

# DonorsChoose

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

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

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

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

# About the DonorsChoose Data Set

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

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. <b>Example:</b> 123456789
<code>project_title</code>		Title of the project. <b>Example:</b> Art Will Make You Feel Good
<code>project_grade_category</code>		Grade level of students for which the project is targeted. One of the following enumerated list of categories: <ul style="list-style-type: none"><li>• Grades K-2</li><li>• Grades 3-5</li><li>• Grades 6-8</li><li>• Grades 9-12</li></ul>
<code>project_subject_categories</code>		One or more (comma-separated) subject categories for the project from the following enumerated list of categories: <ul style="list-style-type: none"><li>• Applied &amp; Technical Education</li><li>• Art &amp; Design</li><li>• Care &amp; Safety</li><li>• Health &amp; Physical Education</li><li>• History &amp; Social Studies</li><li>• Literacy &amp; Language</li><li>• Math &amp; Science</li><li>• Music &amp; Performing Arts</li><li>• Special Education</li></ul>
<code>project_subject_subcategories</code>		One or more (comma-separated) subject subcategories for the project from the following enumerated list of categories: <ul style="list-style-type: none"><li>• Music &amp; Performing Arts</li><li>• Literacy &amp; Language, Math &amp; Science</li></ul>
<code>school_state</code>		State where school is located ( <a href="https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_abbreviations">Two-letter U.S. postal abbreviation</a> ). <b>Example:</b> CA
<code>project_resource_summary</code>		An explanation of the resources needed for the project. <b>Example:</b> My students need hands on literacy materials to enhance their sensory
<code>project_essay_1</code>		First application essay
<code>project_essay_2</code>		Second application essay
<code>project_essay_3</code>		Third application essay
<code>project_essay_4</code>		Fourth application essay
<code>project_submitted_datetime</code>		Datetime when project application was submitted. <b>Example:</b> 2011-01-01 12:43:21
<code>teacher_id</code>		A unique identifier for the teacher of the proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4f1

6/7/2019	Donors_Choose_KNN_Assignment_3	
Feature		Description
		Teacher's title. One of the following enumerated values: <code>teacher_title</code>
<code>teacher_prefix</code>	• • • • • •	
<code>teacher_number_of_previously_posted_projects</code>		Number of project applications previously submitted by the same teacher. <b>Example:</b> 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. <b>Example:</b> <code>p036502</code>
<code>description</code>	Description of the resource. <b>Example:</b> <code>Tenor Saxophone Reeds, Box of 25</code>
<code>quantity</code>	Quantity of the resource required. <b>Example:</b> <code>3</code>
<code>price</code>	Price of the resource required. <b>Example:</b> <code>9.95</code>

**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.

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 [2]:

```

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

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

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

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

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

from tqdm import tqdm
import os

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

```

C:\Users\RASHU TYAGI\Anaconda3\lib\site-packages\smart\_open\ssh.py:34: Use  
 rWarning: paramiko missing, opening SSH/SCP/SFTP paths will be disabled.  
 `pip install paramiko` to suppress  
 warnings.warn('paramiko missing, opening SSH/SCP/SFTP paths will be disa  
 bled. `pip install paramiko` to suppress')

## 1.1 Reading Data

In [3]:

```

project_data = pd.read_csv(r'D:\Rashu Studies\AppliedAICourse\Assignments\Mandatory Ass
ignments\Mandatory Assignment 3 Donors Choose KNN\train_data.csv')
resource_data = pd.read_csv(r'D:\Rashu Studies\AppliedAICourse\Assignments\Mandatory As
signments\Mandatory Assignment 3 Donors Choose KNN\resources.csv')

```

In [4]:

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

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

-----

The attributes of data : ['Unnamed: 0' 'id' 'teacher\_id' 'teacher\_prefix' 'school\_state' 'project\_submitted\_datetime' 'project\_grade\_category' 'project\_subject\_categories' 'project\_subject\_subcategories' 'project\_title' 'project\_essay\_1' 'project\_essay\_2' 'project\_essay\_3' 'project\_essay\_4' 'project\_resource\_summary' 'teacher\_number\_of\_previously\_posted\_projects' 'project\_is\_approved']

In [5]:

```
# 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[5]:

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

◀ ◻ ▶

In [6]:

```
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[6]:

	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

**NOW THE MOST IMPORTANT THING HERE IS THAT YOU SHOULD SPLIT OUR DATA INTO TRAIN AND TEST BEFORE APPLYING ANY FIT TECHNIQUE LIKE BOW OR TFIDF BECAUSE OTHERWISE THERE WILL BE DATA LEAKAGE PROBLEM. ALSO FOR PREPROCESSING LIKE STANDARDIZATION AND NORMALIZATION ALSO WE SHOULD KEEP IN MIND THAT TRAIN TEST SPLIT SHOULD BE DONE BEFORE APPLYING THOSE PREPROCESSING TECHNIQUES**

In [7]:

```
# REFER THIS SOUNDCLOUD LINK : https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf
```

## Train\_Test Split

In [8]:

```
# train test split
# note that here This stratify parameter makes a split so that the proportion of values
in the sample produced will be the same as the proportion of values provided to paramet
er stratify.
#For example, if variable y is a binary categorical variable with values 0 and 1 and th
ere are 25% of zeros and 75% of ones, stratify=y will make sure that your random split
has 25% of 0's and 75% of 1's.

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(project_data, project_data['project
_is_approved'], test_size=0.33, stratify = project_data['project_is_approved'])

# now getting the crossvalidation data from our train data
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, strat
ify=y_train)
```

In [9]:

```
# Now we will be removing the column "project_is_approved" because that is the only one  
which our model needs to predict
```

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

In [ ]:

In [ ]:

In [ ]:

**Now we will do all kind of preprocessing required for the train data ,test data,crossvalidation data separately**

## **FOR TRAIN DATA**

### **Preprocessing of `project\_subject\_categories`**

In [10]:

```
categories = list(X_train['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 replacing all the ' ' (space) with '' (empty) ex: "Math & Science"=> "Math&Science"
            temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_') # we are replacing the & value into
    cat_list.append(temp.strip())

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

from collections import Counter
my_counter = Counter()
for word in X_train['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]))
```

## Preprocessing of project\_subject\_subcategories



In [11]:

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

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

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

X_train['clean_subcategories'] = sub_cat_list
X_train.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 X_train['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]))
```

## Text preprocessing

In [12]:

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

In [13]:

```
X_train.head(2)
```

Out[13]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
8134	135883	p090948	ade6aad520456e6e23490918ef5ea2d4	Ms.	CA
32738	113417	p140603	e37e94882ace7b10b482892d78787fa4	Ms.	CA



In [14]:

```
# printing some random reviews
print(X_train['essay'].values[0])
print("="*50)
print(X_train['essay'].values[150])
print("="*50)
print(X_train['essay'].values[1000])
print("="*50)
print(X_train['essay'].values[20000])
print("="*50)
print(X_train['essay'].values[9999])
print("="*50)
```

As a teacher in a high-poverty school, I have the unique opportunity to work with a group of thirty-three students from varying backgrounds who all contribute unique skills and experiences to our classroom. Our class is made up of a multilingual and global community of learners, with over half the class speaking a language other than English at home. These bright, innovative, spunky students are enthusiastic readers, writers, mathematicians, scientists and historians. Many students do not have access to writing materials at home and would greatly benefit from having access to a rich writing center from which they can easily build the worlds and ideas that are just waiting to be put down on paper. Writing is a critical part of 5th grade. It is the beginning of the 5 paragraph essay, and my students will be doing a lot of work researching different topics and developing their own persuasive essays. They will also have an opportunity to write a memoir at the end of the school year to reflect on their elementary career. Writing time is the best opportunity that my students have to express themselves and their own beliefs. To do this they need a place to put all of this writing, which is why I am requesting composition books for all of my students. They also need an organized center where they can find all of the tools for their writing - from paper trays to pencils bins. Your contributions will greatly influence my students experience with writing this year. To say thank you, every donor will receive a sample of student writing from our narrative unit!nannan

My students thrive on learning new things and sharing what they know! Our school has over 70% of students who receive free or reduced lunch and we are a Title I school that has a high poverty and low income rate. My students come to school having to think about more than just learning and it is my job to make sure that they feel loved and have an environment where there are no limits to their dreams! I have a wide variety of students in my classroom from the language they speak to their learning abilities and disabilities! My children are independent, dreamers, helpers and learners. In first grade, our main focus is to teach our students how to read. There are so many apps that reinforce the reading skills I teach every day! I currently have two iPads in my classroom. One was provided by our county and one was funded through my last Donors Choose project. It is the highlight of their day if they get to work on the iPad that day. I have many struggling readers in my class this year and they get excited to learn when they have the opportunity to work in the iPad. If I had three iPads in my classroom, students would have double the opportunity to enhance their reading skills. Many of my students do not access to the Internet or any technology at home. Having the opportunity to use technology in the classroom is beneficial to each child on so many levels. We use apps such as reading friendly, epic, abcmouse, abcy and more. I am able to tailor the needs to each specific student on the iPad. Your generous donation to our project will improve our First grade classroom by building stronger readers and learners.nannan

Fifth grade has always been my favorite grade. 5th graders personalities are all very special and incredibly unique. Their energy is contagious! They have a wide variety of hobbies, interests, hopes and dreams. However, as different as my kids all are, they do have ONE thing in common.... They all LOVE graphic novels! My class is completely enthralled by them! The novels fly off of the shelves! Discussions about the books are vibrant and lively. Requests for new titles happen daily. My 5th grades had these thoughts and opinions to share about their experiences reading graphic novels. "Graphic novels are the best books ever!" "I like graphic novels because they are interesting and because they are so adventurous. You never know what will happen next!" "They take you to a whole other fantasy world." "Graphic novels are funny and the pictures help to tell the story." "They are easier to read and they make me feel more confident." "The drawings are amazing. The details

are great.\"\\r\\nBased on these glowing reviews, my curiosity was peaked.  
\\r\\nThe kids have converted me also! \\r\\n\\r\\n\\r\\nnnnnn

My students are brand new students - kindergarten students with such an eagerness to learn. We are located on a military installation and have students from all over the country and/or world. My students are diverse and bright. Some of these students do not have access to technology at home and look forward to the time that is given in class to explore technology and all the fun (and academic) things that we can do with it. Our goal is not only to learn, but to have fun doing so. Some of my students do not have access to technology at home and we would like to enhance student learning by adding tablets to our classroom resources. With these tablets, the students will be able to access educational games that are paired with our curriculum and will expand their learning. Students will be able to practice using technology while listening to a book online, reading an e-book, playing math games, looking at pictures of famous presidents, and so much more. The possibilities are endless by adding technology to our classroom and that is what we strive to do - to aim for the stars and then some!

\"Good Morning Mrs. A\" is the greeting I hear every morning from my students. Along with hugs and high fives, it is important they say hello because it sets the tone for their day. By greeting me they know I am here for them. I am here to care for them, teach them, and guide them along the way.

My school is 95% free and reduced. My class is 95% African American children. The students in my class often don't have their basic needs met before they come to school. Clean socks, breakfast, and lotion are some of the things I take care of before we start our learning day. I thrive on seeing their smiling faces and the eagerness they have to learn. Imagine you are on a plane and passengers around you are watching a video or listening to music without headphones. It causes you to be annoyed and distracted. In my classroom each student has their own computer provided by the district. With my school being very transit and low income, most students cannot provide their own headphones.

Students in my classroom create their own personal learning path with their Chromebooks. With access to their own headphones they will be able to complete personal task efficiently. They will not become confused or frustrated by hearing other students programs. They will be able to concentrate and focus on their own task. This will ultimately give them a higher level of personal success.

In [15]:

```
# creating a function named as decontracted which does the job of decontraction

# 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 [16]:

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

My students are brand new students - kindergarten students with such an eagerness to learn. We are located on a military installation and have students from all over the country and/or world. My students are diverse and bright. Some of these students do not have access to technology at home and look forward to the time that is given in class to explore technology and all the fun (and academic) things that we can do with it. Our goal is not only to learn, but to have fun doing so. Some of my students do not have access to technology at home and we would like to enhance student learning by adding tablets to our classroom resources. With these tablets, the students will be able to access educational games that are paired with our curriculum and will expand their learning. Students will be able to practice using technology while listening to a book online, reading an e-book, playing math games, looking at pictures of famous presidents, and so much more. The possibilities are endless by adding technology to our classroom and that is what we strive to do - to aim for the stars and then some!nannan

=====

In [17]:

```
#\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)
```

My students are brand new students - kindergarten students with such an eagerness to learn. We are located on a military installation and have students from all over the country and/or world. My students are diverse and bright. Some of these students do not have access to technology at home and look forward to the time that is given in class to explore technology and all the fun (and academic) things that we can do with it. Our goal is not only to learn, but to have fun doing so. Some of my students do not have access to technology at home and we would like to enhance student learning by adding tablets to our classroom resources. With these tablets, the students will be able to access educational games that are paired with our curriculum and will expand their learning. Students will be able to practice using technology while listening to a book online, reading an e-book, playing math games, looking at pictures of famous presidents, and so much more. The possibilities are endless by adding technology to our classroom and that is what we strive to do - to aim for the stars and then some!nannan

In [18]:

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

My students are brand new students kindergarten students with such an eagerness to learn We are located on a military installation and have students from all over the country and or world My students are diverse and bright Some of these students do not have access to technology at home and look forward to the time that is given in class to explore technology and all the fun and academic things that we can do with it Our goal is not only to learn but to have fun doing so Some of my students do not have access to technology at home and we would like to enhance student learning by adding tablets to our classroom resources With these tablets the students will be able to access educational games that are paired with our curriculum and will expand their learning Students will be able to practice using technology while listening to a book online reading an e book playing math games looking at pictures of famous presidents and so much more The possibilities are endless by adding technology to our classroom and that is what we strive to do to aim for the stars and then some nannan

In [19]:

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# because although they are in this list but they matter a lot because
# they change the meaning of the entire sentence.
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you'r
e", "you've",\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him',
            'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 't
hey', 'them', 'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "th
at'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'ha
d', 'having', 'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as'
, 'until', 'while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through'
, 'during', 'before', 'after',\
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'ov
er', 'under', 'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'an
y', 'both', 'each', 'few', 'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too'
, 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'no
w', 'd', 'll', 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't",
            'doesn', "doesn't", 'hadn',\
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'migh
tn', "mightn't", 'mustn',\
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'w
asn', "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [20]:

```
# Combining all the above preprocessing techniques for all the project essays
from tqdm import tqdm
preprocessed_essays_Train = []
# tqdm is for printing the status bar
for sentence in tqdm(X_train['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\n', ' ')
    sent = sent.replace('\\t', ' ')
    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_Train.append(sent.lower().strip())
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 49041/49041 [00:21<00:00, 2256.99it/s]
```



In [21]:

```
# after preprocesing of project essays
preprocessed_essays_Train[20000]
```

Out[21]:

```
'students brand new students kindergarten students eagerness learn located
military installation students country world students diverse bright stude
nts not access technology home look forward time given class explore techn
ology fun academic things goal not learn fun students not access technolog
y home would like enhance student learning adding tablets classroom resour
ces tablets students able access educational games paired curriculum expan
d learning students able practice using technology listening book online r
eading e book playing math games looking pictures famous presidents much p
ossibilities endless adding technology classroom strive aim stars nannan'
```

## Preprocessing of project\_title

**Now we will simply apply the above preprocessing steps on the project title for the train data as well, as it is also a text feature**

In [22]:

```
# printing some random titles.
print(X_train['project_title'].values[0])
print("="*50)
print(X_train['project_title'].values[150])
print("="*50)
print(X_train['project_title'].values[1000])
print("="*50)
print(X_train['project_title'].values[20000])
print("="*50)
print(X_train['project_title'].values[9999])
print("="*50)
```

```
Help Build Our Writing Center
=====
Growing minds with an iPad
=====
Gotta Get MORE Graphic Novels!
=====
Taking Off with Technology!
=====
What Did It Say? Can You Hear It Now?
=====
```

**We have already written the preprocessing codes for different preprocessing approaches now we will simply use those codes on the project titles**



In [25]:

```
categories = list(X_test['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 replacing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
            temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
    temp = temp.replace('&','_') # we are replacing the & value into
    cat_list.append(temp.strip())

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

from collections import Counter
my_counter = Counter()
for word in X_test['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]))
```

## Preprocessing of project\_subject\_subcategories

In [26]:

```
sub_categories = list(X_test['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())

X_test['clean_subcategories'] = sub_cat_list
X_test.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 X_test['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]))
```

## Text Preprocessing

In [27]:

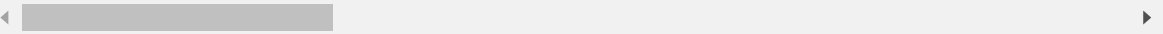
```
# merge two column text dataframe:
X_test["essay"] = X_test["project_essay_1"].map(str) + \
                  X_test["project_essay_2"].map(str) + \
                  X_test["project_essay_3"].map(str) + \
                  X_test["project_essay_4"].map(str)
```

In [28]:

```
X_test.head(5)
```

Out[28]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
36675	158529	p012531	d5418aad0ebd66f05263f9c42773b003	Mrs.	CA
88995	104451	p013643	72cd1d32d22b1d1e98b839e20aff93fb	Mr.	IL
71367	26079	p205185	eef03ef3528ea1ab9fc3b0d7c54c464a	Ms.	MI
53982	110544	p191726	52a0ec26e1010b35fa234518653328fb	Mrs.	AL
108487	128586	p053525	c92897ab11467e64f4ff6dd95897a46e	Ms.	NV



In [29]:

```
# printing some random reviews
print(X_test['essay'].values[0])
print("="*50)
print(X_test['essay'].values[150])
print("="*50)
print(X_test['essay'].values[1000])
print("="*50)
print(X_test['essay'].values[20000])
print("="*50)
print(X_test['essay'].values[9999])
print("="*50)
```

My first grade classroom is home to students who are enthusiastic and excited to learn. Our school is rich in diversity where over 75% of the students are English Language Learners. The students have backgrounds coming from Guatemala, Mexico and West Africa. Our school is also poverty stricken where more than 95% of our students are on free and reduced price lunch.

My goal is to provide my students with engaging, creative and meaningful learning experiences so they can grow into confident, respectful and literate young people. We are a Title 1 school which means funds and resources are extremely limited at school and in my students' homes. Many of my students are not able to afford to buy crayons, markers or even glue in their homes. The opportunities I provide for them in the classroom are that much more meaningful and important to their overall growth.

I am asking for these materials to help increase and foster my students' creativity and imagination. Their families are not able to afford crayons or markers so the use and exploration of art is very limited without these materials. The supplies I am requesting will help my students stay focused on the learning while having fun and engaging experiences.

Coloring improves students' motor skills, stimulates creatively and provides stress relief. My students love to create and learn better through hands on activities. Crayons and markers will be used daily to add brightly colored illustrations to our writing and help us to work on our fine motor skills. We will use the dry erase markers to work on letter and number formation. We will also use the markers to work on writing our phonics sounds.

It is also important for my students to practice organization and responsibility for their materials. I am requesting small containers so the students will have a place to store their crayons. They will be able to keep their containers at their seats allowing them to have their crayons within arms reach.

Thank you for taking the time to consider my project worth your donation. My students will benefit from your generosity!

I teach an amazing group of 3rd graders at a high poverty inner city publi

c school in Indianapolis. \r\n\r\nAll of our students are on the free breakfast and lunch program. Many of our students' parents work odd hours and they barely get to spend time with their children. I have some students being raised by a single parent or a grandparent. No matter what my students are going through at home, they always try their personal best in the classroom. My students know how high my expectations are for each one of them! My students work very hard in the classroom to earn awards. They earn awards for gaining fluency on Reflex Math and they also earn awards for growing on Pivot and Raps180, both district mandated assessments. They have worked so hard and deserve to get awards to show their hard work in beautiful color! I would also like to have color ink to print strategy posters and task cards for my students to use in the classroom. They are always referring to posters to help guide them on work. My students also love working in small groups during reading and math and need some new task cards! They also need laminating pouches to keep their awards and task cards safe!\r\nHaving colorful, vibrant materials for my classroom will help encourage all of my eager students to work on their tasks and enjoy their surroundings.\r\nI will be able to readily create and print materials that will be used for this year and future ones as well.nannan

=====  
My students are three, four and five year olds that attend either my a.m. class or my p.m. pre-kindergarten class. My students are considered at-risk in order to qualify for the pre-kindergarten program at my school. Many of my students receive speech services. My students are bright, energetic, eager to learn children who I know will go far with the right start.\r\n\r\nI teach at a title I school in the suburbs of Chicago. My district provides free breakfast and lunch to all students. My students enjoy listening to stories. Everyday I read at least two books to my students. This component to my listening station will enable my students to listen to even more books daily. My students will learn to independently listen to stories. Students can listen to a story alone or with a partner. They will learn the skills to follow along in the book, using audio and visual cues as to when to turn the pages.\r\n\r\nEarly exposure to reading can create a life long love of books! \r\nThis reading center will help do that. I have numerous books that are age appropriate that will foster independent reading.nannan

=====  
My co-teacher and I have 48 amazing first graders who are curious, eager and excited to be at school. Our classrooms are full of life, love, and laughter. \r\nOur students love to collaborate and work together to solve problems and we are committed to providing a flexible and engaging environment that fits their energy level. As most first graders, our students love to move and are full of energy. We want to preserve the energy and wonder in our students so they will continue to be active learners. Do you remember your first grade classroom? If you were to walk into a first grade classroom today, it is likely you would see the same, or similar environment. \r\nWe want more for our students; we want to provide our students with an active learning environment in which they can fidget when needed and showcase their learning using dry erase tables.\r\nIn an article for The Education Facilities Clearinghouse, author Greg Smith explains, "American schools, designed around a standard learning environment that supports a lecture style of teaching, have remained relatively unchanged for the last 50 years while we have evolved into a society of visual and tactile learners. So what defines these 21st century learning environments? The words we most often used are flexible, agile, and adaptable, words that ultimately mean being able to adjust to new conditions, modify for a new use or purpose, and allow flexibility to engage students in a variety of ways." \r\nWe hope to give our students a learning environment such as the one Smith illustrated in his article. In order to achieve this environment, we need to provide flexible seating options.\r\nThis grant will provide dry erase film for tables and alternative seating for two first grade classes. We en



vision students choosing how and where to learn to meet their individual needs. For example, the wobble stools are designed to allow students who need tactile stimuli the freedom to move, so they can focus on learning. Other options include standing, laying on mats, and comfortable couches.\r\nAs teachers in the 21st century, we have the responsibility to teach and encourage collaboration as we meet students' sensory needs. The dry erase surfaces and mobile furniture will provide unlimited opportunities for students to collaborate with others.nannan

=====

In [30]:

```
# creating a function named as decontracted which does the job of decontraction

# 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 [31]:

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

My students are three, four and five year olds that attend either my a.m. class or my p.m. pre-kindergarten class. My students are considered at-risk in order to qualify for the pre-kindergarten program at my school. Many of my students receive speech services. My students are bright, energetic, eager to learn children who I know will go far with the right start.\r\n\r\nI teach at a title I school in the suburbs of Chicago. My district provides free breakfast and lunch to all students. My students enjoy listening to stories. Everyday I read at least two books to my students. This component to my listening station will enable my students to listen to even more books daily. My students will learn to independently listen to stories. Students can listen to a story alone or with a partner. They will learn the skills to follow along in the book, using audio and visual cues as to when to turn the pages.\r\n\r\n\r\nEarly exposure to reading can create a life long love of books! \r\n\r\nThis reading center will help do that. I have numerous books that are age appropriate that will foster independent reading.nannan

=====

In [32]:

```
#\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)
```

My students are three, four and five year olds that attend either my a.m. class or my p.m. pre-kindergarten class. My students are considered at-risk in order to qualify for the pre-kindergarten program at my school. Many of my students receive speech services. My students are bright, energetic, eager to learn children who I know will go far with the right start. I teach at a title I school in the suburbs of Chicago. My district provides free breakfast and lunch to all students. My students enjoy listening to stories. Everyday I read at least two books to my students. This component to my listening station will enable my students to listen to even more books daily. My students will learn to independently listen to stories. Students can listen to a story alone or with a partner. They will learn the skills to follow along in the book, using audio and visual cues as to when to turn the pages. Early exposure to reading can create a life long love of books! This reading center will help do that. I have numerous books that are age appropriate that will foster independent reading.

In [33]:

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

My students are three four and five year olds that attend either my a m class or my p m pre kindergarten class My students are considered at risk in order to qualify for the pre kindergarten program at my school Many of my students receive speech services My students are bright energetic eager to learn children who I know will go far with the right start I teach at a title I school in the suburbs of Chicago My district provides free breakfast and lunch to all students My students enjoy listening to stories Everyday I read at least two books to my students This component to my listening station will enable my students to listen to even more books daily My students will learn to independently listen to stories Students can listen to a story alone or with a partner They will learn the skills to follow along in the book using audio and visual cues as to when to turn the pages Early exposure to reading can create a life long love of books This reading center will help do that I have numerous books that are age appropriate that will foster independent reading

In [34]:

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# because although they are in this list but they matter a lot because
# they change the meaning of the entire sentence.
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you'r
e", "you've", \
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him',
'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 't
hey', 'them', 'their', \
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "th
at'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'ha
d', 'having', 'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as'
, 'until', 'while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through'
, 'during', 'before', 'after', \
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'ov
er', 'under', 'again', 'further', \
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'an
y', 'both', 'each', 'few', 'more', \
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too'
, 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'no
w', 'd', 'll', 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't",
'doesn', "doesn't", 'hadn', \
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'migh
tn', "mightn't", 'mustn', \
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'w
asn', "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [35]:

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

```
100%|████████████████████████████████████████████████████████████████████████████████| 36052/36052 [00:15<00:00, 2289.59it/s]
```

In [36]:

```
# after preprocessing of project essays
preprocessed_essays_Test[20000]
```

Out[36]:

```
'students three four five year olds attend either class p pre kindergarten
class students considered risk order qualify pre kindergarten program scho
ol many students receive speech services students bright energetic eager l
earn children know go far right start teach title school suburbs chicago d
istrict provides free breakfast lunch students students enjoy listening st
ories everyday read least two books students component listening station e
nable students listen even books daily students learn independently listen
stories students listen story alone partner learn skills follow along book
using audio visual cues turn pages early exposure reading create life long
love books reading center help numerous books age appropriate foster indep
endent reading nannan'
```

## Preprocessing of project\_title

**Now we will simply apply the above preprocessing steps on the project title for the test data as well, as it is also a text feature**

In [37]:

```
# printing some random titles.
print(X_test['project_title'].values[0])
print("="*50)
print(X_test['project_title'].values[150])
print("="*50)
print(X_test['project_title'].values[1000])
print("="*50)
print(X_test['project_title'].values[20000])
print("="*50)
print(X_test['project_title'].values[9999])
print("="*50)
```

Putting Our Hands-On Learning

=====

Let's Get Coloring!

=====

Brighten My Student's School Year! Part 3

=====

Independent Reading In Pre-Kindergarten!

=====

Flexible in First

=====

**We have already written the preprocessing codes for different preprocessing approaches now we will simply use those codes on the project titles**



In [40]:

```
categories = list(X_cv['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 replacing all the ' ' (space) with '' (empty) ex: "Math & Science"=> "Math&Science"
            temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing spaces
    temp = temp.replace('&', '_') # we are replacing the & value into
    cat_list.append(temp.strip())

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

from collections import Counter
my_counter = Counter()
for word in X_test['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]))
```

## Preprocessing of project\_subject\_subcategories

In [41]:

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

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

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

X_cv['clean_subcategories'] = sub_cat_list
X_cv.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 X_cv['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]))
```

## Text Preprocessing

In [42]:

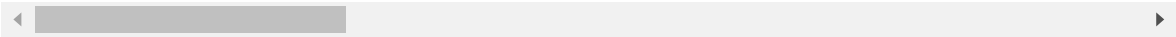
```
# merge two column text dataframe:
X_cv["essay"] = X_cv["project_essay_1"].map(str) + \
                X_cv["project_essay_2"].map(str) + \
                X_cv["project_essay_3"].map(str) + \
                X_cv["project_essay_4"].map(str)
```

In [43]:

```
X_cv.head(5)
```

Out[43]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
83534	147714	p027347	12735fdf8d2d34bc18aed105a86f25f2	Mrs.	WA
74468	87925	p138058	6c80fb44be75778c35db983dd5adf707	Mrs.	NY
47678	31762	p139835	617af30ebdd0b231867955f6ee16d793	Mrs.	GA
16522	133707	p032906	2a518c2b4a6576c3b1367184c649d120	Ms.	CA
71102	69968	p176570	941044a67fae3ddda70fa20234f00e93	Mrs.	NY





In [44]:

```
# printing some random reviews
print(X_cv['essay'].values[0])
print("="*50)
print(X_cv['essay'].values[150])
print("="*50)
print(X_cv['essay'].values[1000])
print("="*50)
print(X_cv['essay'].values[20000])
print("="*50)
print(X_cv['essay'].values[9999])
print("="*50)
```

This is a highly diverse school, over 50% are African American, 5% white, and the remaining 45% are huge mix. This school is situated in one of the most linguistically diverse neighborhoods in the United States, and our school is proof of that. The majority of our students are on the free and reduced lunch program, and many of our families have several adults working in the home (including many of our students). \r\n\r\nWhen describing our school, students and staff alike would say that we are a big family. We love each other and take care of each other. The staff at our school are, in my opinion, the best in the district, we care about our students and we put them first over everything. Also in my opinion, the students here are the BEST young adults in the entire universe. Although we have a strong academic program and a strong basketball program at our school, our arts program has a long way to go. We have some huge strides to make up in terms of equitable resources and it is my goal to grow the arts department and showcase the enormous talent that fills our hallways. For the last decade, students have mostly been exposed to painting with watercolor, drawing with pencil and marker and paper mache sculpture. I want to expose my students to even more exciting and awesome art materials. I plan to show them The Great Wave by Hokusai and have them experiment with woodblock printing with a modern twist. I want to show them how to play with batik and tie dyeing T-shirts and fabrics. I want my students to learn to discover their identities through self portraits with paint markers and paper cutting. I'd like my students to have real, lasting, memorable experiences in my classroom. I hope to inspire them to love art, to be amazed at what they can create, and to provide them with the opportunity to dive head-first into the world of fine art. I believe a good teacher not only sparks a fire in each student, but helps to keep that fire burning by showing students how to see the world in a new way. I hope that my lessons and the experiences my students gain from the projects in my room will inspire them to create art for the rest of their lives. nannan

=====

First grade is truly the age of discovery as school becomes less of a daunting new challenge and more of an exploration into how much students can learn. Our first grade classroom is full of students who are eager to learn. As a classroom that has close to 50% of students who use mostly Spanish (or other languages!) at home, we will be working hard to read and write all while trying to stay true to the culture that makes their families proud. \r\n\r\nWe are a class "on the move" learning through embodied cognition and trying to stay as active as possible. We strive to incorporate movement into every lesson or routine inside of our classroom, because as our bodies are moving, our brains are expanding. We are working hard and playing hard all school year long! Even in first grade, we are preparing students to someday become independent citizens and healthy, enriching members of society. In today's global and technologically-advancing world, even first graders need to learn how to be productive, safe, and creative on the computer. Unfortunately, the need for greater access to technology isn't always met with the increase in school resources and funding. This project will allow my students more time to use computers while not straining my school's current technology. \r\nThese chromebooks will allow my students greater access to the internet, to the apps that help them learn, and to the world of knowledge and wonder that lies just beyond the desktop screen. Students will no longer get frustrated by slow or out-dated technology, but instead will come into the class excited to use chromebooks that are new, up-to-date, and efficient. nannan

=====

Our Transitional Kindergarten Classroom has a morning class and an afternoon class. We have 21 students in each class, so we have 42 students sharing our classroom and materials each day! \r\n\r\nOur elementary school is a Title 1 school where many of our students receive free or reduced lunch. \r\n\r\nMy students come from very diverse backgrounds as well as many students in my class are English Language Learners. \r\n\r\nMy students are n

ow recognizing sight words in books and shout out to each other when they see something they know! We are in search of replacing our library foam mat with a fire resistant rug! Our foam mat did not pass fire code in the recent Fire Marshall visit. \r\n\r\nWith 42 students sharing our classroom and materials each day in our Transitional Kindergarten Classroom, our materials and our environment get lots of love! Other than our main letter rug for circle time, our entire classroom is covered with cold tile! \r\n\r\nIn our classroom library our students retell stories and songs for literacy centers, listen to stories with our listening center station, sit in our library during projector screen time, present puppet shows together, create felt board stories, participate in small group reading time with Ms. G-W or Mrs. Mew, and of course they enjoy getting cozy to read books of their choice during Free Choice time. \r\n\r\nPlease help us make our library a cozy and inviting place to be and read! \r\n\r\nThank you! \r\nMs. G-W, Mrs. Mew, AM TK and PM TKnannan

=====

I work for a Title I school with students from diverse backgrounds. We face many financial challenges as a Title I School, and we value every resource that we have. I do my best to challenge each and every one of my students, because I have students with various academic considerations. Therefore, I work diligently to differentiate and individualize instruction, so that everyone can learn to the best of their ability. \r\n\r\nMy goal is for every student to grow and succeed - academically and socially - and come to school with a curiosity for learning. My students are dynamic and sweet. Our school is community-oriented and supportive. We value partnerships with families, and we look for ways to help our students grow. Our school has created strong partnerships in the community this year to build conscious, young learners. My students value giving back to others in return. I love to create curious, confident learners. Technology integration is one of my greatest passions in life...next to teaching that is. Research shows that students learn best and actually retain more information when technology is incorporated into their every day learning. In my classroom, I do a lot of small group instruction as well as instructional videos and digital portfolios. As a completely portable learning tool, the iPad allows documentation and student sharing to be taken to a whole different level. The iPad will also be linked up to our classroom projector to showcase student work during our share out time as well as being hooked up to a TV near the small group table for quick and easy access to information during reading, writing, and math groups. This makes it even easier to cater my instruction to individual students' needs. nannan

=====

My class is generally made up of overly \"active\" students, mostly boys. At this age, movement is very important. I encourage them to learn through song, dance, and movement. I am hoping to further their ability to have the freedom to move/learn through flexible seating. There is an abundance of research that shows students who are allowed to sit in various types of seating (with the ability to move and not be restrained to traditional seating styles), have increased ability to be academically successful. \r\n\r\nI hope to remove most of my traditional desk/chair type seating from my classroom and replace them with rocker chairs, bean bag seats, bench seats, stability balls, wobble chairs, and floor pillows. \r\n\r\nReplacing traditional desk/chairs with flexible seating will allow my students to move while learning. This type of environment allows my \"active\" students the freedom to move while still working. My class is generally made up of overly \"active\" students, mostly boys. At this age, movement is very important. I encourage them to learn through song, dance, and movement. I am hoping to further their ability to have the freedom to move/learn through flexible seating. There is an abundance of research that shows students who are allowed to sit in various types of seating (with the ability to move and not be restrained to traditional seating styles), have increased ability to be academically successful. \r\n\r\nI hope to remove most of my

traditional desk/chair type seating from my classroom and replace them with rocker chairs, bean bag seats, bench seats, stability balls, wobble chairs, and floor pillows. \r\n\r\nReplacing traditional desk/chairs with flexible seating will allow my students to move while learning. This type of environment allows my \"active\" students the freedom to move while still working.nannan

=====

In [45]:

```
# creating a function named as decontracted which does the job of decontraction

# 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 [46]:

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

I work for a Title I school with students from diverse backgrounds. We face many financial challenges as a Title I School, and we value every resource that we have. I do my best to challenge each and every one of my students, because I have students with various academic considerations. Therefore, I work diligently to differentiate and individualize instruction, so that everyone can learn to the best of their ability. \r\n\r\nMy goal is for every student to grow and succeed - academically and socially - and come to school with a curiosity for learning. My students are dynamic and sweet. Our school is community-oriented and supportive. We value partnerships with families, and we look for ways to help our students grow. Our school has created strong partnerships in the community this year to build conscious, young learners. My students value giving back to others in return. I love to create curious, confident learners. Technology integration is one of my greatest passions in life...next to teaching that is. Research shows that students learn best and actually retain more information when technology is incorporated into their every day learning. In my classroom, I do a lot of small group instruction as well as instructional videos and digital portfolios. As a completely portable learning tool, the iPad allows documentation and student sharing to be taken to a whole different level. The iPad will also be linked up to our classroom projector to showcase student work during our share out time as well as being hooked up to a TV near the small group table for quick and easy access to information during reading, writing, and math groups. This makes it even easier to cater my instruction to individual students' needs.nannan

=====

In [47]:

```
#\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)
```

I work for a Title I school with students from diverse backgrounds. We face many financial challenges as a Title I School, and we value every resource that we have. I do my best to challenge each and every one of my students, because I have students with various academic considerations. Therefore, I work diligently to differentiate and individualize instruction, so that everyone can learn to the best of their ability. My goal is for every student to grow and succeed - academically and socially - and come to school with a curiosity for learning. My students are dynamic and sweet. Our school is community-oriented and supportive. We value partnerships with families, and we look for ways to help our students grow. Our school has created strong partnerships in the community this year to build conscious, young learners. My students value giving back to others in return. I love to create curious, confident learners. Technology integration is one of my greatest passions in life...next to teaching that is. Research shows that students learn best and actually retain more information when technology is incorporated into their every day learning. In my classroom, I do a lot of small group instruction as well as instructional videos and digital portfolios. As a completely portable learning tool, the iPad allows documentation and student sharing to be taken to a whole different level. The iPad will also be linked up to our classroom projector to showcase student work during our share out time as well as being hooked up to a TV near the small group table for quick and easy access to information during reading, writing, and math groups. This makes it even easier to cater my instruction to individual students' needs.

In [48]:

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# because although they are in this list but they matter a lot because
# they change the meaning of the entire sentence.
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you'r
e", "you've", \
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him',
'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 't
hey', 'them', 'their', \
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "th
at'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'ha
d', 'having', 'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as'
, 'until', 'while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through'
, 'during', 'before', 'after', \
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'ov
er', 'under', 'again', 'further', \
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'an
y', 'both', 'each', 'few', 'more', \
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too'
, 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'no
w', 'd', 'll', 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't",
'doesn', "doesn't", 'hadn', \
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'migh
tn', "mightn't", 'mustn', \
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'w
asn', "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [49]:

```
# Combining all the above preprocessing techniques for all the project essays
from tqdm import tqdm
preprocessed_essays_Cv = []
# tqdm is for printing the status bar
for sentence in tqdm(X_cv['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\n', ' ')
    sent = sent.replace('\\t', ' ')
    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_Cv.append(sent.lower().strip())
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
████████████████████████████████████████████████████████████████████████████████| 24155/24155 [00:10<00:00, 2218.92it/s]
```

In [50]:

```
# after preprocessing of project essays
preprocessed_essays_Test[20000]
```

Out[50]:

```
'students three four five year olds attend either class p pre kindergarten
class students considered risk order qualify pre kindergarten program scho
ol many students receive speech services students bright energetic eager l
earn children know go far right start teach title school suburbs chicago d
istrict provides free breakfast lunch students students enjoy listening st
ories everyday read least two books students component listening station e
nable students listen even books daily students learn independently listen
stories students listen story alone partner learn skills follow along book
using audio visual cues turn pages early exposure reading create life long
love books reading center help numerous books age appropriate foster indep
endent reading nannan'
```

## Preprocessing of project\_title

**Now we will simply apply the above preprocessing steps on the project title for the Cross Validation data as well, as it is also a text feature**

In [51]:

```
# printing some random titles.
print(X_cv['project_title'].values[0])
print("="*50)
print(X_cv['project_title'].values[150])
print("="*50)
print(X_cv['project_title'].values[1000])
print("="*50)
print(X_cv['project_title'].values[20000])
print("="*50)
print(X_cv['project_title'].values[9999])
print("="*50)
```

Help Fund Memorable Art Projects!

=====

Crazy for Chromebooks

=====

In Need of a Cozy Place to BE and READ!

=====

21st Century Guided Math and Literacy Rotations

=====

Creating a friendly learning environment for active learners....

=====





In [54]:

```
project_data.columns
```

Out[54]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',  
      'Date', '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_approve  
d'],  
      dtype='object')
```

we are going to consider

- school\_state : categorical data
- clean\_categories : categorical data
- clean\_subcategories : categorical data
- project\_grade\_category : categorical data
- teacher\_prefix : categorical data
- project\_title : text data
- text : text data
- project\_resource\_summary: text data (optinal)
- quantity : numerical (optinal)
- teacher\_number\_of\_previously\_posted\_projects : numerical
- price : numerical

**Now firstly we will be vectorizing the categorical data**

**For vectorizing the categorical data we will be using One Hot Encoding Technique**

**One Hot Encoding Of Project Clean Categories**

In [55]:

```
# we use count vectorizer to convert the values into one hot encoded features

from sklearn.feature_extraction.text import CountVectorizer

vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False,
                             binary=True)

# we will be using the X_train for fitting our model because that is the only data a user knows rest all are for testing purposes
vectorizer.fit(X_train['clean_categories'].values)

print(vectorizer.get_feature_names())

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

['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
```

In [56]:

```
print("Shape of Train data matrix after one hot encoding ",categories_one_hot_train.shape)
print("Shape of Test data matrix after one hot encoding ",categories_one_hot_test.shape)
print("Shape of CV data matrix after one hot encoding ",categories_one_hot_cv.shape)

Shape of Train data matrix after one hot encoding (49041, 9)
Shape of Test data matrix after one hot encoding (36052, 9)
Shape of CV data matrix after one hot encoding (24155, 9)
```

## One Hot Encoding Of Cleaned Project Sub Category

In [57]:

```
# we use count vectorizer to convert the values into one hot encoded features

from sklearn.feature_extraction.text import CountVectorizer

vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)

# we will be using the X_train for fitting our model because that is the only data a user knows rest all are for testing purposes
vectorizer.fit(X_train['clean_subcategories'].values)

print(vectorizer.get_feature_names())

subcategories_one_hot_train = vectorizer.transform(X_train['clean_subcategories'].values)

subcategories_one_hot_test = vectorizer.transform(X_test['clean_subcategories'].values)

subcategories_one_hot_cv = vectorizer.transform(X_cv['clean_subcategories'].values)
```

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

In [58]:

```
print("Shape of Train data matrix after one hot encoding ",subcategories_one_hot_train.shape)
print("Shape of Test data matrix after one hot encoding ",subcategories_one_hot_test.shape)
print("Shape of CV data matrix after one hot encoding ",subcategories_one_hot_cv.shape)
```

```
Shape of Train data matrix after one hot encoding (49041, 30)
Shape of Test data matrix after one hot encoding (36052, 30)
Shape of CV data matrix after one hot encoding (24155, 30)
```

## One hot encoding of teacher prefix

In [59]:

```
mylist_teacher_prefix = list(X_train['teacher_prefix'])
```

In [60]:

```
# We are removing the duplicate values from our list of the teacher prefix  
# Source :- https://www.w3schools.com/python/python\_howto\_remove\_duplicates.asp  
mylist_teacher_prefix_actual_Train = list(dict.fromkeys(mylist_teacher_prefix))
```

In [61]:

```
# removing the nan from the teacher prefix category as there is no such category of teacher which exists  
  
mylist_teacher_prefix_actual_Train = [p for p in mylist_teacher_prefix_actual_Train if str(p) != 'nan']  
mylist_teacher_prefix_actual_Train
```

Out[61]:

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

In [62]:

```
# we use count vectorizer to convert the values into one hot encoded features  
# now we are working on teacher prefix data  
  
from sklearn.feature_extraction.text import CountVectorizer  
vectorizer = CountVectorizer(vocabulary=mylist_teacher_prefix_actual_Train, lowercase=False, binary=True)  
  
# I was getting an error like "np.nan is an invalid document, expected byte or unicode string."  
# below is the solution  
  
# https://stackoverflow.com/questions/39303912/tfidfvectorizer-in-scikit-learn-valueerror-np-nan-is-an-invalid-document  
  
vectorizer.fit(X_train['teacher_prefix'].values.astype('U'))  
print(vectorizer.get_feature_names())  
  
teacher_prefix_one_hot_train = vectorizer.transform(X_train['teacher_prefix'].values.astype('U'))  
teacher_prefix_one_hot_test = vectorizer.transform(X_test['teacher_prefix'].values.astype('U'))  
teacher_prefix_one_hot_cv = vectorizer.transform(X_cv['teacher_prefix'].values.astype('U'))
```

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

In [63]:

```
print("Shape of matrix of Train data after one hot encoding ",teacher_prefix_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding ",teacher_prefix_one_hot_test.shape)
print("Shape of matrix of Cross Validation data after one hot encoding ",teacher_prefix_one_hot_cv.shape)
```

Shape of matrix of Train data after one hot encoding (49041, 5)  
 Shape of matrix of Test data after one hot encoding (36052, 5)  
 Shape of matrix of Cross Validation data after one hot encoding (24155, 5)

## One Hot encoding of project grade category

In [64]:

```
mylist_project_grade_category = list(X_train['project_grade_category'])
```

In [65]:

```
# We are removing the duplicate values from our list of the project grade category
# Source :- https://www.w3schools.com/python/python_howto_remove_duplicates.asp
mylist_project_grade_category_actual = list(dict.fromkeys(mylist_project_grade_category))
```

In [66]:

```
type(mylist_project_grade_category_actual)
print(mylist_project_grade_category_actual[0])

n = len(mylist_project_grade_category_actual)
print(n)

# I already saw by running the code that the word Grades is unnecessarily present in the elements of list hence trying to remove that word
# how to remove a word from a sentence --> https://codescracker.com/python/program/python-program-remove-word-from-sentence.htm
for m in range(0,4,1):
    words = mylist_project_grade_category_actual[m].split()
    mylist_project_grade_category_actual[m] = ''.join([j for j in words if j not in "Grades"])
print(mylist_project_grade_category_actual)
```

Grades 3-5

4

['3-5', 'PreK-2', '6-8', '9-12']

In [67]:

```
# we use count vectorizer to convert the values into one hot encoded features

# now we are working on project grade category data

from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary=mylist_project_grade_category_actual, lowercase
=False, binary=True)

vectorizer.fit(X_train['project_grade_category'].values)
print(vectorizer.get_feature_names())

project_grade_categories_one_hot_train = vectorizer.transform(X_train['project_grade_ca
tegory'].values)
project_grade_categories_one_hot_test = vectorizer.transform(X_test['project_grade_cate
gory'].values)
project_grade_categories_one_hot_cv = vectorizer.transform(X_cv['project_grade_categor
y'].values)

['3-5', 'PreK-2', '6-8', '9-12']
```

In [68]:

```
print("Shape of matrix of Train data after one hot encoding ",project_grade_categories_
one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding ",project_grade_categories_o
ne_hot_test.shape)
print("Shape of matrix of Cross Validation data after one hot encoding ",project_grade_
categories_one_hot_cv.shape)
```

```
Shape of matrix of Train data after one hot encoding (49041, 4)
Shape of matrix of Test data after one hot encoding (36052, 4)
Shape of matrix of Cross Validation data after one hot encoding (24155,
4)
```

## One hot encoding of School States

In [69]:

```
type(list(project_data['school_state']))
```

Out[69]:

```
list
```

In [70]:

```
mylist = list(X_train['school_state'])
```

In [71]:

```
# We are removing the duplicate values from our list of the state codes

# Source :- https://www.w3schools.com/python/python\_howto\_remove\_duplicates.asp

mylist_actual = list(dict.fromkeys(mylist))
```

In [72]:

```
# we use count vectorizer to convert the values into one hot encoded features

# now we are working on school state data

from sklearn.feature_extraction.text import CountVectorizer

vectorizer = CountVectorizer(vocabulary=mylist_actual, lowercase=False, binary=True)

vectorizer.fit(X_train['school_state'].values)

print(vectorizer.get_feature_names())

school_state_categories_one_hot_train = vectorizer.transform(X_train['school_state'].values)
school_state_categories_one_hot_test = vectorizer.transform(X_test['school_state'].values)
school_state_categories_one_hot_cv = vectorizer.transform(X_cv['school_state'].values)

['CA', 'MA', 'TX', 'GA', 'NJ', 'NC', 'AL', 'NY', 'VA', 'PA', 'OK', 'UT',
 'KS', 'OR', 'SC', 'FL', 'AZ', 'CO', 'DE', 'MT', 'WA', 'MN', 'MI', 'NM', 'M
D', 'NV', 'CT', 'AR', 'IL', 'WI', 'IN', 'MS', 'MO', 'LA', 'OH', 'TN', 'I
D', 'WV', 'RI', 'KY', 'AK', 'SD', 'ME', 'WY', 'NH', 'IA', 'NE', 'DC', 'V
T', 'HI', 'ND']
```

In [73]:

```
print("Shape of matrix of Train data after one hot encoding ", school_state_categories_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding ", school_state_categories_one_hot_test.shape)
print("Shape of matrix of Cross Validation data after one hot encoding ", school_state_categories_one_hot_cv.shape)
```

```
Shape of matrix of Train data after one hot encoding (49041, 51)
Shape of matrix of Test data after one hot encoding (36052, 51)
Shape of matrix of Cross Validation data after one hot encoding (24155, 51)
```

**Now as we have done the vectorizing of categorical data now we will be vectorizing the text data using different techniques**

## Vectorizing the text data

### Technique -1 Bag of words(BOW)

## Essays

### Essays Train Data



In [74]:

```
# I am not setting the value of min_df here because i read here  
# https://stackoverflow.com/questions/27697766/understanding-min-df-and-max-df-in-scikit-countvectorizer  
# that min_df helps in better performance but might also give poor clusters hence am trying without it this time  
  
vectorizer = CountVectorizer()  
vectorizer.fit(preprocessed_essays_Train)  
  
essay_bow_train = vectorizer.transform(preprocessed_essays_Train)  
  
print("Shape of matrix after bag of words ", essay_bow_train.shape)
```

Shape of matrix after bag of words (49041, 41074)

## Essay Test Data

In [75]:

```
# I am not setting the value of min_df here because i read here  
# https://stackoverflow.com/questions/27697766/understanding-min-df-and-max-df-in-scikit-countvectorizer  
# that min_df helps in better performance but might also give poor clusters hence am trying without it this time  
  
  
# now note that the below two lines are wrong because we have already trained the  
# model on the training data now using that model we should get the bow representation  
# of test data. After all training from test data only and then checking for its accuracy will ofcourse give  
# good accuracy.  
  
# lines not to be used (I used them but then going through the code realised the mistake)  
  
#vectorizer = CountVectorizer()  
#vectorizer.fit(preprocessed_essays_Test)  
  
essay_bow_test = vectorizer.transform(preprocessed_essays_Test)  
  
print("Shape of matrix after bag of words ", essay_bow_test.shape)
```

Shape of matrix after bag of words (36052, 41074)

## Essay Cross Validation data

In [76]:

```
# I am not setting the value of min_df here because i read here
# https://stackoverflow.com/questions/27697766/understanding-min-df-and-max-df-in-scikit-countvectorizer
# that min_df helps in better performance but might also give poor clusters hence am trying without it this time

# similarly below two lines should not be used

#vectorizer = CountVectorizer()
#vectorizer.fit(preprocessed_essays_Cv)

essay_bow_cv = vectorizer.transform(preprocessed_essays_Cv)

print("Shape of matrix after bag of words ",essay_bow_cv.shape)
```

Shape of matrix after bag of words (24155, 41074)

## Project Title

### Bag of words on Project Title Train data

In [77]:

```
vectorizer.fit(preprocessed_project_titles_Train)
project_title_bow_train = vectorizer.transform(preprocessed_project_titles_Train)
print("Shape of matrix after bag of words ",project_title_bow_train.shape)
```

Shape of matrix after bag of words (49041, 11659)

### Bag of words on Project Title Test data

In [78]:

```
project_title_bow_test = vectorizer.transform(preprocessed_project_titles_Test)
print("Shape of matrix after bag of words ",project_title_bow_test.shape)
```

Shape of matrix after bag of words (36052, 11659)

### Bag of words on Project Title Cross Validation data

In [79]:

```
project_title_bow_cv = vectorizer.transform(preprocessed_project_titles_Cv)
print("Shape of matrix after bag of words ",project_title_bow_cv.shape)
```

Shape of matrix after bag of words (24155, 11659)

## Technique-2 TFIDF

## Essay Data

### Essay Train Data

In [156]:

```
from sklearn.feature_extraction.text import TfidfVectorizer

vectorizer = TfidfVectorizer(max_features = 6000, min_df = 10)
vectorizer.fit(preprocessed_essays_Train)

text_tfidf_train = vectorizer.transform(preprocessed_essays_Train)
print("Shape of matrix after tfidf ",text_tfidf_train.shape)

print(type(text_tfidf_train))

# we are converting a dictionary with word as a key, and the idf as a value

dictionary = dict(zip(vectorizer.get_feature_names(), list(vectorizer.idf_)))
tfidf_words = set(dictionary.keys())
```

Shape of matrix after tfidf (49041, 6000)  
<class 'scipy.sparse.csr.csr\_matrix'>

### Essay Test data

In [157]:

```
text_tfidf_test = vectorizer.transform(preprocessed_essays_Test)
print("Shape of matrix after tfidf ",text_tfidf_test.shape)
```

Shape of matrix after tfidf (36052, 6000)

### Essay Cross Validation Data

In [158]:

```
text_tfidf_cv = vectorizer.transform(preprocessed_essays_Cv)
print("Shape of matrix after tfidf ",text_tfidf_cv.shape)
```

Shape of matrix after tfidf (24155, 6000)

## Project Title

### Project title train data

In [159]:

```
vectorizer = TfidfVectorizer(max_features = 2000, min_df = 10)

vectorizer.fit(preprocessed_project_titles_Train)

project_title_tfidf_train = vectorizer.transform(preprocessed_project_titles_Train)

print("Shape of matrix after tfidf ",project_title_tfidf_train.shape)

# we are converting a dictionary with word as a key, and the idf as a value

dictionary = dict(zip(vectorizer.get_feature_names(), list(vectorizer.idf_)))
tfidf_project_title_words = set(dictionary.keys())
```

Shape of matrix after tfidf (49041, 2000)

## Project title test data

In [160]:

```
project_title_tfidf_test = vectorizer.transform(preprocessed_project_titles_Test)
print("Shape of matrix after tfidf ",project_title_tfidf_test.shape)
```

Shape of matrix after tfidf (36052, 2000)

## Project title cross validation data

In [161]:

```
title_tfidf_cv = vectorizer.transform(preprocessed_project_titles_Cv)
print("Shape of matrix after tfidf ",title_tfidf_cv.shape)
```

Shape of matrix after tfidf (24155, 2000)

# Technique-3 Average Word to Vector

## Using pretrained w2v model in the file glove

In [86]:

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

## Essay Train data

In [87]:

```
# average Word2Vec
# compute average word2vec for each preprocessed essay.

avg_w2v_vectors_train = []; # the avg-w2v for each essay is stored in this list
for sentence in tqdm(preprocessed_essays_Train): # for each essay
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the essay
    for word in sentence.split(): # for each word in a essay
        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]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
████████| 49041/49041 [00:10<00:00, 4615.36it/s]
```

```
49041
300
```

## Essay Test data

In [88]:

```
# average Word2Vec
# compute average word2vec for each preprocessed essay.

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

print(len(avg_w2v_vectors_test))
print(len(avg_w2v_vectors_test[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
████████| 36052/36052 [00:07<00:00, 4621.37it/s]
```

```
36052
300
```

## Essay Cross Validation data

In [89]:

```
# average Word2Vec
# compute average word2vec for each preprocessed essay.

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

print(len(avg_w2v_vectors_cv))
print(len(avg_w2v_vectors_cv[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
██████| 24155/24155 [00:05<00:00, 4598.35it/s]
```

```
24155
300
```

## Project Title train data

In [90]:

```
# average Word2Vec
# compute average word2vec for each preprocessed essay.

avg_w2v_vectors_project_title_train = []; # the avg-w2v for each essay is stored in this list
for sentence in tqdm(preprocessed_project_titles_Train): # for each essay
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the essay
    for word in sentence.split(): # for each word in a essay
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_project_title_train.append(vector)

print(len(avg_w2v_vectors_project_title_train))
print(len(avg_w2v_vectors_project_title_train[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
██████| 49041/49041 [00:00<00:00, 84872.79it/s]
```

```
49041
300
```

## Project Title test data



## Essay train data



In [162]:

```
# tfidf Word2Vec
# computing average word2vec for each Project Title is stored in this list

tfidf_w2v_vectors_text_train = []; # the tfidf-w2v for each essay
for sentence in tqdm(preprocessed_essays_Train): # for each 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 title
    for word in sentence.split(): # for each word in a title
        if (word in glove_words) and (word in dictionary):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sen
            tence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # ge
            tting 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_text_train.append(vector)

print(len(tfidf_w2v_vectors_text_train))
print(len(tfidf_w2v_vectors_text_train[0]))
```

```

0%|
| 0/49041 [00:00<?, ?it/s]

0%||
| 84/49041 [00:00<00:58, 833.63it/s]

0%||
| 163/49041 [00:00<00:59, 819.48it/s]

1%||
| 247/49041 [00:00<00:59, 823.83it/s]

1%||
| 341/49041 [00:00<00:57, 848.86it/s]

1%||
| 432/49041 [00:00<00:56, 865.18it/s]

1%||
| 518/49041 [00:00<00:56, 860.94it/s]

1%||
| 596/49041 [00:00<00:58, 833.01it/s]

1%||
| 674/49041 [00:00<00:59, 815.43it/s]

2%||
| 759/49041 [00:00<00:58, 823.94it/s]

2%||
| 849/49041 [00:01<00:57, 843.23it/s]

2%||
| 932/49041 [00:01<00:57, 834.48it/s]

2%||
| 1017/49041 [00:01<00:57, 836.50it/s]

2%||
| 1102/49041 [00:01<00:57, 840.43it/s]

2%||
| 1186/49041 [00:01<00:57, 836.24it/s]

3%||
| 1270/49041 [00:01<00:58, 823.22it/s]

3%||
| 1353/49041 [00:01<00:58, 811.30it/s]

3%||
| 1434/49041 [00:01<00:59, 794.09it/s]

3%||
| 1516/49041 [00:01<00:59, 800.04it/s]

3%||
| 1599/49041 [00:01<00:58, 806.84it/s]

```

3%|█████  
| 1684/49041 [00:02<00:57, 817.60it/s]

4%|█████  
| 1766/49041 [00:02<00:58, 814.11it/s]

4%|█████  
| 1853/49041 [00:02<00:56, 828.80it/s]

4%|█████  
| 1940/49041 [00:02<00:56, 836.62it/s]

4%|█████  
| 2024/49041 [00:02<00:57, 823.45it/s]

4%|█████  
| 2107/49041 [00:02<00:58, 808.38it/s]

4%|█████  
| 2188/49041 [00:02<00:58, 807.02it/s]

5%|█████  
| 2272/49041 [00:02<00:57, 816.07it/s]

5%|█████  
| 2359/49041 [00:02<00:56, 830.66it/s]

5%|█████  
| 2445/49041 [00:02<00:55, 837.46it/s]

5%|█████  
| 2529/49041 [00:03<00:56, 824.29it/s]

5%|█████  
| 2621/49041 [00:03<00:54, 847.26it/s]

6%|█████  
| 2706/49041 [00:03<00:55, 841.29it/s]

6%|█████  
| 2791/49041 [00:03<00:55, 830.45it/s]

6%|█████  
| 2880/49041 [00:03<00:54, 842.88it/s]

6%|█████  
| 2965/49041 [00:03<00:55, 835.73it/s]

6%|█████  
| 3049/49041 [00:03<00:55, 825.82it/s]

6%|█████  
| 3135/49041 [00:03<00:54, 834.72it/s]

7%|█████  
| 3221/49041 [00:03<00:54, 840.89it/s]

7%|█████  
| 3311/49041 [00:03<00:53, 852.98it/s]

7%|█████

| 3404/49041 [00:04<00:52, 872.87it/s]

7%|███████

| 3493/49041 [00:04<00:51, 877.88it/s]

7%|███████

| 3581/49041 [00:04<00:52, 867.12it/s]

7%|███████

| 3668/49041 [00:04<00:52, 865.83it/s]

8%|███████

| 3755/49041 [00:04<00:53, 850.08it/s]

8%|███████

| 3843/49041 [00:04<00:52, 855.92it/s]

8%|███████

| 3930/49041 [00:04<00:52, 859.35it/s]

8%|███████

| 4017/49041 [00:04<00:52, 860.59it/s]

8%|███████

| 4104/49041 [00:04<00:52, 858.98it/s]

9%|███████

| 4190/49041 [00:05<00:52, 850.85it/s]

9%|███████

| 4279/49041 [00:05<00:52, 860.41it/s]

9%|███████

| 4366/49041 [00:05<00:52, 850.04it/s]

9%|███████

| 4456/49041 [00:05<00:51, 862.68it/s]

9%|███████

| 4543/49041 [00:05<00:52, 841.42it/s]

9%|███████

| 4628/49041 [00:05<00:53, 837.11it/s]

10%|███████

| 4715/49041 [00:05<00:52, 840.02it/s]

10%|███████

| 4800/49041 [00:05<00:53, 833.83it/s]

10%|███████

| 4884/49041 [00:05<00:53, 831.38it/s]

10%|███████

| 4968/49041 [00:05<00:54, 810.29it/s]

10%|███████

| 5057/49041 [00:06<00:52, 831.53it/s]

10%|███████

| 5142/49041 [00:06<00:52, 835.11it/s]

11%|██████████  
| 5232/49041 [00:06<00:51, 851.72it/s]

11%|██████████  
| 5318/49041 [00:06<00:51, 843.63it/s]

11%|██████████  
| 5405/49041 [00:06<00:51, 844.14it/s]

11%|██████████  
| 5490/49041 [00:06<00:52, 828.74it/s]

11%|██████████  
| 5575/49041 [00:06<00:52, 832.16it/s]

12%|██████████  
| 5659/49041 [00:06<00:53, 806.36it/s]

12%|██████████  
| 5748/49041 [00:06<00:52, 823.48it/s]

12%|██████████  
| 5835/49041 [00:06<00:51, 835.15it/s]

12%|██████████  
| 5934/49041 [00:07<00:49, 874.76it/s]

12%|██████████  
| 6023/49041 [00:07<00:50, 844.75it/s]

12%|██████████  
| 6109/49041 [00:07<00:51, 832.64it/s]

13%|██████████  
| 6200/49041 [00:07<00:50, 849.79it/s]

13%|██████████  
| 6286/49041 [00:07<00:51, 836.13it/s]

13%|██████████  
| 6370/49041 [00:07<00:51, 831.54it/s]

13%|██████████  
| 6454/49041 [00:07<00:51, 832.84it/s]

13%|██████████  
| 6538/49041 [00:07<00:52, 807.28it/s]

14%|██████████  
| 6623/49041 [00:07<00:52, 815.43it/s]

14%|██████████  
| 6707/49041 [00:08<00:51, 820.89it/s]

14%|██████████  
| 6795/49041 [00:08<00:50, 836.65it/s]

14%|██████████  
| 6879/49041 [00:08<00:51, 814.66it/s]

14%|██████████  
| 6968/49041 [00:08<00:50, 834.57it/s]

14%|██████████  
| 7052/49041 [00:08<00:50, 824.64it/s]

15%|██████████  
| 7135/49041 [00:08<00:51, 821.33it/s]

15%|██████████  
| 7222/49041 [00:08<00:50, 834.46it/s]

15%|██████████  
| 7306/49041 [00:08<00:50, 827.42it/s]

15%|██████████  
| 7389/49041 [00:08<00:51, 814.25it/s]

15%|██████████  
| 7478/49041 [00:08<00:49, 833.56it/s]

15%|██████████  
| 7564/49041 [00:09<00:49, 838.93it/s]

16%|██████████  
| 7652/49041 [00:09<00:48, 849.35it/s]

16%|██████████  
| 7739/49041 [00:09<00:48, 855.19it/s]

16%|██████████  
| 7833/49041 [00:09<00:46, 877.17it/s]

16%|██████████  
| 7921/49041 [00:09<00:47, 872.86it/s]

16%|██████████  
| 8009/49041 [00:09<00:48, 840.21it/s]

17%|██████████  
| 8095/49041 [00:09<00:48, 844.35it/s]

17%|██████████  
| 8180/49041 [00:09<00:49, 831.62it/s]

17%|██████████  
| 8264/49041 [00:09<00:49, 817.08it/s]

17%|██████████  
| 8346/49041 [00:09<00:50, 811.40it/s]

17%|██████████  
| 8428/49041 [00:10<00:53, 765.56it/s]

17%|██████████  
| 8510/49041 [00:10<00:52, 775.11it/s]

18%|██████████  
| 8591/49041 [00:10<00:51, 782.38it/s]

18%|██████████

| 8680/49041 [00:10<00:49, 811.36it/s]  
18%|██████████  
| 8762/49041 [00:10<00:50, 794.53it/s]  
18%|██████████  
| 8842/49041 [00:10<00:51, 778.19it/s]  
18%|██████████  
| 8921/49041 [00:10<00:52, 766.11it/s]  
18%|██████████  
| 9001/49041 [00:10<00:51, 775.90it/s]  
19%|██████████  
| 9079/49041 [00:10<00:52, 766.31it/s]  
19%|██████████  
| 9165/49041 [00:11<00:50, 790.61it/s]  
19%|██████████  
| 9250/49041 [00:11<00:49, 806.75it/s]  
19%|██████████  
| 9334/49041 [00:11<00:48, 814.64it/s]  
19%|██████████  
| 9421/49041 [00:11<00:47, 826.41it/s]  
19%|██████████  
| 9511/49041 [00:11<00:46, 846.26it/s]  
20%|██████████  
| 9596/49041 [00:11<00:47, 828.27it/s]  
20%|██████████  
| 9690/49041 [00:11<00:45, 857.25it/s]  
20%|██████████  
| 9777/49041 [00:11<00:46, 849.13it/s]  
20%|██████████  
| 9863/49041 [00:11<00:47, 816.61it/s]  
20%|██████████  
| 9946/49041 [00:11<00:49, 795.40it/s]  
20%|██████████  
| 10027/49041 [00:12<00:50, 777.38it/s]  
21%|██████████  
| 10108/49041 [00:12<00:49, 785.21it/s]  
21%|██████████  
| 10187/49041 [00:12<00:49, 778.84it/s]  
21%|██████████  
| 10266/49041 [00:12<00:50, 768.74it/s]  
21%|██████████  
| 10348/49041 [00:12<00:49, 778.66it/s]

21%|███████████████████|  
| 10427/49041 [00:12<00:49, 777.92it/s]

21%|███████████████████|  
| 10510/49041 [00:12<00:48, 788.60it/s]

22%|███████████████████|  
| 10591/49041 [00:12<00:48, 790.96it/s]

22%|███████████████████|  
| 10671/49041 [00:12<00:50, 756.03it/s]

22%|███████████████████|  
| 10748/49041 [00:13<00:50, 756.30it/s]

22%|███████████████████|  
| 10824/49041 [00:13<00:50, 751.68it/s]

22%|███████████████████|  
| 10906/49041 [00:13<00:49, 770.16it/s]

22%|███████████████████|  
| 10988/49041 [00:13<00:48, 783.77it/s]

23%|███████████████████|  
| 11077/49041 [00:13<00:46, 811.25it/s]

23%|███████████████████|  
| 11159/49041 [00:13<00:47, 800.23it/s]

23%|███████████████████|  
| 11240/49041 [00:13<00:47, 798.97it/s]

23%|███████████████████|  
| 11326/49041 [00:13<00:46, 814.65it/s]

23%|███████████████████|  
| 11419/49041 [00:13<00:44, 844.45it/s]

23%|███████████████████|  
| 11504/49041 [00:13<00:45, 829.48it/s]

24%|███████████████████|  
| 11592/49041 [00:14<00:44, 837.45it/s]

24%|███████████████████|  
| 11677/49041 [00:14<00:45, 829.55it/s]

24%|███████████████████|  
| 11761/49041 [00:14<00:45, 823.61it/s]

24%|███████████████████|  
| 11849/49041 [00:14<00:44, 839.10it/s]

24%|███████████████████|  
| 11935/49041 [00:14<00:44, 840.84it/s]

25%|███████████████████|  
| 12026/49041 [00:14<00:43, 856.27it/s]



[illegible]

```

13839/49041 [00:16<00:44, 792.76it/s]
28%|███████████|
| 13919/49041 [00:16<00:44, 792.83it/s]
29%|███████████|
| 13999/49041 [00:16<00:50, 692.04it/s]
29%|███████████|
| 14071/49041 [00:17<00:51, 684.75it/s]
29%|███████████|
| 14148/49041 [00:17<00:49, 704.91it/s]
29%|███████████|
| 14231/49041 [00:17<00:47, 734.60it/s]
29%|███████████|
| 14311/49041 [00:17<00:46, 751.91it/s]
29%|███████████|
| 14397/49041 [00:17<00:44, 776.88it/s]
30%|███████████|
| 14476/49041 [00:17<00:44, 772.77it/s]
30%|███████████|
| 14559/49041 [00:17<00:43, 785.29it/s]
30%|███████████|
| 14640/49041 [00:17<00:43, 791.65it/s]
30%|███████████|
| 14721/49041 [00:17<00:43, 795.22it/s]
30%|███████████|
| 14801/49041 [00:18<00:44, 773.66it/s]
30%|███████████|
| 14883/49041 [00:18<00:43, 783.24it/s]
31%|███████████|
| 14970/49041 [00:18<00:42, 805.47it/s]
31%|███████████|
| 15061/49041 [00:18<00:40, 832.66it/s]
31%|███████████|
| 15145/49041 [00:18<00:41, 808.34it/s]
31%|███████████|
| 15227/49041 [00:18<00:42, 787.11it/s]
31%|███████████|
| 15311/49041 [00:18<00:42, 801.36it/s]
31%|███████████|
| 15398/49041 [00:18<00:41, 819.77it/s]
32%|███████████|
| 15481/49041 [00:18<00:41, 811.90it/s]

```

```

32%|███████████████████████████████|
| 15563/49041 [00:18<00:41, 813.27it/s]

32%|███████████████████████████████|
| 15645/49041 [00:19<00:41, 813.59it/s]

32%|███████████████████████████████|
| 15733/49041 [00:19<00:40, 830.78it/s]

32%|███████████████████████████████|
| 15821/49041 [00:19<00:39, 842.83it/s]

32%|███████████████████████████████|
| 15915/49041 [00:19<00:38, 868.69it/s]

33%|███████████████████████████████|
| 16003/49041 [00:19<00:40, 816.39it/s]

33%|███████████████████████████████|
| 16088/49041 [00:19<00:40, 823.00it/s]

33%|███████████████████████████████|
| 16171/49041 [00:19<00:41, 797.33it/s]

33%|███████████████████████████████|
| 16252/49041 [00:19<00:41, 792.22it/s]

33%|███████████████████████████████|
| 16332/49041 [00:19<00:41, 785.50it/s]

33%|███████████████████████████████|
| 16418/49041 [00:19<00:40, 804.24it/s]

34%|███████████████████████████████|
| 16499/49041 [00:20<00:41, 792.62it/s]

34%|███████████████████████████████|
| 16580/49041 [00:20<00:40, 796.64it/s]

34%|███████████████████████████████|
| 16663/49041 [00:20<00:40, 800.67it/s]

34%|███████████████████████████████|
| 16744/49041 [00:20<00:40, 803.37it/s]

34%|███████████████████████████████|
| 16827/49041 [00:20<00:39, 811.08it/s]

34%|███████████████████████████████|
| 16912/49041 [00:20<00:39, 822.04it/s]

35%|███████████████████████████████|
| 17003/49041 [00:20<00:37, 844.37it/s]

35%|███████████████████████████████|
| 17088/49041 [00:20<00:37, 841.89it/s]

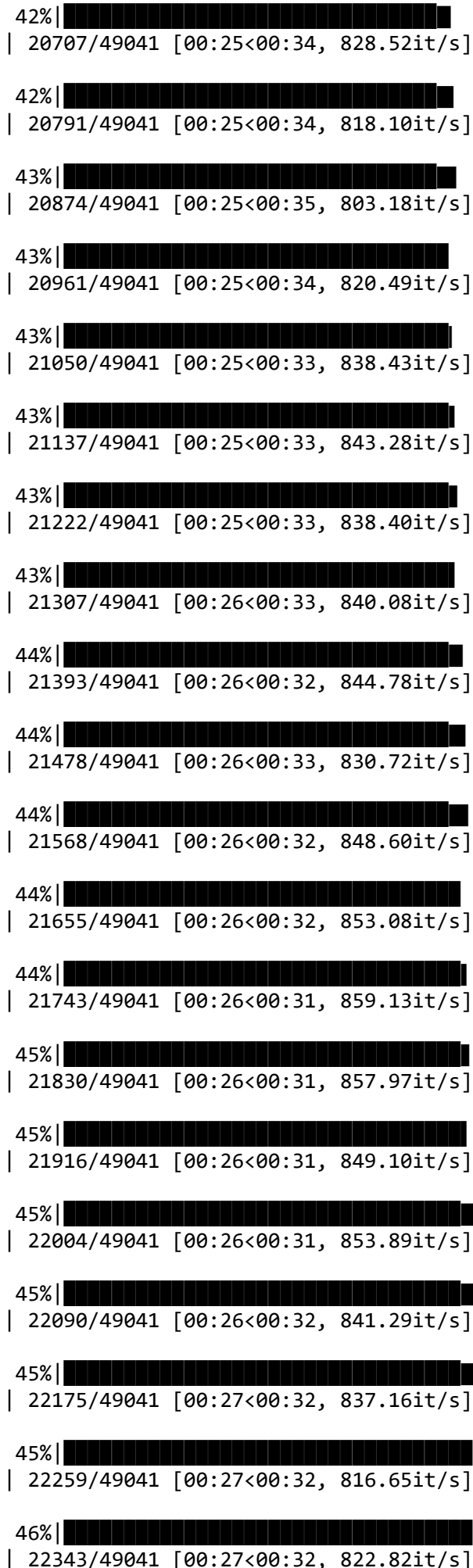
35%|███████████████████████████████|
| 17173/49041 [00:20<00:38, 824.89it/s]

```

[illegible]

```

18945/49041 [00:23<00:37, 803.76it/s]
39%|███████████|
19026/49041 [00:23<00:38, 774.05it/s]
39%|███████████|
19109/49041 [00:23<00:38, 783.93it/s]
39%|███████████|
19188/49041 [00:23<00:38, 779.40it/s]
39%|███████████|
19267/49041 [00:23<00:38, 773.99it/s]
39%|███████████|
19358/49041 [00:23<00:36, 808.75it/s]
40%|███████████|
19441/49041 [00:23<00:36, 813.18it/s]
40%|███████████|
19528/49041 [00:23<00:35, 827.71it/s]
40%|███████████|
19618/49041 [00:23<00:34, 846.46it/s]
40%|███████████|
19703/49041 [00:24<00:36, 814.17it/s]
40%|███████████|
19785/49041 [00:24<00:36, 804.57it/s]
41%|███████████|
19866/49041 [00:24<00:37, 778.98it/s]
41%|███████████|
19945/49041 [00:24<00:38, 748.53it/s]
41%|███████████|
20028/49041 [00:24<00:37, 769.38it/s]
41%|███████████|
20113/49041 [00:24<00:36, 786.55it/s]
41%|███████████|
20193/49041 [00:24<00:37, 778.30it/s]
41%|███████████|
20273/49041 [00:24<00:36, 780.64it/s]
42%|███████████|
20357/49041 [00:24<00:35, 797.02it/s]
42%|███████████|
20444/49041 [00:24<00:35, 815.98it/s]
42%|███████████|
20533/49041 [00:25<00:34, 835.07it/s]
42%|███████████|
20622/49041 [00:25<00:33, 846.65it/s]
```



[illegible]

[illegible]



53% |

```
| 25961/49041 [00:31<00:27, 836.43it/s]
```

53% |

```
| 26048/49041 [00:31<00:27, 844.36it/s]
```

53% |

```
| 26138/49041 [00:31<00:26, 856.07it/s]
```

53% |

```
| 26224/49041 [00:31<00:27, 832.03it/s]
```

54% |

```
| 26312/49041 [00:31<00:26, 844.08it/s]
```

54% |

```
| 26397/49041 [00:32<00:27, 829.22it/s]
```

54% |

```
| 26481/49041 [00:32<00:27, 810.72it/s]
```

54% |

```
| 26570/49041 [00:32<00:27, 831.04it/s]
```

54% |

```
| 26656/49041 [00:32<00:26, 838.29it/s]
```

55% |

```
| 26741/49041 [00:32<00:26, 837.54it/s]
```

55% |

```
| 26826/49041 [00:32<00:26, 839.42it/s]
```

55% |

```
| 26911/49041 [00:32<00:26, 832.33it/s]
```

55% |

```
| 26995/49041 [00:32<00:26, 828.62it/s]
```

55% |

```
| 27088/49041 [00:32<00:25, 852.50it/s]
```

55% |

```
| 27174/49041 [00:33<00:25, 842.88it/s]
```

56% |

```
| 27260/49041 [00:33<00:25, 846.18it/s]
```

56% |

```
| 27345/49041 [00:33<00:26, 823.42it/s]
```

56% |

```
| 27428/49041 [00:33<00:27, 799.84it/s]
```

56% |

```
| 27509/49041 [00:33<00:27, 791.69it/s]
```

56% |

```
| 27593/49041 [00:33<00:26, 803.84it/s]
```

[illegible]

```
| 31184/49041 [00:37<00:21, 829.96it/s]
```



[illegible]

[illegible]

[illegible]

[illegible]



| 39866/49041 [00:48<00:11, 821.65it/s]

81%|

| 39949/49041 [00:48<00:11, 792.94it/s]

82%|

| 40034/49041 [00:48<00:11, 806.34it/s]

82%|

| 40115/49041 [00:48<00:11, 796.26it/s]

82%|

| 40203/49041 [00:48<00:10, 815.66it/s]

82%|

| 40287/49041 [00:48<00:10, 821.05it/s]

82%|

| 40378/49041 [00:48<00:10, 843.80it/s]

83%|

| 40463/49041 [00:49<00:10, 832.73it/s]

83%|

| 40547/49041 [00:49<00:10, 823.24it/s]

83%|

| 40631/49041 [00:49<00:10, 826.48it/s]

83%|

| 40720/49041 [00:49<00:09, 842.75it/s]

83%|

| 40806/49041 [00:49<00:09, 846.02it/s]

83%|

| 40895/49041 [00:49<00:09, 856.94it/s]

84%|

| 40982/49041 [00:49<00:09, 858.95it/s]

84%|

| 41073/49041 [00:49<00:09, 871.82it/s]

84%|

| 41162/49041 [00:49<00:09, 875.31it/s]

84%|

| 41250/49041 [00:49<00:08, 869.61it/s]

84%|

| 41338/49041 [00:50<00:08, 870.73it/s]

84%|

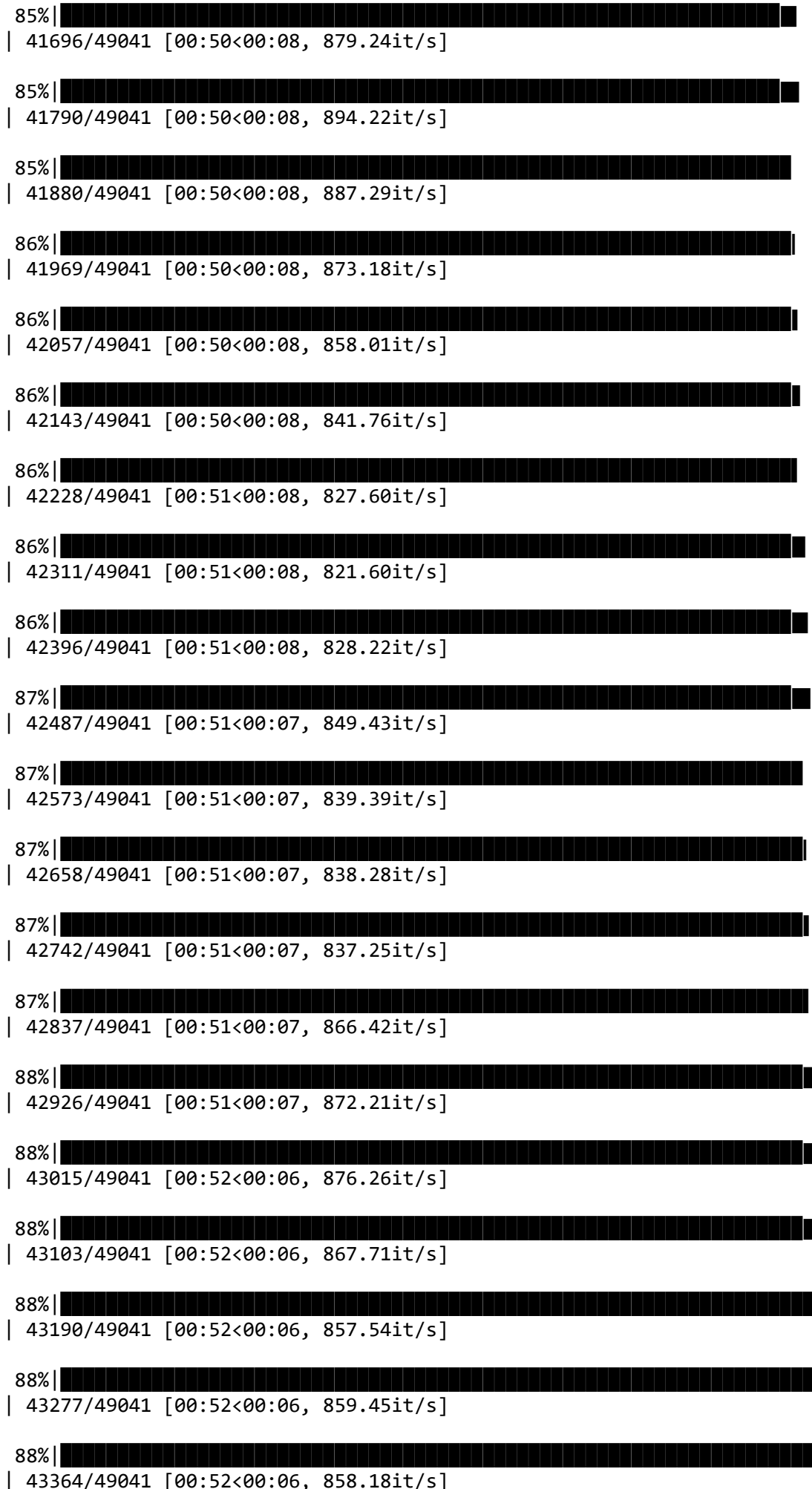
| 41426/49041 [00:50<00:08, 866.47it/s]

85%|

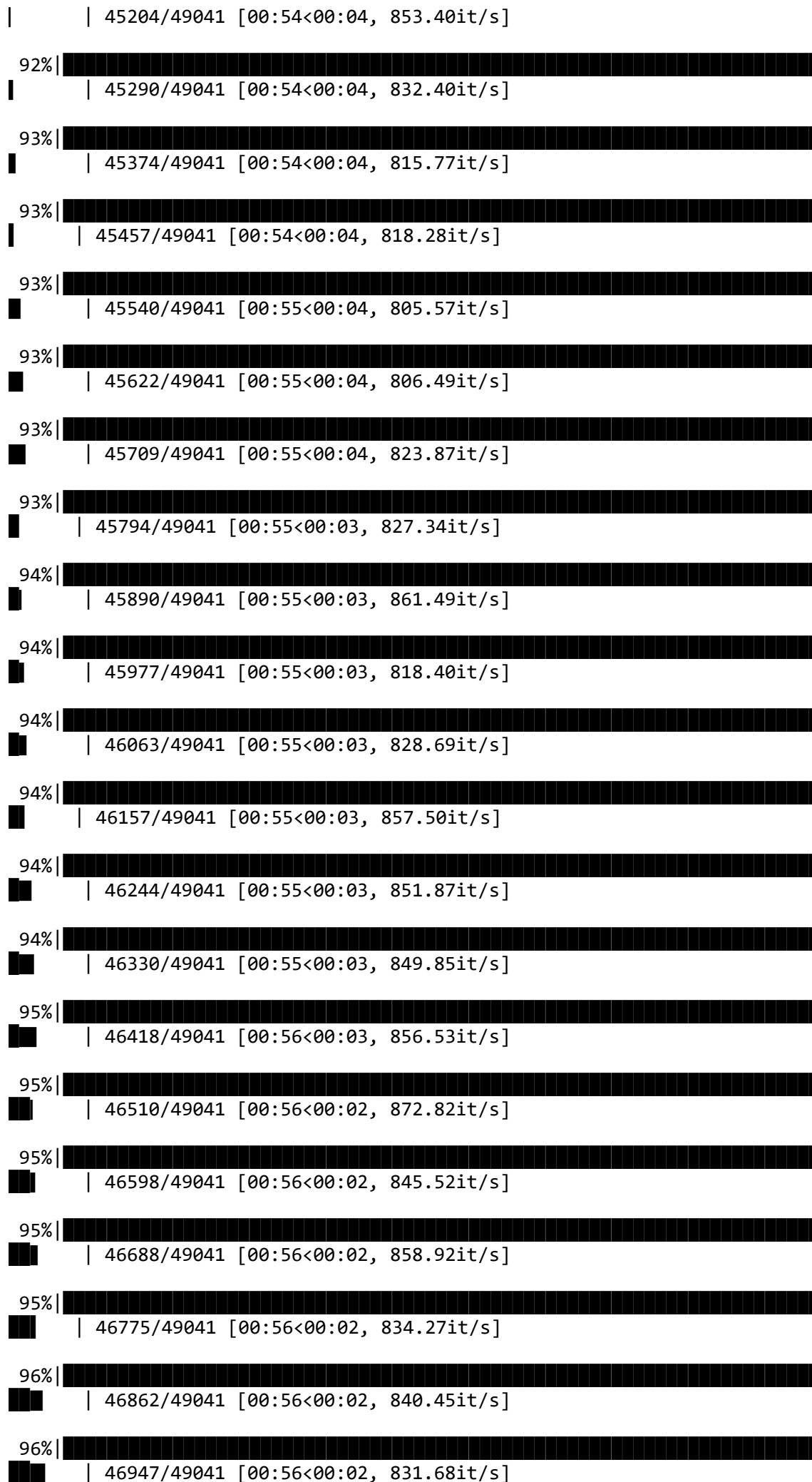
| 41515/49041 [00:50<00:08, 871.52it/s]

85%|

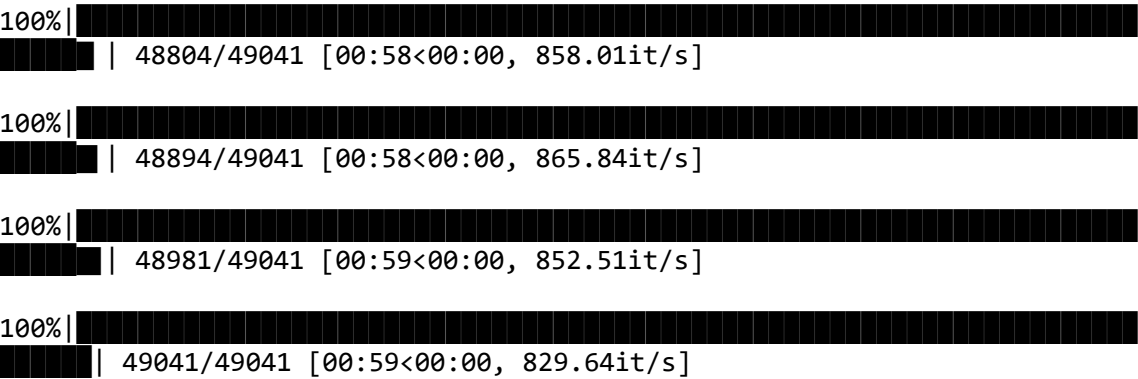
| 41607/49041 [00:50<00:08, 881.13it/s]



92% |



[illegible]



49041  
300

Essay test data

In [163]:

```
# tfidf Word2Vec
# computing average word2vec for each Project Title is stored in this list

tfidf_w2v_vectors_text_test = []; # the tfidf-w2v for each essay
for sentence in tqdm(preprocessed_essays_Test): # for each 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 title
    for word in sentence.split(): # for each word in a title
        if (word in glove_words) and (word in dictionary):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sen
            tence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # ge
            tting 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_text_test.append(vector)

print(len(tfidf_w2v_vectors_text_test))
print(len(tfidf_w2v_vectors_text_test[0]))
```

```

0%|
| 0/36052 [00:00<?, ?it/s]

0%||
| 94/36052 [00:00<00:38, 933.02it/s]

1%|█
| 181/36052 [00:00<00:39, 911.08it/s]

1%|█
| 273/36052 [00:00<00:39, 911.75it/s]

1%|██
| 368/36052 [00:00<00:38, 920.94it/s]

1%|██
| 461/36052 [00:00<00:38, 921.63it/s]

2%|███
| 545/36052 [00:00<00:39, 893.43it/s]

2%|███
| 636/36052 [00:00<00:39, 896.48it/s]

2%|███
| 723/36052 [00:00<00:39, 883.62it/s]

2%|████
| 808/36052 [00:00<00:40, 872.46it/s]

2%|████
| 892/36052 [00:01<00:40, 860.45it/s]

3%|████
| 981/36052 [00:01<00:40, 867.17it/s]

3%|████
| 1066/36052 [00:01<00:41, 849.75it/s]

3%|█████
| 1158/36052 [00:01<00:40, 867.28it/s]

3%|█████
| 1246/36052 [00:01<00:39, 870.31it/s]

4%|█████
| 1333/36052 [00:01<00:40, 853.11it/s]

4%|█████
| 1424/36052 [00:01<00:39, 867.53it/s]

4%|█████
| 1511/36052 [00:01<00:40, 853.65it/s]

4%|█████
| 1597/36052 [00:01<00:40, 843.66it/s]

5%|█████
| 1687/36052 [00:01<00:40, 857.81it/s]

```



5%|██████  
| 1773/36052 [00:02<00:41, 834.22it/s]

5%|██████  
| 1857/36052 [00:02<00:41, 831.72it/s]

5%|██████  
| 1948/36052 [00:02<00:40, 852.00it/s]

6%|██████  
| 2034/36052 [00:02<00:39, 852.53it/s]

6%|██████  
| 2127/36052 [00:02<00:38, 872.53it/s]

6%|██████  
| 2215/36052 [00:02<00:38, 870.20it/s]

6%|██████  
| 2303/36052 [00:02<00:39, 856.00it/s]

7%|██████  
| 2389/36052 [00:02<00:40, 823.46it/s]

7%|██████  
| 2475/36052 [00:02<00:40, 833.80it/s]

7%|██████  
| 2566/36052 [00:02<00:39, 853.54it/s]

7%|██████  
| 2658/36052 [00:03<00:38, 868.19it/s]

8%|██████  
| 2746/36052 [00:03<00:39, 849.71it/s]

8%|██████  
| 2832/36052 [00:03<00:39, 848.48it/s]

8%|██████  
| 2918/36052 [00:03<00:39, 845.06it/s]

8%|██████  
| 3003/36052 [00:03<00:40, 822.67it/s]

9%|██████  
| 3088/36052 [00:03<00:39, 828.86it/s]

9%|██████  
| 3173/36052 [00:03<00:39, 833.28it/s]

9%|██████  
| 3261/36052 [00:03<00:38, 845.05it/s]

9%|██████  
| 3346/36052 [00:03<00:38, 840.71it/s]

10%|██████  
| 3431/36052 [00:04<00:38, 841.65it/s]

10%|██████

| 3517/36052 [00:04<00:38, 844.77it/s]

10%|██████

| 3602/36052 [00:04<00:38, 834.50it/s]

10%|██████

| 3695/36052 [00:04<00:37, 856.90it/s]

10%|██████

| 3785/36052 [00:04<00:37, 867.49it/s]

11%|██████

| 3872/36052 [00:04<00:37, 852.24it/s]

11%|██████

| 3958/36052 [00:04<00:37, 854.55it/s]

11%|██████

| 4045/36052 [00:04<00:37, 858.45it/s]

11%|██████

| 4137/36052 [00:04<00:36, 872.96it/s]

12%|██████

| 4225/36052 [00:04<00:36, 870.57it/s]

12%|██████

| 4313/36052 [00:05<00:36, 863.80it/s]

12%|██████

| 4404/36052 [00:05<00:36, 875.38it/s]

12%|██████

| 4492/36052 [00:05<00:36, 856.95it/s]

13%|██████

| 4578/36052 [00:05<00:37, 845.97it/s]

13%|██████

| 4663/36052 [00:05<00:37, 835.39it/s]

13%|██████

| 4747/36052 [00:05<00:37, 832.47it/s]

13%|██████

| 4831/36052 [00:05<00:37, 827.92it/s]

14%|██████

| 4923/36052 [00:05<00:36, 851.81it/s]

14%|██████

| 5017/36052 [00:05<00:35, 874.69it/s]

14%|██████

| 5108/36052 [00:05<00:35, 883.19it/s]

14%|██████

| 5197/36052 [00:06<00:36, 845.63it/s]

15%|██████

| 5283/36052 [00:06<00:36, 846.71it/s]

15%|██████████  
| 5370/36052 [00:06<00:36, 851.69it/s]

15%|██████████  
| 5458/36052 [00:06<00:35, 858.12it/s]

15%|██████████  
| 5545/36052 [00:06<00:35, 857.25it/s]

16%|██████████  
| 5632/36052 [00:06<00:35, 859.25it/s]

16%|██████████  
| 5719/36052 [00:06<00:35, 853.00it/s]

16%|██████████  
| 5809/36052 [00:06<00:34, 864.75it/s]

16%|██████████  
| 5897/36052 [00:06<00:34, 867.32it/s]

17%|██████████  
| 5984/36052 [00:06<00:35, 858.55it/s]

17%|██████████  
| 6070/36052 [00:07<00:35, 852.10it/s]

17%|██████████  
| 6156/36052 [00:07<00:35, 850.08it/s]

17%|██████████  
| 6243/36052 [00:07<00:34, 854.02it/s]

18%|██████████  
| 6329/36052 [00:07<00:34, 851.42it/s]

18%|██████████  
| 6419/36052 [00:07<00:34, 863.31it/s]

18%|██████████  
| 6506/36052 [00:07<00:34, 859.57it/s]

18%|██████████  
| 6593/36052 [00:07<00:34, 849.56it/s]

19%|██████████  
| 6679/36052 [00:07<00:34, 843.32it/s]

19%|██████████  
| 6769/36052 [00:07<00:34, 855.40it/s]

19%|██████████  
| 6863/36052 [00:07<00:33, 877.27it/s]

19%|██████████  
| 6955/36052 [00:08<00:32, 887.77it/s]

20%|██████████  
| 7044/36052 [00:08<00:33, 860.86it/s]

20%|████████████████████  
| 7135/36052 [00:08<00:33, 873.19it/s]

20%|████████████████████  
| 7223/36052 [00:08<00:33, 870.73it/s]

20%|████████████████████  
| 7313/36052 [00:08<00:32, 877.45it/s]

21%|████████████████████  
| 7405/36052 [00:08<00:32, 885.35it/s]

21%|████████████████████  
| 7495/36052 [00:08<00:32, 885.22it/s]

21%|████████████████████  
| 7584/36052 [00:08<00:32, 877.53it/s]

21%|████████████████████  
| 7672/36052 [00:08<00:32, 873.75it/s]

22%|████████████████████  
| 7761/36052 [00:09<00:32, 875.13it/s]

22%|████████████████████  
| 7849/36052 [00:09<00:33, 851.80it/s]

22%|████████████████████  
| 7940/36052 [00:09<00:32, 867.74it/s]

22%|████████████████████  
| 8027/36052 [00:09<00:32, 867.66it/s]

23%|████████████████████  
| 8114/36052 [00:09<00:32, 865.16it/s]

23%|████████████████████  
| 8201/36052 [00:09<00:32, 844.69it/s]

23%|████████████████████  
| 8286/36052 [00:09<00:32, 841.86it/s]

23%|████████████████████  
| 8371/36052 [00:09<00:33, 835.09it/s]

23%|████████████████████  
| 8458/36052 [00:09<00:32, 841.03it/s]

24%|████████████████████  
| 8543/36052 [00:09<00:33, 829.51it/s]

24%|████████████████████  
| 8627/36052 [00:10<00:33, 828.37it/s]

24%|████████████████████  
| 8710/36052 [00:10<00:33, 822.16it/s]

24%|████████████████████  
| 8796/36052 [00:10<00:32, 826.60it/s]

25%|████████████████████

8879/36052	[00:10<00:33, 812.33it/s]
25%	<div><div></div></div>
8967/36052	[00:10<00:32, 827.46it/s]
25%	<div><div></div></div>
9054/36052	[00:10<00:32, 835.67it/s]
25%	<div><div></div></div>
9141/36052	[00:10<00:31, 841.44it/s]
26%	<div><div></div></div>
9232/36052	[00:10<00:31, 860.52it/s]
26%	<div><div></div></div>
9320/36052	[00:10<00:30, 863.78it/s]
26%	<div><div></div></div>
9407/36052	[00:10<00:31, 858.63it/s]
26%	<div><div></div></div>
9501/36052	[00:11<00:30, 877.33it/s]
27%	<div><div></div></div>
9589/36052	[00:11<00:30, 855.82it/s]
27%	<div><div></div></div>
9675/36052	[00:11<00:31, 850.07it/s]
27%	<div><div></div></div>
9763/36052	[00:11<00:30, 857.00it/s]
27%	<div><div></div></div>
9849/36052	[00:11<00:31, 828.94it/s]
28%	<div><div></div></div>
9933/36052	[00:11<00:31, 830.36it/s]
28%	<div><div></div></div>
10020/36052	[00:11<00:30, 840.08it/s]
28%	<div><div></div></div>
10112/36052	[00:11<00:30, 860.72it/s]
28%	<div><div></div></div>
10204/36052	[00:11<00:29, 875.85it/s]
29%	<div><div></div></div>
10292/36052	[00:12<00:29, 862.35it/s]
29%	<div><div></div></div>
10383/36052	[00:12<00:29, 872.65it/s]
29%	<div><div></div></div>
10471/36052	[00:12<00:29, 860.17it/s]
29%	<div><div></div></div>
10558/36052	[00:12<00:29, 859.74it/s]
30%	<div><div></div></div>
10645/36052	[00:12<00:30, 830.80it/s]

[illegible]

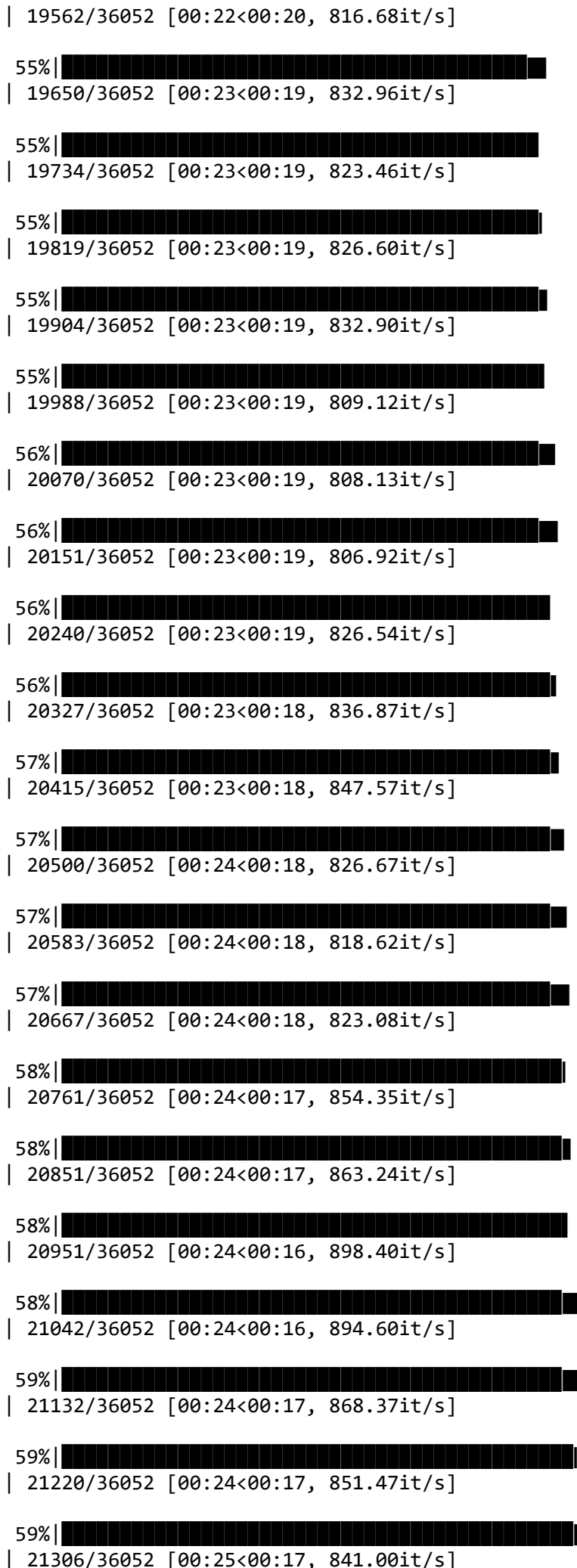
[illegible]

[illegible]



```
49%|███████████████████████████████████████████████████████████████████████████  
| 17753/36052 [00:20<00:21, 854.44it/s]
```

```
| 49%|███████████  
| 17839/36052 [00:20<00:21, 846.63it/s]  
  
50%|███████████  
| 17926/36052 [00:20<00:21, 851.66it/s]  
  
50%|███████████  
| 18012/36052 [00:21<00:21, 844.64it/s]  
  
50%|███████████  
| 18097/36052 [00:21<00:21, 829.61it/s]  
  
50%|███████████  
| 18181/36052 [00:21<00:21, 828.44it/s]  
  
51%|███████████  
| 18268/36052 [00:21<00:21, 837.49it/s]  
  
51%|███████████  
| 18357/36052 [00:21<00:20, 852.47it/s]  
  
51%|███████████  
| 18443/36052 [00:21<00:20, 848.91it/s]  
  
51%|███████████  
| 18528/36052 [00:21<00:20, 839.86it/s]  
  
52%|███████████  
| 18616/36052 [00:21<00:20, 849.65it/s]  
  
52%|███████████  
| 18702/36052 [00:21<00:20, 848.37it/s]  
  
52%|███████████  
| 18787/36052 [00:22<00:21, 820.17it/s]  
  
52%|███████████  
| 18870/36052 [00:22<00:20, 818.88it/s]  
  
53%|███████████  
| 18957/36052 [00:22<00:20, 831.62it/s]  
  
53%|███████████  
| 19043/36052 [00:22<00:20, 839.39it/s]  
  
53%|███████████  
| 19128/36052 [00:22<00:20, 822.65it/s]  
  
53%|███████████  
| 19220/36052 [00:22<00:19, 847.84it/s]  
  
54%|███████████  
| 19306/36052 [00:22<00:20, 813.61it/s]  
  
54%|███████████  
| 19392/36052 [00:22<00:20, 825.25it/s]  
  
54%|███████████  
| 19478/36052 [00:22<00:19, 833.53it/s]
```



```
| 21391/36052 [00:25<00:17, 836.89it/s]
```

```
| 21477/36052 [00:25<00:17, 841.88it/s]
```

```
| 21562/36052 [00:25<00:17, 839.97it/s]
```

```
| 21647/36052 [00:25<00:17, 838.50it/s]
```

```
| 21736/36052 [00:25<00:16, 852.80it/s]
```

```
| 21822/36052 [00:25<00:16, 843.08it/s]
```

```
| 21909/36052 [00:25<00:16, 849.16it/s]
```

```
21994/36052 [00:25<00:16, 837.48it/s]
```

```
22080/36052 [00:25<00:16, 842.23it/s]
```

```
| 22165/36052 [00:26<00:16, 842.70it/s]
```

```
| 22254/36052 [00:26<00:16, 854.55it/s]
```

```
| 22340/36052 [00:26<00:16, 854.30it/s]
```

```
| 22431/36052 [00:26<00:15, 868.47it/s]
```

```
| 22518/36052 [00:26<00:16, 829.98it/s]
```

```
| 22608/36052 [00:26<00:15, 848.87it/s]
```

```
| 22694/36052 [00:26<00:15, 845.33it/s]
```

```
| 22782/36052 [00:26<00:15, 854.92it/s]
```

```
| 22870/36052 [00:26<00:15, 861.61it/s]
```

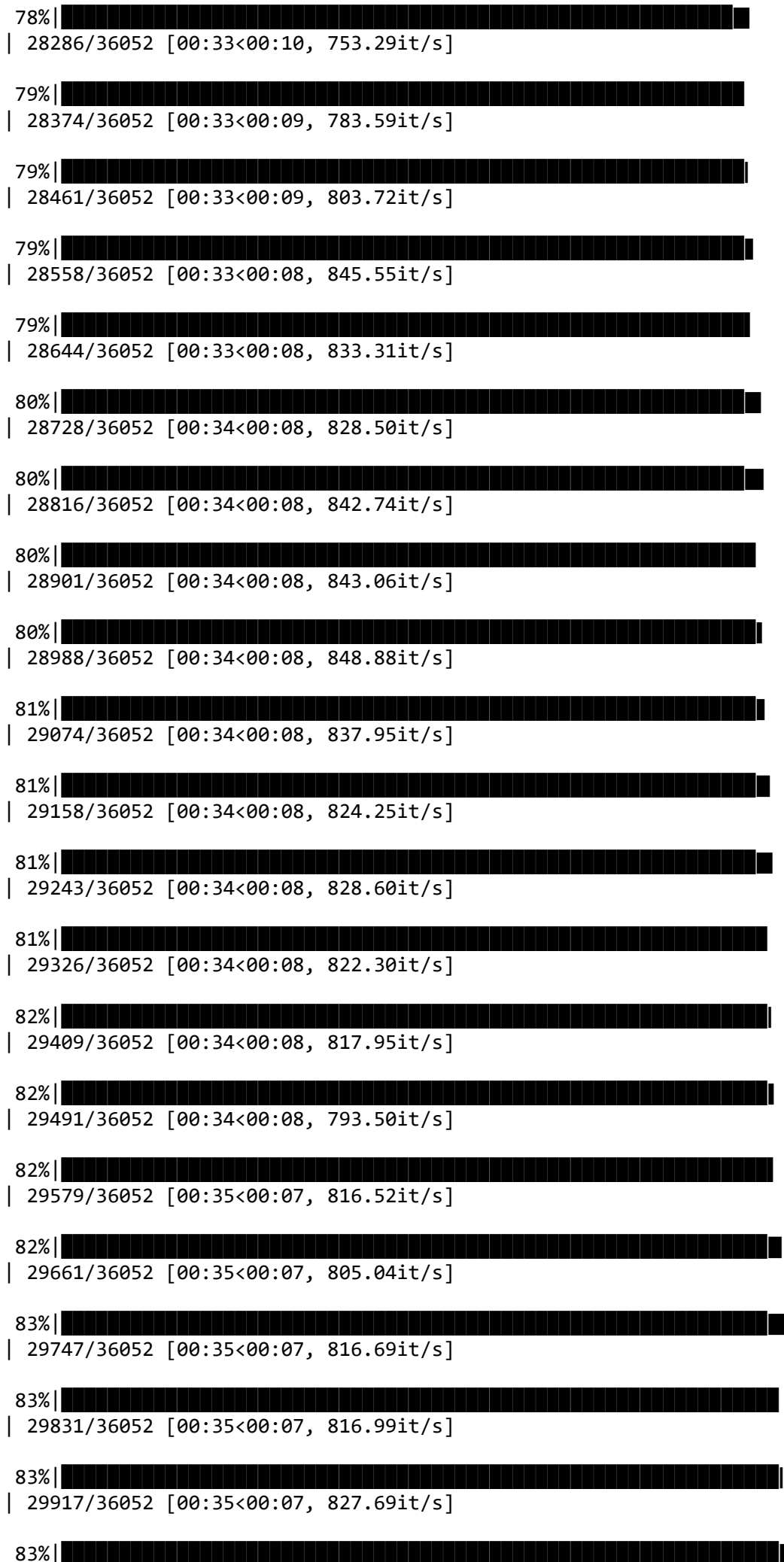
```
| 22957/36052 [00:26<00:15. 859.28it/s]
```

```
| 23044/36052 [00:27<00:15. 854.33it/s]
```

[illegible]

[illegible]

[illegible]





| 30004/36052 [00:35<00:07, 838.23it/s]

83%|

| 30088/36052 [00:35<00:07, 815.06it/s]

84%|

| 30175/36052 [00:35<00:07, 829.06it/s]

84%|

| 30259/36052 [00:35<00:07, 821.00it/s]

84%|

| 30345/36052 [00:36<00:06, 830.56it/s]

84%|

| 30429/36052 [00:36<00:06, 805.33it/s]

85%|

| 30514/36052 [00:36<00:06, 814.10it/s]

85%|

| 30596/36052 [00:36<00:06, 812.74it/s]

85%|

| 30681/36052 [00:36<00:06, 821.56it/s]

85%|

| 30765/36052 [00:36<00:06, 826.11it/s]

86%|

| 30856/36052 [00:36<00:06, 847.80it/s]

86%|

| 30941/36052 [00:36<00:06, 810.39it/s]

86%|

| 31029/36052 [00:36<00:06, 828.36it/s]

86%|

| 31114/36052 [00:36<00:05, 832.81it/s]

87%|

| 31200/36052 [00:37<00:05, 838.84it/s]

87%|

| 31288/36052 [00:37<00:05, 849.80it/s]

87%|

| 31374/36052 [00:37<00:05, 830.15it/s]

87%|

| 31458/36052 [00:37<00:05, 813.09it/s]

87%|

| 31544/36052 [00:37<00:05, 825.91it/s]

88%|

| 31632/36052 [00:37<00:05, 839.51it/s]

88%|

| 31720/36052 [00:37<00:05, 846.75it/s]

```
| 31806/36052 [00:37<00:05, 848.84it/s]
```

```
| 31891/36052 [00:37<00:05, 830.01it/s]
```

```
| 31980/36052 [00:37<00:04, 845.37it/s]
```

```
| 32065/36052 [00:38<00:05, 787.51it/s]
```

```
| 32151/36052 [00:38<00:04, 803.90it/s]
```

```
| 32239/36052 [00:38<00:04, 823.41it/s]
```

```
| 32331/36052 [00:38<00:04, 848.53it/s]
```

```
| 32419/36052 [00:38<00:04, 855.85it/s]
```

```
| 32506/36052 [00:38<00:04, 844.74it/s]
```

```
| 32594/36052 [00:38<00:04, 853.34it/s]
```

```
| 32680/36052 [00:38<00:03, 847.94it/s]
```

```
| 32771/36052 [00:38<00:03, 860.09it/s]
```

```
| 32858/36052 [00:39<00:03, 860.93it/s]
```

```
| 32945/36052 [00:39<00:03, 842.53it/s]
```

```
| 33030/36052 [00:39<00:03, 818.65it/s]
```

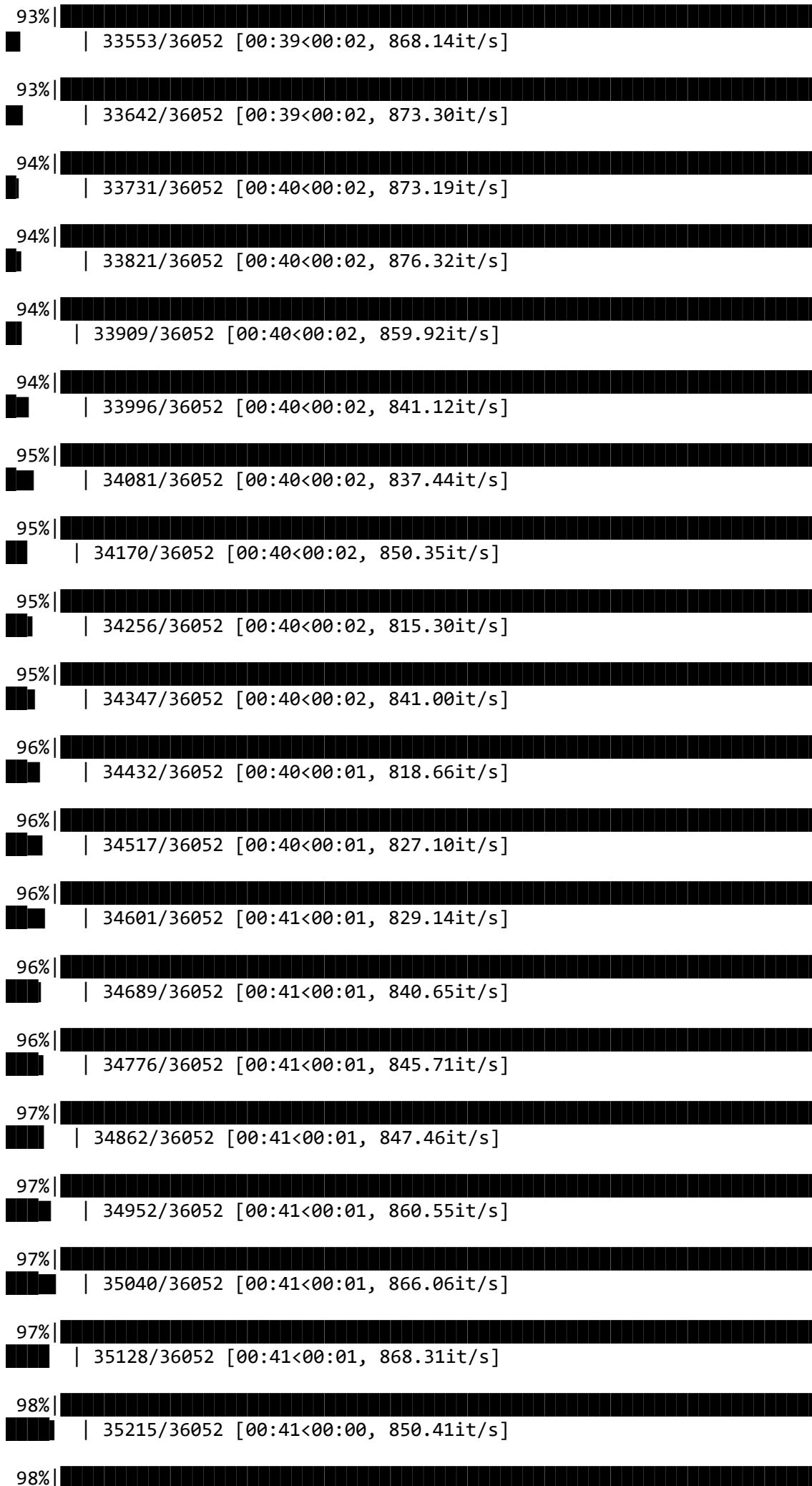
```
| 33113/36052 [00:39<00:03, 820.04it/s]
```

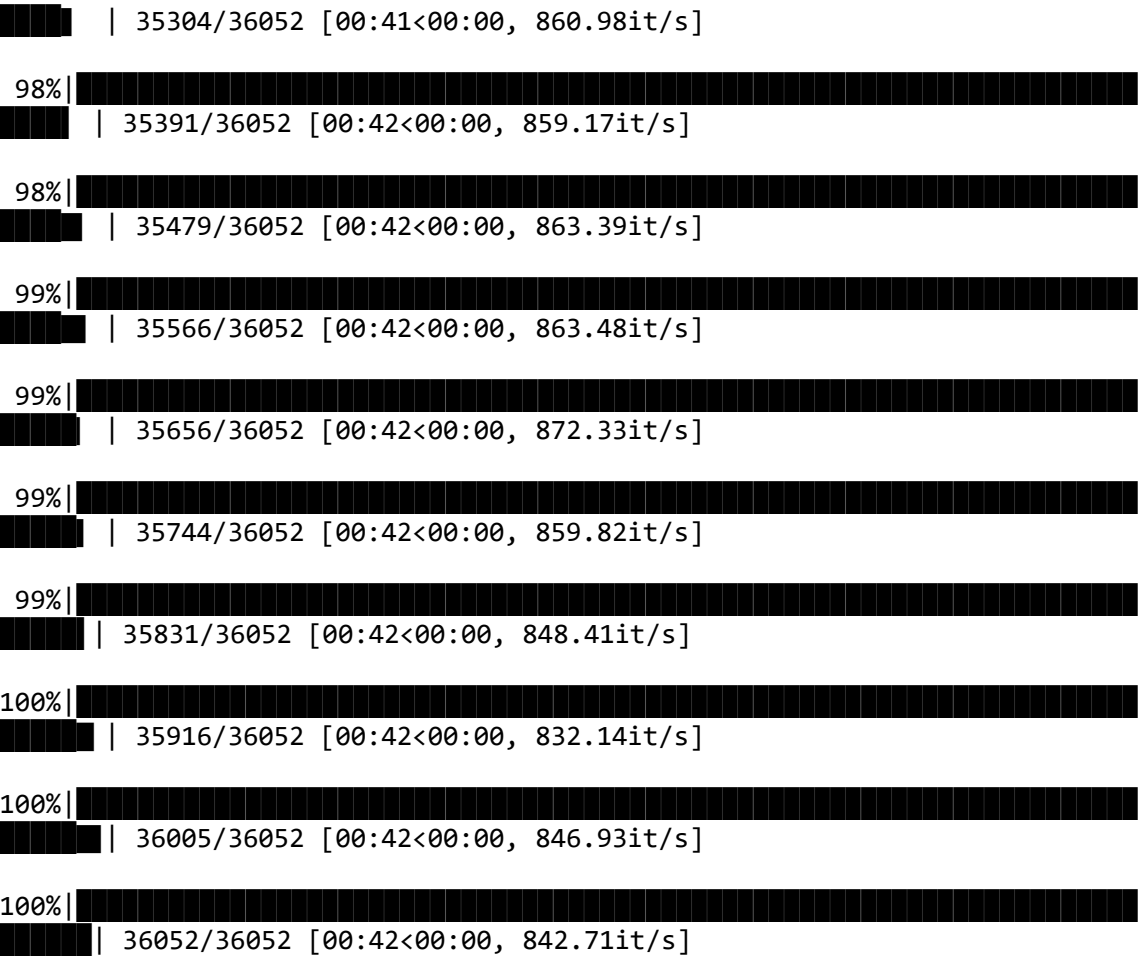
```
| 33196/36052 [00:39<00:03, 819.70it/s]
```

```
33288/36052 [00:39<00:03, 845.58it/s]
```

```
33373/36052 [00:39<00:03, 839.92it/s]
```

```
33461/36052 [00:39<00:03, 849.46it/s]
```





36052  
300

## Essay cross validation

In [164]:

```
# tfidf Word2Vec
# computing average word2vec for each Project Title is stored in this list

tfidf_w2v_vectors_text_cv = []; # the tfidf-w2v for each essay
for sentence in tqdm(preprocessed_essays_Cv): # for each 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 title
    for word in sentence.split(): # for each word in a title
        if (word in glove_words) and (word in dictionary):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sen
            tence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # ge
            tting 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_text_cv.append(vector)

print(len(tfidf_w2v_vectors_text_cv))
print(len(tfidf_w2v_vectors_text_cv[0]))
```

0%|  
| 0/24155 [00:00<?, ?it/s]

0%||  
| 82/24155 [00:00<00:29, 814.07it/s]

1%||  
| 171/24155 [00:00<00:28, 833.73it/s]

1%|■  
| 260/24155 [00:00<00:28, 848.08it/s]

1%|■  
| 342/24155 [00:00<00:28, 837.49it/s]

2%|■  
| 429/24155 [00:00<00:28, 845.26it/s]

2%|■  
| 509/24155 [00:00<00:28, 829.21it/s]

2%|■■  
| 583/24155 [00:00<00:29, 795.86it/s]

3%|■■  
| 657/24155 [00:00<00:30, 771.55it/s]

3%|■■  
| 743/24155 [00:00<00:29, 794.49it/s]

3%|■■■  
| 828/24155 [00:01<00:28, 808.69it/s]

4%|■■■  
| 907/24155 [00:01<00:29, 798.84it/s]

4%|■■■  
| 988/24155 [00:01<00:28, 800.42it/s]

4%|■■■■  
| 1080/24155 [00:01<00:27, 831.24it/s]

5%|■■■■  
| 1163/24155 [00:01<00:28, 812.06it/s]

5%|■■■■■  
| 1256/24155 [00:01<00:27, 842.44it/s]

6%|■■■■■  
| 1341/24155 [00:01<00:28, 795.65it/s]

6%|■■■■■  
| 1436/24155 [00:01<00:27, 834.80it/s]

6%|■■■■■  
| 1521/24155 [00:01<00:27, 830.15it/s]

7%|■■■■■  
| 1610/24155 [00:01<00:26, 845.55it/s]

7%|██████|  
| 1696/24155 [00:02<00:26, 838.11it/s]

7%|██████|  
| 1786/24155 [00:02<00:26, 851.49it/s]

8%|██████|  
| 1874/24155 [00:02<00:25, 858.02it/s]

8%|██████|  
| 1964/24155 [00:02<00:25, 865.92it/s]

8%|██████|  
| 2051/24155 [00:02<00:26, 849.96it/s]

9%|██████|  
| 2137/24155 [00:02<00:25, 851.15it/s]

9%|██████|  
| 2223/24155 [00:02<00:26, 841.95it/s]

10%|██████|  
| 2308/24155 [00:02<00:26, 830.20it/s]

10%|██████|  
| 2395/24155 [00:02<00:26, 835.21it/s]

10%|██████|  
| 2484/24155 [00:02<00:25, 849.08it/s]

11%|██████|  
| 2570/24155 [00:03<00:25, 835.71it/s]

11%|██████|  
| 2654/24155 [00:03<00:25, 827.71it/s]

11%|██████|  
| 2743/24155 [00:03<00:25, 841.39it/s]

12%|██████|  
| 2828/24155 [00:03<00:25, 834.70it/s]

12%|██████|  
| 2912/24155 [00:03<00:25, 831.92it/s]

12%|██████|  
| 2996/24155 [00:03<00:25, 832.52it/s]

13%|██████|  
| 3080/24155 [00:03<00:26, 806.60it/s]

13%|██████|  
| 3168/24155 [00:03<00:25, 823.29it/s]

13%|██████|  
| 3254/24155 [00:03<00:25, 832.21it/s]

14%|██████|  
| 3338/24155 [00:04<00:25, 818.15it/s]

14%|██████|

| 3420/24155 [00:04<00:25, 816.93it/s]

15%|██████████

| 3503/24155 [00:04<00:25, 819.09it/s]

15%|██████████

| 3588/24155 [00:04<00:24, 823.96it/s]

15%|██████████

| 3671/24155 [00:04<00:24, 823.90it/s]

16%|██████████

| 3754/24155 [00:04<00:24, 821.50it/s]

16%|██████████

| 3841/24155 [00:04<00:24, 833.72it/s]

16%|██████████

| 3925/24155 [00:04<00:24, 833.77it/s]

17%|██████████

| 4010/24155 [00:04<00:24, 836.76it/s]

17%|██████████

| 4099/24155 [00:04<00:23, 850.26it/s]

17%|██████████

| 4187/24155 [00:05<00:23, 857.14it/s]

18%|██████████

| 4280/24155 [00:05<00:22, 875.97it/s]

18%|██████████

| 4370/24155 [00:05<00:22, 878.57it/s]

18%|██████████

| 4458/24155 [00:05<00:22, 859.20it/s]

19%|██████████

| 4545/24155 [00:05<00:23, 847.93it/s]

19%|██████████

| 4630/24155 [00:05<00:23, 829.46it/s]

20%|██████████

| 4719/24155 [00:05<00:23, 844.98it/s]

20%|██████████

| 4806/24155 [00:05<00:22, 850.48it/s]

20%|██████████

| 4896/24155 [00:05<00:22, 862.96it/s]

21%|██████████

| 4988/24155 [00:05<00:21, 874.98it/s]

21%|██████████

| 5076/24155 [00:06<00:22, 846.87it/s]

21%|██████████

| 5165/24155 [00:06<00:22, 857.49it/s]



22%|████████████████████|  
| 5251/24155 [00:06<00:22, 843.80it/s]

22%|████████████████████|  
| 5336/24155 [00:06<00:22, 838.82it/s]

22%|████████████████████|  
| 5423/24155 [00:06<00:22, 846.13it/s]

23%|████████████████████|  
| 5513/24155 [00:06<00:21, 857.34it/s]

23%|████████████████████|  
| 5605/24155 [00:06<00:21, 873.40it/s]

24%|████████████████████|  
| 5702/24155 [00:06<00:20, 895.99it/s]

24%|████████████████████|  
| 5792/24155 [00:06<00:21, 869.36it/s]

24%|████████████████████|  
| 5880/24155 [00:07<00:21, 848.11it/s]

25%|████████████████████|  
| 5966/24155 [00:07<00:22, 825.33it/s]

25%|████████████████████|  
| 6054/24155 [00:07<00:21, 839.25it/s]

25%|████████████████████|  
| 6141/24155 [00:07<00:21, 846.50it/s]

26%|████████████████████|  
| 6231/24155 [00:07<00:20, 859.99it/s]

26%|████████████████████|  
| 6318/24155 [00:07<00:21, 843.62it/s]

27%|████████████████████|  
| 6407/24155 [00:07<00:20, 855.28it/s]

27%|████████████████████|  
| 6493/24155 [00:07<00:20, 849.70it/s]

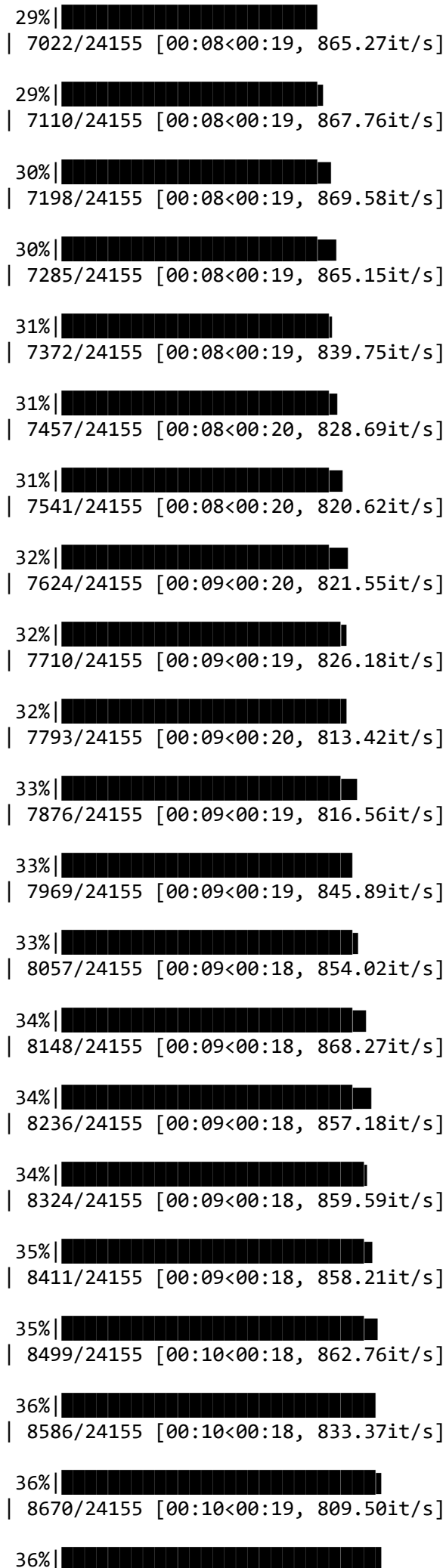
27%|████████████████████|  
| 6583/24155 [00:07<00:20, 859.91it/s]

28%|████████████████████|  
| 6670/24155 [00:07<00:20, 846.01it/s]

28%|████████████████████|  
| 6758/24155 [00:08<00:20, 854.10it/s]

28%|████████████████████|  
| 6844/24155 [00:08<00:20, 853.99it/s]

29%|████████████████████|  
| 6933/24155 [00:08<00:20, 857.69it/s]



43%|██████████  
| 10476/24155 [00:12<00:16, 844.50it/s]

```
51%|███████████████████████████████████████████████████████████████████████████████  
| 12234/24155 [00:14<00:14.851 39it/s]
```

[illegible]

```
| 15807/24155 [00:18<00:09, 841.99it/s]
```

[illegible]

```
| 17669/24155 [00:20<00:07, 849.39it/s]
```

```
| 17756/24155 [00:20<00:07, 853.64it/s]
```

```
| 17842/24155 [00:21<00:07, 851.14it/s]
```

```
| 17928/24155 [00:21<00:07, 846.89it/s]
```

```
| 18013/24155 [00:21<00:07, 833.57it/s]
```

```
| 18097/24155 [00:21<00:07, 823.96it/s]
```

```
| 18180/24155 [00:21<00:07, 821.46it/s]
```

```
| 18263/24155 [00:21<00:07, 805.49it/s]
```

```
18346/24155 [00:21<00:07, 810.96it/s]
```

```
18429/24155 [00:21<00:07, 812.50it/s]
```

```
18514/24155 [00:21<00:06, 821.59it/s]
```

```
| 18607/24155 [00:22<00:06, 849.67it/s]
```

```
| 18693/24155 [00:22<00:06, 850.88it/s]
```

```
| 18779/24155 [00:22<00:06. 846.72it/s]
```

```
| 18864/24155 [00:22<00:06. 840.85it/s]
```

```
| 18952/24155 [00:22<00:06.850.43it/s]
```

```
| 19038/24155 [00:22<00:06.846 40it/s]
```

19127/24155 [00:22<00:05 857 21it/s]

```
| 19216/24155 [00:22<00:05  862 44it/s]
```

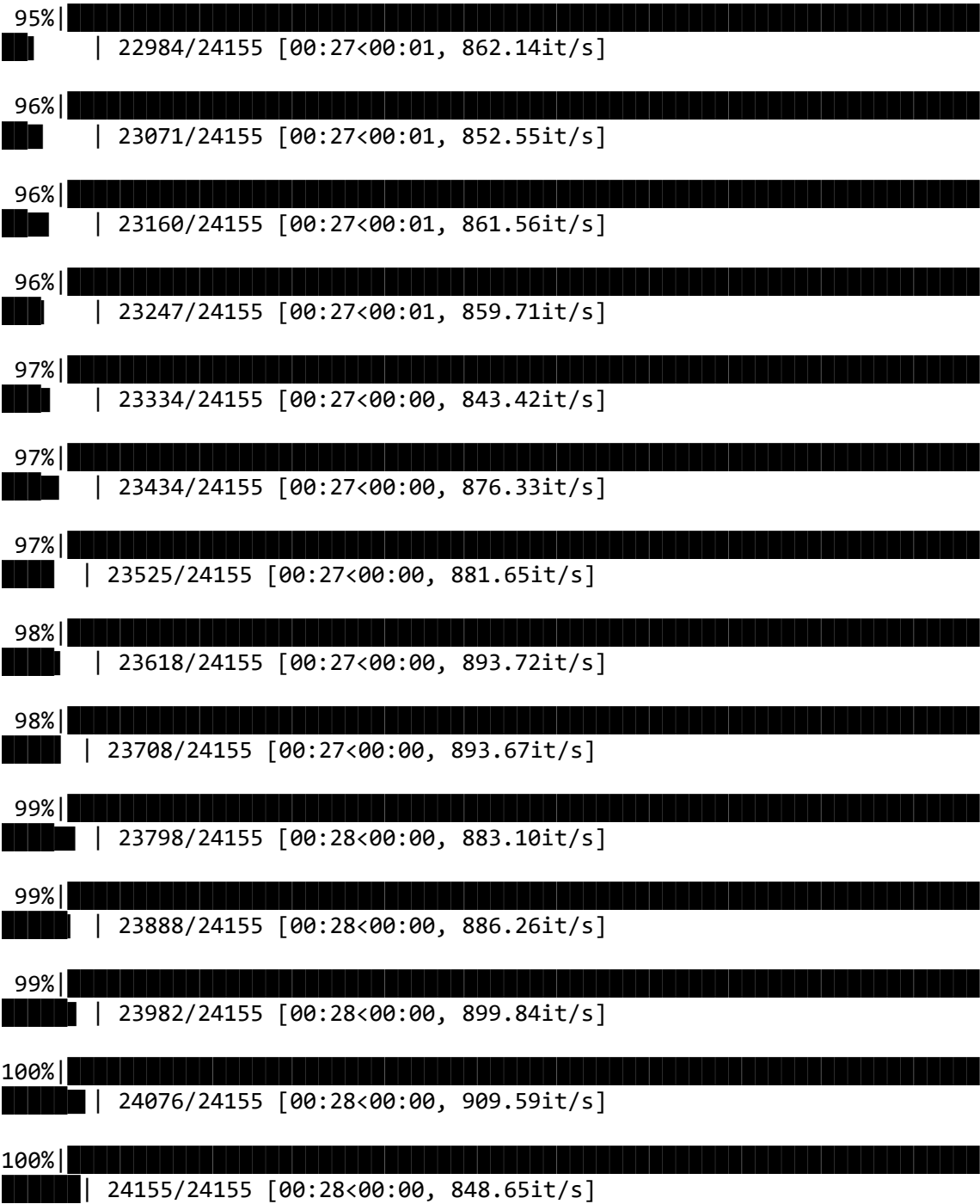
19303/24155 [00:22<00:05 855 27it/s]

80% |



```
21160/24155 [00:24<00:03, 845.76it/s]
```

[illegible]



24155  
300

# Project Title Train data

In [171]:

```
vectorizer = TfidfVectorizer(max_features = 2000, min_df = 10)
vectorizer.fit(preprocessed_project_titles_Train)
project_title_tfidf_train = vectorizer.transform(preprocessed_project_titles_Train)
print("Shape of matrix after tfidf ", project_title_tfidf_train.shape)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(vectorizer.get_feature_names(), list(vectorizer.idf_)))
tfidf_project_title_words = set(dictionary.keys())
```

Shape of matrix after tfidf (49041, 2000)

In [172]:

```
# tfidf Word2Vec
# computing average word2vec for each Project Title is stored in this list

tfidf_w2v_vectors_project_title_train = []; # the tfidf-w2v for each essay
for sentence in tqdm(preprocessed_project_titles_Train): # for each 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 title
    for word in sentence.split(): # for each word in a title
        if (word in glove_words) and (word in dictionary):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sen
            tence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # ge
            tting 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_project_title_train.append(vector)

print(len(tfidf_w2v_vectors_project_title_train))
print(len(tfidf_w2v_vectors_project_title_train[0]))
```

## Project title test data

In [173]:

```

# tfidf Word2Vec
# computing average word2vec for each Project Title is stored in this list

tfidf_w2v_vectors_project_title_test = []; # the tfidf-w2v for each essay
for sentence in tqdm(preprocessed_project_titles_Test): # for each 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 title
    for word in sentence.split(): # for each word in a title
        if (word in glove_words) and (word in dictionary):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sen
            tence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # ge
            tting 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_project_title_test.append(vector)

print(len(tfidf_w2v_vectors_project_title_test))
print(len(tfidf_w2v_vectors_project_title_test[0]))

```

```

0%|
| 0/36052 [00:00<?, ?it/s]

12%|██████████
| 4361/36052 [00:00<00:00, 43282.10it/s]

25%|██████████
| 8858/36052 [00:00<00:00, 43681.59it/s]

37%|██████████
| 13161/36052 [00:00<00:00, 43391.86it/s]

50%|██████████
| 17850/36052 [00:00<00:00, 44293.40it/s]

62%|██████████
| 22426/36052 [00:00<00:00, 44611.28it/s]

75%|██████████
| 27020/36052 [00:00<00:00, 44902.01it/s]

88%|██████████
| 31645/36052 [00:00<00:00, 45204.30it/s]

100%|██████████
| 36052/36052 [00:00<00:00, 44897.84it/s]

36052
300

```

## project title cross validation data

In [174]:

```
# tfidf Word2Vec
# computing average word2vec for each Project Title is stored in this list

tfidf_w2v_vectors_project_title_cv = []; # the tfidf-w2v for each essay
for sentence in tqdm(preprocessed_project_titles_Cv): # for each 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 title
    for word in sentence.split(): # for each word in a title
        if (word in glove_words) and (word in dictionary):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sen
            tence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # ge
            tting 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_project_title_cv.append(vector)

print(len(tfidf_w2v_vectors_project_title_cv))
print(len(tfidf_w2v_vectors_project_title_cv[0]))
```

```
0%|
| 0/24155 [00:00<?, ?it/s]

18%|██████████
| 4459/24155 [00:00<00:00, 44266.46it/s]

37%|██████████
| 8981/24155 [00:00<00:00, 44452.31it/s]

56%|██████████
| 13583/24155 [00:00<00:00, 44815.48it/s]

75%|██████████
| 18224/24155 [00:00<00:00, 45185.54it/s]

94%|██████████
| 22813/24155 [00:00<00:00, 45293.28it/s]

100%|██████████
| 24155/24155 [00:00<00:00, 45101.61it/s]

24155
300
```

## Numerical Features

### vectorizing numerical features



## Price for projects

In [100]:

```
# now firstly we will try to add the price and the quantity of the items required from  
the resource dataframe  
  
price_quantity_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'  
}).reset_index()  
  
price_quantity_data.head(2)
```

Out[100]:

	id	price	quantity
0	p000001	459.56	7
1	p000002	515.89	21

**now we need to join the above dataframe with our train,test,cv data that we already have**

In [101]:

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

**we will be performing the normalization of the numerical data here**

**Normalizing the price data**

In [102]:

```
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))

price_train = normalizer.transform(X_train['price'].values.reshape(-1,1))
price_cv = normalizer.transform(X_cv['price'].values.reshape(-1,1))
price_test = normalizer.transform(X_test['price'].values.reshape(-1,1))

print("After vectorizations")
print(price_train.shape, y_train.shape)
print(price_cv.shape, y_cv.shape)
print(price_test.shape, y_test.shape)
```

After vectorizations  
(49041, 1) (49041,)  
(24155, 1) (24155,)  
(36052, 1) (36052,)

### Normalizing the quantity data

In [103]:

```
from sklearn.preprocessing import Normalizer

normalizer = Normalizer()

# normalizer.fit(X_train['quantity'].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['quantity'].values.reshape(-1,1))

quantity_train = normalizer.transform(X_train['quantity'].values.reshape(-1,1))
quantity_test = normalizer.transform(X_test['quantity'].values.reshape(-1,1))
quantity_cv = normalizer.transform(X_cv['quantity'].values.reshape(-1,1))

print("After vectorizations")
print(quantity_train.shape, y_train.shape)
print(quantity_test.shape, y_test.shape)
print(quantity_cv.shape, y_cv.shape)
```

```
After vectorizations
(49041, 1) (49041,)
(36052, 1) (36052,)
(24155, 1) (24155,)
```

**Normalizing the number of previously posted projects by a teacher**

In [104]:

```
normalizer = Normalizer()

# normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].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))

prev_projects_train = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
prev_projects_cv = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
prev_projects_test = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

print("After converting into vectors form")
print(prev_projects_train.shape, y_train.shape)
print(prev_projects_cv.shape, y_cv.shape)
print(prev_projects_test.shape, y_test.shape)
```

After converting into vectors form

```
(49041, 1) (49041,)
(24155, 1) (24155,)
(36052, 1) (36052,)
```

## Set 1: categorical, numerical features + project\_title(BOW) + preprocessed\_essay (BOW)

Now we need to merge all the numerical vectors(categorical features,text features,numerical features) given above for set-1 which we created using different methods

In [105]:

```
from scipy.sparse import hstack

X_train_merge = hstack((categories_one_hot_train,subcategories_one_hot_train,teacher_prefix_one_hot_train,project_grade_categories_one_hot_train,school_state_categories_one_hot_train,essay_bow_train,project_title_bow_train,price_train,quantity_train,prev_projects_train)).tocsr()
X_test_merge = hstack((categories_one_hot_test,subcategories_one_hot_test,teacher_prefix_one_hot_test,project_grade_categories_one_hot_test,school_state_categories_one_hot_test,essay_bow_test,project_title_bow_test,price_test,quantity_test,prev_projects_test)).tocsr()
X_cv_merge = hstack((categories_one_hot_cv,subcategories_one_hot_cv,teacher_prefix_one_hot_cv,project_grade_categories_one_hot_cv,school_state_categories_one_hot_cv,essay_bow_cv,project_title_bow_cv,price_cv,quantity_cv,prev_projects_cv)).tocsr()
```

In [106]:

```
# this will be our finally created data matrix dimensions  
  
print(X_train_merge.shape, y_train.shape)  
print(X_test_merge.shape, y_test.shape)  
print(X_cv_merge.shape, y_cv.shape)
```

```
(49041, 52835) (49041,)  
(36052, 52835) (36052,)  
(24155, 52835) (24155,)
```

## Applying KNN

**Hyper parameter Tuning (USING THE BASIC FOR LOOP METHOD GRIDSEARCH-CV CAN ALSO BE USED HERE)**

In [108]:

```
import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

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

train_auc = []
cv_auc = []
a = []
b = []

K = [1, 5, 10, 15, 21, 31, 41, 51,]

for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i)
    neigh.fit(X_train_merge, y_train)

    y_train_pred = neigh.predict_proba(X_train_merge)[:,-1]
    y_cv_pred = neigh.predict_proba(X_cv_merge)[:,-1]

    # 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))
    a.append(y_train_pred)
    b.append(y_cv_pred)

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

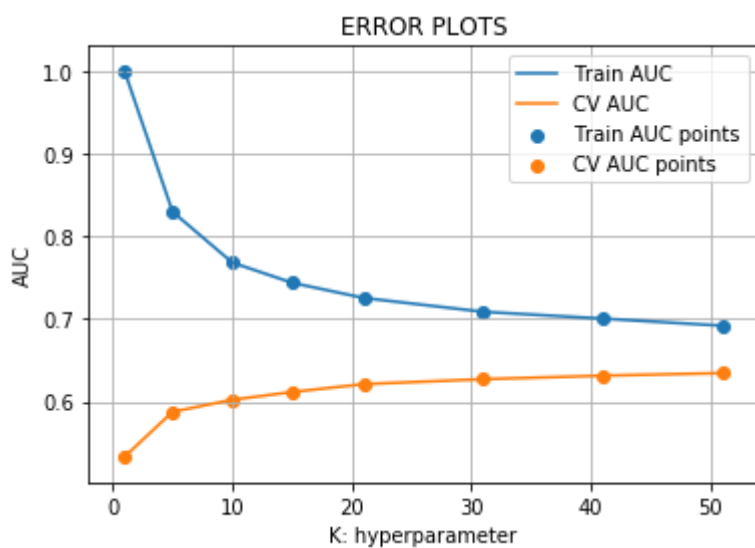
plt.scatter(K, train_auc, label='Train AUC points')
plt.scatter(K, cv_auc, label='CV AUC points')

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

```

0%|
| 0/8 [00:00<?, ?it/s]
12%|██████████
| 1/8 [04:18<30:08, 258.32s/it]
25%|██████████
| 2/8 [09:00<26:33, 265.62s/it]
38%|██████████
| 3/8 [13:41<22:30, 270.11s/it]
50%|██████████
| 4/8 [18:24<18:15, 273.83s/it]
62%|██████████
| 5/8 [23:06<13:49, 276.45s/it]
75%|██████████
| 6/8 [27:43<09:13, 276.68s/it]
88%|██████████
██████████ | 7/8 [32:31<04:39, 279.87s/it]
100%|██████████
██████████ | 8/8 [37:15<00:00, 281.31s/it]

```



FROM ABOVE PLOT WE CAN SAY THAT THE BEST K WILL BE AROUND 28 BECAUSE AFTER THAT THE GRAPH OF CV\_AUC IS ALMOST CONSTANT HENCE LET USE CONSIDER BEST K = 29

In [109]:

```
best_k_set_1 = 31
```

NOW USING THE BEST K VALUE RECEIVED FROM ABOVE WE WILL BE TRAINING OUR MODEL

In [110]:

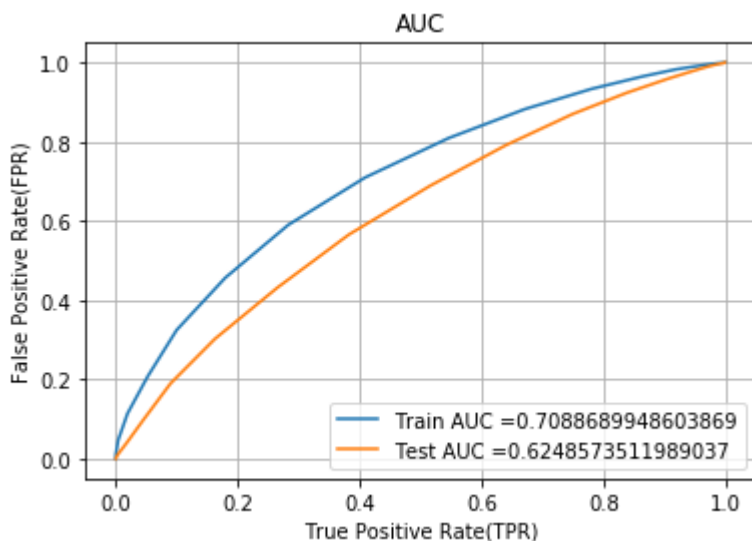
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k_set_1, n_jobs = -4)
neigh.fit(X_train_merge, 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 = neigh.predict_proba(X_train_merge)[: ,1]
y_test_pred = neigh.predict_proba(X_test_merge)[: ,1]

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

plt.plot(train_fpr, train_tpr, label="Train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



In [ ]:

**Creating the confusion matrix for the above train data results**



In [111]:

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

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

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

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

### Creating the confusion matrix for train and test data

In [112]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

```
=====
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24791603938644033 for threshold 0.742
[[ 3374  4052]
 [ 7946 33669]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.24975469520310625 for threshold 0.774
[[ 2644  2815]
 [ 9518 21075]]
```

### Showing the confusion matrix for train data visually

In [113]:

```
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

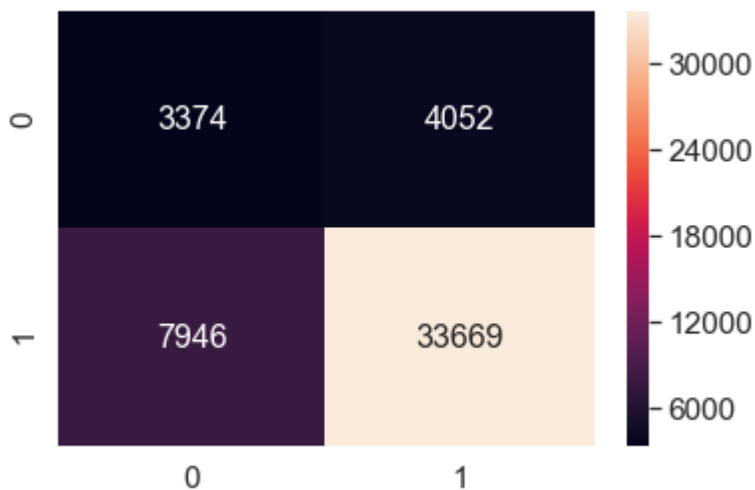
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt

df_cm = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)), range(2),
                     range(2))
plt.figure(figsize = (10,7))
sn.set(font_scale=1.4)#for label size
sn.heatmap(df_cm, annot=True,annot_kws={"size": 16},fmt='g')# font size
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.24791603938644033 for threshold 0.742

Out[113]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x22b949326d8>



**We can see that the True Positive values is the highest which is a good sign also. But the True Negative value is less which is not a good sign**

**Showing the confusion matrix for the test data visually**

In [114]:

```
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

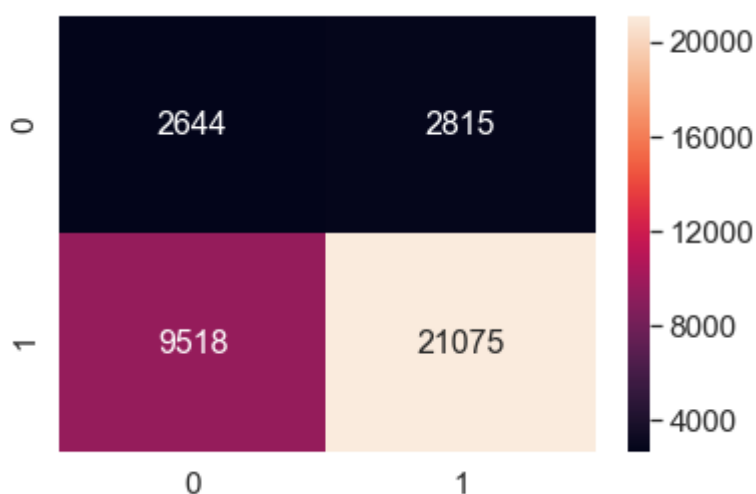
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt

df_cm_test = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)), range(2),
                           range(2))
#plt.figure(figsize = (10,7))
sn.set(font_scale=1.4)#for label size
sn.heatmap(df_cm_test, annot=True, annot_kws={"size": 16}, fmt='g')# font size
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.24975469520310625 for threshold 0.774

Out[114]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x22b94d93080>



From above we can say that we have large number of True positive but we also have a large number of False positive hence we can say that our model is a bit biased towards the positive class that's why it is predicting the positive class so many times

## Set 2 : categorical, numerical features + project\_title(TFIDF) + preprocessed\_essay(TFIDF)

Now we need to merge all the numerical vectors(categorical features, text features, numerical features) given above for set-2 which we created using different methods

In [123]:

```
from scipy.sparse import hstack

X_train_merge_set_2 = hstack((categories_one_hot_train,subcategories_one_hot_train,teacher_prefix_one_hot_train,project_grade_categories_one_hot_train,school_state_categories_one_hot_train,text_tfidf_train,project_title_tfidf_train,price_train,quantity_train,prev_projects_train)).tocsr()
X_test_merge_set_2 = hstack((categories_one_hot_test,subcategories_one_hot_test,teacher_prefix_one_hot_test,project_grade_categories_one_hot_test,school_state_categories_one_hot_test,text_tfidf_test,project_title_tfidf_test,price_test,quantity_test,prev_projects_test)).tocsr()
X_cv_merge_set_2 = hstack((categories_one_hot_cv,subcategories_one_hot_cv,teacher_prefix_one_hot_cv,project_grade_categories_one_hot_cv,school_state_categories_one_hot_cv,text_tfidf_cv,title_tfidf_cv,price_cv,quantity_cv,prev_projects_cv)).tocsr()
```

In [124]:

```
# this will be our finally created data matrix dimensions

print(X_train_merge_set_2.shape, y_train.shape)
print(X_test_merge_set_2.shape, y_test.shape)
print(X_cv_merge_set_2.shape, y_cv.shape)
```

```
(49041, 8102) (49041,)
(36052, 8102) (36052,)
(24155, 8102) (24155,)
```

## Hyper parameter Tuning (USING THE BASIC FOR LOOP METHOD GRIDSEARCH-CV CAN ALSO BE USED HERE)

In [125]:

```
import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

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

train_auc = []
cv_auc = []
a = []
b = []

K = [1, 5, 10, 15, 21, 31, 41, 51,]

for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i)
    neigh.fit(X_train_merge_set_2, y_train)

    y_train_pred = neigh.predict_proba(X_train_merge_set_2)[:,-1]
    y_cv_pred = neigh.predict_proba(X_cv_merge_set_2)[:,-1]

    # 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))
    a.append(y_train_pred)
    b.append(y_cv_pred)

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

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

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

```

0%|
| 0/8 [00:00<?, ?it/s]

12%|███████
| 1/8 [04:18<30:07, 258.16s/it]

25%|██████████
| 2/8 [08:57<26:26, 264.47s/it]

38%|██████████████
| 3/8 [13:40<22:29, 269.96s/it]

50%|██████████████████
| 4/8 [18:17<18:08, 272.11s/it]

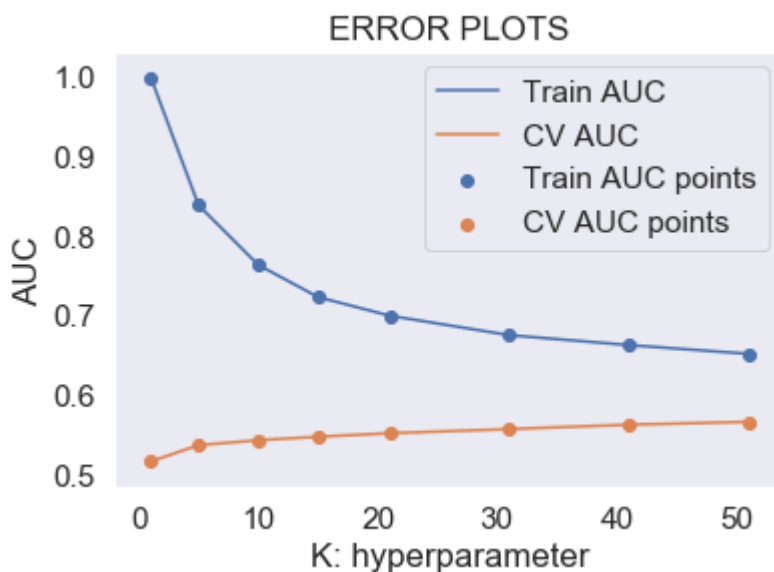
62%|██████████████████████
| 5/8 [22:50<13:36, 272.33s/it]

75%|██████████████████████████
| 6/8 [27:20<09:03, 271.69s/it]

88%|██████████████████████████████
| 7/8 [31:52<04:31, 271.70s/it]

100%|██████████████████████████████████
| 8/8 [36:22<00:00, 271.38s/it]

```



**FROM ABOVE PLOT WE CAN SAY THAT THE BEST K WILL BE AROUND 35 BECAUSE AFTER THAT THE GRAPH OF CV\_AUC IS ALMOST CONSTANT HENCE LET USE CONSIDER BEST K = 35**

In [126]:

```
best_k_set_2 = 41
```

## NOW USING THE BEST K VALUE RECEIVED FROM ABOVE WE WILL BE TRAINING OUR MODEL

In [127]:

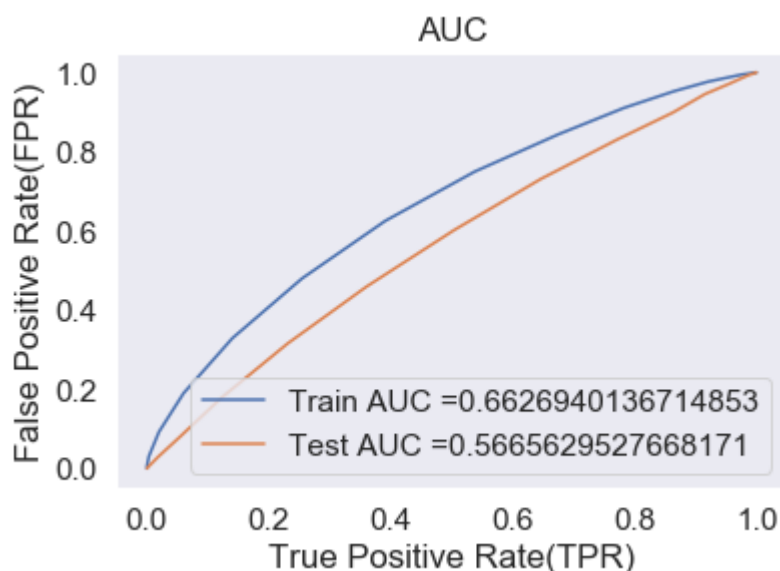
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k_set_2)
neigh.fit(X_train_merge_set_2, 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 = neigh.predict_proba(X_train_merge_set_2)[:,-1]
y_test_pred = neigh.predict_proba(X_test_merge_set_2)[:,-1]

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

plt.plot(train_fpr, train_tpr, label="Train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



**Creating the confusion matrix for the above data results**

In [128]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

```
=====
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24856813276940434 for threshold 0.829
[[ 3432  3994]
 [10419 31196]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.24998589797327114 for threshold 0.854
[[ 2709  2750]
 [12186 18407]]
```

## Showing the train confusion matrix visually



In [129]:

```
# Code for this segment from here -->> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

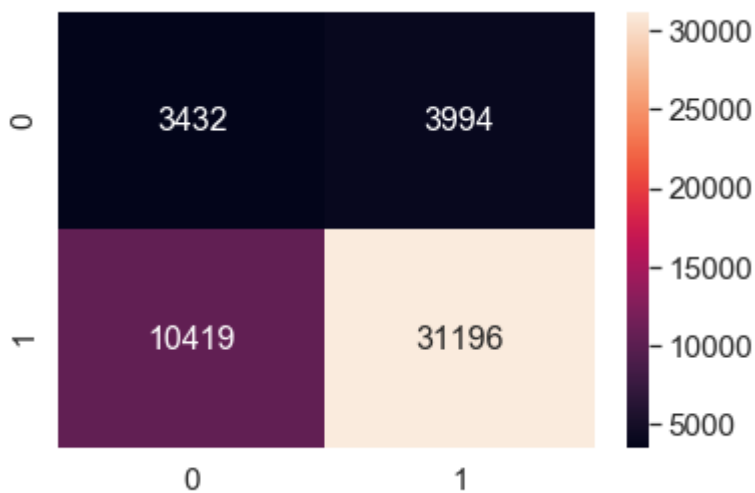
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt

df_cm_train = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)), range(2),
                           range(2))
plt.figure(figsize = (10,7))
sn.set(font_scale=1.4)#for label size
sn.heatmap(df_cm_train, annot=True,annot_kws={"size": 16},fmt = 'g')# font size
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.24856813276940434 for threshold 0.829

Out[129]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x22b949fa320>



**Showing the confusion matrix for the test data visually**

In [130]:

```
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

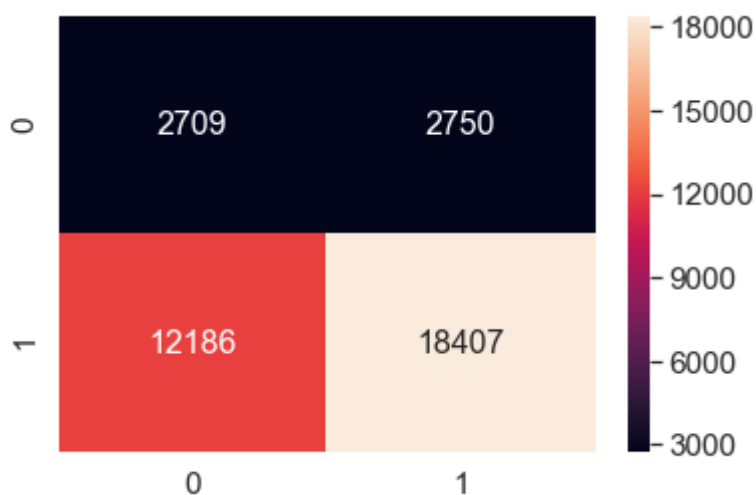
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt

df_cm_train = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds,
test_fpr, test_fpr)))
#plt.figure(figsize = (10,7))
sn.set(font_scale=1.4)#for label size
sn.heatmap(df_cm_train, annot=True,annot_kws={"size": 16},fmt = 'g')# font size
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.24998589797327114 for threshold 0.854

Out[130]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x22b9454ad68>



Here we can see that false positives and true positives are high in number it might also link towards the fact that our model is biased towards the positive class.

In [ ]:

**Set 3 : categorical, numerical features +  
project\_title(AVG W2V) + preprocessed\_essay (AVG  
W2V)**

In [131]:

```
from scipy.sparse import hstack

X_train_merge_set_3 = hstack((categories_one_hot_train, subcategories_one_hot_train, teacher_prefix_one_hot_train, project_grade_categories_one_hot_train, school_state_categories_one_hot_train, avg_w2v_vectors_train, avg_w2v_vectors_project_title_train, price_train, quantity_train, prev_projects_train)).tocsr()
X_test_merge_set_3 = hstack((categories_one_hot_test, subcategories_one_hot_test, teacher_prefix_one_hot_test, project_grade_categories_one_hot_test, school_state_categories_one_hot_test, avg_w2v_vectors_test, avg_w2v_vectors_project_title_test, price_test, quantity_test, prev_projects_test)).tocsr()
X_cv_merge_set_3 = hstack((categories_one_hot_cv, subcategories_one_hot_cv, teacher_prefix_one_hot_cv, project_grade_categories_one_hot_cv, school_state_categories_one_hot_cv, avg_w2v_vectors_cv, avg_w2v_vectors_project_title_cv, price_cv, quantity_cv, prev_projects_cv)).tocsr()
```

In [132]:

```
# this will be our finally created data matrix dimensions
```

```
print(X_train_merge_set_3.shape, y_train.shape)
print(X_test_merge_set_3.shape, y_test.shape)
print(X_cv_merge_set_3.shape, y_cv.shape)
```

```
(49041, 702) (49041,)
(36052, 702) (36052,)
(24155, 702) (24155,)
```

## Hyper parameter Tuning (USING THE BASIC FOR LOOP METHOD GRIDSEARCH-CV CAN ALSO BE USED HERE)

In [133]:

```
import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

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

train_auc = []
cv_auc = []
a = []
b = []

K = [1, 5, 10, 15, 21, 31, 41, 51,]

for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i)
    neigh.fit(X_train_merge_set_3, y_train)

    y_train_pred = neigh.predict_proba(X_train_merge_set_3)[:,-1]
    y_cv_pred = neigh.predict_proba(X_cv_merge_set_3)[:,-1]

    # 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))
    a.append(y_train_pred)
    b.append(y_cv_pred)

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

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

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

```

0%|
| 0/8 [00:00<?, ?it/s]

12%|███████
| 1/8 [1:02:21<7:16:27, 3741.02s/it]

25%|██████████
| 2/8 [2:05:53<6:16:15, 3762.60s/it]

38%|██████████████
| 3/8 [3:10:13<5:15:57, 3791.53s/it]

50%|██████████████████
| 4/8 [4:13:08<4:12:27, 3786.78s/it]

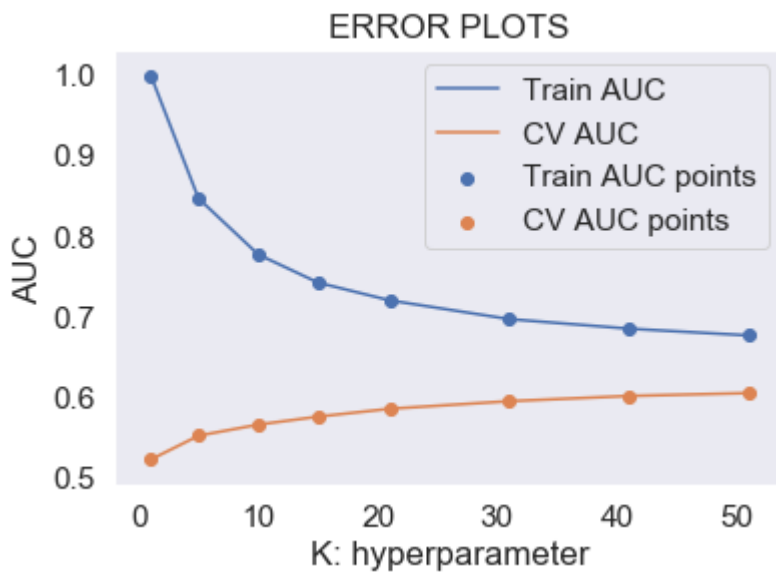
62%|██████████████████████
| 5/8 [10:16:03<7:39:09, 9183.31s/it]

75%|██████████████████████████
| 6/8 [16:44:54<7:27:34, 13427.49s/it]

88%|██████████████████████████████
| 7/8 [18:23:57<3:06:21, 11182.00s/it]

100%|██████████████████████████████████
| 8/8 [19:36:54<00:00, 9140.60s/it]

```



**FROM ABOVE PLOT WE CAN SAY THAT THE BEST K WILL BE AROUND 35 BECAUSE AFTER THAT THE GRAPH OF CV\_AUC IS ALMOST CONSTANT HENCE LET USE CONSIDER BEST K = 35**

In [134]:

```
best_k_set_3 = 35
```

## NOW USING THE BEST K VALUE RECEIVED FROM ABOVE WE WILL BE TRAINING OUR MODEL

In [135]:

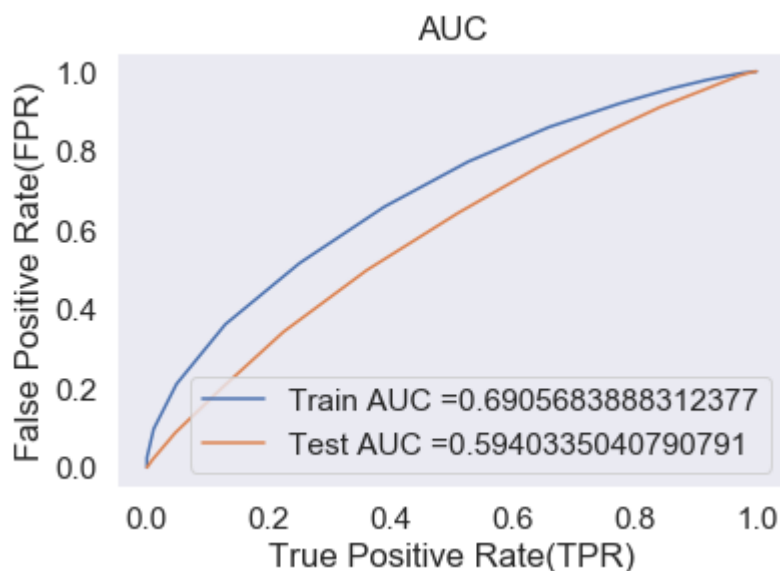
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k_set_3)
neigh.fit(X_train_merge_set_3, 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 = neigh.predict_proba(X_train_merge_set_3)[:,-1]
y_test_pred = neigh.predict_proba(X_test_merge_set_3)[:,-1]

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

plt.plot(train_fpr, train_tpr, label="Train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



**Creating the confusion matrix for the above data results**

In [136]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

```
=====
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24925272201839369 for threshold 0.829
[[ 3510  3916]
 [ 9472 32143]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.24990032945891405 for threshold 0.857
[[ 2675  2784]
 [10949 19644]]
```

## Showing the train confusion matrix visually

In [140]:

```
# Code for this segment from here -->> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
```

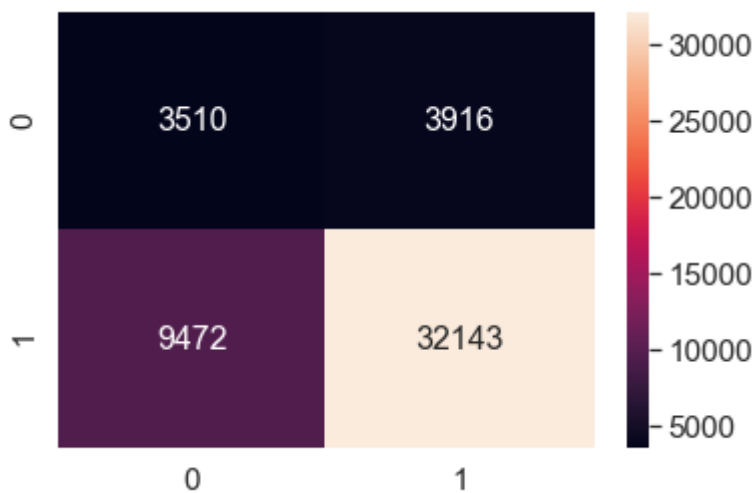
```
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
```

```
df_cm_train = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)), range(2),
                           range(2))
plt.figure(figsize = (10,7))
sn.set(font_scale=1.4)#for label size
sn.heatmap(df_cm_train, annot=True,annot_kws={"size": 16},fmt = 'g')# font size
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.24925272201839369 for threshold 0.829

Out[140]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x22b945ce5f8>



**We can see above clearly that the number of False positives are more than the number of False Negatives**

**Showing the test confusion matrix visually**



In [141]:

```
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

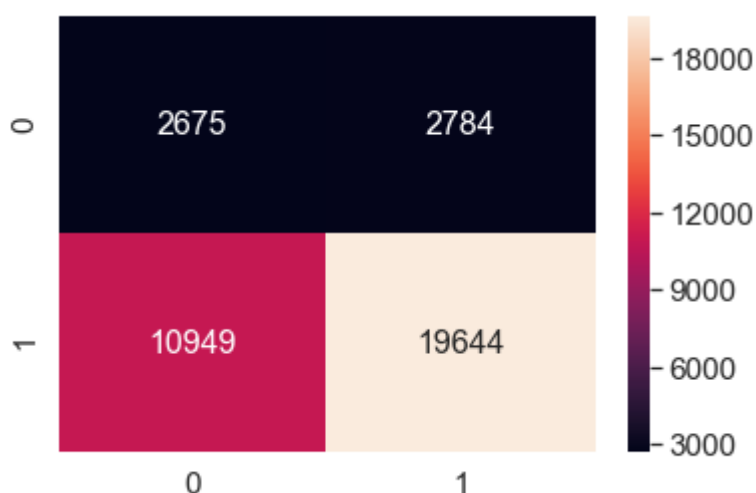
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt

df_cm_train = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds,
test_fpr, test_fpr)))
#plt.figure(figsize = (10,7))
sn.set(font_scale=1.4)#for label size
sn.heatmap(df_cm_train, annot=True,annot_kws={"size": 16},fmt = 'g')# font size
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.24990032945891405 for threshold 0.857

Out[141]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x22be57c8860>



Here we can see that we are getting a lot of False positives and true positives which directly indicate that our model is biased towards the positive class.

## SET-4 Applying KNN brute force on TFIDF W2V

categorical, numerical features + project\_title(TFIDF W2V) + preprocessed\_essay (TFIDF W2V)

In [185]:

```
from scipy.sparse import hstack

X_train_merge_set_4 = hstack((categories_one_hot_train,subcategories_one_hot_train,teacher_prefix_one_hot_train,project_grade_categories_one_hot_train,school_state_categories_one_hot_train,tfidf_w2v_vectors_text_train,tfidf_w2v_vectors_project_title_train,price_train,quantity_train,prev_projects_train)).tocsr()
X_test_merge_set_4 = hstack((categories_one_hot_test,subcategories_one_hot_test,teacher_prefix_one_hot_test,project_grade_categories_one_hot_test,school_state_categories_one_hot_test,tfidf_w2v_vectors_text_test,tfidf_w2v_vectors_project_title_test,price_test,quantity_test,prev_projects_test)).tocsr()
X_cv_merge_set_4 = hstack((categories_one_hot_cv,subcategories_one_hot_cv,teacher_prefix_one_hot_cv,project_grade_categories_one_hot_cv,school_state_categories_one_hot_cv,tfidf_w2v_vectors_text_cv,tfidf_w2v_vectors_project_title_cv,price_cv,quantity_cv,prev_projects_cv)).tocsr()
```

In [186]:

```
# this will be our finally created data matrix dimensions

print(X_train_merge_set_4.shape, y_train.shape)
print(X_test_merge_set_4.shape, y_test.shape)
print(X_cv_merge_set_4.shape,y_cv.shape)
```

```
(49041, 702) (49041,)
(36052, 702) (36052,)
(24155, 702) (24155,)
```

## Hyper parameter Tuning (USING THE BASIC FOR LOOP METHOD GRIDSEARCH-CV CAN ALSO BE USED HERE)

In [187]:

```
import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

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

train_auc = []
cv_auc = []
a = []
b = []

K = [1, 5, 10, 15, 21, 31, 41, 51,]

for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i)
    neigh.fit(X_train_merge_set_4, y_train)

    y_train_pred = neigh.predict_proba(X_train_merge_set_4)[:,-1]
    y_cv_pred = neigh.predict_proba(X_cv_merge_set_4)[:,-1]

    # 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))
    a.append(y_train_pred)
    b.append(y_cv_pred)

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

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

plt.legend()

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

```

0%|
| 0/8 [00:00<?, ?it/s]

12%|███████
| 1/8 [1:02:49<7:19:45, 3769.32s/it]

25%|██████████
| 2/8 [2:06:29<6:18:28, 3784.70s/it]

38%|██████████████
| 3/8 [3:10:13<5:16:22, 3796.41s/it]

50%|██████████████████
| 4/8 [4:11:43<4:10:57, 3764.38s/it]

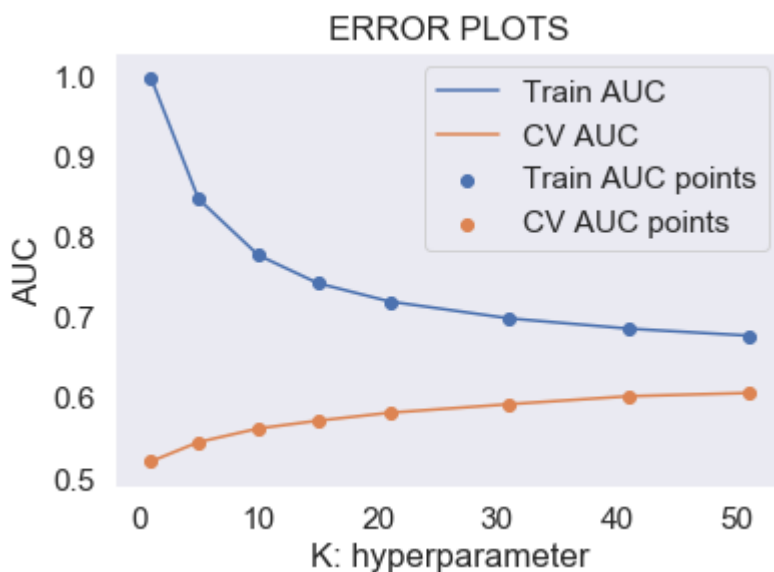
62%|██████████████████████
| 5/8 [5:12:47<3:06:42, 3734.19s/it]

75%|██████████████████████████
| 6/8 [6:16:34<2:05:24, 3762.06s/it]

88%|██████████████████████████████
| 7/8 [7:20:15<1:02:59, 3779.81s/it]

100%|██████████████████████████████████
| 8/8 [8:22:27<00:00, 3765.44s/it]

```



**FROM ABOVE PLOT WE CAN SAY THAT THE BEST K WILL BE AROUND 35 BECAUSE AFTER THAT THE GRAPH OF CV\_AUC IS ALMOST CONSTANT HENCE LET USE CONSIDER BEST K = 35**

In [189]:

```
best_k_set_4 = 53
```

**NOW USING THE BEST K VALUE RECEIVED FROM ABOVE WE WILL BE TRAINING OUR MODEL**

In [190]:

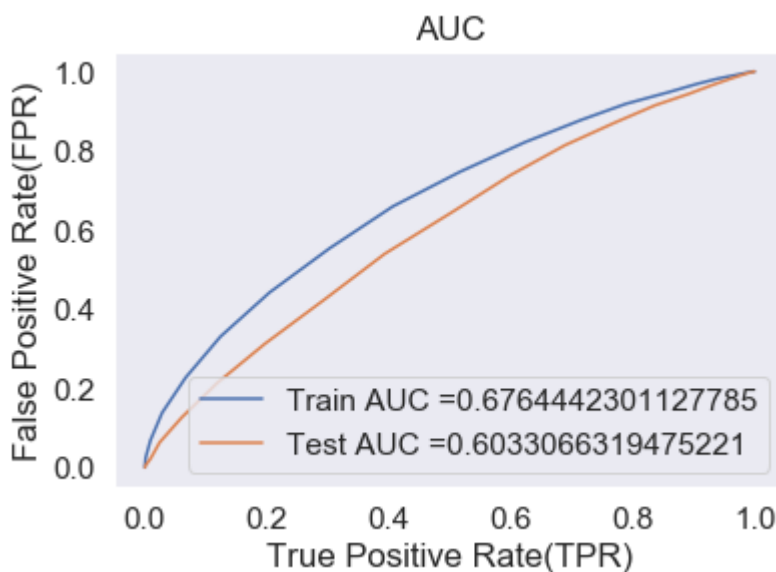
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k_set_4, n_jobs = -1)
neigh.fit(X_train_merge_set_4, 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 = neigh.predict_proba(X_train_merge_set_4)[:,1]
y_test_pred = neigh.predict_proba(X_test_merge_set_4)[:,1]

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

plt.plot(train_fpr, train_tpr, label="Train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



**Creating confusion matrix for the above results**

In [191]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

```
=====
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.2497121069369317 for threshold 0.83
[[ 3587  3839]
 [10545 31070]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.2499830121331791 for threshold 0.849
[[ 2707  2752]
 [10864 19729]]
```

## Showing the train confusion matrix visually

In [192]:

```
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
```

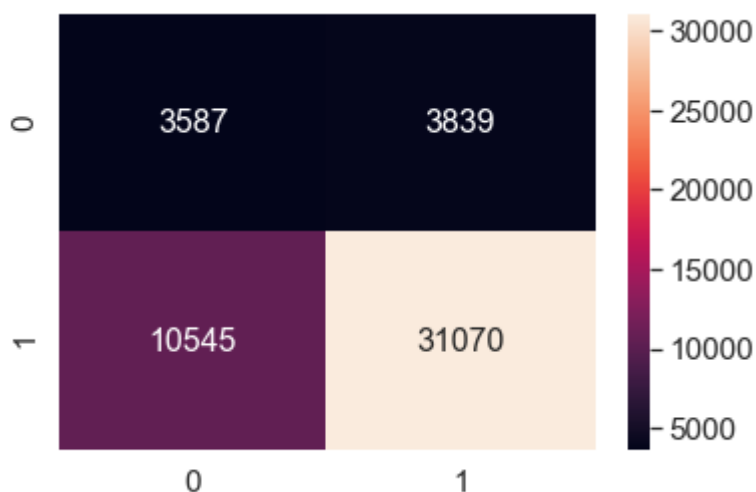
```
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
```

```
df_cm_train = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)), range(2),
                           range(2))
plt.figure(figsize = (10,7))
sn.set(font_scale=1.4)#for label size
sn.heatmap(df_cm_train, annot=True,annot_kws={"size": 16},fmt = 'g')# font size
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.2497121069369317 for threshold 0.83

Out[192]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x22b94440c50>



## Showing the test confusion matrix visually

In [193]:

```
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

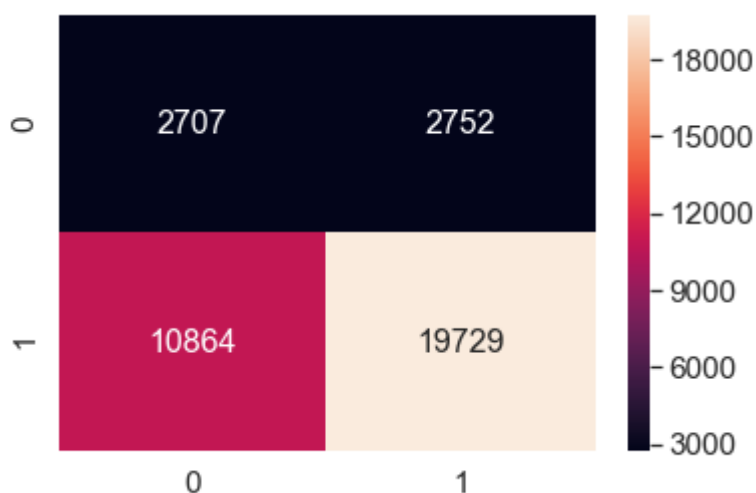
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt

df_cm_train = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds,
test_fpr, test_fpr)))
#plt.figure(figsize = (10,7))
sn.set(font_scale=1.4)#for label size
sn.heatmap(df_cm_train, annot=True,annot_kws={"size": 16},fmt = 'g')# font size
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.2499830121331791 for threshold 0.849

Out[193]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x22be57017b8>



The model has good amount of true positives which is a good sign. But a lot of negative class points are classified as positive class which is not a good sign.

## Select Best K features On Set\_2

In [195]:

```
print(X_train_merge_set_2.shape, y_train.shape)
print(X_test_merge_set_2.shape, y_test.shape)
print(X_cv_merge_set_2.shape, y_cv.shape)
```

```
(49041, 8102) (49041,)
(36052, 8102) (36052,)
(24155, 8102) (24155,)
```

Now above we have 8102 features out of which we want to select the best 2000 features



In [196]:

```
from sklearn.datasets import load_digits
from sklearn.feature_selection import SelectKBest, chi2
```

In [199]:

```
X_train_merge_set_2_best = SelectKBest(chi2, k=2000).fit_transform(X_train_merge_set_2,
y_train)
X_test_merge_set_2_best = SelectKBest(chi2, k=2000).fit_transform(X_test_merge_set_2,y_
test)
X_cv_merge_set_2_best = SelectKBest(chi2, k=2000).fit_transform(X_cv_merge_set_2,y_cv)
```

In [200]:

```
print(X_train_merge_set_2_best.shape)
print(X_test_merge_set_2_best.shape)
print(X_cv_merge_set_2_best.shape)
```

(49041, 2000)

(36052, 2000)

(24155, 2000)

**Hyper parameter Tuning for the best 2000 features  
(USING THE BASIC FOR LOOP METHOD  
GRIDSEARCH-CV CAN ALSO BE USED HERE)**

In [201]:

```
import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

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

train_auc = []
cv_auc = []
a = []
b = []

K = [1, 5, 10, 15, 21, 31, 41, 51]

for i in tqdm(K):
    neigh = KNeighborsClassifier(n_neighbors=i)
    neigh.fit(X_train_merge_set_2_best, y_train)

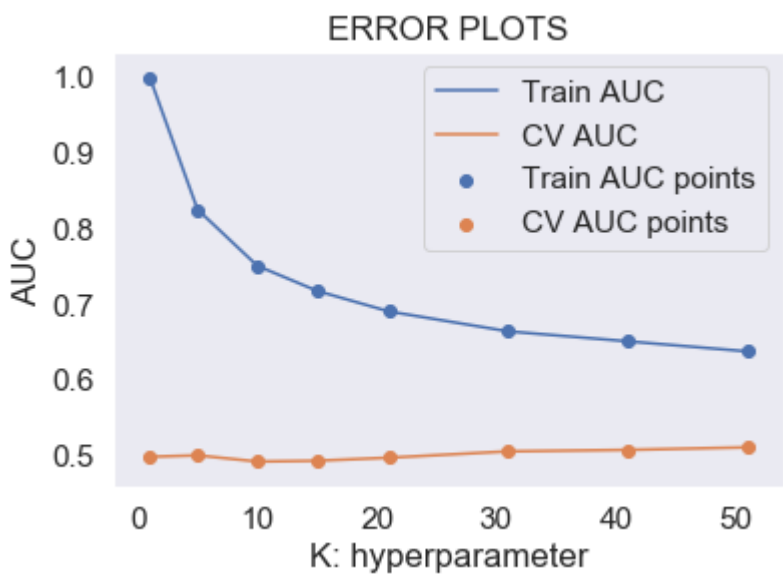
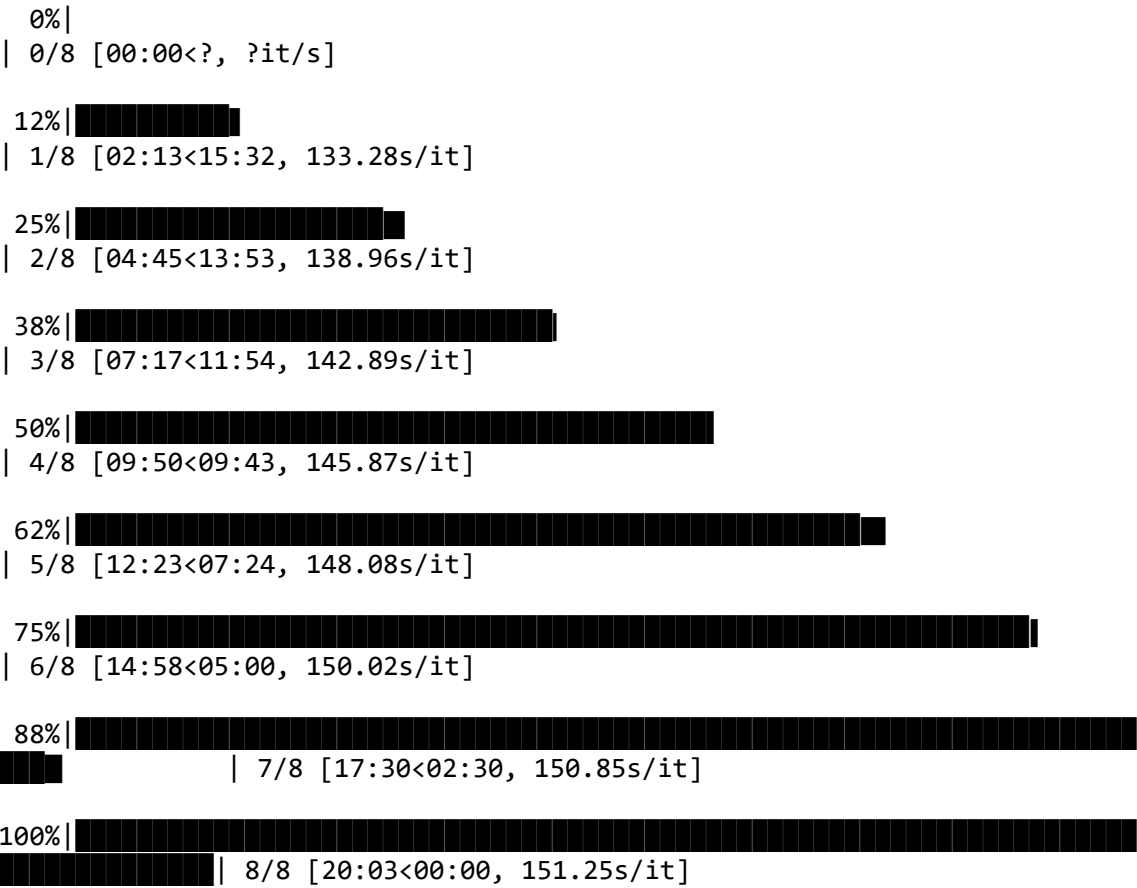
    y_train_pred = neigh.predict_proba(X_train_merge_set_2_best)[:,-1]
    y_cv_pred = neigh.predict_proba(X_cv_merge_set_2_best)[:,-1]

    # 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))
    a.append(y_train_pred)
    b.append(y_cv_pred)

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

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

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



**From above we can easily see that for lower values of k we have large AUC on training data while we have Low AUC of the test data hence we can say that our model falls under the category of overfitting for those lower values of k.**

**Choosing the best k from above plot**

In [202]:

```
best_k_selectKBest = 49
```

**NOW USING THE BEST K VALUE RECEIVED FROM ABOVE WE WILL BE TRAINING OUR MODEL**

In [203]:

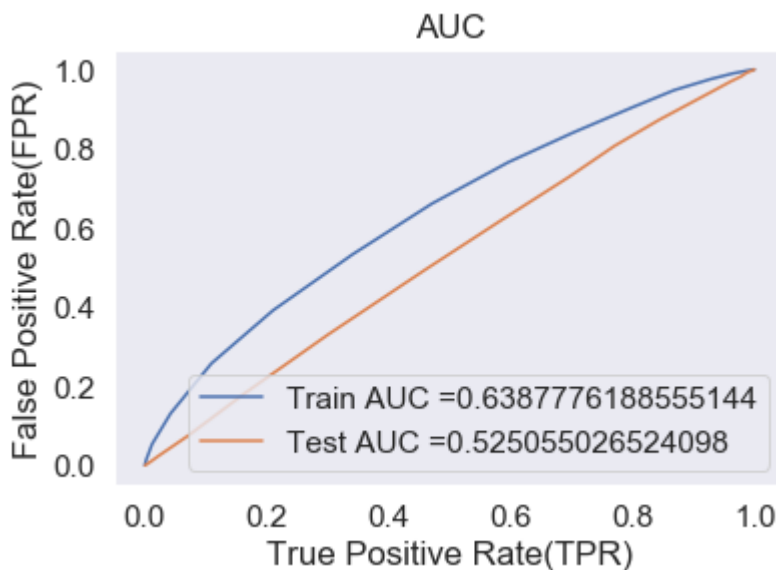
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

neigh = KNeighborsClassifier(n_neighbors=best_k_selectKBest)
neigh.fit(X_train_merge_set_2_best, 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 = neigh.predict_proba(X_train_merge_set_2_best)[: ,1]
y_test_pred = neigh.predict_proba(X_test_merge_set_2_best)[: ,1]

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

plt.plot(train_fpr, train_tpr, label="Train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



**Auc values received here for the test data are the worst it means that working on only 2000 features was not a good idea.**

**Now getting the confusion matrix for the train and test data based on above plots**

In [204]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

```
=====
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.2491539469166972 for threshold 0.837
[[ 3929  3497]
 [14100 27515]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.2483979476844841 for threshold 0.878
[[ 3822  1637]
 [20501 10092]]
```

**visually plotting the confusion matrix for train data**

In [205]:

```
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

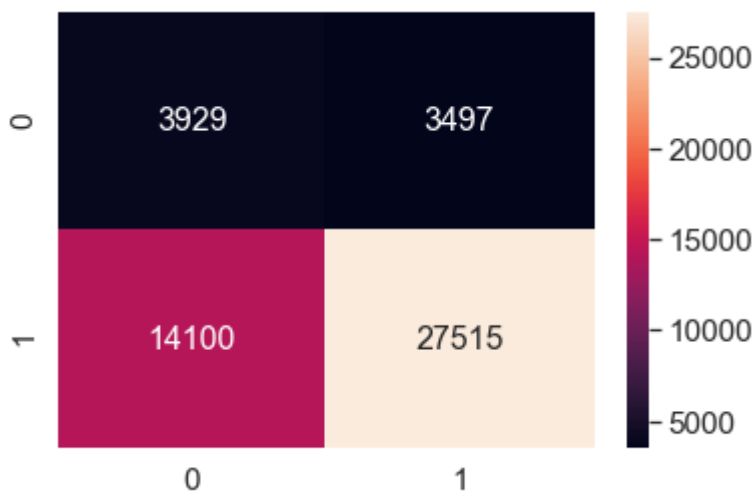
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt

df_cm_train = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)), range(2),
                           range(2))
plt.figure(figsize = (10,7))
sn.set(font_scale=1.4)#for label size
sn.heatmap(df_cm_train, annot=True,annot_kws={"size": 16},fmt = 'g')# font size
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.2491539469166972 for threshold 0.837

Out[205]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x22be59e9208>



## Visually plotting the confusion matrix for test data

In [206]:

```
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
```

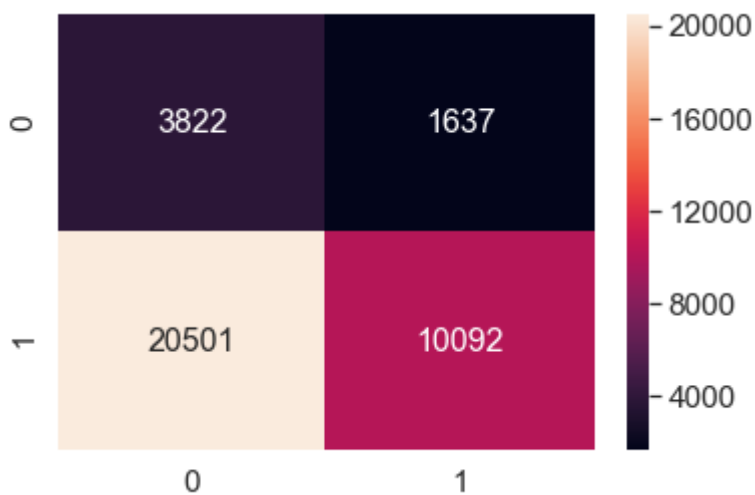
```
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
```

```
df_cm_train = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds,
test_fpr, test_fpr)))
#plt.figure(figsize = (10,7))
sn.set(font_scale=1.4)#for label size
sn.heatmap(df_cm_train, annot=True,annot_kws={"size": 16},fmt = 'g')# font size
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.2483979476844841 for threshold 0.878

Out[206]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x22b9490b5c0>



Cannot find much of good predictions and also the model is giving the worse AUC as compared to previous models which we did in different sets.

## Conclusions



In [207]:

```
# Compare all your models using Prettytable Library
# http://zetcode.com/python/prettytable/

from prettytable import PrettyTable

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

x = PrettyTable()
x.field_names = ["Vectorizer Technique", "Model", "Hyper Parameter K Value", "AUC"]

x.add_row(["BOW", "Brute", 31, 0.62])
x.add_row(["TFIDF", "Brute", 41, 0.56])
x.add_row(["AVG W2V", "Brute", 35, 0.59])
x.add_row(["TFIDF W2V", "Brute", 53, 0.60])

print(x)
```

Vectorizer Technique	Model	Hyper Parameter K Value	AUC
BOW	Brute	31	0.62
TFIDF	Brute	41	0.56
AVG W2V	Brute	35	0.59
TFIDF W2V	Brute	53	0.6

**Looks like k in the range of 31-41 will work good based on the trainings for different vectorizations techniques.**

In [ ]:

In [ ]: