)
print(vectorizer.get_feature_names())

school_state_one_hot_train = vectorizer.transform(X_train['school_state'].values)
school_state_one_hot_cv = vectorizer.transform(X_cv['school_state'].values)
school_state_one_hot_test = vectorizer.transform(X_test['school_state'].values)
print("Shape of matrix after one hot encoding_train ", school_state_one_hot_train.shape)
print("Shape of matrix after one hot encoding_cv ", school_state_one_hot_cv.shape)
print("Shape of matrix after one hot encoding_test ", school_state_one_hot_test.shape)

['AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DC', 'DE', 'FL', 'GA', 'HI', 'IA', 'ID', 'IL', 'IN', 'KS', 'KY', 'LA', 'MA', 'MD', 'ME', 'MI', 'MN', 'MO', 'MS', 'MT', 'NC', 'ND', 'NE', 'NH', 'NJ', 'NM', 'NV', 'NY', 'OH', 'OK', 'OR', 'PA', 'RI', 'SC', 'SD', 'TN', 'TX', 'UT', 'VA', 'VT', 'WA', 'WI', 'WV', 'WY']
Shape of matrix after one hot encoding_train (22445, 51)
Shape of matrix after one hot encoding_cv (11055, 51)
Shape of matrix after one hot encoding_test (16500, 51)
```

```
In [32]: # Replacing Nan's with maximum occurred value: https://stackoverflow.com/a/51053916/8089731
project_data['teacher_prefix'].value_counts().argmax()
project_data.fillna(value=project_data['teacher_prefix'].value_counts().argmax(), axis=1, inplace=True)
```

```
In [33]: # we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(lowercase=False, binary=True)
vectorizer.fit(X_train['teacher_prefix'].values.astype('U'))
print(vectorizer.get_feature_names())

#https://stackoverflow.com/a/39308809/8089731
teacher_prefix_one_hot_train = vectorizer.transform(X_train['teacher_prefix'].values.astype('U'))
teacher_prefix_one_hot_cv = vectorizer.transform(X_cv['teacher_prefix'].values.astype('U'))
teacher_prefix_one_hot_test = vectorizer.transform(X_test['teacher_prefix'].values.astype('U'))
print("Shape of matrix after one hot encodig_train ", teacher_prefix_one_hot_train.shape)
print("Shape of matrix after one hot encodig_cv ", teacher_prefix_one_hot_cv.shape)
print("Shape of matrix after one hot encodig_test ", teacher_prefix_one_hot_test.shape)

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

```
In [34]: print(project_data['clean_project_grade_category'].unique())

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

```
In [35]: # we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
https://stackoverflow.com/a/38161028/8089731
pattern = "(?u)\\b[\\w-]+\\b"
vectorizer = CountVectorizer(token_pattern=pattern, lowercase=False, binary=True)
vectorizer.fit(X_train['clean_project_grade_category'].values)
print(vectorizer.get_feature_names())

#https://stackoverflow.com/a/39308809/8089731
project_grade_category_one_hot_train = vectorizer.transform(X_train['clean_project_grade_category'].values)
project_grade_category_one_hot_cv = vectorizer.transform(X_cv['clean_project_grade_category'].values)
project_grade_category_one_hot_test = vectorizer.transform(X_test['clean_project_grade_category'].values)
print("Shape of matrix after one hot encoding_train ", project_grade_category_one_hot_train.shape)
print("Shape of matrix after one hot encoding_cv ", project_grade_category_one_hot_cv.shape)
print("Shape of matrix after one hot encoding_test ", project_grade_category_one_hot_test.shape)

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

## Vectorizing Numerical features

```

In [36]: # check this one: https://www.youtube.com/watch?v=0H0q0cLn3Z4&t=530s
standardization sklearn: https://scikit-learn.org/stable/modules/generated/s
klearn.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 32
9. ... 399. 287.73 5.5].
Reshape your data either using array.reshape(-1, 1)

price_scalar = StandardScaler()
price_scalar.fit(X_train['price'].values.reshape(-1,1)) # finding the mean and
standard 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 mean and variance.
price_standardized_train = price_scalar.transform(X_train['price'].values.resh
ape(-1, 1))
price_standardized_cv = price_scalar.transform(X_cv['price'].values.reshape(-1
, 1))
price_standardized_test = price_scalar.transform(X_test['price'].values.reshap
e(-1, 1))
print(price_standardized_train.shape)
print(price_standardized_cv.shape)
print(price_standardized_test.shape)

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

```

```
In [37]: # check this one: https://www.youtube.com/watch?v=0H0qOcLn3Z4&t=530s
standardization sklearn: https://scikit-learn.org/stable/modules/generated/s
klearn.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 32
9. ... 399. 287.73 5.5].
Reshape your data either using array.reshape(-1, 1)

quantity_scalar = StandardScaler()
quantity_scalar.fit(X_train['quantity'].values.reshape(-1,1)) # finding the me
an and standard deviation of this data
print(f"Mean : {quantity_scalar.mean_[0]}, Standard deviation : {np.sqrt(qua
ntity_scalar.var_[0])}")

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

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:  
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:  
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:  
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:  
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

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

```
In [38]: # check this one: https://www.youtube.com/watch?v=0H0q0cLn3Z4&t=530s
standardization sklearn: https://scikit-learn.org/stable/modules/generated/s
klearn.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 32
9. ... 399. 287.73 5.5].
Reshape your data either using array.reshape(-1, 1)

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

Now standardize the data with above mean and variance.
teacher_number_of_previously_posted_projects_standardized_train = teacher_numbe
r_of_previously_posted_projects_scalar.transform(X_train['teacher_number_of_p
reviously_posted_projects'].values.reshape(-1, 1))
teacher_number_of_previously_posted_projects_standardized_cv = teacher_number_
of_previously_posted_projects_scalar.transform(X_cv['teacher_number_of_previou
sly_posted_projects'].values.reshape(-1, 1))
teacher_number_of_previously_posted_projects_standardized_test = teacher_numbe
r_of_previously_posted_projects_scalar.transform(X_test['teacher_number_of_pre
viously_posted_projects'].values.reshape(-1, 1))
print(teacher_number_of_previously_posted_projects_standardized_train.shape)
print(teacher_number_of_previously_posted_projects_standardized_cv.shape)
print(teacher_number_of_previously_posted_projects_standardized_test.shape)
```

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:  
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:  
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:  
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:  
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

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

In [39]: *# please write all the code with proper documentation, and proper titles for each subsection*  
*# go through documentations and blogs before you start coding*  
*# first figure out what to do, and then think about how to do.*  
*# reading and understanding error messages will be very much helpfull in debugging your code*  
*# make sure you featurize train and test data separatly*  
  
*# when you plot any graph make sure you use*  
*# a. Title, that describes your plot, this will be very helpful to the reader*  
*# b. Legends if needed*  
*# c. X-axis label*  
*# d. Y-axis label*

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

In [40]: X\_train.head(2)

Out[40]:

|       | Unnamed: 0 | id      | teacher_id                       | teacher_prefix | school_state | proj |
|-------|------------|---------|----------------------------------|----------------|--------------|------|
| 34627 | 38025      | p202072 | c870270453a99e69b43ff65e95face8b | Mrs.           | SC           |      |
| 9670  | 83652      | p039580 | 497ba7027394d6490d2d2b34af2db5fd | Mrs.           | FL           |      |

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

```
In [41]: # We are considering only the words which appeared in at least 10 documents(ro
ws or projects).
vectorizer = CountVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_essays'])
```

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

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

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

```
In [42]: # We are considering only the words which appeared in at least 10 documents(ro
ws or projects).
vectorizer = CountVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_titles'])
```

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

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



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

```
In [43]: from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_essays'])

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

Shape of matrix after tfidf\_text\_train (22445, 8894)

Shape of matrix after tfidf\_text\_cv (11055, 8894)

Shape of matrix after tfidf\_text\_test (16500, 8894)

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

```
In [44]: from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_titles'])

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

Shape of matrix after tfidf\_title\_train (22445, 1249)

Shape of matrix after tfidf\_title\_cv (11055, 1249)

Shape of matrix after tfidf\_title\_test (16500, 1249)

```
In [45]: # 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())
```

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

```
In [46]: # average Word2Vec
compute average word2vec for each review.
avg_w2v_essays_vectors_train = []; # the avg-w2v for each sentence/review is s
tored in this list
for sentence in tqdm1(X_train['preprocessed_essays']): # for each review/sente
nce
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_essays_vectors_train.append(vector)

avg_w2v_essays_vectors_cv = []; # the avg-w2v for each sentence/review is stor
ed in this list
for sentence in tqdm1(X_cv['preprocessed_essays']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_essays_vectors_cv.append(vector)

avg_w2v_essays_vectors_test = []; # the avg-w2v for each sentence/review is st
ored in this list
for sentence in tqdm1(X_test['preprocessed_essays']): # for each review/senten
ce
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_essays_vectors_test.append(vector)

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

22445  
11055  
16500  
300  
300  
300

**Avg W2V on TITLES(Train,cv,test)**

```
In [47]: # average Word2Vec
compute average word2vec for each review.
avg_w2v_titles_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm1(X_train['preprocessed_titles']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_titles_vectors_train.append(vector)

avg_w2v_titles_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm1(X_cv['preprocessed_titles']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_titles_vectors_cv.append(vector)

avg_w2v_titles_vectors_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm1(X_test['preprocessed_titles']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_titles_vectors_test.append(vector)

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

22445  
11055  
16500  
300  
300  
300

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

```

In [48]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['preprocessed_essays'])
we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_
)))
tfidf_words = set(tfidf_model.get_feature_names())

*
average Word2Vec
compute average word2vec for each review.
tfidf_w2v_essays_vectors_train = []; # the avg-w2v for each sentence/review is
stored in this list
for sentence in tqdm1(X_train['preprocessed_essays']): # for each review/sente
nce
 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/revie
w
 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 v
alue((sentence.count(word)/len(sentence.split())))
 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split
())) # getting the tfidf value for each word
 vector += (vec * tf_idf) # calculating tfidf weighted w2v
 tf_idf_weight += tf_idf
 if tf_idf_weight != 0:
 vector /= tf_idf_weight
 tfidf_w2v_essays_vectors_train.append(vector)

average Word2Vec
compute average word2vec for each review.
tfidf_w2v_essays_vectors_cv = []; # the avg-w2v for each sentence/review is st
ored in this list
for sentence in tqdm1(X_cv['preprocessed_essays']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 tf_idf_weight = 0; # num of words with a valid vector in the sentence/revie
w
 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 v
alue((sentence.count(word)/len(sentence.split())))
 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split
())) # getting the tfidf value for each word
 vector += (vec * tf_idf) # calculating tfidf weighted w2v
 tf_idf_weight += tf_idf
 if tf_idf_weight != 0:
 vector /= tf_idf_weight
 tfidf_w2v_essays_vectors_cv.append(vector)

average Word2Vec
compute average word2vec for each review.

```

```

tfidf_w2v_essays_vectors_test = []; # the avg-w2v for each sentence/review is
 stored in this list
for sentence in tqdm1(X_test['preprocessed_essays']): # for each review/senten
 ce
 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/revie
 w
 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 v
 alue((sentence.count(word)/len(sentence.split())))
 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split
 ())) # getting the tfidf value for each word
 vector += (vec * tf_idf) # calculating tfidf weighted w2v
 tf_idf_weight += tf_idf
 if tf_idf_weight != 0:
 vector /= tf_idf_weight
 tfidf_w2v_essays_vectors_test.append(vector)

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

```

```

22445
11055
16500
300
300
300

```

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

```

In [49]: ## S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['preprocessed_essays'])
we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())

average Word2Vec
compute average word2vec for each review.
tfidf_w2v_essays_vectors_train = []; # the avg-w2v for each sentence/review
is stored in this list
for sentence in tqdm1(X_train['preprocessed_essays']): # for each review/sen
tence
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/rev
ew
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.spl
it())) # getting the tfidf value for each word
vector += (vec * tf_idf) # calculating tfidf weighted w2v
tf_idf_weight += tf_idf
if tf_idf_weight != 0:
vector /= tf_idf_weight
tfidf_w2v_essays_vectors_train.append(vector)

average Word2Vec
compute average word2vec for each review.
tfidf_w2v_essays_vectors_cv = []; # the avg-w2v for each sentence/review is
stored in this list
for sentence in tqdm1(X_cv['preprocessed_essays']): # for each review/senten
ce
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/rev
ew
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.spl
it())) # getting the tfidf value for each word
vector += (vec * tf_idf) # calculating tfidf weighted w2v
tf_idf_weight += tf_idf
if tf_idf_weight != 0:
vector /= tf_idf_weight
tfidf_w2v_essays_vectors_cv.append(vector)

average Word2Vec

```



```
compute average word2vec for each review.
tfidf_w2v_essays_vectors_test = []; # the avg-w2v for each sentence/review is
stored in this list
for sentence in tqdm1(X_test['preprocessed_essays']): # for each review/sentence
vector = np.zeros(300) # as word vectors are of zero length
tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
for word in sentence.split(): # for each word in a review/sentence
if (word in glove_words) and (word in tfidf_words):
vec = model[word] # getting the vector for each word
here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
getting the tfidf value for each word
vector += (vec * tf_idf) # calculating tfidf weighted w2v
tf_idf_weight += tf_idf
if tf_idf_weight != 0:
vector /= tf_idf_weight
tfidf_w2v_essays_vectors_test.append(vector)

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

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

```

In [50]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['preprocessed_titles'])
we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_
)))
tfidf_words = set(tfidf_model.get_feature_names())

*
average Word2Vec
compute average word2vec for each review.
tfidf_w2v_titles_vectors_train = []; # the avg-w2v for each sentence/review is
stored in this list
for sentence in tqdm1(X_train['preprocessed_titles']): # for each review/sente
nce
 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/revie
w
 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 v
alue((sentence.count(word)/len(sentence.split())))
 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split
())) # getting the tfidf value for each word
 vector += (vec * tf_idf) # calculating tfidf weighted w2v
 tf_idf_weight += tf_idf
 if tf_idf_weight != 0:
 vector /= tf_idf_weight
 tfidf_w2v_titles_vectors_train.append(vector)

average Word2Vec
compute average word2vec for each review.
tfidf_w2v_titles_vectors_cv = []; # the avg-w2v for each sentence/review is st
ored in this list
for sentence in tqdm1(X_cv['preprocessed_titles']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 tf_idf_weight = 0; # num of words with a valid vector in the sentence/revie
w
 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 v
alue((sentence.count(word)/len(sentence.split())))
 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split
())) # getting the tfidf value for each word
 vector += (vec * tf_idf) # calculating tfidf weighted w2v
 tf_idf_weight += tf_idf
 if tf_idf_weight != 0:
 vector /= tf_idf_weight
 tfidf_w2v_titles_vectors_cv.append(vector)

average Word2Vec
compute average word2vec for each review.

```

```

tfidf_w2v_titles_vectors_test = []; # the avg-w2v for each sentence/review is
 stored in this list
for sentence in tqdm1(X_test['preprocessed_titles']): # for each review/senten
ce
 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/revie
w
 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 v
 alue((sentence.count(word)/len(sentence.split())))
 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split
 ())) # getting the tfidf value for each word
 vector += (vec * tf_idf) # calculating tfidf weighted w2v
 tf_idf_weight += tf_idf
 if tf_idf_weight != 0:
 vector /= tf_idf_weight
 tfidf_w2v_titles_vectors_test.append(vector)

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

```

```

22445
11055
16500
300
300
300

```

In [51]:

```

please write all the code with proper documentation, and proper titles for e
ach 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 debug
ging 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 rea
der
 # b. Legends if needed
 # c. X-axis label
 # d. Y-axis label

```

## 2.4 Applying 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 [ ]:

```
In [56]: # please write all the code with proper documentation, and proper titles for each subsection
go through documentations and blogs before you start coding
first figure out what to do, and then think about how to do.
reading and understanding error messages will be very much helpfull in debugging your code

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

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

```
In [57]: # Please write all the code with proper documentation
```

```
In [58]: # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((categories_one_hot_train, sub_categories_one_hot_train, school_state_one_hot_train, teacher_prefix_one_hot_train
 , project_grade_category_one_hot_train, price_standardized_train,
 , teacher_number_of_previously_posted_projects_standardized_train, text_bow_train, title_bow_train)).tocsr()
X_cr = hstack((categories_one_hot_cv, sub_categories_one_hot_cv, school_state_one_hot_cv, teacher_prefix_one_hot_cv
 , project_grade_category_one_hot_cv, price_standardized_cv, quantity_standardized_cv
 , teacher_number_of_previously_posted_projects_standardized_cv, text_bow_cv, title_bow_cv)).tocsr()
X_te = hstack((categories_one_hot_test, sub_categories_one_hot_test, school_state_one_hot_test, teacher_prefix_one_hot_test
 , project_grade_category_one_hot_test, price_standardized_test, quantity_standardized_test
 , teacher_number_of_previously_posted_projects_standardized_test, text_bow_test, title_bow_test)).tocsr()

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

Final Data matrix on BOW

(22445, 10229) (22445,)

(11055, 10229) (11055,)

(16500, 10229) (16500,)

=====

=====

## 1.1 Method 1: Simple for loop (if you are having memory limitations use this)

```
In [59]: def batch_predict(clf, data):
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs

 y_data_pred = []
 tr_loop = data.shape[0] - data.shape[0]%1000
 # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
 # in this for loop we will iterate until the last 1000 multiplier
 for i in range(0, tr_loop, 1000):
 y_data_pred.extend(clf.predict_proba(data[i:i+1000]))[:,1])
 # we will be predicting for the last data points
 y_data_pred.extend(clf.predict_proba(data[tr_loop:]))[:,1])
 return y_data_pred
```

```

In [101]: 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 non-thresholded measure of decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

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

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

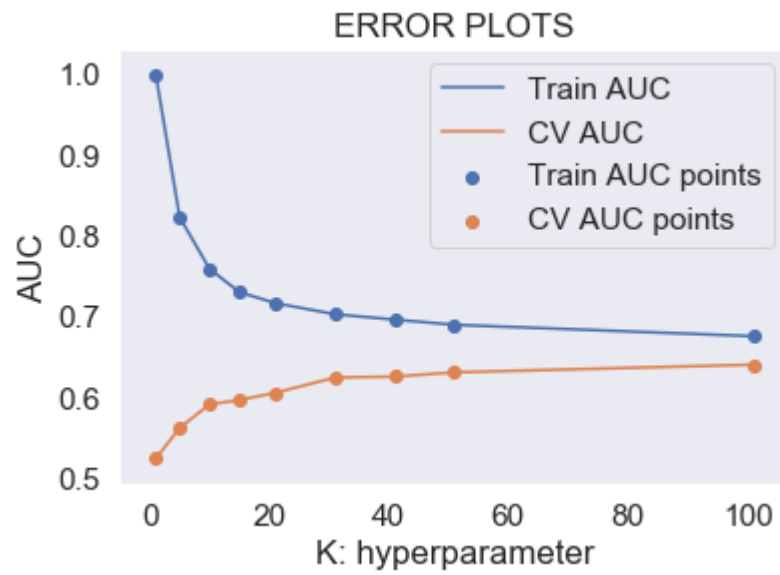
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs
 train_auc.append(roc_auc_score(y_train[:, :], y_train_pred))
 cv_auc.append(roc_auc_score(y_cv[:, :], y_cv_pred))

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

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

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

```



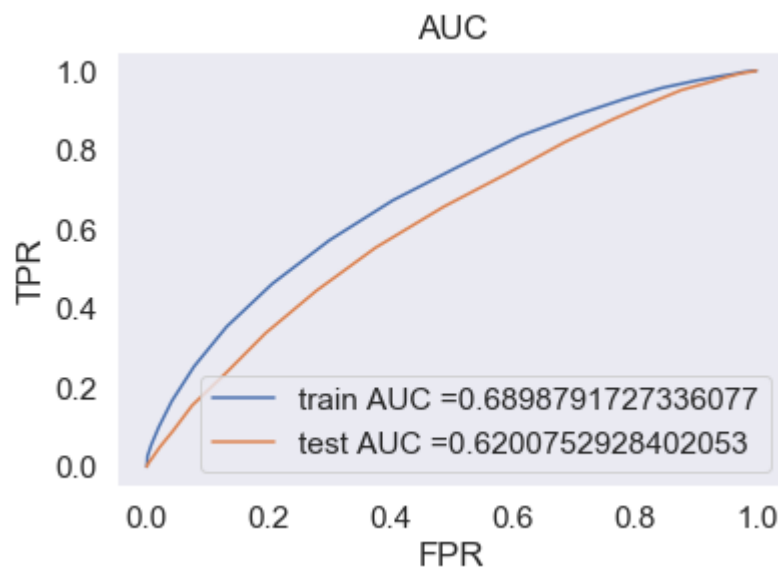
```
In [102]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

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

y_train_pred = batch_predict(neigh, X_tr[:, :])
y_test_pred = batch_predict(neigh, X_te[:, :])

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

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





```
In [103]: # we are writing our own function for predict, with defined threshold
we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

 t = threshold[np.argmax(tpr*(1-fpr))]

 # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

 print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
 predictions = []
 for i in proba:
 if i>=t:
 predictions.append(1)
 else:
 predictions.append(0)
 return predictions
```

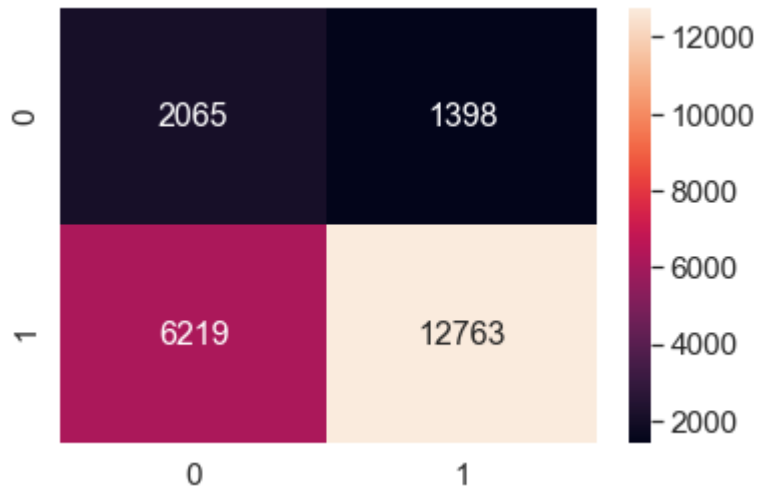
```
In [104]: print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

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

```
In [105]: conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of  $tpr*(1-fpr)$  0.4009390570266228 for threshold 0.784

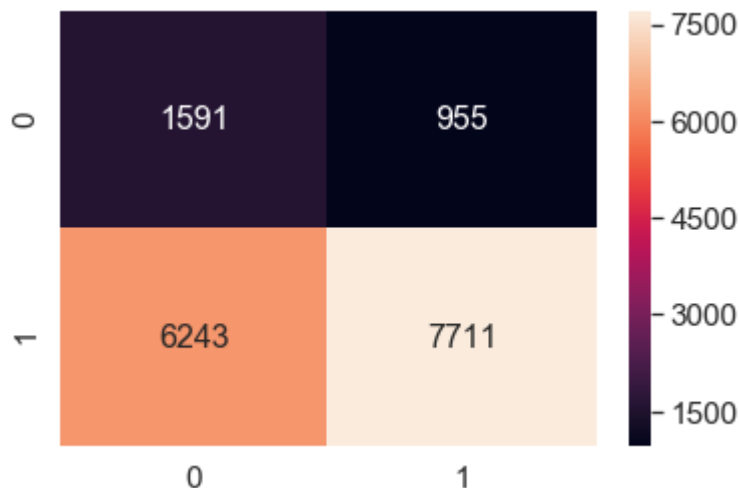
Out[105]: <matplotlib.axes.\_subplots.AxesSubplot at 0x19785933240>



```
In [106]: conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of  $tpr*(1-fpr)$  0.3453216161597511 for threshold 0.804

Out[106]: <matplotlib.axes.\_subplots.AxesSubplot at 0x197853de940>



```
In [107]: print(train_fpr.shape)
print(train_tpr.shape)
print(len(y_train_pred))

(29,)
(29,)
22445
```

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

```
In [108]: # Please write all the code with proper documentation
```

```
In [60]: # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train,teacher_prefix_one_hot_train
 ,project_grade_category_one_hot_train,price_standardized_train,
 quantity_standardized_train,
 teacher_number_of_previously_posted_projects_standardized_train,
 text_tfidf_train,title_tfidf_train)).tocsr()
X_cr = hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv,teacher_prefix_one_hot_cv
 ,project_grade_category_one_hot_cv,price_standardized_cv,quantity_standardized_cv
 ,teacher_number_of_previously_posted_projects_standardized_cv,
 text_tfidf_cv,title_tfidf_cv)).tocsr()
X_te = hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test,teacher_prefix_one_hot_test
 ,project_grade_category_one_hot_test,price_standardized_test,quantity_standardized_test
 ,teacher_number_of_previously_posted_projects_standardized_test,
 text_tfidf_test,title_tfidf_test)).tocsr()

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

Final Data matrix on TFIDF

```
(22445, 10229) (22445,)
(11055, 10229) (11055,)
(16500, 10229) (16500,)
```

```
=====
=====
```

## 1.1 Method 1: Simple for loop (if you are having memory limitations use this)

```
In [61]: def batch_predict(clf, data):
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs

 y_data_pred = []
 tr_loop = data.shape[0] - data.shape[0]%1000
 # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
 # in this for loop we will iterate until the last 1000 multiplier
 for i in range(0, tr_loop, 1000):
 y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
 # we will be predicting for the last data points
 y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
 return y_data_pred
```

```

In [114]: 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 non-thresholded measure of decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

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

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

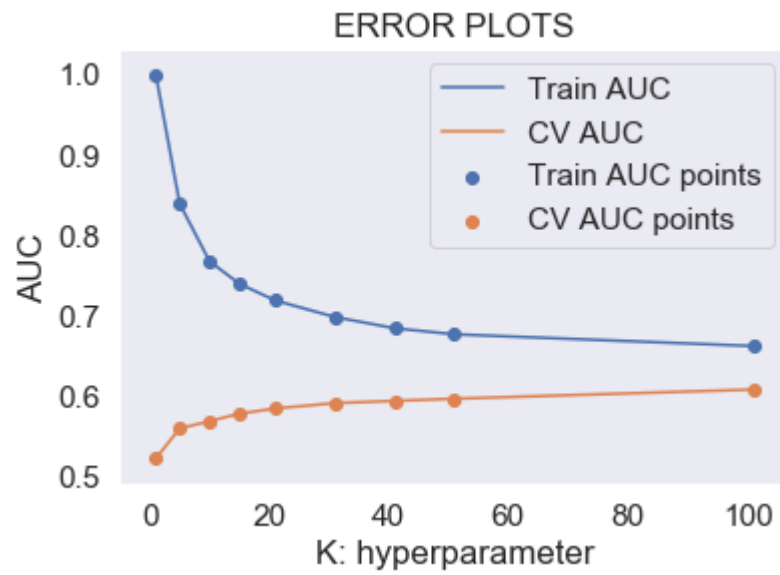
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs
 train_auc.append(roc_auc_score(y_train[:, :], y_train_pred))
 cv_auc.append(roc_auc_score(y_cv[:, :], y_cv_pred))

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

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

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

```



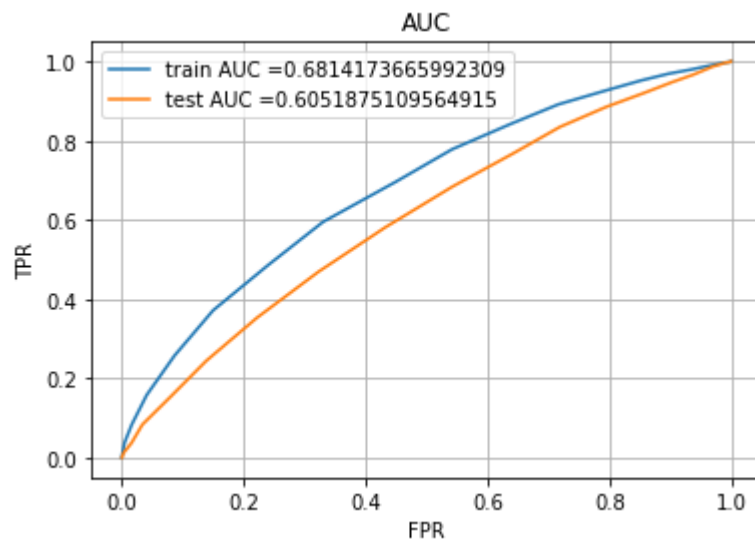
```
In [64]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.neighbors import KNeighborsClassifier

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

y_train_pred = batch_predict(neigh, X_tr[:, :])
y_test_pred = batch_predict(neigh, X_te[:, :])

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

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



```
In [65]: # we are writing our own function for predict, with defined threshold
we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

 t = threshold[np.argmax(tpr*(1-fpr))]

 # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

 print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
 predictions = []
 for i in proba:
 if i>=t:
 predictions.append(1)
 else:
 predictions.append(0)
 return predictions
```

```
In [66]: print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

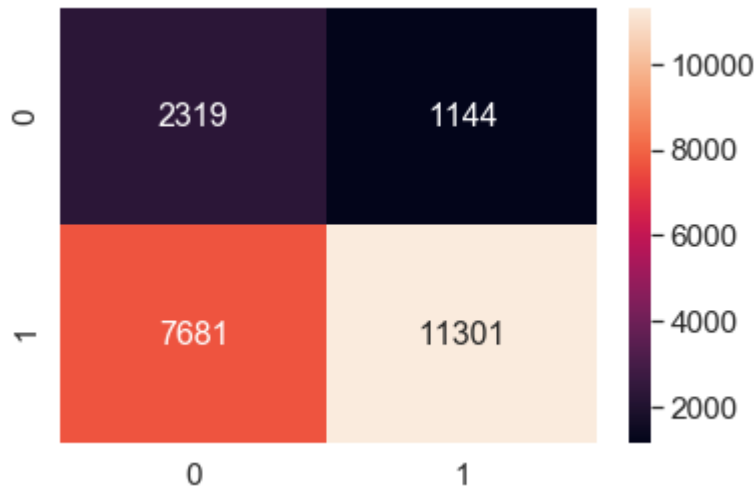
```
=====
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.3986788188746559 for threshold 0.855
[[2319 1144]
 [7681 11301]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.33006483202973835 for threshold 0.855
[[1445 1101]
 [5839 8115]]
```



```
In [67]: conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of  $tpr*(1-fpr)$  0.3986788188746559 for threshold 0.855

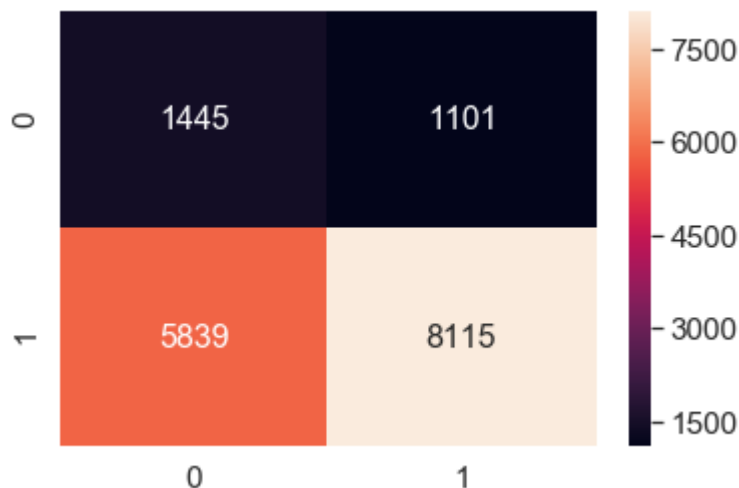
Out[67]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1d709a211d0>



```
In [68]: conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of  $tpr*(1-fpr)$  0.33006483202973835 for threshold 0.855

Out[68]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1d71169c780>



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

```
In [69]: # Please write all the code with proper documentation
```

```
In [70]: # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((categories_one_hot_train, sub_categories_one_hot_train, school_state_one_hot_train, teacher_prefix_one_hot_train
 , project_grade_category_one_hot_train, price_standardized_train,
 quantity_standardized_train,
 teacher_number_of_previously_posted_projects_standardized_train,
 avg_w2v_essays_vectors_train,
 avg_w2v_titles_vectors_train)).tocsr()
X_cr = hstack((categories_one_hot_cv, sub_categories_one_hot_cv, school_state_one_hot_cv, teacher_prefix_one_hot_cv
 , project_grade_category_one_hot_cv, price_standardized_cv, quantity_standardized_cv
 , teacher_number_of_previously_posted_projects_standardized_cv,
 avg_w2v_essays_vectors_cv,
 avg_w2v_titles_vectors_cv)).tocsr()
X_te = hstack((categories_one_hot_test, sub_categories_one_hot_test, school_state_one_hot_test, teacher_prefix_one_hot_test
 , project_grade_category_one_hot_test, price_standardized_test, quantity_standardized_test
 , teacher_number_of_previously_posted_projects_standardized_test,
 avg_w2v_essays_vectors_test,
 avg_w2v_titles_vectors_test)).tocsr()

print("Final Data matrix on AVGW2V")
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
print("="*100)
```

Final Data matrix on AVGW2V

(22445, 702) (22445,)

(11055, 702) (11055,)

(16500, 702) (16500,)

=====

=====

## 1.1 Method 1: Simple for loop (if you are having memory limitations use this)

```
In [71]: def batch_predict(clf, data):
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs

 y_data_pred = []
 tr_loop = data.shape[0] - data.shape[0]%1000
 # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
 # in this for loop we will iterate until the last 1000 multiplier
 for i in range(0, tr_loop, 1000):
 y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
 # we will be predicting for the last data points
 y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
 return y_data_pred
```

```

In [137]: 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 non-thresholded measure of decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

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

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

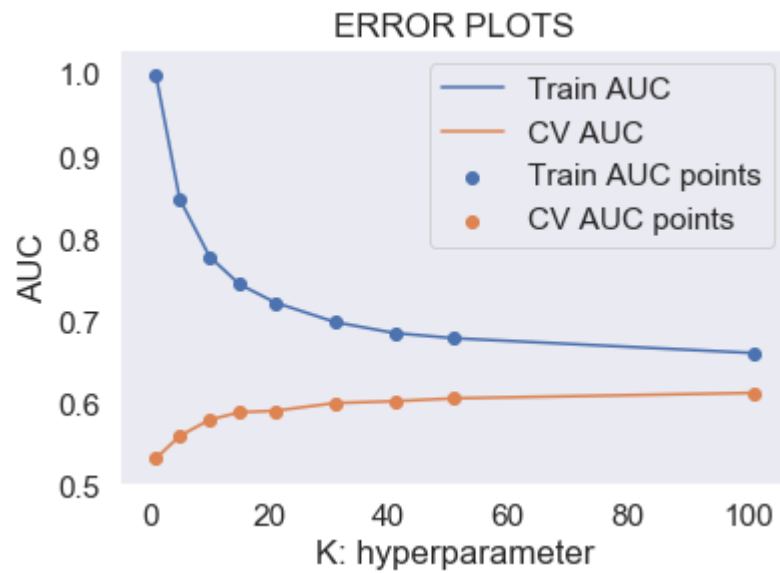
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs
 train_auc.append(roc_auc_score(y_train[:5747],y_train_pred))
 cv_auc.append(roc_auc_score(y_cv[:5747], y_cv_pred))

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

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
print("Execution time: " + str((time.time() - start_time)) + ' ms')

```



Execution time: 981.1687150001526 ms

```

In [72]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.
 #html#sklearn.metrics.roc_curve
 from sklearn.metrics import roc_curve, auc
 from sklearn.neighbors import KNeighborsClassifier

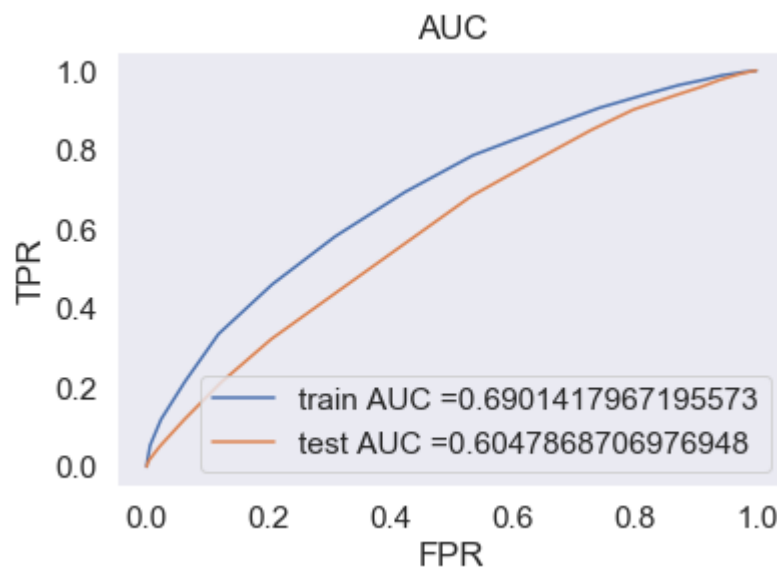
 start_time = time.time()
 neigh = KNeighborsClassifier(n_neighbors=53)
 neigh.fit(X_tr[:, :], y_train[:])
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates
 # of the positive class
 # not the predicted outputs

 y_train_pred = batch_predict(neigh, X_tr[:, :])
 y_test_pred = batch_predict(neigh, X_te[:])

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

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

```



Execution time: 1464.733214378357 ms

In [ ]:

```
In [76]: # we are writing our own function for predict, with defined threshold
we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

 t = threshold[np.argmax(tpr*(1-fpr))]

 # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

 print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
 predictions = []
 for i in proba:
 if i>=t:
 predictions.append(1)
 else:
 predictions.append(0)
 return predictions
```

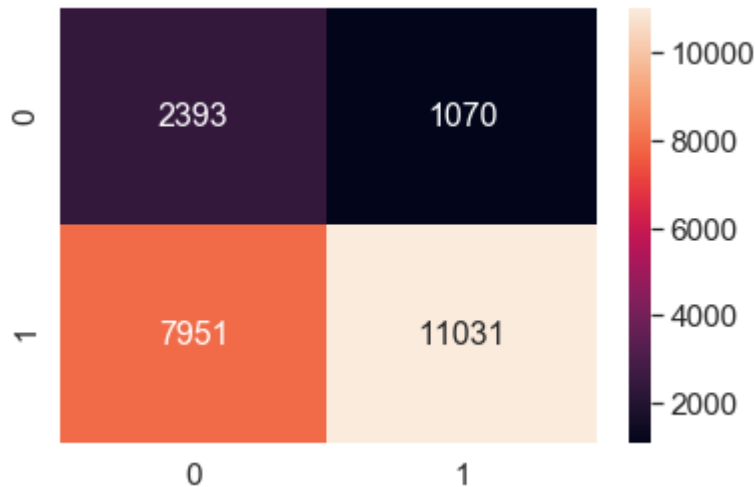
```
In [77]: print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

```
=====
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.40157172168487176 for threshold 0.868
[[2393 1070]
 [7951 11031]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.32514391073531806 for threshold 0.868
[[1453 1093]
 [6004 7950]]
```

```
In [78]: conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.40157172168487176 for threshold 0.868

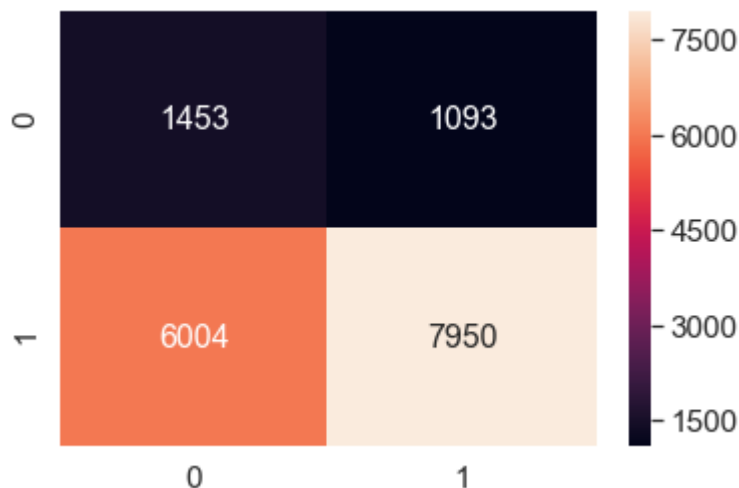
Out[78]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1d712735048>



```
In [79]: conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.32514391073531806 for threshold 0.868

Out[79]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1d711ab6908>



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

```
In [80]: # Please write all the code with proper documentation
```



```
In [81]: # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((categories_one_hot_train, sub_categories_one_hot_train, school_state_one_hot_train, teacher_prefix_one_hot_train
 , project_grade_category_one_hot_train, price_standardized_train,
 , teacher_number_of_previously_posted_projects_standardized_train, tfidf_w2v_essays_vectors_train
 , tfidf_w2v_titles_vectors_train)).tocsr()
X_cr = hstack((categories_one_hot_cv, sub_categories_one_hot_cv, school_state_one_hot_cv, teacher_prefix_one_hot_cv
 , project_grade_category_one_hot_cv, price_standardized_cv, quantity_standardized_cv
 , teacher_number_of_previously_posted_projects_standardized_cv, tfidf_w2v_essays_vectors_cv
 , tfidf_w2v_titles_vectors_cv)).tocsr()
X_te = hstack((categories_one_hot_test, sub_categories_one_hot_test, school_state_one_hot_test, teacher_prefix_one_hot_test
 , project_grade_category_one_hot_test, price_standardized_test, quantity_standardized_test
 , teacher_number_of_previously_posted_projects_standardized_test, tfidf_w2v_essays_vectors_test
 , tfidf_w2v_titles_vectors_test)).tocsr()

print("Final Data matrix on TFIDF W2V")
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
print("="*100)
```

Final Data matrix on TFIDF W2V

(22445, 702) (22445,)

(11055, 702) (11055,)

(16500, 702) (16500,)

=====

## 1.1 Method 1: Simple for loop (if you are having memory limitations use this)

```
In [83]: def batch_predict(clf, data):
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs

 y_data_pred = []
 tr_loop = data.shape[0] - data.shape[0]%1000
 # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
 # in this for loop we will iterate until the last 1000 multiplier
 for i in range(0, tr_loop, 1000):
 y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
 # we will be predicting for the last data points
 y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
 return y_data_pred
```

```

In [150]: import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score
"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

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

"""

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

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

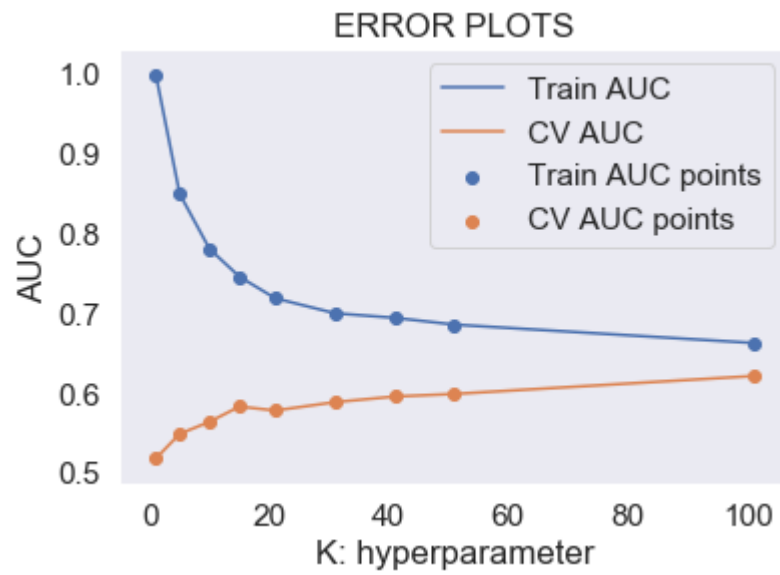
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs
 train_auc.append(roc_auc_score(y_train[:5747],y_train_pred))
 cv_auc.append(roc_auc_score(y_cv[:5747], y_cv_pred))

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

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
print("Execution time: " + str((time.time() - start_time)) + ' ms')

```



Execution time: 996.8478593826294 ms

```

In [84]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.
 #html#sklearn.metrics.roc_curve
 from sklearn.metrics import roc_curve, auc

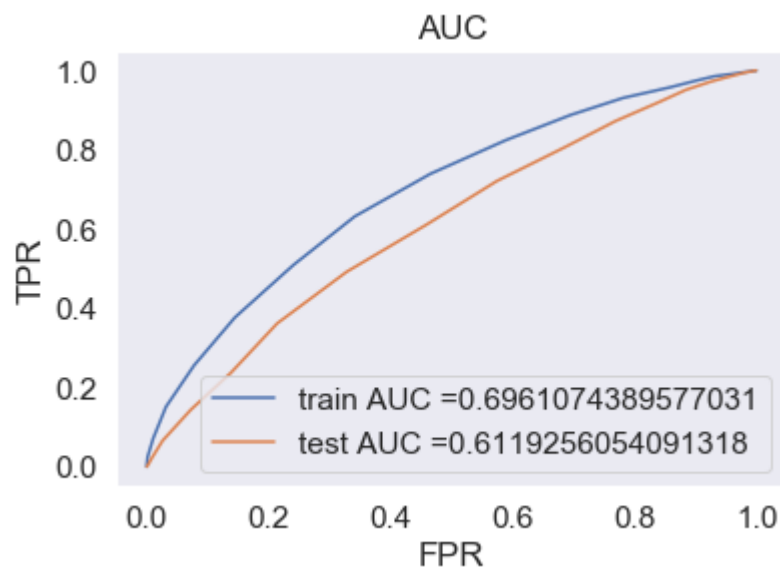
 start_time = time.time()
 neigh = KNeighborsClassifier(n_neighbors=50)
 neigh.fit(X_tr[:, :], y_train[:])
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates
 # of the positive class
 # not the predicted outputs

 y_train_pred = batch_predict(neigh, X_tr[:, :])
 y_test_pred = batch_predict(neigh, X_te[:, :])

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

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

```



Execution time: 1695.6296381950378 ms

```
In [85]: # we are writing our own function for predict, with defined threshold
we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

 t = threshold[np.argmax(tpr*(1-fpr))]

 # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

 print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
 predictions = []
 for i in proba:
 if i>=t:
 predictions.append(1)
 else:
 predictions.append(0)
 return predictions
```

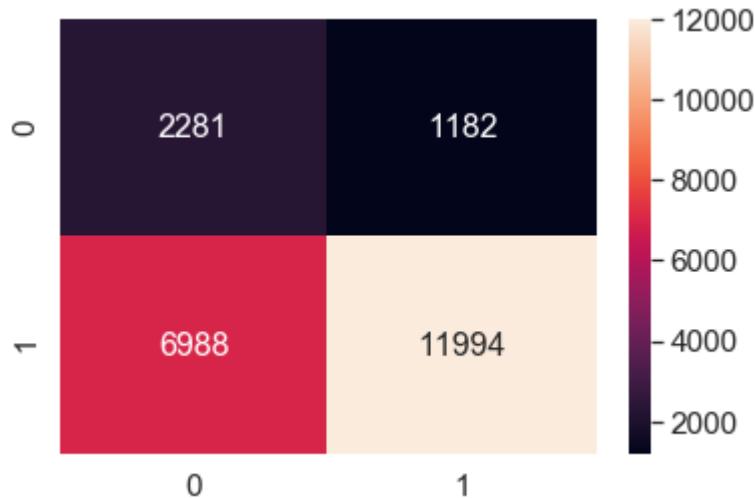
```
In [86]: print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

```
=====
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.41619309361060725 for threshold 0.86
[[2281 1182]
 [6988 11994]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.3307163104988324 for threshold 0.86
[[1380 1166]
 [5440 8514]]
```

```
In [87]: conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```

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

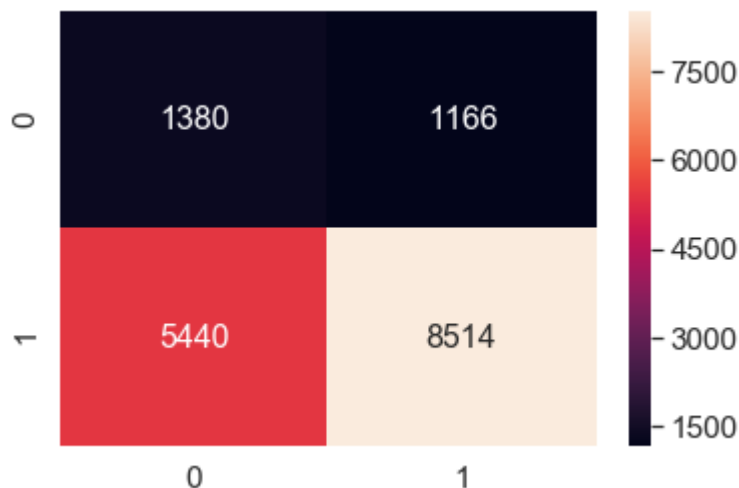
Out[87]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1d711d98c18>



```
In [88]: conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

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

Out[88]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1d712723c18>



## 2.5 Feature selection with `SelectKBest`

```
In [89]: # please write all the code with proper documentation, and proper titles for each subsection
go through documentations and blogs before you start coding
first figure out what to do, and then think about how to do.
reading and understanding error messages will be very much helpfull in debugging your code

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

```
In [52]: # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((categories_one_hot_train, sub_categories_one_hot_train, school_state_one_hot_train, teacher_prefix_one_hot_train,
 , project_grade_category_one_hot_train, price_standardized_train, quantity_standardized_train,
 , teacher_number_of_previously_posted_projects_standardized_train, text_tfidf_train, title_tfidf_train)).tocsr().toarray()
X_cr = hstack((categories_one_hot_cv, sub_categories_one_hot_cv, school_state_one_hot_cv, teacher_prefix_one_hot_cv,
 , project_grade_category_one_hot_cv, price_standardized_cv, quantity_standardized_cv,
 , teacher_number_of_previously_posted_projects_standardized_cv, text_tfidf_cv, title_tfidf_cv)).tocsr().toarray()
X_te = hstack((categories_one_hot_test, sub_categories_one_hot_test, school_state_one_hot_test, teacher_prefix_one_hot_test,
 , project_grade_category_one_hot_test, price_standardized_test, quantity_standardized_test,
 , teacher_number_of_previously_posted_projects_standardized_test, text_tfidf_test, title_tfidf_test)).tocsr().toarray()
```



```

In [53]: #####
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
X_tr = scaler.fit_transform(X_tr,y_train)
X_cr = scaler.transform(X_cr)
X_te = scaler.transform(X_te)
#####
from sklearn.feature_selection import SelectKBest, chi2
t = SelectKBest(chi2,k=2000).fit(X_tr, y_train)
X_tr = t.transform(X_tr)
X_te = t.transform(X_te)
X_cr = t.transform(X_cr)
#####
print("Final Data matrix on TFIDF")
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
print("="*100)

```

Final Data matrix on TFIDF

(22445, 2000) (22445,)

(11055, 2000) (11055,)

(16500, 2000) (16500,)

=====

```

In [93]: # #####
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
X_tr = scaler.fit_transform(X_tr)
X_cr = scaler.fit_transform(X_cr)
X_te = scaler.fit_transform(X_te)
#####
from sklearn.feature_selection import SelectKBest, chi2
t = SelectKBest(chi2,k=2000)
X_tr = t.fit_transform(X_tr,y_train)
X_te = t.fit_transform(X_te,y_test)
X_cr = t.fit_transform(X_cr,y_cv)
#####
print("Final Data matrix on TFIDF")
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
print("="*100)

```

## 1.1 Method 1: Simple for loop (if you are having memory limitations use this)

```
In [56]: def batch_predict(clf, data):
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs

 y_data_pred = []
 tr_loop = data.shape[0] - data.shape[0]%1000
 # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
 # in this for loop we will iterate until the last 1000 multiplier
 for i in range(0, tr_loop, 1000):
 y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
 # we will be predicting for the last data points
 y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
 return y_data_pred
```

```

In [56]: 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 non-thresholded measure of decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

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

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

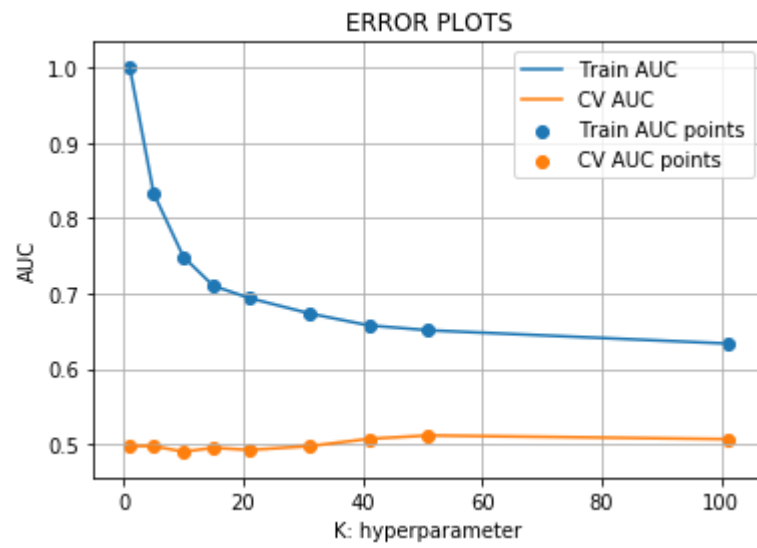
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs
 train_auc.append(roc_auc_score(y_train[:5747],y_train_pred))
 cv_auc.append(roc_auc_score(y_cv[:5747], y_cv_pred))

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

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
print("Execution time: " + str((time.time() - start_time)) + ' ms')

```



Execution time: 2072.0459604263306 ms

```

In [57]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.
 #html#sklearn.metrics.roc_curve
 from sklearn.metrics import roc_curve, auc
 from sklearn.neighbors import KNeighborsClassifier

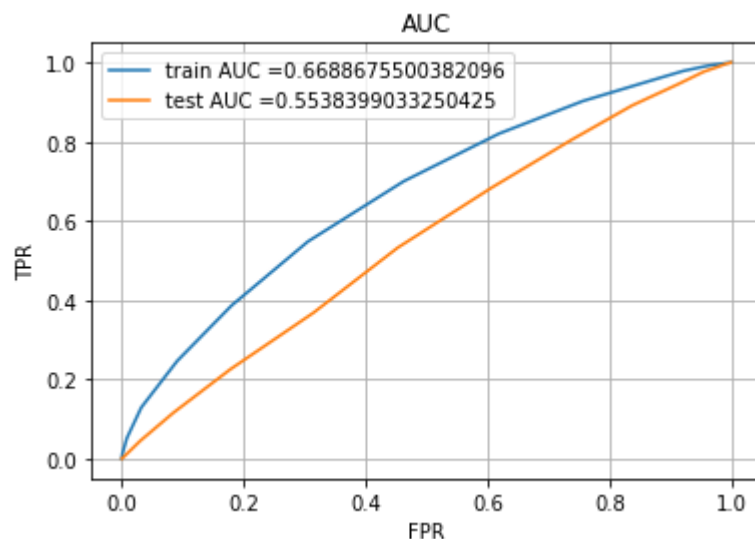
 start_time = time.time()
 neigh = KNeighborsClassifier(n_neighbors=49)
 neigh.fit(X_tr[:, :], y_train[:])
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates
 # of the positive class
 # not the predicted outputs

 y_train_pred = batch_predict(neigh, X_tr[:, :])
 y_test_pred = batch_predict(neigh, X_te[:])

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

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

```



Execution time: 3397.581614255905 ms

```
In [58]: # we are writing our own function for predict, with defined threshold
we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

 t = threshold[np.argmax(tpr*(1-fpr))]

 # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

 print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
 predictions = []
 for i in proba:
 if i>=t:
 predictions.append(1)
 else:
 predictions.append(0)
 return predictions
```

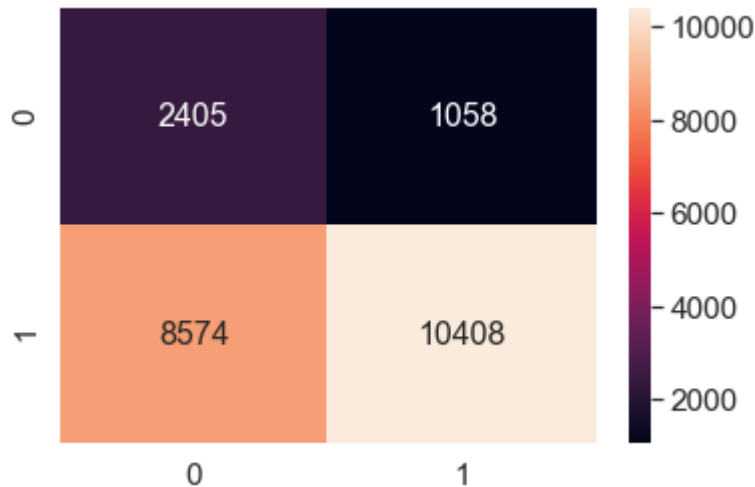
```
In [59]: print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

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

```
In [60]: conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of  $tpr*(1-fpr)$  0.38079207704501006 for threshold 0.878

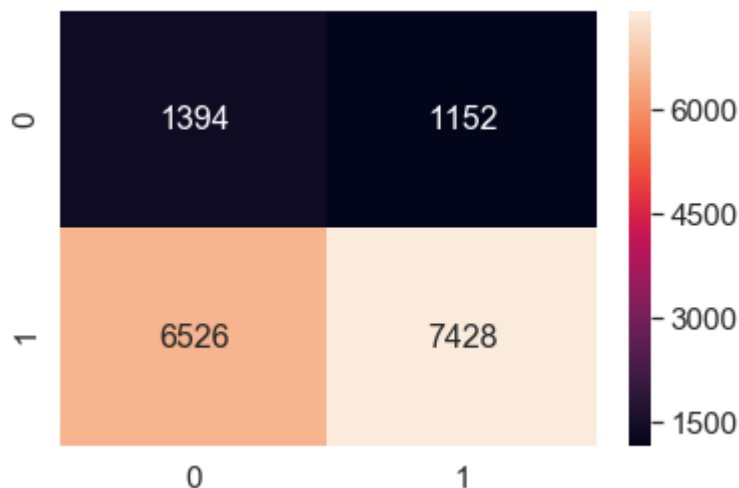
Out[60]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1f589e16a20>



```
In [61]: conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of  $tpr*(1-fpr)$  0.2914590539378573 for threshold 0.878

Out[61]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1f58c005518>



### 3. Conclusions

```
In []: # Please compare all your models using Prettytable Library
```

```
In [1]: from prettytable import PrettyTable
#If you get a ModuleNotFoundError error , install prettytable using: pip3 inst
all prettytable
x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Hyper Parameter", "AUC"]
x.add_row(["BOW", "Brute", 51, 0.62])
x.add_row(["TFIDF", "Brute", 55, 0.60])
x.add_row(["AVG W2V", "Brute", 53, 0.60])
x.add_row(["TFIDF W2V", "Brute", 50, 0.61])
x.add_row(["TFIDF", "Top 2000", 49, 0.55])
print(x)
```

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

```
In []:
```