

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	<ul style="list-style-type: none">••	Title of the project. Examples: Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: Grades PreK-2 Grades 3-5 Grades 6-8 Grades 9-12
<code>project_subject_categories</code>	<ul style="list-style-type: none">•••••••••	One or more (comma-separated) subject categories for the project from the following enumerated list of values: Applied Learning Care & Hunger Health & Sports History & Civics Literacy & Language Math & Science Music & The Arts Special Needs Warmth
	<ul style="list-style-type: none">••	Examples: Music & The Arts Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (Two-letter U.S. postal code (https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes)). Example: WY	
<code>project_subject_subcategories</code>	<ul style="list-style-type: none">••	One or more (comma-separated) subject subcategories for the project. Examples: Literacy Literature & Writing, Social Sciences
<code>project_resource_summary</code>	<ul style="list-style-type: none">•	An explanation of the resources needed for the project. Example: My students need hands on literacy materials to manage sensory needs!</code

Feature	Description
project_essay_1	First application essay*
project_essay_2	Second application essay*
project_essay_3	Third application essay*
project_essay_4	Fourth application essay*
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <div><ul style="list-style-type: none">nanDr.Mr.Mrs.Ms.Teacher.</div>
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. Example: 2

* See the section **Notes on the Essay Data** for more details about these features.

Additionally, the `resources.csv` data set provides more data about the resources required for each project. Each line in this file represents a resource required by a project:

Feature	Description
id	A <code>project_id</code> value from the <code>train.csv</code> file. Example: p036502
description	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
quantity	Quantity of the resource required. Example: 3
price	Price of the resource required. Example: 9.95

Note: Many projects require multiple resources. The `id` value corresponds to a `project_id` in `train.csv`, so you use it as a key to retrieve all resources needed for a project:

The data set contains the following label (the value you will attempt to predict):

Label	Description
project_is_approved	A binary flag indicating whether DonorsChoose approved the project. A value of <code>0</code> indicates the project was not approved, and a value of <code>1</code> indicates the project was approved.

Notes on the Essay Data

Prior to May 17, 2016, the prompts for the essays were as follows:

- `__project_essay_1:__` "Introduce us to your classroom"
- `__project_essay_2:__` "Tell us more about your students"
- `__project_essay_3:__` "Describe how your students will use the materials you're requesting"
- `__project_essay_3:__` "Close by sharing why your project will make a difference"

Starting on May 17, 2016, the number of essays was reduced from 4 to 2, and the prompts for the first 2 essays were changed to the following:

- `__project_essay_1:__` "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- `__project_essay_2:__` "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

For all projects with `project_submitted_datetime` of 2016-05-17 and later, the values of `project_essay_3` and `project_essay_4` will be NaN.

```
In [1]: ▶ 1  ## import all the modules
2  %matplotlib inline
3  import warnings
4  warnings.filterwarnings("ignore")
5
6  import sqlite3
7  import pandas as pd
8  import numpy as np
9  import nltk
10 import string
11 import matplotlib.pyplot as plt
12 import seaborn as sns
13 from sklearn.feature_extraction.text import TfidfTransformer
14 from sklearn.feature_extraction.text import TfidfVectorizer
15 from scipy.sparse import hstack
16 from sklearn.feature_extraction.text import CountVectorizer
17 from sklearn.metrics import confusion_matrix
18 from sklearn import metrics
19 from sklearn.metrics import roc_curve, auc
20 from nltk.stem.porter import PorterStemmer
21 from sklearn.preprocessing import Normalizer
22 from sklearn.linear_model import LogisticRegression
23 from sklearn.metrics import roc_auc_score
24 import re
25 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
26 import string
27 from nltk.corpus import stopwords
28 from nltk.stem import PorterStemmer
29 from nltk.stem.wordnet import WordNetLemmatizer
30 from nltk.sentiment import SentimentIntensityAnalyzer
31 from gensim.models import Word2Vec
32 from gensim.models import KeyedVectors
33 import pickle
34
35 from tqdm import tqdm
36 import os
37
38 # from plotly import plotly
39 # import plotly.offline as offline
40 # import plotly.graph_objs as go
41 # offline.init_notebook_mode()
42 from collections import Counter
```

1.1 Reading the Data

```
In [2]: ▶ 1  project_data = pd.read_csv('train_data.csv')
2  resource_data = pd.read_csv('resources.csv')
```

```
In [3]: 1 print("Number of data points in project train data", project_data.shape)
2 print('-'*50)
3 print("The attributes of data :", project_data.columns.values)
```

Number of data points in project 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 [4]: 1 print("Number of data points in resource train data", resource_data.shape)
2 print(resource_data.columns.values)
3 resource_data.head(2)
```

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

Out[4]:

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

1.1 Preprocessing Categorical Features: project_grade_category

```
In [5]: 1 print("Project grade", project_data['project_grade_category'].value_counts(dropna=False))
2 ## visulaize how project grade looks like
3 print('-'*50)
4 print(project_data['project_grade_category'].values[1000])
5 print(project_data['project_grade_category'].values[1500])
```

Project grade Grades PreK-2 44225
Grades 3-5 37137
Grades 6-8 16923
Grades 9-12 10963
Name: project_grade_category, dtype: int64

Grades 3-5
Grades PreK-2

```
In [6]: 1 # https://stackoverflow.com/questions/36383821/pandas-dataframe-apply-function-to-column-strings-based-on-other-column-value
2 project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('Grades ', '')
3 project_data['project_grade_category'] = project_data['project_grade_category'].str.replace(' ', '_')
4 project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('-', '_')
5 project_data['project_grade_category'] = project_data['project_grade_category'].str.lower()
6 project_data['project_grade_category'].value_counts()
```

```
Out[6]: prek_2      44225
3_5      37137
6_8      16923
9_12     10963
Name: project_grade_category, dtype: int64
```

1.2 Preprocessing Categorical Features: project_subject_category

```
In [7]: 1 catogories = list(project_data['project_subject_categories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
3 # reference from course material : reference_EDA.ipynb
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7 cat_list = []
8 for i in catogories:
9     temp = ""
10    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
11    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
12        if 'The' in j.split(): # this will split each of the catogory based on space "Math & Science"=> "Math", "&", "Science"
13            j=j.replace('The', '') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
14            j = j.replace(' ', '') # we are placeing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
15            temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
16            temp = temp.replace('&', '_') # we are replacing the & value into
17    cat_list.append(temp.strip())
18
19 project_data['clean_categories'] = cat_list
20 project_data.drop(['project_subject_categories'], axis=1, inplace=True)
21 project_data.head(2)
22
23
24 ### maintain a dict that
25 my_counter=Counter()
26
27 for word in project_data['clean_categories'].values:
28     my_counter.update(word.split())
29 cat_dict=dict(my_counter)
30
31 sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
```

1.3 Preprocessing Categorical Features: project_subject_subcategory

```

In [8]: ▶ 1 sub_categories = list(project_data['project_subject_subcategories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
3
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7
8 sub_cat_list = []
9 for i in sub_categories:
10     temp = ""
11     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
12     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
13         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
14             j=j.replace('The', '') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
15             j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
16             temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
17             temp = temp.replace('&', '_')
18     sub_cat_list.append(temp.strip())
19
20 project_data['clean_subcategories'] = sub_cat_list
21 project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
22
23 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
24 my_counter = Counter()
25 for word in project_data['clean_subcategories'].values:
26     my_counter.update(word.split())
27
28 sub_cat_dict = dict(my_counter)
29 sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))

```

1.3 Preprocessing Categorical Features: school_state

```
In [9]: 1 project_data['school_state'].value_counts()
        2 ## Convert it to Lower
        3 project_data['school_state'] = project_data['school_state'].str.lower()
        4 project_data['school_state'].value_counts(dropna=False)
```

```
Out[9]: ca      15388
        tx       7396
        ny       7318
        fl       6185
        nc       5091
        il       4350
        ga       3963
        sc       3936
        mi       3161
        pa       3109
        in       2620
        mo       2576
        oh       2467
        la       2394
        ma       2389
        wa       2334
        ok       2276
        nj       2237
        az       2147
        va       2045
        wi       1827
        al       1762
        ut       1731
        tn       1688
        ct       1663
        md       1514
        nv       1367
        ms       1323
        ky       1304
        or       1242
        mn       1208
        co       1111
        ar       1049
        id        693
        ia        666
        ks        634
        nm        557
        dc        516
        hi        507
        me        505
        wv        503
        nh        348
        ak        345
        de        343
        ne        309
        sd        300
        ri        285
        mt        245
        nd        143
        wy         98
        vt         80
        Name: school_state, dtype: int64
```

1.4 Preprocessing Categorical Features: Teacher_prefix

```
In [10]: 1 print(project_data['teacher_prefix'].value_counts(dropna=False))
2 # try to remove the dots from the teacher prefix and replace nan with mrs.
3 project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('Mrs.')
4 project_data['teacher_prefix']=project_data['teacher_prefix'].str.replace('.', '')
5 project_data['teacher_prefix']=project_data['teacher_prefix'].str.lower()
6 project_data['teacher_prefix']=project_data['teacher_prefix'].str.strip()
```

```
Mrs.      57269
Ms.       38955
Mr.       10648
Teacher   2360
Dr.        13
NaN         3
Name: teacher_prefix, dtype: int64
```

1.5 Combining all the essays

```
In [11]: 1 # merge two column text dataframe:
2 project_data["essay"] = project_data["project_essay_1"].map(str) + \
3                        project_data["project_essay_2"].map(str) + \
4                        project_data["project_essay_3"].map(str) + \
5                        project_data["project_essay_4"].map(str)
```

1.6 Number of Words in the Essay and Title

```
In [12]: 1 source:'''https://www.geeksforgeeks.org/python-program-to-count-words-in-a-sentence/'''
2 words_counter=[]
3 for string in project_data['essay']:
4     res = len(re.findall(r'\w+', string))
5     words_counter.append(res)
6
7 project_data["words_in_essay"] = words_counter
8
9 words_counter=[]
10
11 for string in project_data['project_title']:
12     res = len(re.findall(r'\w+', string))
13     words_counter.append(res)
14 project_data["words_in_title"] = words_counter
```

1.7. Preprocessing Numerical Values: price

```
In [13]: 1 ## calculate the overall count of resources and the total price for each project id
2 price_data=resource_data.groupby('id').agg({'price':'sum','quantity':'sum' })
3
```

```
In [14]: 1 project_data=pd.merge(project_data, price_data, on='id', how='left')
```


1.8 Preprocessing Text Features: project_title

```
In [15]: 1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\ 're", " are", phrase)
12    phrase = re.sub(r"\ 's", " is", phrase)
13    phrase = re.sub(r"\ 'd", " would", phrase)
14    phrase = re.sub(r"\ 'll", " will", phrase)
15    phrase = re.sub(r"\ 't", " not", phrase)
16    phrase = re.sub(r"\ 've", " have", phrase)
17    phrase = re.sub(r"\ 'm", " am", phrase)
18    return phrase
```

```
In [16]: 1 # https://gist.github.com/sebleier/554280
2 # we are removing the words from the stop words List: 'no', 'nor', 'not'
3 stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
4             "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
5             'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', \
6             'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
7             'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
8             'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
9             'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
10            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
11            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \
12            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
13            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
14            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', \
15            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', \
16            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
17            'won', "won't", 'wouldn', "wouldn't"]
```

```
In [17]: 1 print("printing some random reviews")
2 print(9, project_data['project_title'].values[9])
3 print(34, project_data['project_title'].values[34])
4 print(147, project_data['project_title'].values[147])
```

```
printing some random reviews
9 Just For the Love of Reading--\r\nPure Pleasure
34 \"Have A Ball!!!\"
147 Who needs a Chromebook?\r\nWE DO!!
```

```
In [18]: 1 # Combining all the above stundents
2 from tqdm import tqdm
3 def preprocess_text(text_data):
4     preprocessed_text = []
5     # tqdm is for printing the status bar
6     for sentence in tqdm(text_data):
7         sent = decontracted(sentence)
8         sent = sent.replace('\\r', ' ')
9         sent = sent.replace('\\n', ' ')
10        sent = sent.replace('\\\"', ' ')
11        sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12        # https://gist.github.com/sebleier/554280
13        sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
14        preprocessed_text.append(sent.lower().strip())
15    return preprocessed_text
```

```
In [19]: 1 preprocessed_titles = preprocess_text(project_data['project_title'].values)

100%|████████████████████████████████████████| 109248/109248 [00:02<00:00, 36843.35it/s]
```

```
In [20]: 1 print("printing some random reviews")
2 print(9, preprocessed_titles[9])
3 print(34, preprocessed_titles[34])
4 print(147, preprocessed_titles[147])
```

```
printing some random reviews
9 love reading pure pleasure
34 ball
147 needs chromebook
```

1.9 Preprocessing Text Features: essay

In [21]:

```

1 print("printing some random essay")
2 print(9, project_data['essay'].values[9])
3 print('-'*50)
4 print(34, project_data['essay'].values[34])
5 print('-'*50)
6 print(147, project_data['essay'].values[147])

```

printing some random essay

9 Over 95% of my students are on free or reduced lunch. I have a few who are homeless, but despite that, they come to school with an eagerness to learn. My students are inquisitive eager learners who embrace the challenge of not having great books and other resources every day. Many of them are not afforded the opportunity to engage with these big colorful pages of a book on a regular basis at home and they don't travel to the public library. \r\nIt is my duty as a teacher to do all I can to provide each student an opportunity to succeed in every aspect of life. \r\nReading is Fundamental! My students will read these books over and over again while boosting their comprehension skills. These books will be used for read alouds, partner reading and for Independent reading. \r\nThey will engage in reading to build their "Love for Reading" by reading for pure enjoyment. They will be introduced to some new authors as well as some old favorites. I want my students to be ready for the 21st Century and know the pleasure of holding a good hard back book in hand. There's nothing like a good book to read! \r\nMy students will soar in Reading, and more because of your consideration and generous funding contribution. This will help build stamina and prepare for 3rd grade. Thank you so much for reading our proposal!nannan

34 My students mainly come from extremely low-income families, and the majority of them come from homes where both parents work full time. Most of my students are at school from 7:30 am to 6:00 pm (2:30 to 6:00 pm in the after-school program), and they all receive free and reduced meals for breakfast and lunch. \r\n\r\n\r\nI want my students to feel as comfortable in my classroom as they do at home. Many of my students take on multiple roles both at home as well as in school. They are sometimes the caretakers of younger siblings, cooks, babysitters, academics, friends, and most of all, they are developing who they are going to become as adults. I consider it an essential part of my job to model helping others gain knowledge in a positive manner. As a result, I have a community of students who love helping each other in and outside of the classroom. They consistently look for opportunities to support each other's learning in a kind and helpful way. I am excited to be experimenting with alternative seating in my classroom this school year. Studies have shown that giving students the option of where they sit in a classroom increases focus as well as motivation. \r\n\r\nBy allowing students choice in the classroom, they are able to explore and create in a welcoming environment. Alternative classroom seating has been experimented with more frequently in recent years. I believe (along with many others), that every child learns differently. This does not only apply to how multiplication is memorized, or a paper is written, but applies to the space in which they are asked to work. I have had students in the past ask "Can I work in the library? Can I work on the carpet?" My answer was always, "As long as you're learning, you can work wherever you want!" \r\n\r\nWith the yoga balls and the lap-desks, I will be able to increase the options for seating in my classroom and expand its imaginable space.nannan

147 My students are eager to learn and make their mark on the world.\r\n\r\nThey come from a Title 1 school and need extra love.\r\n\r\nMy fourth grade students are in a high poverty area and still come to school every day to get their education. I am trying to make it fun and educational for them so they can get the most out of their schooling. I created a caring environment for the students to bloom! They deserve the best.\r\nThank you!\r\nI am requesting 1 Chromebook to access online interventions, differentiate instruction, and get extra practice. The Chromebook will be used to supplement ELA and math instruction. Students will play ELA and math games that are engaging and fun, as well as participate in assignments online. This in turn will help my students improve their skills. Having a Chromebook in the classroom would not only allow students to use the programs at their own pace, but would ensure more students are getting adequate time to use the programs. The online programs have been especially beneficial to my students with special needs. They are able to work at their level as well as be challenged with some different materials. This is making these students more confident in their abilities.\r\n\r\nThe Chromebook would allow my students to have daily access to computers and increase their computing skills.\r\nThis will change their lives for the better as they become more successful in school. Having access to technology in the classroom would help bridge the achievement gap.nannan

In [22]:

```

1 preprocessed_essays = preprocess_text(project_data['essay'].values)

```

100%|████████████████████████████████████████| 109248/109248 [01:13<00:00, 1484.64it/s]

```
In [23]: 1 print("printing some random essay")
2 print(9, preprocessed_essays[9])
3 print('-'*50)
4 print(34, preprocessed_essays[34])
5 print('-'*50)
6 print(147, preprocessed_essays[147])
7
8 #merge the column in the project_data
9 project_data['processed_essay']=preprocessed_essays
```

printing some random essay

9 95 students free reduced lunch homeless despite come school eagerness learn students inquisitive eager learners embrace challenge not great books resources every day many not afforded opportunity engage big colorful pages book regular basis home not travel public library duty teacher provide student opportunity succeed every aspect life reading fund amental students read books boosting comprehension skills books used read alouds partner reading independent reading engage reading build love reading reading pure enjoyment in troduced new authors well old favorites want students ready 21st century know pleasure holding good hard back book hand nothing like good book read students soar reading consid eration generous funding contribution help build stamina prepare 3rd grade thank much reading proposal nannan

34 students mainly come extremely low income families majority come homes parents work full time students school 7 30 6 00 pm 2 30 6 00 pm school program receive free reduced m eals breakfast lunch want students feel comfortable classroom home many students take multiple roles home well school sometimes caretakers younger siblings cooks babysitters ac ademics friends developing going become adults consider essential part job model helping others gain knowledge positive manner result community students love helping outside cl assroom consistently look opportunities support learning kind helpful way excited experimenting alternative seating classroom school year studies shown giving students option s it classroom increases focus well motivation allowing students choice classroom able explore create welcoming environment alternative classroom seating experimented frequently recent years believe along many others every child learns differently not apply multiplication memorized paper written applies space asked work students past ask work library w ork carpet answer always long learning work wherever want yoga balls lap desks able increase options seating classroom expand imaginable space nannan

147 students eager learn make mark world come title 1 school need extra love fourth grade students high poverty area still come school every day get education trying make fun e ducational get schooling created caring environment students bloom deserve best thank requesting 1 chromebook access online interventions differentiate instruction get extra pr actice chromebook used supplement ela math instruction students play ela math games engaging fun well participate assignments online turn help students improve skills chromeboo k classroom would not allow students use programs pace would ensure students getting adequate time use programs online programs especially beneficial students special needs abl e work level well challenged different materials making students confident abilities chromebook would allow students daily access computers increase computing skills change liv es better become successful school access technology classroom would help bridge achievement gap nannan

1.10 Preprocessing Text Features: Project Title

```
In [24]: 1 print("printing some random project_titles")
2 print(9, project_data['project_title'].values[9])
3 print('-'*50)
4 print(34, project_data['project_title'].values[34])
5 print('-'*50)
6 print(147, project_data['project_title'].values[147])
```

printing some random project_titles

9 Just For the Love of Reading--\r\nPure Pleasure

34 \"Have A Ball!!!\"

147 Who needs a Chromebook?\r\nWE DO!!

```
In [25]: 1 processed_title = preprocess_text(project_data['project_title'].values)
```

100%|████████████████████████████████████████| 109248/109248 [00:05<00:00, 21518.70it/s]

```
In [26]: 1 print("printing some random title")
2 print(9, preprocessed_titles[9])
3 print('-'*50)
4 print(34, preprocessed_titles[34])
5 print('-'*50)
6 print(147, preprocessed_titles[147])
7
8 #merge the column in the project_data
9 project_data['processed_title']=preprocessed_titles
```

```
printing some random title
9 love reading pure pleasure
```

```
-----
34 ball
```

```
-----
147 needs chromebook
```

Creating sentiment columns

```
In [27]: 1 ## craete the sentiment columns using
2
3 neg=[]
4 pos=[]
5 neu=[]
6 compound=[]
7 sentiment_model=SentimentIntensityAnalyzer()
8 for text in tqdm(project_data['processed_essay']):
9     pol_scores = sentiment_model.polarity_scores(text)
10     neg.append(pol_scores['neg'])
11     pos.append(pol_scores['pos'])
12     neu.append(pol_scores['neu'])
13     compound.append(pol_scores['compound'])
14
15 project_data['pos']=pos
16 project_data['neg']=neg
17 project_data['neu']=neu
18 project_data['compound']=compound
```

```
100%|████████████████████████████████████████| 109248/109248 [06:24<00:00, 284.07it/s]
```

2 Train,Test,CV Split

```
In [28]: 1 # train test split using sklearn.model selection
2 from sklearn.model_selection import train_test_split
3 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'], random_state=0)
4 X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train, random_state=0)
```

```
In [29]: 1 ## drop the y labels from splits
2 X_train.drop(['project_is_approved'], axis=1, inplace=True)
3 X_test.drop(['project_is_approved'], axis=1, inplace=True)
4 X_cv.drop(['project_is_approved'], axis=1, inplace=True)
```

In [30]: ▶

1 X_train.head(2)

Out[30]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_category	project_title	project_essay_1	project_essay_2	...	words_in_essay	words_
75742	118221	p186156	f50f55a2b44b65b54f38f03c5df21922	mrs	tx	2017-03-01 16:21:46	9_12	I-Waste: a Multi-Media Art Installation on El...	My students are creative human beings. They a...	It's no secret that the arts are underfunded i...	...	241	
61001	57644	p180433	9e0fb5827f551d7e6966f8b3985e387b	ms	ny	2017-03-09 10:19:06	6_8	Authentic Listening and Speaking Activities fo...	The Hyde Park Central School District is locat...	One area my students struggle the most with is...	...	254	

2 rows × 27 columns

3. VECTORIZING DATA

3.1 One hot encoding on Categorical

In [31]: ▶

1 # we use count vectorizer to convert the values into one hot vectors
2 ## clean categories
3
4 cat_vectorize = CountVectorizer(lowercase=False, binary=True)
5 cat_vectorize.fit(X_train['clean_categories'].values)
6
7 train_categories = cat_vectorize.transform(X_train['clean_categories'].values)
8 test_categories = cat_vectorize.transform(X_test['clean_categories'].values)
9 cv_categories = cat_vectorize.transform(X_cv['clean_categories'].values)
10
11 print(cat_vectorize.get_feature_names())
12 print("Shape of matrix of Train data after one hot encoding ",train_categories.shape)
13 print("Shape of matrix of Test data after one hot encoding ",test_categories.shape)
14 print("Shape of matrix of CV data after one hot encoding ",cv_categories.shape)

['AppliedLearning', 'Care_Hunger', 'Health_Sports', 'History_Civics', 'Literacy_Language', 'Math_Science', 'Music_Arts', 'SpecialNeeds', 'Warmth']
Shape of matrix of Train data after one hot encoding (49041, 9)
Shape of matrix of Test data after one hot encoding (36052, 9)
Shape of matrix of CV data after one hot encoding (24155, 9)

```
In [32]: ▶ 1 # we use count vectorizer to convert the values into one hot vectors
2 ## clean subcategories
3
4 subcat_vectorize = CountVectorizer(lowercase=False, binary=True)
5 subcat_vectorize.fit(X_train['clean_subcategories'].values)
6
7 train_subcategories = subcat_vectorize.transform(X_train['clean_subcategories'].values)
8 test_subcategories = subcat_vectorize.transform(X_test['clean_subcategories'].values)
9 cv_subcategories = subcat_vectorize.transform(X_cv['clean_subcategories'].values)
10
11 print(subcat_vectorize.get_feature_names())
12 print("Shape of matrix of Train data after one hot encoding ",train_subcategories.shape)
13 print("Shape of matrix of Test data after one hot encoding ",test_subcategories.shape)
14 print("Shape of matrix of CV data after one hot encoding ",cv_subcategories.shape)
15
```

```
['AppliedSciences', 'Care_Hunger', 'CharacterEducation', 'Civics_Government', 'College_CareerPrep', 'CommunityService', 'ESL', 'EarlyDevelopment', 'Economics', 'EnvironmentalSc
ience', 'Extracurricular', 'FinancialLiteracy', 'ForeignLanguages', 'Gym_Fitness', 'Health_LifeScience', 'Health_Wellness', 'History_Geography', 'Literacy', 'Literature_Writin
g', 'Mathematics', 'Music', 'NutritionEducation', 'Other', 'ParentInvolvement', 'PerformingArts', 'SocialSciences', 'SpecialNeeds', 'TeamSports', 'VisualArts', 'Warmth']
Shape of matrix of Train data after one hot encoding (49041, 30)
Shape of matrix of Test data after one hot encoding (36052, 30)
Shape of matrix of CV data after one hot encoding (24155, 30)
```

```
In [33]: ▶ 1 # we use count vectorizer to convert the values into one hot vectors
2 ## school state
3
4
5 sklstate_vectorize = CountVectorizer(lowercase=False, binary=True)
6 sklstate_vectorize.fit(X_train['school_state'].values)
7
8 sklstate_train = sklstate_vectorize.transform(X_train['school_state'].values)
9 sklstate_test = sklstate_vectorize.transform(X_test['school_state'].values)
10 sklstate_cv = sklstate_vectorize.transform(X_cv['school_state'].values)
11
12 print(sklstate_vectorize.get_feature_names())
13 print("Shape of matrix of Train data after one hot encoding ",sklstate_train.shape)
14 print("Shape of matrix of Test data after one hot encoding ",sklstate_test.shape)
15 print("Shape of matrix of CV data after one hot encoding ",sklstate_cv.shape)
```

```
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'ks', 'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd',
'ne', 'nh', 'nj', 'nm', 'nv', 'ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv', 'wy']
Shape of matrix 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 CV data after one hot encoding (24155, 51)
```



```
In [34]: 1 # we use count vectorizer to convert the values into one hot vectors
2 ## teacher_prefix
3 from sklearn.feature_extraction.text import CountVectorizer
4
5 teacher_prefix_vectorize = CountVectorizer(lowercase=False, binary=True)
6 teacher_prefix_vectorize.fit(X_train['teacher_prefix'].values)
7
8 teacher_prefix_train = teacher_prefix_vectorize.transform(X_train['teacher_prefix'].values)
9 teacher_prefix_test = teacher_prefix_vectorize.transform(X_test['teacher_prefix'].values)
10 teacher_prefix_cv = teacher_prefix_vectorize.transform(X_cv['teacher_prefix'].values)
11
12 print(teacher_prefix_vectorize.get_feature_names())
13 print("Shape of matrix of Train data after one hot encoding ",teacher_prefix_train.shape)
14 print("Shape of matrix of Test data after one hot encoding ",teacher_prefix_test.shape)
15 print("Shape of matrix of CV data after one hot encoding ",teacher_prefix_cv.shape)
```

```
['dr', 'mr', 'mrs', 'ms', 'teacher']
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 CV data after one hot encoding (24155, 5)
```

```
In [35]: 1 # we use count vectorizer to convert the values into one hot vectors
2 ## project_grade
3
4 proj_grade_vectorize = CountVectorizer(lowercase=False, binary=True)
5 proj_grade_vectorize.fit(X_train['project_grade_category'].values)
6
7 proj_grade_train = proj_grade_vectorize.transform(X_train['project_grade_category'].values)
8 proj_grade_test = proj_grade_vectorize.transform(X_test['project_grade_category'].values)
9 proj_grade_cv = proj_grade_vectorize.transform(X_cv['project_grade_category'].values)
10
11 print(proj_grade_vectorize.get_feature_names())
12 print("Shape of matrix of Train data after one hot encoding ",proj_grade_train.shape)
13 print("Shape of matrix of Test data after one hot encoding ",proj_grade_test.shape)
14 print("Shape of matrix of CV data after one hot encoding ",proj_grade_cv.shape)
```

```
['3_5', '6_8', '9_12', 'prek_2']
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 CV data after one hot encoding (24155, 4)
```

3.2 Vectorizing Text data

3.2.1 BOW on Essay data


```
In [36]: 1  ##Considering the words that appeared in atleast 10 documents
2
3  bow_essay = CountVectorizer(min_df=10,max_features=5000,ngram_range=(1,2))
4  bow_essay.fit(X_train['processed_essay'])
5
6  bow_essay_train = bow_essay.transform(X_train['processed_essay'])
7
8  print("Shape of matrix after one hot encoding ",bow_essay_train.shape)
9
10 ## tranform Test data
11
12 bow_essay_test = bow_essay.transform(X_test['processed_essay'])
13
14 print("Shape of matrix after one hot encoding ",bow_essay_test.shape)
15
16
17 ## Teansform cv data
18
19 bow_essay_cv = bow_essay.transform(X_cv['processed_essay'])
20 print("Shape of matrix after one hot encoding ",bow_essay_cv.shape)
```

```
Shape of matrix after one hot encoding (49041, 5000)
Shape of matrix after one hot encoding (36052, 5000)
Shape of matrix after one hot encoding (24155, 5000)
```

3.2.2 BOW on Title data

```
In [37]: 1  ##Considering the words that appeared in atleast 10 documents
2
3  bow_title = CountVectorizer(min_df=10,max_features=5000,ngram_range=(1,2))
4  bow_title.fit(X_train['processed_title'])
5
6  bow_title_train = bow_title.transform(X_train['processed_title'])
7
8  print("Shape of matrix after one hot encoding ",bow_title_train.shape)
9
10 ## tranform Test data
11
12 bow_title_test = bow_title.transform(X_test['processed_title'])
13
14 print("Shape of matrix after one hot encoding ",bow_title_test.shape)
15
16
17 ## Teansform cv data
18
19 bow_title_cv = bow_title.transform(X_cv['processed_title'])
20 print("Shape of matrix after one hot encoding ",bow_title_cv.shape)
```

```
Shape of matrix after one hot encoding (49041, 3244)
Shape of matrix after one hot encoding (36052, 3244)
Shape of matrix after one hot encoding (24155, 3244)
```

3.2.3 TFIDF on Essay data

```

In [38]: 1  ##Considering the words that appeared in atleast 10 documents
2
3  tfidf_essay = TfidfVectorizer(min_df=10,max_features=5000,ngram_range=(1,2))
4  tfidf_essay.fit(X_train['processed_essay'])
5
6  tfidf_essay_train = tfidf_essay.transform(X_train['processed_essay'])
7
8  print("Shape of matrix after one hot encoding ",tfidf_essay_train.shape)
9
10 ## tranform Test data
11
12 tfidf_essay_test = tfidf_essay.transform(X_test['processed_essay'])
13
14 print("Shape of matrix after one hot encoding ",tfidf_essay_test.shape)
15
16
17 ## Teansform cv data
18
19 tfidf_essay_cv = tfidf_essay.transform(X_cv['processed_essay'])
20 print("Shape of matrix after one hot encoding ",tfidf_essay_cv.shape)

```

```

Shape of matrix after one hot encoding (49041, 5000)
Shape of matrix after one hot encoding (36052, 5000)
Shape of matrix after one hot encoding (24155, 5000)

```

3.2.4 TFIDF on Title data

```

In [39]: 1  ##Considering the words that appeared in atleast 10 documents
2
3  tfidf_title = TfidfVectorizer(min_df=10,max_features=5000,ngram_range=(1,2))
4  tfidf_title.fit(X_train['processed_title'])
5
6  tfidf_title_train = tfidf_title.transform(X_train['processed_title'])
7
8  print("Shape of matrix after one hot encoding ",tfidf_title_train.shape)
9
10 ## tranform Test data
11
12 tfidf_title_test = tfidf_title.transform(X_test['processed_title'])
13
14 print("Shape of matrix after one hot encoding ",tfidf_title_test.shape)
15
16
17 ## Teansform cv data
18
19 tfidf_title_cv = tfidf_title.transform(X_cv['processed_title'])
20 print("Shape of matrix after one hot encoding ",tfidf_title_cv.shape)

```

```

Shape of matrix after one hot encoding (49041, 3244)
Shape of matrix after one hot encoding (36052, 3244)
Shape of matrix after one hot encoding (24155, 3244)

```

3.2.5 Avg W2V on Essay data using Pretrained Models

In [40]:

```
1 # stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
2 # make sure you have the glove_vectors file
3 ## Glove vectors are global vectors for words which has vector every word in 300d .
4 ## for read more :https://nlp.stanford.edu/projects/glove/
5 with open('glove_vectors', 'rb') as f:
6     model = pickle.load(f)
7     glove_words = set(model.keys())
```

In [41]:

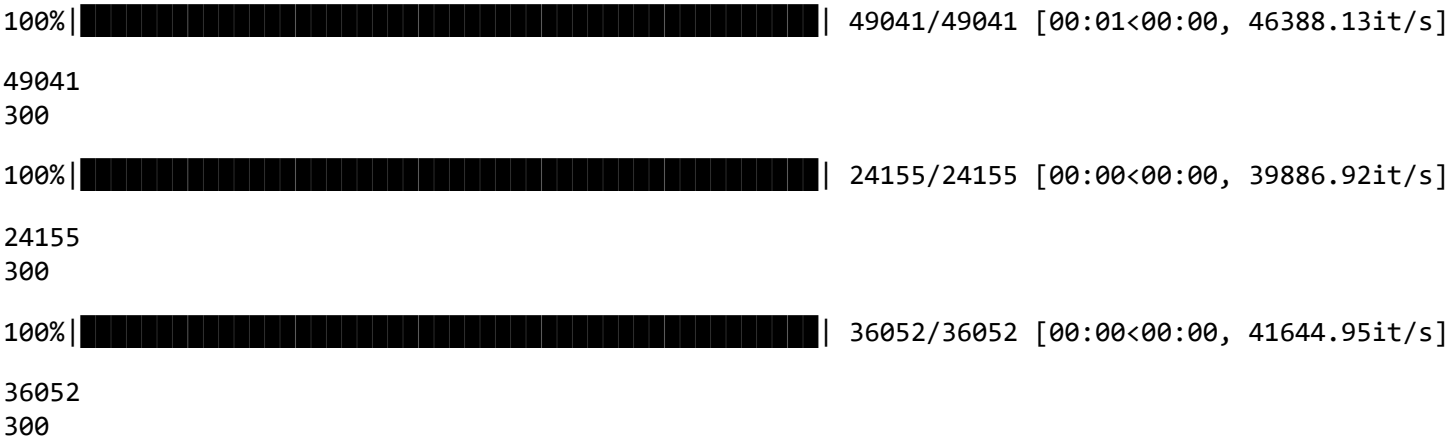
```
1  # average Word2Vec on train
2  # compute average word2vec for each review.
3
4  avg_w2v_vectors_train = [];
5
6  for sentence in tqdm(X_train["processed_essay"]): # for each review/sentence
7      vector = np.zeros(300) # as word vectors are of zero length
8      cnt_words = 0; # num of words with a valid vector in the sentence/review
9      for word in sentence.split(): # for each word in a review/sentence
10         if word in glove_words:
11             vector += model[word]
12             cnt_words += 1
13         if cnt_words != 0:
14             vector /= cnt_words
15         avg_w2v_vectors_train.append(vector)
16
17 print(len(avg_w2v_vectors_train))
18 print(len(avg_w2v_vectors_train[0]))
19
20
21 # average Word2Vec on CV
22 # compute average word2vec for each review.
23
24 avg_w2v_vectors_cv = [];
25
26 for sentence in tqdm(X_cv["processed_essay"]): # for each review/sentence
27     vector = np.zeros(300) # as word vectors are of zero length
28     cnt_words = 0; # num of words with a valid vector in the sentence/review
29     for word in sentence.split(): # for each word in a review/sentence
30         if word in glove_words:
31             vector += model[word]
32             cnt_words += 1
33         if cnt_words != 0:
34             vector /= cnt_words
35         avg_w2v_vectors_cv.append(vector)
36
37 print(len(avg_w2v_vectors_cv))
38 print(len(avg_w2v_vectors_cv[0]))
39
40
41 # average Word2Vec on test
42 # compute average word2vec for each review.
43
44 avg_w2v_vectors_test = [];
45
46 for sentence in tqdm(X_test["processed_essay"]): # for each review/sentence
47     vector = np.zeros(300) # as word vectors are of zero length
48     cnt_words = 0; # num of words with a valid vector in the sentence/review
49     for word in sentence.split(): # for each word in a review/sentence
50         if word in glove_words:
51             vector += model[word]
52             cnt_words += 1
53         if cnt_words != 0:
54             vector /= cnt_words
55         avg_w2v_vectors_test.append(vector)
56
57 print(len(avg_w2v_vectors_test))
58 print(len(avg_w2v_vectors_test[0]))
```

100% | 49041/49041 [00:24<00:00, 1965.82it/s]

36052
300

In [42]:

```
1  # average Word2Vec on train
2  # compute average word2vec for each review.
3
4  avg_w2v_vectors_title_train = [];
5
6  for sentence in tqdm(X_train["processed_title"]): # for each review/sentence
7      vector = np.zeros(300) # as word vectors are of zero length
8      cnt_words = 0; # num of words with a valid vector in the sentence/review
9      for word in sentence.split(): # for each word in a review/sentence
10         if word in glove_words:
11             vector += model[word]
12             cnt_words += 1
13         if cnt_words != 0:
14             vector /= cnt_words
15         avg_w2v_vectors_title_train.append(vector)
16
17 print(len(avg_w2v_vectors_title_train))
18 print(len(avg_w2v_vectors_title_train[0]))
19
20
21 # average Word2Vec on CV
22 # compute average word2vec for each review.
23
24 avg_w2v_vectors_title_cv = [];
25
26 for sentence in tqdm(X_cv["processed_title"]): # for each review/sentence
27     vector = np.zeros(300) # as word vectors are of zero length
28     cnt_words = 0; # num of words with a valid vector in the sentence/review
29     for word in sentence.split(): # for each word in a review/sentence
30         if word in glove_words:
31             vector += model[word]
32             cnt_words += 1
33         if cnt_words != 0:
34             vector /= cnt_words
35         avg_w2v_vectors_title_cv.append(vector)
36
37 print(len(avg_w2v_vectors_title_cv))
38 print(len(avg_w2v_vectors_title_cv[0]))
39
40
41 # average Word2Vec on test
42 # compute average word2vec for each review.
43
44 avg_w2v_vectors_title_test = [];
45
46 for sentence in tqdm(X_test["processed_title"]): # for each review/sentence
47     vector = np.zeros(300) # as word vectors are of zero length
48     cnt_words = 0; # num of words with a valid vector in the sentence/review
49     for word in sentence.split(): # for each word in a review/sentence
50         if word in glove_words:
51             vector += model[word]
52             cnt_words += 1
53         if cnt_words != 0:
54             vector /= cnt_words
55         avg_w2v_vectors_title_test.append(vector)
56
57 print(len(avg_w2v_vectors_title_test))
58 print(len(avg_w2v_vectors_title_test[0]))
```



3.2.7 Weighted tfidf on Essay data using Pretrained Models

In [43]:

```

1  # average Word2Vec on train
2  # compute average word2vec for each review.
3  # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
4
5  tfidf_model = TfidfVectorizer()
6  tfidf_model.fit(X_train["processed_essay"])
7  # we are converting a dictionary with word as a key, and the idf as a value
8  dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
9  tfidf_words = set(tfidf_model.get_feature_names())
10
11 tfidf_w2v_vectors_train = [];
12
13 for sentence in tqdm(X_train["processed_essay"]): # for each review/sentence
14     vector = np.zeros(300) # as word vectors are of zero length
15     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
16     for word in sentence.split(): # for each word in a review/sentence
17         if (word in glove_words) and (word in tfidf_words):
18             vec = model[word] # getting the vector for each word
19             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
20             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
21             vector += (vec * tf_idf) # calculating tfidf weighted w2v
22             tf_idf_weight += tf_idf
23     if tf_idf_weight != 0:
24         vector /= tf_idf_weight
25     tfidf_w2v_vectors_train.append(vector)
26
27 print(len(tfidf_w2v_vectors_train))
28 print(len(tfidf_w2v_vectors_train[0]))
29
30
31 # average Word2Vec on CV
32 # compute average word2vec for each review.
33
34 tfidf_w2v_vectors_cv = [];
35
36 for sentence in tqdm(X_cv["processed_essay"]): # for each review/sentence
37     vector = np.zeros(300) # as word vectors are of zero length
38     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
39     for word in sentence.split(): # for each word in a review/sentence
40         if (word in glove_words) and (word in tfidf_words):
41             vec = model[word] # getting the vector for each word
42             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
43             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
44             vector += (vec * tf_idf) # calculating tfidf weighted w2v
45             tf_idf_weight += tf_idf
46     if tf_idf_weight != 0:
47         vector /= tf_idf_weight
48     tfidf_w2v_vectors_cv.append(vector)
49
50
51 print(len(tfidf_w2v_vectors_cv))
52 print(len(tfidf_w2v_vectors_cv[0]))
53
54
55 # average Word2Vec on test
56 # compute average word2vec for each review.
57
58 tfidf_w2v_vectors_test = [];
59
60 for sentence in tqdm(X_test["processed_essay"]): # for each review/sentence

```


300

localhost:8888/notebooks/Documents/appleidai/LR/Assignment_7_LR.ipynb

In [44]:

```

1  # average Word2Vec on train
2  # compute average word2vec for each review.
3  # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
4
5  tfidf_model = TfidfVectorizer()
6  tfidf_model.fit(X_train["processed_title"])
7  # we are converting a dictionary with word as a key, and the idf as a value
8  dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
9  tfidf_words = set(tfidf_model.get_feature_names())
10
11 tfidf_w2v_vectors_title_train = [];
12
13 for sentence in tqdm(X_train["processed_title"]): # for each review/sentence
14     vector = np.zeros(300) # as word vectors are of zero length
15     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
16     for word in sentence.split(): # for each word in a review/sentence
17         if (word in glove_words) and (word in tfidf_words):
18             vec = model[word] # getting the vector for each word
19             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
20             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
21             vector += (vec * tf_idf) # calculating tfidf weighted w2v
22             tf_idf_weight += tf_idf
23     if tf_idf_weight != 0:
24         vector /= tf_idf_weight
25     tfidf_w2v_vectors_title_train.append(vector)
26
27 print(len(tfidf_w2v_vectors_title_train))
28 print(len(tfidf_w2v_vectors_title_train[0]))
29
30
31 # average Word2Vec on CV
32 # compute average word2vec for each review.
33
34 tfidf_w2v_vectors_title_cv = [];
35
36 for sentence in tqdm(X_cv["processed_title"]): # for each review/sentence
37     vector = np.zeros(300) # as word vectors are of zero length
38     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
39     for word in sentence.split(): # for each word in a review/sentence
40         if (word in glove_words) and (word in tfidf_words):
41             vec = model[word] # getting the vector for each word
42             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
43             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
44             vector += (vec * tf_idf) # calculating tfidf weighted w2v
45             tf_idf_weight += tf_idf
46     if tf_idf_weight != 0:
47         vector /= tf_idf_weight
48     tfidf_w2v_vectors_title_cv.append(vector)
49
50
51 print(len(tfidf_w2v_vectors_title_cv))
52 print(len(tfidf_w2v_vectors_title_cv[0]))
53
54
55 # average Word2Vec on test
56 # compute average word2vec for each review.
57
58 tfidf_w2v_vectors_title_test = [];
59
60 for sentence in tqdm(X_test["processed_title"]): # for each review/sentence

```

100%	<div><div></div></div>	49041/49041 [00:02<00:00, 19513.53it/s]
49041		
300		
100%	<div><div></div></div>	24155/24155 [00:01<00:00, 22377.67it/s]
24155		
300		
100%	<div><div></div></div>	36052/36052 [00:01<00:00, 19622.08it/s]
36052		
300		

4.1 Price

```

In [45]: ▶ 1 normalizer = Normalizer()
2 # normalizer.fit(X_train['price'].values)
3 # this will rise an error Expected 2D array, got 1D array instead:
4 # array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
5 # Reshape your data either using
6 # array.reshape(-1, 1) if your data has a single feature
7 # array.reshape(1, -1) if it contains a single sample.
8 normalizer.fit(X_train['price'].values.reshape(1,-1))
9
10 X_train_price_norm = normalizer.transform(X_train['price'].values.reshape(1,-1))
11 X_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(1,-1))
12 X_test_price_norm = normalizer.transform(X_test['price'].values.reshape(1,-1))
13
14 print("After vectorizations")
15 print(X_train_price_norm.shape, y_train.shape)
16 print(X_cv_price_norm.shape, y_cv.shape)
17 print(X_test_price_norm.shape, y_test.shape)
18 print("="*100)
19
20 ## reshaping
21 X_train_price_norm=X_train_price_norm.reshape(-1,1)
22 X_cv_price_norm=X_cv_price_norm.reshape(-1,1)
23 X_test_price_norm=X_test_price_norm.reshape(-1,1)

```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

=====

4.2 Quantity

```

In [46]: ▶ 1
          2 normalizer = Normalizer()
          3
          4 # normalizer.fit(X_train['price'].values)
          5 # this will rise an error Expected 2D array, got 1D array instead:
          6 # array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
          7 # Reshape your data either using
          8 # array.reshape(-1, 1) if your data has a single feature
          9 # array.reshape(1, -1) if it contains a single sample.
         10
         11 normalizer.fit(X_train['quantity'].values.reshape(1,-1))
         12
         13 quantity_train_norm = normalizer.transform(X_train['quantity'].values.reshape(1,-1))
         14 quantity_cv_norm = normalizer.transform(X_cv['quantity'].values.reshape(1,-1))
         15 quantity_test_norm = normalizer.transform(X_test['quantity'].values.reshape(1,-1))
         16
         17 print("After vectorizations")
         18 print(quantity_train_norm.shape, y_train.shape)
         19 print(quantity_cv_norm.shape, y_cv.shape)
         20 print(quantity_test_norm.shape, y_test.shape)
         21 print("="*100)
         22
         23 ## reshaping
         24 quantity_train_norm=quantity_train_norm.reshape(-1,1)
         25 quantity_cv_norm=quantity_cv_norm.reshape(-1,1)
         26 quantity_test_norm=quantity_test_norm.reshape(-1,1)

```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

=====

4.3 Number of Previously posted projects

```

In [47]: ▶ 1 normalizer = Normalizer()
2
3 # normalizer.fit(X_train['price'].values)
4 # this will rise an error Expected 2D array, got 1D array instead:
5 # array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
6 # Reshape your data either using
7 # array.reshape(-1, 1) if your data has a single feature
8 # array.reshape(1, -1) if it contains a single sample.
9
10 normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
11
12 prev_projects_train_norm = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
13 prev_projects_cv_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
14 prev_projects_test_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
15
16 print("After vectorizations")
17 print(prev_projects_train_norm.shape, y_train.shape)
18 print(prev_projects_cv_norm.shape, y_cv.shape)
19 print(prev_projects_test_norm.shape, y_test.shape)
20 print("=="*100)
21
22 ## reshaping
23 prev_projects_train_norm=prev_projects_train_norm.reshape(-1,1)
24 prev_projects_cv_norm=prev_projects_cv_norm.reshape(-1,1)
25 prev_projects_test_norm=prev_projects_test_norm.reshape(-1,1)

```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

=====

4.4 Title Word counts

```
In [48]: 1 normalizer = Normalizer()
2
3 normalizer.fit(X_train['words_in_title'].values.reshape(1,-1))
4
5 title_word_count_train_norm = normalizer.transform(X_train['words_in_title'].values.reshape(1,-1))
6 title_word_count_cv_norm = normalizer.transform(X_cv['words_in_title'].values.reshape(1,-1))
7 title_word_count_test_norm = normalizer.transform(X_test['words_in_title'].values.reshape(1,-1))
8
9 print("After vectorizations")
10 print(title_word_count_train_norm.shape, y_train.shape)
11 print(title_word_count_cv_norm.shape, y_cv.shape)
12 print(title_word_count_test_norm.shape, y_test.shape)
13 print("="*100)
14
15 ## reshaping
16 title_word_count_train_norm=title_word_count_train_norm.reshape(-1,1)
17 title_word_count_cv_norm=title_word_count_cv_norm.reshape(-1,1)
18 title_word_count_test_norm=title_word_count_test_norm.reshape(-1,1)
```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

=====

4.5 Essay Words Counts

```
In [49]: 1 normalizer = Normalizer()
2
3 normalizer.fit(X_train['words_in_essay'].values.reshape(1,-1))
4
5 essay_word_count_train_norm = normalizer.transform(X_train['words_in_essay'].values.reshape(1,-1))
6 essay_word_count_cv_norm = normalizer.transform(X_cv['words_in_essay'].values.reshape(1,-1))
7 essay_word_count_test_norm = normalizer.transform(X_test['words_in_essay'].values.reshape(1,-1))
8
9 print("After vectorizations")
10 print(essay_word_count_train_norm.shape, y_train.shape)
11 print(essay_word_count_cv_norm.shape, y_cv.shape)
12 print(essay_word_count_test_norm.shape, y_test.shape)
13
14 ## reshaping
15 essay_word_count_train_norm=essay_word_count_train_norm.reshape(-1,1)
16 essay_word_count_cv_norm=essay_word_count_cv_norm.reshape(-1,1)
17 essay_word_count_test_norm=essay_word_count_test_norm.reshape(-1,1)
```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

4.6 Vectorizing sentiment Columns

```

In [50]: 1  ### vectorize pos
          2  Normalize=Normalizer()
          3  Normalize.fit(X_train['pos'].values.reshape(1,-1))
          4  sentiment_pos_train_norm=Normalize.transform(X_train['pos'].values.reshape(1,-1))
          5  sentiment_pos_test_norm=Normalize.transform(X_test['pos'].values.reshape(1,-1))
          6  sentiment_pos_cv_norm=Normalize.transform(X_cv['pos'].values.reshape(1,-1))
          7  sentiment_pos_train_norm=sentiment_pos_train_norm.reshape(-1,1)
          8  sentiment_pos_test_norm=sentiment_pos_test_norm.reshape(-1,1)
          9  sentiment_pos_cv_norm=sentiment_pos_cv_norm.reshape(-1,1)

In [51]: 1  ### vectorize neg
          2  Normalize=Normalizer()
          3  Normalize.fit(X_train['neg'].values.reshape(1,-1))
          4  sentiment_neg_train_norm=Normalize.transform(X_train['neg'].values.reshape(1,-1))
          5  sentiment_neg_test_norm=Normalize.transform(X_test['neg'].values.reshape(1,-1))
          6  sentiment_neg_cv_norm=Normalize.transform(X_cv['neg'].values.reshape(1,-1))
          7  sentiment_neg_train_norm=sentiment_neg_train_norm.reshape(-1,1)
          8  sentiment_neg_test_norm=sentiment_neg_test_norm.reshape(-1,1)
          9  sentiment_neg_cv_norm=sentiment_neg_cv_norm.reshape(-1,1)

In [52]: 1  ### vectorize compound
          2  Normalize=Normalizer()
          3  Normalize.fit(X_train['compound'].values.reshape(1,-1))
          4  sentiment_compound_train_norm=Normalize.transform(X_train['compound'].values.reshape(1,-1))
          5  sentiment_compound_test_norm=Normalize.transform(X_test['compound'].values.reshape(1,-1))
          6  sentiment_compound_cv_norm=Normalize.transform(X_cv['compound'].values.reshape(1,-1))
          7  sentiment_compound_train_norm=sentiment_compound_train_norm.reshape(-1,1)
          8  sentiment_compound_test_norm=sentiment_compound_test_norm.reshape(-1,1)
          9  sentiment_compound_cv_norm=sentiment_compound_cv_norm.reshape(-1,1)

In [53]: 1  ### vectorize neu
          2  Normalize=Normalizer()
          3  Normalize.fit(X_train['neu'].values.reshape(1,-1))
          4  sentiment_neu_train_norm=Normalize.transform(X_train['neu'].values.reshape(1,-1))
          5  sentiment_neu_test_norm=Normalize.transform(X_test['neu'].values.reshape(1,-1))
          6  sentiment_neu_cv_norm=Normalize.transform(X_cv['neu'].values.reshape(1,-1))
          7  sentiment_neu_train_norm=sentiment_neu_train_norm.reshape(-1,1)
          8  sentiment_neu_test_norm=sentiment_neu_test_norm.reshape(-1,1)
          9  sentiment_neu_cv_norm=sentiment_neu_cv_norm.reshape(-1,1)

```

Assignment 7: Logistic Regression

1. [Task-1] Logistic Regression(either SGDClassifier with log loss, or LogisticRegression) on these feature sets

- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_eassay ('BOW with bi-grams' with 'min_df=10' and 'max_features=5000')
- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay ('TFIDF with bi-grams' with 'min_df=10' and 'max_features=5000')
- **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)
- **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

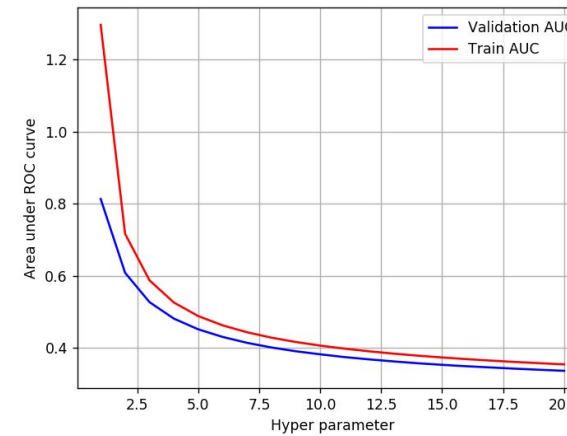
2. Hyper paramter tuning (find best hyper parameters corresponding the algorithm that you choose)

- Find the best hyper parameter which will give the maximum [AUC \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) value
- Find the best hyper paramter using k-fold cross validation or simple cross validation data

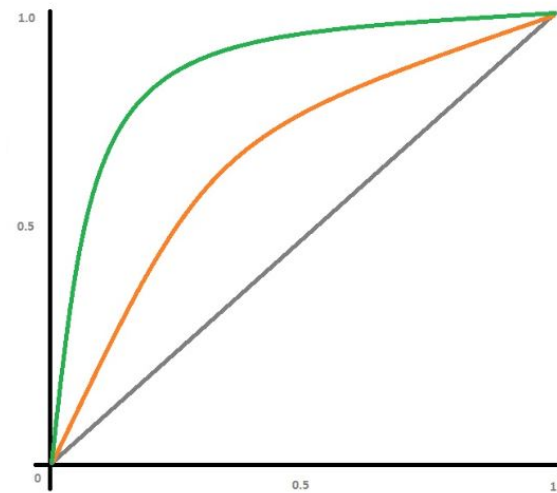
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.



- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.



- Along with plotting ROC curve, you need to print the [confusion matrix](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/>) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).

	Predicted: NO	Predicted: YES
Actual: NO	TN = ??	FP = ??
Actual: YES	FN = ??	TP = ??

(<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)

4. [Task-2] Apply Logistic Regression on the below feature set **Set 5** by finding the best hyper parameter as suggested in step 2 and step 3.

5. Consider these set of features **Set 5** :

- **school_state** : categorical data
- **clean_categories** : categorical data
- **clean_subcategories** : categorical data
- **project_grade_category** :categorical data
- **teacher_prefix** : categorical data
- **quantity** : numerical data
- **teacher_number_of_previously_posted_projects** : numerical data

- **price** : numerical data
- **sentiment score's of each of the essay** : numerical data
- **number of words in the title** : numerical data
- **number of words in the combine essays** : numerical data

And apply the Logistic regression on these features by finding the best hyper paramter as suggested in step 2 and step 3

6. Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library [link \(http://zetcode.com/python/prettytable/\)](http://zetcode.com/python/prettytable/)

Vectorizer	Model	Hyper parameter	AUC
BOW	Brute	7	0.78
TFIDF	Brute	12	0.79
W2V	Brute	10	0.78
TFIDFW2V	Brute	6	0.78

5. SET 1 : categorical, numerical features + project_title(BOW) + preprocessed_eassay (BOW with bi-grams with min_df=10 and max_features=5000)

```
In [54]: ▶ 1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2
3
4 X_tr = hstack((train_categories, train_subcategories,sklstate_train,teacher_prefix_train,
5               proj_grade_train,bow_essay_train,bow_title_train,
6               X_train_price_norm,quantity_train_norm,prev_projects_train_norm,title_word_count_train_norm,
7               essay_word_count_train_norm)).tocsr()
8
9 X_te = hstack((test_categories, test_subcategories,sklstate_test,teacher_prefix_test,
10              proj_grade_test,bow_essay_test,bow_title_test,
11              X_test_price_norm,quantity_test_norm,prev_projects_test_norm,title_word_count_test_norm,
12              essay_word_count_test_norm)).tocsr()
13
14 X_cr = hstack((cv_categories, cv_subcategories,sklstate_cv,teacher_prefix_cv,
15              proj_grade_cv,bow_essay_cv,bow_title_cv,
16              X_cv_price_norm,quantity_cv_norm,prev_projects_cv_norm,title_word_count_cv_norm,
17              essay_word_count_cv_norm)).tocsr()
18
19
20 print(X_tr.shape)
21 print(X_te.shape)
22 print(X_cr.shape)
```

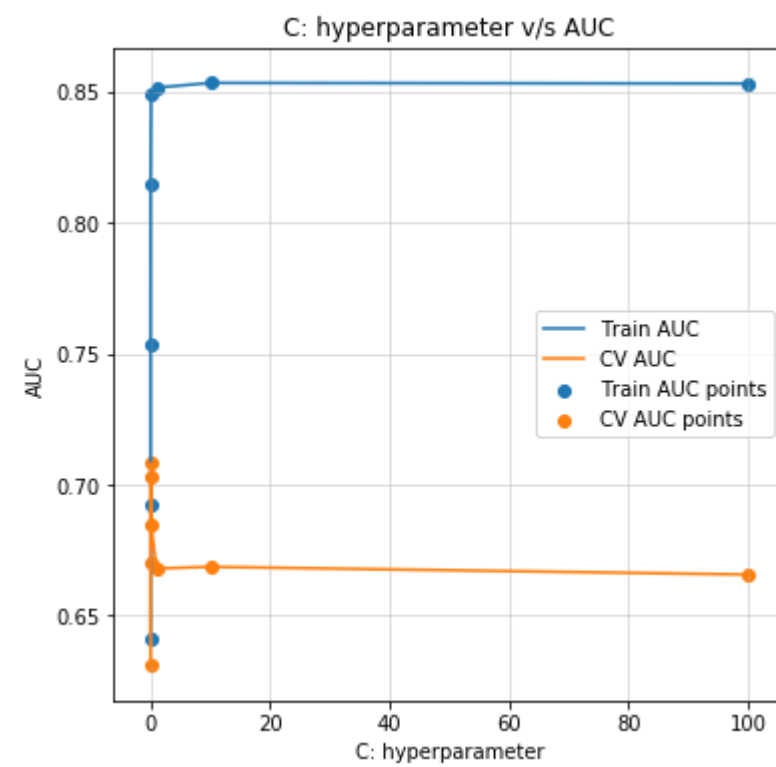
```
(49041, 8348)
(36052, 8348)
(24155, 8348)
```

```
Final Data matrix
(49041, 8348) (49041,)
(24155, 8348) (24155,)
(36052, 8348) (36052,)
=====
```

5.1 Write own function to find which alpha gives maximum auc

[illegible]

```
In [57]: 1 plt.figure(figsize=(6,6))
2         plt.plot(c, train_auc, label='Train AUC')
3         plt.plot(c, cv_auc, label='CV AUC')
4         plt.scatter(c, train_auc, label='Train AUC points')
5         plt.scatter(c, cv_auc, label='CV AUC points')
6
7         plt.legend()
8         plt.xlabel("C: hyperparameter")
9         plt.ylabel("AUC")
10        plt.title("C: hyperparameter v/s AUC")
11        plt.grid(which='major', alpha=0.5)
12        plt.grid(which='minor', alpha=0.2)
13        plt.show()
```



Observations

1. The c values are taken in the range of 10^{-5} to 10^5
2. We can see a steep drop in AUC as the c value increases above 0.1
3. We can see as the c values increase in train data AUC also increases. While it is opposite in the case of cv data implying the case of overfitting .
4. with a c value less than 1 we can see both cv and train data AUC converge.
5. Optimal c value which can be selected in this can be 0.001

5.2 GridSearch CV using K-fold Crossvalidation with k=10

```
In [58]: 1 from sklearn.model_selection import GridSearchCV
2 from scipy.stats import randint as sp_randint
3 from sklearn.model_selection import RandomizedSearchCV
4
5 ## As we are seeing steady decreasing in AUC value in CV as c value increases more than 1 .Hence keep more values less
6 ## than 1
7 parameters={"C" : [0.0001,0.001,0.01,0.1,0.6,1,5,10,20] }
8
9 clf = GridSearchCV(LR,parameters, cv=5, scoring='roc_auc',verbose=1,n_jobs=3,return_train_score=True)
10 clf.fit(X_tr, y_train)
11 results = pd.DataFrame.from_dict(clf.cv_results_)
12
```

Fitting 5 folds for each of 9 candidates, totalling 45 fits

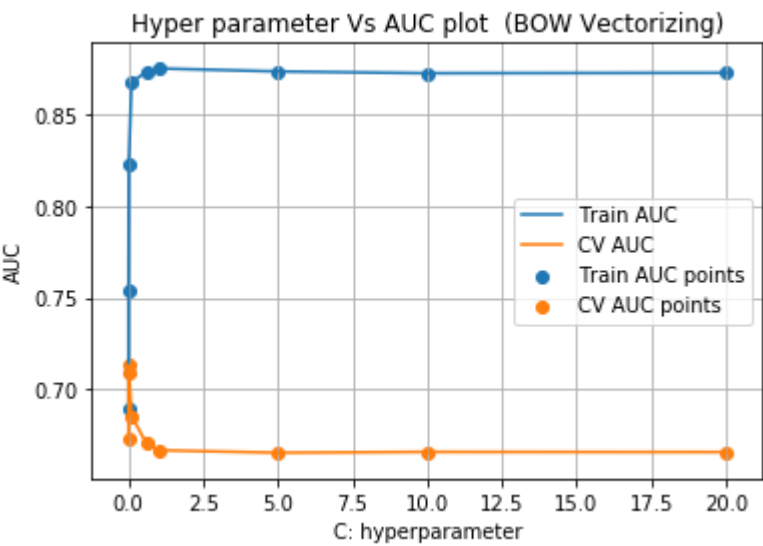
[Parallel(n_jobs=3)]: Using backend LokyBackend with 3 concurrent workers.

[Parallel(n_jobs=3)]: Done 45 out of 45 | elapsed: 59.7s finished

```
In [59]: 1 results = results.sort_values(['param_C'])
2 RS_alphas=results['param_C']
3 train_auc= results['mean_train_score']
4 train_auc_std= results['std_train_score']
5 cv_auc = results['mean_test_score']
6 cv_auc_std= results['std_test_score']
```

In [60]: ▶

```
1 plt.plot(RS_alphas, train_auc, label='Train AUC')
2 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
3 # plt.gca().fill_between(K, train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
4
5 plt.plot(RS_alphas, cv_auc, label='CV AUC')
6 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
7 # plt.gca().fill_between(K, cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
8
9 plt.scatter(RS_alphas, train_auc, label='Train AUC points')
10 plt.scatter(RS_alphas, cv_auc, label='CV AUC points')
11
12 plt.legend()
13 plt.xlabel("C: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("Hyper parameter Vs AUC plot (BOW Vectorizing)")
16 plt.grid()
17 plt.show()
18
19 results.head()
```



Out[60]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	params	split0_test_score	split1_test_score	split2_test_score	split3_test_score	...	mean_test_score	std_test_score	rank_test_score	split0_tra
0	1.193842	0.053759	0.013450	0.001847	0.0001	{'C': 0.0001}	0.676880	0.679327	0.663630	0.675099	...	0.672731	0.005749	4	
1	2.556824	0.103688	0.013152	0.001134	0.001	{'C': 0.001}	0.711388	0.713364	0.703465	0.709601	...	0.708885	0.003504	2	
2	4.170823	0.068766	0.014387	0.001301	0.01	{'C': 0.01}	0.716206	0.717785	0.707984	0.713126	...	0.713403	0.003427	1	
3	4.167730	0.080110	0.014416	0.001293	0.1	{'C': 0.1}	0.691922	0.688252	0.679283	0.686546	...	0.685731	0.004391	3	
4	4.098698	0.033398	0.013387	0.002114	0.6	{'C': 0.6}	0.674811	0.674728	0.665261	0.672970	...	0.670699	0.004303	5	

5 rows × 21 columns

Observations

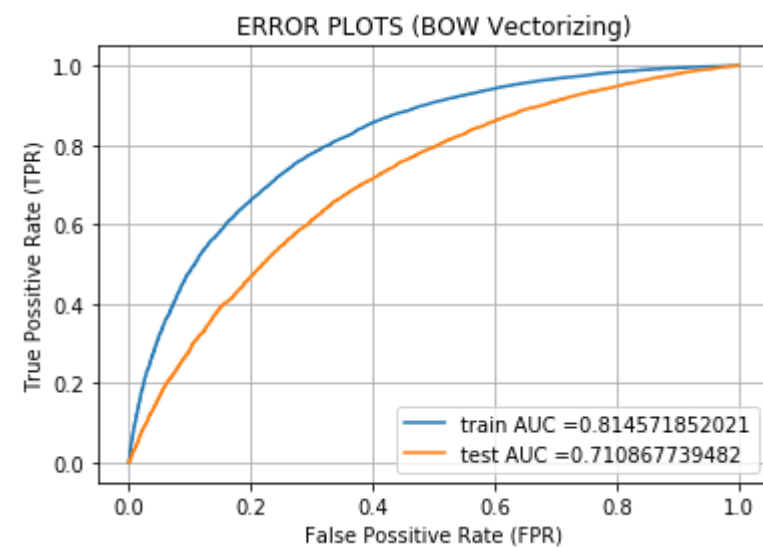
1. Using Grid Search technique we find that as C values increase after 0.001, AUC decreases steeply and maintains a steadiness beyond that.
2. Clearly we see the case of overfit from the graph when C values are beyond 1 (i.e. very high train AUC and low CV AUC)
3. Both Train and CV auc converge below 0.01.

5.3 Train the model using the best Hyperparameter value

```
In [61]: 1 ### https://forums.fast.ai/t/hyperparameter-random-search-interpretation/8591 ---to get the best hyper parameter as a result of Random search
2 best_C = clf.best_params_
3 print('Best Alpha as a result of Grid Search',best_C)
```

Best Alpha as a result of Grid Search {'C': 0.01}

```
In [63]: 1 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2 from sklearn.metrics import roc_curve, auc
3
4 LR=LogisticRegression(penalty="l2",C=best_C['C'],n_jobs=3)
5 LR.fit(X_tr, y_train)
6 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
7 # not the predicted outputs
8
9 y_train_pred = batch_predict(LR, X_tr)
10 y_test_pred = batch_predict(LR, X_te)
11
12 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
13 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
14
15 plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
16 plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
17 plt.legend()
18 plt.xlabel("False Possitive Rate (FPR)")
19 plt.ylabel("True Possitive Rate (TPR)")
20 plt.title("ERROR PLOTS (BOW Vectorizing)")
21 plt.grid()
22 plt.show()
```



Observations :

1. Test AUC found to be 0.78 and Train AUC as 0.81 after training model using best hyperparameter 0.01 and vectorizing text data using Bag Of Words .

5.4 Confusion Matrix

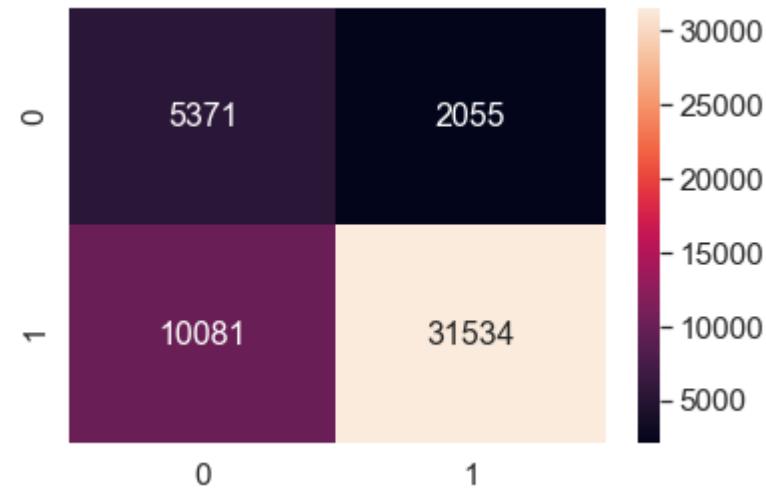
```
In [64]: ► 1  ## Finding best threshold for predictions
2  def best_threshold(thresholds,fpr,tpr):
3      t=thresholds[np.argmax(tpr*(1-fpr))]
4      # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
5      print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
6      return t
7
8  def predict_with_best_t(proba, threshold):
9      predictions = []
10     for i in proba:
11         if i>=threshold:
12             predictions.append(1)
13         else:
14             predictions.append(0)
15     return predictions
```

```
In [65]: ► 1  print("="*100)
2  from sklearn.metrics import confusion_matrix
3  best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
4  print("Train confusion matrix")
5  print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
6  print("Test confusion matrix")
7  print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

```
=====
the maximum value of tpr*(1-fpr) 0.548061596919 for threshold 0.821
Train confusion matrix
[[ 5371  2055]
 [10081 31534]]
Test confusion matrix
[[ 3173  2286]
 [ 8177 22416]]
```

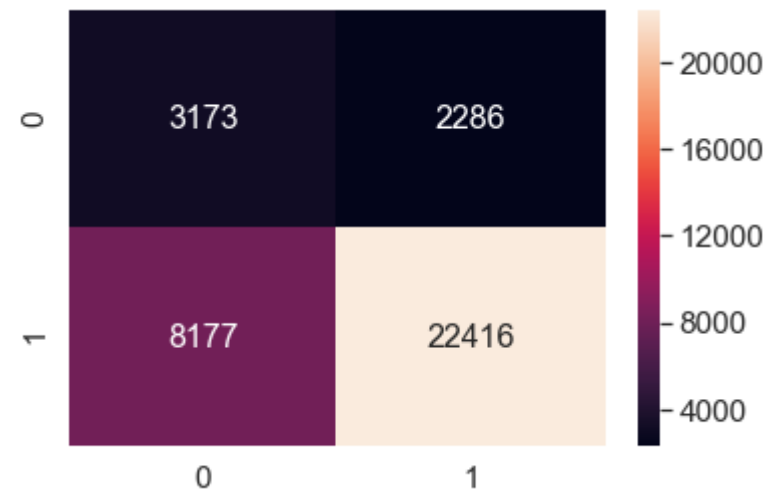
5.4.1 Plotting Confusion Matrix on Train data


```
In [66]: 1 ### PLOT the matrix for Train
2 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3 df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
4 # plt.figure(figsize=(10,7))
5 sns.set(font_scale=1.4) # for label size
6 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7 plt.show()
```



5.4.2 Plotting Confusion Matrix on Test data

```
In [67]: 1 ### PLOT the matrix for Train
2 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3 df_cm = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
4 # plt.figure(figsize=(10,7))
5 sns.set(font_scale=1.4) # for label size
6 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7 plt.show()
```



Observations :

1. We can observe from train and test we are getting majority True positives
2. Least number of data falls in False negative, which refers as least number of projects were incorrectly predicted as not approved in both Test and Train.
3. For a model to perform well we need High True Positive Rate and Low False Positive Rate. From the above our train data has True Positive Rate as 91% and False Positive Rate as 75%
4. In our test data : True Positive Rate as 90% and False Positive Rate as 76.7%.
5. Test data is nearly accurate as train data.

6. SET 2 : categorical, numerical features + project_title(TFIDF) + preprocessed_eassay (TFIDF with bi-grams with min_df=10 and max_features=5000)

```
In [68]: 1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2 from scipy.sparse import hstack
3
4
5 X_tr = hstack((train_categories, train_subcategories, sklstate_train, teacher_prefix_train,
6               proj_grade_train, tfidf_essay_train, tfidf_title_train,
7               X_train_price_norm, quantity_train_norm, prev_projects_train_norm, title_word_count_train_norm,
8               essay_word_count_train_norm)).tocsr()
9
10 X_te = hstack((test_categories, test_subcategories, sklstate_test, teacher_prefix_test,
11               proj_grade_test, tfidf_essay_test, tfidf_title_test,
12               X_test_price_norm, quantity_test_norm, prev_projects_test_norm, title_word_count_test_norm,
13               essay_word_count_test_norm)).tocsr()
14
15 X_cr = hstack((cv_categories, cv_subcategories, sklstate_cv, teacher_prefix_cv,
16               proj_grade_cv, tfidf_essay_cv, tfidf_title_cv,
17               X_cv_price_norm, quantity_cv_norm, prev_projects_cv_norm, title_word_count_cv_norm,
18               essay_word_count_cv_norm)).tocsr()
19
20
21 print(X_tr.shape)
22 print(X_te.shape)
23 print(X_cr.shape)
```

```
(49041, 8348)
(36052, 8348)
(24155, 8348)
```

```
In [69]: 1 print("Final Data matrix")
2 print(X_tr.shape, y_train.shape)
3 print(X_cr.shape, y_cv.shape)
4 print(X_te.shape, y_test.shape)
5 print("="*100)
```

```
Final Data matrix
(49041, 8348) (49041,)
(24155, 8348) (24155,)
(36052, 8348) (36052,)
```

```
=====
```

6.1 Write own function to find which alpha gives maximum auc

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import roc_auc_score
3 ## Lets consider a set of alphas from 10 ** -4 to 10 ** 4
4
5 c = [0.00001,0.0001,0.001,0.01,0.1,1,10,100]
6 train_auc=[]
7 cv_auc=[]
8 ## for each alpha now train the model multiple times and store the cv and train auc for plotting it later.
9 for i in tqdm(c):
10     LR=LogisticRegression(penalty="l2",C=i,n_jobs=3)
11     LR.fit(X_tr,y_train)
12     y_train_pred=batch_predict(LR,X_tr)
13     y_cv_pred=batch_predict(LR,X_cr)
14     train_auc.append(roc_auc_score(y_train,y_train_pred))
15     cv_auc.append(roc_auc_score(y_cv,y_cv_pred))

```

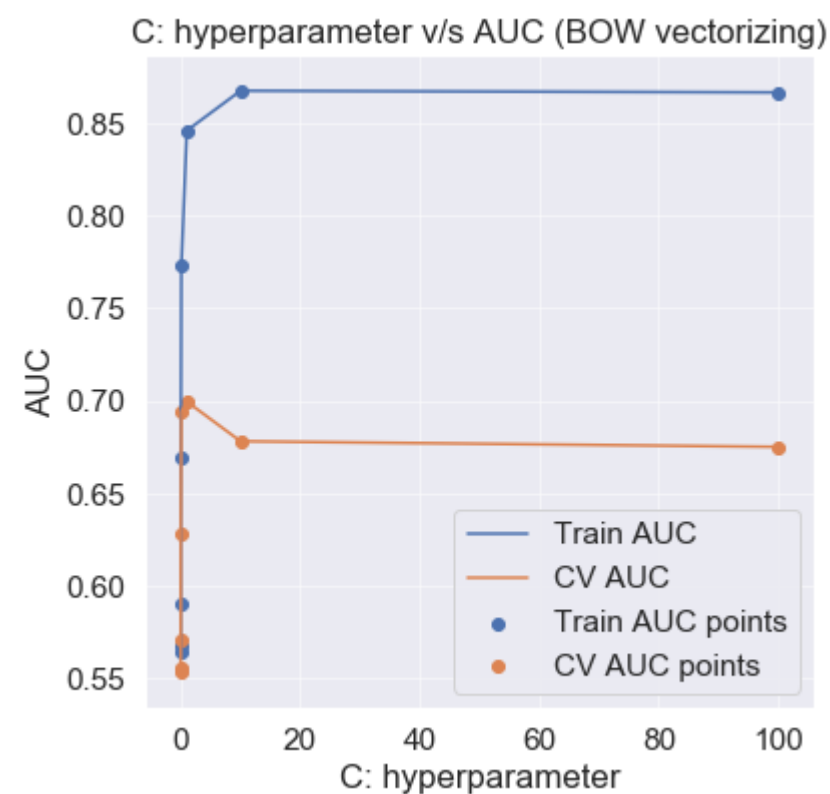
[illegible]

In [72]:

```

1 plt.figure(figsize=(6,6))
2
3
4 plt.plot(c, train_auc, label='Train AUC')
5 plt.plot(c, cv_auc, label='CV AUC')
6
7 plt.scatter(c, train_auc, label='Train AUC points')
8 plt.scatter(c, cv_auc, label='CV AUC points')
9
10
11 plt.legend()
12 plt.xlabel("C: hyperparameter")
13 plt.ylabel("AUC")
14 plt.title("C: hyperparameter v/s AUC (BOW vectorizing)")
15 plt.grid(which='major', alpha=0.5)
16 plt.grid(which='minor', alpha=0.2)
17 plt.show()

```



Observations

1. The c values are taken in the range of 10^{-5} to 10^2
2. We can see a steep drop in AUC as the c value increases above 0.1
3. We can see as the c values increase in train data AUC also increases. While it is opposite in the case of cv data implying the case of overfitting .
4. with a c value less than 1 we can see both cv and train data AUC converge.
5. Optimal c value which can be selected in this can be 0.001 / 0.01 .

6.2 GridSearch CV using K-fold Crossvalidation with k=10

```
In [77]: 1 from sklearn.model_selection import GridSearchCV
2 from scipy.stats import randint as sp_randint
3 from sklearn.model_selection import RandomizedSearchCV
4
5 ## As we are seeing steady decreasing in AUC value in CV as c value increases more than 1 .Hence keep more values less
6 ## than 1
7 parameters={"C" : [0.0001,0.001,0.01,0.1,0.6,1,5,10,20] }
8
9 clf = GridSearchCV(LR, parameters,return_train_score=True,cv=5, scoring='roc_auc',verbose=1,n_jobs=3)
10 clf.fit(X_tr, y_train)
11 results = pd.DataFrame.from_dict(clf.cv_results_)
12
```

Fitting 5 folds for each of 9 candidates, totalling 45 fits

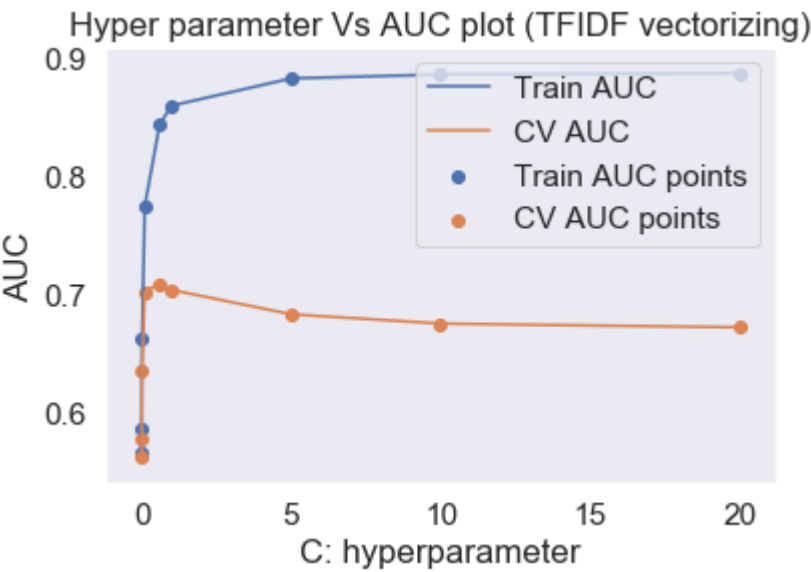
[Parallel(n_jobs=3)]: Using backend LokyBackend with 3 concurrent workers.

[Parallel(n_jobs=3)]: Done 45 out of 45 | elapsed: 51.5s finished

```
In [78]: 1 results = results.sort_values(['param_C'])
2 RS_alphas=results['param_C']
3 train_auc= results['mean_train_score']
4 train_auc_std= results['std_train_score']
5 cv_auc = results['mean_test_score']
6 cv_auc_std= results['std_test_score']
```

In [79]: ▶

```
1 plt.plot(RS_alphas, train_auc, label='Train AUC')
2 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
3 # plt.gca().fill_between(K, train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
4
5 plt.plot(RS_alphas, cv_auc, label='CV AUC')
6 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
7 # plt.gca().fill_between(K, cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
8
9 plt.scatter(RS_alphas, train_auc, label='Train AUC points')
10 plt.scatter(RS_alphas, cv_auc, label='CV AUC points')
11
12 plt.legend()
13 plt.xlabel("C: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("Hyper parameter Vs AUC plot (TFIDF vectorizing)")
16 plt.grid()
17 plt.show()
18
19 results.head()
```



Out[79]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	params	split0_test_score	split1_test_score	split2_test_score	split3_test_score	...	mean_test_score	std_test_score	rank_test_score	split0_tr
0	0.546082	0.027421	0.013334	0.003953	0.0001	{'C': 0.0001}	0.570097	0.558643	0.564393	0.562584	...	0.563047	0.004089	9	
1	0.638051	0.160989	0.013557	0.001385	0.001	{'C': 0.001}	0.588894	0.573337	0.580318	0.578605	...	0.578912	0.005712	8	
2	1.859549	0.119618	0.017579	0.001035	0.01	{'C': 0.01}	0.648677	0.633101	0.636064	0.636344	...	0.636536	0.006696	7	
3	4.323542	0.381404	0.013335	0.000560	0.1	{'C': 0.1}	0.711483	0.701515	0.698013	0.703077	...	0.702781	0.004668	3	
4	4.528854	0.160870	0.014351	0.001719	0.6	{'C': 0.6}	0.716366	0.710137	0.704070	0.706579	...	0.709298	0.004133	1	

5 rows × 21 columns

Observations

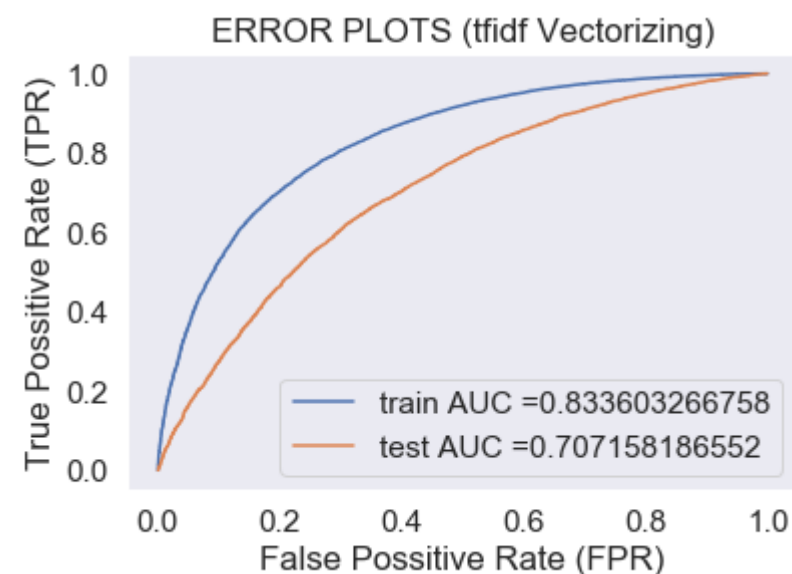
1. Using Grid Search technique we find that as C values increase after 0.001, AUC decreases steeply and maintains a steadiness beyond that.
2. Clearly we see the case of overfit from the graph when C values are beyond 1 (i.e. very high train AUC and low CV AUC)

6.3 Train the model using the best Hyperparameter value

```
In [81]: 1 ### https://forums.fast.ai/t/hyperparameter-random-search-interpretation/8591 ---to get the best hyper parameter as a result of Random search
2 best_C = clf.best_params_
3 print('Best Alpha as a result of Grid Search', best_C)
```

Best Alpha as a result of Grid Search {'C': 0.6}

```
In [83]: 1 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2 from sklearn.metrics import roc_curve, auc
3
4
5 LR=LogisticRegression(penalty="l2", C=0.6, n_jobs=3)
6 LR.fit(X_tr, y_train)
7 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
8 # not the predicted outputs
9
10 y_train_pred = batch_predict(LR, X_tr)
11 y_test_pred = batch_predict(LR, X_te)
12
13 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
14 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
15
16 plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
17 plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
18 plt.legend()
19 plt.xlabel("False Positive Rate (FPR)")
20 plt.ylabel("True Positive Rate (TPR)")
21 plt.title("ERROR PLOTS (tfidf Vectorizing)")
22 plt.grid()
23 plt.show()
```



Observations :

1. Test AUC found to be 0.71 and Train AUC as 0.83 after training model using best hyperparameter 1 and vectorizing text data using TFIDF
2. TFIDF performs better than Bag of Words .

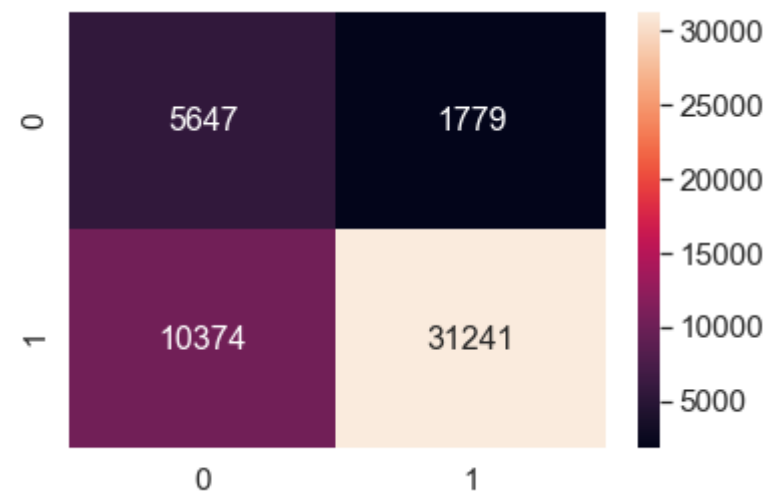
6.4 Confusion Matrix

```
In [84]: 1 print("="*100)
2 from sklearn.metrics import confusion_matrix
3 best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
4 print("Train confusion matrix")
5 print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
6 print("Test confusion matrix")
7 print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

```
=====
the maximum value of tpr*(1-fpr) 0.570870854274 for threshold 0.829
Train confusion matrix
[[ 5647  1779]
 [10374 31241]]
Test confusion matrix
[[ 3189  2270]
 [ 8580 22013]]
```

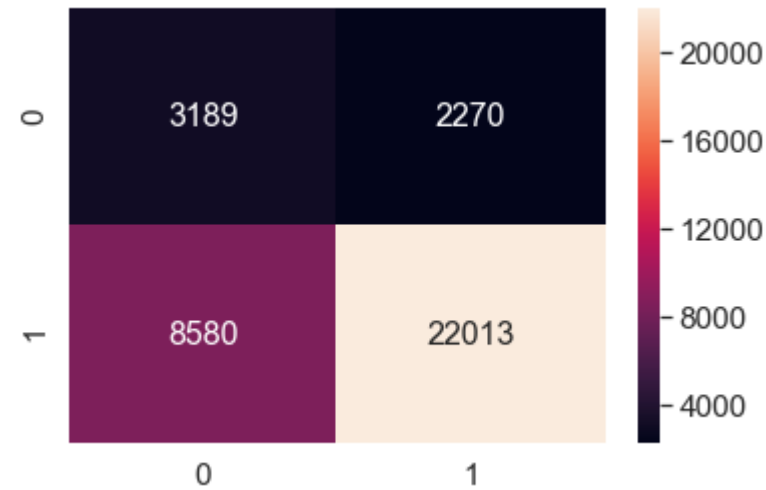
6.4.1 Plotting Confusion Matrix on Train data

```
In [85]: 1 ### PLOT the matrix for Train
2 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3 df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
4 # plt.figure(figsize=(10,7))
5 sns.set(font_scale=1.4) # for label size
6 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7 plt.show()
```



6.4.2 Plotting Confusion Matrix on Test data


```
In [86]: 1 ### PLOT the matrix for Train
2 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3 df_cm = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
4 # plt.figure(figsize=(10,7))
5 sns.set(font_scale=1.4) # for label size
6 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7 plt.show()
```



Observations :

- 1.We can observe from train and test we are getting majority True positives
- 2.Least number of data falls in False negative,which refers as least number of projects were incorrectly predicted as not approved in both Test and Train.
- 3.For a model to perform well we need High True Positive Rate and Low False Positive Rate.From the above our train data has True Positive Rate as 95% and False Positive Rate as 65%
- 4.In our test data : True Positive Rate as 91% and False Positive Rate as 73%.
- 5.Test data AUC is not close to train AUC.
- 6.TFIDF performs better than Bag of words as it has higher % of true positives rate and lower % of False positive rates .

7. SET 3 : categorical, numerical features + project_titleAVG W2V)+ preprocessed_eassay (AVG W2V)

```

In [87]: 1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2 from scipy.sparse import hstack
3
4
5 X_tr = hstack((train_categories, train_subcategories, sklstate_train, teacher_prefix_train,
6               proj_grade_train, avg_w2v_vectors_train, avg_w2v_vectors_title_train,
7               X_train_price_norm, quantity_train_norm, prev_projects_train_norm, title_word_count_train_norm,
8               essay_word_count_train_norm)).tocsr()
9
10 X_te = hstack((test_categories, test_subcategories, sklstate_test, teacher_prefix_test,
11               proj_grade_test, avg_w2v_vectors_test, avg_w2v_vectors_title_test,
12               X_test_price_norm, quantity_test_norm, prev_projects_test_norm, title_word_count_test_norm,
13               essay_word_count_test_norm)).tocsr()
14
15 X_cr = hstack((cv_categories, cv_subcategories, sklstate_cv, teacher_prefix_cv,
16               proj_grade_cv, avg_w2v_vectors_cv, avg_w2v_vectors_title_cv,
17               X_cv_price_norm, quantity_cv_norm, prev_projects_cv_norm, title_word_count_cv_norm,
18               essay_word_count_cv_norm)).tocsr()
19
20
21 print(X_tr.shape)
22 print(X_te.shape)
23 print(X_cr.shape)

```

```

(49041, 704)
(36052, 704)
(24155, 704)

```

```

In [88]: 1 print("Final Data matrix")
2 print(X_tr.shape, y_train.shape)
3 print(X_cr.shape, y_cv.shape)
4 print(X_te.shape, y_test.shape)
5 print("="*100)

```

```

Final Data matrix
(49041, 704) (49041,)
(24155, 704) (24155,)
(36052, 704) (36052,)
=====

```

7.1 Write own function to find which alpha gives maximum auc

In [91]:

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import roc_auc_score
3 ## Lets consider a set of alphas from 10 ** -4 to 10 ** 4
4 def batch_predict(clf, data):
5     # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
6     # not the predicted outputs
7
8     y_data_pred = []
9     tr_loop = data.shape[0] - data.shape[0]%1000
10    # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
11    # in this for loop we will iterate until the last 1000 multiplier
12    for i in range(0, tr_loop, 1000):
13        y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
14    # we will be predicting for the last data points
15    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
16
17    return y_data_pred
18 c = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100]
19 train_auc = []
20 cv_auc = []
21 ## for each alpha now train the model multiple times and store the cv and train auc for plotting it later.
22 for i in tqdm(c):
23     LR = LogisticRegression(penalty="l2", C=i, n_jobs=3)
24     LR.fit(X_tr, y_train)
25     y_train_pred = batch_predict(LR, X_tr)
26     y_cv_pred = batch_predict(LR, X_cr)
27     train_auc.append(roc_auc_score(y_train, y_train_pred))
28     cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

```

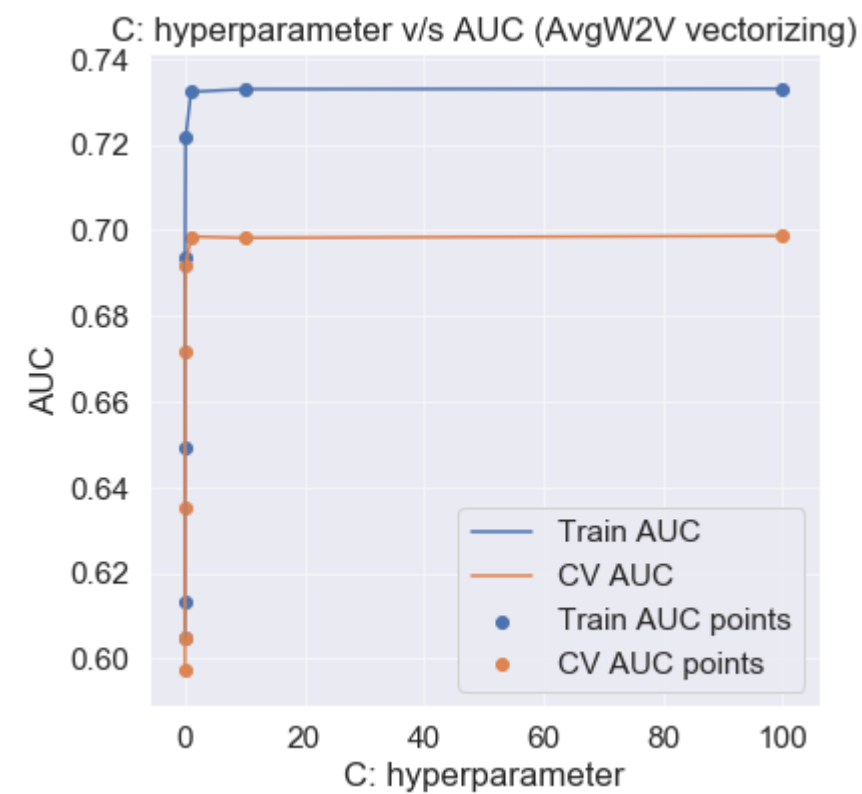
100% | 8/8 [01:57<00:00, 14.67s/it]

In [93]:

```

1 plt.figure(figsize=(6,6))
2
3
4 plt.plot(c, train_auc, label='Train AUC')
5 plt.plot(c, cv_auc, label='CV AUC')
6
7 plt.scatter(c, train_auc, label='Train AUC points')
8 plt.scatter(c, cv_auc, label='CV AUC points')
9
10
11 plt.legend()
12 plt.xlabel("C: hyperparameter")
13 plt.ylabel("AUC")
14 plt.title("C: hyperparameter v/s AUC (AvgW2V vectorizing)")
15 plt.grid(which='major', alpha=0.5)
16 plt.grid(which='minor', alpha=0.2)
17 plt.show()

```



Observations

1. The c values are taken in the range of 10^{-5} to 10^2
2. Model performs well in test as well unlike TFIDF and BOW as we don't see any drop in AUC as C value increases.
4. As C value increases more than 10, we see the curve tends to maintain a constant AUC.
5. Optimal c value which can be selected in this can be 10.

7.2 GridSearch CV using K-fold Crossvalidation with k=10

```
In [98]: 1 from sklearn.model_selection import GridSearchCV
2 from scipy.stats import randint as sp_randint
3 from sklearn.model_selection import RandomizedSearchCV
4 from sklearn.linear_model import LogisticRegression
5 ## As we are seeing steady decreasing in AUC value in CV as c value increases more than 1 .Hence keep more values less
6 ## than 1
7 ## time complexity increases as number of params increases hence reducing the cv to 3
8 ## better to choose parameters below 15 as we saw as number of parameter increases beyond 10 curve maintains steadiness.
9 parameters={"C" : [0.0001,0.001,0.01,0.1,0.6,15] }
10
11 clf = GridSearchCV(LogisticRegression(),parameters,return_train_score=True, cv=3, scoring='roc_auc',verbose=2,n_jobs=3)
12 clf.fit(X_tr, y_train)
13 results = pd.DataFrame.from_dict(clf.cv_results_)
14
```

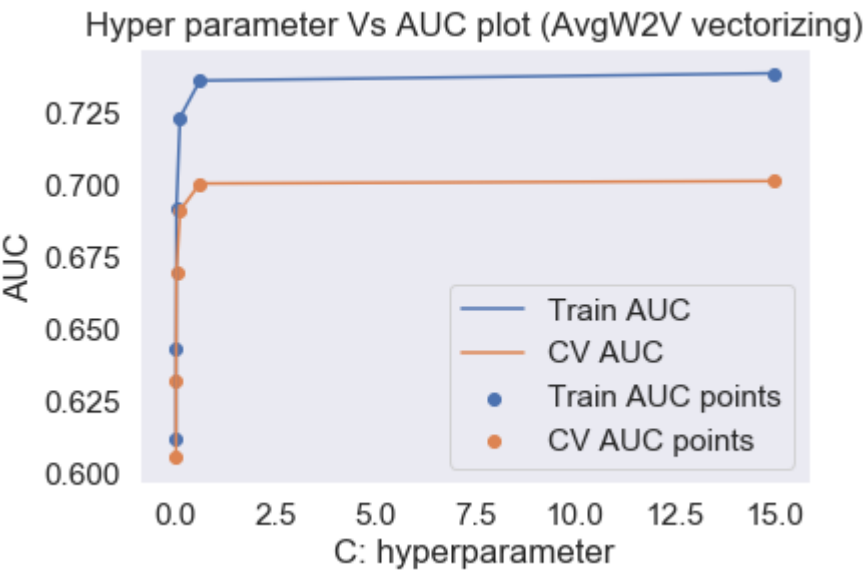
Fitting 3 folds for each of 6 candidates, totalling 18 fits

```
[Parallel(n_jobs=3)]: Using backend LokyBackend with 3 concurrent workers.
[Parallel(n_jobs=3)]: Done 18 out of 18 | elapsed: 57.0s finished
```

```
In [99]: 1 results = results.sort_values(['param_C'])
2 RS_alphas=results['param_C']
3 train_auc= results['mean_train_score']
4 train_auc_std= results['std_train_score']
5 cv_auc = results['mean_test_score']
6 cv_auc_std= results['std_test_score']
```

In [101]: ▶

```
1 plt.plot(RS_alphas, train_auc, label='Train AUC')
2 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
3 # plt.gca().fill_between(K, train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
4
5 plt.plot(RS_alphas, cv_auc, label='CV AUC')
6 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
7 # plt.gca().fill_between(K, cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
8
9 plt.scatter(RS_alphas, train_auc, label='Train AUC points')
10 plt.scatter(RS_alphas, cv_auc, label='CV AUC points')
11
12 plt.legend()
13 plt.xlabel("C: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("Hyper parameter Vs AUC plot (AvgW2V vectorizing)")
16 plt.grid()
17 plt.show()
18
19 results.head()
```



Out[101]:

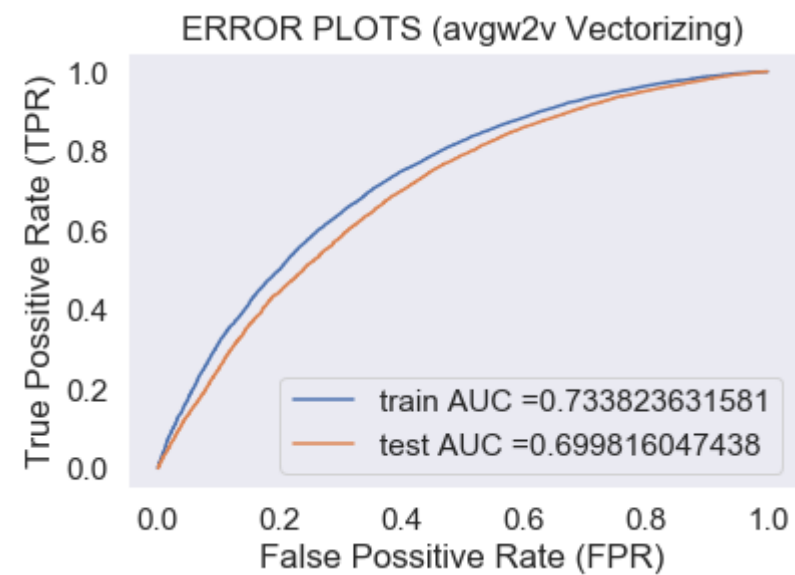
	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	params	split0_test_score	split1_test_score	split2_test_score	mean_test_score	std_test_score	rank_test_score	split0_train_score	split1_train_
0	2.598803	0.076121	0.061490	3.886097e-03	0.0001	{'C': 0.0001}	0.604057	0.614561	0.597704	0.605441	0.006951	6	0.612198	0.6
1	3.588904	0.093584	0.064000	1.325077e-06	0.001	{'C': 0.001}	0.631559	0.639590	0.624728	0.631959	0.006074	5	0.642570	0.6
2	9.496942	0.752247	0.060000	7.867412e-07	0.01	{'C': 0.01}	0.670737	0.673017	0.665133	0.669629	0.003312	4	0.690365	0.6
3	12.769815	0.158370	0.068156	3.191006e-03	0.1	{'C': 0.1}	0.692827	0.690926	0.689305	0.691019	0.001439	3	0.721735	0.7
4	13.127893	0.268059	0.062734	1.547107e-02	0.6	{'C': 0.6}	0.703818	0.700155	0.697212	0.700395	0.002702	2	0.734615	0.7

7.3 Train the model using the best Hyperparameter value

```
In [104]: 1 ### https://forums.fast.ai/t/hyperparameter-random-search-interpretation/8591 ---to get the best hyper parameter as a result of Random search
2 best_C = clf.best_params_
3 print('Best Alpha as a result of Grid Search',best_C)
```

Best Alpha as a result of Grid Search {'C': 15}

```
In [106]: 1 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2 from sklearn.metrics import roc_curve, auc
3
4
5 LR=LogisticRegression(penalty="l2",C=15,n_jobs=3)
6 LR.fit(X_tr, y_train)
7 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
8 # not the predicted outputs
9
10 y_train_pred = batch_predict(LR, X_tr)
11 y_test_pred = batch_predict(LR, X_te)
12
13 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
14 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
15
16 plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
17 plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
18 plt.legend()
19 plt.xlabel("False Possitive Rate (FPR)")
20 plt.ylabel("True Possitive Rate (TPR)")
21 plt.title("ERROR PLOTS (avgw2v Vectorizing)")
22 plt.grid()
23 plt.show()
```



Observations :

1. Test AUC found to be 0.69 and Train AUC as 0.73 after training model using best hyperparameter 1 and vectorizing text data using AvgW2V
2. TFIDF performs better than Bag of Words and AvgW2V.

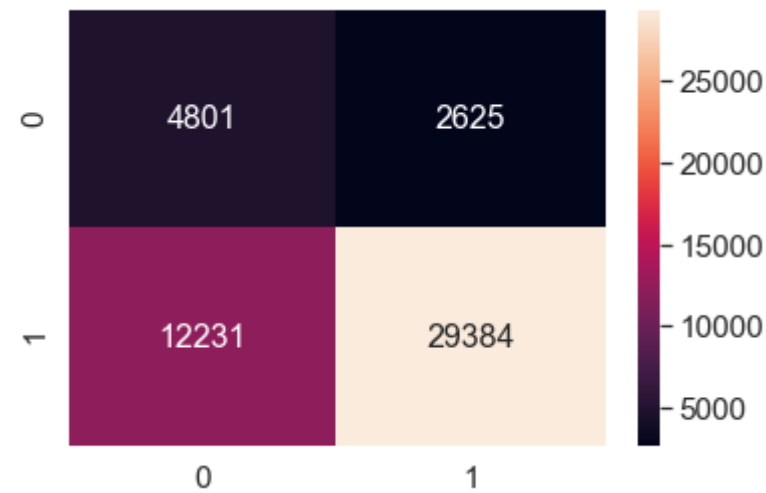
7.4 Confusion Matrix

```
In [107]: 1 print("-"*100)
2 from sklearn.metrics import confusion_matrix
3 best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
4 print("Train confusion matrix")
5 print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
6 print("Test confusion matrix")
7 print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

```
=====
the maximum value of tpr*(1-fpr) 0.456496841971 for threshold 0.836
Train confusion matrix
[[ 4801  2625]
 [12231 29384]]
Test confusion matrix
[[ 3282  2177]
 [ 9210 21383]]
```

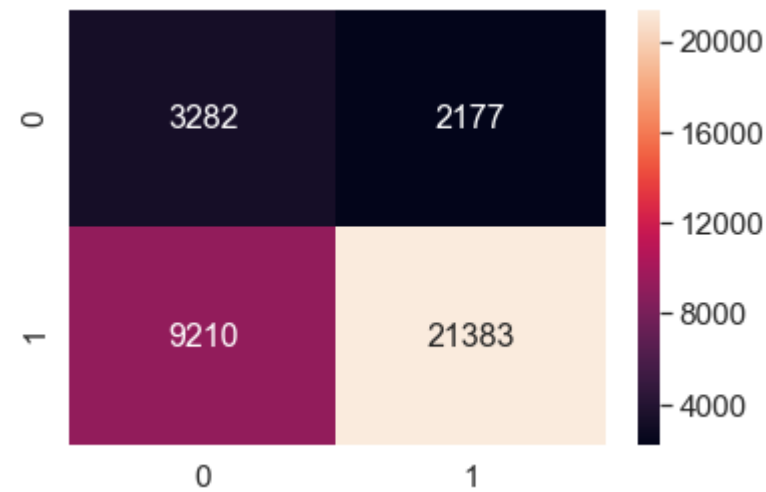
7.4.1 Ploting Confusion Matrix on Train data

```
In [108]: 1 ### PLOT the matrix for Train
2 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3 df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
4 # plt.figure(figsize=(10,7))
5 sns.set(font_scale=1.4) # for label size
6 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7 plt.show()
```



7.4.2 Plotting Confusion Matrix on Test data


```
In [109]: 1 ### PLOT the matrix for Train
2 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3 df_cm = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
4 # plt.figure(figsize=(10,7))
5 sns.set(font_scale=1.4) # for label size
6 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16}, fmt='g') # font size
7 plt.show()
```



Observations :

1. We can observe from train and test we are getting majority True positives
2. Least number of data falls in False negative, which refers as least number of projects were incorrectly predicted as not approved in both Test and Train.
3. For a model to perform well we need High True Positive Rate and Low False Positive Rate. From the above our train data has True Positive Rate as 91% and False Positive Rate as 71.8%
4. In our test data : True Positive Rate as 91% and False Positive Rate as 73.7%.

8. SET 4 : categorical, numerical features + project_title(tfidf W2V)+ preprocessed_eassay (tfidf W2V)

```

In [110]: 1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
          2
          3
          4 X_tr = hstack((train_categories, train_subcategories, sklstate_train, teacher_prefix_train,
          5                   proj_grade_train, tfidf_w2v_vectors_train, tfidf_w2v_vectors_title_train,
          6                   X_train_price_norm, quantity_train_norm, prev_projects_train_norm, title_word_count_train_norm,
          7                   essay_word_count_train_norm)).tocsr()
          8
          9 X_te = hstack((test_categories, test_subcategories, sklstate_test, teacher_prefix_test,
         10                   proj_grade_test, tfidf_w2v_vectors_test, tfidf_w2v_vectors_title_test,
         11                   X_test_price_norm, quantity_test_norm, prev_projects_test_norm, title_word_count_test_norm,
         12                   essay_word_count_test_norm)).tocsr()
         13
         14 X_cr = hstack((cv_categories, cv_subcategories, sklstate_cv, teacher_prefix_cv,
         15                   proj_grade_cv, tfidf_w2v_vectors_cv, tfidf_w2v_vectors_title_cv,
         16                   X_cv_price_norm, quantity_cv_norm, prev_projects_cv_norm, title_word_count_cv_norm,
         17                   essay_word_count_cv_norm)).tocsr()
         18
         19
         20 print(X_tr.shape)
         21 print(X_te.shape)
         22 print(X_cr.shape)

```

```

(49041, 704)
(36052, 704)
(24155, 704)

```

```

In [111]: 1 print("Final Data matrix")
          2 print(X_tr.shape, y_train.shape)
          3 print(X_cr.shape, y_cv.shape)
          4 print(X_te.shape, y_test.shape)
          5 print("="*100)

```

```

Final Data matrix
(49041, 704) (49041,)
(24155, 704) (24155,)
(36052, 704) (36052,)
=====

```

8.1 Write own function to find which alpha gives maximum auc

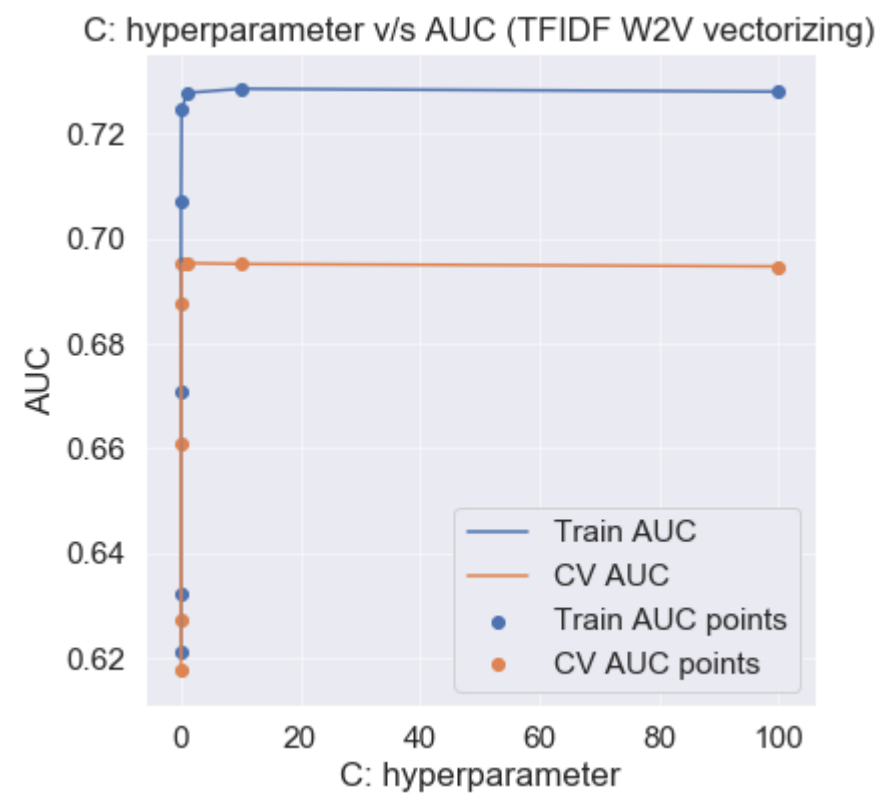
[illegible]

In [115]:

```

1 plt.figure(figsize=(6,6))
2
3
4 plt.plot(c, train_auc, label='Train AUC')
5 plt.plot(c, cv_auc, label='CV AUC')
6
7 plt.scatter(c, train_auc, label='Train AUC points')
8 plt.scatter(c, cv_auc, label='CV AUC points')
9
10
11 plt.legend()
12 plt.xlabel("C: hyperparameter")
13 plt.ylabel("AUC")
14 plt.title("C: hyperparameter v/s AUC (TFIDF W2V vectorizing)")
15 plt.grid(which='major', alpha=0.5)
16 plt.grid(which='minor', alpha=0.2)
17 plt.show()

```



Observations

1. The c values are taken in the range of 10^{-5} to 10^2
2. Model performs well in test as well unlike TFIDF and BOW as we don't see any drop in AUC as C value increases.
3. As C value increases more than 10, we see the curve tends to maintain a constant AUC.
4. Optimal c value which can be selected in this can be 10.

8.2 GridSearch CV using K-fold Crossvalidation with k=10

```
In [118]: 1 from sklearn.model_selection import GridSearchCV
2 from scipy.stats import randint as sp_randint
3 from sklearn.model_selection import RandomizedSearchCV
4
5 ## As we are seeing steady decreasing in AUC value in CV as c value increases more than 1 .Hence keep more values less
6 ## than 1
7 parameters={"C" : [0.0001,0.001,0.01,0.1,0.6,15] }
8
9 clf = GridSearchCV(LogisticRegression(), parameters,return_train_score=True, cv=3, scoring='roc_auc',verbose=1,n_jobs=3)
10 clf.fit(X_tr, y_train)
11 results = pd.DataFrame.from_dict(clf.cv_results_)
12
```

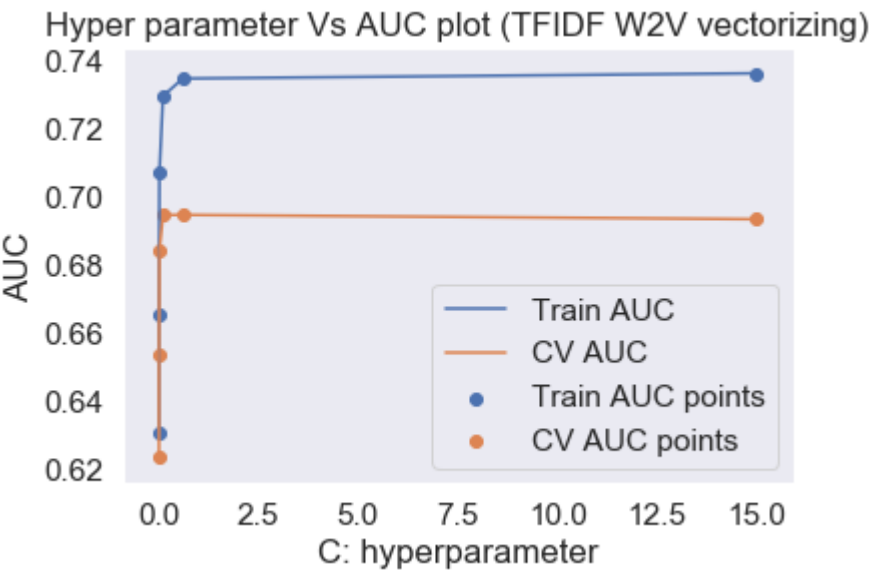
Fitting 3 folds for each of 6 candidates, totalling 18 fits

[Parallel(n_jobs=3)]: Using backend LokyBackend with 3 concurrent workers.
[Parallel(n_jobs=3)]: Done 18 out of 18 | elapsed: 50.0s finished

```
In [119]: 1 results = results.sort_values(['param_C'])
2 RS_alphas=results['param_C']
3 train_auc= results['mean_train_score']
4 train_auc_std= results['std_train_score']
5 cv_auc = results['mean_test_score']
6 cv_auc_std= results['std_test_score']
```

In [121]: ▶

```
1 plt.plot(RS_alphas, train_auc, label='Train AUC')
2 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
3 # plt.gca().fill_between(K, train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
4
5 plt.plot(RS_alphas, cv_auc, label='CV AUC')
6 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
7 # plt.gca().fill_between(K, cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
8
9 plt.scatter(RS_alphas, train_auc, label='Train AUC points')
10 plt.scatter(RS_alphas, cv_auc, label='CV AUC points')
11
12 plt.legend()
13 plt.xlabel("C: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("Hyper parameter Vs AUC plot (TFIDF W2V vectorizing)")
16 plt.grid()
17 plt.show()
18
19 results.head()
```



Out[121]:

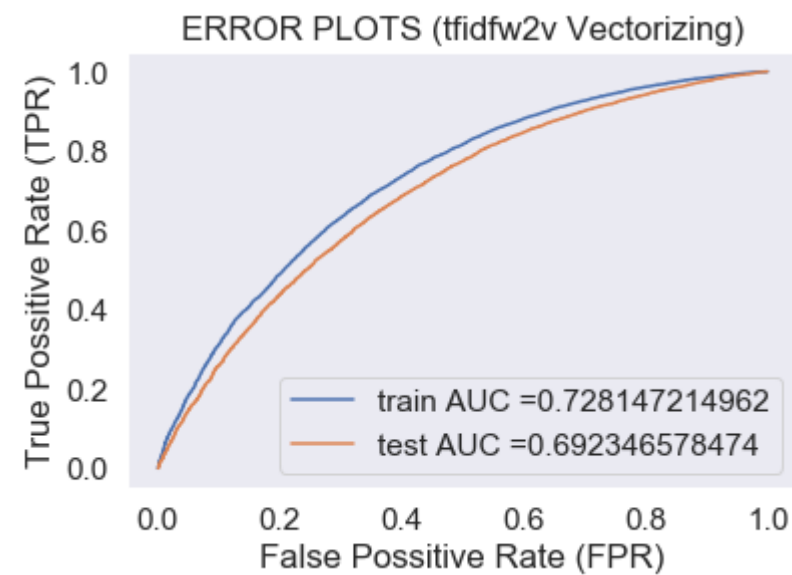
	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	params	split0_test_score	split1_test_score	split2_test_score	mean_test_score	std_test_score	rank_test_score	split0_train_score	split1_train_
0	2.088079	0.083738	0.058178	0.008707	0.0001	{'C': 0.0001}	0.622038	0.634861	0.612612	0.623170	0.009119	6	0.630641	0.6
1	4.141592	0.320524	0.053860	0.016224	0.001	{'C': 0.001}	0.652642	0.662772	0.645173	0.653529	0.007212	5	0.665051	0.6
2	9.961026	0.448239	0.050200	0.000470	0.01	{'C': 0.01}	0.684188	0.688978	0.679134	0.684100	0.004019	4	0.706368	0.7
3	10.649860	0.211833	0.050188	0.000455	0.1	{'C': 0.1}	0.696801	0.696148	0.690633	0.694527	0.002766	2	0.728341	0.7
4	10.340343	0.170759	0.053215	0.004685	0.6	{'C': 0.6}	0.697518	0.695351	0.690891	0.694586	0.002759	1	0.732455	0.7

8.3 Train the model using the best Hyperparameter value

```
In [122]: 1 ### https://forums.fast.ai/t/hyperparameter-random-search-interpretation/8591 ---to get the best hyper parameter as a result of Random search
2 best_C = clf.best_params_
3 print('Best Alpha as a result of Grid Search',best_C)
```

Best Alpha as a result of Grid Search {'C': 0.6}

```
In [124]: 1 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2 from sklearn.metrics import roc_curve, auc
3
4
5 LR=LogisticRegression(penalty="l2",C=best_C['C'],n_jobs=3)
6 LR.fit(X_tr, y_train)
7 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
8 # not the predicted outputs
9
10 y_train_pred = batch_predict(LR, X_tr)
11 y_test_pred = batch_predict(LR, X_te)
12
13 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
14 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
15
16 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
17 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
18 plt.legend()
19 plt.xlabel("False Possitive Rate (FPR)")
20 plt.ylabel("True Possitive Rate (TPR)")
21 plt.title("ERROR PLOTS (tfidf2v Vectorizing)")
22 plt.grid()
23 plt.show()
```



Observations :

1. Test AUC found to be 0.69 and Train AUC as 0.72 after training model using best hyperparameter 0.6 and vectorizing text data using TFIDF

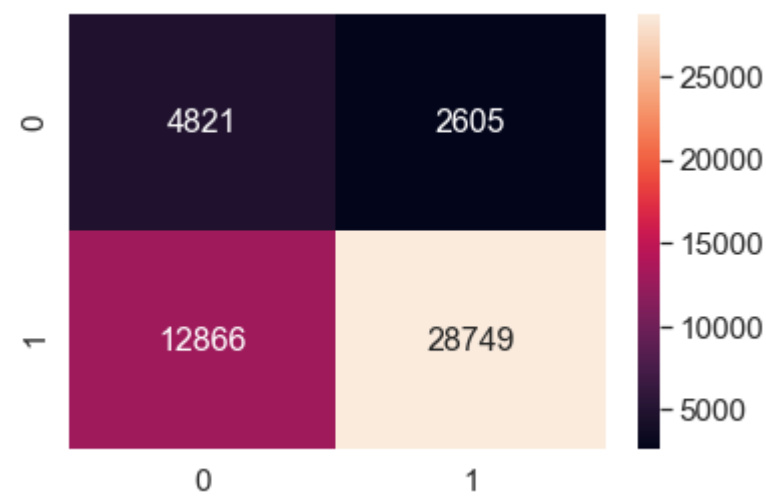
8.4 Confusion Matrix

```
In [125]: 1 print("="*100)
2 from sklearn.metrics import confusion_matrix
3 best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
4 print("Train confusion matrix")
5 print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
6 print("Test confusion matrix")
7 print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

```
=====
the maximum value of tpr*(1-fpr) 0.448492340575 for threshold 0.834
Train confusion matrix
[[ 4821  2605]
 [12866 28749]]
Test confusion matrix
[[ 3300  2159]
 [ 9810 20783]]
```

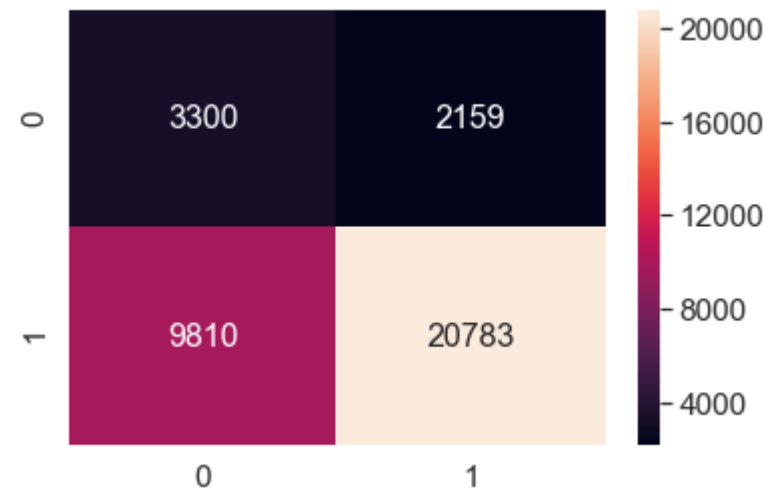
8.4.1 Ploting Confusion Matrix on Train data

```
In [126]: 1 ### PLOT the matrix for Train
2 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3 df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
4 # plt.figure(figsize=(10,7))
5 sns.set(font_scale=1.4) # for label size
6 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7 plt.show()
```



8.4.2 Ploting Confusion Matrix on Test data


```
In [127]: 1 ### PLOT the matrix for Train
2 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3 df_cm = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
4 # plt.figure(figsize=(10,7))
5 sns.set(font_scale=1.4) # for label size
6 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7 plt.show()
```



Observations :

1. We can observe from train and test we are getting majority True positives
2. Least number of data falls in False negative, which refers as least number of projects were incorrectly predicted as not approved in both Test and Train.
3. For a model to perform well we need High True Positive Rate and Low False Positive Rate. From the above our train data has True Positive Rate as 93% and False Positive Rate as 65%
4. In our test data : True Positive Rate as 91% and False Positive Rate as 71%.
5. Test data AUC is not close to train AUC.
6. TFIDF performs better than Bag of words as it has higher % of true positives rate and lower % of False positive rates .

9. Set5 . Apply Logistic Regression on the specified feature set Set 5 by finding the best hyper parameter as suggested in step 2 and step 3

```

In [128]: ▶ 1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2
3 X_tr = hstack((train_categories, train_subcategories,sklstate_train,teacher_prefix_train,
4               proj_grade_train,sentiment_neg_train_norm,
5               X_train_price_norm,quantity_train_norm,
6               prev_projects_train_norm,title_word_count_train_norm,
7               essay_word_count_train_norm)).tocsr()
8
9 X_te = hstack((test_categories, test_subcategories,sklstate_test,teacher_prefix_test,
10              proj_grade_test,sentiment_neg_test_norm,
11              X_test_price_norm,quantity_test_norm,
12              prev_projects_test_norm,title_word_count_test_norm,
13              essay_word_count_test_norm)).tocsr()
14
15 X_cr = hstack((cv_categories, cv_subcategories,sklstate_cv,teacher_prefix_cv,
16              proj_grade_cv,sentiment_neg_cv_norm,
17              X_cv_price_norm,quantity_cv_norm,
18              prev_projects_cv_norm,title_word_count_cv_norm,
19              essay_word_count_cv_norm)).tocsr()
20
21
22 print(X_tr.shape)
23 print(X_te.shape)
24 print(X_cr.shape)

```

```

(49041, 105)
(36052, 105)
(24155, 105)

```

Training model with different c values

```

In [130]: ▶ 1
2 ### Train the model with loops
3 train_auc=[]
4 cv_auc=[]
5 c=[0.00001,0.0001,0.001,0.01,0.1,1,10,20,50,100,500,1000]
6 for i in tqdm(c):
7     LR=LogisticRegression(C=i,penalty='l2',n_jobs=3)
8     LR.fit(X_tr,y_train)
9     train_pred=batch_predict(LR,X_tr)
10    cv_pred=batch_predict(LR,X_cr)
11    train_auc.append(roc_auc_score(y_train,train_pred))
12    cv_auc.append(roc_auc_score(y_cv,cv_pred))

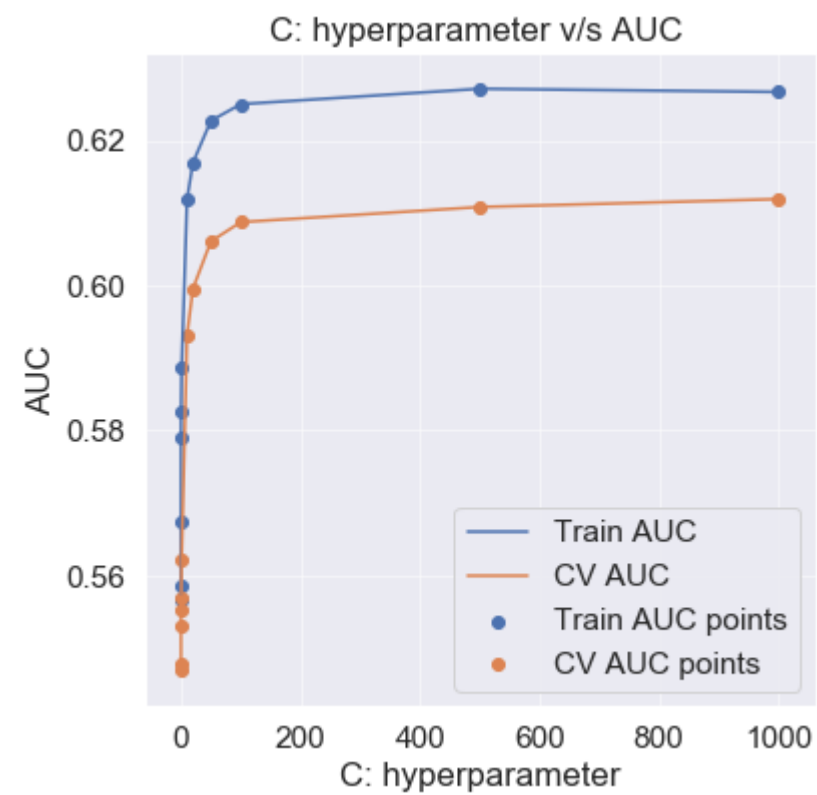
```

```

100%|████████████████████████████████████████████████████████████████████████████████| 12/12 [00:07<00:00, 1.60it/s]

```

```
In [131]: ▶ 1 plt.figure(figsize=(6,6))
2
3 plt.plot(c, train_auc, label='Train AUC')
4 plt.plot(c, cv_auc, label='CV AUC')
5
6 plt.scatter(c, train_auc, label='Train AUC points')
7 plt.scatter(c, cv_auc, label='CV AUC points')
8
9
10 plt.legend()
11 plt.xlabel("C: hyperparameter")
12 plt.ylabel("AUC")
13 plt.title("C: hyperparameter v/s AUC ")
14 plt.grid(which='major', alpha=0.5)
15 plt.grid(which='minor', alpha=0.2)
16 plt.show()
```



Training model using Gridsearchcv

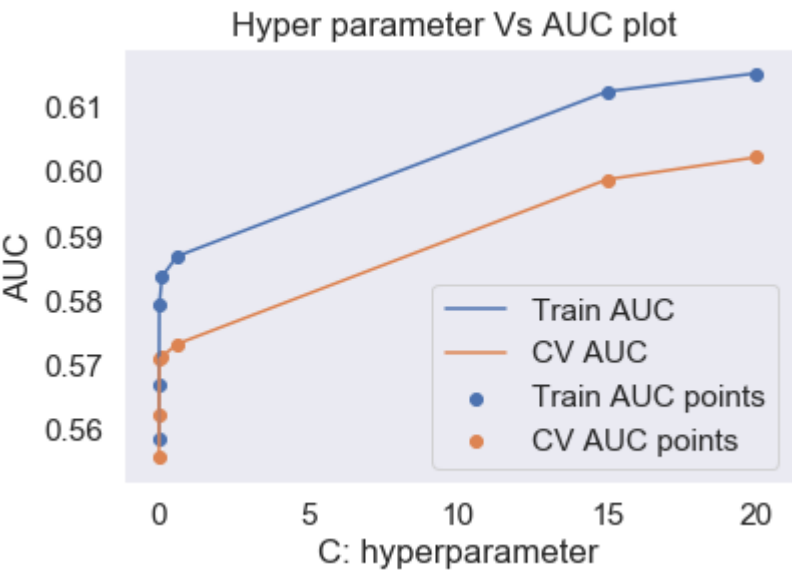
```
In [138]: 1 parameters = {"C" : [0.0001,0.001,0.01,0.1,0.6,15,20] }  
2  
3 grid_model=GridSearchCV(LR,parameters,return_train_score=True,n_jobs=3,scoring='roc_auc',cv=5)  
4 grid_model.fit(X_tr,y_train)  
5
```

```
Out[138]: GridSearchCV(cv=5, estimator=LogisticRegression(C=1000, n_jobs=3), n_jobs=3,  
                  param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 0.6, 15, 20]},  
                  return_train_score=True, scoring='roc_auc')
```

```
In [139]: 1 results=pd.DataFrame.from_dict(grid_model.cv_results_)  
2 RS_alphas=results['param_C']  
3 train_auc=results['mean_train_score']  
4 cv_auc=results['mean_test_score']
```

In [140]: ▶

```
1 plt.plot(RS_alphas, train_auc, label='Train AUC')
2 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
3 # plt.gca().fill_between(K, train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
4
5 plt.plot(RS_alphas, cv_auc, label='CV AUC')
6 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
7 # plt.gca().fill_between(K, cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
8
9 plt.scatter(RS_alphas, train_auc, label='Train AUC points')
10 plt.scatter(RS_alphas, cv_auc, label='CV AUC points')
11
12 plt.legend()
13 plt.xlabel("C: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("Hyper parameter Vs AUC plot ")
16 plt.grid()
17 plt.show()
18
19 results.head()
```



Out[140]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	params	split0_test_score	split1_test_score	split2_test_score	split3_test_score	...	mean_test_score	std_test_score	rank_test_score	split0_tr
0	0.119279	0.001493	0.004589	0.000798	0.0001	{'C': 0.0001}	0.562555	0.550991	0.556968	0.555057	...	0.555516	0.004114	7	
1	0.119480	0.000746	0.005785	0.001162	0.001	{'C': 0.001}	0.572008	0.556316	0.564359	0.561480	...	0.562257	0.005689	6	
2	0.243848	0.039447	0.007187	0.000746	0.01	{'C': 0.01}	0.583984	0.565542	0.572532	0.567999	...	0.570782	0.007219	5	
3	0.545133	0.075107	0.005386	0.001197	0.1	{'C': 0.1}	0.585985	0.565890	0.572392	0.567554	...	0.571192	0.007891	4	
4	0.540406	0.040533	0.004382	0.000482	0.6	{'C': 0.6}	0.588760	0.567409	0.574234	0.569279	...	0.573155	0.008277	3	

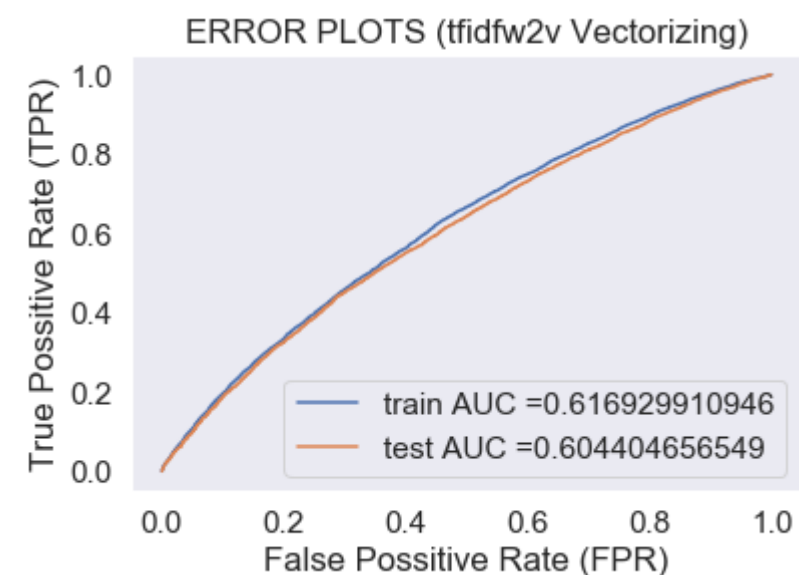
5 rows × 21 columns

```
In [141]: 1 ### https://forums.fast.ai/t/hyperparameter-random-search-interpretation/8591 ---to get the best hyper parameter as a result of Random search
2 best_C = grid_model.best_params_
3 print('Best Alpha as a result of Grid Search',best_C)
```

Best Alpha as a result of Grid Search {'C': 20}

Train the model using the best parameter

```
In [143]: 1 LR = LogisticRegression(C=best_C['C'],penalty='l2',n_jobs=3)
2 LR.fit(X_tr,y_train)
3 # Predict
4 y_train_pred=batch_predict(LR,X_tr)
5 y_test_pred=batch_predict(LR,X_te)
6 # FPR,TPR,thresholds
7 train_fpr,train_tpr,tr_thresholds=roc_curve(y_train,y_train_pred)
8 test_fpr,test_tpr,te_thresholds=roc_curve(y_test,y_test_pred)
9 #
10 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("False Possitive Rate (FPR)")
14 plt.ylabel("True Possitive Rate (TPR)")
15 plt.title("ERROR PLOTS (tfidf2v Vectorizing)")
16 plt.grid()
17 plt.show()
18
```



Plotting Confusion Matrix on Train

```
In [144]: 1 best_t=best_threshold(tr_thresholds,train_fpr,train_tpr)
2 print("Train confusion matrix")
3 print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
4 print("Test confusion matrix")
5 print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
6 df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.342058548506 for threshold 0.844

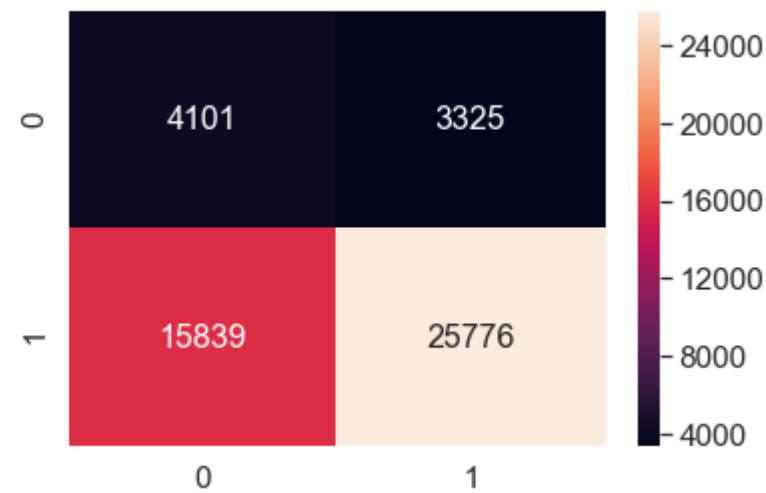
Train confusion matrix

```
[[ 4101  3325]
 [15839 25776]]
```

Test confusion matrix

```
[[ 2964  2495]
 [12110 18483]]
```

```
In [145]: 1 # plt.figure(figsize=(10,7))
2 sns.set(font_scale=1.4) # for label size
3 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
4 plt.show()
```



In [146]: ▶

```
1 df_cm = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
2 # plt.figure(figsize=(10,7))
3 sns.set(font_scale=1.4) # for label size
4 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
5 plt.show()
```



In [148]: ▶

```
1 ##http://zetcode.com/python/prettytable/
2
3 from prettytable import PrettyTable
4
5 x = PrettyTable()
6 x.field_names = ["Vectorizer", "Model", "Hyper Parameter", " Train AUC" ,"Test AUC "]
7
8 x.add_row(["BOW", "Logistic Regression", 0.01, 0.81,0.71])
9 x.add_row(["TFIDF", "Logistic Regression",0.6, 0.83, 0.70])
10 x.add_row(["AvgW2v", "Logistic Regression",15, 0.73, 0.69])
11 x.add_row(["TFIDF w2v", "Logistic Regression",0.6, 0.72, 0.69])
12 x.add_row(["Without text", "Logistic Regression",20, 0.61, 0.60])
13 print(x)
```

Vectorizer	Model	Hyper Parameter	Train AUC	Test AUC
BOW	Logistic Regression	0.01	0.81	0.71
TFIDF	Logistic Regression	0.6	0.83	0.7
AvgW2v	Logistic Regression	15	0.73	0.69
TFIDF w2v	Logistic Regression	0.6	0.72	0.69
Without text	Logistic Regression	20	0.61	0.6