

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
<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
<code>project_submitted_datetime</code>	Datetime when project application was submitted. Example: 2016-04-12:43:56
<code>teacher_id</code>	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c1
<code>teacher_prefix</code>	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> • • • • •
<code>teacher_number_of_previously_posted_projects</code>	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
<code>id</code>	A <code>project_id</code> value from the <code>train.csv</code> file. Example: p036502
<code>description</code>	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
<code>quantity</code>	Quantity of the resource required. Example: 3
<code>price</code>	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
<code>project_is_approved</code>	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 [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
import os

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

1.1 Reading Data

```
In [2]: project_data = pd.read_csv('train_data.csv', nrows = 80000) # Taking 80000 data points
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 (80000, 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 [4]: # 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[4]:

	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 [5]: 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[5]:

	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 [6]: 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 [7]: 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 [8]: # 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 [9]: project_data.head(2)
```

Out[9]:

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

1.4.2.3 Using Pretrained Models: TFIDF weighted W2V

```
In [10]: # 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[79999])
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.

=====

"Why can't I play the drums, Mr. Reyes?" I hear this one often. I want my students to experience "percussion" as opposed to "the drums." We need some percussion instruments in order to help make this dream come to fruition. We have a diverse population with almost entirely "minority students" in attendance, including many English Learners (several who barely speak any English). We have great staff members that love the kids and support them in all aspects of their education. Students' safety is the primary interest at our school. We hope to provide a well-rounded music education where band is "cool" and kids look forward to coming to class. We are finding ourselves in need of instruments to help keep the momentum of our program going. My students will be able to experience more advanced pieces of music with this instrument. I am asking for a high quality instrument so that it may be used and taken care of over time. Eventually, I'd love to start a percussion ensemble. This investment will greatly help the percussion section at my school! My students come from very troubled homes and they deal with stress that people shouldn't have to experience. They know adversity more than the average person and have to overcome it in whatever ways they can manage. Music can be an amazing outlet for these kids. Your generous support could help change a life for a young troubled youth. Thank you for your help!

=====

I have over 80 students on my ESL roster that need my services during the school year. They come in at varying levels of English acquisition, so differentiation is very important in my classroom. Technology is only becoming more important, so giving my students this chance to be hands-on is wonderful. My school has over 740 students; about 80% are of a minority group and roughly 60% of them are English Language Learners. There are more than 7 languages represented, some mainstream like Spanish or Vietnamese, and some indigenous like Mam. The neighborhood where my school is located is very low income, with some families being in poverty or homeless. With all of these things working against them, you will not meet a group of students who are more dedicated to learning. They are excited to try new things. My school works very hard to give the students as much access to print (good literature) and hands-on activities as we can. If my project is funded, my students will be able to use the iPad Minis during class time in several different ways. They will use them to look up wo

rds they don't understand, or to find synonyms and antonyms to help broaden their vocabularies. There are so many apps available that help English Language Learners practice and apply what they are learning in their classrooms. My students will be able to use these apps as hands-on practice when our lessons are complete, or as a reward for a job well done. As the majority of our students are low income, being able to work with technology in school can help them to become more future-ready. This project will make a difference in my classroom because it gives my students another avenue for learning. It gets technology into their hands that they don't have access to normally, making them more future-ready. It will give them a hands-on learning tool to make the material more meaningful to them. It will give them a sense of ownership over the material and be an encouragement for learning a new language.

=====

In my classroom, I have the opportunity to make a difference in the lives of children who face challenges associated with limited English language proficiency and socioeconomic status. Most of the students at my school walk into our classrooms after only being in the country for a few days-Imagine walking into a room and seeing a book for the first time.. or a pencil. A smiling teacher is talking to you in a language you don't understand-Many of our other students are homeless or being raised in a single parent household and receive a free lunch based on their socioeconomic status. These are the students I teach every day.\r\nDespite these challenges, I believe that every single student can learn, grow and accomplish amazing things. I can not control what happens outside of the classroom, however, I can certainly control what happens inside the classroom. I will do everything in my power to give my students the best education and most positive and nurturing learning environment. I am hopeful to inspire our earliest learners to continue on a path to academic excellence. My job is to not only teach my students how to read, but to also get them excited about it! A major part in motivating the students is allowing them to choose where they sit in the classroom as they read-no one wants to sit at their desk every day! The problem with this, however, is our current classroom library has two pillows and one bean bag. This makes the option of getting comfy while reading around the room very limited. \r\nDonations for this project will allow more of my students the option to be comfortable while they read. My students will be able to sit on the soft pillows not only in the library, but also around the room. My students will also enjoy the soft plush animals to snuggle as they read. With your help, I can transform our reading time into something every student will enjoy!nannan

=====

My first graders are eager to learn about the world around them. They come to school each day full of enthusiasm and genuinely love learning. \r\n\r\nOur diverse class includes students from a variety of cultural and economic backgrounds. Many come from homes where parents can't afford or simply don't know the importance of books, so it is important to me to provide an environment that is rich in literature so that students learn to love reading. I want my students to be lifelong learners, and reading is the best way! I have used these magazines in the past, and kids absolutely LOVE them!! The topics are of high interest for children and always correspond to real world issues that are important for kids to learn. The subscription also includes online resources such as videos, printable worksheets, and skill-based games. \r\n\r\nThese materials will expose students to rigorous and interesting nonfiction text that will spark their curiosity about the world around them. The topics allow me to teach the nonfiction text standards using interesting materials. They always lead to engaging discussions and inspire students to find additional information about the various topics.nannan

=====

```
In [11]: # 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 [12]: sent = decontracted(project_data['essay'].values[20000])
print(sent)
print("="*50)
```

In my classroom, I have the opportunity to make a difference in the lives of children who face challenges associated with limited English language proficiency and socioeconomic status. Most of the students at my school walk into our classrooms after only being in the country for a few days-Imagine walking into a room and seeing a book for the first time.. or a pencil. A smiling teacher is talking to you in a language you do not understand-Many of our other students are homeless or being raised in a single parent household and receive a free lunch based on their socioeconomic status. These are the students I teach every day.\r\nDespite these challenges, I believe that every single student can learn, grow and accomplish amazing things. I can not control what happens outside of the classroom, however, I can certainly control what happens inside the classroom. I will do everything in my power to give my students the best education and most positive and nurturing learning environment. I am hopeful to inspire our earliest learners to continue on a path to academic excellence.My job is to not only teach my students how to read, but to also get them excited about it! A major part in motivating the students is allowing them to choose where they sit in the classroom as they read-no one wants to sit at their desk every day! The problem with this, however, is our current classroom library has two pillows and one bean bag. This makes the option of getting comfy while reading around the room very limited. \r\nDonations for this project will allow more of my students the option to be comfortable while they read. My students will be able to sit on the soft pillows not only in the library, but also around the room. My students will also enjoy the soft plush animals to snuggle as they read. With your help, I can transform our reading time into something every student will enjoy!nannan

=====

```
In [13]: # \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)
```

In my classroom, I have the opportunity to make a difference in the lives of children who face challenges associated with limited English language proficiency and socioeconomic status. Most of the students at my school walk into our classrooms after only being in the country for a few days-Imagine walking into a room and seeing a book for the first time.. or a pencil. A smiling teacher is talking to you in a language you do not understand-Many of our other students are homeless or being raised in a single parent household and receive a free lunch based on their socioeconomic status. These are the students I teach every day. Despite these challenges, I believe that every single student can learn, grow and accomplish amazing things. I can not control what happens outside of the classroom, however, I can certainly control what happens inside the classroom. I will do everything in my power to give my students the best education and most positive and nurturing learning environment. I am hopeful to inspire our earliest learners to continue on a path to academic excellence. My job is to not only teach my students how to read, but to also get them excited about it! A major part in motivating the students is allowing them to choose where they sit in the classroom as they read-no one wants to sit at their desk every day! The problem with this, however, is our current classroom library has two pillows and one bean bag. This makes the option of getting comfy while reading around the room very limited. Donations for this project will allow more of my students the option to be comfortable while they read. My students will be able to sit on the soft pillows not only in the library, but also around the room. My students will also enjoy the soft plush animals to snuggle as they read. With your help, I can transform our reading time into something every student will enjoy!nannan

```
In [14]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

In my classroom I have the opportunity to make a difference in the lives of children who face challenges associated with limited English language proficiency and socioeconomic status. Most of the students at my school walk into our classrooms after only being in the country for a few days. Imagine walking into a room and seeing a book for the first time or a pencil. A smiling teacher is talking to you in a language you do not understand. Many of our other students are homeless or being raised in a single parent household and receive a free lunch based on their socioeconomic status. These are the students I teach every day. Despite these challenges, I believe that every single student can learn, grow, and accomplish amazing things. I can not control what happens outside of the classroom, however, I can certainly control what happens inside the classroom. I will do everything in my power to give my students the best education and most positive and nurturing learning environment. I am hopeful to inspire our earliest learners to continue on a path to academic excellence. My job is to not only teach my students how to read but to also get them excited about it. A major part in motivating the students is allowing them to choose where they sit in the classroom as they read. No one wants to sit at their desk every day. The problem with this, however, is our current classroom library has two pillows and one bean bag. This makes the option of getting comfy while reading around the room very limited. Donations for this project will allow more of my students the option to be comfortable while they read. My students will be able to sit on the soft pillows not only in the library but also around the room. My students will also enjoy the soft plush animals to snuggle as they read. With your help, I can transform our reading time into something every student will enjoy.

```
In [15]: # 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 [16]: # 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('\n', ' ')
    sent = sent.replace('\t', ' ')
    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%|████████████████████████████████████████████████████████████████████████████████|
80000/80000 [00:38<00:00, 2087.22it/s]
```



```
In [17]: # after preprocessing
preprocessed_essays[20000]
```

```
Out[17]: 'classroom opportunity make difference lives children face challenges associa
ted limited english language proficiency socioeconomic status students school
walk classrooms country days imagine walking room seeing book first time penc
il smiling teacher talking language not understand many students homeless rai
sed single parent household receive free lunch based socioeconomic status stu
dents teach every day despite challenges believe every single student learn g
row accomplish amazing things not control happens outside classroom however c
ertainly control happens inside classroom everything power give students best
education positive nurturing learning environment hopeful inspire earliest le
arners continue path academic excellence job not teach students read also get
excited major part motivating students allowing choose sit classroom read no
one wants sit desk every day problem however current classroom library two pi
llows one bean bag makes option getting comfy reading around room limited don
ations project allow students option comfortable read students able sit soft
pillows not library also around room students also enjoy soft plush animals s
nuggle read help transform reading time something every student enjoy nannan'
```

1.4 Preprocessing of `project_title`

```
In [18]: # similarly you can preprocess the titles also
def preprocess_text_func(text_data):
    sent = decontracted(text_data)
    sent = sent.replace('\r', ' ')
    sent = sent.replace('\t', ' ')
    sent = sent.replace('\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    return sent.lower()
```

```
In [19]: preprocessed_titles = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_title'].values):
    preprocessed_titles.append(preprocess_text_func(sentence))
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
80000/80000 [00:01<00:00, 48061.09it/s]
```

```
In [20]: print(preprocessed_titles[5:12])
```

```
['breakout box ignite engagement', '21st century learning multimedia', 'ipad
learners', 'a flexible classroom flexible minds', 'make powerful movies', 'ro
bots taking 2nd grade', 'time kids to learn about science']
```

```
In [21]: print(project_data["project_title"].values[5:12])

['Breakout Box to Ignite Engagement!'
 '21st Century Learning with Multimedia' 'iPad for Learners'
 'A flexible classroom for flexible minds!' 'Make Powerful Movies!!'
 'Robots are Taking over 2nd Grade'
 'Time for Kids....To Learn About Science and more!']
```

1.5 Preparing data for models

```
In [22]: project_data.columns

Out[22]: Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
               'Date', 'project_grade_category', 'project_title', 'project_essay_1',
               'project_essay_2', 'project_essay_3', 'project_essay_4',
               'project_resource_summary',
               'teacher_number_of_previously_posted_projects', 'project_is_approved',
               'clean_categories', 'clean_subcategories', 'essay'],
              dtype='object')
```

we are going to consider

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

1.5.1 Vectorizing Categorical data

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

```
In [23]: # 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 (80000, 9)
```

```
In [24]: # we use count vectorizer to convert the values into one
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
sub_categories_one_hot = vectorizer.fit_transform(project_data['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encoding ", sub_categories_one_hot.shape)
```

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement',
'Civics_Government', 'Extracurricular', 'ForeignLanguages', 'NutritionEducation',
'Warmth', 'Care_Hunger', 'SocialSciences', 'PerformingArts', 'CharacterEducation',
'TeamSports', 'Other', 'College_CareerPrep', 'History_Geography',
'Music', 'Health_LifeScience', 'EarlyDevelopment', 'ESL', 'Gym_Fitness', 'EnvironmentalScience',
'VisualArts', 'Health_Wellness', 'AppliedSciences', 'SpecialNeeds', 'Literature_Writing',
'Mathematics', 'Literacy']
Shape of matrix after one hot encoding (80000, 30)
```

```
In [25]: # you can do the similar thing with state, teacher_prefix and project_grade_category also
# function to perform one hot encoding
def perform_one_hot_encoding(listdata, category, fillnan_value=""):
    vectorizer = CountVectorizer(vocabulary=listdata, lowercase=False, binary=True)
    vectorizer.fit(project_data[category].fillna(fillnan_value).values)
    print(vectorizer.get_feature_names())
    print("="*50)
    return vectorizer.transform(project_data[category].fillna(fillnan_value).values)
```

```
In [26]: # One hot encoding for school state
countries_list = sorted(project_data["school_state"].value_counts().keys())
school_state_one_hot = perform_one_hot_encoding(countries_list, "school_state")
print("Shape of matrix after one hot encoding ", school_state_one_hot.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 (80000, 51)
```

```
In [27]: # Project_Grade_Category - replacing hyphens, spaces with Underscores
project_data['project_grade_category'] = project_data['project_grade_category']
.map({'Grades PreK-2': 'Grades_PreK_2',

'Grades 6-8' : 'Grades_6_8',

'Grades 3-5' : 'Grades_3_5',

'Grades 9-12' : 'Grades_9_12'})
project_data['teacher_prefix'] = project_data['teacher_prefix'].map({'Mrs.':
'Mrs', 'Ms.': 'Ms', 'Mr.' : 'Mr',

'Teacher'

: 'Teacher', 'Dr.' : 'Dr'})
```

```
In [28]: project_data['project_grade_category'].head(5)
```

```
Out[28]: 55660    Grades_PreK_2
76127      Grades_3_5
51140    Grades_PreK_2
473      Grades_PreK_2
41558      Grades_3_5
Name: project_grade_category, dtype: object
```

```
In [29]: project_data['teacher_prefix'].head(5)
```

```
Out[29]: 55660    Mrs
76127      Ms
51140    Mrs
473      Mrs
41558    Mrs
Name: teacher_prefix, dtype: object
```

```
In [30]: project_data['teacher_prefix'].isnull().values.any() # since there are null va
lues replacing them with most common value Mrs
```

```
Out[30]: True
```

```
In [31]: project_data["teacher_prefix"].fillna("Mrs", inplace=True)
```

```
In [32]: project_data['teacher_prefix'].isnull().values.any() # No more Null Values
```

```
Out[32]: False
```

```
In [33]: # One hot encoding for teacher_prefix
teacher_prefix_list = sorted(project_data["teacher_prefix"].value_counts().keys())
print(teacher_prefix_list)
teacher_prefix_one_hot = perform_one_hot_encoding(teacher_prefix_list, "teacher_prefix", "Mrs.")
print("Shape of matrix after one hot encoding ", teacher_prefix_one_hot.shape)

['Dr', 'Mr', 'Mrs', 'Ms', 'Teacher']
['Dr', 'Mr', 'Mrs', 'Ms', 'Teacher']
=====
Shape of matrix after one hot encoding (80000, 5)
```

```
In [34]: # One hot encoding for project_grade_category
grade_list = sorted(project_data["project_grade_category"].value_counts().keys())
grade_one_hot = perform_one_hot_encoding(grade_list, "project_grade_category")
print("Shape of matrix after one hot encoding ", grade_one_hot.shape)

['Grades_3_5', 'Grades_6_8', 'Grades_9_12', 'Grades_PreK_2']
=====
Shape of matrix after one hot encoding (80000, 4)
```

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

```
In [35]: # We are considering only the words which appeared in at least 10 documents(ros or projects).
vectorizer = CountVectorizer(min_df=10)
text_bow = vectorizer.fit_transform(preprocessed_essays)
print("Shape of matrix after one hot encoding ", text_bow.shape)

Shape of matrix after one hot encoding (80000, 14627)
```

```
In [36]: # you can vectorize the title also
# before you vectorize the title make sure you preprocess it
vectorizer_titles = CountVectorizer(min_df=10)
text_bow_titles = vectorizer_titles.fit_transform(preprocessed_titles)
print("Shape of matrix after one hot encoding ", text_bow_titles.shape)

Shape of matrix after one hot encoding (80000, 2740)
```

1.5.2.2 TFIDF vectorizer

```
In [37]: 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 encoding ",text_tfidf.shape)
```

Shape of matrix after one hot encoding (80000, 14627)

```
In [38]: # TFIDF Vectorizer for Preprocessed Title
vectorizer_titles = TfidfVectorizer(min_df=10)
text_tfidf_titles = vectorizer_titles.fit_transform(preprocessed_titles)
print("Shape of matrix after one hot encoding ",text_tfidf_titles.shape)
```

Shape of matrix after one hot encoding (80000, 2740)

1.5.2.3 Using Pretrained Models: Avg W2V

```

In [ ]: '''
# 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)

```


1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

```
In [42]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(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())
```

```
In [43]: # average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in th
is list
for sentence in tqdm(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_vectors.append(vector)

print(len(tfidf_w2v_vectors))
print(len(tfidf_w2v_vectors[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 80000/80000 [01:58<00:00, 672.90it/s]
```

```
80000
300
```

```
In [44]: # Similarly you can vectorize for title also
tfidf_w2v_vectors_titles = []; # the avg-w2v for each project_title is stored
in this list
for sentence in tqdm(preprocessed_titles): # for each project_title
    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_vectors_titles.append(vector)

print(len(tfidf_w2v_vectors_titles))
print(len(tfidf_w2v_vectors_titles[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
80000/80000 [00:01<00:00, 40167.22it/s]
```

```
80000
```

```
300
```

1.5.3 Vectorizing Numerical features

```
In [45]: price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'
}).reset_index()
project_data = pd.merge(project_data, price_data, on='id', how='left')
```

```
In [46]: # 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))
```

Mean : 299.16610437500003, Standard deviation : 375.7800454521539

```
In [47]: price_standardized
```

```
Out[47]: array([[ 1.13333292],
 [-0.22921947],
 [ 0.07939191],
 ...,
 [-0.08142025],
 [ 0.26567109],
 [-0.78148403]])
```

```
In [48]: # Vectorizing teacher_number_of_previously_posted_projects
teacher_number_of_previously_posted_projects_scalar = StandardScaler()
teacher_number_of_previously_posted_projects_scalar.fit(project_data['teacher_number_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_projects_scalar.var_[0])}")

# Now standardize the data with above mean and variance.
teacher_number_of_previously_posted_projects_standardized = teacher_number_of_previously_posted_projects_scalar.transform(project_data['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
```

Mean : 11.214825, Standard deviation : 27.961699254862445

1.5.4 Merging all the above features

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

```
In [49]: # Categorical
print(school_state_one_hot.shape)
print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(teacher_prefix_one_hot.shape)
print(grade_one_hot.shape)
print(text_bow_titles.shape)
print(text_bow.shape)
# Numerical
print(price_standardized.shape)
print(teacher_number_of_previously_posted_projects_standardized.shape)

(80000, 51)
(80000, 9)
(80000, 30)
(80000, 5)
(80000, 4)
(80000, 2740)
(80000, 14627)
(80000, 1)
(80000, 1)
```

```
In [50]: # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
# with the same hstack function we are concatinating a sparse matrix and a den
se matirx :)
X = hstack((school_state_one_hot, categories_one_hot, sub_categories_one_hot, t
eachar_prefix_one_hot,
           grade_one_hot, text_bow_titles, text_bow, price_standardized,
           teacher_number_of_previously_posted_projects_standardized))
X.shape
```

```
Out[50]: (80000, 17468)
```

```
In [51]: X.__class__
```

```
Out[51]: scipy.sparse.coo.coo_matrix
```

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

In [60]: `project_data.head(3)`

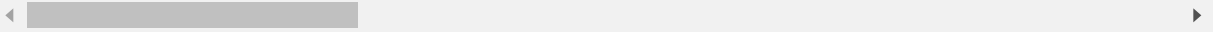
Out[60]:

|   | Unnamed: 0 | id      | teacher_id                       | teacher_prefix | school_state | Date                |
|---|------------|---------|----------------------------------|----------------|--------------|---------------------|
| 0 | 8393       | p205479 | 2bf07ba08945e5d8b2a3f269b2b3cfe5 | Mrs            | CA           | 2016-04-27 00:27:36 |
| 1 | 37728      | p043609 | 3f60494c61921b3b43ab61bdde2904df | Ms             | UT           | 2016-04-27 00:31:25 |
| 2 | 74477      | p189804 | 4a97f3a390bfe21b99cf5e2b81981c73 | Mrs            | CA           | 2016-04-27 00:46:53 |

```
In [61]: y = project_data['project_is_approved'].values
X = project_data.drop(['project_is_approved'], axis=1)
X.head(1)
```

Out[61]:

| Unnamed:<br>0 |      | id      | teacher_id                       | teacher_prefix | school_state | Date                |
|---------------|------|---------|----------------------------------|----------------|--------------|---------------------|
| 0             | 8393 | p205479 | 2bf07ba08945e5d8b2a3f269b2b3cfe5 | Mrs            | CA           | 2016-04-27 00:27:36 |



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

```
In [63]: # train test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify=y)
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33)
```

```
In [64]: print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)
```

```
(35912, 19) (35912,)
(17688, 19) (17688,)
(26400, 19) (26400,)
```

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

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

```
In [71]: # School State
vectorizer = CountVectorizer()
vectorizer.fit(X_train['school_state'].values) # fit has to happen only on tra
in data

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

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

After vectorizations

```
(35912, 51) (35912,)
(17688, 51) (17688,)
(26400, 51) (26400,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'i
a', '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']
=====
=====
```

```
In [72]: # teacher_prefix
vectorizer = CountVectorizer()
vectorizer.fit(X_train['teacher_prefix'].values) # fit has to happen only on t
rain data

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

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

After vectorizations

```
(35912, 5) (35912,)
(17688, 5) (17688,)
(26400, 5) (26400,)
['dr', 'mr', 'mrs', 'ms', 'teacher']
=====
=====
```



```
In [73]: # project_grade_category
vectorizer = CountVectorizer()
vectorizer.fit(X_train['project_grade_category'].values) # fit has to happen o
nly on train data

we use the fitted CountVectorizer to convert the text to vector
X_train_grade_ohe = vectorizer.transform(X_train['project_grade_category'].val
ues)
X_cv_grade_ohe = vectorizer.transform(X_cv['project_grade_category'].values)
X_test_grade_ohe = vectorizer.transform(X_test['project_grade_category'].value
s)

print("After vectorizations")
print(X_train_grade_ohe.shape, y_train.shape)
print(X_cv_grade_ohe.shape, y_cv.shape)
print(X_test_grade_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

```
After vectorizations
(35912, 4) (35912,)
(17688, 4) (17688,)
(26400, 4) (26400,)
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
=====
=====
```

```
In [77]: # categories
vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_categories'].values) # fit has to happen only on
train data

we use the fitted CountVectorizer to convert the text to vector
X_train_category_ohe = vectorizer.transform(X_train['clean_categories'].values
)
X_cv_category_ohe = vectorizer.transform(X_cv['clean_categories'].values)
X_test_category_ohe = vectorizer.transform(X_test['clean_categories'].values)

print("After vectorizations")
print(X_train_category_ohe.shape, y_train.shape)
print(X_cv_category_ohe.shape, y_cv.shape)
print(X_test_category_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

```
After vectorizations
(35912, 9) (35912,)
(17688, 9) (17688,)
(26400, 9) (26400,)
['appliedlearning', 'care_hunger', 'health_sports', 'history_civics', 'litera
cy_language', 'math_science', 'music_arts', 'specialneeds', 'warmth']
=====
=====
```

```
In [78]: # sub categories
vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_subcategories'].values) # fit has to happen only
on train data

we use the fitted CountVectorizer to convert the text to vector
X_train_subcategory_ohe = vectorizer.transform(X_train['clean_subcategories'].
values)
X_cv_subcategory_ohe = vectorizer.transform(X_cv['clean_subcategories'].values
)
X_test_subcategory_ohe = vectorizer.transform(X_test['clean_subcategories'].va
lues)

print("After vectorizations")
print(X_train_subcategory_ohe.shape, y_train.shape)
print(X_cv_subcategory_ohe.shape, y_cv.shape)
print(X_test_subcategory_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

```
After vectorizations
(35912, 30) (35912,)
(17688, 30) (17688,)
(26400, 30) (26400,)
['appliedsciences', 'care_hunger', 'charactereducation', 'civics_government',
'college_careerprep', 'communityservice', 'earlydevelopment', 'economics', 'e
nvironmentalscience', 'esl', 'extracurricular', 'financialliteracy', 'foreign
languages', 'gym_fitness', 'health_lifescience', 'health_wellness', 'history_
geography', 'literacy', 'literature_writing', 'mathematics', 'music', 'nutrit
ioneducation', 'other', 'parentinvolvement', 'performingarts', 'socialscience
s', 'specialneeds', 'teamsports', 'visualarts', 'warmth']
=====
=====
```

### Encoding numerical features

```
In [81]: from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['price'].values)
this will rise an error Expected 2D array, got 1D array instead:
array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
Reshape your data either using
array.reshape(-1, 1) if your data has a single feature
array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['price'].values.reshape(-1,1))

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

print("After vectorizations")
print(X_train_price_norm.shape, y_train.shape)
print(X_cv_price_norm.shape, y_cv.shape)
print(X_test_price_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

(35912, 1) (35912,)

(17688, 1) (17688,)

(26400, 1) (26400,)

=====

=====

```
In [82]: # teacher previously posted projects
normalizer = Normalizer()
normalizer.fit(X_train['price'].values)
this will rise an error Expected 2D array, got 1D array instead:
array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
Reshape your data either using
array.reshape(-1, 1) if your data has a single feature
array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.
reshape(-1,1))

X_train_teach_prev_norm = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
X_cv_teach_prev_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
X_test_teach_prev_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

print("After vectorizations")
print(X_train_teach_prev_norm.shape, y_train.shape)
print(X_cv_teach_prev_norm.shape, y_cv.shape)
print(X_test_teach_prev_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

(35912, 1) (35912,)

(17688, 1) (17688,)

(26400, 1) (26400,)

=====

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

```
In [65]: vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['essay'].values) # fit has to happen only on train data
```

```
Out[65]: CountVectorizer(analyzer='word', binary=False, decode_error='strict',
dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
lowercase=True, max_df=1.0, max_features=5000, min_df=10,
ngram_range=(1, 4), preprocessor=None, stop_words=None,
strip_accents=None, token_pattern='(?u)\\b\\w\\w+\\b',
tokenizer=None, vocabulary=None)
```

```
In [66]: # we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vectorizer.transform(X_train['essay'].values)
X_cv_essay_bow = vectorizer.transform(X_cv['essay'].values)
X_test_essay_bow = vectorizer.transform(X_test['essay'].values)
```

```
In [67]: print("After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
print(X_cv_essay_bow.shape, y_cv.shape)
print(X_test_essay_bow.shape, y_test.shape)
print("="*100)
```

```
After vectorizations
(35912, 5000) (35912,)
(17688, 5000) (17688,)
(26400, 5000) (26400,)
```

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

```
In [68]: # Preprocessing project_title
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['project_title'].values) # fit has to happen only on train data
```

```
Out[68]: CountVectorizer(analyzer='word', binary=False, decode_error='strict',
 dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
 lowercase=True, max_df=1.0, max_features=5000, min_df=10,
 ngram_range=(1, 4), preprocessor=None, stop_words=None,
 strip_accents=None, token_pattern='(?u)\\b\\w\\w+\\b',
 tokenizer=None, vocabulary=None)
```

```
In [69]: # we use the fitted CountVectorizer to convert the text to vector
X_train_pj_title_bow = vectorizer.transform(X_train['project_title'].values)
X_cv_pj_title_bow = vectorizer.transform(X_cv['project_title'].values)
X_test_pj_title_bow = vectorizer.transform(X_test['project_title'].values)
```

```
In [70]: print("After vectorizations")
print(X_train_pj_title_bow.shape, y_train.shape)
print(X_cv_pj_title_bow.shape, y_cv.shape)
print(X_test_pj_title_bow.shape, y_test.shape)
print("="*100)
```

```
After vectorizations
(35912, 4272) (35912,)
(17688, 4272) (17688,)
(26400, 4272) (26400,)
```

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

## 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

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

```
In [83]: # concatenating all the features
merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((X_train_essay_bow, X_train_state_ohe, X_train_teacher_ohe,
 X_train_grade_ohe, X_train_price_norm, X_train_category_ohe,
 X_train_subcategory_ohe, X_train_teach_prev_norm,
 X_train_pj_title_bow)).tocsr()

X_cr = hstack((X_cv_essay_bow, X_cv_state_ohe, X_cv_teacher_ohe,
 X_cv_grade_ohe, X_cv_category_ohe, X_cv_subcategory_ohe,
 X_cv_price_norm, X_cv_teach_prev_norm, X_cv_pj_title_bow)).tocsr()

X_te = hstack((X_test_essay_bow, X_test_state_ohe, X_test_teacher_ohe,
 X_test_grade_ohe, X_test_category_ohe, X_test_subcategory_ohe,
 X_test_price_norm, X_test_teach_prev_norm,
 X_test_pj_title_bow)).tocsr()
```

```
In [84]: print("Final Data matrix - for set 1")
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 - for set 1
(35912, 9373) (35912,)
(17688, 9373) (17688,)
(26400, 9373) (26400,)
=====
=====
```

```
In [85]: # Since there was memory errors while trying with RandomizedSearchCV, trying the for loop approach
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 tr_loop will be 49041 - 49041%1000 = 49000
 # in this for loop we will iterate until the last 1000 multiplier
 for i in range(0, tr_loop, 1000):
 y_data_pred.extend(clf.predict_proba(data[i:i+1000]))[:,1])
 # we will be predicting for the last data points
 if data.shape[0]%1000 != 0:
 y_data_pred.extend(clf.predict_proba(data[tr_loop:]))[:,1])

 return y_data_pred
```

```
In [86]: import matplotlib.pyplot as plt
 from sklearn.neighbors import KNeighborsClassifier
 from sklearn.metrics import roc_auc_score
```

```
In [87]: train_auc = []
cv_auc = []
K = [3, 15, 25, 51, 101]
for i in tqdm(K):
 neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
 neigh.fit(X_tr, y_train)

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

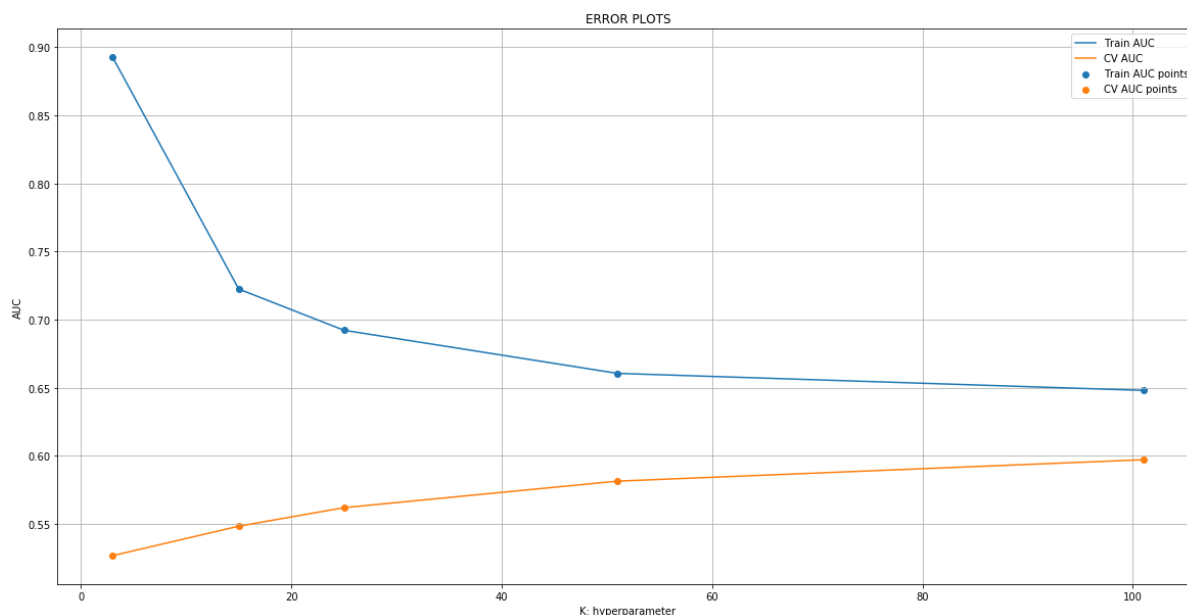
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
 # not the predicted outputs
 train_auc.append(roc_auc_score(y_train, y_train_pred))
 cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

0%|
| 0/5 [00:00<?, ?it/s]
20%|
| 1/5 [01:36<06:24, 96.09s/it]
40%|
| 2/5 [03:18<04:53, 97.99s/it]
60%|
| 3/5 [05:01<03:18, 99.38s/it]
80%|
| 4/5 [06:44<01:40, 100.46s/it]
100%|
| 5/5 [08:27<00:00, 101.43s/it]
```

```
In [88]: plt.figure(figsize=(20,10))
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 [89]: # from the above error plot, we have max AUC for CV with least difference between Train and CV AUCs at K=105
best_k = 105
```

## Testing the performance on test data, plotting ROC Curves

```
In [90]: from sklearn.metrics import roc_curve, auc

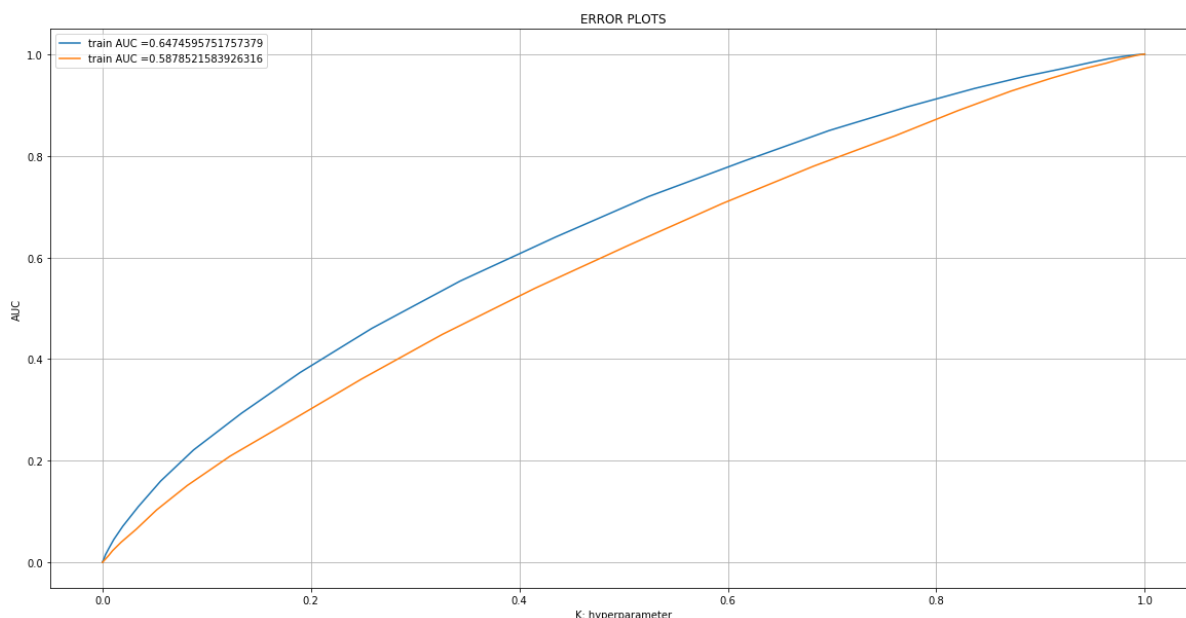
neigh = KNeighborsClassifier(n_neighbors=best_k, n_jobs=-1)
neigh.fit(X_tr, y_train)
roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
not the predicted outputs

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

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



```
In [91]: plt.figure(figsize=(20,10))
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



The Blue plot in the Graph is Train AUC and Orange Plot is Test AUC. Initially the label was accidentally copy pasted as Train AUC for both. Since redrawing plot would require to re run entire ipynb and it takes up a lot of time. so documenting the information as a markdown.

```
In [92]: # we are writing our own function for predict, with defined threshold
we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
 t = threshold[np.argmax(tpr*(1-fpr))]
 # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
 print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
 return t

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

```
In [93]: print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

```
=====
=====
the maximum value of tpr*(1-fpr) 0.36356411719545345 for threshold 0.838
Train confusion matrix
[[3599 1880]
 [13589 16844]]
Test confusion matrix
[[2342 1660]
 [10321 12077]]
```

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

```
In [96]: # preprocessing TFIDF of Text Essays and Project Titles
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10, ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train["essay"].values)
X_train_essay_tfidf = vectorizer.transform(X_train['essay'].values)
X_cv_essay_tfidf = vectorizer.transform(X_cv['essay'].values)
X_test_essay_tfidf = vectorizer.transform(X_test['essay'].values)

print("Shape of Datamatrix after TFIDF Vectorization")
print(X_train_essay_tfidf.shape, y_train.shape)
print(X_cv_essay_tfidf.shape, y_cv.shape)
print(X_test_essay_tfidf.shape, y_test.shape)
print("="*100)
```

```
Shape of Datamatrix after TFIDF Vectorization
(35912, 5000) (35912,)
(17688, 5000) (17688,)
(26400, 5000) (26400,)
```

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

```
In [97]: # Similarly you can vectorize for title also
vectorizer_titles = TfidfVectorizer(min_df=10, ngram_range=(1,4), max_features=5000)
vectorizer_titles.fit(X_train["project_title"])

X_train_pj_title_tfidf = vectorizer.transform(X_train['project_title'].values)
X_cv_pj_title_tfidf = vectorizer.transform(X_cv['project_title'].values)
X_test_pj_title_tfidf = vectorizer.transform(X_test['project_title'].values)

print("Shape of Datamatrix after TFIDF Vectorization")
print(X_train_pj_title_tfidf.shape, y_train.shape)
print(X_cv_pj_title_tfidf.shape, y_cv.shape)
print(X_test_pj_title_tfidf.shape, y_test.shape)
print("="*100)
```

Shape of Datamatrix after TFIDF Vectorization

(35912, 5000) (35912,)

(17688, 5000) (17688,)

(26400, 5000) (26400,)

=====

```
In [98]: # Concatinating all the features for Set 2

X_tr = hstack((X_train_essay_tfidf, X_train_state_oh, X_train_teacher_oh,
 X_train_grade_oh, X_train_price_norm, X_train_category_oh,
 X_train_subcategory_oh, X_train_teach_prev_norm,
 X_train_pj_title_tfidf)).tocsr()

X_cr = hstack((X_cv_essay_tfidf, X_cv_state_oh, X_cv_teacher_oh,
 X_cv_grade_oh, X_cv_category_oh, X_cv_subcategory_oh,
 X_cv_price_norm, X_cv_teach_prev_norm, X_cv_pj_title_tfidf)).to
csr()

X_te = hstack((X_test_essay_tfidf, X_test_state_oh, X_test_teacher_oh,
 X_test_grade_oh, X_test_category_oh, X_test_subcategory_oh,
 X_test_price_norm, X_test_teach_prev_norm,
 X_test_pj_title_tfidf)).tocsr()
```

```
In [99]: train_auc = []
cv_auc = []
K = [3, 15, 25, 51, 101]
for i in tqdm(K):
 neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
 neigh.fit(X_tr, y_train)

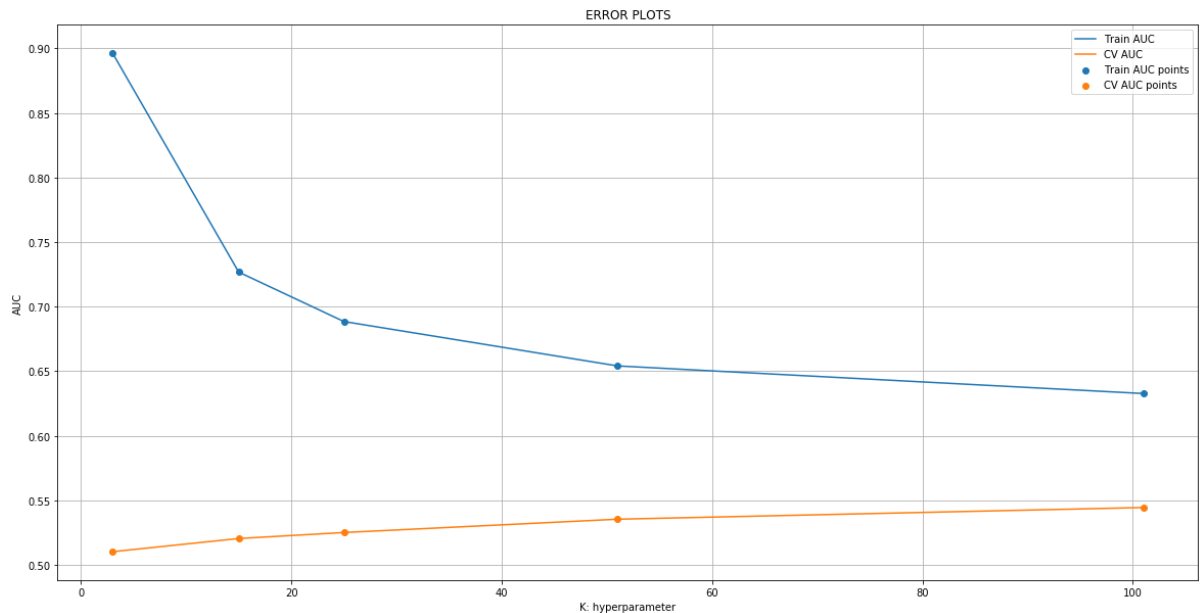
 y_train_pred = batch_predict(neigh, X_tr)
 y_cv_pred = batch_predict(neigh, X_cr)

 # 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))
```

```
In [100]: plt.figure(figsize=(20,10))
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()
```

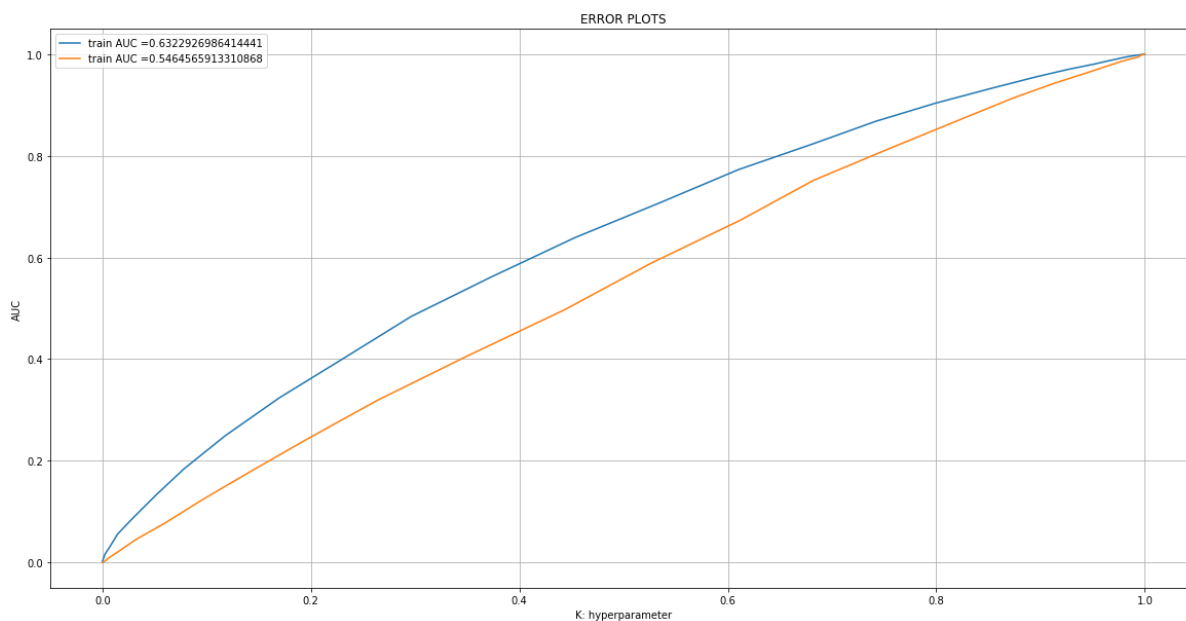


```
In [101]: best_k = 105
Testing the performance of the model on test data, plotting ROC curves
neigh = KNeighborsClassifier(n_neighbors=best_k, n_jobs=-1)
neigh.fit(X_tr, y_train)

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

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

```
In [102]: plt.figure(figsize=(20,10))
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



The Blue plot in the Graph is Train AUC and Orange Plot is Test AUC. Initially the label was accidentally copy pasted as Train AUC for both. Since redrawing plot would require to re run entire ipynb and it takes up a lot of time. so documenting the information as a markdown.

```
In [103]: # Confusion matrix
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))

=====
=====
the maximum value of tpr*(1-fpr) 0.35208438007015214 for threshold 0.857
Train confusion matrix
[[3437 2042]
 [13352 17081]]
Test confusion matrix
[[1898 2104]
 [9227 13171]]
```

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

```
In [104]: # Please write all the code with proper documentation
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 [105]: # average Word2Vec
compute average word2vec for each review.
avg_w2v_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['essay'].values): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_vectors_train.append(vector)

print(len(avg_w2v_vectors_train))
print(len(avg_w2v_vectors_train[0]))
print(avg_w2v_vectors_train[0])
```



```

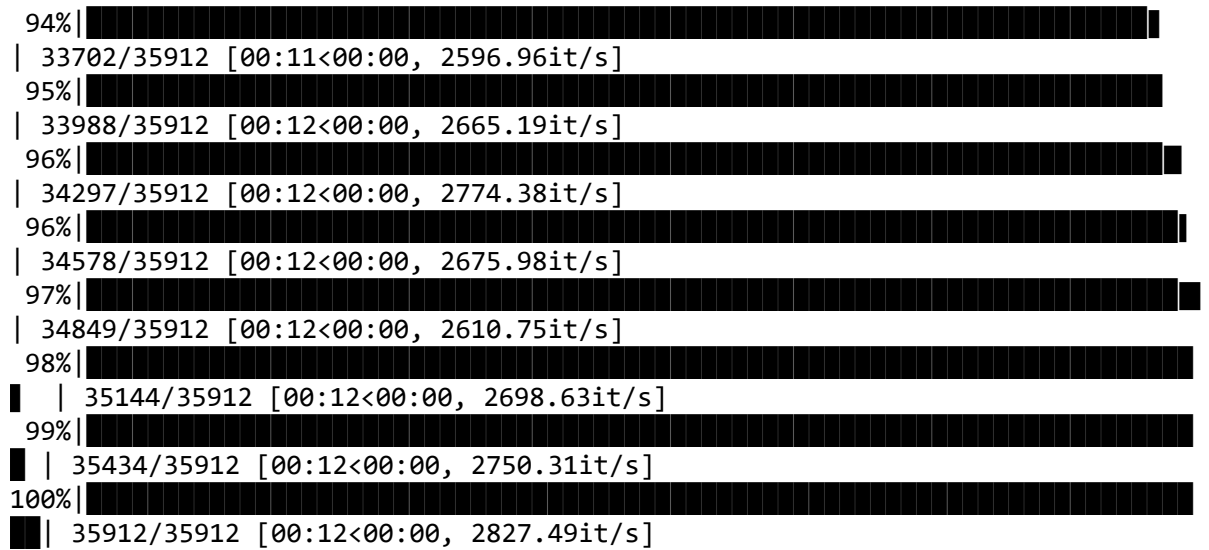
0%|
| 0/35912 [00:00<?, ?it/s]
1%|█
| 304/35912 [00:00<00:11, 3017.95it/s]
2%|██
| 588/35912 [00:00<00:11, 2955.51it/s]
2%|████
| 876/35912 [00:00<00:11, 2925.92it/s]
3%|█████
| 1185/35912 [00:00<00:11, 2967.03it/s]
4%|██████
| 1478/35912 [00:00<00:11, 2949.30it/s]
5%|███████
| 1786/35912 [00:00<00:11, 2980.99it/s]
6%|████████
| 2069/35912 [00:00<00:11, 2927.37it/s]
7%|█████████
| 2380/35912 [00:00<00:11, 2973.47it/s]
7%|██████████
| 2661/35912 [00:00<00:11, 2845.31it/s]
8%|███████████
| 2975/35912 [00:01<00:11, 2921.77it/s]
9%|████████████
| 3289/35912 [00:01<00:10, 2977.64it/s]
10%|█████████████
| 3589/35912 [00:01<00:10, 2977.96it/s]
11%|██████████████
| 3884/35912 [00:01<00:10, 2954.11it/s]
12%|███████████████
| 4178/35912 [00:01<00:11, 2784.89it/s]
12%|████████████████
| 4457/35912 [00:01<00:11, 2763.87it/s]
13%|█████████████████
| 4765/35912 [00:01<00:10, 2845.89it/s]
14%|██████████████████
| 5085/35912 [00:01<00:10, 2937.69it/s]
15%|███████████████████
| 5410/35912 [00:01<00:10, 3018.73it/s]
16%|████████████████████
| 5714/35912 [00:01<00:10, 2825.43it/s]
17%|█████████████████████
| 6007/35912 [00:02<00:10, 2849.92it/s]
18%|██████████████████████
| 6303/35912 [00:02<00:10, 2875.95it/s]
18%|███████████████████████
| 6593/35912 [00:02<00:10, 2762.07it/s]
19%|████████████████████████
| 6893/35912 [00:02<00:10, 2823.56it/s]
20%|█████████████████████████
| 7196/35912 [00:02<00:09, 2876.47it/s]
21%|██████████████████████████
| 7501/35912 [00:02<00:09, 2920.28it/s]
22%|███████████████████████████
| 7800/35912 [00:02<00:09, 2934.53it/s]
23%|████████████████████████████
| 8095/35912 [00:02<00:09, 2932.74it/s]
23%|█████████████████████████████

```



[illegible]





35912

300

```
[-3.75512795e-02 -6.69488513e-02 -6.70803363e-02 -1.90723491e-01
 4.33536981e-02 -4.46185099e-02 -3.70475594e+00 1.76297780e-01
 6.31845660e-03 -2.25834307e-01 1.01603749e-01 6.71636127e-03
 3.48734495e-02 -8.40747873e-02 -1.04109690e-01 -6.75093146e-02
 -9.43774835e-02 -7.41148665e-02 1.05182125e-01 9.59291340e-02
 2.52670566e-02 -4.15442519e-02 -6.39903320e-02 9.43120042e-02
 -5.90282906e-02 -7.06461320e-02 6.07404236e-02 -9.37109480e-02
 -7.13501285e-02 -9.25726127e-02 -2.86350256e-01 -4.88240014e-02
 6.97527967e-02 3.46846363e-02 -7.01084848e-02 -7.31722137e-02
 -5.38548307e-02 -1.13977837e-01 1.75220575e-02 -7.33068595e-02
 -6.97333679e-02 1.49121928e-01 -5.89289241e-03 -2.05862483e-01
 -5.39752434e-02 -8.60852127e-02 7.63792505e-02 -1.01604432e-01
 -5.63580118e-02 -1.25356231e-02 4.37574858e-03 4.46424500e-02
 -3.00707406e-02 -4.03732665e-02 6.80637217e-02 -9.11545091e-02
 8.31546783e-02 -2.76295835e-02 -3.92472906e-02 8.98411288e-02
 -3.75046005e-02 -4.32509764e-02 7.75953461e-02 -5.26954420e-02
 3.56176792e-03 1.19512247e-01 5.98629509e-02 8.99909561e-02
 2.07070452e-01 -1.77531633e-01 -1.46071150e-01 3.66391651e-02
 -2.79509458e-02 -7.51180635e-02 -3.10207371e-02 -2.23594016e-01
 1.25378812e-01 3.30216623e-02 7.55462792e-02 -1.01349446e-01
 3.69506439e-02 -4.14687118e-01 -6.68423505e-02 -1.28211963e-01
 -2.76120401e-02 5.80354334e-02 8.18738108e-02 -9.31786991e-02
 9.92280844e-02 -5.41338708e-02 5.15280594e-02 -5.00081396e-02
 -3.32537236e-02 5.64134986e-02 2.91459538e-02 -2.82199116e-01
 -2.51290434e+00 -6.20086571e-02 1.23330066e-01 1.13975784e-01
 -1.09734866e-01 5.42183345e-02 2.74902544e-01 -2.18736425e-02
 -9.30018113e-03 -6.85015759e-03 7.03599340e-02 -1.39569992e-01
 1.52875708e-03 -2.32907156e-02 -6.03078500e-02 6.99941005e-03
 5.77007557e-02 1.85541224e-01 -7.88257231e-02 7.99913747e-02
 -1.00275681e-01 -8.50787684e-03 1.17561448e-01 5.42727670e-02
 -8.98184340e-02 2.21527170e-02 4.74494210e-02 -1.85058015e-01
 3.79218074e-02 -1.81355519e-02 4.37706354e-02 -5.77079245e-02
 7.01073726e-03 1.13459167e-01 -1.51163398e-02 5.51013995e-02
 -8.17969623e-02 -1.15731642e-01 8.39074829e-02 -3.59646696e-02
 1.03642450e-01 -3.99207349e-02 2.39044960e-01 3.59695305e-01
 1.31053569e-01 4.37903283e-02 8.22113660e-02 4.06554434e-03
 -6.46671292e-03 -3.38558616e-02 9.73816792e-02 -7.05552262e-02
 3.21097126e-01 1.03204479e-01 -7.17701604e-03 -5.84411000e-02
 3.94897877e-02 -9.15246811e-02 1.07563754e-01 -7.09258915e-02
 2.14623821e-02 4.15949281e-02 -1.46675448e-01 -2.90320335e-02
 8.26662104e-02 -8.53794514e-02 -3.91979717e-02 -7.01792995e-02
 -2.39609097e-02 3.54544764e-02 -7.54533198e-02 1.27745723e-01
 1.24336272e-01 -1.03204526e-01 -3.97040623e-02 -2.38604896e-02
 -8.77995146e-02 -1.40369784e-01 -1.14014311e-01 1.91901149e-01
 -1.04500453e-01 -2.22308863e-02 -8.49270698e-02 -7.62908456e-02
 1.19993331e-01 1.57618745e-01 -4.11730561e-02 -6.75506980e-02
 1.61752132e-02 -1.13839944e-01 -2.31425892e-02 2.74699910e-02
 1.13403331e-01 3.36503462e-03 -1.08718588e-02 -7.37779118e-02
 -9.34916019e-03 1.01493253e-02 -1.11434184e-02 -1.27577170e-01
 6.30833887e-02 6.94061448e-02 9.47414604e-02 -1.78464753e-02
 1.54618331e-01 -6.59983879e-02 -5.03701236e-02 5.30649009e-02
 -4.28139481e-02 8.20074540e-02 3.36096502e-02 -1.35731005e-01
 2.29840204e-01 -3.82775307e-02 8.40885778e-02 -3.06141199e-02
 -1.07436777e-01 -1.43445985e-01 -6.16349788e-02 -9.09813509e-02
 -9.85752830e-02 -7.56855292e-02 -1.74586554e-02 6.06722302e-02
```

|                 |                 |                 |                 |
|-----------------|-----------------|-----------------|-----------------|
| -9.78945961e-02 | -1.13928922e-01 | -4.61330797e-02 | -1.80887821e-02 |
| -2.60652231e+00 | 9.17834104e-02  | -6.90932764e-02 | 2.74259265e-02  |
| 1.33784651e-02  | -1.22021278e-01 | 7.16237646e-02  | 7.87142259e-02  |
| -1.79656509e-02 | -2.52493118e-02 | -1.09611346e-01 | 1.86658514e-01  |
| -4.44758121e-02 | -3.36083245e-02 | -1.01186310e-01 | 3.02425453e-02  |
| -1.13465377e-01 | 3.82199332e-02  | -2.19400804e-01 | 1.19033893e-01  |
| -1.68650943e-02 | -3.45514575e-02 | -3.26362642e-02 | -4.12791164e-02 |
| -5.13733425e-02 | 5.39199042e-02  | 1.51183019e-04  | -5.22306712e-02 |
| 7.28870057e-02  | -5.81570783e-02 | 1.01209083e-01  | 7.96570038e-03  |
| 7.70483802e-02  | -5.84984816e-02 | 4.63630439e-02  | -1.08076597e-02 |
| 1.18586374e-01  | -6.39099319e-02 | -3.24825458e-02 | -5.67872764e-02 |
| 8.10332840e-02  | -8.76950491e-02 | -2.64661931e-01 | -3.06185264e-02 |
| 1.43698059e-01  | 2.25140673e-02  | -3.73186382e-02 | -8.39512039e-02 |
| -2.13309127e-01 | 5.59564858e-02  | 5.24485278e-02  | 6.97427901e-02  |
| 6.31198585e-02  | 1.07468608e-02  | -3.39265542e-02 | -2.89662286e-02 |
| 3.00834358e-01  | -3.05820868e-02 | 2.70612792e-02  | 1.47800315e-01  |
| 1.00594321e-02  | 1.29999731e-01  | 4.63337080e-03  | 6.27311689e-02  |
| -2.71931377e-02 | -3.35124858e-02 | -2.49807858e-03 | -9.85138981e-02 |
| -8.56981675e-02 | -7.69237292e-02 | -1.47365100e-02 | 1.20339675e-02  |
| -4.93174528e-02 | 3.68590708e-02  | 1.01812524e-01  | 1.16092321e-02] |

```
In [106]: avg_w2v_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['essay'].values): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_vectors_cv.append(vector)
```



```

0%|
| 0/17688 [00:00<?, ?it/s]
2%|█
| 301/17688 [00:00<00:05, 2988.19it/s]
3%|██
| 601/17688 [00:00<00:05, 2985.19it/s]
5%|████
| 907/17688 [00:00<00:05, 3000.78it/s]
7%|█████
| 1212/17688 [00:00<00:05, 3008.86it/s]
9%|██████
| 1515/17688 [00:00<00:05, 3008.61it/s]
10%|███████
| 1813/17688 [00:00<00:05, 2993.36it/s]
12%|████████
| 2120/17688 [00:00<00:05, 3009.47it/s]
14%|█████████
| 2413/17688 [00:00<00:05, 2978.52it/s]
15%|██████████
| 2717/17688 [00:00<00:05, 2990.25it/s]
17%|███████████
| 3020/17688 [00:01<00:04, 2995.56it/s]
19%|████████████
| 3331/17688 [00:01<00:04, 3022.55it/s]
21%|█████████████
| 3636/17688 [00:01<00:04, 3024.14it/s]
22%|██████████████
| 3942/17688 [00:01<00:04, 3028.23it/s]
24%|███████████████
| 4255/17688 [00:01<00:04, 3051.53it/s]
26%|████████████████
| 4565/17688 [00:01<00:04, 3059.28it/s]
28%|█████████████████
| 4870/17688 [00:01<00:04, 3049.79it/s]
29%|██████████████████
| 5175/17688 [00:01<00:04, 3007.27it/s]
31%|███████████████████
| 5476/17688 [00:01<00:04, 2914.56it/s]
33%|████████████████████
| 5769/17688 [00:01<00:04, 2912.81it/s]
34%|█████████████████████
| 6077/17688 [00:02<00:03, 2954.80it/s]
36%|██████████████████████
| 6387/17688 [00:02<00:03, 2990.58it/s]
38%|███████████████████████
| 6699/17688 [00:02<00:03, 3021.82it/s]
40%|████████████████████████
| 7018/17688 [00:02<00:03, 3063.93it/s]
41%|█████████████████████████
| 7332/17688 [00:02<00:03, 3079.73it/s]
43%|██████████████████████████
| 7641/17688 [00:02<00:03, 2903.14it/s]
45%|███████████████████████████
| 7947/17688 [00:02<00:03, 2942.26it/s]
47%|████████████████████████████
| 8259/17688 [00:02<00:03, 2987.14it/s]
48%|█████████████████████████████

```

```
| 17317/17688 [00:05<00:00, 3075.67it/s]
```



```
In [107]: avg_w2v_vectors_test = []; # the avg-w2v for each sentence/review is stored in
 this list
 for sentence in tqdm(X_test['essay'].values): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_vectors_test.append(vector)
```

```

0%|
| 0/26400 [00:00<?, ?it/s]
1%|█
| 307/26400 [00:00<00:08, 3047.72it/s]
2%|██
| 581/26400 [00:00<00:08, 2941.45it/s]
3%|███
| 880/26400 [00:00<00:08, 2949.46it/s]
4%|████
| 1176/26400 [00:00<00:08, 2946.17it/s]
6%|█████
| 1476/26400 [00:00<00:08, 2955.71it/s]
7%|██████
| 1780/26400 [00:00<00:08, 2974.11it/s]
8%|███████
| 2071/26400 [00:00<00:08, 2948.03it/s]
9%|████████
| 2370/26400 [00:00<00:08, 2954.08it/s]
10%|█████████
| 2680/26400 [00:00<00:07, 2990.06it/s]
11%|██████████
| 2977/26400 [00:01<00:07, 2977.46it/s]
12%|███████████
| 3285/26400 [00:01<00:07, 3001.08it/s]
14%|███████████
| 3590/26400 [00:01<00:07, 3009.07it/s]
15%|████████████
| 3893/26400 [00:01<00:07, 3008.75it/s]
16%|█████████████
| 4192/26400 [00:01<00:07, 2987.55it/s]
17%|█████████████
| 4495/26400 [00:01<00:07, 2993.66it/s]
18%|█████████████
| 4806/26400 [00:01<00:07, 3021.19it/s]
19%|█████████████
| 5114/26400 [00:01<00:07, 3031.87it/s]
21%|█████████████
| 5435/26400 [00:01<00:06, 3076.89it/s]
22%|█████████████
| 5756/26400 [00:01<00:06, 3108.89it/s]
23%|█████████████
| 6072/26400 [00:02<00:06, 3117.32it/s]
24%|█████████████
| 6384/26400 [00:02<00:06, 3056.70it/s]
25%|█████████████
| 6693/26400 [00:02<00:06, 3059.96it/s]
27%|█████████████
| 7015/26400 [00:02<00:06, 3099.73it/s]
28%|█████████████
| 7326/26400 [00:02<00:06, 2997.81it/s]
29%|█████████████
| 7640/26400 [00:02<00:06, 3032.66it/s]
30%|█████████████
| 7953/26400 [00:02<00:06, 3054.68it/s]
31%|█████████████
| 8260/26400 [00:02<00:05, 3052.33it/s]
32%|█████████████

```

```
8571/26400 [00:02<00:05, 3063.04it/s]
34%|███████████|
| 8888/26400 [00:02<00:05, 3087.75it/s]
35%|███████████|
| 9202/26400 [00:03<00:05, 3096.53it/s]
36%|███████████|
| 9512/26400 [00:03<00:05, 3027.48it/s]
37%|███████████|
| 9829/26400 [00:03<00:05, 3062.33it/s]
38%|███████████|
| 10144/26400 [00:03<00:05, 3081.65it/s]
40%|███████████|
| 10461/26400 [00:03<00:05, 3100.97it/s]
41%|███████████|
| 10772/26400 [00:03<00:05, 3096.78it/s]
42%|███████████|
| 11087/26400 [00:03<00:04, 3105.95it/s]
43%|███████████|
| 11399/26400 [00:03<00:04, 3103.37it/s]
44%|███████████|
| 11711/26400 [00:03<00:04, 3101.57it/s]
46%|███████████|
| 12031/26400 [00:03<00:04, 3123.76it/s]
47%|███████████|
| 12344/26400 [00:04<00:04, 3082.05it/s]
48%|███████████|
| 12653/26400 [00:04<00:04, 3077.51it/s]
49%|███████████|
| 12961/26400 [00:04<00:04, 3044.46it/s]
50%|███████████|
| 13273/26400 [00:04<00:04, 3060.15it/s]
51%|███████████|
| 13583/26400 [00:04<00:04, 3065.32it/s]
53%|███████████|
| 13907/26400 [00:04<00:04, 3109.16it/s]
54%|███████████|
| 14221/26400 [00:04<00:03, 3111.59it/s]
55%|███████████|
| 14539/26400 [00:04<00:03, 3125.05it/s]
56%|███████████|
| 14859/26400 [00:04<00:03, 3140.40it/s]
57%|███████████|
| 15176/26400 [00:04<00:03, 3142.38it/s]
59%|███████████|
| 15491/26400 [00:05<00:03, 3100.83it/s]
60%|███████████|
| 15802/26400 [00:05<00:03, 3096.80it/s]
61%|███████████|
| 16112/26400 [00:05<00:03, 3072.65it/s]
62%|███████████|
| 16420/26400 [00:05<00:03, 3067.95it/s]
63%|███████████|
| 16727/26400 [00:05<00:03, 3052.93it/s]
65%|███████████|
| 17033/26400 [00:05<00:03, 3012.47it/s]
66%|███████████|
| 17338/26400 [00:05<00:03, 3017.07it/s]
```

| Percentage | Value       | Time          | Throughput   |
|------------|-------------|---------------|--------------|
| 67%        | 17655/26400 | [00:05<00:02, | 3054.91it/s] |
| 68%        | 17967/26400 | [00:05<00:02, | 3067.53it/s] |
| 69%        | 18282/26400 | [00:05<00:02, | 3084.88it/s] |
| 70%        | 18604/26400 | [00:06<00:02, | 3117.88it/s] |
| 72%        | 18922/26400 | [00:06<00:02, | 3129.49it/s] |
| 73%        | 19236/26400 | [00:06<00:02, | 3107.29it/s] |
| 74%        | 19549/26400 | [00:06<00:02, | 3107.29it/s] |
| 75%        | 19860/26400 | [00:06<00:02, | 3101.32it/s] |
| 76%        | 20171/26400 | [00:06<00:02, | 3087.94it/s] |
| 78%        | 20486/26400 | [00:06<00:01, | 3099.60it/s] |
| 79%        | 20804/26400 | [00:06<00:01, | 3116.58it/s] |
| 80%        | 21130/26400 | [00:06<00:01, | 3151.57it/s] |
| 81%        | 21450/26400 | [00:06<00:01, | 3159.09it/s] |
| 82%        | 21767/26400 | [00:07<00:01, | 3155.45it/s] |
| 84%        | 22083/26400 | [00:07<00:01, | 3140.55it/s] |
| 85%        | 22398/26400 | [00:07<00:01, | 3136.51it/s] |
| 86%        | 22712/26400 | [00:07<00:01, | 3093.80it/s] |
| 87%        | 23022/26400 | [00:07<00:01, | 3079.71it/s] |
| 88%        | 23331/26400 | [00:07<00:00, | 3076.07it/s] |
| 90%        | 23652/26400 | [00:07<00:00, | 3108.45it/s] |
| 91%        | 23978/26400 | [00:07<00:00, | 3145.74it/s] |
| 92%        | 24296/26400 | [00:07<00:00, | 3149.09it/s] |
| 93%        | 24612/26400 | [00:08<00:00, | 3136.14it/s] |
| 94%        | 24926/26400 | [00:08<00:00, | 3130.44it/s] |
| 96%        | 25249/26400 | [00:08<00:00, | 3152.80it/s] |
| 97%        | 25565/26400 | [00:08<00:00, | 3148.17it/s] |
| 98%        | 25880/26400 | [00:08<00:00, | 3113.95it/s] |
| 100%       | 26400/26400 | [00:08<00:00, | 3074.05it/s] |

```
In [108]: # avg w2v for project_titles
avg_w2v_vectors_pj_title_train = []; # the avg-w2v for each sentence/review is
stored in this list
for sentence in tqdm(X_train['project_title'].values): # for each review/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_vectors_pj_title_train.append(vector)

print(len(avg_w2v_vectors_pj_title_train))
print(len(avg_w2v_vectors_pj_title_train[0]))
print(avg_w2v_vectors_pj_title_train[0])
```



```
0%|
| 0/35912 [00:00<?, ?it/s]
41%|
| 14736/35912 [00:00<00:00, 146291.46it/s]
100%|
| 35912/35912 [00:00<00:00, 148181.44it/s]
```

35912

300

```

[-6.3004e-02 3.1239e-02 -1.8579e-01 4.4635e-02 -1.3728e-01 2.0852e-01
-3.0928e+00 5.1131e-01 3.0234e-02 -5.0480e-01 1.8275e-01 -1.7641e-01
-5.2801e-02 1.8643e-01 -1.5230e-01 -5.1885e-02 1.0524e-01 1.1579e-02
 2.4091e-01 1.0702e-02 6.3483e-03 -5.2369e-03 -9.6554e-02 -6.0413e-02
-9.9636e-02 -4.6387e-02 -4.1514e-01 1.3550e-01 6.4773e-02 -1.5405e-01
-9.3040e-01 4.1329e-01 -6.9841e-02 -1.2174e-02 2.4677e-01 -3.4287e-01
 3.9454e-01 -2.0151e-01 3.9742e-01 -5.3066e-02 -4.7727e-01 -4.1655e-01
 3.5602e-01 2.1111e-01 1.0663e-01 -4.1312e-01 8.2616e-03 -3.8526e-01
-1.6343e-01 6.6813e-01 -1.9065e-01 3.3698e-01 -4.3189e-01 5.8042e-02
 5.5317e-01 3.4748e-01 6.4496e-02 2.2694e-03 2.0072e-01 2.2098e-01
-3.0341e-01 -2.8312e-01 -3.1363e-01 5.3067e-01 -2.5296e-02 -1.5465e-01
 3.6343e-02 -4.6435e-01 2.3240e-01 1.8428e-01 -1.2815e-01 -1.3027e-02
-1.9626e-01 -3.1234e-01 -7.1010e-02 -7.8771e-01 -2.9466e-01 -1.7576e-01
 6.7903e-02 -3.9181e-01 2.6900e-01 -6.8174e-01 2.5190e-01 -2.4059e-01
 2.7325e-01 -1.4917e-01 8.3279e-01 -3.0103e-01 -2.8582e-02 -3.3734e-01
-2.9910e-01 -2.0228e-01 1.4101e-01 2.3135e-03 -2.7270e-01 -5.7234e-01
-3.2172e+00 -4.6713e-01 -2.0993e-01 -4.8806e-01 4.9618e-01 7.8517e-02
 6.6969e-01 1.0868e-01 1.9829e-01 1.5625e-01 2.7715e-01 7.0318e-01
-2.9784e-01 2.5801e-01 -3.4829e-01 -4.3414e-01 9.9575e-02 4.7150e-01
-3.1241e-01 4.5226e-01 7.6259e-01 3.6508e-01 2.5221e-01 -1.2493e-01
 9.8863e-02 6.2592e-03 -2.9063e-01 3.1175e-01 3.8153e-01 4.8195e-01
 1.0216e-01 -1.0257e-01 -2.5863e-02 -3.7701e-01 8.7649e-02 -1.3177e-01
-1.1675e-01 4.5153e-01 1.4577e-01 1.9500e-01 1.2384e-01 9.8835e-02
 5.3697e-01 4.9282e-01 6.9783e-02 7.0670e-02 -1.4149e-01 5.8837e-02
 6.3967e-02 -1.9756e-01 -3.6390e-01 1.0592e-01 1.6831e-01 6.1402e-01
-8.0258e-02 3.4805e-01 -1.1233e-01 -1.4224e-01 8.6569e-03 1.5688e-01
 3.5244e-01 2.2691e-01 -5.0237e-02 -2.4736e-01 -2.8841e-01 -8.3235e-02
-1.5517e-01 -3.5509e-01 -1.4373e-01 2.6160e-01 1.8030e-01 -3.8211e-02
-3.8005e-01 -8.3100e-01 -3.3748e-01 3.7304e-02 -1.9378e-01 -5.4925e-01
-7.3450e-02 1.3146e-01 -1.0388e+00 3.5216e-01 5.0045e-01 -1.7853e-01
 9.6374e-02 -9.7410e-02 6.3969e-01 -8.9913e-02 1.7378e-01 -9.7138e-02
-2.9316e-01 -1.9741e-02 -2.8150e-01 7.4473e-01 -2.2485e-01 1.9390e-01
 2.9590e-02 -4.0950e-02 7.4386e-02 -8.6166e-02 2.5382e-01 1.1339e-01
 1.1157e-01 7.2497e-01 -6.2336e-01 -7.8740e-02 -1.9173e-01 1.7641e-01
 6.7506e-02 -3.6622e-01 -5.8263e-02 1.7013e-01 5.0655e-01 1.4760e-01
 3.5870e-01 -2.2899e-01 1.1184e-01 -2.1720e-02 -2.2123e-01 -3.1221e-01
-6.3867e-02 -2.2743e-01 -4.0240e-02 4.5925e-01 -4.4871e-02 4.3306e-02
-3.1833e-01 1.7021e-01 -3.5768e+00 3.8211e-01 4.8462e-02 -1.0886e-01
-1.6740e-01 1.1626e-01 3.4908e-02 2.2989e-01 2.4380e-01 3.4586e-01
 1.5139e-01 6.6616e-01 -4.0853e-01 3.7251e-01 -6.9017e-01 -2.1990e-01
 4.4070e-01 -1.8349e-01 -3.1073e-01 1.8793e-01 -3.3471e-01 3.3364e-02
 3.1918e-02 -9.5115e-01 -1.7393e-01 -1.2559e-01 1.0029e-02 -1.4839e-02
 9.8882e-02 -1.7765e-01 1.9029e-02 -3.3748e-01 -4.4808e-01 -1.8664e-01
 7.2643e-02 1.9137e-01 9.5844e-01 -3.0384e-03 3.5579e-02 -1.8228e-01
-5.5415e-01 -3.2349e-01 -3.7264e-01 3.3125e-01 5.4911e-01 1.5412e-01
 2.1926e-01 -1.5290e-01 -1.9009e-01 3.4024e-02 3.3934e-01 2.8025e-01
 1.0679e-01 2.9297e-01 -5.1538e-02 -4.7361e-02 -2.1209e-01 2.5042e-01
-9.6005e-02 1.8085e-01 1.4923e-01 4.3730e-01 -3.0703e-01 3.2686e-01
 9.3145e-02 2.9420e-01 -2.5635e-01 -3.2371e-01 -2.8559e-01 -5.0285e-01
-6.1455e-02 -3.5492e-01 -2.7389e-01 -1.5301e-01 5.3436e-01 -1.3498e-01]

```

```
In [109]: avg_w2v_vectors_pj_title_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['project_title'].values): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_vectors_pj_title_cv.append(vector)
```

```
0%|
| 0/17688 [00:00<?, ?it/s]
100%|
| 17688/17688 [00:00<00:00, 146571.93it/s]
```

```
In [110]: avg_w2v_vectors_pj_title_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['project_title'].values): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 cnt_words = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if word in glove_words:
 vector += model[word]
 cnt_words += 1
 if cnt_words != 0:
 vector /= cnt_words
 avg_w2v_vectors_pj_title_test.append(vector)
```

```
0%|
| 0/26400 [00:00<?, ?it/s]
100%|
| 26400/26400 [00:00<00:00, 153010.79it/s]
```

```

In [111]: X_tr = hstack((avg_w2v_vectors_train, X_train_state_ohe, X_train_teacher_ohe,
 X_train_grade_ohe, X_train_category_ohe,
 X_train_subcategory_ohe, X_train_price_norm,
 X_train_teach_prev_norm, avg_w2v_vectors_pj_title_train)).tocsr()
()

X_cr = hstack((avg_w2v_vectors_cv, X_cv_state_ohe, X_cv_teacher_ohe,
 X_cv_grade_ohe, X_cv_category_ohe,
 X_cv_subcategory_ohe, X_cv_price_norm,
 X_cv_teach_prev_norm, avg_w2v_vectors_pj_title_cv)).tocsr()

X_te = hstack((avg_w2v_vectors_test, X_test_state_ohe, X_test_teacher_ohe,
 X_test_grade_ohe, X_test_category_ohe,
 X_test_subcategory_ohe, X_test_price_norm,
 X_test_teach_prev_norm, avg_w2v_vectors_pj_title_test)).tocsr()

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

```

Final Data matrix

(35912, 701) (35912,)

(17688, 701) (17688,)

(26400, 701) (26400,)

=====

=====

```
In [112]: train_auc = []
cv_auc = []
K = [3, 15, 25, 51, 101]
for i in tqdm(K):
 neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
 neigh.fit(X_tr, y_train)

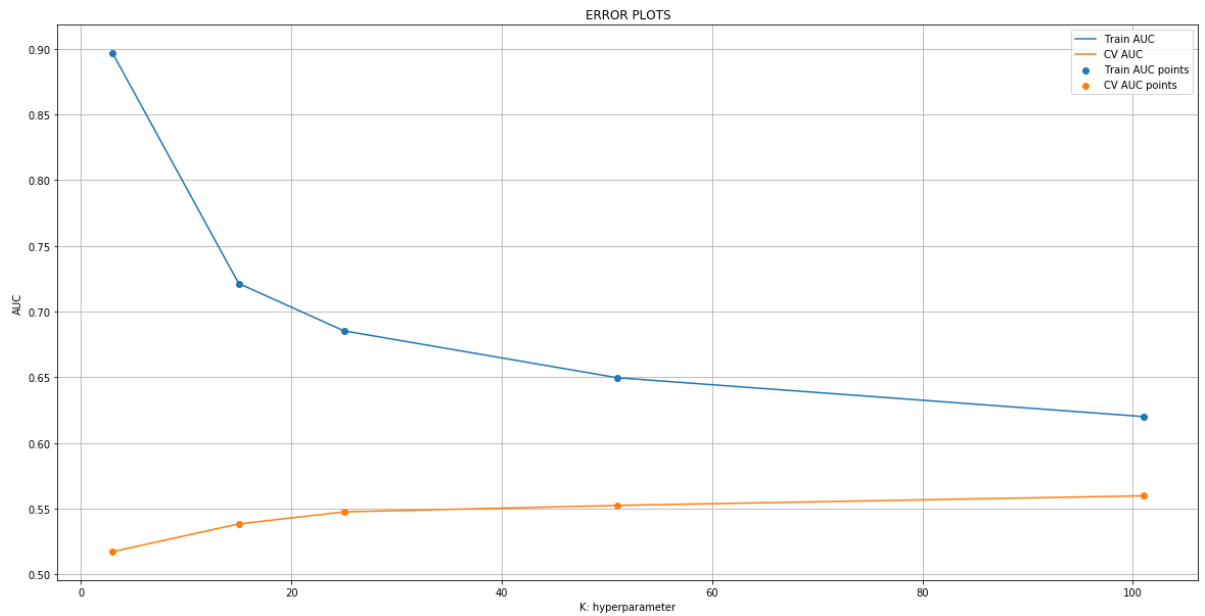
 y_train_pred = batch_predict(neigh, X_tr)
 y_cv_pred = batch_predict(neigh, X_cr)

 # 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))
```

```
In [113]: plt.figure(figsize=(20,10))
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()
```



```

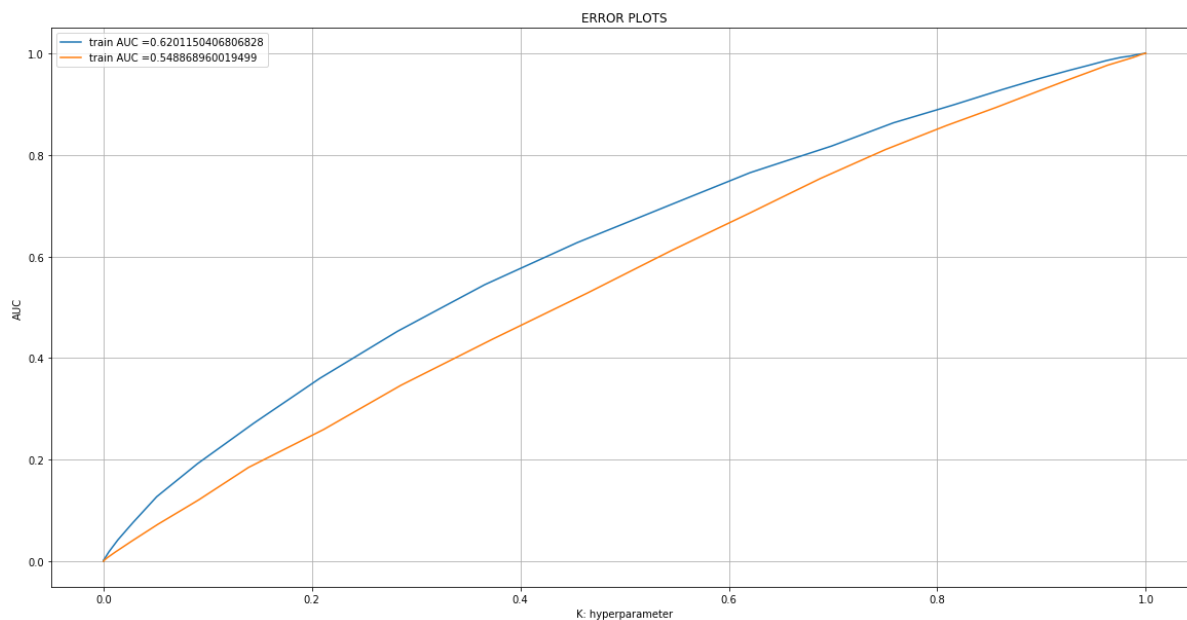
In [114]: best_k = 101
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k, n_jobs=-1)
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.figure(figsize=(20,10))
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



The Blue plot in the Graph is Train AUC and Orange Plot is Test AUC. Initially the label was accidentally copy pasted as Train AUC for both. Since redrawing plot would require to re run entire ipynb and it takes up a lot of time. so documenting the information as a markdown.

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

```
In [115]: # Please write all the code with proper documentation
preprocessing project_title and essay with TFIDF W2V Vectorization
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['essay'])
we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_
)))
tfidf_words = set(tfidf_model.get_feature_names())
```



```

In [116]: # average Word2Vec
compute average word2vec for each review.
tfidf_w2v_vectors_train = []; # the avg-w2v for each sentence/review is stored
in this list
for sentence in tqdm(X_train['essay']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if (word in glove_words) and (word in tfidf_words):
 vec = model[word] # getting the vector for each word
 # here we are multiplying idf value(dictionary[word]) and the tf value
 # ((sentence.count(word)/len(sentence.split())))
 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split
 ())) # getting the tfidf value for each word
 vector += (vec * tf_idf) # calculating tfidf weighted w2v
 tf_idf_weight += tf_idf
 if tf_idf_weight != 0:
 vector /= tf_idf_weight
 tfidf_w2v_vectors_train.append(vector)

print(len(tfidf_w2v_vectors_train))
print(len(tfidf_w2v_vectors_train[0]))

```

```

0%|
| 0/35912 [00:00<?, ?it/s]
0%|
| 28/35912 [00:00<02:09, 277.97it/s]
0%|
| 56/35912 [00:00<02:09, 276.33it/s]
0%||
| 85/35912 [00:00<02:08, 278.89it/s]
0%||
| 116/35912 [00:00<02:05, 284.60it/s]
0%||
| 145/35912 [00:00<02:06, 283.08it/s]
0%||
| 172/35912 [00:00<02:08, 277.54it/s]
1%|
| 200/35912 [00:00<02:08, 277.67it/s]
1%|
| 228/35912 [00:00<02:08, 276.92it/s]
1%|
| 260/35912 [00:00<02:04, 285.71it/s]
1%|
| 291/35912 [00:01<02:02, 291.17it/s]
1%|
| 320/35912 [00:01<02:03, 288.45it/s]
1%|
| 349/35912 [00:01<02:03, 288.28it/s]
1%|
| 378/35912 [00:01<02:07, 278.24it/s]
1%|
| 407/35912 [00:01<02:06, 280.25it/s]
1%|
| 436/35912 [00:01<02:05, 281.69it/s]
1%|
| 465/35912 [00:01<02:10, 272.37it/s]
1%|
| 495/35912 [00:01<02:07, 277.22it/s]
1%|
| 523/35912 [00:01<02:19, 253.39it/s]
2%|
| 553/35912 [00:01<02:13, 264.57it/s]
2%|
| 588/35912 [00:02<02:03, 284.97it/s]
2%|
| 618/35912 [00:02<02:16, 258.22it/s]
2%|
| 646/35912 [00:02<02:14, 263.10it/s]
2%|
| 675/35912 [00:02<02:11, 268.58it/s]
2%|
| 703/35912 [00:02<02:10, 270.55it/s]
2%|
| 734/35912 [00:02<02:06, 278.46it/s]
2%|
| 767/35912 [00:02<02:00, 290.82it/s]
2%|
| 799/35912 [00:02<01:59, 294.28it/s]
2%|

```

| 830/35912 [00:02<01:59, 294.80it/s]  
2%|██████████  
| 861/35912 [00:03<01:58, 296.86it/s]  
2%|██████████  
| 891/35912 [00:03<02:00, 291.11it/s]  
3%|██████████  
| 922/35912 [00:03<02:00, 289.30it/s]  
3%|██████████  
| 952/35912 [00:03<01:59, 291.81it/s]  
3%|██████████  
| 982/35912 [00:03<01:58, 293.59it/s]  
3%|██████████  
| 1012/35912 [00:03<02:01, 286.42it/s]  
3%|██████████  
| 1043/35912 [00:03<01:59, 291.68it/s]  
3%|██████████  
| 1075/35912 [00:03<01:56, 298.17it/s]  
3%|██████████  
| 1107/35912 [00:03<01:55, 302.07it/s]  
3%|██████████  
| 1140/35912 [00:03<01:52, 308.44it/s]  
3%|██████████  
| 1171/35912 [00:04<01:54, 304.61it/s]  
3%|██████████  
| 1202/35912 [00:04<01:58, 293.43it/s]  
3%|██████████  
| 1232/35912 [00:04<02:05, 275.31it/s]  
4%|██████████  
| 1260/35912 [00:04<02:07, 272.09it/s]  
4%|██████████  
| 1293/35912 [00:04<02:01, 285.18it/s]  
4%|██████████  
| 1322/35912 [00:04<02:04, 277.00it/s]  
4%|██████████  
| 1352/35912 [00:04<02:02, 282.14it/s]  
4%|██████████  
| 1383/35912 [00:04<01:59, 288.56it/s]  
4%|██████████  
| 1413/35912 [00:04<02:03, 278.33it/s]  
4%|██████████  
| 1442/35912 [00:05<02:04, 276.33it/s]  
4%|██████████  
| 1472/35912 [00:05<02:01, 282.43it/s]  
4%|██████████  
| 1505/35912 [00:05<01:57, 292.29it/s]  
4%|██████████  
| 1540/35912 [00:05<01:52, 306.11it/s]  
4%|██████████  
| 1571/35912 [00:05<01:53, 303.88it/s]  
4%|██████████  
| 1603/35912 [00:05<01:51, 307.02it/s]  
5%|██████████  
| 1634/35912 [00:05<01:54, 300.13it/s]  
5%|██████████  
| 1665/35912 [00:05<01:54, 298.03it/s]  
5%|██████████  
| 1695/35912 [00:05<01:59, 286.06it/s]

5%|██████|  
| 1726/35912 [00:06<01:56, 292.24it/s]  
5%|██████|  
| 1756/35912 [00:06<02:10, 261.00it/s]  
5%|██████|  
| 1787/35912 [00:06<02:04, 273.46it/s]  
5%|██████|  
| 1816/35912 [00:06<02:09, 264.03it/s]  
5%|██████|  
| 1847/35912 [00:06<02:03, 275.04it/s]  
5%|██████|  
| 1875/35912 [00:06<02:04, 274.31it/s]  
5%|██████|  
| 1907/35912 [00:06<01:58, 286.02it/s]  
5%|██████|  
| 1936/35912 [00:06<01:58, 286.58it/s]  
5%|██████|  
| 1972/35912 [00:06<01:51, 303.92it/s]  
6%|██████|  
| 2003/35912 [00:07<01:56, 291.34it/s]  
6%|██████|  
| 2035/35912 [00:07<01:53, 298.77it/s]  
6%|██████|  
| 2066/35912 [00:07<01:53, 298.80it/s]  
6%|██████|  
| 2097/35912 [00:07<01:54, 294.57it/s]  
6%|██████|  
| 2130/35912 [00:07<01:51, 303.76it/s]  
6%|██████|  
| 2161/35912 [00:07<01:59, 282.50it/s]  
6%|██████|  
| 2197/35912 [00:07<01:52, 300.70it/s]  
6%|██████|  
| 2228/35912 [00:07<01:51, 302.78it/s]  
6%|██████|  
| 2260/35912 [00:07<01:49, 306.22it/s]  
6%|██████|  
| 2291/35912 [00:07<01:49, 305.77it/s]  
6%|██████|  
| 2322/35912 [00:08<01:53, 295.02it/s]  
7%|██████|  
| 2357/35912 [00:08<01:48, 308.20it/s]  
7%|██████|  
| 2389/35912 [00:08<01:58, 282.26it/s]  
7%|██████|  
| 2419/35912 [00:08<01:56, 286.75it/s]  
7%|██████|  
| 2454/35912 [00:08<01:50, 302.62it/s]  
7%|██████|  
| 2485/35912 [00:08<01:51, 300.61it/s]  
7%|██████|  
| 2516/35912 [00:08<01:55, 290.01it/s]  
7%|██████|  
| 2546/35912 [00:08<01:55, 288.11it/s]  
7%|██████|  
| 2582/35912 [00:08<01:50, 302.80it/s]  
7%|██████|

| 2613/35912 [00:09<01:52, 297.30it/s]  
7%|██████  
| 2644/35912 [00:09<01:54, 289.45it/s]  
7%|██████  
| 2674/35912 [00:09<01:58, 280.48it/s]  
8%|██████  
| 2706/35912 [00:09<01:54, 289.12it/s]  
8%|██████  
| 2736/35912 [00:09<01:59, 277.93it/s]  
8%|██████  
| 2767/35912 [00:09<01:56, 285.47it/s]  
8%|██████  
| 2796/35912 [00:09<01:58, 279.58it/s]  
8%|██████  
| 2825/35912 [00:09<01:59, 277.98it/s]  
8%|██████  
| 2857/35912 [00:09<01:54, 288.81it/s]  
8%|██████  
| 2887/35912 [00:10<01:59, 276.20it/s]  
8%|██████  
| 2916/35912 [00:10<01:58, 279.61it/s]  
8%|██████  
| 2945/35912 [00:10<02:04, 265.80it/s]  
8%|██████  
| 2972/35912 [00:10<02:04, 264.91it/s]  
8%|██████  
| 3002/35912 [00:10<02:00, 273.99it/s]  
8%|██████  
| 3030/35912 [00:10<02:01, 269.62it/s]  
9%|██████  
| 3058/35912 [00:10<02:05, 261.43it/s]  
9%|██████  
| 3092/35912 [00:10<01:57, 279.02it/s]  
9%|██████  
| 3125/35912 [00:10<01:52, 292.01it/s]  
9%|██████  
| 3158/35912 [00:10<01:49, 298.58it/s]  
9%|██████  
| 3189/35912 [00:11<01:53, 287.08it/s]  
9%|██████  
| 3219/35912 [00:11<01:52, 289.39it/s]  
9%|██████  
| 3249/35912 [00:11<01:56, 279.66it/s]  
9%|██████  
| 3278/35912 [00:11<02:00, 271.03it/s]  
9%|██████  
| 3306/35912 [00:11<02:02, 266.85it/s]  
9%|██████  
| 3333/35912 [00:11<02:06, 256.57it/s]  
9%|██████  
| 3359/35912 [00:11<02:08, 252.55it/s]  
9%|██████  
| 3390/35912 [00:11<02:02, 265.54it/s]  
10%|██████  
| 3417/35912 [00:11<02:05, 258.66it/s]  
10%|██████  
| 3445/35912 [00:12<02:03, 263.42it/s]

10%|██████████  
| 3472/35912 [00:12<02:10, 249.42it/s]  
10%|██████████  
| 3499/35912 [00:12<02:07, 254.01it/s]  
10%|██████████  
| 3533/35912 [00:12<02:00, 269.82it/s]  
10%|██████████  
| 3563/35912 [00:12<01:57, 276.10it/s]  
10%|██████████  
| 3593/35912 [00:12<01:54, 281.51it/s]  
10%|██████████  
| 3625/35912 [00:12<01:51, 289.10it/s]  
10%|██████████  
| 3655/35912 [00:12<01:51, 289.14it/s]  
10%|██████████  
| 3685/35912 [00:12<01:53, 283.45it/s]  
10%|██████████  
| 3714/35912 [00:13<01:58, 272.75it/s]  
10%|██████████  
| 3751/35912 [00:13<01:49, 294.87it/s]  
11%|██████████  
| 3782/35912 [00:13<01:52, 284.67it/s]  
11%|██████████  
| 3811/35912 [00:13<01:52, 284.79it/s]  
11%|██████████  
| 3840/35912 [00:13<01:54, 279.94it/s]  
11%|██████████  
| 3869/35912 [00:13<01:55, 277.43it/s]  
11%|██████████  
| 3897/35912 [00:13<02:03, 258.42it/s]  
11%|██████████  
| 3924/35912 [00:13<02:04, 257.52it/s]  
11%|██████████  
| 3954/35912 [00:13<01:59, 268.41it/s]  
11%|██████████  
| 3983/35912 [00:14<01:57, 272.44it/s]  
11%|██████████  
| 4015/35912 [00:14<01:52, 283.09it/s]  
11%|██████████  
| 4044/35912 [00:14<01:52, 283.68it/s]  
11%|██████████  
| 4073/35912 [00:14<01:53, 281.62it/s]  
11%|██████████  
| 4105/35912 [00:14<01:50, 288.40it/s]  
12%|██████████  
| 4134/35912 [00:14<01:55, 275.94it/s]  
12%|██████████  
| 4168/35912 [00:14<01:49, 290.43it/s]  
12%|██████████  
| 4199/35912 [00:14<01:47, 295.42it/s]  
12%|██████████  
| 4229/35912 [00:14<01:48, 291.83it/s]  
12%|██████████  
| 4259/35912 [00:14<01:51, 284.44it/s]  
12%|██████████  
| 4293/35912 [00:15<01:46, 296.19it/s]  
12%|██████████

| 4323/35912 [00:15<01:47, 293.21it/s]  
12%|██████████  
| 4353/35912 [00:15<01:50, 285.36it/s]  
12%|██████████  
| 4388/35912 [00:15<01:44, 301.53it/s]  
12%|██████████  
| 4419/35912 [00:15<01:46, 296.42it/s]  
12%|██████████  
| 4453/35912 [00:15<01:42, 307.66it/s]  
12%|██████████  
| 4485/35912 [00:15<01:47, 292.75it/s]  
13%|██████████  
| 4517/35912 [00:15<01:45, 296.48it/s]  
13%|██████████  
| 4550/35912 [00:15<01:43, 303.50it/s]  
13%|██████████  
| 4581/35912 [00:16<01:46, 295.21it/s]  
13%|██████████  
| 4615/35912 [00:16<01:42, 305.92it/s]  
13%|██████████  
| 4646/35912 [00:16<01:45, 296.81it/s]  
13%|██████████  
| 4676/35912 [00:16<01:46, 294.47it/s]  
13%|██████████  
| 4706/35912 [00:16<01:46, 292.05it/s]  
13%|██████████  
| 4736/35912 [00:16<01:53, 273.70it/s]  
13%|██████████  
| 4770/35912 [00:16<01:47, 288.70it/s]  
13%|██████████  
| 4800/35912 [00:16<01:47, 289.69it/s]  
13%|██████████  
| 4830/35912 [00:16<01:50, 281.43it/s]  
14%|██████████  
| 4860/35912 [00:17<01:48, 285.34it/s]  
14%|██████████  
| 4889/35912 [00:17<01:50, 279.50it/s]  
14%|██████████  
| 4918/35912 [00:17<01:51, 277.93it/s]  
14%|██████████  
| 4951/35912 [00:17<01:46, 290.41it/s]  
14%|██████████  
| 4985/35912 [00:17<01:42, 301.49it/s]  
14%|██████████  
| 5019/35912 [00:17<01:39, 308.93it/s]  
14%|██████████  
| 5051/35912 [00:17<01:45, 293.55it/s]  
14%|██████████  
| 5083/35912 [00:17<01:42, 300.39it/s]  
14%|██████████  
| 5114/35912 [00:17<01:44, 294.82it/s]  
14%|██████████  
| 5145/35912 [00:17<01:43, 296.02it/s]  
14%|██████████  
| 5181/35912 [00:18<01:38, 310.49it/s]  
15%|██████████  
| 5213/35912 [00:18<01:38, 311.70it/s]

15%|██████████|  
| 5246/35912 [00:18<01:37, 314.50it/s]  
15%|██████████|  
| 5278/35912 [00:18<01:37, 313.60it/s]  
15%|██████████|  
| 5310/35912 [00:18<01:39, 306.69it/s]  
15%|██████████|  
| 5341/35912 [00:18<01:45, 289.01it/s]  
15%|██████████|  
| 5371/35912 [00:18<01:45, 289.09it/s]  
15%|██████████|  
| 5401/35912 [00:18<01:51, 273.37it/s]  
15%|██████████|  
| 5432/35912 [00:18<01:48, 282.08it/s]  
15%|██████████|  
| 5461/35912 [00:19<01:51, 274.16it/s]  
15%|██████████|  
| 5489/35912 [00:19<01:50, 275.30it/s]  
15%|██████████|  
| 5524/35912 [00:19<01:43, 292.85it/s]  
15%|██████████|  
| 5556/35912 [00:19<01:42, 296.56it/s]  
16%|██████████|  
| 5589/35912 [00:19<01:39, 305.24it/s]  
16%|██████████|  
| 5620/35912 [00:19<01:42, 296.36it/s]  
16%|██████████|  
| 5650/35912 [00:19<01:44, 288.26it/s]  
16%|██████████|  
| 5683/35912 [00:19<01:41, 299.04it/s]  
16%|██████████|  
| 5714/35912 [00:19<01:40, 300.72it/s]  
16%|██████████|  
| 5745/35912 [00:20<01:41, 296.72it/s]  
16%|██████████|  
| 5778/35912 [00:20<01:38, 304.52it/s]  
16%|██████████|  
| 5809/35912 [00:20<01:40, 299.31it/s]  
16%|██████████|  
| 5840/35912 [00:20<01:40, 300.04it/s]  
16%|██████████|  
| 5871/35912 [00:20<01:39, 302.31it/s]  
16%|██████████|  
| 5903/35912 [00:20<01:37, 306.77it/s]  
17%|██████████|  
| 5934/35912 [00:20<01:38, 304.35it/s]  
17%|██████████|  
| 5965/35912 [00:20<01:40, 298.31it/s]  
17%|██████████|  
| 5996/35912 [00:20<01:39, 301.10it/s]  
17%|██████████|  
| 6027/35912 [00:20<01:40, 296.99it/s]  
17%|██████████|  
| 6057/35912 [00:21<01:48, 275.97it/s]  
17%|██████████|  
| 6085/35912 [00:21<01:50, 269.40it/s]  
17%|██████████|



| 6116/35912 [00:21<01:48, 275.40it/s]  
17%|██████████  
| 6148/35912 [00:21<01:44, 285.32it/s]  
17%|██████████  
| 6178/35912 [00:21<01:43, 286.48it/s]  
17%|██████████  
| 6207/35912 [00:21<01:44, 284.38it/s]  
17%|██████████  
| 6237/35912 [00:21<01:43, 287.46it/s]  
17%|██████████  
| 6268/35912 [00:21<01:41, 292.43it/s]  
18%|██████████  
| 6298/35912 [00:21<01:40, 294.03it/s]  
18%|██████████  
| 6328/35912 [00:22<01:42, 290.01it/s]  
18%|██████████  
| 6358/35912 [00:22<01:42, 288.97it/s]  
18%|██████████  
| 6387/35912 [00:22<01:43, 285.25it/s]  
18%|██████████  
| 6416/35912 [00:22<01:43, 285.19it/s]  
18%|██████████  
| 6448/35912 [00:22<01:40, 292.61it/s]  
18%|██████████  
| 6478/35912 [00:22<01:43, 284.96it/s]  
18%|██████████  
| 6514/35912 [00:22<01:36, 303.41it/s]  
18%|██████████  
| 6548/35912 [00:22<01:34, 309.48it/s]  
18%|██████████  
| 6580/35912 [00:22<01:38, 298.83it/s]  
18%|██████████  
| 6613/35912 [00:22<01:35, 306.07it/s]  
19%|██████████  
| 6644/35912 [00:23<01:36, 302.98it/s]  
19%|██████████  
| 6677/35912 [00:23<01:35, 304.83it/s]  
19%|██████████  
| 6713/35912 [00:23<01:32, 316.39it/s]  
19%|██████████  
| 6745/35912 [00:23<01:36, 302.44it/s]  
19%|██████████  
| 6776/35912 [00:23<01:38, 297.04it/s]  
19%|██████████  
| 6809/35912 [00:23<01:35, 304.75it/s]  
19%|██████████  
| 6844/35912 [00:23<01:32, 315.57it/s]  
19%|██████████  
| 6876/35912 [00:23<01:31, 316.19it/s]  
19%|██████████  
| 6908/35912 [00:23<01:40, 289.24it/s]  
19%|██████████  
| 6944/35912 [00:24<01:35, 302.17it/s]  
19%|██████████  
| 6975/35912 [00:24<01:38, 294.33it/s]  
20%|██████████  
| 7005/35912 [00:24<01:44, 275.86it/s]

20%|████████████████████|  
| 7034/35912 [00:24<01:47, 269.28it/s]  
20%|████████████████████|  
| 7071/35912 [00:24<01:38, 291.34it/s]  
20%|████████████████████|  
| 7101/35912 [00:24<01:38, 291.55it/s]  
20%|████████████████████|  
| 7132/35912 [00:24<01:37, 296.23it/s]  
20%|████████████████████|  
| 7163/35912 [00:24<01:37, 293.65it/s]  
20%|████████████████████|  
| 7193/35912 [00:24<01:39, 287.28it/s]  
20%|████████████████████|  
| 7226/35912 [00:25<01:36, 295.90it/s]  
20%|████████████████████|  
| 7256/35912 [00:25<01:39, 287.96it/s]  
20%|████████████████████|  
| 7288/35912 [00:25<01:38, 289.85it/s]  
20%|████████████████████|  
| 7318/35912 [00:25<01:38, 289.67it/s]  
20%|████████████████████|  
| 7351/35912 [00:25<01:35, 297.65it/s]  
21%|████████████████████|  
| 7381/35912 [00:25<01:36, 295.09it/s]  
21%|████████████████████|  
| 7412/35912 [00:25<01:35, 298.77it/s]  
21%|████████████████████|  
| 7442/35912 [00:25<01:36, 294.98it/s]  
21%|████████████████████|  
| 7472/35912 [00:25<01:38, 289.00it/s]  
21%|████████████████████|  
| 7508/35912 [00:25<01:32, 306.60it/s]  
21%|████████████████████|  
| 7540/35912 [00:26<01:35, 296.94it/s]  
21%|████████████████████|  
| 7570/35912 [00:26<01:37, 289.48it/s]  
21%|████████████████████|  
| 7600/35912 [00:26<01:38, 286.10it/s]  
21%|████████████████████|  
| 7629/35912 [00:26<01:39, 284.12it/s]  
21%|████████████████████|  
| 7658/35912 [00:26<01:39, 282.75it/s]  
21%|████████████████████|  
| 7687/35912 [00:26<01:41, 279.34it/s]  
21%|████████████████████|  
| 7716/35912 [00:26<01:44, 270.83it/s]  
22%|████████████████████|  
| 7750/35912 [00:26<01:37, 287.91it/s]  
22%|████████████████████|  
| 7782/35912 [00:26<01:34, 296.24it/s]  
22%|████████████████████|  
| 7812/35912 [00:27<01:38, 284.10it/s]  
22%|████████████████████|  
| 7841/35912 [00:27<01:41, 275.50it/s]  
22%|████████████████████|  
| 7869/35912 [00:27<01:42, 273.01it/s]  
22%|████████████████████|

| 7900/35912 [00:27<01:39, 281.05it/s]  
22%|████████████████████  
| 7929/35912 [00:27<01:40, 278.19it/s]  
22%|████████████████████  
| 7957/35912 [00:27<01:45, 265.50it/s]  
22%|████████████████████  
| 7992/35912 [00:27<01:37, 285.02it/s]  
22%|████████████████████  
| 8022/35912 [00:27<01:38, 283.04it/s]  
22%|████████████████████  
| 8053/35912 [00:27<01:36, 289.22it/s]  
23%|████████████████████  
| 8089/35912 [00:27<01:30, 305.99it/s]  
23%|████████████████████  
| 8121/35912 [00:28<01:30, 307.63it/s]  
23%|████████████████████  
| 8153/35912 [00:28<01:33, 295.96it/s]  
23%|████████████████████  
| 8185/35912 [00:28<01:31, 302.16it/s]  
23%|████████████████████  
| 8217/35912 [00:28<01:31, 303.18it/s]  
23%|████████████████████  
| 8248/35912 [00:28<01:34, 293.33it/s]  
23%|████████████████████  
| 8278/35912 [00:28<01:34, 292.94it/s]  
23%|████████████████████  
| 8308/35912 [00:28<01:35, 289.29it/s]  
23%|████████████████████  
| 8338/35912 [00:28<01:35, 289.27it/s]  
23%|████████████████████  
| 8368/35912 [00:28<01:40, 274.23it/s]  
23%|████████████████████  
| 8400/35912 [00:29<01:36, 285.96it/s]  
23%|████████████████████  
| 8435/35912 [00:29<01:31, 301.22it/s]  
24%|████████████████████  
| 8469/35912 [00:29<01:29, 307.06it/s]  
24%|████████████████████  
| 8501/35912 [00:29<01:36, 284.59it/s]  
24%|████████████████████  
| 8533/35912 [00:29<01:33, 292.16it/s]  
24%|████████████████████  
| 8563/35912 [00:29<01:34, 289.59it/s]  
24%|████████████████████  
| 8593/35912 [00:29<01:34, 287.82it/s]  
24%|████████████████████  
| 8623/35912 [00:29<01:34, 289.90it/s]  
24%|████████████████████  
| 8658/35912 [00:29<01:29, 303.49it/s]  
24%|████████████████████  
| 8689/35912 [00:30<01:35, 286.23it/s]  
24%|████████████████████  
| 8719/35912 [00:30<01:34, 286.31it/s]  
24%|████████████████████  
| 8749/35912 [00:30<01:34, 287.18it/s]  
24%|████████████████████  
| 8778/35912 [00:30<01:39, 272.81it/s]



10505/35912 [00:36&lt;01:34



```

11397/35912 [00:39<01:28, 277.00it/s]
32%|███████████|
| 11428/35912 [00:39<01:26, 282.44it/s]
32%|███████████|
| 11457/35912 [00:39<01:26, 282.40it/s]
32%|███████████|
| 11486/35912 [00:39<01:30, 270.55it/s]
32%|███████████|
| 11518/35912 [00:40<01:27, 280.19it/s]
32%|███████████|
| 11552/35912 [00:40<01:23, 292.20it/s]
32%|███████████|
| 11584/35912 [00:40<01:21, 299.40it/s]
32%|███████████|
| 11617/35912 [00:40<01:20, 302.29it/s]
32%|███████████|
| 11648/35912 [00:40<01:25, 283.92it/s]
33%|███████████|
| 11684/35912 [00:40<01:20, 301.06it/s]
33%|███████████|
| 11716/35912 [00:40<01:19, 303.26it/s]
33%|███████████|
| 11750/35912 [00:40<01:17, 311.93it/s]
33%|███████████|
| 11782/35912 [00:40<01:23, 289.84it/s]
33%|███████████|
| 11812/35912 [00:41<01:25, 280.73it/s]
33%|███████████|
| 11842/35912 [00:41<01:24, 284.03it/s]
33%|███████████|
| 11871/35912 [00:41<01:25, 282.69it/s]
33%|███████████|
| 11905/35912 [00:41<01:20, 296.40it/s]
33%|███████████|
| 11939/35912 [00:41<01:18, 305.17it/s]
33%|███████████|
| 11971/35912 [00:41<01:17, 308.79it/s]
33%|███████████|
| 12003/35912 [00:41<01:19, 300.07it/s]
34%|███████████|
| 12034/35912 [00:41<01:24, 284.10it/s]
34%|███████████|
| 12064/35912 [00:41<01:23, 287.26it/s]
34%|███████████|
| 12094/35912 [00:41<01:22, 290.35it/s]
34%|███████████|
| 12125/35912 [00:42<01:22, 289.58it/s]
34%|███████████|
| 12155/35912 [00:42<01:24, 281.35it/s]
34%|███████████|
| 12184/35912 [00:42<01:32, 257.68it/s]
34%|███████████|
| 12213/35912 [00:42<01:29, 263.88it/s]
34%|███████████|
| 12242/35912 [00:42<01:27, 269.15it/s]
34%|███████████|
| 12270/35912 [00:42<01:27, 271.73it/s]

```





[illegible]

|                                       |            |
|---------------------------------------|------------|
| 39%                                   | ██████████ |
| 14062/35912 [00:48<01:12, 301.12it/s] |            |
| 39%                                   | ██████████ |
| 14093/35912 [00:48<01:15, 287.92it/s] |            |
| 39%                                   | ██████████ |
| 14124/35912 [00:48<01:14, 291.12it/s] |            |
| 39%                                   | ██████████ |
| 14160/35912 [00:49<01:11, 305.92it/s] |            |
| 40%                                   | ██████████ |
| 14191/35912 [00:49<01:13, 294.28it/s] |            |
| 40%                                   | ██████████ |
| 14224/35912 [00:49<01:11, 302.71it/s] |            |
| 40%                                   | ██████████ |
| 14255/35912 [00:49<01:12, 298.94it/s] |            |
| 40%                                   | ██████████ |
| 14286/35912 [00:49<01:12, 298.92it/s] |            |
| 40%                                   | ██████████ |
| 14317/35912 [00:49<01:13, 292.16it/s] |            |
| 40%                                   | ██████████ |
| 14355/35912 [00:49<01:08, 312.59it/s] |            |
| 40%                                   | ██████████ |
| 14387/35912 [00:49<01:09, 308.66it/s] |            |
| 40%                                   | ██████████ |
| 14419/35912 [00:49<01:15, 284.79it/s] |            |
| 40%                                   | ██████████ |
| 14449/35912 [00:50<01:17, 277.40it/s] |            |
| 40%                                   | ██████████ |
| 14478/35912 [00:50<01:16, 278.85it/s] |            |
| 40%                                   | ██████████ |
| 14510/35912 [00:50<01:14, 287.91it/s] |            |
| 40%                                   | ██████████ |
| 14540/35912 [00:50<01:16, 277.92it/s] |            |
| 41%                                   | ██████████ |
| 14569/35912 [00:50<01:15, 280.84it/s] |            |
| 41%                                   | ██████████ |
| 14599/35912 [00:50<01:15, 284.11it/s] |            |
| 41%                                   | ██████████ |
| 14628/35912 [00:50<01:16, 277.86it/s] |            |
| 41%                                   | ██████████ |
| 14659/35912 [00:50<01:14, 284.64it/s] |            |
| 41%                                   | ██████████ |
| 14688/35912 [00:50<01:20, 264.56it/s] |            |
| 41%                                   | ██████████ |
| 14715/35912 [00:51<01:20, 261.74it/s] |            |
| 41%                                   | ██████████ |
| 14742/35912 [00:51<01:22, 256.11it/s] |            |
| 41%                                   | ██████████ |
| 14775/35912 [00:51<01:17, 271.36it/s] |            |
| 41%                                   | ██████████ |
| 14805/35912 [00:51<01:15, 278.80it/s] |            |
| 41%                                   | ██████████ |
| 14834/35912 [00:51<01:15, 277.44it/s] |            |
| 41%                                   | ██████████ |
| 14868/35912 [00:51<01:11, 292.34it/s] |            |
| 41%                                   | ██████████ |
| 14901/35912 [00:51<01:09, 301.27it/s] |            |
| 42%                                   | ██████████ |

```
14932/35912 [00:51<01:11, 292.90it/s]
42%|███████████|
| 14965/35912 [00:51<01:09, 300.86it/s]
42%|███████████|
| 14998/35912 [00:52<01:08, 307.56it/s]
42%|███████████|
| 15030/35912 [00:52<01:07, 308.73it/s]
42%|███████████|
| 15062/35912 [00:52<01:08, 303.41it/s]
42%|███████████|
| 15096/35912 [00:52<01:06, 312.90it/s]
42%|███████████|
| 15128/35912 [00:52<01:08, 305.35it/s]
42%|███████████|
| 15159/35912 [00:52<01:13, 281.15it/s]
42%|███████████|
| 15193/35912 [00:52<01:10, 295.21it/s]
42%|███████████|
| 15224/35912 [00:52<01:12, 284.11it/s]
42%|███████████|
| 15254/35912 [00:52<01:12, 285.63it/s]
43%|███████████|
| 15286/35912 [00:52<01:10, 294.55it/s]
43%|███████████|
| 15316/35912 [00:53<01:09, 295.52it/s]
43%|███████████|
| 15346/35912 [00:53<01:09, 295.31it/s]
43%|███████████|
| 15376/35912 [00:53<01:09, 294.32it/s]
43%|███████████|
| 15411/35912 [00:53<01:07, 304.48it/s]
43%|███████████|
| 15442/35912 [00:53<01:09, 295.01it/s]
43%|███████████|
| 15472/35912 [00:53<01:11, 284.93it/s]
43%|███████████|
| 15503/35912 [00:53<01:10, 290.59it/s]
43%|███████████|
| 15533/35912 [00:53<01:10, 287.69it/s]
43%|███████████|
| 15563/35912 [00:53<01:10, 290.65it/s]
43%|███████████|
| 15593/35912 [00:54<01:12, 282.06it/s]
44%|███████████|
| 15622/35912 [00:54<01:12, 278.09it/s]
44%|███████████|
| 15654/35912 [00:54<01:10, 286.57it/s]
44%|███████████|
| 15683/35912 [00:54<01:11, 284.44it/s]
44%|███████████|
| 15712/35912 [00:54<01:14, 271.09it/s]
44%|███████████|
| 15741/35912 [00:54<01:13, 275.93it/s]
44%|███████████|
| 15770/35912 [00:54<01:12, 278.60it/s]
44%|███████████|
| 15804/35912 [00:54<01:08, 293.25it/s]
```

```
| 44%| ██████████
| 15837/35912 [00:54<01:06, 302.77it/s]
44%| ██████████
| 15872/35912 [00:54<01:03, 314.93it/s]
44%| ██████████
| 15904/35912 [00:55<01:08, 291.64it/s]
44%| ██████████
| 15934/35912 [00:55<01:08, 290.91it/s]
44%| ██████████
| 15966/35912 [00:55<01:07, 297.63it/s]
45%| ██████████
| 15997/35912 [00:55<01:08, 292.12it/s]
45%| ██████████
| 16028/35912 [00:55<01:07, 293.28it/s]
45%| ██████████
| 16058/35912 [00:55<01:09, 285.40it/s]
45%| ██████████
| 16089/35912 [00:55<01:07, 291.76it/s]
45%| ██████████
| 16123/35912 [00:55<01:05, 303.29it/s]
45%| ██████████
| 16155/35912 [00:55<01:04, 306.60it/s]
45%| ██████████
| 16186/35912 [00:56<01:05, 301.59it/s]
45%| ██████████
| 16217/35912 [00:56<01:08, 286.62it/s]
45%| ██████████
| 16246/35912 [00:56<01:10, 278.75it/s]
45%| ██████████
| 16278/35912 [00:56<01:08, 288.61it/s]
45%| ██████████
| 16308/35912 [00:56<01:11, 276.07it/s]
46%| ██████████
| 16340/35912 [00:56<01:08, 286.59it/s]
46%| ██████████
| 16376/35912 [00:56<01:04, 302.40it/s]
46%| ██████████
| 16407/35912 [00:56<01:07, 287.14it/s]
46%| ██████████
| 16437/35912 [00:56<01:07, 288.59it/s]
46%| ██████████
| 16467/35912 [00:57<01:12, 267.95it/s]
46%| ██████████
| 16495/35912 [00:57<01:12, 269.32it/s]
46%| ██████████
| 16527/35912 [00:57<01:08, 282.21it/s]
46%| ██████████
| 16562/35912 [00:57<01:05, 296.78it/s]
46%| ██████████
| 16593/35912 [00:57<01:05, 294.03it/s]
46%| ██████████
| 16625/35912 [00:57<01:04, 299.90it/s]
46%| ██████████
| 16657/35912 [00:57<01:03, 302.44it/s]
46%| ██████████
| 16688/35912 [00:57<01:06, 289.57it/s]
47%| ██████████
```

[illegible]

[illegible]

95/180

```
| 54%| ██████████
| 19346/35912 [01:07<00:58, 280.87it/s]
54%| ██████████
| 19375/35912 [01:07<00:58, 281.30it/s]
54%| ██████████
| 19404/35912 [01:07<00:58, 280.78it/s]
54%| ██████████
| 19433/35912 [01:07<01:00, 271.02it/s]
54%| ██████████
| 19461/35912 [01:07<01:01, 266.84it/s]
54%| ██████████
| 19493/35912 [01:07<00:58, 278.83it/s]
54%| ██████████
| 19526/35912 [01:07<00:56, 291.87it/s]
54%| ██████████
| 19561/35912 [01:07<00:54, 302.62it/s]
55%| ██████████
| 19594/35912 [01:07<00:52, 309.70it/s]
55%| ██████████
| 19626/35912 [01:08<00:55, 295.66it/s]
55%| ██████████
| 19659/35912 [01:08<00:53, 302.90it/s]
55%| ██████████
| 19691/35912 [01:08<00:52, 306.30it/s]
55%| ██████████
| 19722/35912 [01:08<00:55, 289.59it/s]
55%| ██████████
| 19752/35912 [01:08<00:56, 285.37it/s]
55%| ██████████
| 19782/35912 [01:08<00:56, 286.52it/s]
55%| ██████████
| 19811/35912 [01:08<00:58, 275.51it/s]
55%| ██████████
| 19846/35912 [01:08<00:55, 291.56it/s]
55%| ██████████
| 19876/35912 [01:08<00:54, 291.71it/s]
55%| ██████████
| 19907/35912 [01:09<00:54, 293.82it/s]
56%| ██████████
| 19941/35912 [01:09<00:52, 305.69it/s]
56%| ██████████
| 19972/35912 [01:09<00:52, 302.74it/s]
56%| ██████████
| 20003/35912 [01:09<00:56, 283.42it/s]
56%| ██████████
| 20032/35912 [01:09<01:01, 258.89it/s]
56%| ██████████
| 20066/35912 [01:09<00:56, 278.34it/s]
56%| ██████████
| 20099/35912 [01:09<00:54, 289.17it/s]
56%| ██████████
| 20129/35912 [01:09<00:55, 285.11it/s]
56%| ██████████
| 20160/35912 [01:09<00:54, 291.54it/s]
56%| ██████████
| 20196/35912 [01:09<00:51, 307.81it/s]
56%| ██████████
```



97/180

```

59%|██████████|
| 21142/35912 [01:13<00:49, 297.90it/s]
59%|██████████|
| 21172/35912 [01:13<00:49, 297.87it/s]
59%|██████████|
| 21202/35912 [01:13<00:50, 290.11it/s]
59%|██████████|
| 21232/35912 [01:13<00:51, 286.53it/s]
59%|██████████|
| 21267/35912 [01:13<00:48, 301.66it/s]
59%|██████████|
| 21298/35912 [01:13<00:49, 294.83it/s]
59%|██████████|
| 21328/35912 [01:13<00:50, 286.43it/s]
59%|██████████|
| 21363/35912 [01:13<00:48, 301.58it/s]
60%|██████████|
| 21394/35912 [01:14<00:47, 302.52it/s]
60%|██████████|
| 21425/35912 [01:14<00:50, 288.81it/s]
60%|██████████|
| 21455/35912 [01:14<00:49, 291.46it/s]
60%|██████████|
| 21487/35912 [01:14<00:48, 296.37it/s]
60%|██████████|
| 21517/35912 [01:14<00:52, 273.34it/s]
60%|██████████|
| 21547/35912 [01:14<00:51, 280.25it/s]
60%|██████████|
| 21576/35912 [01:14<00:51, 276.07it/s]
60%|██████████|
| 21606/35912 [01:14<00:50, 281.46it/s]
60%|██████████|
| 21635/35912 [01:14<00:50, 283.36it/s]
60%|██████████|
| 21665/35912 [01:15<00:50, 281.89it/s]
60%|██████████|
| 21694/35912 [01:15<00:50, 282.84it/s]
60%|██████████|
| 21724/35912 [01:15<00:49, 284.73it/s]
61%|██████████|
| 21755/35912 [01:15<00:48, 290.45it/s]
61%|██████████|
| 21787/35912 [01:15<00:47, 298.11it/s]
61%|██████████|
| 21820/35912 [01:15<00:46, 304.70it/s]
61%|██████████|
| 21851/35912 [01:15<00:47, 296.85it/s]
61%|██████████|
| 21885/35912 [01:15<00:46, 303.06it/s]
61%|██████████|
| 21918/35912 [01:15<00:45, 310.03it/s]
61%|██████████|
| 21950/35912 [01:15<00:44, 312.28it/s]
61%|██████████|
| 21983/35912 [01:16<00:44, 309.62it/s]
61%|██████████|

```

99/180

```
64%|██████████|
| 22915/35912 [01:19<00:46, 279.61it/s]
64%|██████████|
| 22945/35912 [01:19<00:45, 284.03it/s]
64%|██████████|
| 22975/35912 [01:19<00:44, 288.03it/s]
64%|██████████|
| 23004/35912 [01:19<00:45, 285.45it/s]
64%|██████████|
| 23033/35912 [01:19<00:45, 283.67it/s]
64%|██████████|
| 23070/35912 [01:19<00:42, 303.72it/s]
64%|██████████|
| 23101/35912 [01:19<00:43, 296.20it/s]
64%|██████████|
| 23131/35912 [01:20<00:43, 294.94it/s]
64%|██████████|
| 23161/35912 [01:20<00:45, 281.67it/s]
65%|██████████|
| 23195/35912 [01:20<00:43, 294.09it/s]
65%|██████████|
| 23226/35912 [01:20<00:42, 296.35it/s]
65%|██████████|
| 23256/35912 [01:20<00:45, 280.20it/s]
65%|██████████|
| 23289/35912 [01:20<00:43, 291.37it/s]
65%|██████████|
| 23319/35912 [01:20<00:42, 293.28it/s]
65%|██████████|
| 23349/35912 [01:20<00:42, 293.76it/s]
65%|██████████|
| 23379/35912 [01:20<00:42, 294.11it/s]
65%|██████████|
| 23409/35912 [01:20<00:42, 295.21it/s]
65%|██████████|
| 23439/35912 [01:21<00:43, 285.05it/s]
65%|██████████|
| 23470/35912 [01:21<00:42, 290.68it/s]
65%|██████████|
| 23500/35912 [01:21<00:42, 290.25it/s]
66%|██████████|
| 23530/35912 [01:21<00:42, 289.11it/s]
66%|██████████|
| 23561/35912 [01:21<00:42, 293.62it/s]
66%|██████████|
| 23592/35912 [01:21<00:41, 296.02it/s]
66%|██████████|
| 23622/35912 [01:21<00:41, 294.82it/s]
66%|██████████|
| 23653/35912 [01:21<00:41, 297.72it/s]
66%|██████████|
| 23683/35912 [01:21<00:41, 294.26it/s]
66%|██████████|
| 23715/35912 [01:22<00:40, 300.91it/s]
66%|██████████|
| 23746/35912 [01:22<00:40, 297.69it/s]
66%|██████████|
```

```

23776/35912 [01:22<00:41, 290.84it/s]
66%|███████████|
| 23806/35912 [01:22<00:44, 274.46it/s]
66%|███████████|
| 23835/35912 [01:22<00:43, 276.77it/s]
66%|███████████|
| 23863/35912 [01:22<00:43, 274.69it/s]
67%|███████████|
| 23891/35912 [01:22<00:43, 274.05it/s]
67%|███████████|
| 23922/35912 [01:22<00:42, 283.36it/s]
67%|███████████|
| 23951/35912 [01:22<00:42, 284.70it/s]
67%|███████████|
| 23981/35912 [01:22<00:41, 287.69it/s]
67%|███████████|
| 24010/35912 [01:23<00:42, 280.26it/s]
67%|███████████|
| 24040/35912 [01:23<00:41, 282.90it/s]
67%|███████████|
| 24074/35912 [01:23<00:39, 296.55it/s]
67%|███████████|
| 24104/35912 [01:23<00:40, 293.46it/s]
67%|███████████|
| 24134/35912 [01:23<00:39, 294.76it/s]
67%|███████████|
| 24164/35912 [01:23<00:39, 293.94it/s]
67%|███████████|
| 24194/35912 [01:23<00:40, 292.51it/s]
67%|███████████|
| 24227/35912 [01:23<00:38, 300.58it/s]
68%|███████████|
| 24258/35912 [01:23<00:38, 300.06it/s]
68%|███████████|
| 24289/35912 [01:24<00:39, 292.10it/s]
68%|███████████|
| 24322/35912 [01:24<00:38, 299.46it/s]
68%|███████████|
| 24355/35912 [01:24<00:37, 307.38it/s]
68%|███████████|
| 24387/35912 [01:24<00:37, 305.96it/s]
68%|███████████|
| 24418/35912 [01:24<00:37, 304.66it/s]
68%|███████████|
| 24451/35912 [01:24<00:36, 311.23it/s]
68%|███████████|
| 24483/35912 [01:24<00:38, 294.19it/s]
68%|███████████|
| 24513/35912 [01:24<00:38, 295.27it/s]
68%|███████████|
| 24543/35912 [01:24<00:38, 291.71it/s]
68%|███████████|
| 24577/35912 [01:24<00:37, 303.30it/s]
69%|███████████|
| 24608/35912 [01:25<00:38, 290.93it/s]
69%|███████████|
| 24638/35912 [01:25<00:38, 291.27it/s]

```

```
69%|███████████|
| 24668/35912 [01:25<00:39, 286.50it/s]
69%|███████████|
| 24698/35912 [01:25<00:39, 286.49it/s]
69%|███████████|
| 24727/35912 [01:25<00:41, 271.64it/s]
69%|███████████|
| 24760/35912 [01:25<00:39, 285.57it/s]
69%|███████████|
| 24791/35912 [01:25<00:38, 291.06it/s]
69%|███████████|
| 24821/35912 [01:25<00:38, 288.84it/s]
69%|███████████|
| 24851/35912 [01:25<00:39, 278.52it/s]
69%|███████████|
| 24880/35912 [01:26<00:39, 281.27it/s]
69%|███████████|
| 24911/35912 [01:26<00:38, 288.72it/s]
69%|███████████|
| 24946/35912 [01:26<00:36, 300.24it/s]
70%|███████████|
| 24977/35912 [01:26<00:37, 292.22it/s]
70%|███████████|
| 25007/35912 [01:26<00:37, 289.63it/s]
70%|███████████|
| 25038/35912 [01:26<00:36, 294.84it/s]
70%|███████████|
| 25071/35912 [01:26<00:35, 303.96it/s]
70%|███████████|
| 25102/35912 [01:26<00:35, 301.54it/s]
70%|███████████|
| 25134/35912 [01:26<00:35, 305.33it/s]
70%|███████████|
| 25165/35912 [01:26<00:35, 298.99it/s]
70%|███████████|
| 25196/35912 [01:27<00:35, 301.56it/s]
70%|███████████|
| 25227/35912 [01:27<00:36, 293.09it/s]
70%|███████████|
| 25257/35912 [01:27<00:36, 291.93it/s]
70%|███████████|
| 25288/35912 [01:27<00:36, 294.81it/s]
71%|███████████|
| 25318/35912 [01:27<00:36, 293.12it/s]
71%|███████████|
| 25348/35912 [01:27<00:37, 282.88it/s]
71%|███████████|
| 25377/35912 [01:27<00:37, 282.71it/s]
71%|███████████|
| 25410/35912 [01:27<00:35, 294.04it/s]
71%|███████████|
| 25440/35912 [01:27<00:35, 295.17it/s]
71%|███████████|
| 25470/35912 [01:28<00:36, 286.65it/s]
71%|███████████|
| 25499/35912 [01:28<00:36, 282.81it/s]
71%|███████████|
```

[illegible]

[illegible]



[illegible]



[illegible]



|             |                           |     |  |
|-------------|---------------------------|-----|--|
| 30821/35912 | [01:46<00:17, 290.52it/s] | 86% |  |
| 30851/35912 | [01:46<00:17, 288.46it/s] | 86% |  |
| 30887/35912 | [01:46<00:16, 304.63it/s] | 86% |  |
| 30918/35912 | [01:46<00:17, 292.61it/s] | 86% |  |
| 30948/35912 | [01:46<00:17, 290.75it/s] | 86% |  |
| 30978/35912 | [01:47<00:17, 278.98it/s] | 86% |  |
| 31011/35912 | [01:47<00:16, 291.20it/s] | 86% |  |
| 31041/35912 | [01:47<00:17, 277.74it/s] | 87% |  |
| 31071/35912 | [01:47<00:17, 277.98it/s] | 87% |  |
| 31100/35912 | [01:47<00:17, 269.18it/s] | 87% |  |
| 31132/35912 | [01:47<00:17, 278.43it/s] | 87% |  |
| 31161/35912 | [01:47<00:16, 281.20it/s] | 87% |  |
| 31192/35912 | [01:47<00:16, 283.93it/s] | 87% |  |
| 31221/35912 | [01:47<00:17, 271.53it/s] | 87% |  |
| 31249/35912 | [01:48<00:17, 267.95it/s] | 87% |  |
| 31276/35912 | [01:48<00:17, 261.76it/s] | 87% |  |
| 31307/35912 | [01:48<00:16, 274.04it/s] | 87% |  |
| 31336/35912 | [01:48<00:16, 272.58it/s] | 87% |  |
| 31370/35912 | [01:48<00:15, 289.28it/s] | 87% |  |
| 31401/35912 | [01:48<00:15, 294.59it/s] | 88% |  |
| 31431/35912 | [01:48<00:16, 279.10it/s] | 88% |  |
| 31462/35912 | [01:48<00:15, 287.12it/s] | 88% |  |
| 31495/35912 | [01:48<00:14, 298.17it/s] | 88% |  |
| 31526/35912 | [01:48<00:14, 296.67it/s] | 88% |  |
| 31559/35912 | [01:49<00:14, 303.64it/s] | 88% |  |
| 31590/35912 | [01:49<00:14, 296.99it/s] | 88% |  |
| 31622/35912 | [01:49<00:14, 302.03it/s] | 88% |  |
| 31653/35912 | [01:49<00:14, 296.79it/s] | 88% |  |
| 31684/35912 | [01:49<00:14, 299.13it/s] |     |  |



111/180

[illegible]



[illegible]



```
In [117]: # average Word2Vec
compute average word2vec for each review.
tfidf_w2v_vectors_cv = []; # the avg-w2v for each sentence/review is stored in
this list
for sentence in tqdm(X_cv['essay']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if (word in glove_words) and (word in tfidf_words):
 vec = model[word] # getting the vector for each word
 # here we are multiplying idf value(dictionary[word]) and the tf value
 # (sentence.count(word)/len(sentence.split()))
 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split
 ())) # getting the tfidf value for each word
 vector += (vec * tf_idf) # calculating tfidf weighted w2v
 tf_idf_weight += tf_idf
 if tf_idf_weight != 0:
 vector /= tf_idf_weight
 tfidf_w2v_vectors_cv.append(vector)

print(len(tfidf_w2v_vectors_cv))
print(len(tfidf_w2v_vectors_cv[0]))
```

```

0%|
| 0/17688 [00:00<?, ?it/s]
0%||
| 30/17688 [00:00<00:59, 294.90it/s]
0%||
| 60/17688 [00:00<00:59, 294.04it/s]
1%||
| 90/17688 [00:00<00:59, 294.30it/s]
1%||
| 121/17688 [00:00<00:58, 298.21it/s]
1%||
| 146/17688 [00:00<01:02, 281.21it/s]
1%||
| 173/17688 [00:00<01:03, 275.43it/s]
1%||
| 199/17688 [00:00<01:04, 269.16it/s]
1%||
| 232/17688 [00:00<01:01, 284.38it/s]
1%||
| 263/17688 [00:00<01:00, 290.20it/s]
2%||
| 295/17688 [00:01<00:58, 297.93it/s]
2%||
| 328/17688 [00:01<00:56, 306.25it/s]
2%||
| 360/17688 [00:01<00:56, 307.81it/s]
2%||
| 391/17688 [00:01<00:58, 294.66it/s]
2%||
| 423/17688 [00:01<00:57, 301.21it/s]
3%||
| 454/17688 [00:01<00:57, 300.51it/s]
3%||
| 485/17688 [00:01<01:00, 285.94it/s]
3%||
| 516/17688 [00:01<00:58, 291.33it/s]
3%||
| 546/17688 [00:01<01:03, 268.15it/s]
3%||
| 579/17688 [00:01<01:00, 282.86it/s]
3%||
| 612/17688 [00:02<00:57, 294.94it/s]
4%||
| 643/17688 [00:02<01:01, 278.60it/s]
4%||
| 675/17688 [00:02<00:59, 286.95it/s]
4%||
| 705/17688 [00:02<01:01, 278.06it/s]
4%||
| 738/17688 [00:02<00:58, 291.27it/s]
4%||
| 768/17688 [00:02<00:57, 292.35it/s]
5%||
| 798/17688 [00:02<00:57, 291.41it/s]
5%||
| 832/17688 [00:02<00:55, 302.25it/s]
5%||

```

| 863/17688 [00:02<00:55, 302.99it/s]  
5%|██████  
| 896/17688 [00:03<00:54, 308.24it/s]  
5%|██████  
| 928/17688 [00:03<00:53, 311.02it/s]  
5%|██████  
| 960/17688 [00:03<00:53, 311.17it/s]  
6%|██████  
| 994/17688 [00:03<00:52, 315.97it/s]  
6%|██████  
| 1026/17688 [00:03<00:54, 305.63it/s]  
6%|██████  
| 1057/17688 [00:03<00:58, 282.08it/s]  
6%|██████  
| 1088/17688 [00:03<00:57, 286.92it/s]  
6%|██████  
| 1120/17688 [00:03<00:56, 295.50it/s]  
7%|██████  
| 1150/17688 [00:03<00:57, 285.24it/s]  
7%|██████  
| 1184/17688 [00:04<00:55, 296.80it/s]  
7%|██████  
| 1216/17688 [00:04<00:55, 298.54it/s]  
7%|██████  
| 1247/17688 [00:04<00:56, 291.09it/s]  
7%|██████  
| 1279/17688 [00:04<00:55, 297.76it/s]  
7%|██████  
| 1314/17688 [00:04<00:52, 310.28it/s]  
8%|██████  
| 1346/17688 [00:04<00:52, 308.86it/s]  
8%|██████  
| 1378/17688 [00:04<00:58, 280.42it/s]  
8%|██████  
| 1407/17688 [00:04<00:58, 277.76it/s]  
8%|██████  
| 1440/17688 [00:04<00:55, 291.05it/s]  
8%|██████  
| 1470/17688 [00:05<00:55, 292.19it/s]  
8%|██████  
| 1502/17688 [00:05<00:54, 298.55it/s]  
9%|██████  
| 1533/17688 [00:05<00:54, 297.80it/s]  
9%|██████  
| 1563/17688 [00:05<00:54, 296.90it/s]  
9%|██████  
| 1593/17688 [00:05<00:56, 283.75it/s]  
9%|██████  
| 1625/17688 [00:05<00:54, 293.14it/s]  
9%|██████  
| 1655/17688 [00:05<00:54, 291.96it/s]  
10%|██████  
| 1685/17688 [00:05<00:55, 287.79it/s]  
10%|██████  
| 1715/17688 [00:05<00:55, 289.89it/s]  
10%|██████  
| 1745/17688 [00:05<00:54, 290.53it/s]

```

10%|██████████|
| 1775/17688 [00:06<00:57, 278.07it/s]
10%|██████████|
| 1805/17688 [00:06<00:55, 283.72it/s]
10%|██████████|
| 1834/17688 [00:06<00:55, 284.12it/s]
11%|██████████|
| 1871/17688 [00:06<00:51, 304.84it/s]
11%|██████████|
| 1904/17688 [00:06<00:52, 302.78it/s]
11%|██████████|
| 1935/17688 [00:06<00:52, 300.72it/s]
11%|██████████|
| 1973/17688 [00:06<00:49, 318.60it/s]
11%|██████████|
| 2006/17688 [00:06<00:50, 310.40it/s]
12%|██████████|
| 2039/17688 [00:06<00:49, 315.37it/s]
12%|██████████|
| 2071/17688 [00:07<00:50, 306.99it/s]
12%|██████████|
| 2103/17688 [00:07<00:50, 309.22it/s]
12%|██████████|
| 2135/17688 [00:07<00:51, 301.18it/s]
12%|██████████|
| 2166/17688 [00:07<00:53, 292.84it/s]
12%|██████████|
| 2196/17688 [00:07<00:54, 286.74it/s]
13%|██████████|
| 2225/17688 [00:07<00:55, 276.44it/s]
13%|██████████|
| 2257/17688 [00:07<00:53, 286.11it/s]
13%|██████████|
| 2288/17688 [00:07<00:52, 291.45it/s]
13%|██████████|
| 2318/17688 [00:07<00:55, 277.12it/s]
13%|██████████|
| 2346/17688 [00:07<00:57, 267.08it/s]
13%|██████████|
| 2376/17688 [00:08<00:56, 268.96it/s]
14%|██████████|
| 2404/17688 [00:08<00:56, 270.82it/s]
14%|██████████|
| 2435/17688 [00:08<00:54, 279.42it/s]
14%|██████████|
| 2467/17688 [00:08<00:52, 289.89it/s]
14%|██████████|
| 2500/17688 [00:08<00:50, 298.64it/s]
14%|██████████|
| 2533/17688 [00:08<00:49, 305.92it/s]
14%|██████████|
| 2564/17688 [00:08<00:52, 289.35it/s]
15%|██████████|
| 2594/17688 [00:08<00:52, 287.66it/s]
15%|██████████|
| 2623/17688 [00:08<00:53, 283.52it/s]
15%|██████████|

```

| 2652/17688 [00:09<00:53, 279.05it/s]  
15%|██████████  
| 2681/17688 [00:09<00:54, 274.49it/s]  
15%|██████████  
| 2711/17688 [00:09<00:53, 279.53it/s]  
16%|██████████  
| 2743/17688 [00:09<00:51, 287.64it/s]  
16%|██████████  
| 2776/17688 [00:09<00:50, 293.80it/s]  
16%|██████████  
| 2806/17688 [00:09<00:53, 279.38it/s]  
16%|██████████  
| 2835/17688 [00:09<00:53, 275.47it/s]  
16%|██████████  
| 2863/17688 [00:09<00:54, 272.20it/s]  
16%|██████████  
| 2891/17688 [00:09<00:55, 267.64it/s]  
17%|██████████  
| 2921/17688 [00:10<00:53, 275.27it/s]  
17%|██████████  
| 2954/17688 [00:10<00:51, 287.62it/s]  
17%|██████████  
| 2987/17688 [00:10<00:49, 296.94it/s]  
17%|██████████  
| 3018/17688 [00:10<00:48, 300.10it/s]  
17%|██████████  
| 3050/17688 [00:10<00:48, 303.44it/s]  
17%|██████████  
| 3085/17688 [00:10<00:46, 313.74it/s]  
18%|██████████  
| 3119/17688 [00:10<00:45, 317.83it/s]  
18%|██████████  
| 3151/17688 [00:10<00:46, 312.22it/s]  
18%|██████████  
| 3183/17688 [00:10<00:47, 305.76it/s]  
18%|██████████  
| 3214/17688 [00:10<00:47, 303.66it/s]  
18%|██████████  
| 3246/17688 [00:11<00:47, 305.10it/s]  
19%|██████████  
| 3279/17688 [00:11<00:46, 310.64it/s]  
19%|██████████  
| 3312/17688 [00:11<00:45, 315.55it/s]  
19%|██████████  
| 3344/17688 [00:11<00:47, 299.37it/s]  
19%|██████████  
| 3375/17688 [00:11<00:48, 297.50it/s]  
19%|██████████  
| 3405/17688 [00:11<00:48, 295.84it/s]  
19%|██████████  
| 3435/17688 [00:11<00:49, 290.39it/s]  
20%|██████████  
| 3465/17688 [00:11<00:49, 287.58it/s]  
20%|██████████  
| 3495/17688 [00:11<00:49, 288.90it/s]  
20%|██████████  
| 3524/17688 [00:12<00:49, 286.05it/s]

20%|██████████  
| 3553/17688 [00:12<00:50, 279.17it/s]  
20%|██████████  
| 3581/17688 [00:12<00:53, 261.66it/s]  
20%|██████████  
| 3611/17688 [00:12<00:51, 270.80it/s]  
21%|██████████  
| 3644/17688 [00:12<00:49, 285.68it/s]  
21%|██████████  
| 3673/17688 [00:12<00:49, 280.54it/s]  
21%|██████████  
| 3704/17688 [00:12<00:48, 288.18it/s]  
21%|██████████  
| 3734/17688 [00:12<00:48, 287.65it/s]  
21%|██████████  
| 3763/17688 [00:12<00:48, 286.89it/s]  
21%|██████████  
| 3792/17688 [00:12<00:48, 285.50it/s]  
22%|██████████  
| 3821/17688 [00:13<00:49, 279.61it/s]  
22%|██████████  
| 3851/17688 [00:13<00:48, 283.23it/s]  
22%|██████████  
| 3883/17688 [00:13<00:47, 291.95it/s]  
22%|██████████  
| 3913/17688 [00:13<00:48, 285.33it/s]  
22%|██████████  
| 3942/17688 [00:13<00:48, 286.10it/s]  
22%|██████████  
| 3971/17688 [00:13<00:49, 276.81it/s]  
23%|██████████  
| 4000/17688 [00:13<00:49, 279.22it/s]  
23%|██████████  
| 4029/17688 [00:13<00:50, 269.26it/s]  
23%|██████████  
| 4058/17688 [00:13<00:49, 274.59it/s]  
23%|██████████  
| 4094/17688 [00:14<00:46, 290.82it/s]  
23%|██████████  
| 4124/17688 [00:14<00:46, 292.89it/s]  
23%|██████████  
| 4154/17688 [00:14<00:48, 278.03it/s]  
24%|██████████  
| 4184/17688 [00:14<00:47, 282.09it/s]  
24%|██████████  
| 4218/17688 [00:14<00:45, 294.40it/s]  
24%|██████████  
| 4248/17688 [00:14<00:46, 291.13it/s]  
24%|██████████  
| 4280/17688 [00:14<00:45, 297.79it/s]  
24%|██████████  
| 4314/17688 [00:14<00:44, 302.13it/s]  
25%|██████████  
| 4345/17688 [00:14<00:43, 303.79it/s]  
25%|██████████  
| 4376/17688 [00:14<00:44, 302.30it/s]  
25%|██████████



2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108, 2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121, 2122, 2123, 2124, 2125, 2126, 2127, 2128, 2129, 2130, 2131, 2132, 2133, 2134, 2135, 2136, 2137, 2138, 2139, 2140, 2141, 2142, 2143, 2144, 2145, 2146, 2147, 2148, 2149, 2150, 2151, 2152, 2153, 2154, 2155, 2156, 2157, 2158, 2159, 2160, 2161, 2162, 2163, 2164, 2165, 2166, 2167, 2168, 2169, 2170, 2171, 2172, 2173, 2174, 2175, 2176, 2177, 2178, 2179, 2180, 2181, 2182, 2183, 2184, 2185, 2186, 2187, 2188, 2189, 2190, 2191, 2192, 2193, 2194, 2195, 2196, 2197, 2198, 2199, 2200, 2201, 2202, 2203, 2204, 2205, 2206, 2207, 2208, 2209, 2210, 2211, 2212, 2213, 2214, 2215, 2216, 2217, 2218, 2219, 2220, 2221, 2222, 2223, 2224, 2225, 2226, 2227, 2228, 2229, 2230, 2231, 2232, 2233, 2234, 2235, 2236, 2237, 2238, 2239, 2240, 2241, 2242, 2243, 2244, 2245, 2246, 2247, 2248, 2249, 2250, 2251, 2252, 2253, 2254, 2255, 2256, 2257, 2258, 2259, 2260, 2261, 2262, 2263, 2264, 2265, 2266, 2267, 2268, 2269, 2270, 2271, 2272, 2273, 2274, 2275, 2276, 2277, 2278, 2279, 2280, 2281, 2282, 2283, 2284, 2285, 2286, 2287, 2288, 2289, 2290, 2291, 2292, 2293, 2294, 2295, 2296, 2297, 2298, 2299, 2300, 2301, 2302, 2303, 2304, 2305, 2306, 2307, 2308, 2309, 2310, 2311, 2312, 2313, 2314, 2315, 2316, 2317, 2318, 2319, 2320, 2321, 2322, 2323, 2324, 2325, 2326, 2327, 2328, 2329, 2330, 2331, 2332, 2333, 2334, 2335, 2336, 2337, 2338, 2339, 2340, 2341, 2342, 2343, 2344, 2345, 2346, 2347, 2348, 2349, 2350, 2351, 2352, 2353, 2354, 2355, 2356, 2357, 2358, 2359, 2360, 2361, 2362, 2363, 2364, 2365, 2366, 2367, 2368, 2369, 2370, 2371, 2372, 2373, 2374, 2375, 2376, 2377, 2378, 2379, 2380, 2381, 2382, 2383, 2384, 2385, 2386, 2387, 2388, 2389, 2390, 2391, 2392, 2393, 2394, 2395, 2396, 2397, 2398, 2399, 2400, 2401, 2402, 2403, 2404, 2405, 2406, 2407, 2408, 2409, 2410, 2411, 2412, 2413, 2414, 2415, 2416, 2417, 2418, 2419, 2420, 2421, 2422, 2423, 2424, 2425, 2426, 2427, 2428, 2429, 2430, 2431, 2432, 2433, 2434, 2435, 2436, 2437, 2438, 2439, 2440, 2441, 2442, 2443, 2444, 2445, 2446, 2447, 2448, 2449, 2450, 2451, 2452, 2453, 2454, 2455, 2456, 2457, 2458, 2459, 2460, 2461, 2462, 2463, 2464, 2465, 2466, 2467, 2468, 2469, 2470, 2471, 2472, 2473, 2474, 2475, 2476, 2477, 2478, 2479, 2480, 2481, 2482, 2483, 2484, 2485, 2486, 2487, 2488, 2489, 2490, 2491, 2492, 2493, 2494, 2495, 2496, 2497, 2498, 2499, 2500, 2501, 2502, 2503, 2504, 2505, 2506, 2507, 2508, 2509, 2510, 2511, 2512, 2513, 2514, 2515, 2516, 2517, 2518, 2519, 2520, 2521, 2522, 2523, 2524, 2525, 2526, 2527, 2528, 2529, 2530, 2531, 2532, 2533, 2534, 2535, 2536, 2537, 2538, 2539, 2540, 2541, 2542, 2543, 2544, 2545, 2546, 2547, 2548, 2549, 2550, 2551, 2552, 2553, 2554, 2555, 2556, 2557, 2558, 2559, 2560, 2561, 2562, 2563, 2564, 2565, 2566, 2567, 2568, 2569, 2570, 2571, 2572, 2573, 2574, 2575, 2576, 2577, 2578, 2579, 2580, 2581, 2582, 2583, 2584, 2585, 2586, 2587, 2588, 2589, 2590, 2591, 2592, 2593, 2594, 2595, 2596, 2597, 2598, 2599, 2600, 2601, 2602, 2603, 2604, 2605, 2606, 2607, 2608, 2609, 2610, 2611, 2612, 2613, 2614, 2615, 2616, 2617, 2618, 2619, 2620, 2621, 2622, 2623, 2624, 2625, 2626, 2627, 2628, 2629, 2630, 2631, 2632, 2633, 2634, 2635, 2636, 2637, 2638, 2639, 2640, 2641, 2642, 2643, 2644, 2645, 2646, 2647, 2648, 2649, 2650, 2651, 2652, 2653, 2654, 2655, 2656, 2657, 2658, 2659, 2660, 2661, 2662, 2663, 2664, 2665, 2666, 2667, 2668, 2669, 2670, 2671, 2672, 2673, 2674, 2675, 2676, 2677, 2678, 2679, 2680, 2681, 2682, 2683, 2684, 2685, 2686, 26

```

30%|███████████|
| 5305/17688 [00:18<00:42, 291.80it/s]
30%|███████████|
| 5335/17688 [00:18<00:44, 275.06it/s]
30%|███████████|
| 5365/17688 [00:18<00:43, 280.72it/s]
31%|███████████|
| 5395/17688 [00:18<00:43, 284.03it/s]
31%|███████████|
| 5424/17688 [00:18<00:44, 277.82it/s]
31%|███████████|
| 5454/17688 [00:18<00:43, 283.53it/s]
31%|███████████|
| 5486/17688 [00:18<00:41, 291.38it/s]
31%|███████████|
| 5516/17688 [00:18<00:43, 279.39it/s]
31%|███████████|
| 5545/17688 [00:19<00:43, 281.07it/s]
32%|███████████|
| 5579/17688 [00:19<00:40, 295.92it/s]
32%|███████████|
| 5609/17688 [00:19<00:40, 294.75it/s]
32%|███████████|
| 5643/17688 [00:19<00:39, 306.38it/s]
32%|███████████|
| 5674/17688 [00:19<00:40, 294.60it/s]
32%|███████████|
| 5704/17688 [00:19<00:40, 295.56it/s]
32%|███████████|
| 5734/17688 [00:19<00:40, 295.36it/s]
33%|███████████|
| 5764/17688 [00:19<00:41, 286.78it/s]
33%|███████████|
| 5794/17688 [00:19<00:41, 290.01it/s]
33%|███████████|
| 5825/17688 [00:19<00:40, 295.10it/s]
33%|███████████|
| 5855/17688 [00:20<00:40, 289.10it/s]
33%|███████████|
| 5891/17688 [00:20<00:38, 305.90it/s]
33%|███████████|
| 5924/17688 [00:20<00:38, 308.61it/s]
34%|███████████|
| 5958/17688 [00:20<00:37, 315.00it/s]
34%|███████████|
| 5990/17688 [00:20<00:38, 301.55it/s]
34%|███████████|
| 6021/17688 [00:20<00:39, 298.99it/s]
34%|███████████|
| 6055/17688 [00:20<00:37, 307.94it/s]
34%|███████████|
| 6086/17688 [00:20<00:38, 303.37it/s]
35%|███████████|
| 6117/17688 [00:20<00:37, 304.67it/s]
35%|███████████|
| 6148/17688 [00:21<00:38, 296.84it/s]
35%|███████████|

```

```
6178/17688 [00:21<00:39, 288.58it/s]
35%|███████████|
| 6207/17688 [00:21<00:39, 287.52it/s]
35%|███████████|
| 6239/17688 [00:21<00:38, 294.32it/s]
35%|███████████|
| 6270/17688 [00:21<00:38, 297.37it/s]
36%|███████████|
| 6301/17688 [00:21<00:37, 300.41it/s]
36%|███████████|
| 6335/17688 [00:21<00:36, 307.29it/s]
36%|███████████|
| 6366/17688 [00:21<00:37, 301.18it/s]
36%|███████████|
| 6397/17688 [00:21<00:39, 286.37it/s]
36%|███████████|
| 6430/17688 [00:21<00:38, 296.01it/s]
37%|███████████|
| 6460/17688 [00:22<00:38, 291.38it/s]
37%|███████████|
| 6490/17688 [00:22<00:39, 283.34it/s]
37%|███████████|
| 6519/17688 [00:22<00:40, 277.36it/s]
37%|███████████|
| 6547/17688 [00:22<00:40, 277.54it/s]
37%|███████████|
| 6579/17688 [00:22<00:38, 288.47it/s]
37%|███████████|
| 6609/17688 [00:22<00:38, 289.54it/s]
38%|███████████|
| 6641/17688 [00:22<00:37, 297.44it/s]
38%|███████████|
| 6671/17688 [00:22<00:38, 289.82it/s]
38%|███████████|
| 6705/17688 [00:22<00:36, 301.05it/s]
38%|███████████|
| 6737/17688 [00:23<00:36, 304.11it/s]
38%|███████████|
| 6768/17688 [00:23<00:35, 303.41it/s]
38%|███████████|
| 6799/17688 [00:23<00:35, 302.92it/s]
39%|███████████|
| 6830/17688 [00:23<00:36, 294.82it/s]
39%|███████████|
| 6860/17688 [00:23<00:37, 291.42it/s]
39%|███████████|
| 6890/17688 [00:23<00:38, 278.62it/s]
39%|███████████|
| 6923/17688 [00:23<00:37, 287.92it/s]
39%|███████████|
| 6952/17688 [00:23<00:37, 286.21it/s]
39%|███████████|
| 6983/17688 [00:23<00:37, 289.09it/s]
40%|███████████|
| 7013/17688 [00:24<00:36, 290.81it/s]
40%|███████████|
| 7043/17688 [00:24<00:36, 288.66it/s]
```

[illegible]



| Iteration  | Time (s)    | Throughput (it/s) |
|------------|-------------|-------------------|
| 8840/17688 | 00:30<00:30 | 290.87            |
| 8870/17688 | 00:30<00:30 | 292.05            |
| 8900/17688 | 00:30<00:30 | 284.60            |
| 8929/17688 | 00:30<00:31 | 278.21            |
| 8961/17688 | 00:30<00:30 | 287.42            |
| 8990/17688 | 00:30<00:32 | 269.20            |
| 9019/17688 | 00:30<00:31 | 273.77            |
| 9050/17688 | 00:31<00:30 | 281.59            |
| 9079/17688 | 00:31<00:31 | 270.05            |
| 9109/17688 | 00:31<00:30 | 277.05            |
| 9139/17688 | 00:31<00:30 | 282.17            |
| 9168/17688 | 00:31<00:32 | 264.50            |
| 9198/17688 | 00:31<00:31 | 273.69            |
| 9226/17688 | 00:31<00:31 | 271.77            |
| 9257/17688 | 00:31<00:30 | 279.37            |
| 9286/17688 | 00:31<00:30 | 277.83            |
| 9314/17688 | 00:31<00:30 | 277.87            |
| 9344/17688 | 00:32<00:29 | 280.40            |
| 9374/17688 | 00:32<00:29 | 282.99            |
| 9403/17688 | 00:32<00:29 | 277.92            |
| 9431/17688 | 00:32<00:30 | 273.87            |
| 9462/17688 | 00:32<00:29 | 279.40            |
| 9496/17688 | 00:32<00:28 | 290.86            |
| 9526/17688 | 00:32<00:28 | 288.69            |
| 9556/17688 | 00:32<00:27 | 290.53            |
| 9592/17688 | 00:32<00:26 | 306.24            |
| 9623/17688 | 00:33<00:27 | 288.74            |
| 9653/17688 | 00:33<00:28 | 286.41            |

127/180

```
| 60%| ██████████
| 10585/17688 [00:36<00:25, 282.33it/s]
60%| ██████████
| 10615/17688 [00:36<00:24, 285.17it/s]
60%| ██████████
| 10647/17688 [00:36<00:23, 293.40it/s]
60%| ██████████
| 10684/17688 [00:36<00:22, 312.25it/s]
61%| ██████████
| 10716/17688 [00:36<00:22, 308.43it/s]
61%| ██████████
| 10748/17688 [00:36<00:24, 288.49it/s]
61%| ██████████
| 10778/17688 [00:37<00:24, 287.06it/s]
61%| ██████████
| 10809/17688 [00:37<00:23, 291.32it/s]
61%| ██████████
| 10839/17688 [00:37<00:24, 284.91it/s]
61%| ██████████
| 10868/17688 [00:37<00:24, 272.92it/s]
62%| ██████████
| 10896/17688 [00:37<00:24, 272.02it/s]
62%| ██████████
| 10927/17688 [00:37<00:24, 281.08it/s]
62%| ██████████
| 10956/17688 [00:37<00:23, 281.44it/s]
62%| ██████████
| 10988/17688 [00:37<00:23, 289.05it/s]
62%| ██████████
| 11024/17688 [00:37<00:21, 306.64it/s]
63%| ██████████
| 11056/17688 [00:37<00:22, 297.79it/s]
63%| ██████████
| 11087/17688 [00:38<00:22, 298.12it/s]
63%| ██████████
| 11118/17688 [00:38<00:22, 296.64it/s]
63%| ██████████
| 11148/17688 [00:38<00:22, 296.99it/s]
63%| ██████████
| 11178/17688 [00:38<00:23, 280.60it/s]
63%| ██████████
| 11208/17688 [00:38<00:22, 284.74it/s]
64%| ██████████
| 11239/17688 [00:38<00:22, 290.46it/s]
64%| ██████████
| 11269/17688 [00:38<00:22, 280.36it/s]
64%| ██████████
| 11303/17688 [00:38<00:21, 293.08it/s]
64%| ██████████
| 11333/17688 [00:38<00:21, 293.63it/s]
64%| ██████████
| 11363/17688 [00:39<00:21, 291.45it/s]
64%| ██████████
| 11393/17688 [00:39<00:21, 292.45it/s]
65%| ██████████
| 11423/17688 [00:39<00:21, 289.82it/s]
65%| ██████████
```



129/180

```
| 70%| ██████████
| 12352/17688 [00:42<00:18, 292.91it/s]
70%| ██████████
| 12384/17688 [00:42<00:17, 298.25it/s]
70%| ██████████
| 12414/17688 [00:42<00:18, 292.87it/s]
70%| ██████████
| 12450/17688 [00:42<00:17, 308.08it/s]
71%| ██████████
| 12482/17688 [00:42<00:16, 307.32it/s]
71%| ██████████
| 12513/17688 [00:42<00:17, 290.23it/s]
71%| ██████████
| 12543/17688 [00:43<00:17, 289.09it/s]
71%| ██████████
| 12575/17688 [00:43<00:17, 296.29it/s]
71%| ██████████
| 12605/17688 [00:43<00:17, 290.73it/s]
71%| ██████████
| 12635/17688 [00:43<00:17, 289.44it/s]
72%| ██████████
| 12669/17688 [00:43<00:16, 301.56it/s]
72%| ██████████
| 12700/17688 [00:43<00:16, 294.76it/s]
72%| ██████████
| 12730/17688 [00:43<00:18, 273.87it/s]
72%| ██████████
| 12758/17688 [00:43<00:18, 265.71it/s]
72%| ██████████
| 12786/17688 [00:43<00:18, 268.50it/s]
72%| ██████████
| 12818/17688 [00:44<00:17, 279.37it/s]
73%| ██████████
| 12847/17688 [00:44<00:17, 276.22it/s]
73%| ██████████
| 12875/17688 [00:44<00:18, 261.31it/s]
73%| ██████████
| 12904/17688 [00:44<00:17, 268.02it/s]
73%| ██████████
| 12932/17688 [00:44<00:18, 261.83it/s]
73%| ██████████
| 12961/17688 [00:44<00:17, 269.14it/s]
73%| ██████████
| 12996/17688 [00:44<00:16, 288.66it/s]
74%| ██████████
| 13026/17688 [00:44<00:16, 290.51it/s]
74%| ██████████
| 13056/17688 [00:44<00:15, 292.66it/s]
74%| ██████████
| 13088/17688 [00:44<00:15, 297.25it/s]
74%| ██████████
| 13118/17688 [00:45<00:17, 268.00it/s]
74%| ██████████
| 13152/17688 [00:45<00:15, 284.23it/s]
75%| ██████████
| 13182/17688 [00:45<00:15, 281.70it/s]
75%| ██████████
```

```

13211/17688 [00:45<00:15, 282.70it/s]
75%|
13243/17688 [00:45<00:15, 289.98it/s]
75%|
13280/17688 [00:45<00:14, 304.95it/s]
75%|
13311/17688 [00:45<00:15, 281.68it/s]
75%|
13342/17688 [00:45<00:15, 286.63it/s]
76%|
13372/17688 [00:46<00:15, 280.17it/s]
76%|
13401/17688 [00:46<00:15, 278.39it/s]
76%|
13430/17688 [00:46<00:15, 276.36it/s]
76%|
13458/17688 [00:46<00:15, 276.03it/s]
76%|
13486/17688 [00:46<00:15, 276.60it/s]
76%|
13516/17688 [00:46<00:14, 282.65it/s]
77%|
13546/17688 [00:46<00:14, 286.22it/s]
77%|
13579/17688 [00:46<00:14, 292.76it/s]
77%|
13609/17688 [00:46<00:14, 288.33it/s]
77%|
13638/17688 [00:46<00:14, 284.82it/s]
77%|
13667/17688 [00:47<00:14, 284.06it/s]
77%|
13696/17688 [00:47<00:14, 275.47it/s]
78%|
13726/17688 [00:47<00:14, 281.03it/s]
78%|
13756/17688 [00:47<00:13, 285.86it/s]
78%|
13785/17688 [00:47<00:13, 286.47it/s]
78%|
13820/17688 [00:47<00:12, 301.61it/s]
78%|
13851/17688 [00:47<00:12, 299.05it/s]
78%|
13882/17688 [00:47<00:13, 288.19it/s]
79%|
13914/17688 [00:47<00:12, 295.62it/s]
79%|
13944/17688 [00:48<00:13, 282.90it/s]
79%|
13973/17688 [00:48<00:13, 281.90it/s]
79%|
14006/17688 [00:48<00:12, 294.22it/s]
79%|
14037/17688 [00:48<00:12, 296.44it/s]
80%|
14067/17688 [00:48<00:12, 288.31it/s]

```

[illegible]

|             |                           |     |  |
|-------------|---------------------------|-----|--|
| 14981/17688 | [00:51<00:09, 290.17it/s] | 85% |  |
| 15011/17688 | [00:51<00:09, 280.17it/s] | 85% |  |
| 15042/17688 | [00:51<00:09, 287.89it/s] | 85% |  |
| 15071/17688 | [00:51<00:09, 278.81it/s] | 85% |  |
| 15100/17688 | [00:51<00:09, 280.66it/s] | 86% |  |
| 15133/17688 | [00:52<00:08, 293.27it/s] | 86% |  |
| 15163/17688 | [00:52<00:08, 294.62it/s] | 86% |  |
| 15194/17688 | [00:52<00:08, 297.58it/s] | 86% |  |
| 15224/17688 | [00:52<00:08, 297.65it/s] | 86% |  |
| 15255/17688 | [00:52<00:08, 299.74it/s] | 86% |  |
| 15286/17688 | [00:52<00:08, 298.62it/s] | 87% |  |
| 15317/17688 | [00:52<00:07, 298.70it/s] | 87% |  |
| 15350/17688 | [00:52<00:07, 305.96it/s] | 87% |  |
| 15381/17688 | [00:52<00:07, 290.19it/s] | 87% |  |
| 15411/17688 | [00:53<00:07, 289.07it/s] | 87% |  |
| 15442/17688 | [00:53<00:07, 294.43it/s] | 87% |  |
| 15472/17688 | [00:53<00:08, 274.42it/s] | 88% |  |
| 15501/17688 | [00:53<00:07, 275.95it/s] | 88% |  |
| 15532/17688 | [00:53<00:07, 284.00it/s] | 88% |  |
| 15561/17688 | [00:53<00:07, 281.02it/s] | 88% |  |
| 15597/17688 | [00:53<00:06, 300.27it/s] | 88% |  |
| 15628/17688 | [00:53<00:06, 301.59it/s] | 89% |  |
| 15663/17688 | [00:53<00:06, 311.52it/s] | 89% |  |
| 15695/17688 | [00:53<00:06, 311.52it/s] | 89% |  |
| 15727/17688 | [00:54<00:06, 292.76it/s] | 89% |  |
| 15757/17688 | [00:54<00:06, 292.54it/s] | 89% |  |
| 15787/17688 | [00:54<00:06, 286.54it/s] | 89% |  |
| 15816/17688 | [00:54<00:06, 280.31it/s] | 90% |  |
| 15851/17688 | [00:54<00:06, 297.56it/s] |     |  |



```
17630/17688 [01:00<00:00, 284.37it/s]
```

```
100%|██|
██████| 17659/17688 [01:00<00:00, 284.58it/s]
100%|██|
██████| 17688/17688 [01:00<00:00, 290.69it/s]

17688
300
```



```

In [118]: # average Word2Vec
compute average word2vec for each review.
tfidf_w2v_vectors_test = []; # the avg-w2v for each sentence/review is stored
in this list
for sentence in tqdm(X_test['essay']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if (word in glove_words) and (word in tfidf_words):
 vec = model[word] # getting the vector for each word
 # here we are multiplying idf value(dictionary[word]) and the tf value
 # (sentence.count(word)/len(sentence.split()))
 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
 # getting the tfidf value for each word
 vector += (vec * tf_idf) # calculating tfidf weighted w2v
 tf_idf_weight += tf_idf
 if tf_idf_weight != 0:
 vector /= tf_idf_weight
 tfidf_w2v_vectors_test.append(vector)

print(len(tfidf_w2v_vectors_test))
print(len(tfidf_w2v_vectors_test[0]))

```

```

0%|
| 0/26400 [00:00<?, ?it/s]
0%|
| 27/26400 [00:00<01:39, 265.44it/s]
0%||
| 55/26400 [00:00<01:38, 267.54it/s]
0%||
| 89/26400 [00:00<01:32, 285.29it/s]
0%||
| 118/26400 [00:00<01:32, 285.22it/s]
1%||
| 145/26400 [00:00<01:34, 278.12it/s]
1%||
| 172/26400 [00:00<01:35, 274.18it/s]
1%||
| 207/26400 [00:00<01:30, 290.52it/s]
1%|█
| 235/26400 [00:00<01:31, 286.64it/s]
1%|█
| 265/26400 [00:00<01:31, 286.59it/s]
1%|█
| 300/26400 [00:01<01:26, 301.71it/s]
1%|█
| 330/26400 [00:01<01:30, 287.60it/s]
1%|█
| 362/26400 [00:01<01:28, 293.57it/s]
1%|█
| 395/26400 [00:01<01:26, 299.72it/s]
2%|█
| 425/26400 [00:01<01:28, 293.89it/s]
2%|█
| 455/26400 [00:01<01:30, 288.27it/s]
2%|█
| 484/26400 [00:01<01:31, 283.94it/s]
2%|█
| 515/26400 [00:01<01:29, 288.26it/s]
2%|█
| 546/26400 [00:01<01:29, 288.93it/s]
2%|█
| 575/26400 [00:01<01:31, 281.90it/s]
2%|█
| 609/26400 [00:02<01:27, 295.79it/s]
2%|█
| 639/26400 [00:02<01:27, 294.66it/s]
3%|█
| 669/26400 [00:02<01:28, 289.62it/s]
3%|█
| 699/26400 [00:02<01:33, 275.20it/s]
3%|█
| 733/26400 [00:02<01:28, 290.60it/s]
3%|█
| 763/26400 [00:02<01:27, 292.73it/s]
3%|█
| 793/26400 [00:02<01:29, 285.84it/s]
3%|█
| 822/26400 [00:02<01:36, 263.85it/s]
3%|█

```

| 852/26400 [00:02<01:33, 272.46it/s]  
3%|██████████  
| 881/26400 [00:03<01:32, 276.12it/s]  
3%|██████████  
| 909/26400 [00:03<01:32, 275.05it/s]  
4%|██████████  
| 937/26400 [00:03<01:34, 269.56it/s]  
4%|██████████  
| 967/26400 [00:03<01:32, 275.93it/s]  
4%|██████████  
| 996/26400 [00:03<01:31, 277.81it/s]  
4%|██████████  
| 1024/26400 [00:03<01:35, 266.77it/s]  
4%|██████████  
| 1053/26400 [00:03<01:32, 272.78it/s]  
4%|██████████  
| 1081/26400 [00:03<01:32, 272.72it/s]  
4%|██████████  
| 1110/26400 [00:03<01:31, 277.08it/s]  
4%|██████████  
| 1142/26400 [00:04<01:27, 287.37it/s]  
4%|██████████  
| 1171/26400 [00:04<01:30, 277.64it/s]  
5%|██████████  
| 1207/26400 [00:04<01:25, 293.91it/s]  
5%|██████████  
| 1237/26400 [00:04<01:27, 286.66it/s]  
5%|██████████  
| 1268/26400 [00:04<01:26, 290.19it/s]  
5%|██████████  
| 1298/26400 [00:04<01:33, 269.66it/s]  
5%|██████████  
| 1326/26400 [00:04<01:32, 271.31it/s]  
5%|██████████  
| 1357/26400 [00:04<01:29, 279.78it/s]  
5%|██████████  
| 1394/26400 [00:04<01:23, 299.14it/s]  
5%|██████████  
| 1425/26400 [00:04<01:26, 289.05it/s]  
6%|██████████  
| 1455/26400 [00:05<01:28, 281.79it/s]  
6%|██████████  
| 1484/26400 [00:05<01:27, 283.57it/s]  
6%|██████████  
| 1513/26400 [00:05<01:28, 282.39it/s]  
6%|██████████  
| 1542/26400 [00:05<01:27, 284.02it/s]  
6%|██████████  
| 1571/26400 [00:05<01:27, 282.68it/s]  
6%|██████████  
| 1601/26400 [00:05<01:26, 286.24it/s]  
6%|██████████  
| 1630/26400 [00:05<01:26, 286.73it/s]  
6%|██████████  
| 1659/26400 [00:05<01:29, 277.23it/s]  
6%|██████████  
| 1687/26400 [00:05<01:30, 274.20it/s]

```

6%|██████|
| 1715/26400 [00:06<01:30, 273.71it/s]
7%|██████|
| 1746/26400 [00:06<01:27, 281.56it/s]
7%|██████|
| 1783/26400 [00:06<01:21, 302.03it/s]
7%|██████|
| 1814/26400 [00:06<01:22, 297.62it/s]
7%|██████|
| 1850/26400 [00:06<01:18, 312.52it/s]
7%|██████|
| 1882/26400 [00:06<01:18, 311.31it/s]
7%|██████|
| 1914/26400 [00:06<01:19, 308.67it/s]
7%|██████|
| 1946/26400 [00:06<01:21, 299.11it/s]
7%|██████|
| 1977/26400 [00:06<01:23, 293.14it/s]
8%|██████|
| 2007/26400 [00:06<01:24, 289.39it/s]
8%|██████|
| 2037/26400 [00:07<01:23, 290.22it/s]
8%|██████|
| 2067/26400 [00:07<01:31, 266.78it/s]
8%|██████|
| 2095/26400 [00:07<01:33, 259.56it/s]
8%|██████|
| 2125/26400 [00:07<01:30, 268.52it/s]
8%|██████|
| 2154/26400 [00:07<01:30, 267.25it/s]
8%|██████|
| 2181/26400 [00:07<01:30, 266.70it/s]
8%|██████|
| 2212/26400 [00:07<01:27, 277.07it/s]
8%|██████|
| 2242/26400 [00:07<01:25, 282.99it/s]
9%|██████|
| 2275/26400 [00:07<01:21, 295.04it/s]
9%|██████|
| 2305/26400 [00:08<01:22, 291.57it/s]
9%|██████|
| 2335/26400 [00:08<01:22, 292.56it/s]
9%|██████|
| 2365/26400 [00:08<01:24, 283.32it/s]
9%|██████|
| 2394/26400 [00:08<01:25, 279.75it/s]
9%|██████|
| 2424/26400 [00:08<01:24, 284.13it/s]
9%|██████|
| 2454/26400 [00:08<01:23, 288.10it/s]
9%|██████|
| 2488/26400 [00:08<01:19, 300.54it/s]
10%|██████|
| 2519/26400 [00:08<01:20, 294.92it/s]
10%|██████|
| 2549/26400 [00:08<01:27, 272.47it/s]
10%|██████|

```

| 2579/26400 [00:09<01:25, 279.59it/s]  
10%|██████████  
| 2609/26400 [00:09<01:23, 283.24it/s]  
10%|██████████  
| 2638/26400 [00:09<01:23, 283.79it/s]  
10%|██████████  
| 2667/26400 [00:09<01:24, 281.69it/s]  
10%|██████████  
| 2697/26400 [00:09<01:23, 285.53it/s]  
10%|██████████  
| 2726/26400 [00:09<01:23, 283.71it/s]  
10%|██████████  
| 2755/26400 [00:09<01:23, 282.47it/s]  
11%|██████████  
| 2787/26400 [00:09<01:22, 285.18it/s]  
11%|██████████  
| 2817/26400 [00:09<01:21, 288.01it/s]  
11%|██████████  
| 2851/26400 [00:09<01:18, 299.70it/s]  
11%|██████████  
| 2882/26400 [00:10<01:27, 269.12it/s]  
11%|██████████  
| 2914/26400 [00:10<01:23, 281.31it/s]  
11%|██████████  
| 2943/26400 [00:10<01:24, 276.00it/s]  
11%|██████████  
| 2977/26400 [00:10<01:20, 291.22it/s]  
11%|██████████  
| 3009/26400 [00:10<01:19, 294.57it/s]  
12%|██████████  
| 3041/26400 [00:10<01:17, 300.29it/s]  
12%|██████████  
| 3072/26400 [00:10<01:18, 299.00it/s]  
12%|██████████  
| 3103/26400 [00:10<01:19, 291.40it/s]  
12%|██████████  
| 3133/26400 [00:10<01:20, 287.41it/s]  
12%|██████████  
| 3162/26400 [00:11<01:25, 272.21it/s]  
12%|██████████  
| 3194/26400 [00:11<01:21, 284.42it/s]  
12%|██████████  
| 3223/26400 [00:11<01:21, 283.79it/s]  
12%|██████████  
| 3252/26400 [00:11<01:21, 285.01it/s]  
12%|██████████  
| 3281/26400 [00:11<01:23, 278.47it/s]  
13%|██████████  
| 3309/26400 [00:11<01:24, 273.42it/s]  
13%|██████████  
| 3340/26400 [00:11<01:22, 280.61it/s]  
13%|██████████  
| 3370/26400 [00:11<01:20, 285.56it/s]  
13%|██████████  
| 3399/26400 [00:11<01:21, 283.74it/s]  
13%|██████████  
| 3428/26400 [00:12<01:23, 275.27it/s]

13%|██████████|  
| 3458/26400 [00:12<01:21, 281.67it/s]  
13%|██████████|  
| 3488/26400 [00:12<01:20, 284.70it/s]  
13%|██████████|  
| 3517/26400 [00:12<01:21, 280.69it/s]  
13%|██████████|  
| 3548/26400 [00:12<01:19, 287.49it/s]  
14%|██████████|  
| 3577/26400 [00:12<01:19, 286.76it/s]  
14%|██████████|  
| 3608/26400 [00:12<01:18, 291.11it/s]  
14%|██████████|  
| 3639/26400 [00:12<01:18, 291.74it/s]  
14%|██████████|  
| 3672/26400 [00:12<01:17, 293.62it/s]  
14%|██████████|  
| 3702/26400 [00:12<01:20, 280.83it/s]  
14%|██████████|  
| 3731/26400 [00:13<01:22, 274.11it/s]  
14%|██████████|  
| 3762/26400 [00:13<01:20, 282.63it/s]  
14%|██████████|  
| 3795/26400 [00:13<01:17, 292.43it/s]  
14%|██████████|  
| 3825/26400 [00:13<01:18, 286.47it/s]  
15%|██████████|  
| 3854/26400 [00:13<01:21, 277.85it/s]  
15%|██████████|  
| 3891/26400 [00:13<01:15, 299.03it/s]  
15%|██████████|  
| 3922/26400 [00:13<01:21, 277.37it/s]  
15%|██████████|  
| 3951/26400 [00:13<01:20, 280.45it/s]  
15%|██████████|  
| 3980/26400 [00:13<01:19, 282.64it/s]  
15%|██████████|  
| 4009/26400 [00:14<01:21, 273.76it/s]  
15%|██████████|  
| 4037/26400 [00:14<01:22, 270.25it/s]  
15%|██████████|  
| 4065/26400 [00:14<01:24, 265.56it/s]  
16%|██████████|  
| 4094/26400 [00:14<01:22, 271.89it/s]  
16%|██████████|  
| 4127/26400 [00:14<01:18, 282.11it/s]  
16%|██████████|  
| 4160/26400 [00:14<01:15, 294.37it/s]  
16%|██████████|  
| 4192/26400 [00:14<01:15, 295.18it/s]  
16%|██████████|  
| 4222/26400 [00:14<01:16, 289.98it/s]  
16%|██████████|  
| 4252/26400 [00:14<01:17, 286.44it/s]  
16%|██████████|  
| 4281/26400 [00:15<01:20, 275.46it/s]  
16%|██████████|

| 4311/26400 [00:15<01:18, 281.02it/s]  
16%|██████████  
| 4346/26400 [00:15<01:14, 296.61it/s]  
17%|██████████  
| 4376/26400 [00:15<01:15, 291.79it/s]  
17%|██████████  
| 4406/26400 [00:15<01:14, 293.57it/s]  
17%|██████████  
| 4436/26400 [00:15<01:15, 290.56it/s]  
17%|██████████  
| 4471/26400 [00:15<01:11, 304.78it/s]  
17%|██████████  
| 4502/26400 [00:15<01:14, 294.38it/s]  
17%|██████████  
| 4532/26400 [00:15<01:16, 286.95it/s]  
17%|██████████  
| 4561/26400 [00:15<01:16, 285.54it/s]  
17%|██████████  
| 4592/26400 [00:16<01:14, 291.04it/s]  
18%|██████████  
| 4623/26400 [00:16<01:14, 294.18it/s]  
18%|██████████  
| 4653/26400 [00:16<01:17, 281.98it/s]  
18%|██████████  
| 4684/26400 [00:16<01:15, 287.64it/s]  
18%|██████████  
| 4714/26400 [00:16<01:15, 288.11it/s]  
18%|██████████  
| 4744/26400 [00:16<01:14, 290.12it/s]  
18%|██████████  
| 4779/26400 [00:16<01:11, 303.65it/s]  
18%|██████████  
| 4811/26400 [00:16<01:10, 307.72it/s]  
18%|██████████  
| 4842/26400 [00:16<01:10, 305.92it/s]  
18%|██████████  
| 4873/26400 [00:16<01:11, 302.88it/s]  
19%|██████████  
| 4904/26400 [00:17<01:11, 299.92it/s]  
19%|██████████  
| 4936/26400 [00:17<01:11, 299.07it/s]  
19%|██████████  
| 4966/26400 [00:17<01:13, 290.05it/s]  
19%|██████████  
| 4996/26400 [00:17<01:16, 278.53it/s]  
19%|██████████  
| 5027/26400 [00:17<01:14, 285.12it/s]  
19%|██████████  
| 5061/26400 [00:17<01:11, 299.05it/s]  
19%|██████████  
| 5092/26400 [00:17<01:14, 285.02it/s]  
19%|██████████  
| 5121/26400 [00:17<01:17, 275.32it/s]  
20%|██████████  
| 5156/26400 [00:17<01:12, 293.61it/s]  
20%|██████████  
| 5190/26400 [00:18<01:09, 303.09it/s]

20%|██████████  
| 5222/26400 [00:18<01:09, 306.44it/s]  
20%|██████████  
| 5253/26400 [00:18<01:11, 293.78it/s]  
20%|██████████  
| 5283/26400 [00:18<01:12, 290.71it/s]  
20%|██████████  
| 5313/26400 [00:18<01:13, 286.94it/s]  
20%|██████████  
| 5342/26400 [00:18<01:13, 286.38it/s]  
20%|██████████  
| 5376/26400 [00:18<01:10, 299.23it/s]  
20%|██████████  
| 5407/26400 [00:18<01:10, 295.71it/s]  
21%|██████████  
| 5441/26400 [00:18<01:08, 306.30it/s]  
21%|██████████  
| 5474/26400 [00:19<01:06, 312.39it/s]  
21%|██████████  
| 5513/26400 [00:19<01:03, 330.77it/s]  
21%|██████████  
| 5547/26400 [00:19<01:03, 329.87it/s]  
21%|██████████  
| 5582/26400 [00:19<01:02, 334.00it/s]  
21%|██████████  
| 5616/26400 [00:19<01:08, 304.53it/s]  
21%|██████████  
| 5648/26400 [00:19<01:09, 300.56it/s]  
22%|██████████  
| 5687/26400 [00:19<01:04, 320.59it/s]  
22%|██████████  
| 5720/26400 [00:19<01:05, 314.39it/s]  
22%|██████████  
| 5752/26400 [00:19<01:07, 304.59it/s]  
22%|██████████  
| 5783/26400 [00:20<01:08, 301.09it/s]  
22%|██████████  
| 5814/26400 [00:20<01:08, 298.68it/s]  
22%|██████████  
| 5847/26400 [00:20<01:07, 305.97it/s]  
22%|██████████  
| 5878/26400 [00:20<01:07, 304.70it/s]  
22%|██████████  
| 5909/26400 [00:20<01:07, 302.03it/s]  
22%|██████████  
| 5940/26400 [00:20<01:07, 301.09it/s]  
23%|██████████  
| 5971/26400 [00:20<01:08, 297.83it/s]  
23%|██████████  
| 6003/26400 [00:20<01:07, 300.96it/s]  
23%|██████████  
| 6035/26400 [00:20<01:06, 304.91it/s]  
23%|██████████  
| 6067/26400 [00:20<01:06, 306.86it/s]  
23%|██████████  
| 6098/26400 [00:21<01:06, 304.42it/s]  
23%|██████████



26% | ██████████

27%|████████████████████|  
| 7013/26400 [00:24<01:07, 288.85it/s]  
27%|████████████████████|  
| 7043/26400 [00:24<01:06, 289.79it/s]  
27%|████████████████████|  
| 7073/26400 [00:24<01:07, 287.14it/s]  
27%|████████████████████|  
| 7104/26400 [00:24<01:06, 292.20it/s]  
27%|████████████████████|  
| 7134/26400 [00:24<01:06, 291.31it/s]  
27%|████████████████████|  
| 7168/26400 [00:24<01:03, 300.57it/s]  
27%|████████████████████|  
| 7199/26400 [00:24<01:06, 290.80it/s]  
27%|████████████████████|  
| 7230/26400 [00:24<01:04, 295.68it/s]  
28%|████████████████████|  
| 7260/26400 [00:25<01:05, 290.32it/s]  
28%|████████████████████|  
| 7290/26400 [00:25<01:05, 291.68it/s]  
28%|████████████████████|  
| 7320/26400 [00:25<01:07, 281.93it/s]  
28%|████████████████████|  
| 7349/26400 [00:25<01:08, 276.42it/s]  
28%|████████████████████|  
| 7381/26400 [00:25<01:06, 287.63it/s]  
28%|████████████████████|  
| 7414/26400 [00:25<01:04, 295.36it/s]  
28%|████████████████████|  
| 7447/26400 [00:25<01:02, 304.35it/s]  
28%|████████████████████|  
| 7478/26400 [00:25<01:05, 290.79it/s]  
28%|████████████████████|  
| 7508/26400 [00:25<01:06, 285.36it/s]  
29%|████████████████████|  
| 7539/26400 [00:26<01:04, 290.91it/s]  
29%|████████████████████|  
| 7569/26400 [00:26<01:05, 289.55it/s]  
29%|████████████████████|  
| 7601/26400 [00:26<01:03, 297.47it/s]  
29%|████████████████████|  
| 7631/26400 [00:26<01:05, 284.87it/s]  
29%|████████████████████|  
| 7660/26400 [00:26<01:06, 280.02it/s]  
29%|████████████████████|  
| 7692/26400 [00:26<01:05, 286.45it/s]  
29%|████████████████████|  
| 7721/26400 [00:26<01:06, 279.44it/s]  
29%|████████████████████|  
| 7750/26400 [00:26<01:06, 281.93it/s]  
29%|████████████████████|  
| 7779/26400 [00:26<01:06, 278.00it/s]  
30%|████████████████████|  
| 7807/26400 [00:27<01:07, 274.72it/s]  
30%|████████████████████|  
| 7840/26400 [00:27<01:04, 287.95it/s]  
30%|████████████████████|

8735/26400 [00:30&lt;00:58, 3



```
9608/26400 [00:33<01:01, 272.19it/s]
36%|███████████|
| 9636/26400 [00:33<01:02, 269.95it/s]
37%|███████████|
| 9669/26400 [00:33<00:58, 284.26it/s]
37%|███████████|
| 9698/26400 [00:33<00:58, 284.51it/s]
37%|███████████|
| 9731/26400 [00:33<00:56, 295.40it/s]
37%|███████████|
| 9761/26400 [00:33<00:56, 294.39it/s]
37%|███████████|
| 9791/26400 [00:33<00:56, 294.54it/s]
37%|███████████|
| 9821/26400 [00:34<00:56, 294.65it/s]
37%|███████████|
| 9853/26400 [00:34<00:55, 297.84it/s]
37%|███████████|
| 9886/26400 [00:34<00:53, 306.19it/s]
38%|███████████|
| 9917/26400 [00:34<00:54, 304.85it/s]
38%|███████████|
| 9948/26400 [00:34<00:55, 297.81it/s]
38%|███████████|
| 9978/26400 [00:34<00:59, 277.23it/s]
38%|███████████|
| 10009/26400 [00:34<00:57, 285.73it/s]
38%|███████████|
| 10040/26400 [00:34<00:57, 284.77it/s]
38%|███████████|
| 10069/26400 [00:34<00:57, 284.86it/s]
38%|███████████|
| 10100/26400 [00:35<00:56, 290.55it/s]
38%|███████████|
| 10132/26400 [00:35<00:54, 297.35it/s]
38%|███████████|
| 10162/26400 [00:35<00:54, 297.50it/s]
39%|███████████|
| 10192/26400 [00:35<00:55, 293.25it/s]
39%|███████████|
| 10222/26400 [00:35<00:57, 281.38it/s]
39%|███████████|
| 10255/26400 [00:35<00:55, 291.49it/s]
39%|███████████|
| 10285/26400 [00:35<00:56, 287.47it/s]
39%|███████████|
| 10320/26400 [00:35<00:53, 301.61it/s]
39%|███████████|
| 10351/26400 [00:35<00:54, 293.13it/s]
39%|███████████|
| 10383/26400 [00:35<00:53, 296.75it/s]
39%|███████████|
| 10414/26400 [00:36<00:53, 299.97it/s]
40%|███████████|
| 10445/26400 [00:36<00:53, 298.78it/s]
40%|███████████|
| 10475/26400 [00:36<00:54, 292.40it/s]
```

|                                       |            |
|---------------------------------------|------------|
| 40%                                   | ██████████ |
| 10505/26400 [00:36<00:55, 285.63it/s] |            |
| 40%                                   | ██████████ |
| 10534/26400 [00:36<00:56, 281.32it/s] |            |
| 40%                                   | ██████████ |
| 10563/26400 [00:36<00:55, 283.26it/s] |            |
| 40%                                   | ██████████ |
| 10592/26400 [00:36<00:57, 274.17it/s] |            |
| 40%                                   | ██████████ |
| 10621/26400 [00:36<00:56, 278.15it/s] |            |
| 40%                                   | ██████████ |
| 10652/26400 [00:36<00:55, 285.62it/s] |            |
| 40%                                   | ██████████ |
| 10681/26400 [00:37<00:54, 286.30it/s] |            |
| 41%                                   | ██████████ |
| 10710/26400 [00:37<00:54, 285.92it/s] |            |
| 41%                                   | ██████████ |
| 10743/26400 [00:37<00:52, 296.48it/s] |            |
| 41%                                   | ██████████ |
| 10773/26400 [00:37<00:56, 274.91it/s] |            |
| 41%                                   | ██████████ |
| 10801/26400 [00:37<00:56, 274.20it/s] |            |
| 41%                                   | ██████████ |
| 10829/26400 [00:37<00:58, 268.22it/s] |            |
| 41%                                   | ██████████ |
| 10860/26400 [00:37<00:56, 277.48it/s] |            |
| 41%                                   | ██████████ |
| 10888/26400 [00:37<00:57, 271.19it/s] |            |
| 41%                                   | ██████████ |
| 10916/26400 [00:37<00:56, 272.39it/s] |            |
| 41%                                   | ██████████ |
| 10949/26400 [00:37<00:53, 286.15it/s] |            |
| 42%                                   | ██████████ |
| 10979/26400 [00:38<00:54, 283.82it/s] |            |
| 42%                                   | ██████████ |
| 11010/26400 [00:38<00:52, 290.60it/s] |            |
| 42%                                   | ██████████ |
| 11041/26400 [00:38<00:52, 294.70it/s] |            |
| 42%                                   | ██████████ |
| 11073/26400 [00:38<00:51, 297.88it/s] |            |
| 42%                                   | ██████████ |
| 11105/26400 [00:38<00:50, 301.84it/s] |            |
| 42%                                   | ██████████ |
| 11138/26400 [00:38<00:49, 306.56it/s] |            |
| 42%                                   | ██████████ |
| 11169/26400 [00:38<00:50, 300.68it/s] |            |
| 42%                                   | ██████████ |
| 11200/26400 [00:38<00:52, 290.87it/s] |            |
| 43%                                   | ██████████ |
| 11233/26400 [00:38<00:50, 300.99it/s] |            |
| 43%                                   | ██████████ |
| 11264/26400 [00:39<00:53, 283.90it/s] |            |
| 43%                                   | ██████████ |
| 11298/26400 [00:39<00:50, 298.11it/s] |            |
| 43%                                   | ██████████ |
| 11329/26400 [00:39<00:50, 300.06it/s] |            |
| 43%                                   | ██████████ |

[illegible]

[illegible]



153/180



155/180















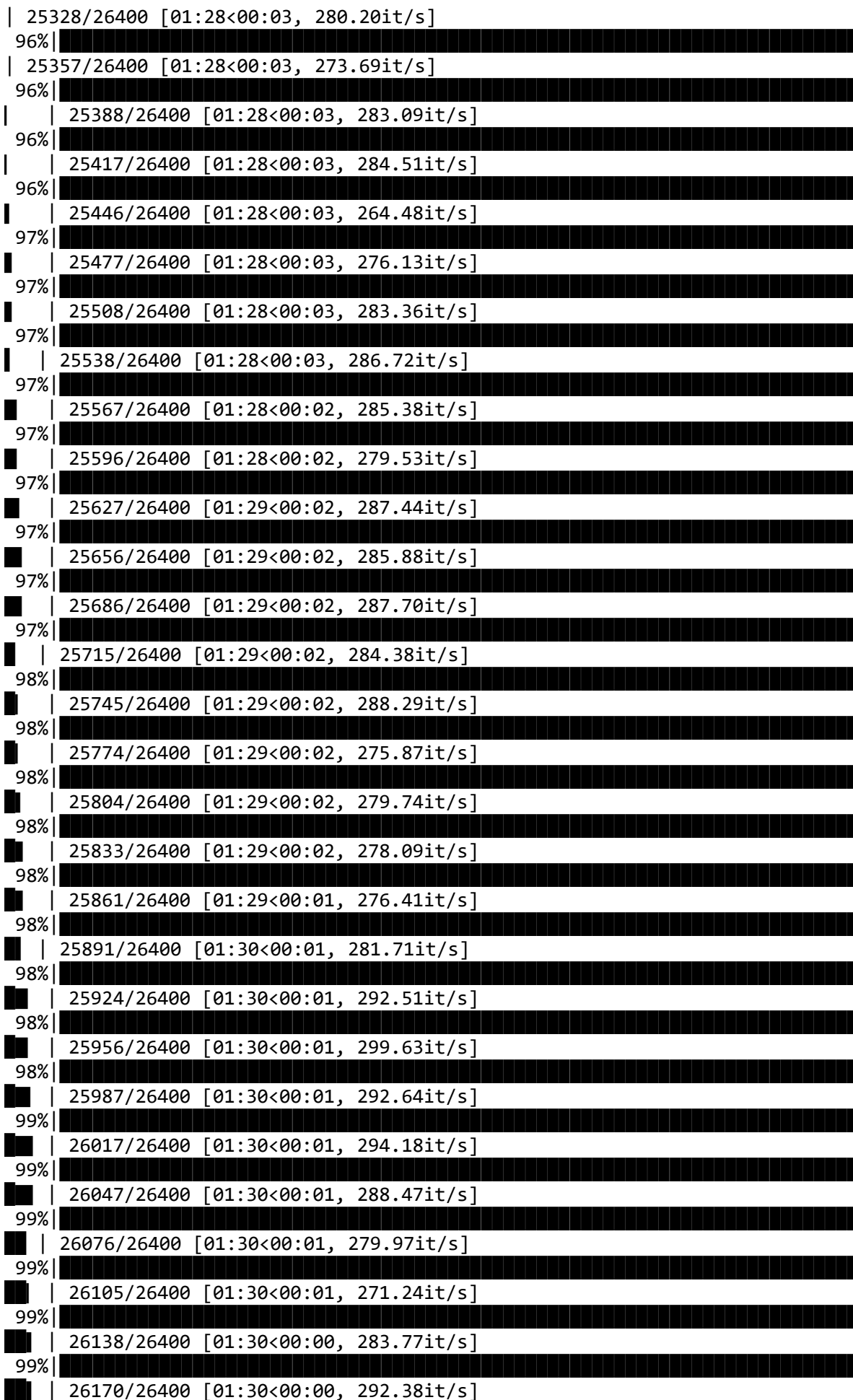
















```
In [120]: # average Word2Vec
compute average word2vec for each review.
tfidf_w2v_vectors_pj_title_train = []; # the avg-w2v for each sentence/review
is stored in this list
for sentence in tqdm(X_train['project_title']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
 for word in sentence.split(): # for each word in a review/sentence
 if (word in glove_words) and (word in tfidf_words):
 vec = model[word] # getting the vector for each word
 # here we are multiplying idf value(dictionary[word]) and the tf value
 # ((sentence.count(word)/len(sentence.split())))
 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split
 ())) # getting the tfidf value for each word
 vector += (vec * tf_idf) # calculating tfidf weighted w2v
 tf_idf_weight += tf_idf
 if tf_idf_weight != 0:
 vector /= tf_idf_weight
 tfidf_w2v_vectors_pj_title_train.append(vector)

print(len(tfidf_w2v_vectors_pj_title_train))
print(len(tfidf_w2v_vectors_pj_title_train[0]))
```

```
0%|
| 0/35912 [00:00<?, ?it/s]
28%|██████████
| 10034/35912 [00:00<00:00, 99611.70it/s]
55%|██████████
| 19722/35912 [00:00<00:00, 98555.89it/s]
100%|██████████
█| 35912/35912 [00:00<00:00, 97318.98it/s]
```

35912

300



```

In [122]: # average Word2Vec
compute average word2vec for each review.
tfidf_w2v_vectors_pj_title_test = []; # the avg-w2v for each sentence/review i
s stored in this list
for sentence in tqdm(X_test['project_title']): # for each review/sentence
 vector = np.zeros(300) # as word vectors are of zero length
 tf_idf_weight = 0; # num of words with a valid vector in the sentence/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_vectors_pj_title_test.append(vector)

print(len(tfidf_w2v_vectors_pj_title_test))
print(len(tfidf_w2v_vectors_pj_title_test[0]))

```

```

0%|
| 0/26400 [00:00<?, ?it/s]
32%|██
| 8398/26400 [00:00<00:00, 83384.07it/s]
100%|██
█| 26400/26400 [00:00<00:00, 93872.79it/s]

```

```

26400
300

```

```

In [123]: # Concatinating all the features
X_tr = hstack((tfidf_w2v_vectors_train, X_train_state_ohe, X_train_teacher_ohe
,
 X_train_grade_ohe, X_train_category_ohe,
 X_train_subcategory_ohe, X_train_price_norm,
 X_train_teach_prev_norm, tfidf_w2v_vectors_pj_title_train)).toc
sr()

X_cr = hstack((tfidf_w2v_vectors_cv, X_cv_state_ohe, X_cv_teacher_ohe,
 X_cv_grade_ohe,X_cv_category_ohe,
 X_cv_subcategory_ohe, X_cv_price_norm,
 X_cv_teach_prev_norm, tfidf_w2v_vectors_pj_title_cv)).tocsr()

X_te = hstack((tfidf_w2v_vectors_test, X_test_state_ohe, X_test_teacher_ohe,
 X_test_grade_ohe, X_test_category_ohe,
 X_test_subcategory_ohe, X_test_price_norm,
 X_test_teach_prev_norm, tfidf_w2v_vectors_pj_title_test)).tocsr
()

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

```

Final Data matrix

(35912, 701) (35912,)

(17688, 701) (17688,)

(26400, 701) (26400,)

=====

```

In [124]: # Hyper parameter Tuning
train_auc = []
cv_auc = []
K = [3, 15, 25, 51, 101]
for i in tqdm(K):
 neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
 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 e
 # stimates 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))

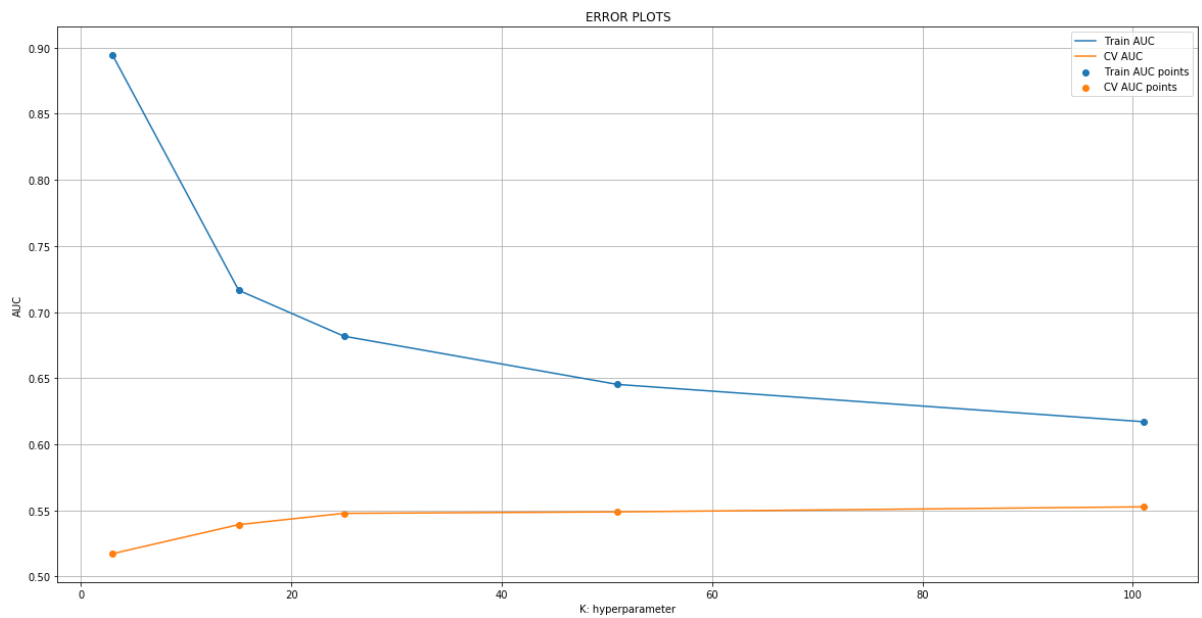
0%|
| 0/5 [00:00<?, ?it/s]
20%|
| 1/5 [07:01<28:04, 421.02s/it]
40%|
| 2/5 [14:10<21:10, 423.59s/it]
60%|
| 3/5 [21:23<14:12, 426.24s/it]
80%|
| 4/5 [28:34<07:07, 427.76s/it]
100%|
| 5/5 [35:44<00:00, 428.98s/it]

```

```
In [125]: plt.figure(figsize=(20,10))
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()
```



```

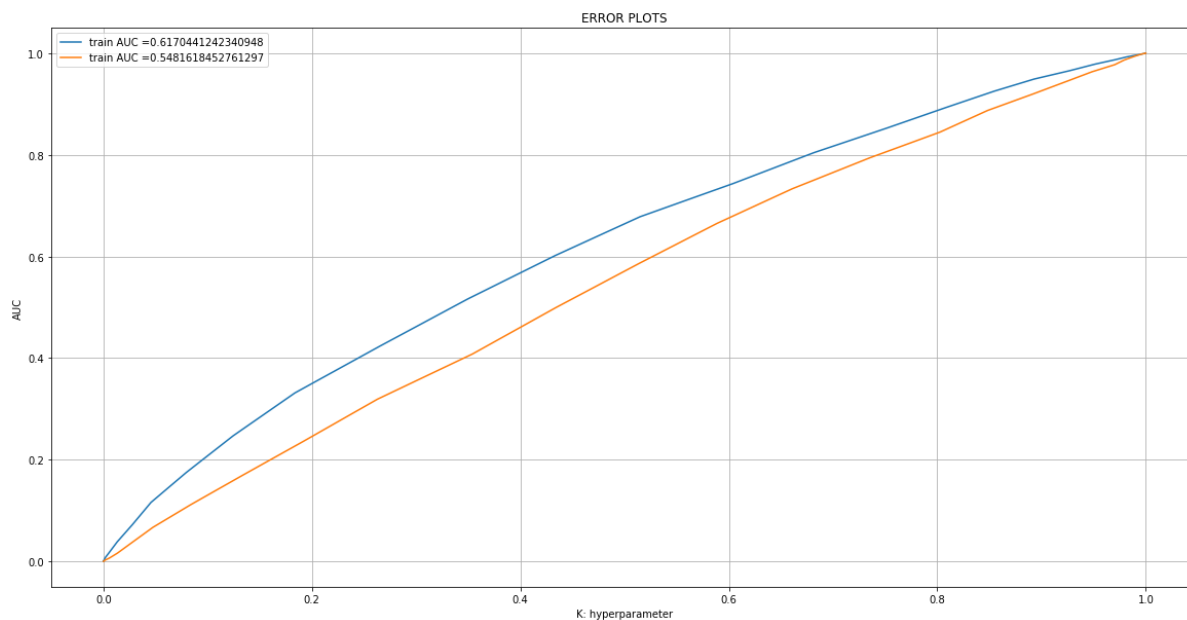
In [126]: best_k = 101
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k, n_jobs=-1)
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.figure(figsize=(20,10))
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



The Blue plot in the Graph is Train AUC and Orange Plot is Test AUC. Initially the label was accidentally copy pasted as Train AUC for both. Since redrawing plot would require to re run entire ipynb and it takes up a lot of time. so documenting the information as a markdown.

## 2.5 Feature selection with `SelectKBest`

```
In [127]: # 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

Performing the Concatination of Features for Set-2

X_tr = hstack((X_train_essay_tfidf, X_train_state_ohe, X_train_teacher_ohe,
 X_train_grade_ohe, X_train_price_norm, X_train_category_ohe,
 X_train_subcategory_ohe, X_train_teach_prev_norm,
 X_train_pj_title_tfidf)).tocsr()

X_cr = hstack((X_cv_essay_tfidf, X_cv_state_ohe, X_cv_teacher_ohe,
 X_cv_grade_ohe, X_cv_category_ohe, X_cv_subcategory_ohe,
 X_cv_price_norm, X_cv_teach_prev_norm, X_cv_pj_title_tfidf)).to
csr()

X_te = hstack((X_test_essay_tfidf, X_test_state_ohe, X_test_teacher_ohe,
 X_test_grade_ohe, X_test_category_ohe, X_test_subcategory_ohe,
 X_test_price_norm, X_test_teach_prev_norm,
 X_test_pj_title_tfidf)).tocsr()
```

```
In [128]: # importing necessary packages for Univariate feature selection
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2
```

```
In [129]: print ("shape before selecting top features")
print (X_tr.shape)
print (X_cr.shape)
print (X_te.shape)
```

```
shape before selecting top features
(35912, 10101)
(17688, 10101)
(26400, 10101)
```

```
In [132]: # Selecting top 2000 features
X_tr_new = SelectKBest(chi2, k=2000).fit_transform(X_tr, y_train)
print (X_tr_new.shape)

(35912, 2000)
```



```
In [134]: X_cr_new = SelectKBest(chi2, k=2000).fit_transform(X_cr, y_cv)
print (X_cr_new.shape)
X_test_new = SelectKBest(chi2, k=2000).fit_transform(X_te, y_test)
print (X_test_new.shape)
```

```
(17688, 2000)
(26400, 2000)
```

```
In [135]: # Hyperparameter Tuning and Error plots with Best 2000 features
train_auc = []
cv_auc = []
K = [3, 15, 25, 51, 101]
for i in tqdm(K):
 neigh = KNeighborsClassifier(n_neighbors=i, n_jobs=-1)
 neigh.fit(X_tr_new, y_train)

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

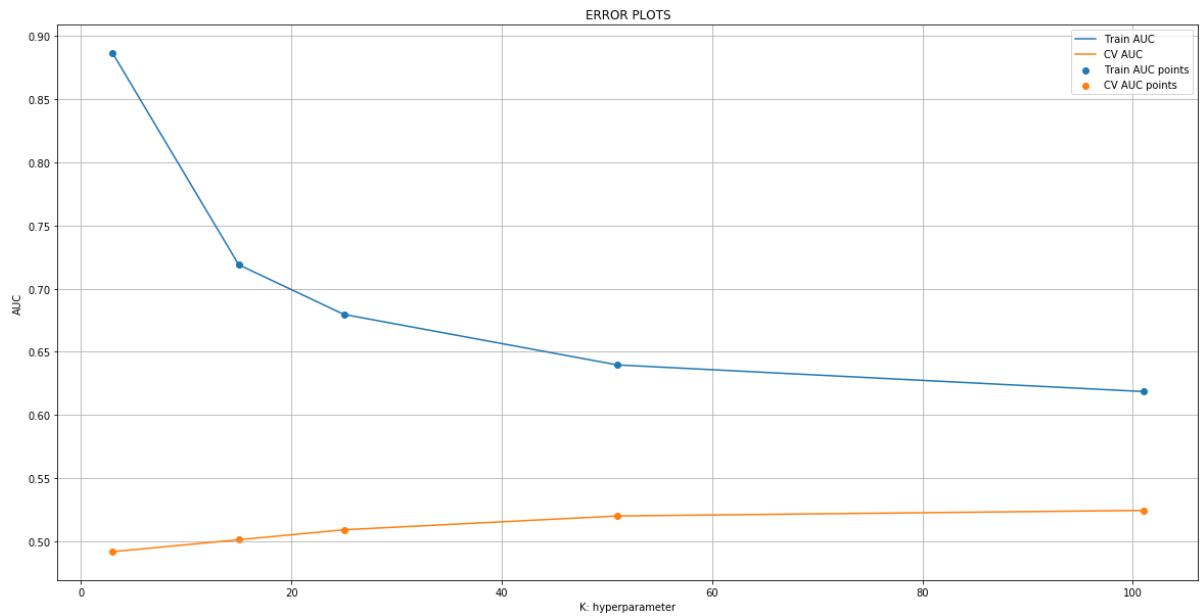
 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability e
 # estimates of the positive class
 # not the predicted outputs
 train_auc.append(roc_auc_score(y_train, y_train_pred))
 cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
```

```
0%|
| 0/5 [00:00<?, ?it/s]
20%|██████████
| 1/5 [00:50<03:22, 50.67s/it]
40%|██████████
| 2/5 [01:45<02:35, 51.86s/it]
60%|██████████
| 3/5 [02:39<01:45, 52.59s/it]
80%|██████████
| 4/5 [03:33<00:53, 53.09s/it]
100%|██████████
██████████ | 5/5 [04:29<00:00, 53.83s/it]
```

```
In [136]: plt.figure(figsize=(20,10))
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()
```



```

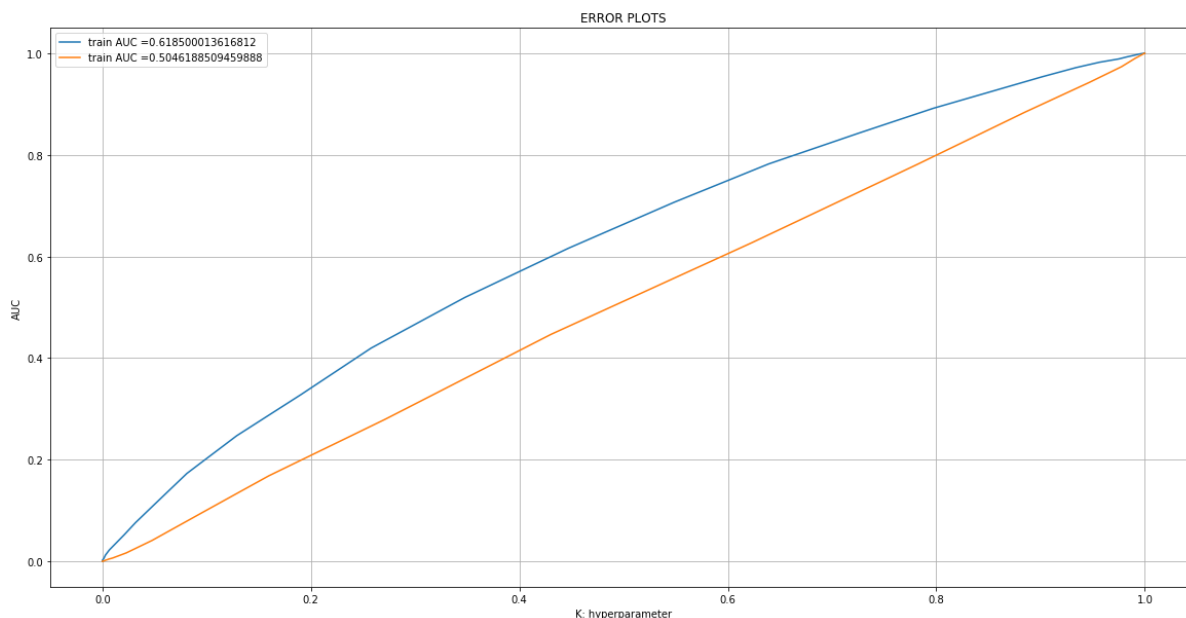
In [138]: best_k = 101
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k, n_jobs=-1)
neigh.fit(X_tr_new, 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_new)
y_test_pred = batch_predict(neigh, X_test_new)

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.figure(figsize=(20,10))
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



The Blue plot in the Graph is Train AUC and Orange Plot is Test AUC. Initially the label was accidentally copy pasted as Train AUC for both. Since redrawing plot would require to re run entire ipynb and it takes up a lot of time. so documenting the information as a markdown.

### 3. Conclusions

```
In [141]: # Please compare all your models using Prettytable library
from prettytable import PrettyTable

#If you get a ModuleNotFoundError error , install prettytable using: pip3 inst
all prettytable

x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Hyperparameter", "Train AUC", "Test A
UC"]

x.add_row(["BOW", "Brute", 105, 0.647, 0.587])
x.add_row(["TFIDF", "Brute", 105, 0.6322, 0.546])
x.add_row(["W2V", "Brute", 101, 0.620, 0.548])
x.add_row(["TFIDF W2V", "Brute", 101, 0.617, 0.548])
print(x)
```

| Vectorizer | Model | Hyperparameter | Train AUC | Test AUC |
|------------|-------|----------------|-----------|----------|
| BOW        | Brute | 105            | 0.647     | 0.587    |
| TFIDF      | Brute | 105            | 0.6322    | 0.546    |
| W2V        | Brute | 101            | 0.62      | 0.548    |
| TFIDF W2V  | Brute | 101            | 0.617     | 0.548    |

**Summary: There was not a lot of difference in the AUC scores with all the features selected and top 2000 feature selection. With SelectionKBest the Train AUC was around 0.61 and Test AUC was around 0.50**