

# Donors choose data analysis

## Table of contents

- [1. About dataset](#)
- [2. Preparing data for analysis - importing libraries, reading data...](#)
- [3. Univariate analysis](#)
- [4. Pre-processing data](#)
- [5. Vectorizing all features - preparing data for classification and modelling](#)
- [6. Vectorizing data using t-SNE](#)
- [7. Classification & Modelling Using K-NN
  - \[7.1 Classification using k-NN\\(k-fold cross validation\\)
    - \\[7.1.1 Classification using k-NN\\\(k-fold cross validation\\\) on imbalanced data\\]\\(#\\)
    - \\[7.2 Classification using k-NN\\\(k-fold cross validation & feature selection\\\)
      - \\\[7.2.1 Classification using k-NN\\\\(k-fold cross validation & feature\\\\\_selection\\\\) on imbalanced data\\\]\\\(#\\\)
      - \\\[7.2.2 Classification using k-NN\\\\(k-fold cross validation & feature\\\\\_selection\\\\) on balanced data\\\]\\\(#\\\)\\]\\(#\\)
    - \\[7.3 Results of analysis using k-NN\\]\\(#\\)
    - \\[7.4 Conclusions of analysis using k-NN\\]\\(#\\)\]\(#\)](#)

## Little History about Data Set

Founded in 2000 by a high school teacher in the Bronx, DonorsChoose.org empowers public school teachers from across the country to request much-needed materials and experiences for their students. At any given time, there are thousands of classroom requests that can be brought to life with a gift of any amount.

## Answers to What and Why Questions on Data Set

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

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

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

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

## About the DonorsChoose Data Set

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

Feature	Description
<code>project_id</code>	A unique identifier for the proposed project. <b>Example:</b> p036502
<code>project_title</code>	Title of the project. <b>Examples:</b> <ul style="list-style-type: none"><li>• Art Will Make You Happy</li><li>• First Grade Fun</li></ul>

Feature	Description
<b>project_grade_category</b>	<p>Grade level of students for which the project is targeted. One of the following enumerated values:</p> <ul style="list-style-type: none"> <li>• Grades PreK-2</li> <li>• Grades 3-5</li> <li>• Grades 6-8</li> <li>• Grades 9-12</li> </ul>
<b>project_subject_categories</b>	<p>One or more (comma-separated) subject categories for the project from the following enumerated list of values:</p> <ul style="list-style-type: none"> <li>• Applied Learning</li> <li>• Care &amp; Hunger</li> <li>• Health &amp; Sports</li> <li>• History &amp; Civics</li> <li>• Literacy &amp; Language</li> <li>• Math &amp; Science</li> <li>• Music &amp; The Arts</li> <li>• Special Needs</li> <li>• Warmth</li> </ul> <p><b>Examples:</b></p> <ul style="list-style-type: none"> <li>• Music &amp; The Arts</li> <li>• Literacy &amp; Language, Math &amp; Science</li> </ul>
<b>school_state</b>	<p>State where school is located (<a href="#">Two-letter U.S. postal code</a>). <b>Example:</b> WY</p>

Feature	Description
<b>project_subject_subcategories</b>	<p>One or more (comma-separated) subcategories for the project. <b>Example:</b></p> <ul style="list-style-type: none"> <li>• Literacy</li> <li>• Literature &amp; Writing, Social Sciences</li> </ul>
<b>project_resource_summary</b>	<p>An explanation of the resources needed for the project. <b>Example:</b></p> <ul style="list-style-type: none"> <li>• My students need hands-on literacy materials to manage sensory needs!</li> </ul>
<b>project_essay_1</b>	First application essay*
<b>project_essay_2</b>	Second application essay*
<b>project_essay_3</b>	Third application essay*
<b>project_essay_4</b>	Fourth application essay*
<b>project_submitted_datetime</b>	Datetime when project application was submitted. <b>Example:</b> 2016-04-28 12:43:56.245
<b>teacher_id</b>	A unique identifier for the teacher of proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4ff1c

Feature	Description
<b>teacher_prefix</b>	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> <li>• nan</li> <li>• Dr.</li> <li>• Mr.</li> <li>• Mrs.</li> <li>• Ms.</li> <li>• Teacher.</li> </ul>
<b>teacher_number_of_previously_posted_projects</b>	Number of project applications previously submitted by the same teacher. <b>Example:</b> 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
<b>id</b>	A project_id value from the train.csv file. <b>Example:</b> p036502
<b>description</b>	Description of the resource. <b>Example:</b> Tenor Saxophone Reeds, Box of 25
<b>quantity</b>	Quantity of the resource required. <b>Example:</b> 3
<b>price</b>	Price of the resource required. <b>Example:</b> 9.95

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

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

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

## Notes on the Essay Data

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

- \_\_project\_essay\_1:\_\_ "Introduce us to your classroom"
- \_\_project\_essay\_2:\_\_ "Tell us more about your students"
- \_\_project\_essay\_3:\_\_ "Describe how your students will use the materials you're requesting"
- \_\_project\_essay\_4:\_\_ "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.

## Importing required libraries

In [1]: `import pandas as pd;`

In [2]: `# numpy for easy numerical computations`

```
import numpy as np
# pandas for dataframes and filterings
import pandas as pd
# sqlite3 library for performing operations on sqlite file
import sqlite3
# matplotlib for plotting graphs
import matplotlib.pyplot as plt
# seaborn library for easy plotting
import seaborn as sns
# warnings library for specific settings
import warnings
# regularlanguage for regex operations
import re
# For loading precomputed models
import pickle

# For loading files from google drive
from google.colab import drive
# For working with files in google drive
drive.mount('/content/drive')
# tqdm for tracking progress of loops
from tqdm import tqdm_notebook as tqdm
# For creating dictionary of words
from collections import Counter
# For creating BagOfWords Model
from sklearn.feature_extraction.text import CountVectorizer
# For creating TfIdfModel
from sklearn.feature_extraction.text import TfidfVectorizer
# For standardizing values
from sklearn.preprocessing import StandardScaler
# For merging sparse matrices along row direction
from scipy.sparse import hstack
# For merging sparse matrices along column direction
from scipy.sparse import vstack
# For calculating TSNE values
from sklearn.manifold import TSNE
# For calculating the accuracy score on cross validate data
from sklearn.metrics import accuracy_score
# For performing the k-fold cross validation
```

```
from sklearn.model_selection import cross_val_score
# For splitting the data set into test and train data
from sklearn import model_selection
# KNeighbors classifier for classification
from sklearn.neighbors import KNeighborsClassifier
# For creating samples for making dataset balanced
from sklearn.utils import resample
# For shuffling the dataframes
from sklearn.utils import shuffle
# For calculating roc_curve parameters
from sklearn.metrics import roc_curve
# For calculating auc value
from sklearn.metrics import auc
# For displaying results in table format
from prettytable import PrettyTable
# For generating confusion matrix
from sklearn.metrics import confusion_matrix
# For selecting most useful features
from sklearn.feature_selection import SelectKBest, f_classif

warnings.filterwarnings('ignore')
```

Go to this URL in a browser: [https://accounts.google.com/o/oauth2/auth?client\\_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect\\_uri=urn%3Aietf%3Awg%3Aoauth%3A2.0%3Aoob&scope=email%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdocs.test%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fpeopleapi.readonly&response\\_type=code](https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3Aietf%3Awg%3Aoauth%3A2.0%3Aoob&scope=email%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdocs.test%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fpeopleapi.readonly&response_type=code)

Enter your authorization code:

.....

Mounted at /content/drive

## Reading and Storing Data

```
In [0]: projectsData = pd.read_csv('drive/My Drive/train_data.csv');
resourcesData = pd.read_csv('drive/My Drive/resources.csv');
```

In [4]: `projectsData.head(3)`

Out[4]:

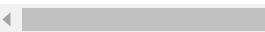
	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bcfd1151f324dd63a	Mr.	FL
2	21895	p182444	3465aa82da834c0582ebd0ef8040ca0	Ms.	AZ

In [5]: `projectsData.tail(3)`

Out[5]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_
109245	143653	p155633	cdbfd04aa041dc6739e9e576b1fb1478	Mrs.	NJ

	Unnamed: 0	id	teacher_id	teacher_prefix	school_
109246	164599	p206114	6d5675dbfafa1371f0e2f6f1b716fe2d	Mrs.	NY
109247	128381	p191189	ca25d5573f2bd2660f7850a886395927	Ms.	VA



In [6]: `resourcesData.head(3)`

Out[6]:

	id	description	quantity	price
0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.00
1	p069063	Bouncy Bands for Desks (Blue support pipes)	3	14.95
2	p069063	Cory Stories: A Kid's Book About Living With Adhd	1	8.45

In [7]: `resourcesData.tail(3)`

Out[7]:

	id	description	quantity	price
1541269	p031981	Black Electrical Tape (GIANT 3 PACK) Each Roll...	6	8.99
1541270	p031981	Flormoon DC Motor Mini Electric Motor 0.5-3V 1...	2	8.14
1541271	p031981	WAYLLSHINE 6PCS 2 x 1.5V AAA Battery Spring Cl...	2	7.39

## Helper functions and classes

```
In [0]: def equalsBorder(numberOfEqualSigns):
    """
    This function prints passed number of equal signs
    """
    print("=". * numberOfEqualSigns);
```

```
In [0]: # Citation link: https://stackoverflow.com/questions/8924173/how-do-i-print-bold-text-in-python
class color:
    PURPLE = '\x1b[95m'
    CYAN = '\x1b[96m'
    DARKCYAN = '\x1b[36m'
    BLUE = '\x1b[94m'
    GREEN = '\x1b[92m'
    YELLOW = '\x1b[93m'
    RED = '\x1b[91m'
    BOLD = '\x1b[1m'
    UNDERLINE = '\x1b[4m'
    END = '\x1b[0m'
```

```
In [0]: def printStyle(text, style):
    "This function prints text with the style passed to it"
    print(style + text + color.END);
```

## Shapes of projects data and resources data

```
In [11]: printStyle("Number of data points in projects data: {}".format(projectsData.shape[0]), color.BOLD);
printStyle("Number of attributes in projects data:{}".format(projectsData.shape[1]), color.BOLD);
equalsBorder(60);
printStyle("Number of data points in resources data: {}".format(resourcesData.shape[0]), color.BOLD);
```

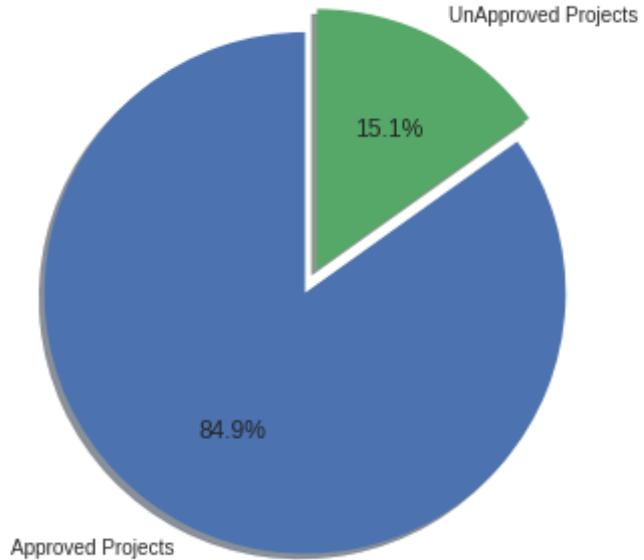
```
printStyle("Number of attributes in resources data: {}".format(resource  
sData.shape[1])), color.BOLD);  
  
Number of data points in projects data: 109248  
Number of attributes in projects data:17  
=====  
Number of data points in resources data: 1541272  
Number of attributes in resources data: 4
```

## Univariate data analysis

In [12]:

```
approvedProjects = projectsData[projectsData.project_is_approved == 1].  
shape[0];  
unApprovedProjects = projectsData[projectsData.project_is_approved == 0]  
.shape[0];  
totalProjects = projectsData.shape[0];  
print("Number of projects approved for funding: {}, ({})".format(approv  
edProjects, (approvedProjects / totalProjects) * 100));  
print("Number of projects not approved for funding: {}, ({})".format(un  
ApprovedProjects, (unApprovedProjects / totalProjects) * 100));  
# Pie chart representation  
# Citation: https://matplotlib.org/gallery/pie\_and\_polar\_charts/pie\_fea  
tures.html  
labels = ["Approved Projects", "UnApproved Projects"];  
explode = (0, 0.1);  
sizes = [approvedProjects, unApprovedProjects];  
figure, ax = plt.subplots();  
ax.pie(sizes, labels = labels, explode = explode, autopct = '%1.1f%%',  
shadow = True, startangle = 90);  
ax.axis('equal');  
plt.rcParams['figure.figsize'] = (7, 7);  
plt.show();
```

```
Number of projects approved for funding: 92706, (84.85830404217927)  
Number of projects not approved for funding: 16542, (15.1416959578073  
9)
```



### Observation:

1. There are more number of approved projects compared to rejected projects. So this is an imbalanced dataset.

## Univariate Analysis : 'school\_state'

### Project proposal percentage in different states

```
In [13]: groupedByStatesData = pd.DataFrame(projectsData.groupby(['school_state'])['project_is_approved'].apply(np.mean)).reset_index();
groupedByStatesData.columns = ['state_code', 'number_of_proposals'];
groupedByStatesData = groupedByStatesData.sort_values(by=['number_of_proposals'], ascending = True);
printStyle("5 States with lowest percentage of project approvals:", col
```

```
or.BOLD);
equalsBorder(60);
groupedByStatesData.head(5)
```

**5 States with lowest percentage of project approvals:**

---

Out[13]:

	state_code	number_of_proposals
46	VT	0.800000
7	DC	0.802326
43	TX	0.813142
26	MT	0.816327
18	LA	0.831245

```
In [14]: printStyle("5 states with highest percentage of project approvals: ", color.BOLD);
equalsBorder(60);
groupedByStatesData.tail(5).iloc[::-1]
```

**5 states with highest percentage of project approvals:**

---

Out[14]:

	state_code	number_of_proposals
8	DE	0.897959
28	ND	0.888112
47	WA	0.876178
35	OH	0.875152
30	NH	0.873563

```
In [0]: def univariateBarPlots(data, col1, col2 = 'project_is_approved', orient
```

```

        ation = 'vertical', plot = True):
    groupedData = data.groupby(col1);
    # Count number of zeros in dataframe python: https://stackoverflow.
    com/a/51540521/4084039
    tempData = pd.DataFrame(groupedData[col2].agg(lambda x: x.eq(1).sum
    ()).reset_index());
    tempData['total'] = pd.DataFrame(groupedData[col2].agg({'total': 'c
    ount'})).reset_index()['total'];
    tempData['approval_rate'] = pd.DataFrame(groupedData[col2].agg({'ap
    proval_rate': 'mean'})).reset_index()['approval_rate'];
    tempData.sort_values(by=['total'], inplace = True, ascending = Fals
    e);
    tempDataWithTotalAndCol2 = tempData[['total', col2, col1]]
    if plot:
        if(orientation == 'vertical'):
            tempDataWithTotalAndCol2.plot(x = col1, align= 'center', ki
            nd = 'bar', title = "Number of projects approved vs rejected", figsize
            = (20, 6), stacked = True, rot = 0);
        else:
            tempDataWithTotalAndCol2.plot(x = col1, align= 'center', ki
            nd = 'barh', title = "Number of projects approved vs rejected", width =
            0.8, figsize = (23, 20), stacked = True);
    return tempData;

```

In [16]:

```

statesCharacteristicsData = univariateBarPlots(projectsData, 'school_st
ate', 'project_is_approved', orientation = 'vertical');
printStyle("Top 5 states with high project proposals", color.BOLD)
equalsBorder(60);
statesCharacteristicsData.head(5)

```

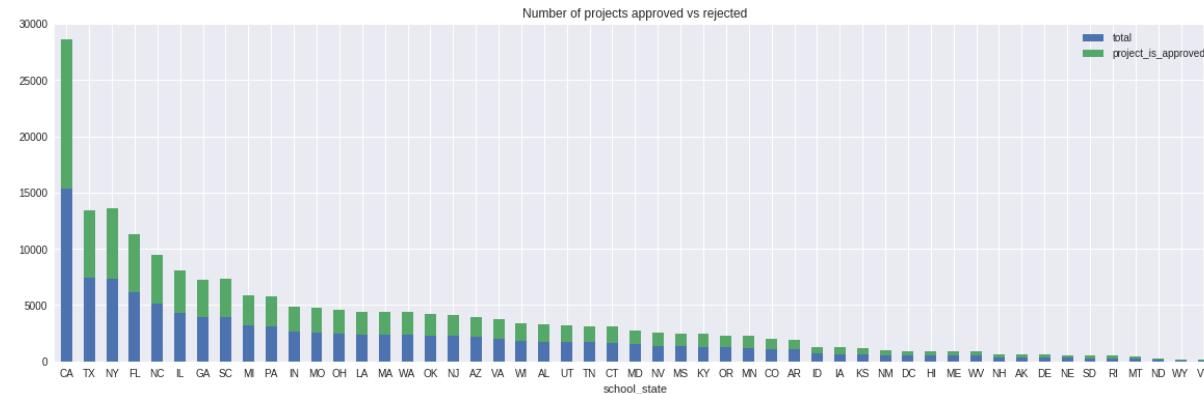
**Top 5 states with high project proposals**

---

Out[16]:

	school_state	project_is_approved	total	approval_rate
4	CA	13205	15388	0.858136
43	TX	6014	7396	0.813142

	<b>school_state</b>	<b>project_is_approved</b>	<b>total</b>	<b>approval_rate</b>
<b>34</b>	NY	6291	7318	0.859661
<b>9</b>	FL	5144	6185	0.831690
<b>27</b>	NC	4353	5091	0.855038



```
In [17]: printStyle("Top 5 states with least project proposals", color.BOLD)
equalsBorder(60);
statesCharacteristicsData.tail(5)
```

### Top 5 states with least project proposals

---

Out[17]:

	<b>school_state</b>	<b>project_is_approved</b>	<b>total</b>	<b>approval_rate</b>
<b>39</b>	RI	243	285	0.852632
<b>26</b>	MT	200	245	0.816327
<b>28</b>	ND	127	143	0.888112
<b>50</b>	WY	82	98	0.836735
<b>46</b>	VT	64	80	0.800000

## Observation:

1. Highest number of project proposals are from CA(California) and it was almost about 16000 projects
2. Every state has more than 80% approval rate.

## Univariate Analysis: teacher\_prefix

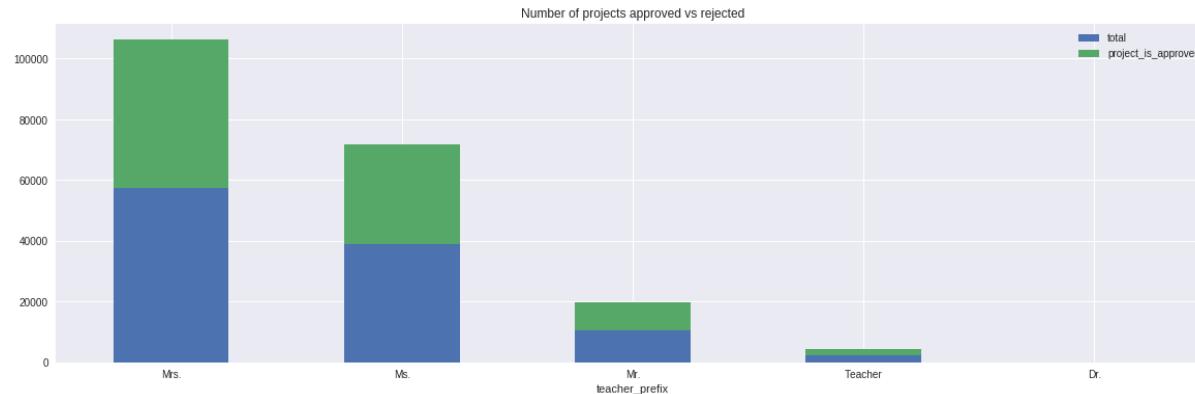
```
In [18]: teacherPrefixCharacteristicsData = univariateBarPlots(projectsData, 'teacher_prefix', 'project_is_approved', orientation = 'vertical', plot = True);
printStyle("Project proposals characteristics based on types of persons", color.BOLD);
equalsBorder(60);
teacherPrefixCharacteristicsData
```

Project proposals characteristics based on types of persons

---

Out[18]:

	teacher_prefix	project_is_approved	total	approval_rate
2	Mrs.	48997	57269	0.855559
3	Ms.	32860	38955	0.843537
1	Mr.	8960	10648	0.841473
4	Teacher	1877	2360	0.795339
0	Dr.	9	13	0.692308



### Observation:

1. When compared to others Dr.'s have proposed very less number of projects.
2. Women have proposed more number of projects than men.

### Univariate Analysis: project\_grade

```
In [19]: gradeCharacteristicsData = univariateBarPlots(projectsData, 'project_grade_category', 'project_is_approved', orientation = 'vertical', plot = True);
printStyle("Project proposal characteristics based on grades", color.BOLD);
equalsBorder(60);
gradeCharacteristicsData
```

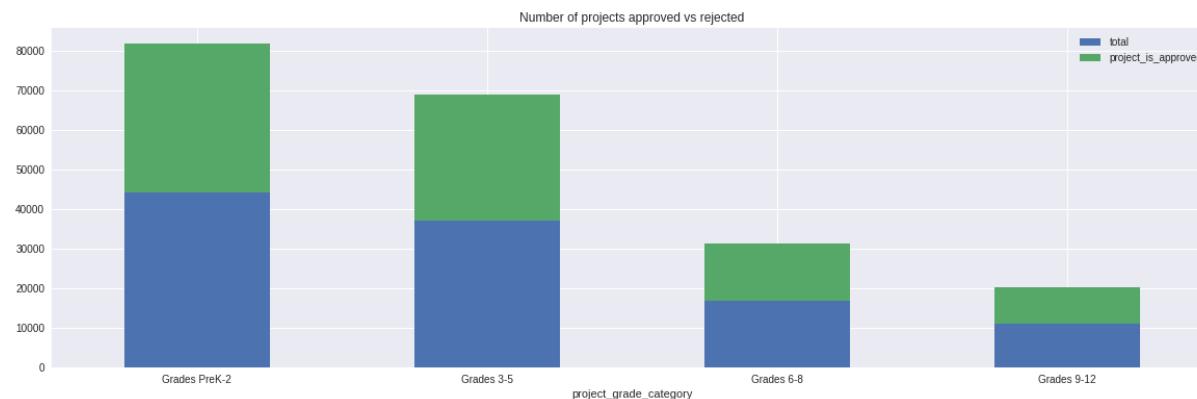
#### Project proposal characteristics based on grades

---

Out[19]:

	project_grade_category	project_is_approved	total	approval_rate
3	Grades PreK-2	37536	44225	0.848751
0	Grades 3-5	31729	37137	0.854377

	project_grade_category	project_is_approved	total	approval_rate
1	Grades 6-8	14258	16923	0.842522
2	Grades 9-12	9183	10963	0.837636



### Observation:

1. Most number of projects proposed are for students less than grade-5 (for primary school students) which means that children are being taught with project oriented teaching which is great.

### Univariate Analysis: project\_subject\_categories

```
In [0]: # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all whitespace-in-a-string-in-python
def cleanCategories(subjectCategories):
```

```

cleanedCategories = []
for subjectCategory in tqdm(subjectCategories):
    tempCategory = ""
    for category in subjectCategory.split(","):
        if 'The' in category.split(): # this will split each of the
            category based on space "Math & Science"=> "Math", "&", "Science"
            category = category.replace('The','') # if we have the
            words "The" we are going to replace it with ''(i.e removing 'The')
            category = category.replace(' ', '') # we are placing all t
            he ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
            tempCategory += category.strip()+" "# abc ".strip() will r
eturn "abc", remove the trailing spaces
        tempCategory = tempCategory.replace('&', '_')
    cleanedCategories.append(tempCategory)
return cleanedCategories

```

```

In [21]: # projectDataWithCleanedCategories = pd.DataFrame(projectsData);
subjectCategories = list(projectsData.project_subject_categories);
cleanedCategories = cleanCategories(subjectCategories);
printStyle("Sample categories: ", color.BOLD);
equalsBorder(60);
print(subjectCategories[0:5]);
equalsBorder(60);
printStyle("Sample cleaned categories: ", color.BOLD);
equalsBorder(60);
print(cleanedCategories[0:5]);
projectsData['cleaned_categories'] = cleanedCategories;
projectsData.head(5)

```

**Sample categories:**

```
=====
['Literacy & Language', 'History & Civics, Health & Sports', 'Health &
Sports', 'Literacy & Language, Math & Science', 'Math & Science']
=====
```

**Sample cleaned categories:**

```
=====
['Literacy_Language ', 'History_Civics_Health_Sports ', 'Health_Sports
', 'Literacy_Language_Math_Science ', 'Math_Science ']
```

In [21]:

In [21]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bcfd1151f324dd63a	Mr.	FL
2	21895	p182444	3465aa82da834c0582ebd0ef8040ca0	Ms.	AZ
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX

In [22]:

```
categoriesCharacteristicsData = univariateBarPlots(projectsData, 'cleaned_categories', 'project_is_approved', orientation = 'horizontal', plot = True);
print("Project proposals characteristics based on subject categories");
```

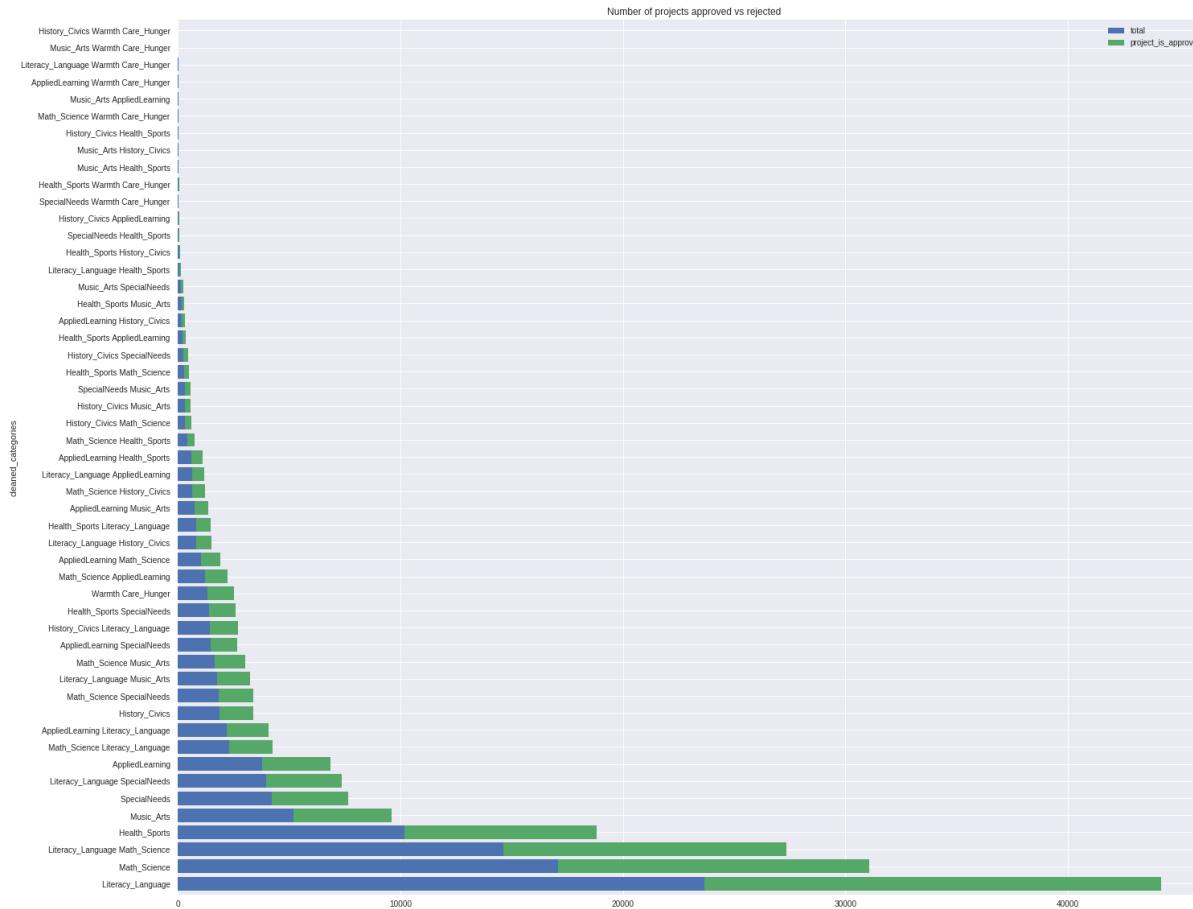
```
equalsBorder(60);  
categoriesCharacteristicsData.head(5)
```

Project proposals characteristics based on subject categories

---

Out[22]:

	cleaned_categories	project_is_approved	total	approval_rate
24	Literacy_Language	20520	23655	0.867470
32	Math_Science	13991	17072	0.819529
28	Literacy_Language Math_Science	12725	14636	0.869432
8	Health_Sports	8640	10177	0.848973
40	Music_Arts	4429	5180	0.855019



```
In [23]: # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
categoriesCounter = Counter()
for subjectCategory in projectsData.cleaned_categories.values:
    categoriesCounter.update(subjectCategory.split());
categoriesCounter
```

Out[23]: Counter({'AppliedLearning': 12135,  
 'Care\_Hunger': 1388,  
 'Health\_Sports': 14223,  
 'History\_Civics': 5914,

```
'Literacy_Language': 52239,  
'Math_Science': 41421,  
'Music_Arts': 10293,  
'SpecialNeeds': 13642,  
'Warmth': 1388})
```

In [24]: *# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039*

```
categoriesDictionary = dict(categoriesCounter);  
sortedCategoriesDictionary = dict(sorted(categoriesDictionary.items(),  
key = lambda keyValue: keyValue[1]));  
sortedCategoriesData = pd.DataFrame.from_dict(sortedCategoriesDictionary,  
orient='index');  
sortedCategoriesData.columns = ['subject_categories'];  
printStyle("Number of projects by Subject Categories: ", color.BOLD);  
equalsBorder(60);  
sortedCategoriesData
```

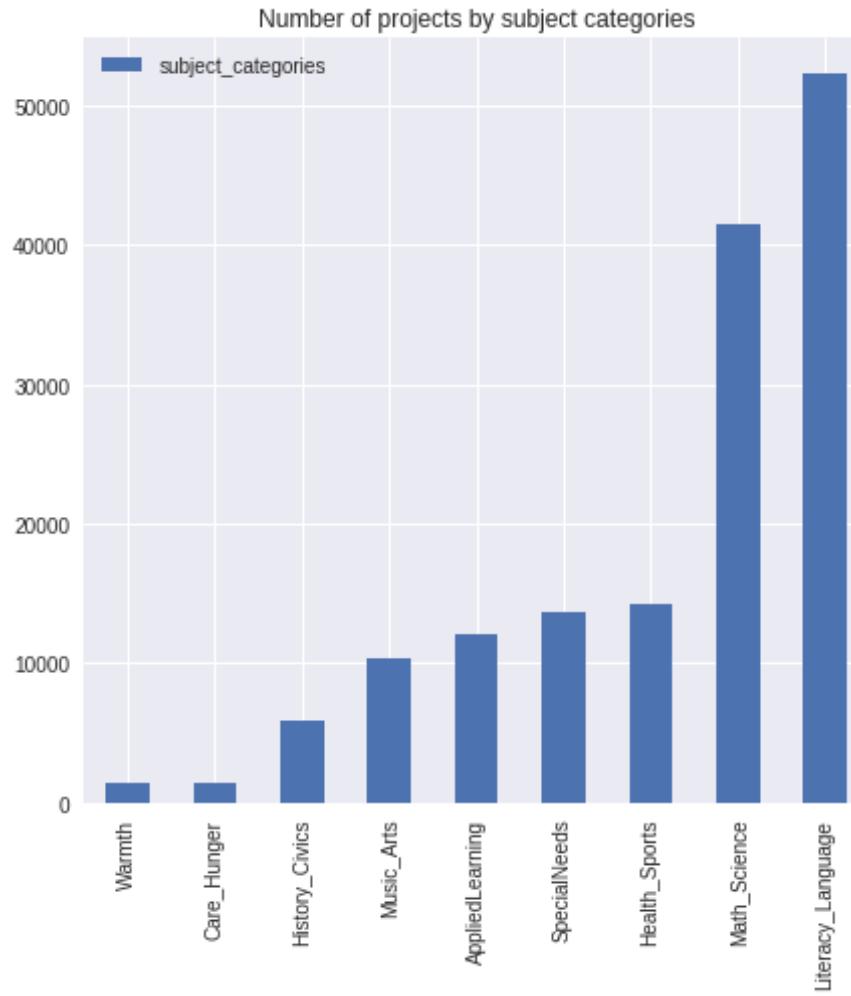
#### Number of projects by Subject Categories:

---

Out[24]:

	subject_categories
Warmth	1388
Care_Hunger	1388
History_Civics	5914
Music_Arts	10293
AppliedLearning	12135
SpecialNeeds	13642
Health_Sports	14223
Math_Science	41421
Literacy_Language	52239

```
In [25]: sortedCategoriesData.plot(kind = 'bar', title = 'Number of projects by subject categories');
```



### Observation:

1. Many number of projects proposed belong to multiple subject categories.
2. When compared to others literacy\_language & math\_science have large number of project proposals.

## Univariate Analysis: project\_subject\_subcategories

```
In [26]: subjectSubCategories = projectsData.project_subject_subcategories;
cleanedSubCategories = cleanCategories(subjectSubCategories);
printStyle("Sample subject sub categories: ", color.BOLD);
equalsBorder(70);
print(subjectSubCategories[0:5]);
equalsBorder(70);
printStyle("Sample cleaned subject sub categories: ", color.BOLD);
equalsBorder(70);
print(cleanedSubCategories[0:5]);
projectsData['cleaned_sub_categories'] = cleanedSubCategories;
```

**Sample subject sub categories:**

```
=====
0          ESL, Literacy
1  Civics & Government, Team Sports
2    Health & Wellness, Team Sports
3          Literacy, Mathematics
4          Mathematics
Name: project_subject_subcategories, dtype: object
=====
```

**Sample cleaned subject sub categories:**

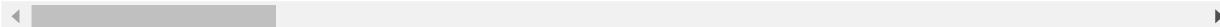
```
=====
['ESL Literacy ', 'Civics_Government TeamSports ', 'Health_Wellness Tea
mSports ', 'Literacy Mathematics ', 'Mathematics ']
```

```
In [27]: projectsData.head(5)
```

Out[27]:

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

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bc1151f324dd63a	Mr.	FL
2	21895	p182444	3465aa82da834c0582ebd0ef8040ca0	Ms.	AZ
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX



```
In [28]: subCategoriesCharacteristicsData = univariateBarPlots(projectsData, 'cleaned_sub_categories', 'project_is_approved', plot = False);
print("Project proposals characteristics based on subject sub categories");
equalsBorder(60);
subCategoriesCharacteristicsData.head(5)
```

## Project proposals characteristics based on subject sub categories

---

Out[28]:

	cleaned_sub_categories	project_is_approved	total	approval_rate
317	Literacy	8371	9486	0.882458
319	Literacy Mathematics	7260	8325	0.872072
331	Literature_Writing Mathematics	5140	5923	0.867803
318	Literacy Literature_Writing	4823	5571	0.865733
342	Mathematics	4385	5379	0.815207

In [29]: *# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039*

```
subjectsSubCategoriesCounter = Counter();
for subCategory in projectsData.cleaned_sub_categories:
    subjectsSubCategoriesCounter.update(subCategory.split());
subjectsSubCategoriesCounter
```

Out[29]: Counter({'AppliedSciences': 10816,
 'Care\_Hunger': 1388,
 'CharacterEducation': 2065,
 'Civics\_Government': 815,
 'College\_CareerPrep': 2568,
 'CommunityService': 441,
 'ESL': 4367,
 'EarlyDevelopment': 4254,
 'Economics': 269,
 'EnvironmentalScience': 5591,
 'Extracurricular': 810,
 'FinancialLiteracy': 568,
 'ForeignLanguages': 890,
 'Gym\_Fitness': 4509,
 'Health\_LifeScience': 4235,
 'Health\_Wellness': 10234,
 'History\_Geography': 3171,

```
'Literacy': 33700,
'Literature_Writing': 22179,
'Mathematics': 28074,
'Music': 3145,
'NutritionEducation': 1355,
'Other': 2372,
'ParentInvolvement': 677,
'PerformingArts': 1961,
'SocialSciences': 1920,
'SpecialNeeds': 13642,
'TeamSports': 2192,
'VisualArts': 6278,
'Warmth': 1388})
```

```
In [30]: # dict sort by value python: https://stackoverflow.com/a/613218/4084039
dictionarySubCategories = dict(subjectsSubCategoriesCounter);
sortedDictionarySubCategories = dict(sorted(dictionarySubCategories.items(), key = lambda keyValue: keyValue[1]));
sortedSubCategoriesData = pd.DataFrame.from_dict(sortedDictionarySubCategories, orient = 'index');
sortedSubCategoriesData.columns = ['subject_sub_categories']
sortedSubCategoriesData.plot(kind = 'bar', title = "Number of projects by subject sub categories");
printStyle("Number of projects sorted by subject sub categories: ", color.BOLD);
equalsBorder(70);
sortedSubCategoriesData
```

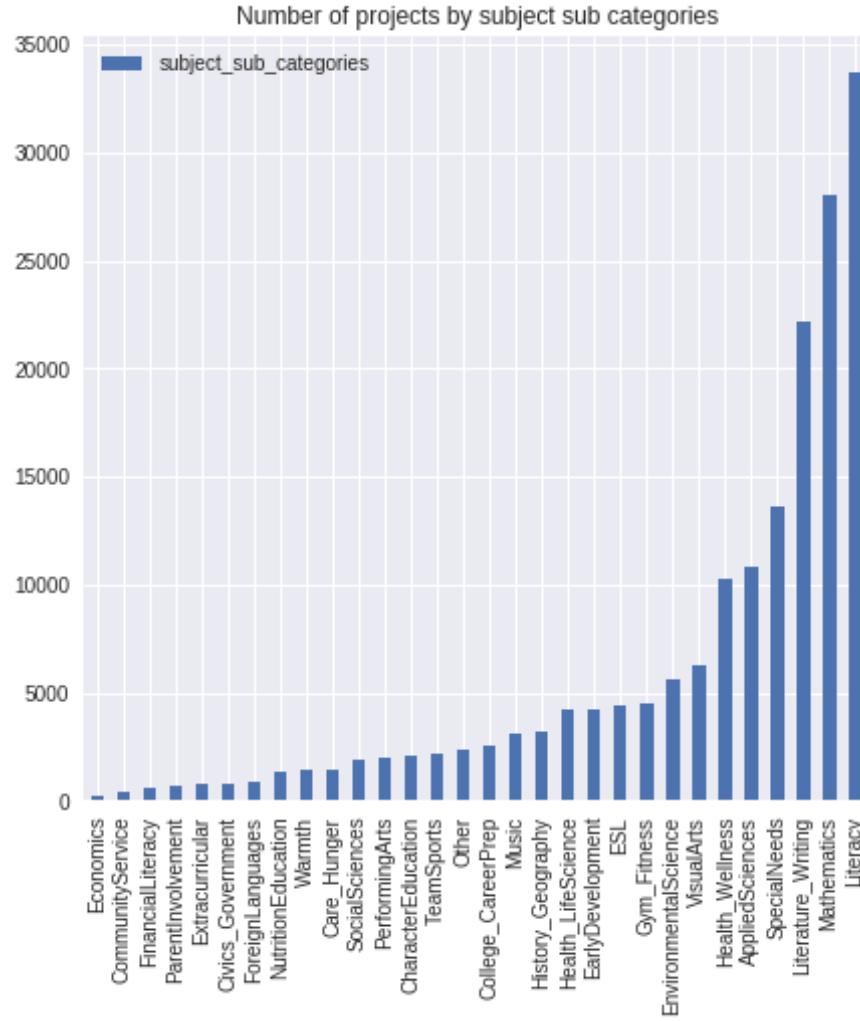
**Number of projects sorted by subject sub categories:**

Out[30]:

	subject_sub_categories
Economics	269
CommunityService	441
FinancialLiteracy	568

	subject_sub_categories
<b>ParentInvolvement</b>	677
<b>Extracurricular</b>	810
<b>Civics_Government</b>	815
<b>ForeignLanguages</b>	890
<b>NutritionEducation</b>	1355
<b>Warmth</b>	1388
<b>Care_Hunger</b>	1388
<b>SocialSciences</b>	1920
<b>PerformingArts</b>	1961
<b>CharacterEducation</b>	2065
<b>TeamSports</b>	2192
<b>Other</b>	2372
<b>College_CareerPrep</b>	2568
<b>Music</b>	3145
<b>History_Geography</b>	3171
<b>Health_LifeScience</b>	4235
<b>EarlyDevelopment</b>	4254
<b>ESL</b>	4367
<b>Gym_Fitness</b>	4509
<b>EnvironmentalScience</b>	5591
<b>VisualArts</b>	6278
<b>Health_Wellness</b>	10234

	<b>subject_sub_categories</b>
<b>AppliedSciences</b>	10816
<b>SpecialNeeds</b>	13642
<b>Literature_Writing</b>	22179
<b>Mathematics</b>	28074
<b>Literacy</b>	33700



## Observation:

1. There are more number of subject subcategories than subject categories.
2. Even more number of projects proposed belong to multiple subject sub categories.

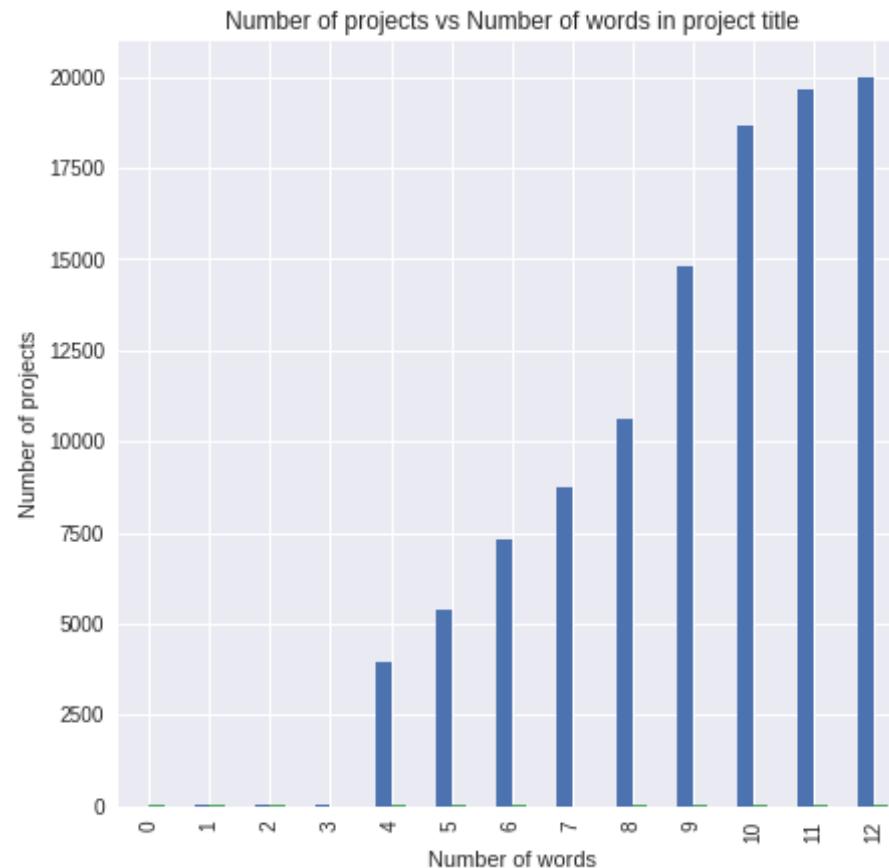
## Univariate Analysis : project\_title

```
In [31]: #How to calculate number of words in a string in DataFrame: https://stackoverflow.com/a/37483537/4084039
wordCounts = projectsData['project_title'].str.split().apply(len).value_
_counts();
dictionaryWordCounts = dict(wordCounts);
dictionaryWordCounts = dict(sorted(dictionaryWordCounts.items(), key =
lambda kv: kv[1]));
wordCountsData = pd.DataFrame.from_dict({'number_of_words': list(dictionaryWordCounts.keys()), 'number_of_projects': list(dictionaryWordCounts.values())}).sort_values(by = ['number_of_projects']);
wordCountsData.plot(kind = 'bar', title = "Number of projects vs Number
of words in project title", legend = False);
plt.xlabel('Number of words');
plt.ylabel('Number of projects');
wordCountsData
```

Out[31]:

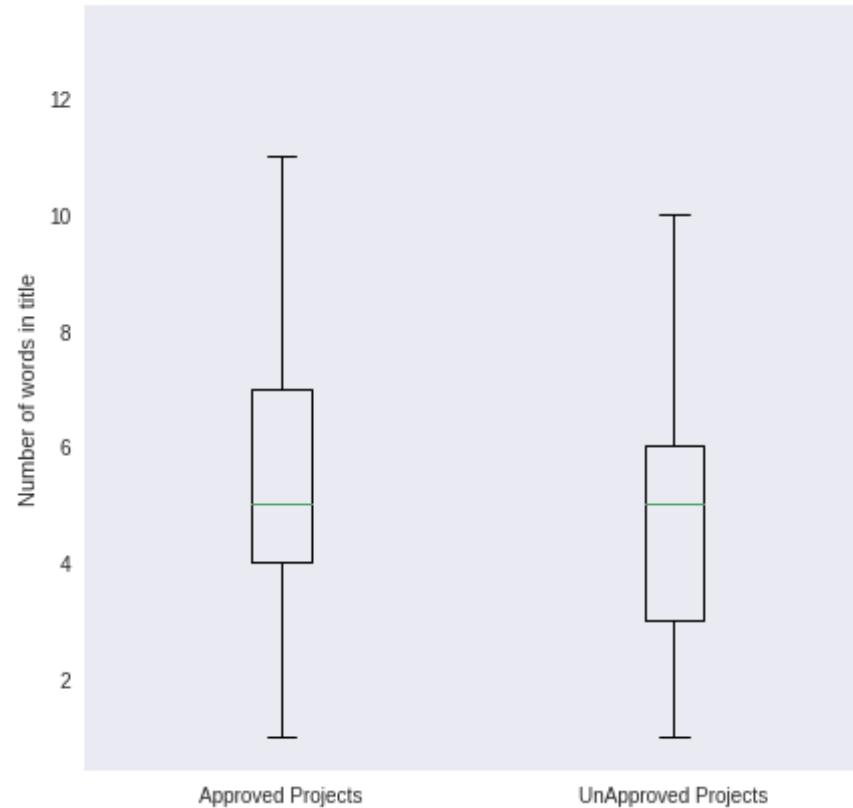
	number_of_projects	number_of_words
0	1	13
1	11	12
2	30	11
3	31	1
4	3968	10
5	5383	9
6	7289	8
7	8733	2
8	10631	7
9	14824	6
10	18691	3

	number_of_projects	number_of_words
11	19677	5
12	19979	4

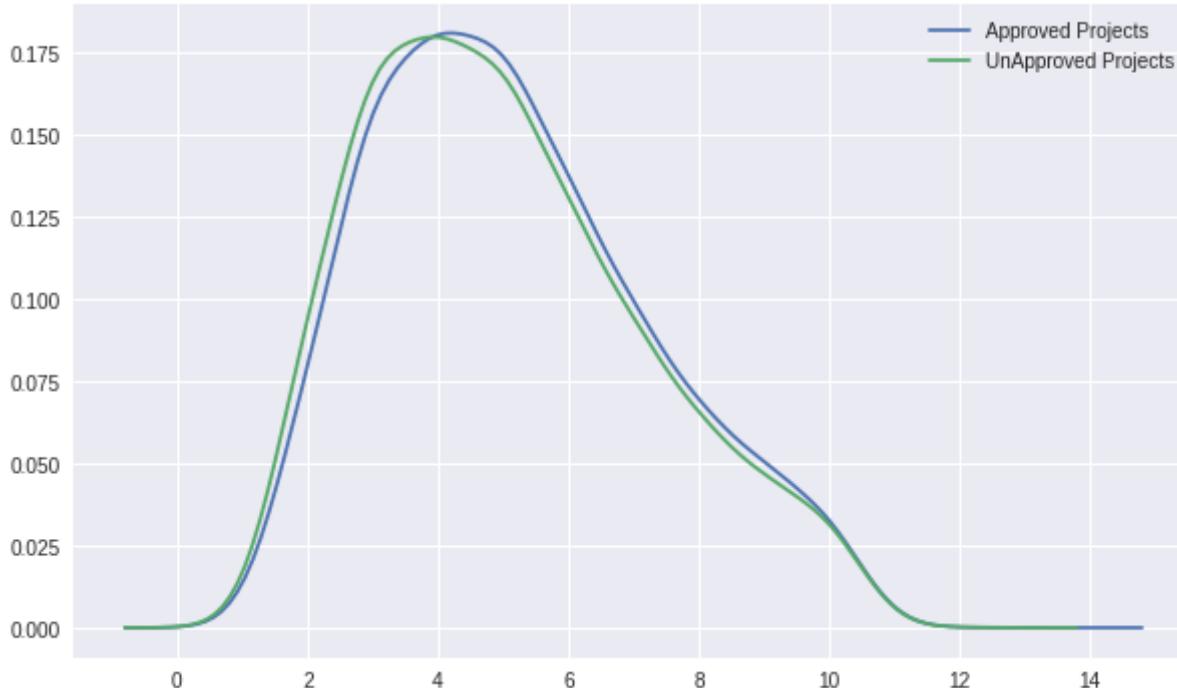


```
In [32]: approvedNumber0fProjects = projectsData[projectsData.project_is_approved == 1]['project_title'].str.split().apply(len);
approvedNumber0fProjects = approvedNumber0fProjects.values
unApprovedNumber0fProjects = projectsData[projectsData.project_is_approved == 0]['project_title'].str.split().apply(len);
unApprovedNumber0fProjects = unApprovedNumber0fProjects.values
```

```
plt.boxplot([approvedNumber0fProjects, unApprovedNumber0fProjects]);
plt.grid();
plt.xticks([1, 2], ['Approved Projects', 'UnApproved Projects']);
plt.ylabel('Number of words in title');
plt.show();
```



```
In [33]: plt.figure(figsize = (10, 6));
sbrn.kdeplot(approvedNumber0fProjects, label = "Approved Projects", bw
= 0.6);
sbrn.kdeplot(unApprovedNumber0fProjects, label = "UnApproved Projects",
bw = 0.6);
plt.legend();
plt.show();
```



### Observations:

1. Most of the approved projects have between 4 to 8 number of words in their project\_title.
2. Most of the rejected projects have between 3 to 6 number of words in their project\_title.

### Univariate Analysis: project\_essay\_1,2,3,4

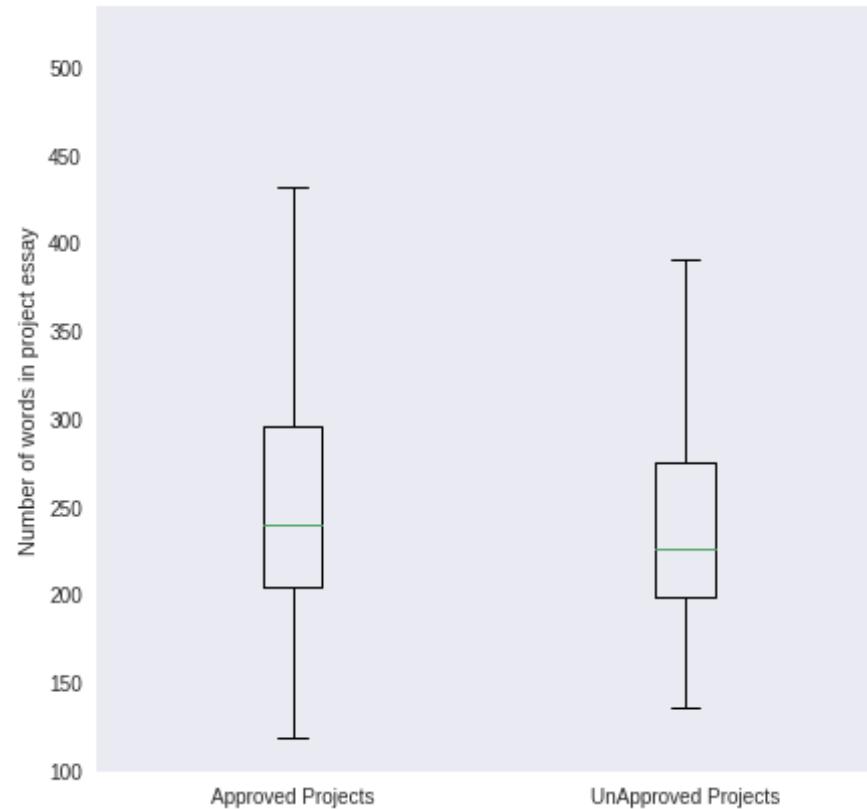
```
In [34]: projectsData['project_essay'] = projectsData['project_essay_1'].map(str) + projectsData['project_essay_2'].map(str) + \\ projectsData['project_essay_3'].map(str) + projectsData['project_essay_4'].map(str); projectsData.head(5)
```

Out[34]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bcfd1151f324dd63a	Mr.	FL
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ
3	45	p246581	f3cb9bffba169bef1a77b243e620b60	Mrs.	KY
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX

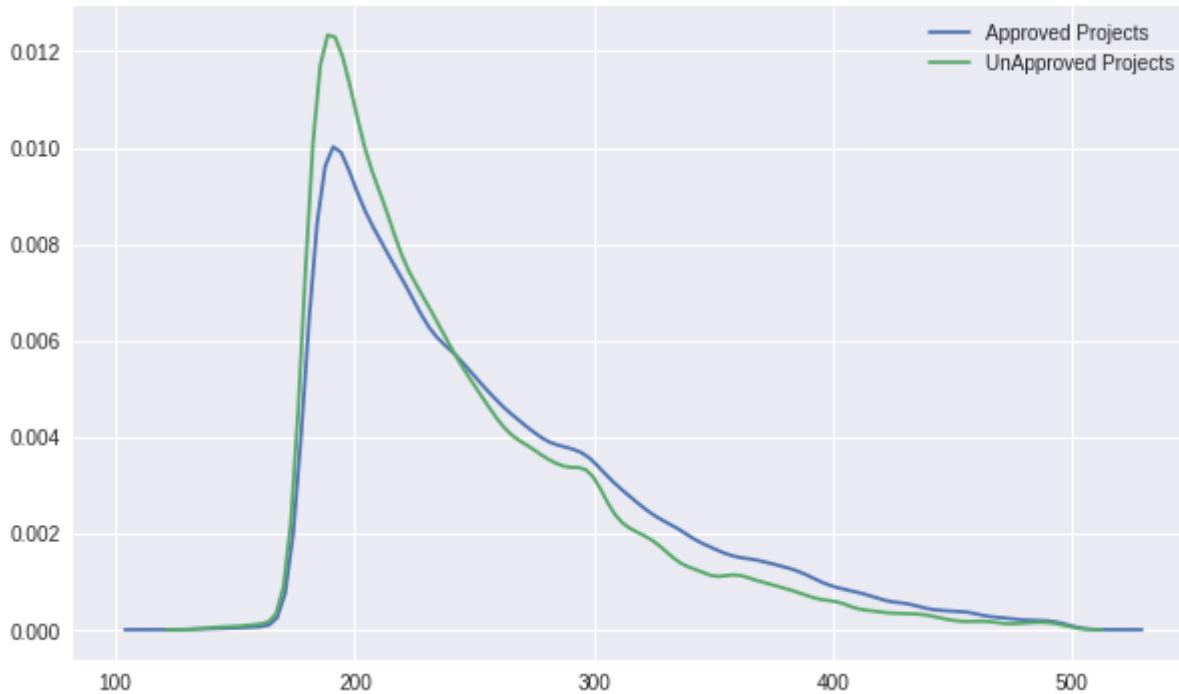
```
In [35]: approvedNumber0fProjects = projectsData[projectsData.project_is_approve  
d == 1]['project_essay'].str.split().apply(len);  
approvedNumber0fProjects = approvedNumber0fProjects.values  
unApprovedNumber0fProjects = projectsData[projectsData.project_is_appro
```

```
ved == 0]['project_essay'].str.split().apply(len);
unApprovedNumber0fProjects = unApprovedNumber0fProjects.values
plt.boxplot([approvedNumber0fProjects, unApprovedNumber0fProjects]);
plt.grid();
plt.xticks([1, 2], ['Approved Projects', 'UnApproved Projects']);
plt.ylabel('Number of words in project essay');
plt.show();
```



```
In [36]: plt.figure(figsize = (10, 6));
sbrn.kdeplot(approvedNumber0fProjects, label = "Approved Projects", bw = 5);
sbrn.kdeplot(unApprovedNumber0fProjects, label = "UnApproved Projects", bw = 5);
```

```
plt.legend();
plt.show();
```



### Observation:

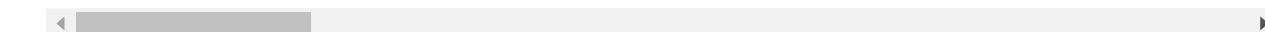
1. The approved and rejected projects overlap largely when plotted based on number of words in project\_essay. So we cannot predict any observation which will be useful for classification.

### Univariate Analysis: price

In [37]: `projectsData.head(5)`

Out[37]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bcfd1151f324dd63a	Mr.	FL
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX



In [38]: `resourcesData.head(5)`

Out[38]:

	id	description	quantity	price

	<b>id</b>	<b>description</b>	<b>quantity</b>	<b>price</b>
<b>0</b>	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.00
<b>1</b>	p069063	Bouncy Bands for Desks (Blue support pipes)	3	14.95
<b>2</b>	p069063	Cory Stories: A Kid's Book About Living With Adhd	1	8.45
<b>3</b>	p069063	Dixon Ticonderoga Wood-Cased #2 HB Pencils, Bo...	2	13.59
<b>4</b>	p069063	EDUCATIONAL INSIGHTS FLUORESCENT LIGHT FILTERS...	3	24.95

```
In [39]: # https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframe
          s-indexes-for-all-groups-in-one-step
priceAndQuantityData = resourcesData.groupby('id').agg({'price': 'sum',
    'quantity': 'sum'}).reset_index();
priceAndQuantityData.head(5)
```

Out[39]:

	<b>id</b>	<b>price</b>	<b>quantity</b>
<b>0</b>	p000001	459.56	7
<b>1</b>	p000002	515.89	21
<b>2</b>	p000003	298.97	4
<b>3</b>	p000004	1113.69	98
<b>4</b>	p000005	485.99	8

```
In [40]: projectsData.shape
```

Out[40]: (109248, 20)

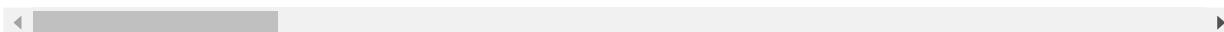
```
In [41]: projectsData = pd.merge(projectsData, priceAndQuantityData, on = 'id',
                                how = 'left');
print(projectsData.shape);
projectsData.head(3)
```

(109248, 22)

Out[41]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bcfd1151f324dd63a	Mr.	FL
2	21895	p182444	3465aaaf82da834c0582ebd0ef8040ca0	Ms.	AZ

3 rows × 22 columns



In [42]: `projectsData[projectsData['id'] == 'p253737']`

Out[42]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	per
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2

1 rows × 22 columns

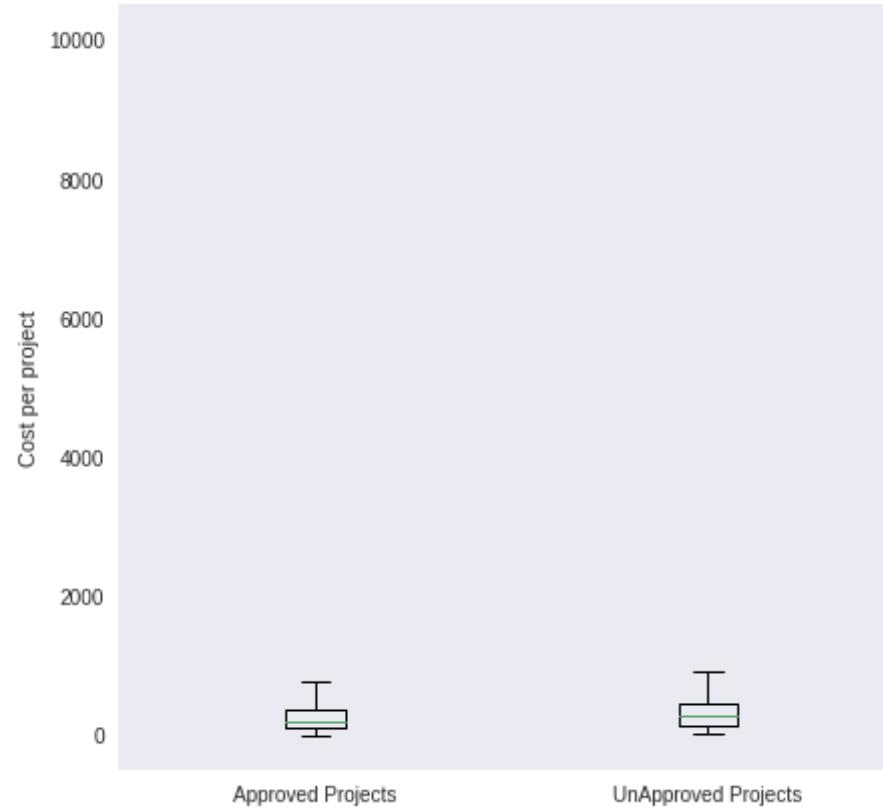


In [43]: `priceAndQuantityData[priceAndQuantityData['id'] == 'p253737']`

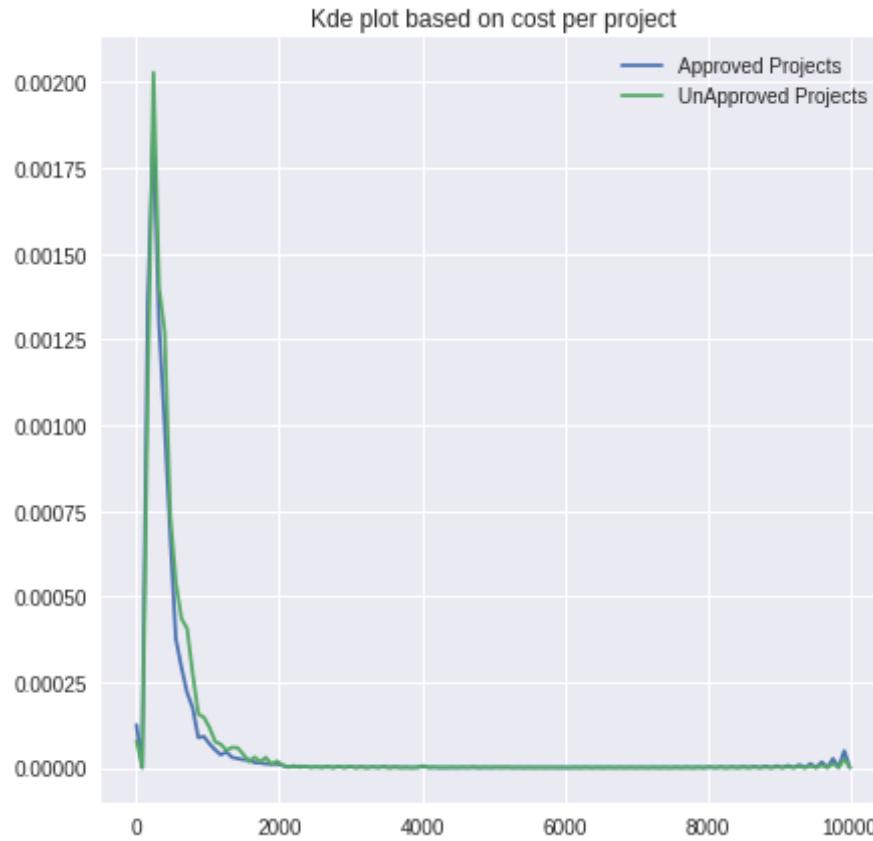
Out[43]:

	<b>id</b>	<b>price</b>	<b>quantity</b>
253736	p253737	154.6	23

In [44]: `approvedProjectsPrice = projectsData[projectsData['project_is_approved'] == 1].price;  
unApprovedProjectsPrice = projectsData[projectsData['project_is_approved'] == 0].price;  
plt.boxplot([approvedProjectsPrice, unApprovedProjectsPrice]);  
plt.grid();  
plt.xticks([1, 2], ['Approved Projects', 'UnApproved Projects']);  
plt.ylabel('Cost per project');  
plt.show();`



```
In [45]: plt.title("Kde plot based on cost per project");
sbrn.kdeplot(approvedProjectsPrice, label = "Approved Projects", bw =
0.6);
sbrn.kdeplot(unApprovedProjectsPrice, label = "UnApproved Projects", bw
= 0.6);
plt.legend();
plt.show();
```



```
In [46]: pricePercentilesApproved = [round(np.percentile(approvedProjectsPrice, percentile), 3) for percentile in np.arange(0, 100, 5)];  
pricePercentilesUnApproved = [round(np.percentile(unApprovedProjectsPrice, percentile), 3) for percentile in np.arange(0, 100, 5)];  
percentileValuePricesData = pd.DataFrame({'Percentile': np.arange(0, 100, 5), 'Approved projects': pricePercentilesApproved, 'UnApproved Projects': pricePercentilesUnApproved});  
percentileValuePricesData
```

Out[46]:

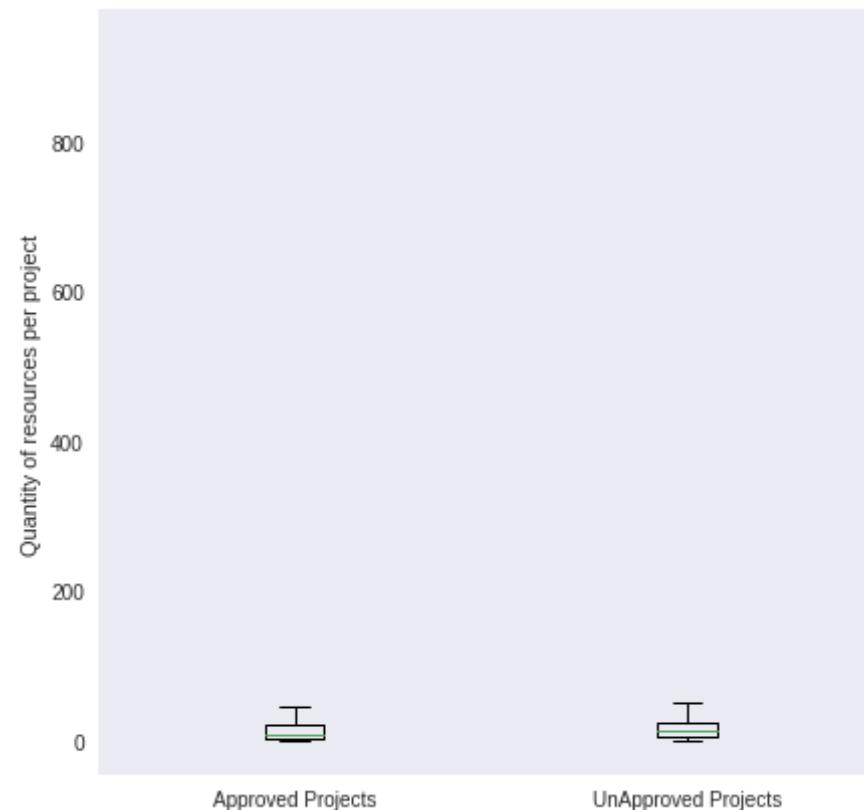
	Approved projects	Percentile	UnApproved Projects
0	0.660	0	1.970

	<b>Approved projects</b>	<b>Percentile</b>	<b>UnApproved Projects</b>
<b>1</b>	13.590	5	41.900
<b>2</b>	33.880	10	73.670
<b>3</b>	58.000	15	99.109
<b>4</b>	77.380	20	118.560
<b>5</b>	99.950	25	140.892
<b>6</b>	116.680	30	162.230
<b>7</b>	137.232	35	184.014
<b>8</b>	157.000	40	208.632
<b>9</b>	178.265	45	235.106
<b>10</b>	198.990	50	263.145
<b>11</b>	223.990	55	292.610
<b>12</b>	255.630	60	325.144
<b>13</b>	285.412	65	362.390
<b>14</b>	321.225	70	399.990
<b>15</b>	366.075	75	449.945
<b>16</b>	411.670	80	519.282
<b>17</b>	479.000	85	618.276
<b>18</b>	593.110	90	739.356
<b>19</b>	801.598	95	992.486

### **Observation:**

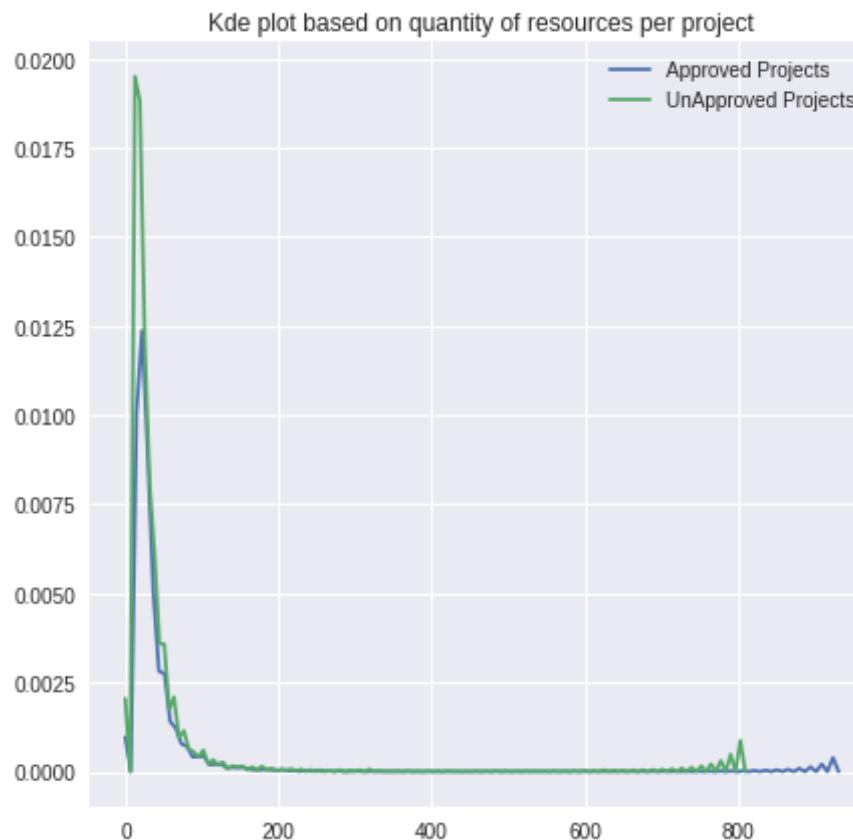
1. Most of the projects proposed are of less cost.

```
In [47]: approvedProjectsQuantity = projectsData[projectsData['project_is_approved'] == 1].quantity;
unApprovedProjectsQuantity = projectsData[projectsData['project_is_approved'] == 0].quantity;
plt.boxplot([approvedProjectsQuantity, unApprovedProjectsQuantity]);
plt.grid();
plt.xticks([1, 2], ['Approved Projects', 'UnApproved Projects']);
plt.ylabel('Quantity of resources per project');
plt.show();
```



```
In [48]: plt.title("Kde plot based on quantity of resources per project");
sbrn.kdeplot(approvedProjectsQuantity, label = "Approved Projects", bw = 0.6);
```

```
sbrn.kdeplot(unApprovedProjectsQuantity, label = "UnApproved Projects",  
    bw = 0.6);  
plt.legend();  
plt.show();
```



```
In [49]: quantityPercentilesApproved = [round(np.percentile(approvedProjectsQuan  
tity, percentile), 3) for percentile in np.arange(0, 100, 5)];  
quantityPercentilesUnApproved = [round(np.percentile(unApprovedProjects  
Quantity, percentile), 3) for percentile in np.arange(0, 100, 5)];  
percentileValueQuantitiesData = pd.DataFrame({'Percentile': np.arange(0  
, 100, 5), 'Approved projects': quantityPercentilesApproved, 'UnApprove  
d Projects': quantityPercentilesUnApproved});  
percentileValueQuantitiesData
```

Out[49]:

	Approved projects	Percentile	UnApproved Projects
0	1.0	0	1.0
1	1.0	5	2.0
2	1.0	10	3.0
3	2.0	15	4.0
4	3.0	20	5.0
5	3.0	25	6.0
6	4.0	30	7.0
7	5.0	35	8.0
8	6.0	40	9.0
9	7.0	45	10.0
10	8.0	50	12.0
11	10.0	55	13.0
12	11.0	60	15.0
13	14.0	65	18.0
14	16.0	70	20.0
15	20.0	75	24.0
16	25.0	80	29.0
17	30.0	85	35.0
18	38.0	90	45.0
19	56.0	95	63.0

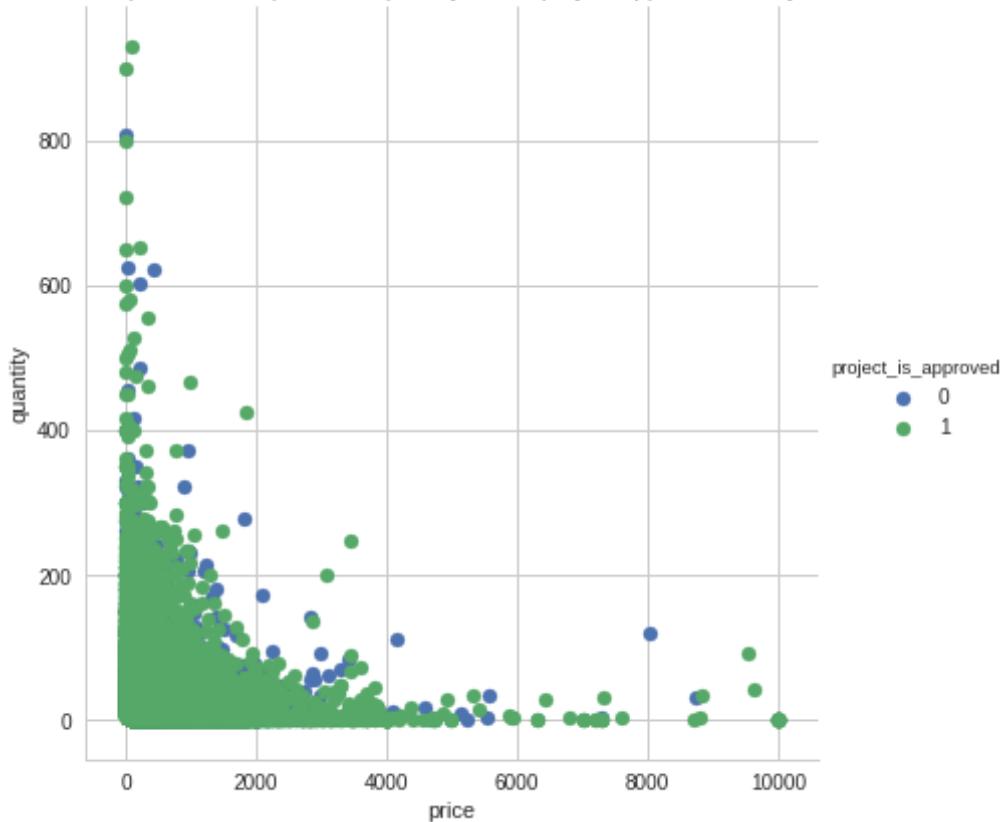
In [50]: `sbrn.set_style('whitegrid');`

```

sbrn.FacetGrid(projectsData, hue = 'project_is_approved', size = 6) \
    .map(plt.scatter, 'price', 'quantity') \
    .add_legend();
plt.title("Scatter plot between price and quantity based project approval and rejection");
plt.show();

```

Scatter plot between price and quantity based project approval and rejection



### Observation:

1. When plotted scatter plot between approved and rejected projects based on price and quantity there is huge overlap. So the projects approval is not actually depending on

price and quantity resources of the project.

### Univariate Analysis: teacher\_number\_of\_previously\_posted\_projects

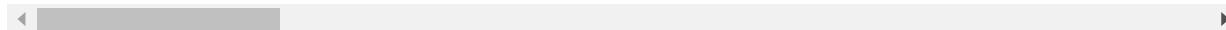
In [51]: `projectsData.head(5)`

Out[51]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bcfd1151f324dd63a	Mr.	FL
2	21895	p182444	3465aa82da834c0582ebd0ef8040ca0	Ms.	AZ
3	45	p246581	f3cb9bffba169bef1a77b243e620b60	Mrs.	KY

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX

5 rows × 22 columns



```
In [52]: previouslyPostedApprovedNumberData = projectsData.groupby('teacher_number_of_previously_posted_projects')['project_is_approved'].agg(lambda x: x.eq(1).sum()).reset_index();
previouslyPostedRejectedNumberData = projectsData.groupby('teacher_number_of_previously_posted_projects')['project_is_approved'].agg(lambda x: x.eq(0).sum()).reset_index();
print("Total number of projects approved: ", len(projectsData[projectsData['project_is_approved'] == 1]));
print("Total number of projects rejected: ", len(projectsData[projectsData['project_is_approved'] == 0]));
print("Number of projects approved categorized by previously_posted: ",
      previouslyPostedApprovedNumberData['project_is_approved'].sum());
print("Number of projects rejected categorized by previously_posted: ",
      previouslyPostedRejectedNumberData['project_is_approved'].sum());
previouslyPostedNumberData = pd.merge(previouslyPostedApprovedNumberData, previouslyPostedRejectedNumberData, on = 'teacher_number_of_previously_posted_projects', how = 'inner');
previouslyPostedNumberData.head(5)
```

Total number of projects approved: 92706

Total number of projects rejected: 16542

Number of projects approved categorized by previously\_posted: 92706

Number of projects rejected categorized by previously\_posted: 16542

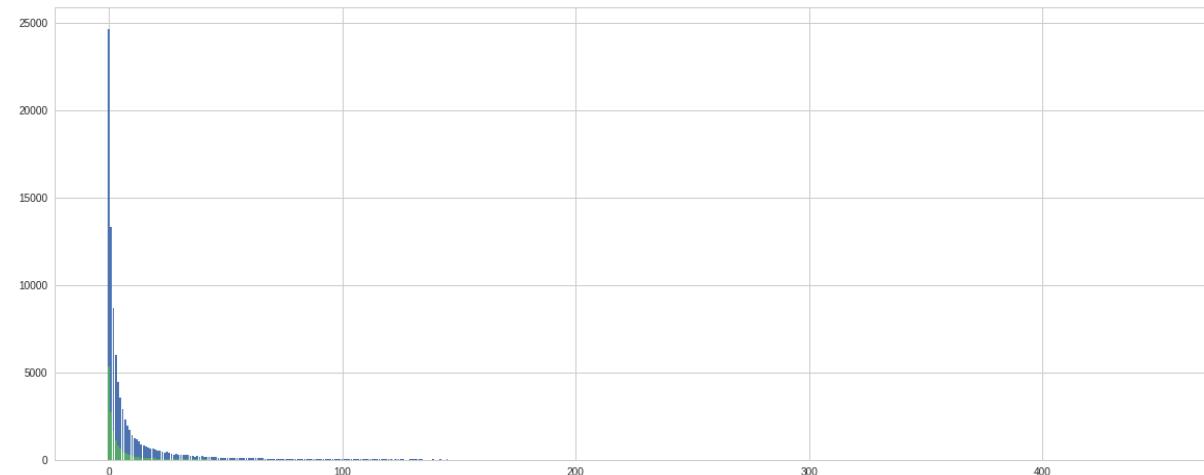
Out[52]:

teacher_number_of_previously_posted_projects	project_is_approved_x	project_is_ap

	teacher_number_of_previously_posted_projects	project_is_approved_x	project_is_ap
0	0	24652	5362
1	1	13329	2729
2	2	8705	1645
3	3	5997	1113
4	4	4452	814

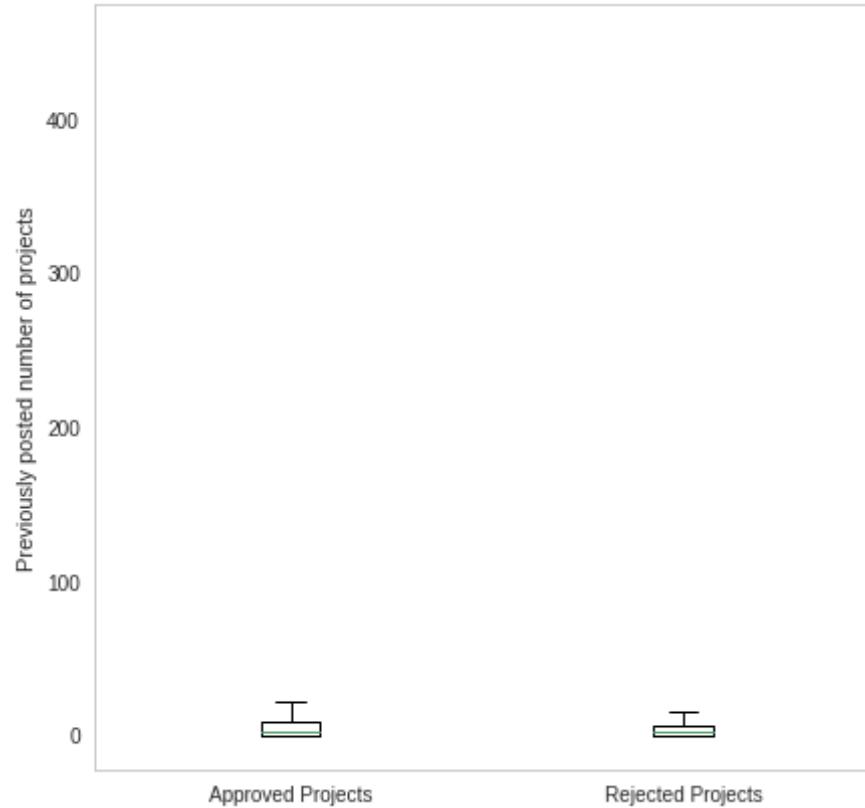


```
In [53]: plt.figure(figsize = (20, 8));
plt.bar(PreviouslyPostedNumberData.teacher_number_of_previously_posted_
projects, PreviouslyPostedNumberData.project_is_approved_x);
plt.bar(PreviouslyPostedNumberData.teacher_number_of_previously_posted_
projects, PreviouslyPostedNumberData.project_is_approved_y);
plt.show();
```

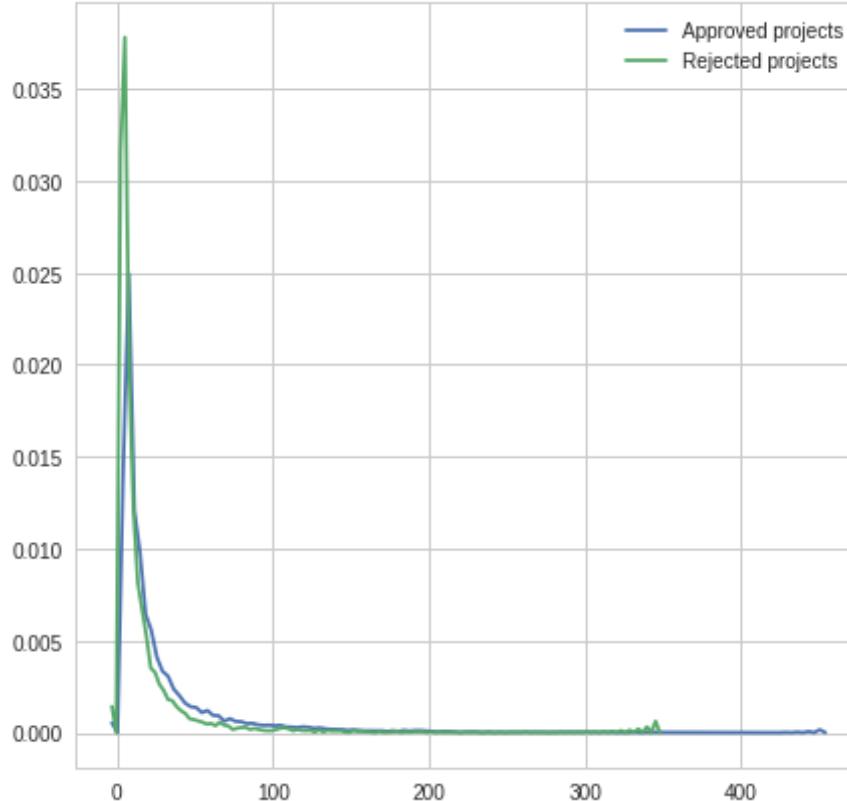


```
In [54]: previouslyPostedApprovedData = projectsData[projectsData['project_is_ap
proved'] == 1].teacher_number_of_previously_posted_projects;
previouslyPostedRejectedData = projectsData[projectsData['project_is_ap
proved'] == 0].teacher_number_of_previously_posted_projects;
```

```
plt.boxplot([previouslyPostedApprovedData, previouslyPostedRejectedData]);
plt.grid();
plt.xticks([1, 2], ['Approved Projects', 'Rejected Projects']);
plt.ylabel('Previously posted number of projects');
plt.show();
```



```
In [55]: sbrn.kdeplot(previouslyPostedApprovedData, label = "Approved projects",
                      bw = 1);
sbrn.kdeplot(previouslyPostedRejectedData, label = "Rejected projects",
                      bw = 1);
plt.show();
```



### Observation:

1. Most of the projects approved and rejected are with less number of teacher\_number\_of\_previously\_posted\_projects. So the approval is not much depending on how many number of projects proposed by teacher previously.

```
In [0]: def stringContainsNumbers(string):
    return any([character.isdigit() for character in string])
```

```
In [57]: numericResourceApprovedData = projectsData[(projectsData['project_resource_summary'].apply(stringContainsNumbers) == True) & (projectsData['project_is_approved'] == 1)]
```

```
textResourceApprovedData = projectsData[(projectsData['project_resource_summary'].apply(stringContainsNumbers) == False) & (projectsData['project_is_approved'] == 1)]
numericResourceRejectedData = projectsData[(projectsData['project_resource_summary'].apply(stringContainsNumbers) == True) & (projectsData['project_is_approved'] == 0)]
textResourceRejectedData = projectsData[(projectsData['project_resource_summary'].apply(stringContainsNumbers) == False) & (projectsData['project_is_approved'] == 0)]
print("Checking whether numbers in resource summary will be useful for project approval?");
equalsBorder(70);
print("Number of approved projects with numbers in resource summary: ", numericResourceApprovedData.shape[0]);
print("Number of rejected projects with numbers in resource summary: ", numericResourceRejectedData.shape[0]);
print("Number of approved projects without numbers in resource summary: ", textResourceApprovedData.shape[0]);
print("Number of rejected projects without numbers in resource summary: ", textResourceRejectedData.shape[0]);
```

Checking whether numbers in resource summary will be useful for project approval?

---

```
=====
Number of approved projects with numbers in resource summary: 14090
Number of rejected projects with numbers in resource summary: 1666
Number of approved projects without numbers in resource summary: 78616
Number of rejected projects without numbers in resource summary: 14876
```

## Observation:

1. The rejection rate of project is less when projects resource summary has numbers in it.
2. Even the number of projects approved without numbers in resource summary is high which means that the classification does not actually depends on whether resource summary contains numerical digits or not.

## Conclusion of univariate analysis:

1. There is huge overlap of approved and rejected projects when taken for all single features. So, this project cannot be classified using single features.
2. project\_title is some what better in text type of feature because of less overlap than others.
3. The project approval is not depending on resources cost, but the probability of project rejection is more when resources cost is more.

## Preprocessing data

In [0]:

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# All stopwords that are needed to be removed in the text
stopWords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
    "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
    'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'it's', 'itself', 'they', 'them', 'their', \
    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
    'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
    'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
    'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
    'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \
    'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
```

```

's', 't', 'can', 'will', 'just', 'don', "don't", 'should',
"should've", 'now', 'd', 'll', 'm', 'o', 're', \
    've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't",
'didn', "didn't", 'doesn', "doesn't", 'hadn',\
    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "is
n't", 'ma', 'mightn', "mightn't", 'mustn',\
    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn',
"shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
    'won', "won't", 'wouldn', "wouldn't"]);
def preProcessingWithAndWithoutStopWords(texts):
"""
    This function takes list of texts and returns preprocessed list of
texts one with
stop words and one without stopwords.
"""

# Variable for storing preprocessed text with stop words
preProcessedTextsWithStopWords = [];
# Variable for storing preprocessed text without stop words
preProcessedTextsWithoutStopWords = [];

# Looping over list of texts for performing pre processing
for text in tqdm(texts, total = len(texts)):
    # Removing all links in the text
    text = re.sub(r"http\S+", "", text);

    # Removing all html tags in the text
    text = re.sub(r"<\w+>", "", text);
    text = re.sub(r"<\w+>", "", text);

    # https://stackoverflow.com/a/47091490/4084039
    # Replacing all below words with adverbs
    text = re.sub(r"won't", "will not", text)
    text = re.sub(r"can't", "can not", text)
    text = re.sub(r"n't", " not", text)
    text = re.sub(r'\re', " are", text)
    text = re.sub(r'\s', " is", text)
    text = re.sub(r'\d', " would", text)
    text = re.sub(r'\ll', " will", text)
    text = re.sub(r'\t", " not", text)

```

```

text = re.sub(r"\'ve", " have", text)
text = re.sub(r"\'m", " am", text)

# Removing backslash symbols in text
text = text.replace('\\r', ' ');
text = text.replace('\\n', ' ');
text = text.replace('\\', ' ');

# Removing all special characters of text
text = re.sub(r"[^a-zA-Z0-9]+", " ", text);

# Converting whole review text into lower case
text = text.lower();

# adding this preprocessed text without stopwords to list
preProcessedTextsWithStopWords.append(text);

# removing stop words from text
textWithoutStopWords = ' '.join([word for word in text.split()
if word not in stopWords]);
# adding this preprocessed text without stopwords to list
preProcessedTextsWithoutStopWords.append(textWithoutStopWords);

return [preProcessedTextsWithStopWords, preProcessedTextsWithoutStopWords];

```

In [59]:

```

texts = [projectsData['project_essay'].values[0]]
preProcessedTextsWithStopWords, preProcessedTextsWithoutStopWords = pre
ProcessingWithAndWithoutStopWords(texts);
print("Example project essay without pre-processing: ");
equalsBorder(70);
print(texts);
equalsBorder(70);
print("Example project essay with stop words and pre-processing: ");
equalsBorder(70);
print(preProcessedTextsWithStopWords);
equalsBorder(70);
print("Example project essay without stop words and pre-processing: ");

```

```
equalsBorder(70);  
print(preProcessedTextsWithoutStopWords);
```

Example project essay without pre-processing:

```
===== ['My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrant s, and native-born Americans bringing the gift of language to our schoo l. \\r\\n\\r\\n We have over 24 languages represented in our English Le arner program with students at every level of mastery. We also have ov er 40 countries represented with the families within our school. Each student brings a wealth of knowledge and experiences to us that open ou r eyes to new cultures, beliefs, and respect.\\\"The limits of your lang uage are the limits of your world.\\\"-Ludwig Wittgenstein Our English learner\\'s have a strong support system at home that begs for more reso urces. Many times our parents are learning to read and speak English a long side of their children. Sometimes this creates barriers for paren ts to be able to help their child learn phonetics, letter recognition, and other reading skills.\\r\\n\\r\\nBy providing these dvd\\'s and play ers, students are able to continue their mastery of the English languag e even if no one at home is able to assist. All families with students within the Level 1 proficiency status, will be offered to be a part o f this program. These educational videos will be specially chosen by t he English Learner Teacher and will be sent home regularly to watch. T he videos are to help the child develop early reading skills.\\r\\n\\r \\nParents that do not have access to a dvd player will have the opport unity to check out a dvd player to use for the year. The plan is to us e these videos and educational dvd\\'s for the years to come for other E L students.\\r\\nnannan']
```

Example project essay with stop words and pre-processing:

```
===== ['my students are english learners that are working on english as their second or third languages we are a melting pot of refugees immigrants a nd native born americans bringing the gift of language to our school we have over 24 languages represented in our english learner program with students at every level of mastery we also have over 40 countries repre sented with the families within our school each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures b eliefs and respect the limits of your language are the limits of your w
```

effects and respect the limits of your language are the limits of your w

orld ludwig wittgenstein our english learner is have a strong support s ystem at home that begs for more resources many times our parents are l earning to read and speak english along side of their children sometime s this creates barriers for parents to be able to help their child learn phonetics letter recognition and other reading skills by providing th ese dvd is and players students are able to continue their mastery of t he english language even if no one at home is able to assist all famili es with students within the level 1 proficiency status will be a offere d to be a part of this program these educational videos will be special ly chosen by the english learner teacher and will be sent home regularl y to watch the videos are to help the child develop early reading skill s parents that do not have access to a dvd player will have the opportu nity to check out a dvd player to use for the year the plan is to use t hese videos and educational dvd is for the years to come for other el s tudents nannan']

=====

Example project essay without stop words and pre-processing:

=====

```
['students english learners working english second third languages melt ing pot refugees immigrants native born americans bringing gift languag e school 24 languages represented english learner program students ever y level mastery also 40 countries represented families within school st udent brings wealth knowledge experiences us open eyes new cultures bel iefs respect limits language limits world ludwig wittgenstein english l earner strong support system home begs resources many times parents lea rning read speak english along side children sometimes creates barriers parents able help child learn phonetics letter recognition reading skil ls providing dvd players students able continue mastery english languag e even no one home able assist families students within level 1 profici ency status offered part program educational videos specially chosen en glish learner teacher sent home regularly watch videos help child devel op early reading skills parents not access dvd player opportunity check dvd player use year plan use videos educational dvd years come el stude nts nannan']
```

```
In [60]: projectEssays = projectsData['project_essay'];
preProcessedEssaysWithStopWords, preProcessedEssaysWithoutStopWords = p reProcessingWithAndWithoutStopWords(projectEssays);
```

```
In [61]: preProcessedEssaysWithoutStopWords[0:3]
```

```
Out[61]: ['students english learners working english second third languages melt  
ing pot refugees immigrants native born americans bringing gift languag  
e school 24 languages represented english learner program students ever  
y level mastery also 40 countries represented families within school st  
udent brings wealth knowledge experiences us open eyes new cultures bel  
iefs respect limits language limits world ludwig wittgenstein english l  
earner strong support system home begs resources many times parents lea  
rning read speak english along side children sometimes creates barriers  
parents able help child learn phonetics letter recognition reading skil  
ls providing dvd players students able continue mastery english languag  
e even no one home able assist families students within level 1 profici  
ency status offered part program educational videos specially chosen en  
glish learner teacher sent home regularly watch videos help child devel  
op early reading skills parents not access dvd player opportunity check  
dvd player use year plan use videos educational dvd years come el stude  
nts nannan',  
'students arrive school eager learn polite generous strive best know e  
ducation succeed life help improve lives school focuses families low in  
comes tries give student education deserve not much students use materi  
als given best projector need school crucial academic improvement stude  
nts technology continues grow many resources internet teachers use grow  
th students however school limited resources particularly technology wi  
thout disadvantage one things could really help classrooms projector pr  
ojector not crucial instruction also growth students projector show pre  
sentations documentaries photos historical land sites math problems muc  
h projector make teaching learning easier also targeting different type  
s learners classrooms auditory visual kinesthetic etc nannan',  
'true champions not always ones win guts mia hamm quote best describes  
students cholla middle school approach playing sports especially girls  
boys soccer teams teams made 7th 8th grade students not opportunity pla  
y organized sport due family financial difficulties teach title one mid  
dle school urban neighborhood 74 students qualify free reduced lunch ma  
ny come activity sport opportunity poor homes students love participate  
sports learn new skills apart team atmosphere school lacks funding meet  
students needs concerned lack exposure not prepare participating sports
```

teams high school end school year goal provide students opportunity learn variety soccer skills positive qualities person actively participates team students campus come school knowing face uphill battle comes participating organized sports players would thrive field confidence appropriate soccer equipment play soccer best abilities students experience helpful person part team teaches positive supportive encouraging others students using soccer equipment practice games daily basis learn practice necessary skills develop strong soccer team experience create opportunity students learn part team positive contribution teammates students get opportunity learn practice variety soccer skills use skills game access type experience nearly impossible without soccer equipment students players utilize practice games nannan']

In [62]: projectTitles = projectsData['project\_title'];  
preProcessedProjectTitlesWithStopWords, preProcessedProjectTitlesWithoutStopWords = preProcessingWithAndWithoutStopWords(projectTitles);  
preProcessedProjectTitlesWithoutStopWords[0:5]

Out[62]: ['educational support english learners home',  
'wanted projector hungry learners',  
'soccer equipment awesome middle school students',  
'techie kindergarteners',  
'interactive math tools']

## Preparing data for classification and modelling

In [0]: pd.DataFrame(projectsData.columns, columns = ['All features in projects data'])

Out[0]:

	All features in projects data
0	Unnamed: 0
1	id
2	teacher_id

All features in projects data	
3	teacher_prefix
4	school_state
5	project_submitted_datetime
6	project_grade_category
7	project_subject_categories
8	project_subject_subcategories
9	project_title
10	project_essay_1
11	project_essay_2
12	project_essay_3
13	project_essay_4
14	project_resource_summary
15	teacher_number_of_previously_posted_projects
16	project_is_approved
17	cleaned_categories
18	cleaned_sub_categories
19	project_essay
20	price
21	quantity

### Useful features:

Here we will consider only below features for classification and we can ignore the other features

**Categorical data:**

1. **school\_state** - categorical data
2. **project\_grade\_category** - categorical data
3. **cleaned\_categories** - categorical data
4. **cleaned\_sub\_categories** - categorical data
5. **teacher\_prefix** - categorical data

**Text data:**

1. **project\_resource\_summary** - text data
2. **project\_title** - text data
3. **project\_resource\_summary** - text data

**Numerical data:**

1. **teacher\_number\_of\_previously\_posted\_projects** - numerical data
2. **price** - numerical data
3. **quantity** - numerical data

## Vectorizing categorical data

### 1. Vectorizing **cleaned\_categories**(**project\_subject\_categories cleaned**) - One Hot Encoding

```
In [0]: # Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique cleaned_categories  
subjectsCategoriesVectorizer = CountVectorizer(vocabulary = list(sorted CategoriesDictionary.keys()), lowercase = False, binary = True);  
# Fitting CountVectorizer with cleaned_categories values
```

```
subjectsCategoriesVectorizer.fit(projectsData['cleaned_categories'].values);
# Vectorizing categories using one-hot-encoding
categoriesVectors = subjectsCategoriesVectorizer.transform(projectsData
['cleaned_categories'].values);
```

```
In [0]: print("Features used in vectorizing categories: ");
equalsBorder(70);
print(subjectsCategoriesVectorizer.get_feature_names());
equalsBorder(70);
print("Shape of cleaned_categories matrix after vectorization(one-hot-e
ncoding): ", categoriesVectors.shape);
equalsBorder(70);
print("Sample vectors of categories: ");
equalsBorder(70);
print(categoriesVectors[0:4])
```

Features used in vectorizing categories:

```
=====
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearn
ing', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Langua
ge']
```

```
=====
Shape of cleaned_categories matrix after vectorization(one-hot-encodin
g): (109248, 9)
```

=====

Sample vectors of categories:

```
=====
(0, 8)      1
(1, 2)      1
(1, 6)      1
(2, 6)      1
(3, 7)      1
(3, 8)      1
```

## 2. Vectorizing cleaned\_sub\_categories(project\_subject\_sub\_categories cleaned) - One Hot Encoding

```
In [0]: # Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique cleaned_sub_categories
subjectsSubCategoriesVectorizer = CountVectorizer(vocabulary = list(sortedDictionarySubCategories.keys()), lowercase = False, binary = True);
# Fitting CountVectorizer with cleaned_sub_categories values
subjectsSubCategoriesVectorizer.fit(projectsData['cleaned_sub_categories'].values);
# Vectorizing sub categories using one-hot-encoding
subCategoriesVectors = subjectsSubCategoriesVectorizer.transform(projectsData['cleaned_sub_categories'].values);
```

```
In [0]: print("Features used in vectorizing subject sub categories: ");
equalsBorder(70);
print(subjectsSubCategoriesVectorizer.get_feature_names());
equalsBorder(70);
print("Shape of cleaned_categories matrix after vectorization(one-hot-encoding): ", subCategoriesVectors.shape);
equalsBorder(70);
print("Sample vectors of categories: ");
equalsBorder(70);
print(subCategoriesVectors[0:4])
```

Features used in vectorizing subject sub categories:

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

```
=====
Shape of cleaned_categories matrix after vectorization(one-hot-encoding): (109248, 30)
```

```
=====
Sample vectors of categories:
```

```
(0, 20)      1
(0, 29)      1
(1, 5)       1
(1, 13)      1
(2, 13)      1
(2, 24)      1
(3, 28)      1
(3, 29)      1
```

### 3. Vectorizing teacher\_prefix - One Hot Encoding

```
In [0]: def giveCounter(data):
    counter = Counter();
    for dataValue in data:
        counter.update(str(dataValue).split());
    return counter
```

```
In [0]: giveCounter(projectsData['teacher_prefix'].values)
```

```
Out[0]: Counter({'Mrs.': 57269,
                  'Mr.': 10648,
                  'Ms.': 38955,
                  'Teacher': 2360,
                  'nan': 3,
                  'Dr.': 13})
```

```
In [0]: projectsData = projectsData.dropna(subset = ['teacher_prefix']);
projectsData.shape
```

```
Out[0]: (109245, 22)
```

```
In [0]: teacherPrefixDictionary = dict(giveCounter(projectsData['teacher_prefix'].values));
# Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique teacher_prefix
teacherPrefixVectorizer = CountVectorizer(vocabulary = list(teacherPrefixDictionary.keys()), lowercase = False, binary = True);
```

```
# Fitting CountVectorizer with teacher_prefix values
teacherPrefixVectorizer.fit(projectsData['teacher_prefix'].values);
# Vectorizing teacher_prefix using one-hot-encoding
teacherPrefixVectors = teacherPrefixVectorizer.transform(projectsData[
    'teacher_prefix'].values);
```

```
In [0]: print("Features used in vectorizing teacher_prefix: ");
equalsBorder(70);
print(teacherPrefixVectorizer.get_feature_names());
equalsBorder(70);
print("Shape of teacher_prefix matrix after vectorization(one-hot-encoding): ", teacherPrefixVectors.shape);
equalsBorder(70);
print("Sample vectors of teacher_prefix: ");
equalsBorder(70);
print(teacherPrefixVectors[0:100]);
```

Features used in vectorizing teacher\_prefix:

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

Shape of teacher\_prefix matrix after vectorization(one-hot-encoding):  
(109245, 5)

Sample vectors of teacher\_prefix:

```
=====
(27, 3)      1
(75, 3)      1
(82, 3)      1
(88, 3)      1
=====
```

```
In [0]: teacherPrefixes = [prefix.replace('.', '') for prefix in projectsData[
    'teacher_prefix'].values];
teacherPrefixes[0:5]
```

Out[0]: ['Mrs', 'Mr', 'Ms', 'Mrs', 'Mrs']

```
In [0]: projectsData['teacher_prefix'] = teacherPrefixes;
projectsData.head(3)
```

Out[0]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs	IN
1	140945	p258326	897464ce9ddc600bcfd1151f324dd63a	Mr	FL
2	21895	p182444	3465aa82da834c0582ebd0ef8040ca0	Ms	AZ

3 rows × 22 columns

In [0]:

```
teacherPrefixDictionary = dict(giveCounter(projectsData['teacher_prefix'].values));
# Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique teacher_prefix
teacherPrefixVectorizer = CountVectorizer(vocabulary = list(teacherPrefixDictionary.keys()), lowercase = False, binary = True);
# Fitting CountVectorizer with teacher_prefix values
teacherPrefixVectorizer.fit(projectsData['teacher_prefix'].values);
# Vectorizing teacher_prefix using one-hot-encoding
teacherPrefixVectors = teacherPrefixVectorizer.transform(projectsData['teacher_prefix'].values);
```

```
In [0]: print("Features used in vectorizing teacher_prefix: ");
equalsBorder(70);
print(teacherPrefixVectorizer.get_feature_names());
equalsBorder(70);
print("Shape of teacher_prefix matrix after vectorization(one-hot-encoding): ", teacherPrefixVectors.shape);
equalsBorder(70);
print("Sample vectors of teacher_prefix: ");
equalsBorder(70);
print(teacherPrefixVectors[0:4]);
```

Features used in vectorizing teacher\_prefix:

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

Shape of teacher\_prefix matrix after vectorization(one-hot-encoding):  
(109245, 5)

Sample vectors of teacher\_prefix:

```
=====
(0, 0)      1
(1, 1)      1
(2, 2)      1
(3, 0)      1
=====
```

#### 4. Vectorizing school\_state - One Hot Encoding

```
In [0]: schoolStateDictionary = dict(giveCounter(projectsData['school_state'].values));
# Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique school states
schoolStateVectorizer = CountVectorizer(vocabulary = list(schoolStateDictionary.keys()), lowercase = False, binary = True);
# Fitting CountVectorizer with school_state values
schoolStateVectorizer.fit(projectsData['school_state'].values);
# Vectorizing school_state using one-hot-encoding
schoolStateVectors = schoolStateVectorizer.transform(projectsData['school_state'].values);
```

```
In [0]: print("Features used in vectorizing school_state: ");
equalsBorder(70);
print(schoolStateVectorizer.get_feature_names());
equalsBorder(70);
print("Shape of school_state matrix after vectorization(one-hot-encoding): ", schoolStateVectors.shape);
equalsBorder(70);
print("Sample vectors of school_state: ");
equalsBorder(70);
print(schoolStateVectors[0:4]);
```

Features used in vectorizing school\_state:  
=====

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

Shape of school\_state matrix after vectorization(one-hot-encoding): (109245, 51)  
=====

Sample vectors of school\_state:  
=====

(0, 0) 1  
(1, 1) 1  
(2, 2) 1  
(3, 3) 1

## 5. Vectorizing project\_grade\_category - One Hot Encoding

```
In [0]: giveCounter(projectsData['project_grade_category'])
```

Out[0]: Counter({'Grades': 109245,  
 'PreK-2': 44225,  
 '6-8': 16923,

```
'3-5': 37135,  
'9-12': 10962})
```

```
In [0]: cleanedGrades = []  
for grade in projectsData['project_grade_category'].values:  
    grade = grade.replace(' ', '';  
    grade = grade.replace('-', 'to');  
    cleanedGrades.append(grade);  
cleanedGrades[0:4]
```

```
Out[0]: ['GradesPreKto2', 'Grades6to8', 'Grades6to8', 'GradesPreKto2']
```

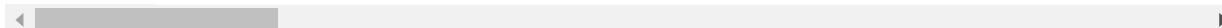
```
In [0]: projectsData['project_grade_category'] = cleanedGrades  
projectsData.head(4)
```

```
Out[0]:
```

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs	IN
1	140945	p258326	897464ce9ddc600bcfd1151f324dd63a	Mr	FL
2	21895	p182444	3465aa82da834c0582ebd0ef8040ca0	Ms	AZ

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs	KY

4 rows × 22 columns



```
In [0]: projectGradeDictionary = dict(giveCounter(projectsData['project_grade_category']));
# Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique project grade categories
projectGradeVectorizer = CountVectorizer(vocabulary = list(projectGradeDictionary.keys()), lowercase = False, binary = True);
# Fitting CountVectorizer with project_grade_category values
projectGradeVectorizer.fit(projectsData['project_grade_category'].values);
# Vectorizing project_grade_category using one-hot-encoding
projectGradeVectors = projectGradeVectorizer.transform(projectsData['project_grade_category'].values);
```

```
In [0]: print("Features used in vectorizing project_grade_category: ");
equalsBorder(70);
print(projectGradeVectorizer.get_feature_names());
equalsBorder(70);
print("Shape of school_state matrix after vectorization(one-hot-encoding): ", projectGradeVectors.shape);
equalsBorder(70);
print("Sample vectors of school_state: ");
equalsBorder(70);
print(projectGradeVectors[0:4]);
```

Features used in vectorizing project\_grade\_category:

=====

['GradesPreKto2', 'Grades6to8', 'Grades3to5', 'Grades9to12']

```
=====
Shape of school_state matrix after vectorization(one-hot-encoding): (1
09245, 4)
=====
Sample vectors of school_state:
=====
(0, 0)      1
(1, 1)      1
(2, 1)      1
(3, 0)      1
```

```
In [0]: projectsDataSub = projectsData[0:40000];
preProcessedEssaysWithoutStopWordsSub = preProcessedEssaysWithoutStopWo
rds[0:40000];
preProcessedProjectTitlesWithoutStopWordsSub = preProcessedProjectTitle
sWithoutStopWords[0:40000];
```

## Vectorizing Text Data

### Bag of Words

#### 1. Vectorizing project\_essay

```
In [0]: # Initializing countvectorizer for bag of words vectorization of prepro
cessed project essays
bowEssayVectorizer = CountVectorizer(min_df = 10);
# Transforming the preprocessed essays to bag of words vectors
bowEssayModel = bowEssayVectorizer.fit_transform(preProcessedEssaysWith
outStopWordsSub);
```

```
In [0]: print("Some of the Features used in vectorizing preprocessed essays: "
);
equalsBorder(70);
print(bowEssayVectorizer.get_feature_names()[-40:]);
```

```
equalsBorder(70);
print("Shape of preprocessed essay matrix after vectorization: ", bowEs
sayModel.shape);
equalsBorder(70);
print("Sample bag-of-words vector of preprocessed essay: ");
equalsBorder(70);
print(bowEssayModel[0])
```

Some of the Features used in vectorizing preprocessed essays:

```
=====
['yeats', 'yell', 'yelling', 'yellow', 'yemen', 'yes', 'yesterday', 'ye
t', 'yield', 'yields', 'yoga', 'york', 'younannan', 'young', 'younger',
'youngest', 'youngsters', 'youth', 'youthful', 'youths', 'youtube', 'yu
mmy', 'zeal', 'zearn', 'zen', 'zenergy', 'zero', 'zest', 'zip', 'ziplo
c', 'zippers', 'zipping', 'zone', 'zoned', 'zones', 'zoo', 'zoom', 'zoo
ming', 'zoos', 'zumba']
```

```
=====
Shape of preprocessed essay matrix after vectorization: (40000, 11077)
```

```
=====
Sample bag-of-words vector of preprocessed essay:
```

```
=====
(0, 6533)    1
(0, 3306)    1
(0, 1981)    1
(0, 11036)   1
(0, 7347)    1
(0, 11029)   1
(0, 10530)   2
(0, 1734)    1
(0, 6855)    1
(0, 7374)    2
(0, 232)     1
(0, 6687)    1
(0, 3211)    1
(0, 2805)    1
(0, 10766)   1
(0, 8133)    1
(0, 8803)    1
(0, 9831)    1
(0, 1794)    1
```

```
(0, 9237)      1
(0, 10639)     3
(0, 3274)      2
(0, 7068)      1
(0, 6798)      1
(0, 9399)      1
:
:
(0, 6123)      2
(0, 5785)      2
(0, 3613)      1
(0, 7703)      2
(0, 5732)      3
(0, 8269)      2
(0, 67)         1
(0, 8670)      2
(0, 5664)      3
(0, 4383)      1
(0, 1339)      1
(0, 553)        1
(0, 1248)      1
(0, 6549)      1
(0, 5003)      1
(0, 8116)      1
(0, 7501)      1
(0, 6207)      1
(0, 5665)      2
(0, 9968)      1
(0, 8736)      1
(0, 10964)     1
(0, 5733)      1
(0, 3449)      7
(0, 9553)      5
```

## 2. Vectorizing project\_title

```
In [0]: # Initializing countvectorizer for bag of words vectorization of preprocessed project titles
bowTitleVectorizer = CountVectorizer(min_df = 10);
```

```
# Transforming the preprocessed project titles to bag of words vectors
bowTitleModel = bowTitleVectorizer.fit_transform(preProcessedProjectTitlesWithoutStopWordsSub);
```

```
In [0]: print("Some of the Features used in vectorizing preprocessed titles: ")
        );
        equalsBorder(70);
        print(bowTitleVectorizer.get_feature_names()[-40:]);
        equalsBorder(70);
        print("Shape of preprocessed title matrix after vectorization: ", bowTitleModel.shape);
        equalsBorder(70);
        print("Sample bag-of-words vector of preprocessed title: ");
        equalsBorder(70);
        print(bowTitleModel[0])
```

Some of the Features used in vectorizing preprocessed titles:

```
=====
['wireless', 'wise', 'wish', 'within', 'without', 'wizards', 'wo', 'wobble',
 'wobbles', 'wobbling', 'wobbly', 'wonder', 'wonderful', 'wonders',
 'word', 'words', 'work', 'workers', 'working', 'works', 'workshop',
 'world', 'worlds', 'worms', 'worth', 'would', 'wow', 'write', 'writer',
 'writers', 'writing', 'ye', 'year', 'yearbook', 'yes', 'yoga', 'young',
 'youth', 'zone', 'zoom']
```

Shape of preprocessed title matrix after vectorization: (40000, 1774)

Sample bag-of-words vector of preprocessed title:

```
=====
(0, 766)      1
(0, 906)      1
(0, 514)      1
(0, 1553)     1
(0, 483)      1
```

## Tf-Idf Vectorization

## 1. Vectorizing project\_essay

```
In [0]: # Intializing tfidf vectorizer for tf-idf vectorization of preprocessed  
# project essays  
tfIdfEssayVectorizer = TfidfVectorizer(min_df = 10);  
# Transforming the preprocessed project essays to tf-idf vectors  
tfIdfEssayModel = tfIdfEssayVectorizer.fit_transform(preProcessedEssays  
WithoutStopWordsSub);
```

```
In [0]: print("Some of the Features used in tf-idf vectorizing preprocessed ess  
ays: ");  
equalsBorder(70);  
print(tfIdfEssayVectorizer.get_feature_names()[-40:]);  
equalsBorder(70);  
print("Shape of preprocessed title matrix after tf-idf vectorization: "  
, tfIdfEssayModel.shape);  
equalsBorder(70);  
print("Sample Tf-Idf vector of preprocessed essay: ");  
equalsBorder(70);  
print(tfIdfEssayModel[0])
```

Some of the Features used in tf-idf vectorizing preprocessed essays:

```
===== ['yeats', 'yell', 'yelling', 'yellow', 'yemen', 'yes', 'yesterday', 'ye  
t', 'yield', 'yields', 'yoga', 'york', 'younannan', 'young', 'younger',  
'youngest', 'youngsters', 'youth', 'youthful', 'youths', 'youtube', 'yu  
mmy', 'zeal', 'zearn', 'zen', 'zenergy', 'zero', 'zest', 'zip', 'ziplo  
c', 'zippers', 'zipping', 'zone', 'zoned', 'zones', 'zoo', 'zoom', 'zoo  
ming', 'zoos', 'zumba'] =====
```

```
===== Shape of preprocessed title matrix after tf-idf vectorization: (40000,  
11077) =====
```

```
===== Sample Tf-Idf vector of preprocessed essay:
```

```
===== (0, 9553)      0.07732161197654648  
(0, 3449)      0.2978137199079083  
(0, 5733)      0.03611311825070974  
(0, 10964)     0.03819325396356506 =====
```

(0, 8736)	0.04966730436190034
(0, 9968)	0.05933894161734909
(0, 5665)	0.13189136979245247
(0, 6207)	0.09909858268088724
(0, 7501)	0.09797369103397546
(0, 8116)	0.09716121418147701
(0, 5003)	0.09174889764250635
(0, 6549)	0.07739523816315956
(0, 1248)	0.09041771504928811
(0, 553)	0.09502243963232913
(0, 1339)	0.07922532406820633
(0, 4383)	0.08387324724715874
(0, 5664)	0.12052414724469786
(0, 8670)	0.03565737676523101
(0, 67)	0.0797508795755641
(0, 8269)	0.18440093271700464
(0, 5732)	0.23244852084297085
(0, 7703)	0.0932371184396508
(0, 3613)	0.033250154942777416
(0, 5785)	0.08336998078832462
(0, 6123)	0.18451571587493337
:	:
(0, 9399)	0.0680639151319745
(0, 6798)	0.08632328546640713
(0, 7068)	0.046135007257522224
(0, 3274)	0.10489683635458984
(0, 10639)	0.2063461965343629
(0, 9237)	0.1100116652395096
(0, 1794)	0.07900547931629058
(0, 9831)	0.03792376194008962
(0, 8803)	0.09740047454864696
(0, 8133)	0.09001501053091984
(0, 10766)	0.07024528926492071
(0, 2805)	0.05089165427462248
(0, 3211)	0.06222851802675729
(0, 6687)	0.022226920710368445
(0, 232)	0.040248356980164615
(0, 7374)	0.1846309297399045
(0, 6855)	0.03799907965204156

```
(0, 1734)      0.07743897673831124
(0, 10530)     0.05491069896079749
(0, 11029)     0.030886589234837624
(0, 7347)      0.06268239285732621
(0, 11036)     0.04610937510882687
(0, 1981)       0.02654012905964554
(0, 3306)       0.1031894334469226
(0, 6533)       0.016043824658976313
```

## 2. Vectorizing project\_title

```
In [0]: # Intializing tfidf vectorizer for tf-idf vectorization of preprocessed
# project titles
tfIdfTitleVectorizer = TfidfVectorizer(min_df = 10);
# Transforming the preprocessed project titles to tf-idf vectors
tfIdfTitleModel = tfIdfTitleVectorizer.fit_transform(preProcessedProjec
tTitlesWithoutStopWordsSub);
```

```
In [0]: print("Some of the Features used in tf-idf vectorizing preprocessed tit
les: ");
equalsBorder(70);
print(tfIdfTitleVectorizer.get_feature_names()[-40:]);
equalsBorder(70);
print("Shape of preprocessed title matrix after tf-idf vectorization: "
, tfIdfTitleModel.shape);
equalsBorder(70);
print("Sample Tf-Idf vector of preprocessed title: ");
equalsBorder(70);
print(tfIdfTitleModel[0])
```

Some of the Features used in tf-idf vectorizing preprocessed titles:

```
=====
['wireless', 'wise', 'wish', 'within', 'without', 'wizards', 'wo', 'wob
ble', 'wobbles', 'wobbling', 'wobbly', 'wonder', 'wonderful', 'wonder
s', 'word', 'words', 'work', 'workers', 'working', 'works', 'workshop',
'world', 'worlds', 'worms', 'worth', 'would', 'wow', 'write', 'writer',
'writers', 'writing', 'ye', 'year', 'yearbook', 'yes', 'yoga', 'young',
'youth', 'zone', 'zoom']
```

```
=====
Shape of preprocessed title matrix after tf-idf vectorization: (40000,
1774)
=====
Sample Tf-Idf vector of preprocessed title:
=====
(0, 483)      0.5356140846908081
(0, 1553)     0.4441059196924978
(0, 514)       0.4615835742389133
(0, 906)       0.3400969810242112
(0, 766)       0.4326223894644794
```

## Average Word2Vector Vectorization

```
In [0]: # stronging variables into pickle files python: http://www.jessicayung.
com/how-to-use-pickle-to-save-and-load-variables-in-python/
# We should have glove_vectors file for creating below model
with open('glove_vectors', 'rb') as f:
    gloveModel = pickle.load(f)
    gloveWords = set(gloveModel.keys())
```

```
In [0]: print("Glove vector of sample word: ");
equalsBorder(70);
print(gloveModel['technology']);
equalsBorder(70);
print("Shape of glove vector: ", gloveModel['technology'].shape);
```

Glove vector of sample word:

```
=====
[-0.26078 -0.36898 -0.022831  0.21666  0.16672 -0.20268
 -3.1219   0.33057  0.71512   0.28874  0.074368 -0.033203
  0.23783  0.21052  0.076562   0.13007 -0.31706 -0.45888
 -0.45463 -0.13191  0.49761   0.072704  0.16811  0.18846
 -0.16688 -0.21973  0.08575   -0.19577 -0.2101   -0.32436
 -0.56336  0.077996 -0.22758   -0.66569  0.14824  0.038945
  0.50881 -0.1352   0.49966   -0.4401   -0.022335 -0.22744
  0.22086  0.21865  0.36647   0.30495 -0.16565  0.038759]
```

0.28108	-0.2167	0.12453	0.65401	0.34584	-0.2557
-0.046363	-0.31111	-0.020936	-0.17122	-0.77114	0.29289
-0.14625	0.39541	-0.078938	0.051127	0.15076	0.085126
0.183	-0.06755	0.26312	0.0087276	0.0066415	0.37033
0.03496	-0.12627	-0.052626	-0.34897	0.14672	0.14799
-0.21821	-0.042785	0.2661	-1.1105	0.31789	0.27278
0.054468	-0.27458	0.42732	-0.44101	-0.19302	-0.32948
0.61501	-0.22301	-0.36354	-0.34983	-0.16125	-0.17195
-3.363	0.45146	-0.13753	0.31107	0.2061	0.33063
0.45879	0.24256	0.042342	0.074837	-0.12869	0.12066
0.42843	-0.4704	-0.18937	0.32685	0.26079	0.20518
-0.18432	-0.47658	0.69193	0.18731	-0.12516	0.35447
-0.1969	-0.58981	-0.88914	0.5176	0.13177	-0.078557
0.032963	-0.19411	0.15109	0.10547	-0.1113	-0.61533
0.0948	-0.3393	-0.20071	-0.30197	0.29531	0.28017
0.16049	0.25294	-0.44266	-0.39412	0.13486	0.25178
-0.044114	1.1519	0.32234	-0.34323	-0.10713	-0.15616
0.031206	0.46636	-0.52761	-0.39296	-0.068424	-0.04072
0.41508	-0.34564	0.71001	-0.364	0.2996	0.032281
0.34035	0.23452	0.78342	0.48045	-0.1609	0.40102
-0.071795	-0.16531	0.082153	0.52065	0.24194	0.17113
0.33552	-0.15725	-0.38984	0.59337	-0.19388	-0.39864
-0.47901	1.0835	0.24473	0.41309	0.64952	0.46846
0.024386	-0.72087	-0.095061	0.10095	-0.025229	0.29435
-0.57696	0.53166	-0.0058338	-0.3304	0.19661	-0.085206
0.34225	0.56262	0.19924	-0.027111	-0.44567	0.17266
0.20887	-0.40702	0.63954	0.50708	-0.31862	-0.39602
-0.1714	-0.040006	-0.45077	-0.32482	-0.0316	0.54908
-0.1121	0.12951	-0.33577	-0.52768	-0.44592	-0.45388
0.66145	0.33023	-1.9089	0.5318	0.21626	-0.13152
0.48258	0.68028	-0.84115	-0.51165	0.40017	0.17233
-0.033749	0.045275	0.37398	-0.18252	0.19877	0.1511
0.029803	0.16657	-0.12987	-0.50489	0.55311	-0.22504
0.13085	-0.78459	0.36481	-0.27472	0.031805	0.53052
-0.20078	0.46392	-0.63554	0.040289	-0.19142	-0.0097011
0.068084	-0.10602	0.25567	0.096125	-0.10046	0.15016
-0.26733	-0.26494	0.057888	0.062678	-0.11596	0.28115
0.25375	-0.17954	0.20615	0.24189	0.062696	0.27719
-0.42601	-0.28619	-0.44697	-0.082253	-0.73415	-0.20675

```
-0.60289 -0.06728 0.15666 -0.042614 0.41368 -0.17367  
-0.54012 0.23883 0.23075 0.13608 -0.058634 -0.089705  
0.18469 0.023634 0.16178 0.23384 0.24267 0.091846 ]  
=====  
Shape of glove vector: (300,)
```

```
In [0]: def getWord2VecVectors(texts):  
    word2VecTextsVectors = [];  
    for preProcessedText in tqdm(texts):  
        word2VecTextVector = np.zeros(300);  
        numberOfWorksInText = 0;  
        for word in preProcessedText.split():  
            if word in gloveWords:  
                word2VecTextVector += gloveModel[word];  
                numberOfWorksInText += 1;  
        if numberOfWorksInText != 0:  
            word2VecTextVector = word2VecTextVector / numberOfWorksInText;  
        word2VecTextsVectors.append(word2VecTextVector);  
    return word2VecTextsVectors;
```

## 1. Vectorizing project\_essay

```
In [0]: word2VecEssaysVectors = getWord2VecVectors(preProcessedEssaysWithoutStopWords);
```

```
In [0]: print("Shape of Word2Vec vectorization matrix of essays: {},{}".format(len(word2VecEssaysVectors), len(word2VecEssaysVectors[0])));  
equalsBorder(70);  
print("Sample essay: ");  
equalsBorder(70);  
print(preProcessedEssaysWithoutStopWords[0]);  
equalsBorder(70);  
print("Word2Vec vector of sample essay: ");  
equalsBorder(70);  
print(word2VecEssaysVectors[0]);
```

```
Shape of Word2Vec vectorization matrix of essays: 109248,300
```

```
=====
```

```
Sample essay:
```

```
=====
```

```
students english learners working english second third languages meltin  
g pot refugees immigrants native born americans bringing gift language  
school 24 languages represented english learner program students every  
level mastery also 40 countries represented families within school stud  
ent brings wealth knowledge experiences us open eyes new cultures belie  
fs respect limits language limits world ludwig wittgenstein english lea  
rner strong support system home begs resources many times parents learn  
ing read speak english along side children sometimes creates barriers p  
arents able help child learn phonetics letter recognition reading skill  
s providing dvd players students able continue mastery english language  
even no one home able assist families students within level 1 proficien  
cy status offered part program educational videos specially chosen engl  
ish learner teacher sent home regularly watch videos help child develop  
early reading skills parents not access dvd player opportunity check dv  
d player use year plan use videos educational dvd years come el student  
s nannan
```

```
=====
```

```
Word2Vec vector of sample essay:
```

```
=====
```

```
[-1.40030644e-02 8.78995685e-02 3.50108161e-02 -5.90358980e-03  
-5.93166809e-02 -6.21039893e-02 -2.96711248e+00 9.45840302e-02  
-8.18737785e-03 4.46964161e-02 -7.64722101e-02 6.97099444e-02  
8.44441262e-02 -1.22974138e-01 -3.55310208e-02 -8.90947154e-02  
1.20959579e-01 -1.21977699e-01 4.61334597e-02 -3.33640832e-02  
1.24900557e-01 7.18837631e-02 -6.14885114e-02 -2.67269047e-02  
6.82086621e-02 -3.60263034e-02 1.17172255e-01 -1.17868631e-01  
-1.13467710e-01 -9.25920168e-02 -2.42461725e-01 -7.92963658e-02  
3.52513154e-03 1.79752468e-01 -4.69217812e-02 -3.56593007e-02  
-7.95331477e-03 -6.71107383e-04 -1.80828067e-02 -1.16224805e-02  
-3.69645852e-02 1.61287176e-01 -1.75201329e-01 -6.02256376e-02  
1.48811886e-02 -9.00106181e-02 7.72160490e-02 7.42989819e-02  
-1.02682389e-02 -1.33311658e-01 -2.82030537e-02 -7.71051879e-03  
7.33988450e-02 3.54095087e-02 -5.80719597e-03 -8.70242758e-02  
-3.57117638e-02 2.78475651e-02 -1.54957291e-01 -3.24157495e-02  
-5.93266570e-02 -8.80254174e-02 2.18914318e-01 -1.22730395e-02
```

-1.05831485e-01	1.53985730e-01	7.15618933e-02	-3.97147470e-02
1.47169116e-01	-4.50476644e-03	-1.49678829e-01	5.52201396e-02
3.04915879e-02	-6.24086617e-02	-7.68483134e-02	-7.50149195e-02
-1.07105068e-01	-2.69954530e-02	1.28067340e-01	-3.42946330e-02
4.24139667e-02	-4.49685043e-01	1.52793905e-01	-9.06178181e-02
-6.67951510e-02	-2.72063766e-02	7.37261792e-02	-8.64977130e-02
1.64616877e-01	4.86745523e-02	-4.44542828e-02	-3.04823530e-02
2.63897436e-02	-6.59345034e-02	-5.21813664e-02	-7.45015886e-02
-2.21975948e+00	8.57858456e-02	7.73778584e-02	1.14644799e-01
-1.50536483e-01	-5.17326940e-02	3.23826117e-02	-1.15700542e-01
7.15651973e-02	9.15412617e-02	5.41334631e-02	-1.25451318e-01
2.80941483e-02	-3.95890262e-02	-1.67010497e-02	1.74708879e-02
4.58374505e-02	2.56664910e-01	3.74891134e-02	3.00990497e-02
-2.18904765e-01	9.37672966e-02	9.99403436e-02	5.26255996e-02
-6.67958718e-02	5.97650946e-02	4.14311192e-02	-6.85917603e-02
1.72453235e-02	1.02485026e-01	3.02940430e-02	9.59998859e-03
1.96364913e-02	1.22438477e-01	7.98410557e-02	1.92611322e-02
6.44085906e-03	4.94252148e-03	-5.36137718e-03	-1.17976934e-01
1.77991634e-01	-2.51954819e-02	8.02478188e-02	2.29125079e-01
3.79080403e-02	1.22892819e-02	7.19621470e-02	-9.25031570e-02
-8.86571674e-02	-4.74898563e-02	1.68688409e-02	-1.15134901e-01
1.76528904e-01	-6.30485141e-02	-4.99678329e-02	-1.00350507e-01
1.25089302e-02	-4.08706114e-02	4.50565289e-02	2.49286074e-02
-1.29713758e-03	-3.21404376e-02	-2.52972249e-02	-9.63531510e-02
8.42448993e-04	-7.29482953e-03	-3.77497893e-02	-9.35034987e-02
-3.45719793e-02	7.15921796e-02	-1.29330935e-01	1.28508101e-02
4.24846988e-02	-8.43078228e-02	4.79772134e-02	-3.05753799e-02
-3.03772013e-02	-2.10572558e-01	-1.05464289e-03	5.18230436e-02
-4.39921874e-02	5.29591584e-02	-1.08551689e-01	2.88053128e-02
-4.88957058e-02	2.31962381e-01	-2.90986193e-02	-2.83725755e-02
-6.80350899e-02	-6.99966387e-02	-6.80414679e-02	-7.63552362e-02
-1.59287859e-02	-2.59947651e-03	-7.81848121e-03	-1.14299579e-01
-2.02054698e-02	1.21184430e-03	2.59984919e-02	-7.64172013e-02
9.47882617e-03	-5.71751181e-02	1.25667972e-01	-4.60388139e-02
5.51296403e-02	-6.73280980e-02	-2.06862389e-02	1.12049165e-01
-7.63451436e-02	4.71124027e-02	6.32404235e-02	-2.13828034e-02
1.24239236e-01	5.08985235e-02	2.05136711e-03	1.45916498e-02
4.25123886e-02	-9.41766832e-02	-3.08569389e-02	-2.57995470e-02
-3.53808765e-02	-7.16000389e-02	1.35426121e-02	4.57596799e-02

```
-1.85721693e-01 -6.62042523e-02 -1.45448285e-01 5.50366758e-02
-2.09367026e+00 1.23479489e-01 -1.46630889e-01 -8.86940765e-02
-7.32806463e-02 -1.48629733e-01 3.23867248e-03 7.08553181e-02
1.10315906e-02 -2.35431879e-02 -7.69633283e-02 -1.13640894e-01
9.96301846e-02 -5.70585054e-02 -5.45997987e-04 9.42995174e-02
-1.40422433e-01 -5.03571812e-04 -2.50305216e-01 3.79384141e-02
-6.44086637e-02 -1.53146188e-02 -2.55858274e-02 -1.10195376e-01
1.62183899e-02 -1.61929591e-02 2.03421993e-02 1.21424534e-01
5.02740463e-02 2.37900799e-02 9.07398322e-02 1.57962685e-02
3.73036075e-02 -8.14876248e-02 1.37349395e-01 -8.17880913e-02
9.27907812e-02 6.76093826e-03 -5.22928389e-02 6.02994188e-02
8.28096711e-03 -1.05344042e-01 -1.02705751e-01 2.45275938e-02
-1.18970611e-02 9.86759282e-02 -1.92870134e-02 9.71936577e-03
-1.40249490e-01 1.61314103e-01 -4.55344879e-02 2.21929812e-02
9.54108215e-02 -1.25028370e-02 2.89625007e-02 1.65818081e-02
-2.34467852e-02 -7.88610081e-02 3.34242148e-03 4.43269879e-02
-4.08419376e-02 6.06990416e-02 2.33916564e-02 -1.02773899e-02
9.21596550e-02 9.90483805e-02 7.50525638e-03 -4.07725570e-03
-6.93980047e-02 -3.50341946e-02 -8.79849597e-02 -4.10474223e-02
4.55004698e-03 2.27073689e-01 1.37340472e-01 4.43856114e-02]
```

## 2. Vectorizing project\_title

```
In [0]: word2VecTitlesVectors = getWord2VecVectors(preProcessedProjectTitlesWithoutStopWords);
```

```
In [0]: print("Shape of Word2Vec vectorization matrix of project titles: {}, {}".format(len(word2VecTitlesVectors), len(word2VecTitlesVectors[0])));
equalsBorder(70);
print("Sample title: ");
equalsBorder(70);
print(preProcessedProjectTitlesWithoutStopWords[0]);
equalsBorder(70);
print("Word2Vec vector of sample title: ");
equalsBorder(70);
print(word2VecTitlesVectors[0]);
```

```
Shape of Word2Vec vectorization matrix of project titles: 109248, 300
=====
Sample title:
=====
educational support english learners home
=====
Word2Vec vector of sample title:
=====
[-4.1285000e-02 4.4970000e-02 1.4283080e-01 1.9901860e-02
 -8.4519200e-02 -4.3207400e-01 -2.8496800e+00 -2.2953320e-01
 2.1736960e-01 3.4239600e-01 -7.5568200e-02 1.8077600e-01
 1.3998316e-01 -1.6401800e-01 -2.9812820e-01 -2.5030200e-01
 2.0420960e-01 -1.6882720e-01 6.5439800e-02 -1.6061000e-01
 2.2179020e-01 2.9944900e-01 2.7358000e-02 -8.8528800e-02
 1.5856400e-01 6.2905000e-02 2.0427440e-01 -1.9312560e-01
 -9.2904600e-02 -2.2050020e-01 -5.7761060e-01 -1.2101294e-01
 1.6846980e-01 2.8212460e-01 -1.8210120e-01 1.7754000e-02
 1.4805200e-01 4.1059000e-02 3.1145000e-02 -9.5658000e-02
 -9.6840000e-03 2.4896520e-01 -2.5047440e-01 7.7859000e-02
 -3.7512000e-03 -2.7071920e-01 2.5586200e-02 2.3205600e-01
 1.0154800e-01 -5.2259200e-01 -1.3211440e-01 1.1908300e-01
 2.7147196e-01 5.6135400e-02 -5.3140200e-02 -1.4937160e-01
 -1.0488160e-01 1.2059600e-01 -1.2639620e-01 -1.4316640e-01
 -2.2147600e-01 -1.9137800e-01 1.6595340e-01 -5.6078000e-02
 3.9884400e-02 1.0854760e-01 1.5552920e-01 7.8204600e-02
 9.5928000e-02 -6.2156000e-03 -1.1407312e-01 3.6862800e-02
 -8.7530020e-02 -4.7668000e-02 -2.3264200e-01 -6.1687200e-02
 -3.1690916e-01 -1.1851380e-01 1.4931240e-01 -7.7857200e-02
 1.8634840e-01 -4.6202100e-01 2.7096800e-01 -3.0512800e-02
 -2.1226400e-01 -1.5356200e-02 1.0844260e-01 -8.2669200e-02
 2.8918600e-01 1.3372960e-01 -8.3522800e-02 4.6474200e-02
 2.0703580e-01 -2.1937640e-01 -1.0252400e-01 -2.5177000e-01
 -2.8408000e+00 1.6622880e-01 1.1216234e-01 2.0837920e-01
 -1.5711600e-01 -1.9159400e-01 -1.4992160e-01 -2.7392820e-01
 3.4989140e-01 1.3991600e-01 1.6275200e-01 1.3887200e-01
 1.8212760e-01 -3.2218600e-02 4.3172000e-02 1.8323640e-01
 1.2295780e-01 4.4706600e-01 2.1688400e-02 -3.8988200e-02
 -3.2467400e-01 3.8389160e-01 -1.4416560e-01 1.1117380e-01
 -1.6218300e-01 1.3871928e-01 1.4305240e-01 -7.6173200e-02
```

8.9476800e-02	2.6043820e-01	5.1114000e-02	1.0619800e-01
1.5968840e-01	1.0530680e-01	8.6300000e-02	1.4667260e-01
1.2320460e-02	-6.6124620e-02	-1.1017760e-01	-1.5091940e-01
2.1297280e-01	-3.2808520e-01	1.4493194e-01	2.1848680e-01
-4.1809800e-03	8.5340000e-02	-1.2410789e-01	-2.2308140e-01
8.8026000e-02	1.9555000e-01	-3.7981400e-02	-1.7720080e-01
3.4328600e-01	-3.7459600e-01	-1.7268200e-01	-2.1554400e-01
-1.1533400e-01	9.9680000e-02	-1.9032980e-01	8.6249800e-02
7.6682200e-02	-9.1090380e-02	-9.3714000e-02	-1.7333260e-01
8.6429960e-02	-6.7933600e-02	-8.6470600e-02	-2.2431600e-01
-2.8319800e-01	1.0138200e-01	-2.8114320e-01	-1.1168240e-01
2.1770560e-02	-1.3971160e-01	2.1795080e-01	-1.1995600e-01
-1.3166600e-02	-3.4848260e-01	-3.0102000e-02	2.3396200e-02
2.8840000e-02	2.8763000e-01	-2.3679600e-02	1.1806440e-01
-3.2261460e-01	2.2622920e-01	1.9506400e-02	1.4363200e-01
-1.3668380e-01	-1.0521880e-01	-3.9385400e-03	-4.6388000e-02
-7.7493780e-02	-2.4700800e-02	-5.2006200e-02	-2.6299360e-01
-2.5607520e-01	2.1704520e-01	5.6336000e-02	-6.3474400e-02
-1.0400400e-01	-1.7901000e-01	2.0326180e-01	-2.8708740e-01
1.0132000e-01	-1.6278080e-01	1.2441440e-01	3.2699820e-01
-4.8321600e-02	-3.6052800e-02	2.2539620e-01	-8.2764000e-03
3.1087258e-01	2.4090500e-01	-9.9590000e-02	1.2362460e-01
1.7440000e-03	-1.6117280e-01	7.4570000e-02	3.1281120e-02
-1.1758000e-02	-1.8464800e-02	-2.0872020e-01	-3.9510000e-03
-5.7714400e-01	-1.8090080e-01	-2.8288200e-01	-2.4662120e-01
-1.8806540e+00	4.4765400e-01	-2.9412700e-01	-1.7280000e-02
-3.1931600e-01	-1.9190500e-01	-1.1642000e-02	1.7475600e-01
1.3068840e-01	1.1943000e-01	-1.7219524e-01	1.9224000e-02
2.2620000e-01	-1.0821980e-01	1.3789060e-01	2.6989320e-01
-2.4364960e-01	-1.3650800e-01	-3.0984180e-01	-3.9546200e-02
-1.1410800e-01	-6.6744640e-02	1.6330620e-01	-4.0601000e-01
9.3793000e-02	-8.3026800e-02	9.0567600e-02	3.1595600e-01
1.6786620e-01	1.0099860e-01	3.5043600e-02	6.6221200e-02
-3.5907800e-02	-2.4589760e-01	2.6006800e-01	-8.0637000e-02
1.5359624e-01	-1.1078680e-01	-5.6956400e-02	2.2253080e-01
3.5808000e-02	-1.8873860e-01	-2.5032660e-01	3.6167400e-02
-2.2424700e-01	2.7863640e-01	2.2622600e-02	1.3753300e-01
-2.3369620e-01	2.8058040e-01	5.0818000e-02	-3.4805800e-02
1.7916600e-01	-7.5374000e-02	7.1228900e-02	1.7556000e-01

```
-5.8004120e-01 -2.0522500e-01 -1.3367960e-01 1.3656000e-02
-2.9052200e-02 1.3698600e-02 1.1746340e-01 -2.3288400e-02
2.7706200e-01 1.6106000e-01 -2.0183340e-01 5.7781800e-02
-2.0954400e-01 -1.4111260e-02 -3.1186860e-01 -2.9536360e-02
-1.7226500e-01 3.5709400e-01 2.9448200e-01 8.5600000e-05]
```

## Tf-Idf Weighted Word2Vec Vectorization

### 1. Vectorizing project\_essay

```
In [0]: # Initializing tfidf vectorizer
tfIdfEssayTempVectorizer = TfidfVectorizer();
# Vectorizing preprocessed essays using tfidf vectorizer initialized above
tfIdfEssayTempVectorizer.fit(preProcessedEssaysWithoutStopWords);
# Saving dictionary in which each word is key and it's idf is value
tfIdfEssayDictionary = dict(zip(tfIdfEssayTempVectorizer.get_feature_names(), list(tfIdfEssayTempVectorizer.idf_)));
# Creating set of all unique words used by tfidf vectorizer
tfIdfEssayWords = set(tfIdfEssayTempVectorizer.get_feature_names());
```

```
In [0]: # Creating list to save tf-idf weighted vectors of essays
tfIdfWeightedWord2VecEssaysVectors = [];
# Iterating over each essay
for essay in tqdm(preProcessedEssaysWithoutStopWords):
    # Sum of tf-idf values of all words in a particular essay
    cumulativeSumTfIdfWeightOfEssay = 0;
    # Tf-Idf weighted word2vec vector of a particular essay
    tfIdfWeightedWord2VecEssayVector = np.zeros(300);
    # Splitting essay into list of words
    splittedEssay = essay.split();
    # Iterating over each word
    for word in splittedEssay:
        # Checking if word is in glove words and set of words used by t
        fIdf essay vectorizer
        if (word in gloveWords) and (word in tfIdfEssayWords):
```

```

        # Tf-Idf value of particular word in essay
        tfIdfValueWord = tfIdfEssayDictionary[word] * (essay.count(
word) / len(splittedEssay));
        # Making tf-idf weighted word2vec
        tfIdfWeightedWord2VecEssayVector += tfIdfValueWord * gloveM
odel[word];
        # Summing tf-idf weight of word to cumulative sum
        cumulativeSumTfIdfWeightOfEssay += tfIdfValueWord;
    if cumulativeSumTfIdfWeightOfEssay != 0:
        # Taking average of sum of vectors with tf-idf cumulative sum
        tfIdfWeightedWord2VecEssayVector = tfIdfWeightedWord2VecEssayVe
ctor / cumulativeSumTfIdfWeightOfEssay;
        # Appending the above calculated tf-idf weighted vector of particul
ar essay to list of vectors of essays
        tfIdfWeightedWord2VecEssaysVectors.append(tfIdfWeightedWord2VecEssa
yVector);

```

In [0]:

```

print("Shape of Tf-Idf weighted Word2Vec vectorization matrix of projec
t essays: {}, {}".format(len(tfIdfWeightedWord2VecEssaysVectors), len(t
fIdfWeightedWord2VecEssaysVectors[0])));
equalsBorder(70);
print("Sample Essay: ");
equalsBorder(70);
print(preProcessedEssaysWithoutStopWords[0]);
equalsBorder(70);
print("Tf-Idf Weighted Word2Vec vector of sample essay: ");
equalsBorder(70);
print(tfIdfWeightedWord2VecEssaysVectors[0]);

```

Shape of Tf-Idf weighted Word2Vec vectorization matrix of project essay
s: 109248, 300

=====

Sample Essay:

=====

students english learners working english second third languages meltin
g pot refugees immigrants native born americans bringing gift language
school 24 languages represented english learner program students every
level mastery also 40 countries represented families within school stud

ent brings wealth knowledge experiences us open eyes new cultures belie  
fs respect limits language limits world ludwig wittgenstein english lea  
rner strong support system home begs resources many times parents learn  
ing read speak english along side children sometimes creates barriers p  
arents able help child learn phonetics letter recognition reading skill  
s providing dvd players students able continue mastery english language  
even no one home able assist families students within level 1 proficien  
cy status offered part program educational videos specially chosen engl  
ish learner teacher sent home regularly watch videos help child develop  
early reading skills parents not access dvd player opportunity check dv  
d player use year plan use videos educational dvd years come el student  
s nannan

=====

Tf-Idf Weighted Word2Vec vector of sample essay:

=====

```
[-5.37582850e-02 7.68689598e-02 7.85741822e-02 4.38958976e-02
 -8.56874440e-02 -1.20832331e-01 -2.68120986e+00 7.17018732e-02
 1.03799206e-04 -5.17255299e-03 -2.67529751e-02 7.40185988e-02
 1.36881934e-01 -8.62706493e-02 -6.35020145e-02 -8.44084597e-02
 1.27523921e-01 -1.77105602e-01 3.68451284e-02 -5.74471880e-02
 1.86477259e-01 9.28786009e-02 -9.73137896e-02 -1.15230456e-02
 4.41962185e-02 -9.32894883e-02 1.11912943e-01 -1.17540961e-01
 -1.22150893e-01 -9.14028838e-02 -1.73918944e-01 -4.54143189e-02
 -7.82036060e-02 3.05617633e-01 -8.71850266e-02 6.31466708e-03
 1.15683161e-01 1.71477594e-02 -5.52983597e-02 9.08989585e-02
 -3.89808292e-04 1.97696142e-01 -4.08078376e-01 -5.39990199e-02
 -1.20129600e-02 -1.12456389e-01 2.92046345e-02 1.37924729e-01
 2.83465620e-02 -2.26817169e-01 -2.29639267e-02 6.94257143e-03
 5.80535394e-02 2.86454339e-02 -7.51508216e-02 -6.21569354e-02
 -1.41805544e-01 2.78707358e-02 -1.63165999e-01 -1.29716251e-01
 -5.67625355e-02 -8.59507500e-02 3.54019902e-01 -4.96274469e-02
 -6.88414062e-02 1.58623510e-01 1.24798600e-01 4.29711440e-02
 7.82814323e-02 -1.73260116e-02 -1.23679491e-01 1.47617250e-01
 4.27083617e-02 -1.16531047e-01 -1.27122530e-01 -5.93638332e-03
 -1.99224414e-01 -8.66160391e-02 2.47701354e-01 1.61218205e-02
 3.56880345e-02 -3.71320273e-01 2.65501745e-01 -4.56454865e-02
 -7.85433814e-02 -5.99177835e-02 4.42212779e-02 -8.20739267e-02
 2.14031939e-01 2.42131497e-02 -1.34069697e-01 7.15871686e-03
 4.00667270e-02 -6.75881497e-02 -7.07967357e-02 -2.15984749e-02
```

-2.09734597e+00	1.02300477e-01	6.61169899e-02	5.70146517e-02
-1.91302495e-01	-1.38114014e-01	-1.10709961e-01	-1.66994098e-01
9.17800823e-02	1.35327093e-01	2.20333244e-02	-3.83844831e-02
2.57206511e-02	-5.54503565e-02	-3.41973653e-03	1.99777588e-02
4.85050396e-02	2.13190534e-01	4.64281665e-02	6.51171751e-02
-5.80015838e-02	1.19900386e-01	1.18803830e-01	7.05550873e-02
-1.87330886e-01	1.41219129e-01	1.33569574e-01	1.00530000e-01
4.14498415e-02	1.39860952e-01	-7.95709830e-02	9.70242332e-02
1.07442882e-01	9.00794808e-02	7.47745032e-02	4.18772282e-02
-7.10347826e-03	-7.62379756e-03	-7.31715828e-02	-1.16370646e-01
2.82271708e-01	-5.30885621e-02	4.51472249e-02	2.61376253e-01
1.29080066e-02	3.96843846e-02	1.04430681e-01	-1.30495811e-01
-1.17999239e-01	-1.02810089e-01	-6.52713784e-02	-1.81350799e-01
1.55415740e-01	-4.43517889e-02	-8.34350788e-02	-1.31445407e-01
-8.87524029e-02	-1.15321245e-02	8.67587067e-03	3.55646447e-02
-4.32365925e-02	2.44285859e-03	2.73165854e-02	-1.91651165e-01
6.70942750e-03	1.45533103e-02	-5.95191056e-02	-9.78336553e-02
-4.61200683e-02	1.04017495e-02	-1.68129330e-01	-5.53455289e-02
-1.95353920e-02	-3.24088827e-03	9.94121739e-02	-2.20584067e-02
1.36190091e-02	-3.13014669e-01	4.46748268e-02	6.11251996e-02
-5.59088914e-02	8.07071841e-02	-7.80920682e-02	1.05535003e-02
-8.49705076e-02	1.87800458e-01	-5.53305425e-02	-4.05296946e-02
-1.68105655e-02	-9.64697267e-02	-1.00114054e-01	-1.25303984e-01
-6.77861115e-02	1.38106300e-02	4.97948787e-02	-1.04414463e-01
3.12147536e-03	-2.46650333e-02	1.56250756e-02	-3.41987984e-02
2.90197738e-02	-1.30795750e-01	1.71425098e-01	-1.33199913e-01
-4.35452619e-02	-1.52841321e-01	3.37717104e-02	2.11400042e-01
-1.08493100e-01	6.64905827e-02	4.45687503e-02	-3.38898797e-03
1.47302984e-01	3.10931848e-02	6.94873935e-03	-3.79090162e-02
3.97055902e-02	-3.12563998e-02	2.99815273e-02	-9.30892230e-03
-3.37192802e-02	-7.79667288e-02	4.20509297e-02	4.33535394e-02
-2.38238094e-01	-4.11188300e-02	-1.93930088e-01	1.15012485e-01
-2.14605373e+00	1.36975648e-01	-1.79026305e-01	-1.42630498e-01
-1.37558424e-01	-1.55433436e-01	-6.96701214e-02	1.05328488e-01
3.43486342e-02	-2.37676310e-03	-6.80980842e-02	-1.92470331e-01
1.54727348e-01	-7.47455695e-02	-1.58054203e-02	3.33369549e-02
-1.70510752e-01	-5.74331307e-02	-2.38994456e-01	5.64188931e-02
-8.55051184e-02	-5.52984572e-02	-5.00408589e-02	-6.81572658e-02
5.15848477e-03	-3.58487773e-02	7.00056842e-02	1.33127170e-01

```
5.57938159e-02 1.03106840e-01 4.18598320e-02 -2.78162076e-03
8.83131944e-02 -1.31482831e-01 1.34875022e-01 -8.31772344e-02
1.62319378e-01 9.25839856e-02 -7.07548194e-02 1.74355644e-01
1.53106818e-02 -1.74504449e-01 -5.39158255e-02 -1.16968555e-02
-1.37824311e-01 1.07713713e-01 4.48548015e-02 1.07272158e-01
-1.59084558e-01 1.94342786e-01 -4.73514319e-02 -4.87250503e-02
2.82023483e-02 -4.18474756e-02 8.04397595e-02 -3.34005484e-02
-1.00808502e-01 -1.15380334e-01 7.05894205e-02 2.92052920e-02
-5.72604859e-02 -7.39274088e-03 1.44106517e-02 -2.64282237e-02
2.31512689e-01 1.50161666e-01 -5.21462274e-02 -1.00796916e-02
-4.47392305e-02 4.83958092e-02 -2.21927272e-01 -9.69846899e-02
-5.91211767e-03 2.52508756e-01 1.08677704e-01 5.05047869e-02]
```

## 2. Vectorizing project\_title

```
In [0]: # Initializing tfidf vectorizer
tfIdfTitleTempVectorizer = TfidfVectorizer();
# Vectorizing preprocessed titles using tfidf vectorizer initialized above
tfIdfTitleTempVectorizer.fit(preProcessedProjectTitlesWithoutStopWords);
# Saving dictionary in which each word is key and it's idf is value
tfIdfTitleDictionary = dict(zip(tfIdfTitleTempVectorizer.get_feature_names(),
                                list(tfIdfTitleTempVectorizer.idf_)));
# Creating set of all unique words used by tfidf vectorizer
tfIdfTitleWords = set(tfIdfTitleTempVectorizer.get_feature_names());
```

```
In [0]: # Creating list to save tf-idf weighted vectors of project titles
tfIdfWeightedWord2VecTitlesVectors = [];
# Iterating over each title
for title in tqdm(preProcessedProjectTitlesWithoutStopWords):
    # Sum of tf-idf values of all words in a particular project title
    cumulativeSumTfIdfWeightOfTitle = 0;
    # Tf-Idf weighted word2vec vector of a particular project title
    tfIdfWeightedWord2VecTitleVector = np.zeros(300);
    # Splitting title into list of words
    splittedTitle = title.split();
```

```

# Iterating over each word
for word in splittedTitle:
    # Checking if word is in glove words and set of words used by t
    fIdf title vectorizer
        if (word in gloveWords) and (word in tfIdfTitleWords):
            # Tf-Idf value of particular word in title
            tfIdfValueWord = tfIdfTitleDictionary[word] * (title.count(
word) / len(splittedTitle));
            # Making tf-idf weighted word2vec
            tfIdfWeightedWord2VecTitleVector += tfIdfValueWord * gloveM
odel[word];
            # Summing tf-idf weight of word to cumulative sum
            cumulativeSumTfIdfWeightOfTitle += tfIdfValueWord;
        if cumulativeSumTfIdfWeightOfTitle != 0:
            # Taking average of sum of vectors with tf-idf cumulative sum
            tfIdfWeightedWord2VecTitleVector = tfIdfWeightedWord2VecTitleVe
ctor / cumulativeSumTfIdfWeightOfTitle;
            # Appending the above calculated tf-idf weighted vector of particul
ar title to list of vectors of project titles
            tfIdfWeightedWord2VecTitlesVectors.append(tfIdfWeightedWord2VecTitl
eVector);

```

In [0]:

```

print("Shape of Tf-Idf weighted Word2Vec vectorization matrix of projec
t titles: {}, {}".format(len(tfIdfWeightedWord2VecTitlesVectors), len(t
fIdfWeightedWord2VecTitlesVectors[0])));
equalsBorder(70);
print("Sample Title: ");
equalsBorder(70);
print(preProcessedProjectTitlesWithoutStopWords[0]);
equalsBorder(70);
print("Tf-Idf Weighted Word2Vec vector of sample title: ");
equalsBorder(70);
print(tfIdfWeightedWord2VecTitlesVectors[0]);

```

Shape of Tf-Idf weighted Word2Vec vectorization matrix of project title  
s: 109248, 300  
=====

Sample Title:  
=====

educational support english learners home

Tf-Idf Weighted Word2Vec vector of sample title:

```
[ -3.23904891e-02  5.58064810e-02  1.32666911e-01  3.84227573e-02
 -6.71984492e-02 -4.30940397e-01 -2.84607947e+00 -2.45905055e-01
  1.96794858e-01  3.19604663e-01 -6.12568872e-02  1.59218099e-01
  1.25129027e-01 -1.67580327e-01 -2.82644062e-01 -2.47555536e-01
  2.18304104e-01 -1.57431101e-01  7.66481545e-02 -1.61436633e-01
  2.38451267e-01  2.86712258e-01  2.70730890e-02 -9.74962294e-02
  1.67511144e-01  7.18131102e-02  1.82846112e-01 -1.96778087e-01
 -8.19948978e-02 -2.25877630e-01 -5.54573752e-01 -1.28462870e-01
  1.61012606e-01  2.94412658e-01 -1.63196910e-01 -1.23217523e-02
  1.37466355e-01  4.45437696e-02  4.65691769e-02 -1.17867965e-01
 -2.41502151e-03  2.24350668e-01 -2.51274676e-01  8.29431360e-02
 -1.65996673e-02 -2.47747576e-01  1.45110611e-03  2.37117949e-01
  9.71345150e-02 -5.13516477e-01 -1.40296688e-01  1.42775548e-01
  2.89949805e-01  6.49771690e-02 -3.41581088e-02 -1.58076306e-01
 -1.07731741e-01  7.59015357e-02 -1.21511682e-01 -1.16519972e-01
 -2.27321940e-01 -1.63525257e-01  1.80860125e-01 -4.17314689e-02
  4.60171896e-02  1.00024674e-01  1.54588362e-01  8.25394911e-02
  7.45768118e-02 -1.80240543e-02 -1.22956246e-01 -4.97450371e-03
 -8.06577406e-02 -5.00614538e-02 -2.15836210e-01 -5.89271531e-02
 -3.26363335e-01 -1.32706775e-01  1.61236199e-01 -1.25038790e-01
  1.96493846e-01 -4.95095193e-01  2.34765396e-01 -4.44646606e-02
 -2.04266125e-01 -3.21415735e-02  8.48111983e-02 -7.27603472e-02
  2.79183660e-01  1.18968262e-01 -7.43300594e-02  6.34587771e-02
  1.99863053e-01 -2.13382053e-01 -1.01221319e-01 -2.49884070e-01
 -2.92249478e+00  1.60273141e-01  7.74579728e-02  1.85323805e-01
 -1.33255909e-01 -2.00013519e-01 -1.31974722e-01 -2.62288530e-01
  3.54852941e-01  1.18537924e-01  1.62207829e-01  1.24436802e-01
  1.98867481e-01 -4.87526944e-03  3.00886908e-02  2.09330567e-01
  1.17189984e-01  3.94887340e-01  2.52941492e-02 -5.13348554e-02
 -2.91140828e-01  4.06939567e-01 -1.70319175e-01  1.17651155e-01
 -1.66813086e-01  1.53049826e-01  1.41255472e-01 -8.10785736e-02
  9.57549943e-02  2.73610111e-01  5.85622995e-02  7.91410001e-02
  1.47619459e-01  9.75521835e-02  6.74487028e-02  1.53125504e-01
  2.02791106e-02 -5.59403852e-02 -1.02109913e-01 -1.22913427e-01
  1.99873969e-01 -3.21872719e-01  1.38343165e-01  2.17196179e-01
```

4.95201760e-03	8.52128333e-02	-1.45880901e-01	-2.10862397e-01
1.20343357e-01	2.15598061e-01	-1.14038072e-02	-1.72172799e-01
3.24157324e-01	-3.82818101e-01	-1.87580283e-01	-2.00827204e-01
-1.41863370e-01	9.63016678e-02	-2.01659119e-01	6.74342164e-02
7.12185747e-02	-1.04314039e-01	-9.08169483e-02	-1.63495605e-01
9.68230169e-02	-5.01176209e-02	-8.34015616e-02	-1.88998660e-01
-2.84065057e-01	1.16975197e-01	-2.80836800e-01	-9.33191327e-02
3.79583269e-02	-1.22755412e-01	2.30408258e-01	-1.31968890e-01
9.72824714e-03	-3.44272546e-01	-2.09522211e-03	2.45944018e-02
2.94077607e-02	2.67568157e-01	-2.69460269e-02	1.25412311e-01
-3.47031083e-01	2.09328241e-01	1.25385338e-02	1.55654760e-01
-1.41368915e-01	-1.01749781e-01	-4.77312036e-04	-4.82325465e-02
-7.15727478e-02	-3.63658602e-02	-4.33504397e-02	-2.71410315e-01
-2.40079853e-01	2.01171435e-01	6.39005674e-02	-4.86787485e-02
-1.48623863e-01	-1.72130906e-01	1.97761227e-01	-3.13043504e-01
1.07772898e-01	-1.54518908e-01	1.31855435e-01	3.39703669e-01
-4.51652340e-02	-4.05998340e-02	2.03610454e-01	8.84982054e-03
3.05974297e-01	2.54736700e-01	-1.06925907e-01	1.27066655e-01
-1.88835779e-02	-1.56632041e-01	8.45142200e-02	5.70681135e-02
1.01119358e-02	-6.62387316e-03	-2.18552410e-01	1.20985419e-02
-5.54006219e-01	-1.72367117e-01	-2.90325016e-01	-2.34816399e-01
-1.94243114e+00	4.36715446e-01	-2.80713863e-01	-6.33991309e-03
-2.90035778e-01	-1.98732349e-01	2.96737137e-02	1.50873684e-01
1.16943997e-01	1.39741722e-01	-1.82238609e-01	4.09714520e-02
2.37176600e-01	-1.24515116e-01	1.41648743e-01	2.64206287e-01
-2.40551078e-01	-1.40415333e-01	-2.92432371e-01	-3.03761027e-02
-9.90320454e-02	-8.43648662e-02	1.81116706e-01	-4.05719699e-01
1.22898740e-01	-8.80109292e-02	1.09543672e-01	2.96110858e-01
1.85027885e-01	9.14976115e-02	9.63416424e-03	5.50340717e-02
-2.59328007e-02	-2.43942768e-01	2.54260096e-01	-1.03280950e-01
1.56799018e-01	-9.58635926e-02	-4.31948365e-02	2.01228907e-01
5.20033765e-02	-2.08030399e-01	-2.49149283e-01	3.11752465e-02
-2.39410711e-01	2.54421815e-01	3.50420005e-02	1.31625993e-01
-2.19027956e-01	2.75093693e-01	4.31276229e-02	-6.89266192e-02
1.80694153e-01	-9.77254221e-02	6.52789959e-02	1.81468103e-01
-5.79288980e-01	-1.91501478e-01	-1.43298895e-01	1.56769073e-02
-2.28584041e-02	-7.96762354e-03	1.38764109e-01	-2.67804890e-02
3.02808634e-01	1.63688874e-01	-1.98263925e-01	8.94007093e-02

```
-2.01132765e-01 8.29230669e-03 -3.17426319e-01 -4.07929287e-02
-1.63872993e-01 3.69860278e-01 2.90009047e-01 4.56005599e-02]
```

## Vectorizing numerical features

### 1. Vectorizing price

```
In [0]: # Standardizing the price data using StandardScaler(Uses mean and std f
or standardization)
priceScaler = StandardScaler();
priceScaler.fit(projectsData['price'].values.reshape(-1, 1));
priceStandardized = priceScaler.transform(projectsData['price'].values.
reshape(-1, 1));
```

```
In [0]: print("Shape of standardized matrix of prices: ", priceStandardized.sha
pe);
equalsBorder(70);
print("Sample original prices: ");
equalsBorder(70);
print(projectsData['price'].values[0:5]);
print("Sample standardized prices: ");
equalsBorder(70);
print(priceStandardized[0:5]);
```

```
Shape of standardized matrix of prices: (109245, 1)
```

```
=====
Sample original prices:
```

```
=====
[154.6 299. 516.85 232.9 67.98]
```

```
=====
Sample standardized prices:
```

```
=====
[[ -0.39052147]
 [ 0.00240752]
 [ 0.5952024 ]
 [-0.17745817]
 [-0.62622444]]
```

## 2. Vectorizing quantity

```
In [0]: # Standardizing the quantity data using StandardScaler(Uses mean and std for standardization)
quantityScaler = StandardScaler();
quantityScaler.fit(projectsData['quantity'].values.reshape(-1, 1));
quantityStandardized = quantityScaler.transform(projectsData['quantity'].values.reshape(-1, 1));
```

```
In [0]: print("Shape of standardized matrix of quantities: ", quantityStandardized.shape);
equalsBorder(70);
print("Sample original quantities: ");
equalsBorder(70);
print(projectsData['quantity'].values[0:5]);
print("Sample standardized quantities: ");
equalsBorder(70);
print(quantityStandardized[0:5]);
```

Shape of standardized matrix of quantities: (109245, 1)

=====

Sample original quantities:

=====

[23 1 22 4 4]

Sample standardized quantities:

=====

```
[[ 0.23045805]
 [-0.6097785 ]
 [ 0.19226548]
 [-0.49520079]
 [-0.49520079]]
```

## 3. Vectorizing teacher\_number\_of\_previously\_posted\_projects

```
In [0]: # Standardizing the teacher_number_of_previously_posted_projects data using StandardScaler(Uses mean and std for standardization)
```

```
previouslyPostedScaler = StandardScaler();
previouslyPostedScaler.fit(projectsData['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1));
previouslyPostedStandardized = previouslyPostedScaler.transform(projectsData['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1));
```

```
In [0]: print("Shape of standardized matrix of teacher_number_of_previously_posted_projects: ", previouslyPostedStandardized.shape);
equalsBorder(70);
print("Sample original quantities:");
equalsBorder(70);
print(projectsData['teacher_number_of_previously_posted_projects'].values[0:5]);
print("Sample standardized teacher_number_of_previously_posted_projects:");
equalsBorder(70);
print(previouslyPostedStandardized[0:5]);
```

```
Shape of standardized matrix of teacher_number_of_previously_posted_projects: (109245, 1)
=====
Sample original quantities:
=====
[0 7 1 4 1]
Sample standardized teacher_number_of_previously_posted_projects:
=====
[[ -0.40153083
  [-0.14952695]
  [-0.36553028]
  [-0.25752861]
  [-0.36553028]]]
```

## Taking 6k points(to avoid memory errors)

```
In [0]: numberOfPoints = 6000;
# Categorical data
```

```
categoriesVectorsSub = categoriesVectors[0:numberOfPoints];
subCategoriesVectorsSub = subCategoriesVectors[0:numberOfPoints];
teacherPrefixVectorsSub = teacherPrefixVectors[0:numberOfPoints];
schoolStateVectorsSub = schoolStateVectors[0:numberOfPoints];
projectGradeVectorsSub = projectGradeVectors[0:numberOfPoints];

# Text data
bowEssayModelSub = bowEssayModel[0:numberOfPoints];
bowTitleModelSub = bowTitleModel[0:numberOfPoints];
tfIdfEssayModelSub = tfIdfEssayModel[0:numberOfPoints];
tfIdfTitleModelSub = tfIdfTitleModel[0:numberOfPoints];
word2VecEssaysVectorsSub = word2VecEssaysVectors[0:numberOfPoints];
word2VecTitlesVectorsSub = word2VecTitlesVectors[0:numberOfPoints];
tfIdfWeightedWord2VecEssaysVectorsSub = tfIdfWeightedWord2VecEssaysVectors[0:numberOfPoints];
tfIdfWeightedWord2VecTitlesVectorsSub = tfIdfWeightedWord2VecTitlesVectors[0:numberOfPoints];

# Numerical data
priceStandardizedSub = priceStandardized[0:numberOfPoints];
quantityStandardizedSub = quantityStandardized[0:numberOfPoints];
previouslyPostedStandardizedSub = previouslyPostedStandardized[0:numberOfPoints];
```

```
In [0]: classesDataSub = projectsData['project_is_approved'][0:numberOfPoints].values
```

```
In [0]: classesDataSub.shape
```

```
Out[0]: (6000,)
```

## Data Visualization using T-SNE

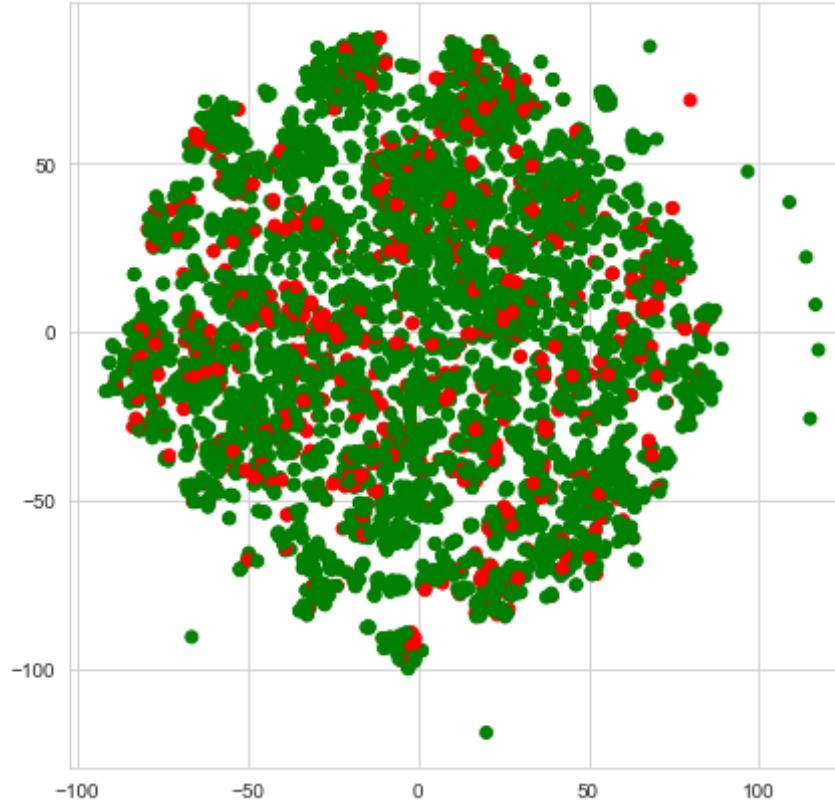
Classification using data merged with bag of words vectorized title and all considered categorical, numerical features

```
In [0]: bowTitleAndOthers = hstack((bowTitleModelSub, categoriesVectorsSub, subCategoriesVectorsSub, teacherPrefixVectorsSub, schoolStateVectorsSub, projectGradeVectorsSub, priceStandardizedSub, previouslyPostedStandardizedSub));
bowTitleAndOthers.shape
```

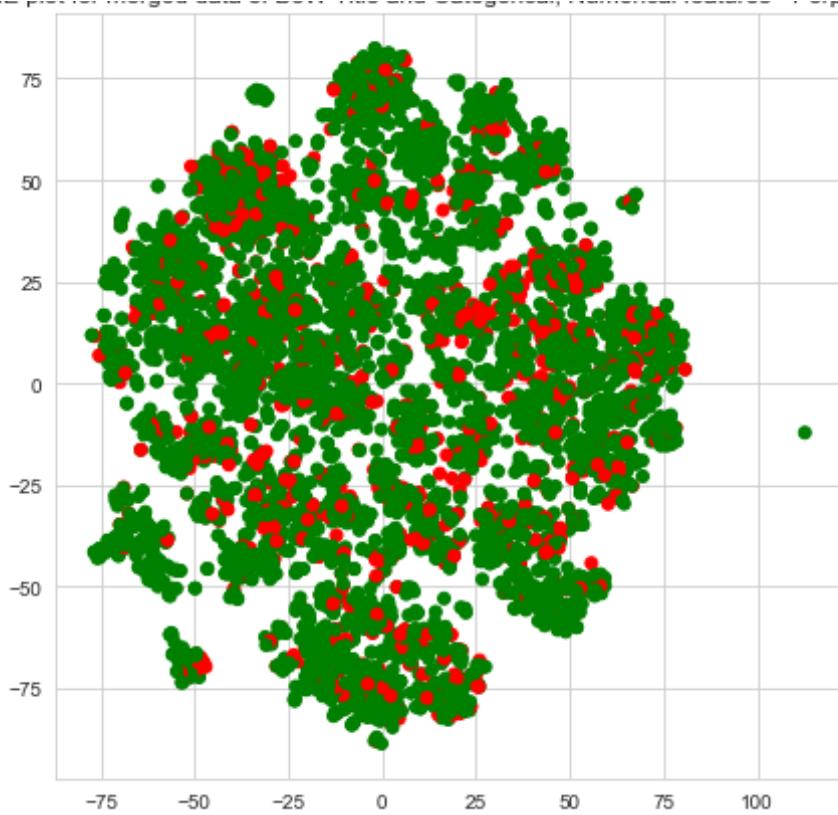
```
Out[0]: (6000, 1875)
```

```
In [0]: perplexityValues = [5, 10, 30, 50, 80, 100]
for perplexityValue in perplexityValues:
    tsne = TSNE(n_components = 2, perplexity = perplexityValue, learning_rate = 200);
    bowTitleAndOthersEmbedded = tsne.fit_transform(bowTitleAndOthers.toarray());
    bowTitleAndOthersTsneData = np.hstack((bowTitleAndOthersEmbedded, classesDataSub.reshape(-1, 1)));
    bowTitleAndOthersTsneDataFrame = pd.DataFrame(bowTitleAndOthersTsneData, columns = ['Dimension1', 'Dimension2', 'Class']);
    colors = {0.0:'red', 1.0:'green'}
    plt.title("TSNE plot for merged data of BoW Title and Categorical, Numerical features - Perplexity({})".format(perplexityValue));
    plt.scatter(bowTitleAndOthersTsneDataFrame['Dimension1'], bowTitleAndOthersTsneDataFrame['Dimension2'], c = bowTitleAndOthersTsneDataFrame['Class'].apply(lambda x: colors[x]));
    plt.show();
```

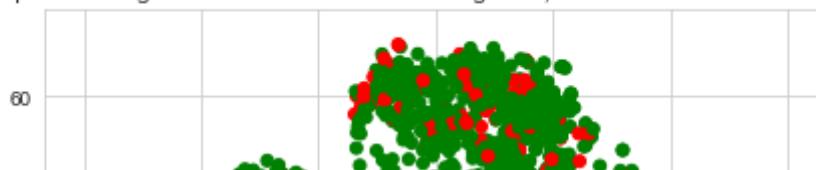
TSNE plot for merged data of BoW Title and Categorical, Numerical features - Perplexity(5)

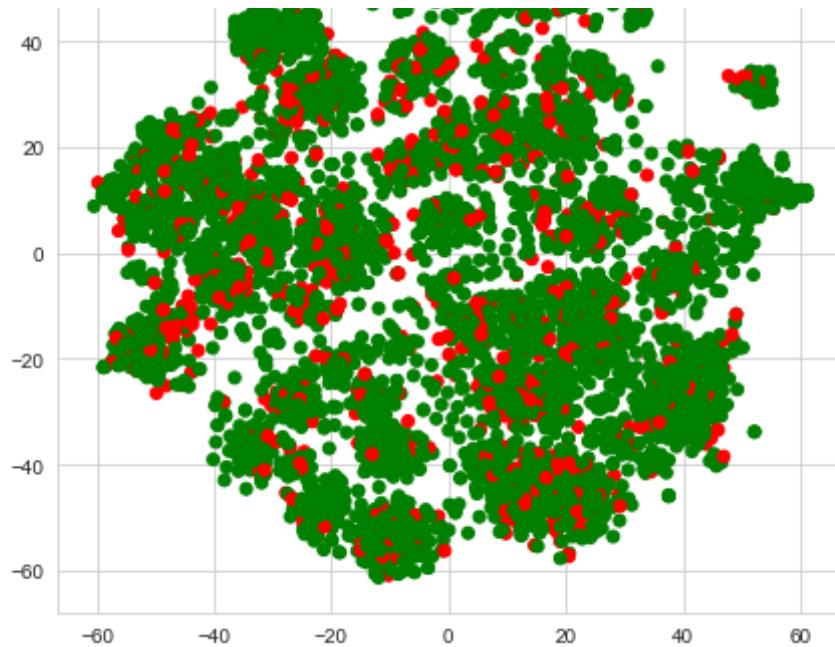


TSNE plot for merged data of BoW Title and Categorical, Numerical features - Perplexity(10)

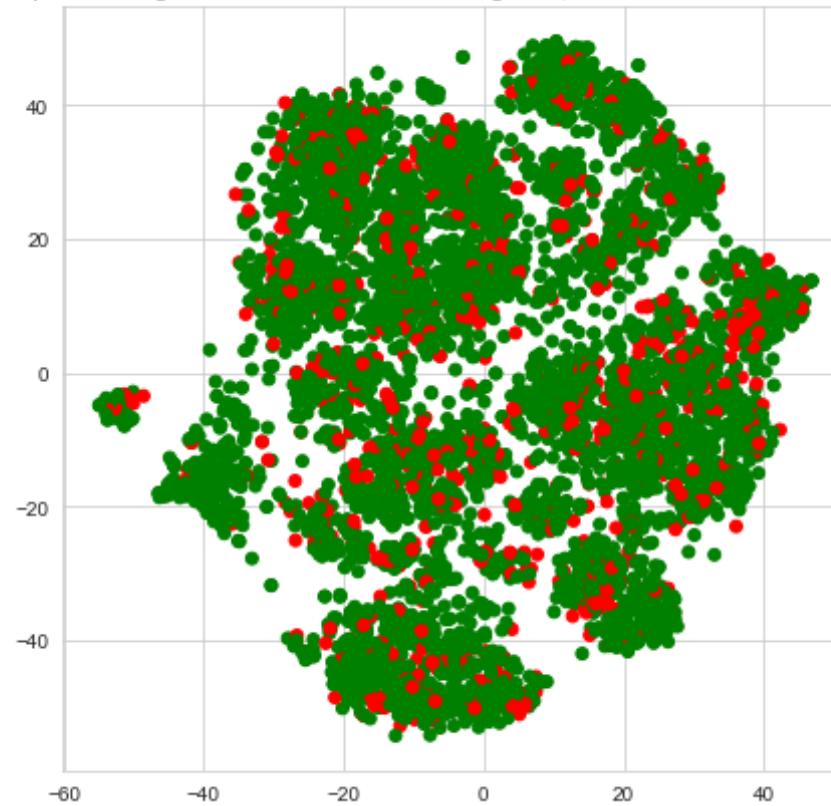


TSNE plot for merged data of BoW Title and Categorical, Numerical features - Perplexity(30)

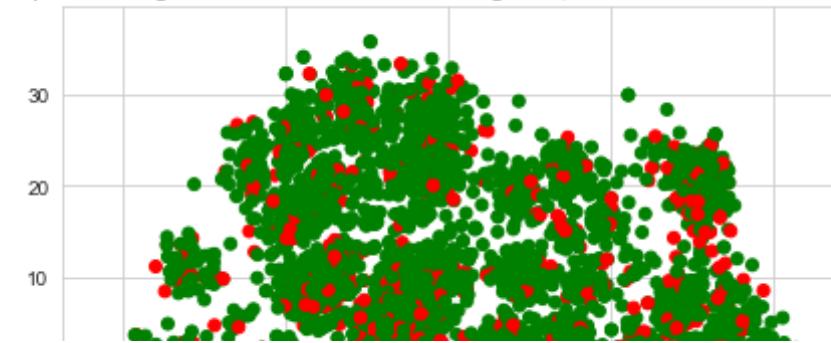


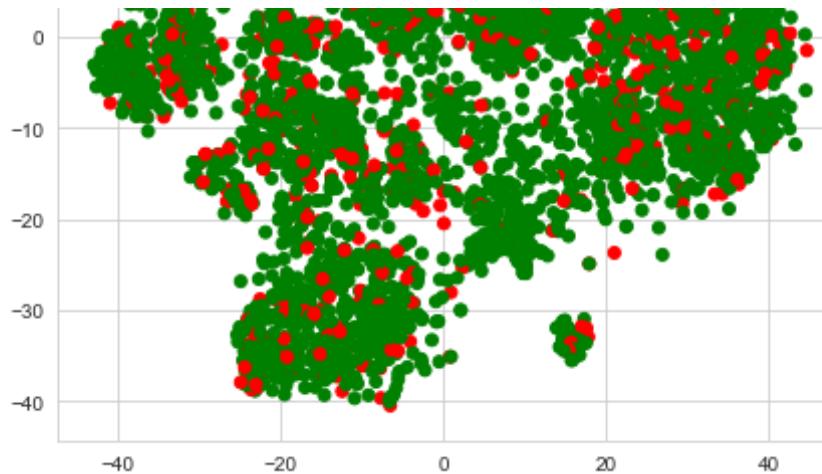


TSNE plot for merged data of BoW Title and Categorical, Numerical features - Perplexity(50)

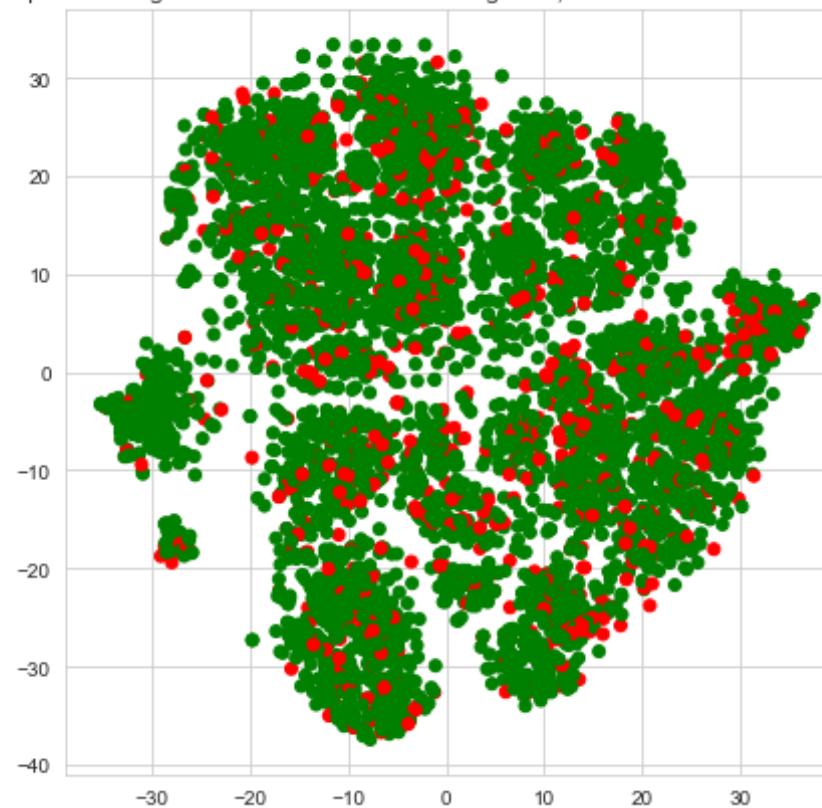


TSNE plot for merged data of BoW Title and Categorical, Numerical features - Perplexity(80)





TSNE plot for merged data of BoW Title and Categorical, Numerical features - Perplexity(100)



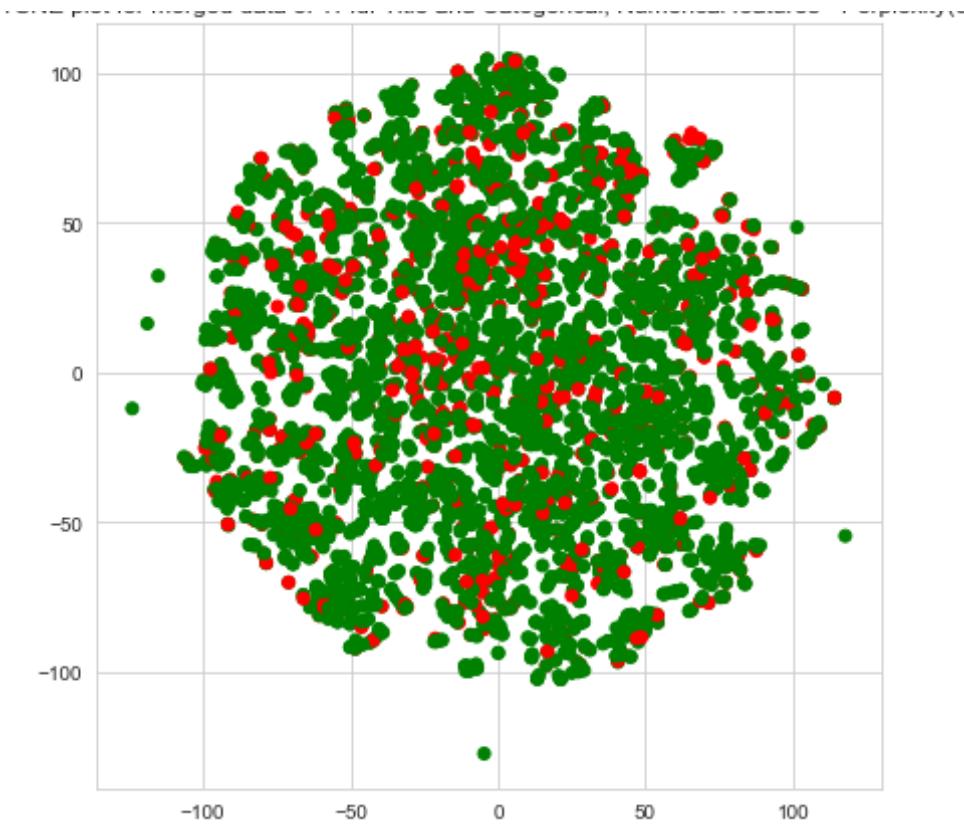
## Classification using data merged with Tf-Idf vectorized title and all considered categorical, numerical features

```
In [0]: tfIdfTitleAndOthers = hstack((tfIdfTitleModelSub, categoriesVectorsSub,
    subCategoriesVectorsSub, teacherPrefixVectorsSub, schoolStateVectorsSu
b, projectGradeVectorsSub, priceStandardizedSub, previouslyPostedStanda
rdizedSub));
tfIdfTitleAndOthers.shape
```

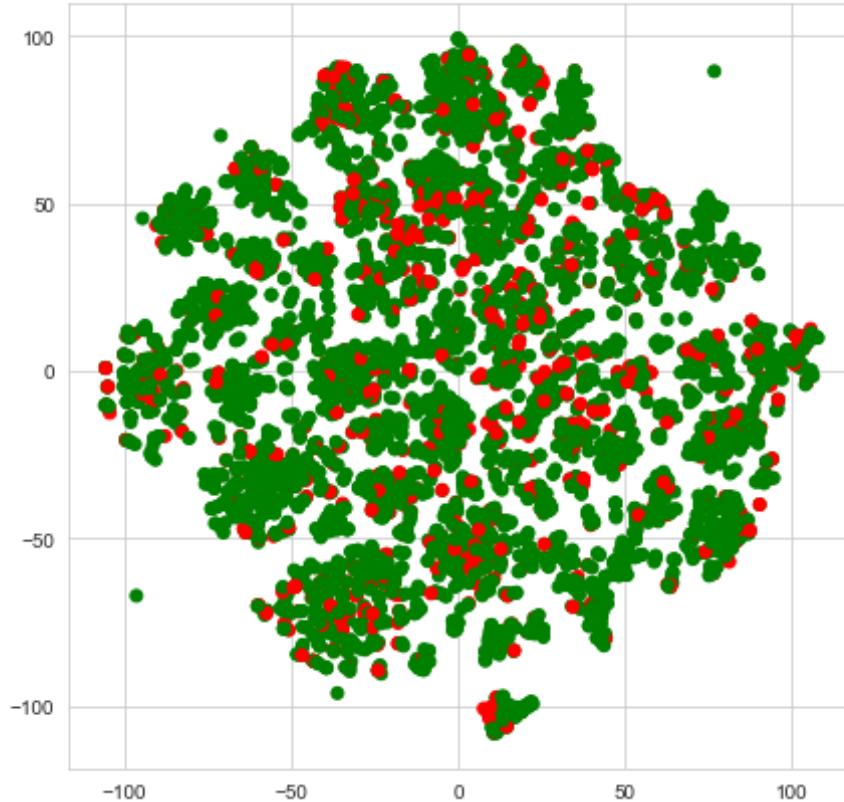
```
Out[0]: (6000, 1875)
```

```
In [0]: perplexityValues = [5, 10, 30, 50, 80, 100]
for perplexityValue in perplexityValues:
    tsne = TSNE(n_components = 2, perplexity = perplexityValue, learnin
g_rate = 200);
    tfIdfTitleAndOthersEmbedded = tsne.fit_transform(tfIdfTitleAndOther
s.toarray());
    tfIdfTitleAndOthersTsneData = np.hstack((tfIdfTitleAndOthersEmbedde
d, classesDataSub.reshape(-1, 1)));
    tfIdfTitleAndOthersTsneDataFrame = pd.DataFrame(tfIdfTitleAndOthers
TsneData, columns = ['Dimension1', 'Dimension2', 'Class']);
    colors = {0.0:'red', 1.0:'green'}
    plt.title("TSNE plot for merged data of Tf-Idf Title and Categorica
l, Numerical features - Perplexity({})".format(perplexityValue));
    plt.scatter(tfIdfTitleAndOthersTsneDataFrame['Dimension1'], tfIdfTi
tleAndOthersTsneDataFrame['Dimension2'], c = tfIdfTitleAndOthersTsneDat
aFrame['Class'].apply(lambda x: colors[x]));
    plt.show();
```

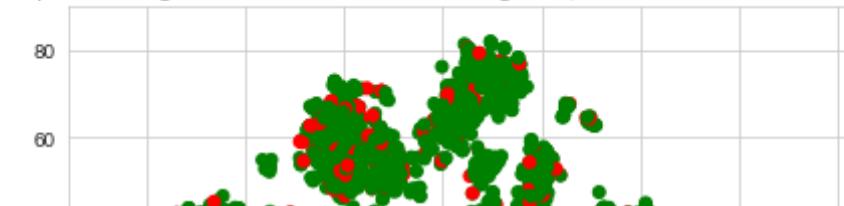
TSNE plot for meraed data of Tf-Idf Title and Categorical. Numerical features - Perplexity(5)

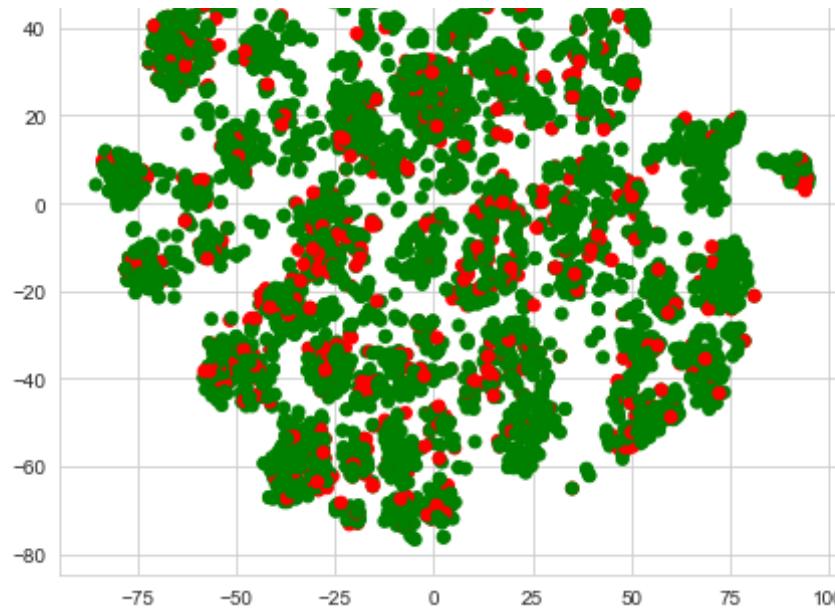


TSNE plot for merged data of Tf-Idf Title and Categorical, Numerical features - Perplexity(10)

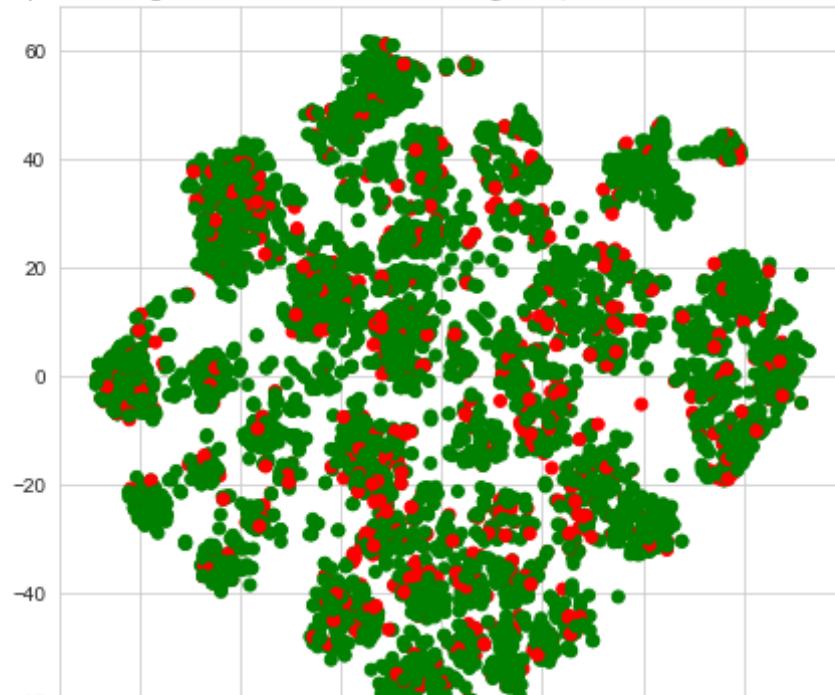


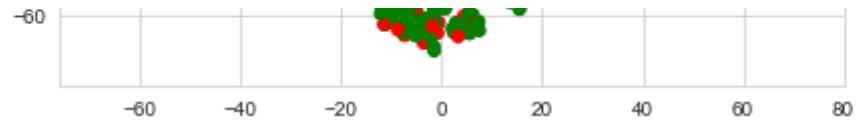
TSNE plot for merged data of Tf-Idf Title and Categorical, Numerical features - Perplexity(30)



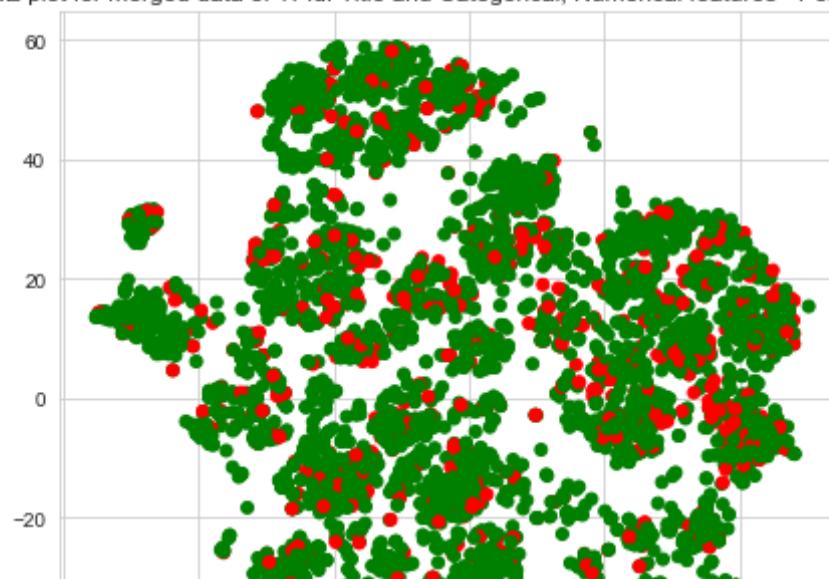


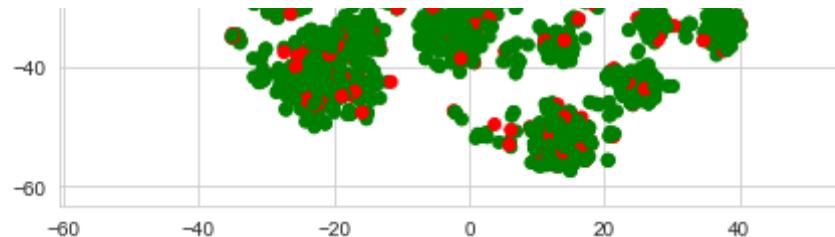
TSNE plot for merged data of Tf-Idf Title and Categorical, Numerical features - Perplexity(50)



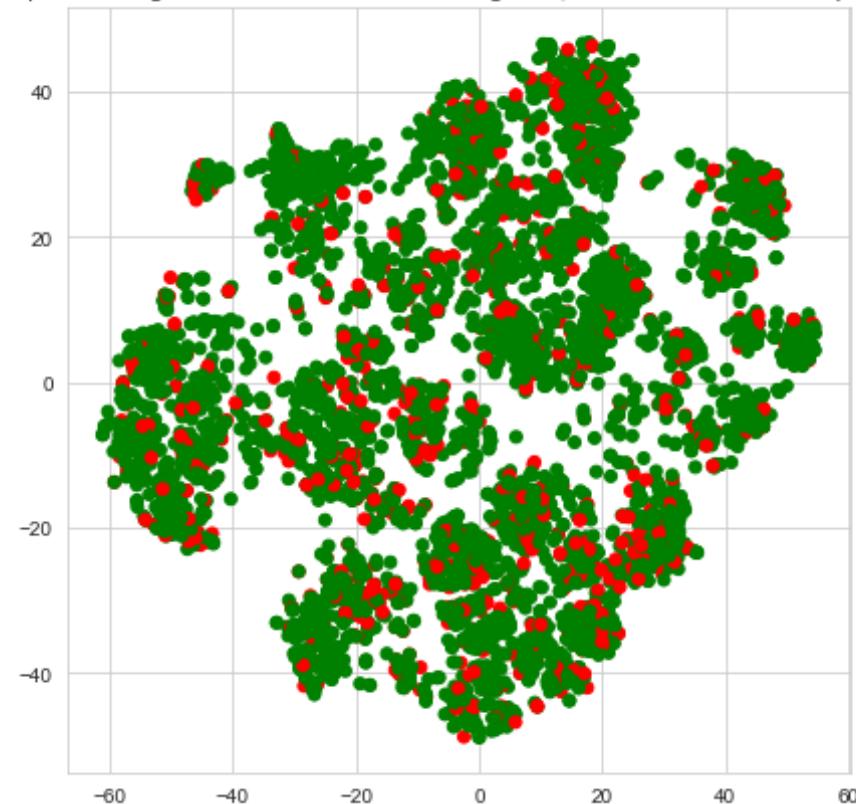


TSNE plot for merged data of Tf-Idf Title and Categorical, Numerical features - Perplexity(80)





TSNE plot for merged data of Tf-Idf Title and Categorical, Numerical features - Perplexity(100)



Classification using data merged with Average Word2Vec vectorized title and all considered categorical, numerical features

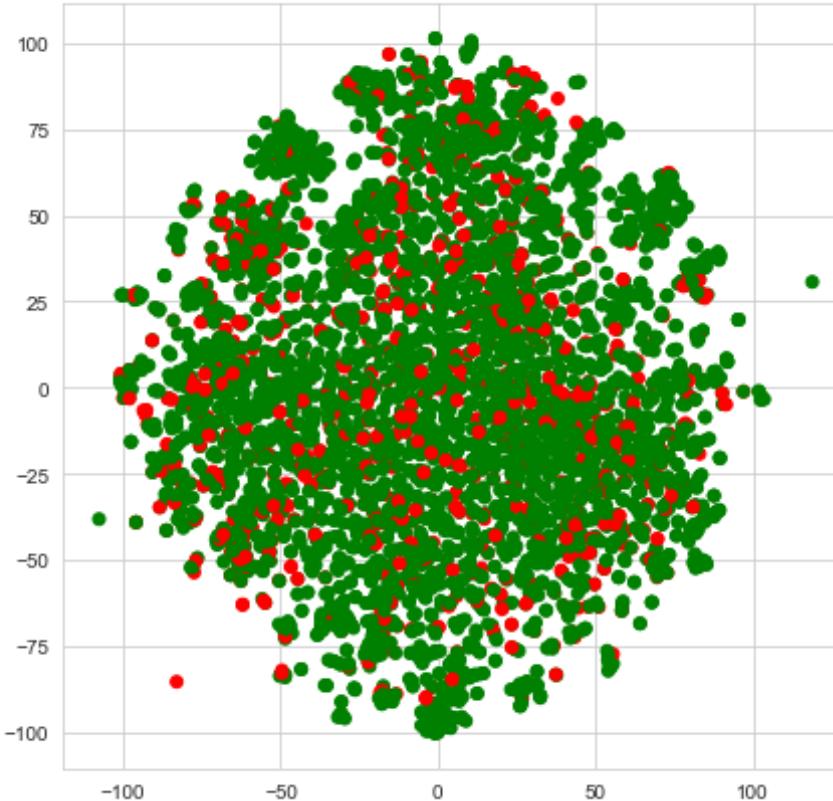
In [0]: word2VecTitleAndOthers = hstack((word2VecTitlesVectorsSub, categoriesVe

```
ctorsSub, subCategoriesVectorsSub, teacherPrefixVectorsSub, schoolState  
VectorsSub, projectGradeVectorsSub, priceStandardizedSub, previouslyPos  
tedStandardizedSub));  
word2VecTitleAndOthers.shape
```

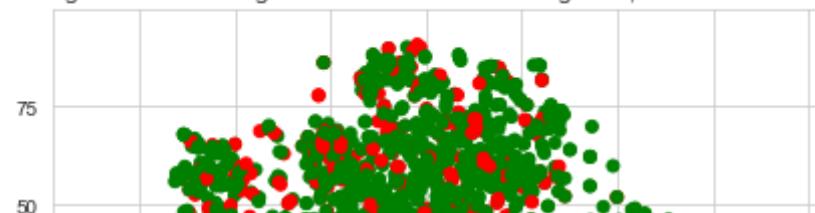
Out[0]: (6000, 401)

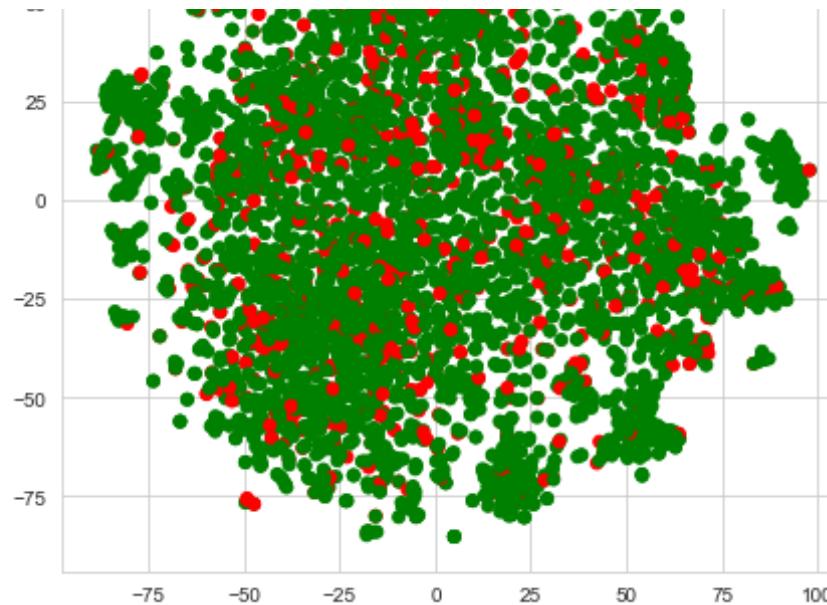
```
In [0]: perplexityValues = [5, 10, 30, 50, 80, 100]  
for perplexityValue in perplexityValues:  
    tsne = TSNE(n_components = 2, perplexity = perplexityValue, learnin  
g_rate = 200);  
    word2VecTitleAndOthersEmbedded = tsne.fit_transform(word2VecTitleAn  
dOthers.toarray());  
    word2VecTitleAndOthersTsneData = np.hstack((word2VecTitleAndOthersE  
mbedded, classesDataSub.reshape(-1, 1)));  
    word2VecTitleAndOthersTsneDataFrame = pd.DataFrame(word2VecTitleAnd  
OthersTsneData, columns = ['Dimension1', 'Dimension2', 'Class']);  
    colors = {0.0:'red', 1.0:'green'}  
    plt.title("TSNE plot for merged data of Average Word2Vec Title and  
Categorical, Numerical features - Perplexity({})".format(perplexityVal  
ue));  
    plt.scatter(word2VecTitleAndOthersTsneDataFrame['Dimension1'], word  
2VecTitleAndOthersTsneDataFrame['Dimension2'], c = word2VecTitleAndOthe  
rsTsneDataFrame['Class'].apply(lambda x: colors[x]));  
    plt.show();
```

TSNE plot for merged data of Average Word2Vec Title and Categorical, Numerical features - Perplexity(5)

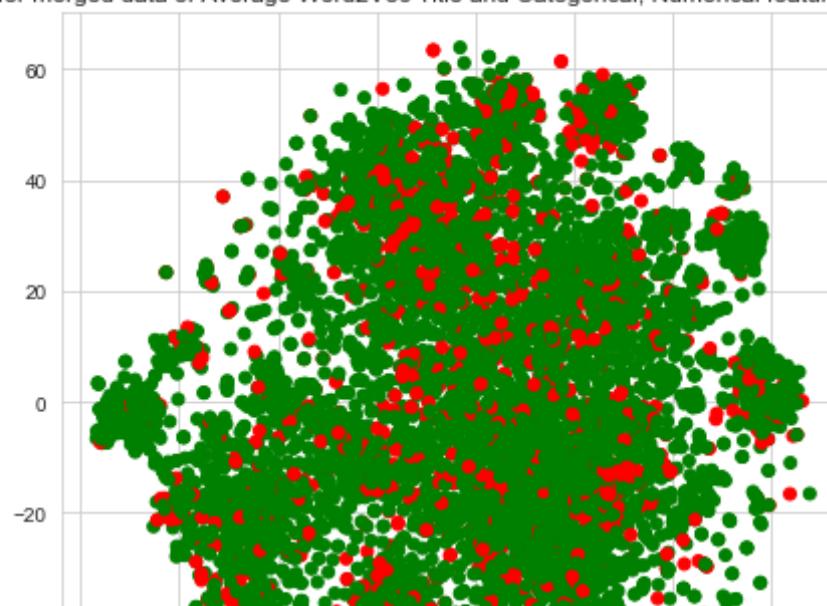


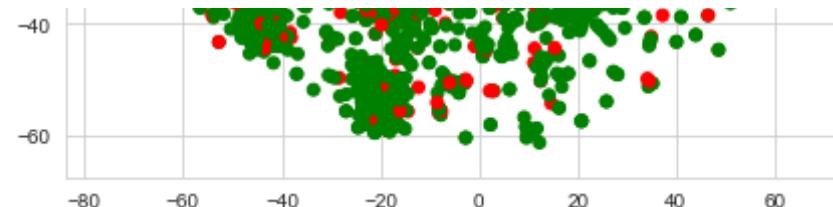
TSNE plot for merged data of Average Word2Vec Title and Categorical, Numerical features - Perplexity(10)



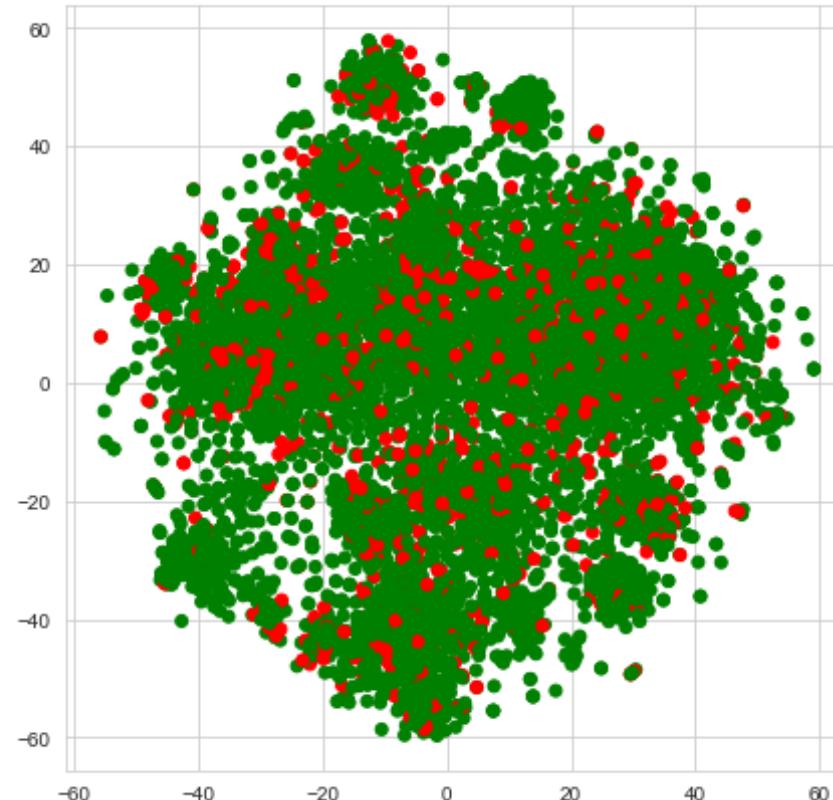


TSNE plot for merged data of Average Word2Vec Title and Categorical, Numerical features - Perplexity(30)



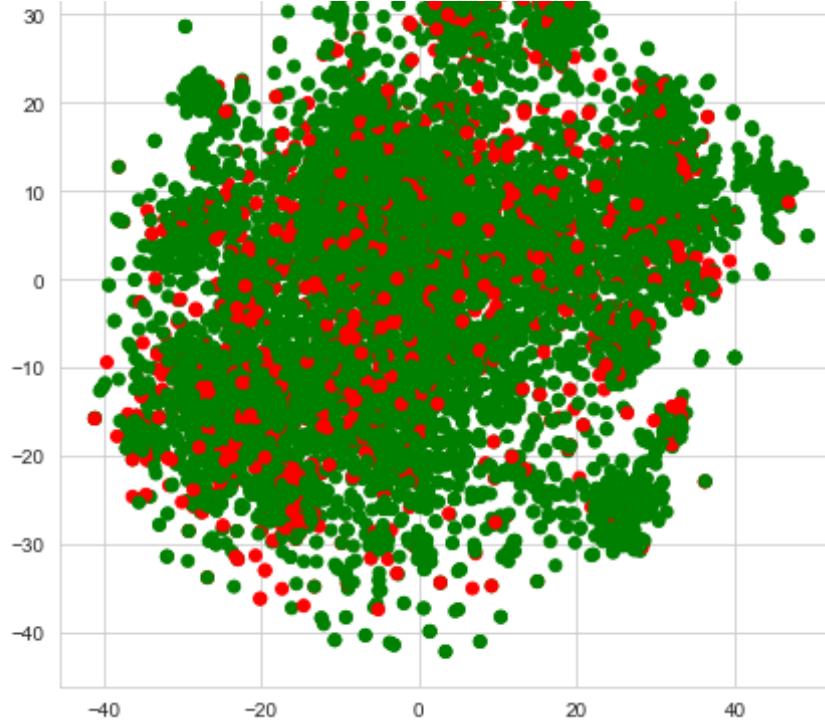


TSNE plot for merged data of Average Word2Vec Title and Categorical, Numerical features - Perplexity(50)

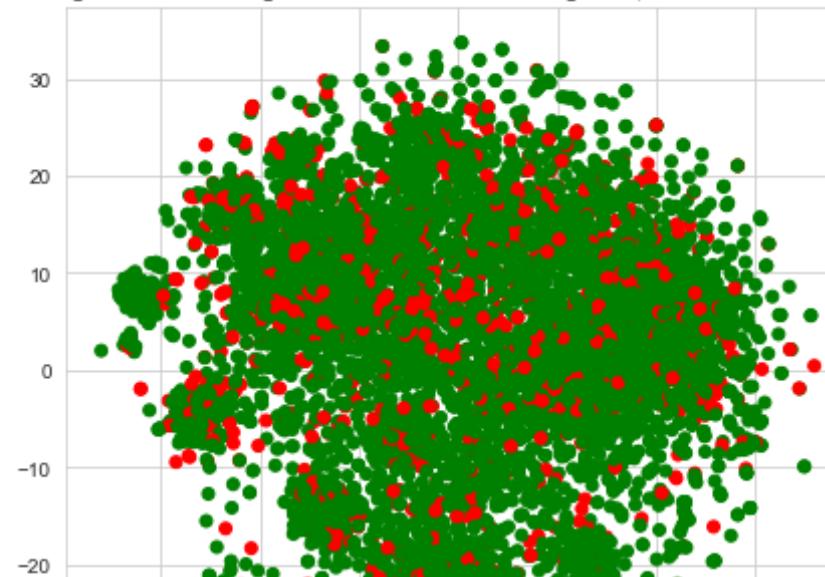


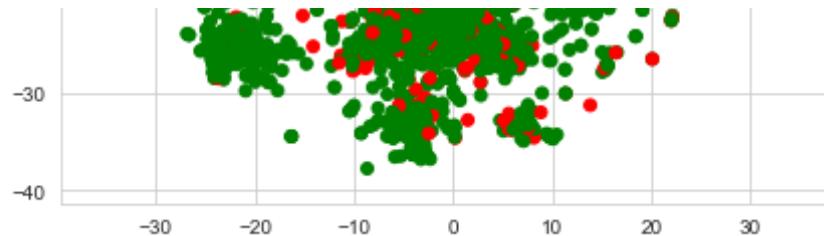
TSNE plot for merged data of Average Word2Vec Title and Categorical, Numerical features - Perplexity(80)





TSNE plot for merged data of Average Word2Vec Title and Categorical, Numerical features - Perplexity(100)





Classification using data merged with Tf-idf Weighted Word2Vec vectorized title and all considered categorical, numerical features

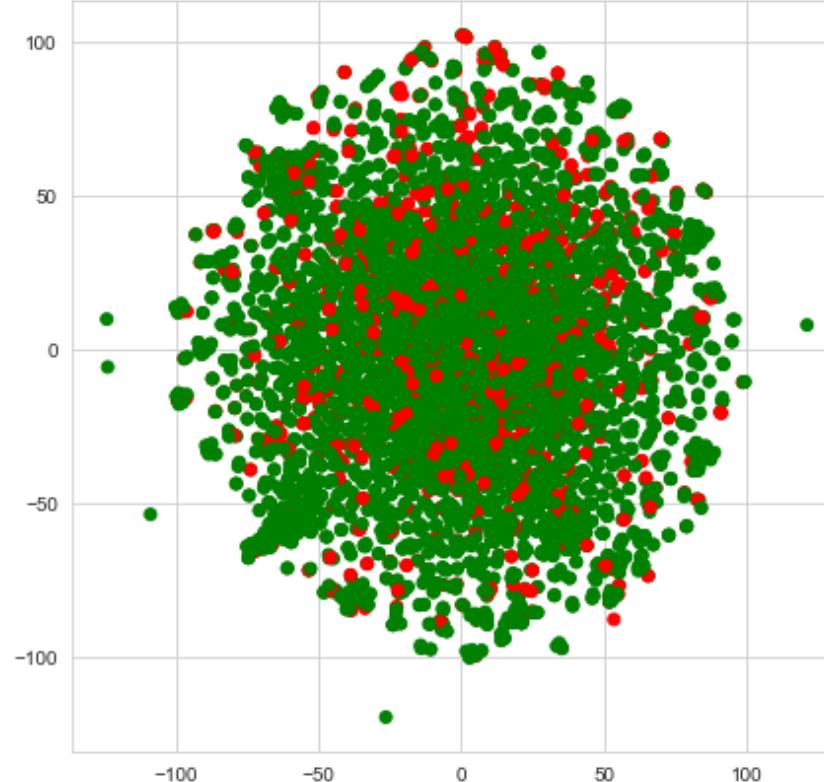
```
In [0]: tfIdfWeightedWord2VecTitleAndOthers = hstack((tfIdfWeightedWord2VecTitlesVectorsSub, categoriesVectorsSub, subCategoriesVectorsSub, teacherPrefixVectorsSub, schoolStateVectorsSub, projectGradeVectorsSub, priceStandardizedSub, previouslyPostedStandardizedSub));
tfIdfWeightedWord2VecTitleAndOthers.shape
```

Out[0]: (6000, 401)

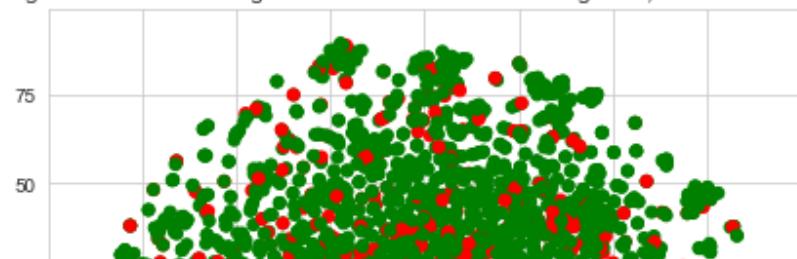
```
In [0]: perplexityValues = [5, 10, 30, 50, 80, 100]
for perplexityValue in perplexityValues:
    tsne = TSNE(n_components = 2, perplexity = perplexityValue, learning_rate = 200);
    tfIdfWeightedWord2VecTitleAndOthersEmbedded = tsne.fit_transform(tfIdfWeightedWord2VecTitleAndOthers.toarray());
    tfIdfWeightedWord2VecTitleAndOthersTsneData = np.hstack((tfIdfWeightedWord2VecTitleAndOthersEmbedded, classesDataSub.reshape(-1, 1)));
    tfIdfWeightedWord2VecTitleAndOthersTsneDataFrame = pd.DataFrame(tfIdfWeightedWord2VecTitleAndOthersTsneData, columns = ['Dimension1', 'Dimension2', 'Class']);
    colors = {0.0:'red', 1.0:'green'}
    plt.title("TSNE plot for merged data of Tf-Idf Weighted Word2Vec Title and Categorical, Numerical features - Perplexity({})".format(perplexityValue));
    plt.scatter(tfIdfWeightedWord2VecTitleAndOthersTsneDataFrame['Dimension1'], tfIdfWeightedWord2VecTitleAndOthersTsneDataFrame['Dimension2'], c = tfIdfWeightedWord2VecTitleAndOthersTsneDataFrame['Class'].apply(
```

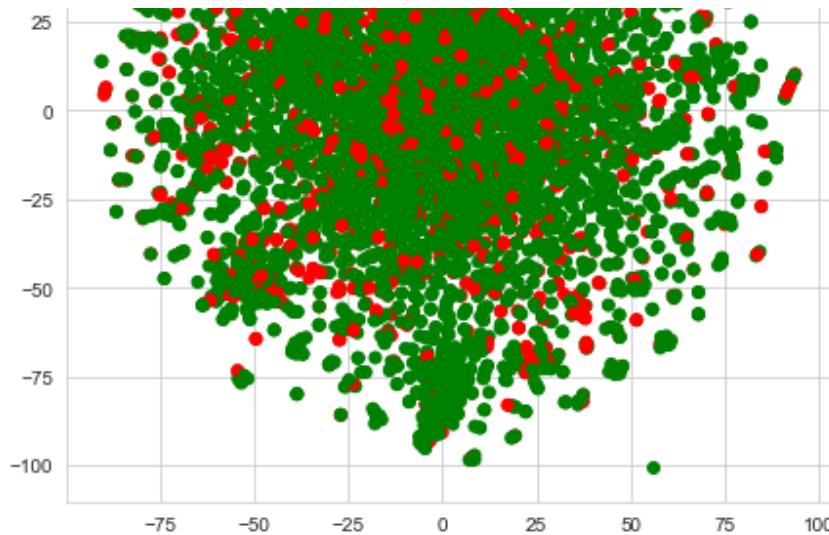
```
lambda x: colors[x]));
plt.show();
```

TSNE plot for merged data of Tf-Idf Weighted Word2Vec Title and Categorical, Numerical features - Perplexity(5)

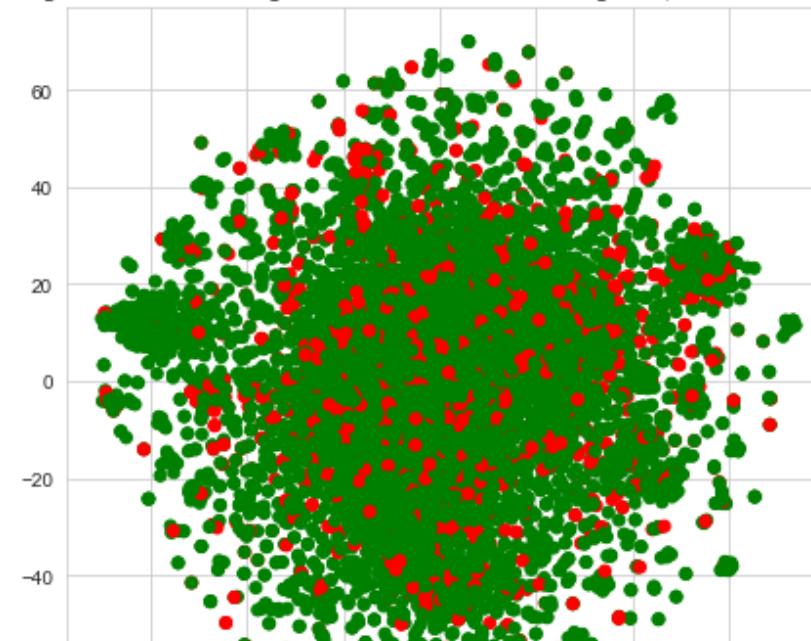


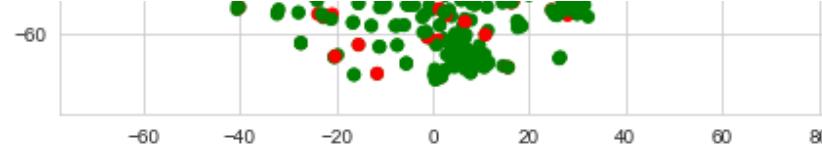
TSNE plot for merged data of Tf-Idf Weighted Word2Vec Title and Categorical, Numerical features - Perplexity(10)



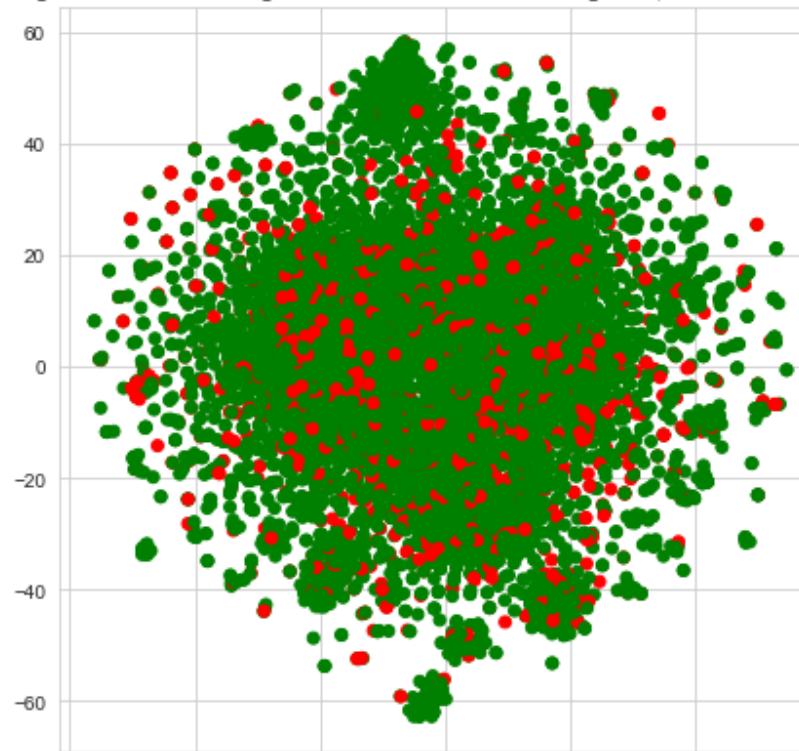


TSNE plot for merged data of Tf-Idf Weighted Word2Vec Title and Categorical, Numerical features - Perplexity(30)



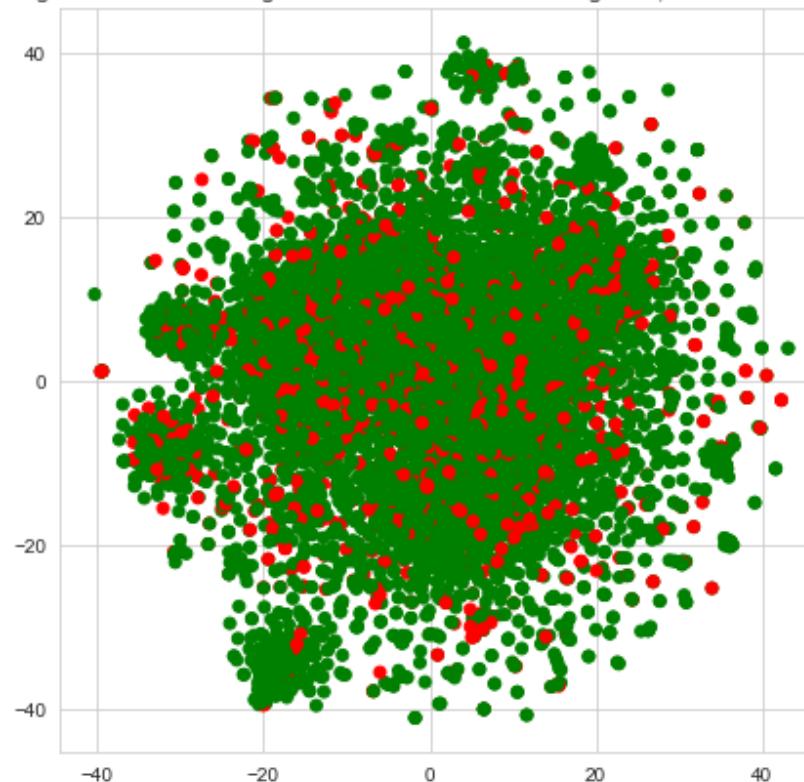


TSNE plot for merged data of Tf-Idf Weighted Word2Vec Title and Categorical, Numerical features - Perplexity(50)



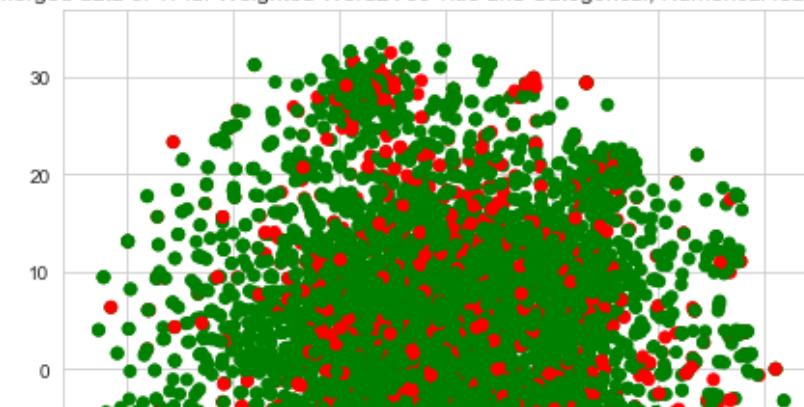
-60 -40 -20 0 20 40 60

TSNE plot for merged data of Tf-Idf Weighted Word2Vec Title and Categorical, Numerical features - Perplexity(80)

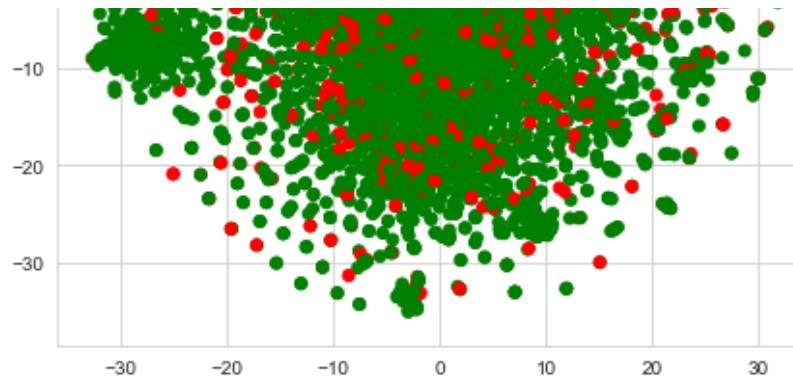


-40 -20 0 20 40

TSNE plot for merged data of Tf-Idf Weighted Word2Vec Title and Categorical, Numerical features - Perplexity(100)



-40 -20 0 20 40



Classification using data merged with all vectorizations of project\_title and with all considered categorical, numerical features

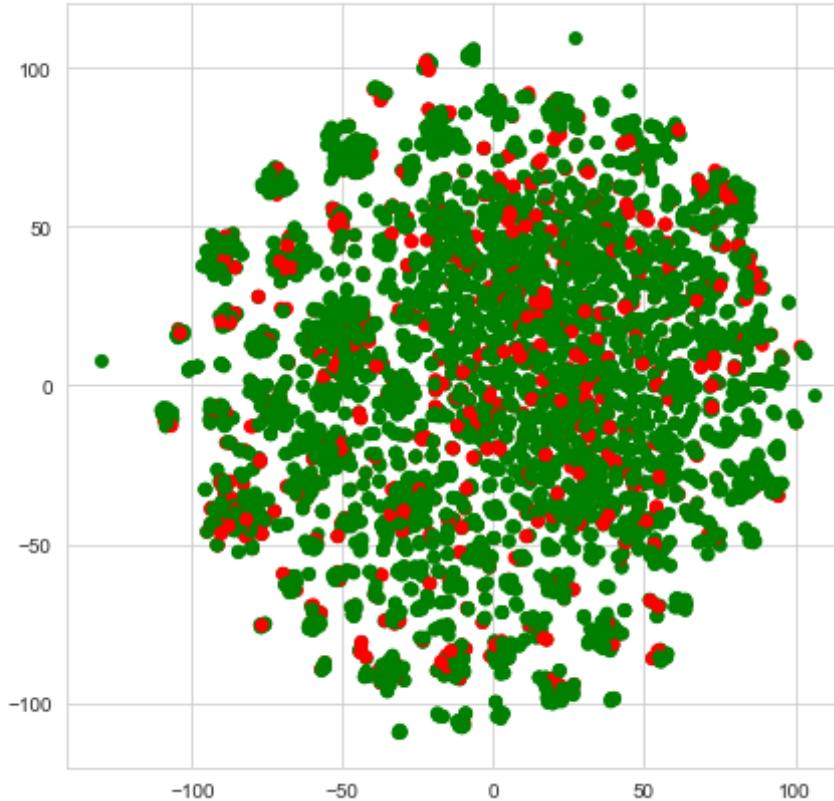
```
In [0]: allFeatures = hstack((bowTitleModelSub, tfIdfTitleModelSub, word2VecTitlesVectorsSub, tfIdfWeightedWord2VecTitlesVectorsSub, categoriesVectorsSub, subCategoriesVectorsSub, teacherPrefixVectorsSub, schoolStateVectorsSub, projectGradeVectorsSub, priceStandardizedSub, previouslyPostedStandardizedSub))
print(allFeatures.shape)

(6000, 4249)
```

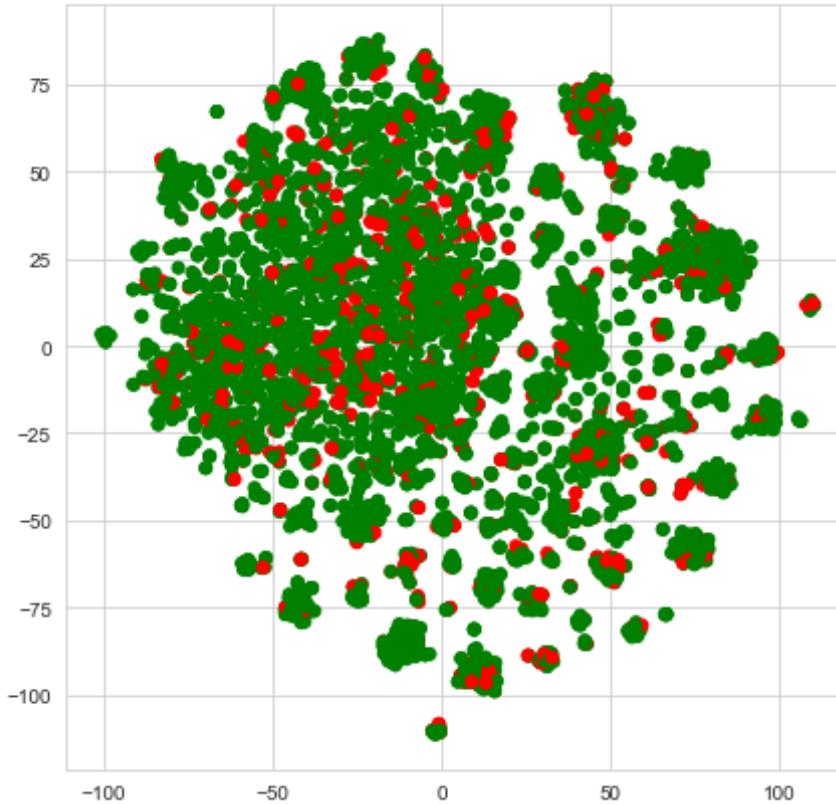
```
In [0]: perplexityValues = [5, 10, 30, 50, 80, 100]
for perplexityValue in perplexityValues:
    tsne = TSNE(n_components = 2, perplexity = perplexityValue, learning_rate = 200);
    allFeaturesEmbedded = tsne.fit_transform(allFeatures.toarray());
    allFeaturesTsneData = np.hstack((allFeaturesEmbedded, classesDataSub.reshape(-1, 1)));
    allFeaturesTsneDataFrame = pd.DataFrame(allFeaturesTsneData, columns = ['Dimension1', 'Dimension2', 'Class']);
    colors = {0.0:'red', 1.0:'green'}
    plt.title("TSNE plot for merged data of all vectorized features and Categorical, Numerical features - Perplexity({})".format(perplexityValue));
    plt.scatter(allFeaturesTsneDataFrame['Dimension1'], allFeaturesTsneDataFrame['Dimension2'], c=colors[allFeaturesTsneDataFrame['Class']])
```

```
 DataFrame['Dimension2'], c = allFeaturesTsnsDataFrame['Class'].apply(lambda x: colors[x]));
plt.show();
```

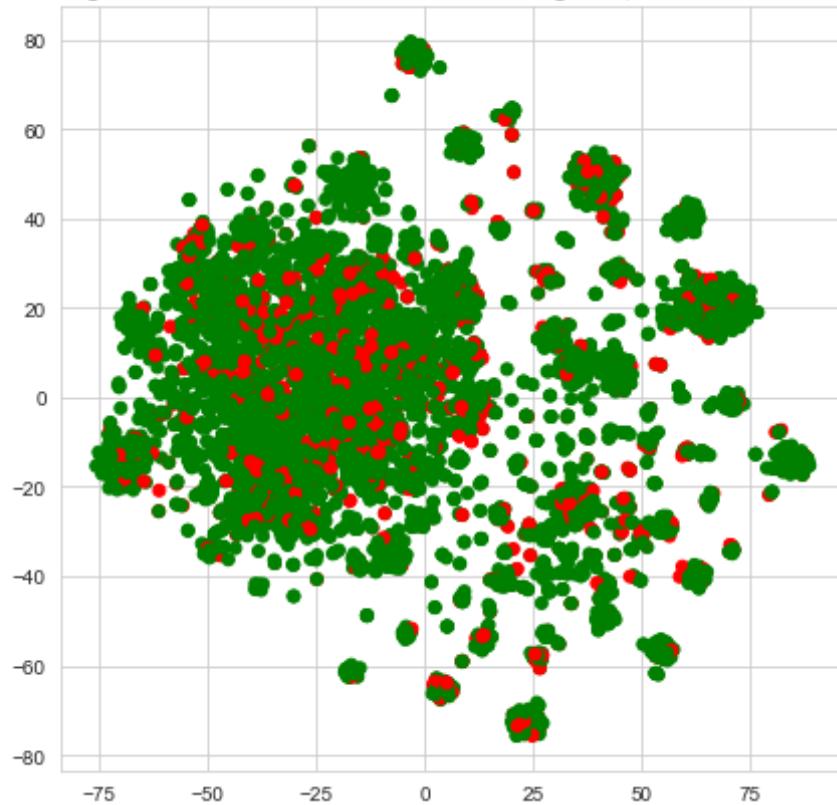
TSNE plot for merged data of all vectorized features and Categorical, Numerical features - Perplexity(5)



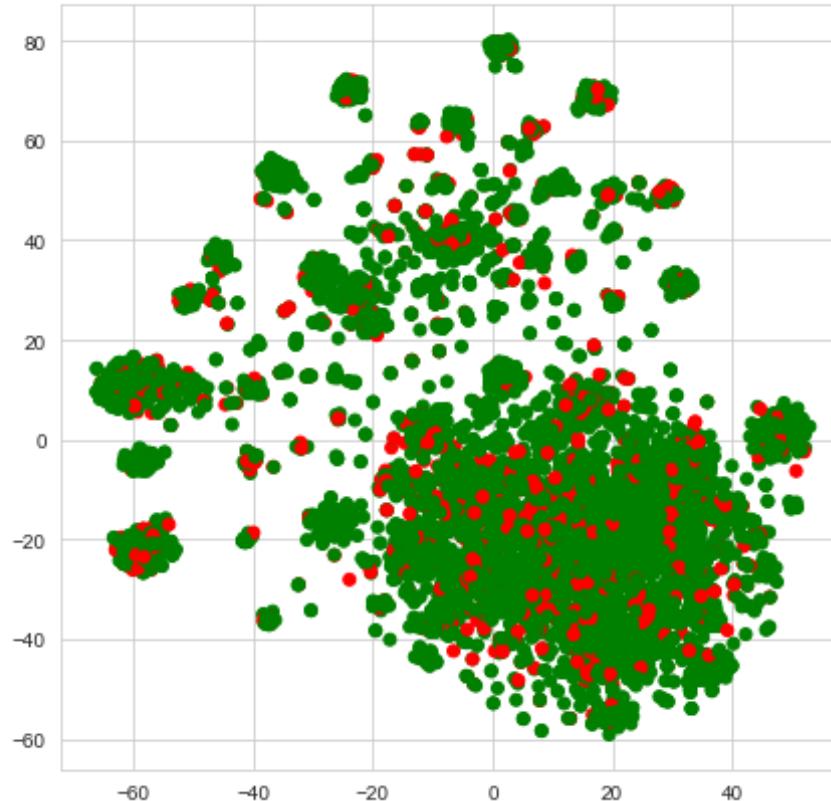
TSNE plot for merged data of all vectorized features and Categorical, Numerical features - Perplexity(10)



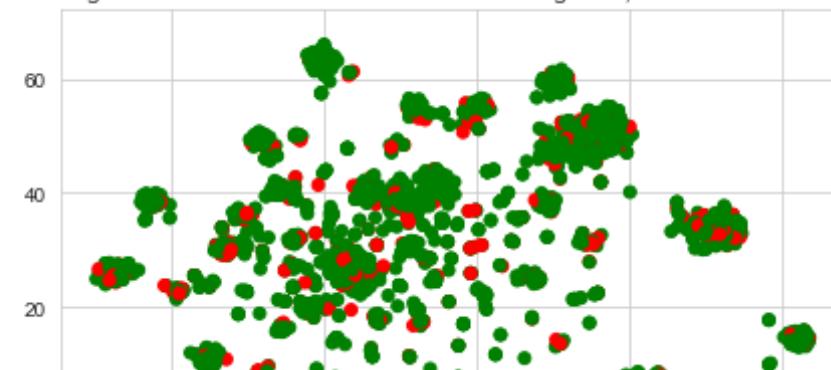
TSNE plot for merged data of all vectorized features and Categorical, Numerical features - Perplexity(30)

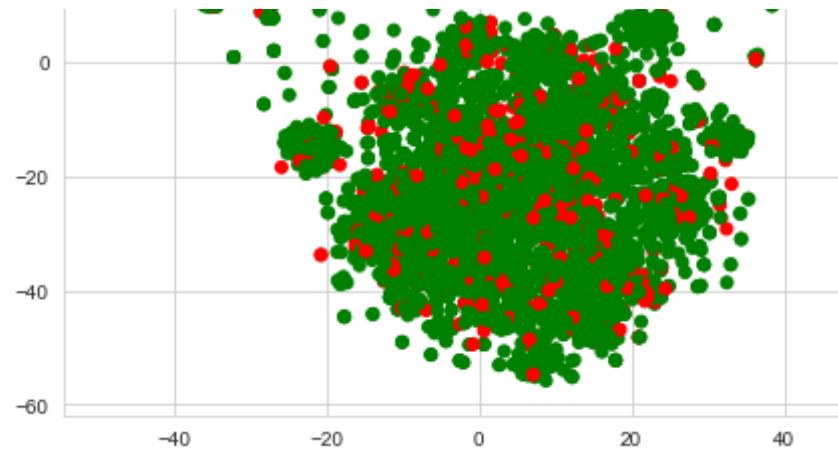


TSNE plot for merged data of all vectorized features and Categorical, Numerical features - Perplexity(50)

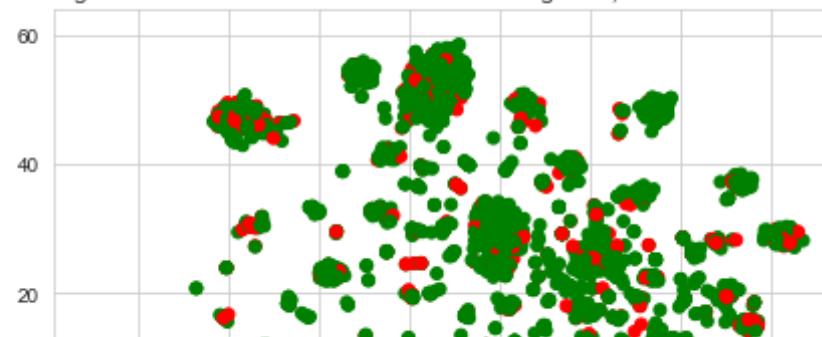


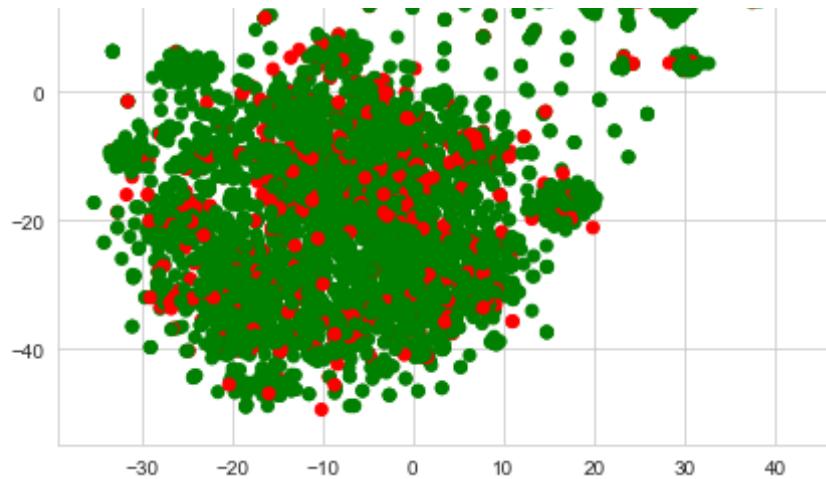
TSNE plot for merged data of all vectorized features and Categorical, Numerical features - Perplexity(80)





TSNE plot for merged data of all vectorized features and Categorical, Numerical features - Perplexity(100)





### Conclusion about data visualization using t-sne:

1. Bag of Words, Tf-Idf are better than word2vec vectorizations because of forming some small group of clusters with less overlap of overall data when compared to others.
2. Higher perplexity values seems better in data visualization because of less overlap of data than others.
3. None of the techniques are useful for classification because of huge overlap of data.
4. It is not seperable problem in 2-dimensions but it may be seperable in higher dimensions.

### Classification & Modelling using K-NN(With 50,000 data points)

#### Classification of data using K-NN(k-fold cross validation)

## Splitting Data(Only training and test)

```
In [63]: projectsData = projectsData.dropna(subset = ['teacher_prefix']);
projectsData.shape
```

```
Out[63]: (109245, 22)
```

```
In [64]: classesData = projectsData['project_is_approved']
print(classesData.shape)
```

```
(109245,)
```

```
In [0]: trainingData, testData, classesTraining, classesTest = model_selection.
train_test_split(projectsData[0:50000], classesData[0:50000], test_size
= 0.3, random_state = 0, stratify = classesData[0:50000]);
```

```
In [66]: print("Shapes of splitted data: ");
equalsBorder(70);

print("testData shape: ", testData.shape);
print("classesTest: ", classesTest.shape);
print("trainingData shape: ", trainingData.shape);
print("classesTraining shape: ", classesTraining.shape);
```

```
Shapes of splitted data:
```

```
=====
testData shape: (15000, 22)
classesTest: (15000,)
trainingData shape: (35000, 22)
classesTraining shape: (35000,)
```

```
In [67]: print("Number of negative points: ", trainingData[trainingData['project
_is_approved'] == 0].shape);
print("Number of positive points: ", trainingData[trainingData['project
_is_approved'] == 1].shape);
```

```
Number of negative points: (5400, 22)
Number of positive points: (29600, 22)
```

## Balancing Data

**Note: Instead of displaying whole vectorization process for balanced and imbalanced data, we have simply disabled below cell while performing analysis on imbalanced data and enabled while performing analysis on balanced data**

```
In [130]: negativeData = trainingData[trainingData['project_is_approved'] == 0];
positiveData = trainingData[trainingData['project_is_approved'] == 1];
negativeDataBalanced = resample(negativeData, replace = True, n_samples
= 29600, random_state = 44);
trainingData = pd.concat([positiveData, negativeDataBalanced]);
trainingData = shuffle(trainingData);
classesTraining = trainingData['project_is_approved'];
print("Testing whether data is balanced:");
equalsBorder(60);
print("Number of positive points: ", trainingData[trainingData['project
_is_approved'] == 1].shape);
print("Number of negative points: ", trainingData[trainingData['project
_is_approved'] == 0].shape);
```

Testing whether data is balanced:

```
=====
Number of positive points: (29600, 22)
Number of negative points: (29600, 22)
```

## Vectorizing categorical data

### 1. Vectorizing cleaned\_categories(project\_subject\_categories cleaned) - One Hot Encoding

```
In [0]: # Using CountVectorizer for performing one-hot-encoding by setting voca
bulary as list of all unique cleaned_categories
subjectsCategoriesVectorizer = CountVectorizer(vocabulary = list(sorted
```

```
CategoriesDictionary.keys(), lowercase = False, binary = True);  
# Fitting CountVectorizer with cleaned_categories values  
subjectsCategoriesVectorizer.fit(trainingData['cleaned_categories'].values);  
# Vectorizing categories using one-hot-encoding  
categoriesVectors = subjectsCategoriesVectorizer.transform(trainingData['cleaned_categories'].values);
```

```
In [132]: print("Features used in vectorizing categories: ");  
equalsBorder(70);  
print(subjectsCategoriesVectorizer.get_feature_names());  
equalsBorder(70);  
print("Shape of cleaned_categories matrix after vectorization(one-hot-encoding): ", categoriesVectors.shape);  
equalsBorder(70);  
print("Sample vectors of categories: ");  
equalsBorder(70);  
print(categoriesVectors[0:4])
```

Features used in vectorizing categories:

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

```
===== Shape of cleaned_categories matrix after vectorization(one-hot-encoding): (59200, 9) =====
```

Sample vectors of categories:

```
===== (0, 5) 1  
(1, 7) 1  
(1, 8) 1  
(2, 8) 1  
(3, 4) 1 =====
```

## 2. Vectorizing cleaned\_sub\_categories(project\_subject\_sub\_categories cleaned) -

## One Hot Encoding

```
In [0]: # Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique cleaned_sub_categories
subjectsSubCategoriesVectorizer = CountVectorizer(vocabulary = list(sortedDictionarySubCategories.keys()), lowercase = False, binary = True);
# Fitting CountVectorizer with cleaned_sub_categories values
subjectsSubCategoriesVectorizer.fit(trainingData['cleaned_sub_categories'].values);
# Vectorizing sub categories using one-hot-encoding
subCategoriesVectors = subjectsSubCategoriesVectorizer.transform(trainingData['cleaned_sub_categories'].values);
```

```
In [134]: print("Features used in vectorizing subject sub categories: ");
equalsBorder(70);
print(subjectsSubCategoriesVectorizer.get_feature_names());
equalsBorder(70);
print("Shape of cleaned_categories matrix after vectorization(one-hot-encoding): ", subCategoriesVectors.shape);
equalsBorder(70);
print("Sample vectors of categories: ");
equalsBorder(70);
print(subCategoriesVectors[0:4])
```

Features used in vectorizing subject sub categories:

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

```
=====
Shape of cleaned_categories matrix after vectorization(one-hot-encoding): (59200, 30)
```

=====
Sample vectors of categories:

```
=====
(0, 26)      1
(1, 27)      1
(1, 28)      1
(2, 29)      1
(3, 12)      1
(3, 19)      1
```

### 3. Vectorizing teacher\_prefix - One Hot Encoding

```
In [0]: def giveCounter(data):
    counter = Counter();
    for dataValue in data:
        counter.update(str(dataValue).split());
    return counter
```

```
In [136]: giveCounter(trainingData['teacher_prefix'].values)
```

```
Out[136]: Counter({'Dr': 4, 'Mr': 5788, 'Mrs': 30558, 'Ms': 21468, 'Teacher': 1382})
```

```
In [0]: teacherPrefixDictionary = dict(giveCounter(trainingData['teacher_prefix'].values));
# Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique teacher_prefix
teacherPrefixVectorizer = CountVectorizer(vocabulary = list(teacherPrefixDictionary.keys()), lowercase = False, binary = True);
# Fitting CountVectorizer with teacher_prefix values
teacherPrefixVectorizer.fit(trainingData['teacher_prefix'].values);
# Vectorizing teacher_prefix using one-hot-encoding
teacherPrefixVectors = teacherPrefixVectorizer.transform(trainingData['teacher_prefix'].values);
```

```
In [138]: print("Features used in vectorizing teacher_prefix: ");
equalsBorder(70);
print(teacherPrefixVectorizer.get_feature_names());
```

```
equalsBorder(70);
print("Shape of teacher_prefix matrix after vectorization(one-hot-encoding): ", teacherPrefixVectors.shape);
equalsBorder(70);
print("Sample vectors of teacher_prefix: ");
equalsBorder(70);
print(teacherPrefixVectors[0:100]);
```

Features used in vectorizing teacher\_prefix:

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

```
=====
Shape of teacher_prefix matrix after vectorization(one-hot-encoding):
(59200, 5)
```

```
=====
Sample vectors of teacher_prefix:
```

```
=====
(0, 0)      1
(1, 0)      1
(2, 0)      1
(3, 0)      1
(4, 1)      1
(5, 0)      1
(6, 2)      1
(7, 1)      1
(8, 0)      1
(9, 0)      1
(10, 0)     1
(11, 0)     1
(12, 0)     1
(13, 0)     1
(14, 2)     1
(15, 0)     1
(16, 2)     1
(17, 0)     1
(18, 0)     1
(19, 2)     1
(20, 0)     1
(21, 0)     1
(22, 3)     1
...       .
```

```
(23, 1)      1
(24, 2)      1
:
(75, 1)      1
(76, 0)      1
(77, 0)      1
(78, 0)      1
(79, 0)      1
(80, 0)      1
(81, 0)      1
(82, 2)      1
(83, 0)      1
(84, 2)      1
(85, 2)      1
(86, 2)      1
(87, 2)      1
(88, 0)      1
(89, 1)      1
(90, 0)      1
(91, 2)      1
(92, 0)      1
(93, 0)      1
(94, 2)      1
(95, 0)      1
(96, 0)      1
(97, 2)      1
(98, 0)      1
(99, 2)      1
```

```
In [139]: teacherPrefixes = [prefix.replace('.', '') for prefix in trainingData['teacher_prefix'].values];
teacherPrefixes[0:5]
```

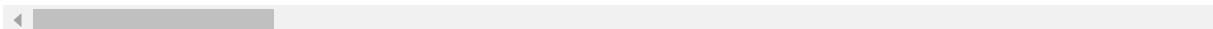
```
Out[139]: ['Mrs', 'Mrs', 'Mrs', 'Mrs', 'Mr']
```

```
In [140]: trainingData['teacher_prefix'] = teacherPrefixes;
trainingData.head(3)
```

```
Out[140]:
```

	Unnamed: 0	id		teacher_id	teacher_prefix	school_s
30339	89181	p191763	346fbe801e5dfcff1733b72fb1fe0da8	Mrs	PA	
11325	152869	p062725	e36201dc361b7cf959475d2ca552a86a	Mrs	IN	
48520	147928	p165461	bf1238aa96041677c85c2316b8d626f5	Mrs	IL	

3 rows × 22 columns



```
In [0]: teacherPrefixDictionary = dict(giveCounter(trainingData['teacher_prefix'].values));
# Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique teacher_prefix
teacherPrefixVectorizer = CountVectorizer(vocabulary = list(teacherPrefixDictionary.keys()), lowercase = False, binary = True);
# Fitting CountVectorizer with teacher_prefix values
teacherPrefixVectorizer.fit(trainingData['teacher_prefix'].values);
# Vectorizing teacher_prefix using one-hot-encoding
teacherPrefixVectors = teacherPrefixVectorizer.transform(trainingData['teacher_prefix'].values);
```

```
In [142]: print("Features used in vectorizing teacher_prefix: ");
equalsBorder(70);
print(teacherPrefixVectorizer.get_feature_names());
```

```

equalsBorder(70);
print("Shape of teacher_prefix matrix after vectorization(one-hot-encoding): ", teacherPrefixVectors.shape);
equalsBorder(70);
print("Sample vectors of teacher_prefix: ");
equalsBorder(70);
print(teacherPrefixVectors[0:4]);

Features used in vectorizing teacher_prefix:
=====
['Mrs', 'Mr', 'Ms', 'Teacher', 'Dr']
=====
Shape of teacher_prefix matrix after vectorization(one-hot-encoding):
(59200, 5)
=====
Sample vectors of teacher_prefix:
=====
(0, 0)      1
(1, 0)      1
(2, 0)      1
(3, 0)      1

```

#### 4. Vectorizing school\_state - One Hot Encoding

In [0]:

```

schoolStateDictionary = dict(giveCounter(trainingData['school_state'].values));
# Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique school states
schoolStateVectorizer = CountVectorizer(vocabulary = list(schoolStateDictionary.keys()), lowercase = False, binary = True);
# Fitting CountVectorizer with school_state values
schoolStateVectorizer.fit(trainingData['school_state'].values);
# Vectorizing school_state using one-hot-encoding
schoolStateVectors = schoolStateVectorizer.transform(trainingData['school_state'].values);

```

In [144]:

```
print("Features used in vectorizing school_state: ");
```

```
equalsBorder(70);
print(schoolStateVectorizer.get_feature_names());
equalsBorder(70);
print("Shape of school_state matrix after vectorization(one-hot-encoding): ", schoolStateVectors.shape);
equalsBorder(70);
print("Sample vectors of school_state: ");
equalsBorder(70);
print(schoolStateVectors[0:4]);
```

Features used in vectorizing school\_state:

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

```
=====
Shape of school_state matrix after vectorization(one-hot-encoding): (59200, 51)
```

Sample vectors of school\_state:

```
=====
(0, 0)      1
(1, 1)      1
(2, 2)      1
(3, 3)      1
```

## 5. Vectorizing project\_grade\_category - One Hot Encoding

```
In [145]: giveCounter(trainingData['project_grade_category'])
```

```
Out[145]: Counter({'Grades3to5': 19774,
                    'Grades6to8': 9073,
                    'Grades9to12': 6090,
                    'GradesPreKto2': 24263})
```

```
In [146]: cleanedGrades = []
```

```
for grade in trainingData['project_grade_category'].values:  
    grade = grade.replace(' ', '' );  
    grade = grade.replace('-', 'to' );  
    cleanedGrades.append(grade);  
cleanedGrades[0:4]
```

Out[146]: ['GradesPreKto2', 'GradesPreKto2', 'GradesPreKto2', 'Grades3to5']

In [147]: trainingData['project\_grade\_category'] = cleanedGrades  
trainingData.head(4)

Out[147]:

	Unnamed: 0	id		teacher_id	teacher_prefix	school_s
30339	89181	p191763	346fbe801e5dfcff1733b72fb1fe0da8	Mrs	PA	
11325	152869	p062725	e36201dc361b7cf959475d2ca552a86a	Mrs	IN	
48520	147928	p165461	bf1238aa96041677c85c2316b8d626f5	Mrs	IL	
5962	117756	p258343	817cecd4322151e0bf9c9b740d01303e	Mrs	SC	

4 rows × 22 columns

```
In [0]: projectGradeDictionary = dict(giveCounter(trainingData['project_grade_category'].values));
# Using CountVectorizer for performing one-hot-encoding by setting vocabulary as list of all unique project grade categories
projectGradeVectorizer = CountVectorizer(vocabulary = list(projectGradeDictionary.keys()), lowercase = False, binary = True);
# Fitting CountVectorizer with project_grade_category values
projectGradeVectorizer.fit(trainingData['project_grade_category'].values);
# Vectorizing project_grade_category using one-hot-encoding
projectGradeVectors = projectGradeVectorizer.transform(trainingData['project_grade_category'].values);
```

```
In [149]: print("Features used in vectorizing project_grade_category: ");
equalsBorder(70);
print(projectGradeVectorizer.get_feature_names());
equalsBorder(70);
print("Shape of school_state matrix after vectorization(one-hot-encoding): ", projectGradeVectors.shape);
equalsBorder(70);
print("Sample vectors of school_state: ");
equalsBorder(70);
print(projectGradeVectors[0:4]);
```

Features used in vectorizing project\_grade\_category:

```
=====
['GradesPreKto2', 'Grades3to5', 'Grades6to8', 'Grades9to12']
=====
```

```
=====
Shape of school_state matrix after vectorization(one-hot-encoding): (5
9200, 4)
=====
```

Sample vectors of school\_state:

```
=====
(0, 0)      1
(1, 0)      1
(2, 0)      1
(3, 1)      1
=====
```

```
In [150]: preProcessedEssaysWithStopWords, preProcessedEssaysWithoutStopWords = p  
reProcessingWithAndWithoutStopWords(trainingData['project_essay']);  
preProcessedProjectTitlesWithStopWords, preProcessedProjectTitlesWithoutStopWords = preProcessingWithAndWithoutStopWords(trainingData['project  
_title']);
```

## Vectorizing Text Data

### Bag of Words

#### 1. Vectorizing project\_essay

```
In [0]: # Initializing countvectorizer for bag of words vectorization of prepro  
cessed project essays  
bowEssayVectorizer = CountVectorizer(min_df = 10);  
# Transforming the preprocessed essays to bag of words vectors  
bowEssayModel = bowEssayVectorizer.fit_transform(preProcessedEssaysWithoutStopWords);
```

```
In [152]: print("Some of the Features used in vectorizing preprocessed essays: "  
);  
equalsBorder(70);  
print(bowEssayVectorizer.get_feature_names()[-40:]);  
equalsBorder(70);  
print("Shape of preprocessed essay matrix after vectorization: ", bowEs  
sayModel.shape);  
equalsBorder(70);  
print("Sample bag-of-words vector of preprocessed essay: ");  
equalsBorder(70);  
print(bowEssayModel[0])
```

Some of the Features used in vectorizing preprocessed essays:

```
=====
```

```
['yes', 'yesterday', 'yesterdays', 'yesteryear', 'yet', 'yield', 'yield  
s', 'yo', 'yoga', 'york', 'younannan', 'young', 'younger', 'youngest',  
'youngsters', 'youth', 'youthful', 'youths', 'youtube', 'yummy', 'zao',  
'zeal', 'zearn', 'zen', 'zenergy', 'zero', 'zest', 'zip', 'ziploc', 'zi  
plock', 'zone', 'zoned', 'zones', 'zoo', 'zoob', 'zoobs', 'zoology', 'z  
oom', 'zoos', 'zumba']
```

```
=====
```

```
Shape of preprocessed essay matrix after vectorization: (59200, 13266)
```

```
=====
```

```
Sample bag-of-words vector of preprocessed essay:
```

```
=====
```

```
(0, 7795)    1  
(0, 3664)    1  
(0, 2632)    1  
(0, 11879)   1  
(0, 10802)   1  
(0, 1901)    1  
(0, 3213)    1  
(0, 3850)    1  
(0, 9229)    1  
(0, 12689)   1  
(0, 12786)   1  
(0, 6424)    1  
(0, 4265)    1  
(0, 12693)   1  
(0, 1866)    1  
(0, 8698)    1  
(0, 12631)   1  
(0, 4403)    1  
(0, 2382)    1  
(0, 7891)    3  
(0, 8461)    1  
(0, 11924)   1  
(0, 6990)    1  
(0, 13081)   1  
(0, 11997)   1  
:      :
```

```
(0, 5323)      1
(0, 11697)     1
(0, 11745)     1
(0, 1193)      1
(0, 4193)      1
(0, 10919)     1
(0, 275)       1
(0, 369)       1
(0, 3413)      1
(0, 5951)      1
(0, 10581)     1
(0, 7606)      1
(0, 190)        1
(0, 10554)     1
(0, 3597)      1
(0, 10365)     2
(0, 12600)     1
(0, 6028)      1
(0, 7104)      2
(0, 11430)     8
(0, 13140)     1
(0, 5949)      1
(0, 6759)      4
(0, 11071)     3
(0, 11763)     2
```

## 2. Vectorizing project\_title

```
In [0]: # Initializing countvectorizer for bag of words vectorization of preprocessed project titles
bowTitleVectorizer = CountVectorizer(min_df = 10);
# Transforming the preprocessed project titles to bag of words vectors
bowTitleModel = bowTitleVectorizer.fit_transform(preProcessedProjectTitlesWithoutStopWords);
```

```
In [154]: print("Some of the Features used in vectorizing preprocessed titles: "
);
```

```
equalsBorder(70);
print(bowTitleVectorizer.get_feature_names()[-40:]);
equalsBorder(70);
print("Shape of preprocessed title matrix after vectorization: ", bowTitleModel.shape);
equalsBorder(70);
print("Sample bag-of-words vector of preprocessed title: ");
equalsBorder(70);
print(bowTitleModel[0])
```

Some of the Features used in vectorizing preprocessed titles:

```
=====
['wobbly', 'wolff', 'wolfpack', 'wonder', 'wonderful', 'wonders', 'wood
winds', 'woodworking', 'word', 'words', 'work', 'workers', 'working',
'workout', 'works', 'worksheets', 'workshop', 'world', 'worlds', 'worm
s', 'worth', 'would', 'wow', 'write', 'writers', 'writing', 'written',
'ye', 'year', 'yearbook', 'years', 'yes', 'yet', 'yoga', 'young', 'youn
gest', 'youth', 'zearn', 'zen', 'zone']
```

```
=====
Shape of preprocessed title matrix after vectorization: (59200, 2412)
```

```
=====
Sample bag-of-words vector of preprocessed title:
```

```
=====
(0, 2104)      1
(0, 1490)      1
(0, 976)       1
(0, 456)        1
(0, 35)         1
(0, 2081)      1
(0, 1052)      1
```

## Tf-Idf Vectorization

### 1. Vectorizing project\_essay

```
In [0]: # Intializing tfidf vectorizer for tf-idf vectorization of preprocessed
          project essays
```

```
tfIdfEssayVectorizer = TfidfVectorizer(min_df = 10);
# Transforming the preprocessed project essays to tf-idf vectors
tfIdfEssayModel = tfIdfEssayVectorizer.fit_transform(preProcessedEssays
WithoutStopWords);
```

```
In [156]: print("Some of the Features used in tf-idf vectorizing preprocessed ess
ays: ");
equalsBorder(70);
print(tfIdfEssayVectorizer.get_feature_names()[-40:]);
equalsBorder(70);
print("Shape of preprocessed title matrix after tf-idf vectorization: "
, tfIdfEssayModel.shape);
equalsBorder(70);
print("Sample Tf-Idf vector of preprocessed essay: ");
equalsBorder(70);
print(tfIdfEssayModel[0])
```

```
Some of the Features used in tf-idf vectorizing preprocessed essays:
=====
['yes', 'yesterday', 'yesterdays', 'yesteryear', 'yet', 'yield', 'yield
s', 'yo', 'yoga', 'york', 'younannan', 'young', 'younger', 'youngest',
'youngsters', 'youth', 'youthful', 'youths', 'youtube', 'yummy', 'zao',
'zeal', 'zearn', 'zen', 'zenergy', 'zero', 'zest', 'zip', 'ziploc', 'zi
plock', 'zone', 'zoned', 'zones', 'zoo', 'zoob', 'zoobs', 'zoology', 'z
oom', 'zoos', 'zumba']
```

```
=====
Shape of preprocessed title matrix after tf-idf vectorization: (59200,
13266)
```

```
=====
Sample Tf-Idf vector of preprocessed essay:
```

```
=====
(0, 11763)    0.09082140050225367
(0, 11071)    0.28489866932225605
(0, 6759)     0.19741278813483798
(0, 5949)     0.12607955973748775
(0, 13140)    0.046688303927799456
(0, 11430)    0.1505360449350122
(0, 7104)     0.0935696559357237
(0, 6028)     0.05126206641829252
```

(0, 12600)	0.07486849272169774
(0, 10365)	0.04362999278168776
(0, 3597)	0.05824813588618048
(0, 10554)	0.09831899936499679
(0, 190)	0.0993377753154874
(0, 7606)	0.11184664293018881
(0, 10581)	0.10249356739557147
(0, 5951)	0.1086476967808625
(0, 3413)	0.09492860719237482
(0, 369)	0.08445272131310896
(0, 275)	0.0546272200470455
(0, 10919)	0.05803104905129563
(0, 4193)	0.09417373660840721
(0, 1193)	0.06510837178172403
(0, 11745)	0.08549962201736003
(0, 11697)	0.04885047783733419
(0, 5323)	0.10473695743155043
:	:
(0, 11997)	0.10489656052237249
(0, 13081)	0.09571531189487342
(0, 6990)	0.1375723500394351
(0, 11924)	0.07641164401255328
(0, 8461)	0.07702040208340973
(0, 7891)	0.1135240006962371
(0, 2382)	0.03259768549428528
(0, 4403)	0.08146827604182014
(0, 12631)	0.10054190196958324
(0, 8698)	0.09180457254556136
(0, 1866)	0.09752494388919074
(0, 12693)	0.16600781752172658
(0, 4265)	0.08315905885954136
(0, 6424)	0.07675605646321514
(0, 12786)	0.075540848566564
(0, 12689)	0.10442178967481681
(0, 9229)	0.14925448343314707
(0, 3850)	0.11003036534807484
(0, 3213)	0.1419212962420623
(0, 1901)	0.1708182609164655
(0, 10802)	0.03867293494193706

```
(0, 11879)      0.06993402095362239
(0, 2632)       0.11037115085040936
(0, 3664)       0.08633413322941873
(0, 7795)       0.019491473369624587
```

## 2. Vectorizing project\_title

```
In [0]: # Intializing tfidf vectorizer for tf-idf vectorization of preprocessed
# project titles
tfIdfTitleVectorizer = TfidfVectorizer(min_df = 10);
# Transforming the preprocessed project titles to tf-idf vectors
tfIdfTitleModel = tfIdfTitleVectorizer.fit_transform(preProcessedProjec
tTitlesWithoutStopWords);
```

```
In [158]: print("Some of the Features used in tf-idf vectorizing preprocessed tit
les: ");
equalsBorder(70);
print(tfIdfTitleVectorizer.get_feature_names()[-40:]);
equalsBorder(70);
print("Shape of preprocessed title matrix after tf-idf vectorization: "
, tfIdfTitleModel.shape);
equalsBorder(70);
print("Sample Tf-Idf vector of preprocessed title: ");
equalsBorder(70);
print(tfIdfTitleModel[0])
```

Some of the Features used in tf-idf vectorizing preprocessed titles:

```
=====
['wobbly', 'wolff', 'wolfpack', 'wonder', 'wonderful', 'wonders', 'wood
winds', 'woodworking', 'word', 'words', 'work', 'workers', 'working',
'workout', 'works', 'worksheets', 'workshop', 'world', 'worlds', 'worm
s', 'worth', 'would', 'wow', 'write', 'writers', 'writing', 'written',
'ye', 'year', 'yearbook', 'years', 'yes', 'yet', 'yoga', 'young', 'you
gest', 'youth', 'zearn', 'zen', 'zone']
```

```
=====
Shape of preprocessed title matrix after tf-idf vectorization: (59200,
2412)
```

Sample Tf-Idf vector of preprocessed title:

```
=====
```

```
(0, 1052)      0.2616760041846661
(0, 2081)      0.25208670560757857
(0, 35)        0.5103968927851699
(0, 456)       0.45692389661844385
(0, 976)       0.4698377009397147
(0, 1490)      0.31726335734686106
(0, 2104)      0.2780108326752883
```

## Average Word2Vector Vectorization

```
In [0]: # storing variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
# We should have glove_vectors file for creating below model
with open('drive/My Drive/glove_vectors', 'rb') as f:
    gloveModel = pickle.load(f)
    gloveWords = set(gloveModel.keys())
```

```
In [160]: print("Glove vector of sample word: ");
equalsBorder(70);
print(gloveModel['technology']);
equalsBorder(70);
print("Shape of glove vector: ", gloveModel['technology'].shape);
```

Glove vector of sample word:

```
=====
```

```
[-0.26078   -0.36898   -0.022831   0.21666   0.16672   -0.20268
 -3.1219    0.33057    0.71512    0.28874   0.074368  -0.033203
  0.23783   0.21052    0.076562   0.13007   -0.31706   -0.45888
 -0.45463   -0.13191   0.49761    0.072704  0.16811    0.18846
 -0.16688   -0.21973   0.08575    -0.19577  -0.2101    -0.32436
 -0.56336   0.077996   -0.22758   -0.66569  0.14824    0.038945
  0.50881   -0.1352    0.49966   -0.4401   -0.022335  -0.22744
  0.22086   0.21865    0.36647    0.30495  -0.16565   0.038759
  0.28108   -0.2167    0.12453    0.65401   0.34584   -0.2557
 -0.046363  -0.31111   -0.020936  -0.17122  -0.77114   0.29289]
```

-0.14625	0.39541	-0.078938	0.051127	0.15076	0.085126
0.183	-0.06755	0.26312	0.0087276	0.0066415	0.37033
0.03496	-0.12627	-0.052626	-0.34897	0.14672	0.14799
-0.21821	-0.042785	0.2661	-1.1105	0.31789	0.27278
0.054468	-0.27458	0.42732	-0.44101	-0.19302	-0.32948
0.61501	-0.22301	-0.36354	-0.34983	-0.16125	-0.17195
-3.363	0.45146	-0.13753	0.31107	0.2061	0.33063
0.45879	0.24256	0.042342	0.074837	-0.12869	0.12066
0.42843	-0.4704	-0.18937	0.32685	0.26079	0.20518
-0.18432	-0.47658	0.69193	0.18731	-0.12516	0.35447
-0.1969	-0.58981	-0.88914	0.5176	0.13177	-0.078557
0.032963	-0.19411	0.15109	0.10547	-0.1113	-0.61533
0.0948	-0.3393	-0.20071	-0.30197	0.29531	0.28017
0.16049	0.25294	-0.44266	-0.39412	0.13486	0.25178
-0.044114	1.1519	0.32234	-0.34323	-0.10713	-0.15616
0.031206	0.46636	-0.52761	-0.39296	-0.068424	-0.04072
0.41508	-0.34564	0.71001	-0.364	0.2996	0.032281
0.34035	0.23452	0.78342	0.48045	-0.1609	0.40102
-0.071795	-0.16531	0.082153	0.52065	0.24194	0.17113
0.33552	-0.15725	-0.38984	0.59337	-0.19388	-0.39864
-0.47901	1.0835	0.24473	0.41309	0.64952	0.46846
0.024386	-0.72087	-0.095061	0.10095	-0.025229	0.29435
-0.57696	0.53166	-0.0058338	-0.3304	0.19661	-0.085206
0.34225	0.56262	0.19924	-0.027111	-0.44567	0.17266
0.20887	-0.40702	0.63954	0.50708	-0.31862	-0.39602
-0.1714	-0.040006	-0.45077	-0.32482	-0.0316	0.54908
-0.1121	0.12951	-0.33577	-0.52768	-0.44592	-0.45388
0.66145	0.33023	-1.9089	0.5318	0.21626	-0.13152
0.48258	0.68028	-0.84115	-0.51165	0.40017	0.17233
-0.033749	0.045275	0.37398	-0.18252	0.19877	0.1511
0.029803	0.16657	-0.12987	-0.50489	0.55311	-0.22504
0.13085	-0.78459	0.36481	-0.27472	0.031805	0.53052
-0.20078	0.46392	-0.63554	0.040289	-0.19142	-0.0097011
0.068084	-0.10602	0.25567	0.096125	-0.10046	0.15016
-0.26733	-0.26494	0.057888	0.062678	-0.11596	0.28115
0.25375	-0.17954	0.20615	0.24189	0.062696	0.27719
-0.42601	-0.28619	-0.44697	-0.082253	-0.73415	-0.20675
-0.60289	-0.06728	0.15666	-0.042614	0.41368	-0.17367
-0.54012	0.23883	0.23075	0.13608	-0.058634	-0.089705

```
 0.18469  0.023634  0.16178  0.23384  0.24267  0.091846 ]  
=====  
Shape of glove vector: (300,)
```

```
In [0]: def getWord2VecVectors(texts):  
    word2VecTextsVectors = [];  
    for preProcessedText in tqdm(texts):  
        word2VecTextVector = np.zeros(300);  
        number0fWordsInText = 0;  
        for word in preProcessedText.split():  
            if word in gloveWords:  
                word2VecTextVector += gloveModel[word];  
                number0fWordsInText += 1;  
        if number0fWordsInText != 0:  
            word2VecTextVector = word2VecTextVector / number0fWordsInTe  
xt;  
    word2VecTextsVectors.append(word2VecTextVector);  
    return word2VecTextsVectors;
```

## 1. Vectorizing project\_essay

```
In [162]: word2VecEssaysVectors = getWord2VecVectors(preProcessedEssaysWithoutSto  
pwords);
```

```
In [163]: print("Shape of Word2Vec vectorization matrix of essays: {},{}".format(  
len(word2VecEssaysVectors), len(word2VecEssaysVectors[0])));  
equalsBorder(70);  
print("Sample essay: ");  
equalsBorder(70);  
print(preProcessedEssaysWithoutStopWords[0]);  
equalsBorder(70);  
print("Word2Vec vector of sample essay: ");  
equalsBorder(70);  
print(word2VecEssaysVectors[0]);
```

Shape of Word2Vec vectorization matrix of essays: 59200,300

```
=====
Sample essay:
=====
```

teacher speech language impaired working students low income urban school district service 70 students moderate severe speech language impairments students difficulty across academic social environments basic tasks take granted including saying good morning participating classroom discussions require practice coaching encouragement students thrive routine love coming school work hard every day overcome disabilities students work various speech language goals including receptive expressive language articulation voice fluency pragmatics achieve functional communication students require low tech materials learn others benefit multi sensory materials including computers web based applications art materials budgets tight teacher wish lists thing past new materials hard come students expand communication utilizing new picture cards verbs opposites associations etc ipads visual verbal prompting easel delivery carryover new skills thank consideration donating students nannan

```
=====
Word2Vec vector of sample essay:
=====
```

```
[ 2.64139000e-02 -1.03656772e-02  4.66371407e-02 -5.42995854e-02
 -6.79666269e-02 -3.90240976e-02 -2.83907585e+00  1.29914374e-01
  1.32283618e-02  1.29509931e-01 -2.15022927e-02 -1.98975837e-03
  2.46369553e-02 -1.36138271e-01 -9.67092301e-02 -8.38304553e-02
  5.31572141e-02 -2.05765813e-02  7.37989593e-02 -1.18184545e-02
 -4.65190593e-02  4.74208179e-02  4.30130146e-02 -3.75941965e-02
 -1.84983048e-02 -7.65586528e-02  1.47973675e-01 -1.45888289e-01
 -8.77171512e-02 -1.05967648e-01 -2.70467412e-01 -6.90247057e-02
  3.15854870e-02  6.89639837e-02 -1.11974002e-01 -7.61549496e-02
  1.36052902e-02 -4.12286220e-02 -6.10443057e-02 -3.92171528e-02
 -3.32440353e-02  1.19208419e-01 -2.40512626e-02 -1.36502031e-01
  2.51629984e-02 -4.09203569e-02  4.60451175e-02  5.42080811e-02
 -3.94468545e-02 -5.86441659e-02 -4.67677837e-02 -3.40983683e-02
  1.71027691e-02  5.64662846e-02  2.11189593e-02 -6.48145285e-02
 -1.96077049e-02  8.39192008e-03 -1.13449745e-01 -3.16621713e-02
 -6.14107821e-02 -1.07176959e-02  9.13451794e-02 -9.90678147e-02
 -7.80341715e-02  5.34884089e-02  2.94555802e-02 -8.82515781e-02
  1.06843683e-01 -6.96431301e-03 -9.88358886e-02 -3.63625797e-02
  4.29165813e-02 -7.52991805e-02 -9.53029512e-03 -2.07781465e-01
```

-5.25286637e-02	-2.41586179e-02	3.90960740e-02	1.12981992e-02
4.96276041e-02	-4.12602374e-01	5.84048293e-02	-9.41574992e-02
-7.66106472e-02	-3.35250569e-02	9.25303585e-02	-3.88805756e-02
1.57611137e-01	1.90392927e-02	-4.24900252e-02	1.21933000e-02
2.45959260e-02	-2.73997602e-02	-1.04229167e-01	-6.91625602e-02
-2.22957634e+00	1.14720215e-01	3.60796081e-02	2.04151993e-01
-8.44915976e-02	-7.69686624e-02	6.09835691e-02	-6.82776618e-02
1.17024415e-01	4.46428861e-02	6.75109951e-02	-1.05180073e-01
8.47961594e-02	8.49469647e-02	-5.65097561e-02	-2.87090618e-02
5.97661098e-02	1.98515136e-01	1.35192611e-01	1.69008333e-02
-2.17186560e-01	2.50708423e-02	7.09105654e-02	-3.08169820e-02
-5.58603171e-02	1.94908706e-02	-6.96867854e-03	-4.55770488e-02
5.64406293e-02	-3.33409984e-02	-1.07356358e-02	-1.99374691e-02
9.42609512e-03	7.83613854e-02	1.67982277e-01	-3.23530797e-02
-1.42924520e-02	-3.94636626e-02	-1.18903322e-01	-6.03470138e-02
1.11492237e-01	7.45746341e-04	9.79713967e-02	1.70162667e-01
7.36360610e-02	4.00126350e-02	4.16237317e-02	-3.74737699e-02
-5.60050787e-02	4.09956886e-02	1.25715348e-01	-7.91884585e-02
1.69255969e-01	-8.75312016e-02	-1.03477382e-03	-9.12928276e-02
2.64354691e-03	3.34730185e-02	-1.17833616e-02	-3.65696724e-02
7.00483065e-02	-5.20841008e-02	-3.42186990e-02	-2.39804894e-02
3.02448780e-02	4.53391480e-02	4.79764421e-02	1.44444528e-02
-4.39220122e-02	-4.87527642e-02	-7.04699213e-02	1.76224935e-02
-7.11332862e-03	-7.97237268e-02	2.40269919e-02	-6.78273577e-03
3.88495935e-04	-1.22577192e-01	5.37204155e-02	5.34761285e-02
-5.50476547e-02	3.50927911e-02	-1.55167502e-01	-4.79386260e-04
-6.03507480e-02	2.14006959e-01	9.48389634e-02	-5.20955854e-03
-1.25615432e-01	5.96900491e-02	-6.18750065e-03	-4.96629098e-02
-6.35718098e-02	3.25734593e-02	1.00759813e-02	-3.72842748e-02
-1.03469643e-01	-2.83488301e-02	-1.88673642e-02	-6.88466415e-02
-6.68359659e-02	1.03524108e-01	3.80108528e-02	9.07260157e-02
1.82775292e-01	-3.14985675e-02	-1.55638049e-02	3.42808797e-02
-9.53908211e-02	-4.98048870e-02	8.88515715e-02	-1.28009524e-01
9.85152309e-02	-1.74774163e-02	-2.28406740e-02	4.73470894e-02
3.84675638e-02	-1.01562789e-01	-1.73566026e-01	6.16483008e-02
-7.51423715e-02	-8.02684545e-02	-3.34895244e-02	-2.60461073e-02
-1.68603872e-01	-1.08821807e-01	-1.67570749e-01	-2.31221008e-02
-1.63629854e+00	2.02827619e-01	-1.00571059e-01	2.11941382e-03
-7.44693203e-02	-1.14754733e-01	-3.06726667e-02	-7.40186366e-02

```
4.30041463e-04 -1.27752250e-01 -2.80107309e-02 -4.03755798e-02
5.17559390e-02 -1.12606712e-01 3.35728707e-02 1.57177450e-01
-6.84731659e-02 3.28213642e-02 -1.92280428e-01 -7.46220699e-02
-5.33523681e-02 5.21797203e-02 2.51429211e-02 -1.18179419e-01
-1.92825610e-02 -7.89913553e-02 -3.00810392e-02 1.84475211e-01
4.84087216e-02 -6.99062772e-02 6.06970862e-02 1.67683992e-02
2.01109618e-02 -9.01677073e-02 2.03833184e-01 -1.65582949e-01
-1.13755122e-03 8.09234593e-02 6.52108598e-02 1.02856195e-02
3.08906122e-02 -1.21981317e-01 -1.08018216e-01 -1.30018583e-01
3.44232108e-02 4.44549048e-02 -3.85859707e-02 5.88687398e-02
-1.84054433e-01 1.25702990e-01 -2.75518690e-02 8.32098276e-02
6.39334008e-02 -6.62163512e-03 -4.15211463e-02 9.53861628e-02
-6.21571260e-02 -6.06478967e-02 -1.88637552e-02 8.26255171e-02
-4.80093325e-02 7.38454138e-02 1.27526285e-02 3.58887057e-02
5.93906057e-02 8.73966748e-02 6.15277367e-02 -3.16783920e-02
2.36692943e-02 -1.26624980e-01 4.03968154e-02 -2.89239341e-02
-8.11545041e-02 2.03916707e-01 1.08410884e-01 1.73225081e-02]
```

## 2. Vectorizing project\_title

```
In [164]: word2VecTitlesVectors = getWord2VecVectors(preProcessedProjectTitlesWithoutStopWords);
```

```
In [165]: print("Shape of Word2Vec vectorization matrix of project titles: {}, {}".format(len(word2VecTitlesVectors), len(word2VecTitlesVectors[0])));
equalsBorder(70);
print("Sample title: ");
equalsBorder(70);
print(preProcessedProjectTitlesWithoutStopWords[0]);
equalsBorder(70);
print("Word2Vec vector of sample title: ");
equalsBorder(70);
print(word2VecTitlesVectors[0]);
```

Shape of Word2Vec vectorization matrix of project titles: 59200, 300

=====

Sample title:

=====

help students achieve communication goals new supplies

=====

Word2Vec vector of sample title:

=====

```
[ 3.58599286e-02  1.84358429e-01  6.24455714e-02 -4.71412857e-02
 -4.48928571e-02 -3.12258571e-02 -3.26542857e+00  1.08774429e-01
 1.79830571e-01  3.02331429e-01 -2.21750857e-01  7.33338571e-02
 -1.08487857e-01 -2.82982857e-01 -2.58511714e-01  4.09127143e-02
 1.28077143e-02 -1.74408714e-01  1.07378143e-01 -1.66132571e-01
 3.04367857e-02 -1.55811557e-01  8.17727143e-02 -4.15785714e-02
 7.53098571e-02 -1.30265143e-01  3.10358571e-01  1.70962857e-02
 -1.91833857e-01 -2.08737429e-01 -3.92974714e-01 -1.93696571e-01
 -5.93123714e-02 -9.18287143e-02  1.56044286e-01 -5.49844286e-02
 -6.48722857e-02 -1.32022857e-01 -1.24002857e-01 -1.59692286e-01
 -2.00012857e-02  7.93137143e-02  1.68705429e-01 -1.16606571e-01
 -1.13671429e-02 -5.66514286e-02 -9.34860000e-03  1.09194429e-01
 -1.16490714e-01 -1.37149429e-01  8.49185714e-02 -1.22370143e-01
 -1.77967143e-02  9.87665714e-02  1.32881286e-01  5.49012857e-02
 -6.34840000e-02 -6.55558000e-02 -9.76285714e-02  1.26923429e-01
 -1.41306429e-01 -7.08714286e-03 -3.72057143e-02 -1.34334286e-01
 -3.53578571e-02  2.20397143e-01  4.36980000e-02 -1.31079514e-01
 -9.26285714e-03 -8.88985714e-02 -2.82517857e-01 -1.08293286e-01
 -5.36587143e-02 -1.99059714e-01 -2.21442857e-01 -3.77005429e-01
 -2.26962857e-02 -2.13760000e-02 -4.97612857e-02 -2.12871429e-02
 1.94804000e-01 -4.29568571e-01  1.90542857e-02  1.83928571e-02
 -1.16381000e-01 -4.79917143e-02  1.00564000e-01 -7.50254286e-02
 3.04007000e-01  4.16568571e-02  9.75822857e-02 -1.57494714e-01
 -6.35042857e-02  4.98613571e-02  1.16806000e-01 -1.91141429e-01
 -2.75040000e+00 -9.30058714e-02  1.06755714e-01  9.83262857e-02
 -1.15591714e-01 -9.52262857e-02  3.04851429e-01  3.83200000e-02
 3.64157143e-01  8.05828571e-02 -6.37357143e-02  1.24780714e-01
 6.84160000e-02  8.13745714e-02 -1.90891286e-01 -1.39602429e-01
 5.04157143e-03  3.10949000e-01  1.50242857e-03 -2.54110286e-01
 -3.98340000e-01  1.74777143e-01  3.75097143e-02  2.43855714e-02
 -5.69917143e-02 -6.05336857e-02 -1.83347033e-01 -1.11100714e-01
 2.97132857e-02  1.14363286e-01  7.84602857e-02  9.98542857e-03
 7.39672429e-02  1.26285714e-02  1.16320857e-01 -5.76894286e-02
```

4.80794286e-02	6.60450000e-02	-1.07607571e-01	-3.58200000e-02
-4.58138571e-02	1.21713714e-01	1.04291857e-01	3.04672857e-01
-2.86100000e-02	-9.56550000e-03	3.54234714e-02	-1.08172657e-01
1.21583857e-01	6.55270000e-02	5.95418571e-02	6.53442857e-03
1.39536143e-01	4.05408571e-02	-5.29578571e-02	-4.02490000e-02
2.80677286e-01	2.32010857e-02	-1.73784729e-01	1.00637000e-01
2.10214857e-01	-9.22588571e-02	1.33228857e-01	-4.27796143e-01
6.24685714e-02	1.45815714e-02	-1.24963000e-01	-1.73338000e-01
-2.54835429e-01	6.37280000e-02	-1.53612143e-01	1.39938429e-01
1.92028571e-01	-2.43279571e-01	-8.26884286e-02	1.21495857e-01
3.73642857e-03	-2.74055571e-01	-4.94607143e-02	-1.99414286e-02
-1.54208157e-01	3.45785714e-02	-2.18685714e-03	4.89657857e-02
4.43497143e-02	2.04480000e-01	-7.55635714e-02	9.97154286e-02
-2.29658571e-02	-1.64465714e-02	6.46419000e-02	8.00474286e-02
-5.45580000e-02	6.18284286e-02	-1.38109286e-01	1.27303714e-01
-1.37338143e-01	-1.69194286e-01	1.27940857e-01	-2.49912286e-01
-4.00922286e-02	2.03992086e-01	-1.01642857e-03	2.63515857e-01
1.85248714e-01	-1.61614286e-01	-1.58317714e-01	1.62107586e-01
-2.13297843e-01	-5.03420000e-02	1.70956157e-01	-2.51525714e-02
8.55785714e-02	-7.32261429e-02	-7.20721429e-02	-3.68628571e-02
-5.77283929e-02	-2.77523857e-01	-2.39722754e-01	9.46920000e-02
3.23488571e-02	-2.28255714e-01	-5.84176714e-02	7.78871429e-02
-4.18634429e-01	-1.34635714e-01	-2.34114286e-01	-1.86807429e-01
-1.88791714e+00	3.09474000e-01	-1.17988286e-01	6.71065714e-02
-9.89038571e-02	-5.58550429e-02	4.26802857e-02	-1.89420000e-01
-1.29985714e-01	2.60400000e-02	2.12255714e-02	3.14854286e-01
1.18105143e-01	-9.22777143e-02	9.96915714e-02	9.84124286e-02
1.24285714e-02	-1.70854286e-02	-3.35268143e-01	-1.84982714e-01
-1.20007857e-01	-1.13500000e-02	3.82493571e-02	-2.23834143e-01
-9.67078571e-02	-2.58624286e-01	-4.37072857e-02	1.72095714e-02
1.24809429e-01	-2.02338714e-01	1.38150857e-01	-7.75595714e-02
1.33574714e-01	5.40310000e-02	7.69847143e-02	-1.55639857e-01
1.99719143e-01	6.88692286e-02	1.23702857e-02	2.16670000e-02
-1.27467143e-02	-4.77471429e-02	8.51634286e-02	-5.61835714e-02
2.93782857e-02	4.57400000e-03	-1.18510286e-01	-5.27137143e-02
-3.98847571e-02	8.63601429e-02	1.14924000e-01	1.90435143e-01
-5.40362857e-02	-2.01485286e-01	1.70955714e-02	1.35353143e-01
-2.77114129e-01	-2.40557143e-03	2.23442857e-03	-1.77624286e-02
9.80580000e-02	2.88576000e-01	1.60625714e-02	8.40617143e-02

```
-2.90109429e-02 8.09714286e-03 1.20169143e-01 -1.48409200e-01
-3.96180000e-02 -1.29316286e-01 -5.06684286e-02 -1.22183286e-01
-3.34631429e-01 3.06598857e-01 3.45759429e-01 -1.37428043e-01]
```

## Tf-Idf Weighted Word2Vec Vectorization

### 1. Vectorizing project\_essay

```
In [0]: # Initializing tfidf vectorizer
tfIdfEssayTempVectorizer = TfidfVectorizer();
# Vectorizing preprocessed essays using tfidf vectorizer initialized above
tfIdfEssayTempVectorizer.fit(preProcessedEssaysWithoutStopWords);
# Saving dictionary in which each word is key and it's idf is value
tfIdfEssayDictionary = dict(zip(tfIdfEssayTempVectorizer.get_feature_names(), list(tfIdfEssayTempVectorizer.idf_)));
# Creating set of all unique words used by tfidf vectorizer
tfIdfEssayWords = set(tfIdfEssayTempVectorizer.get_feature_names());
```

```
In [167]: # Creating list to save tf-idf weighted vectors of essays
tfIdfWeightedWord2VecEssaysVectors = [];
# Iterating over each essay
for essay in tqdm(preProcessedEssaysWithoutStopWords):
    # Sum of tf-idf values of all words in a particular essay
    cumulativeSumTfIdfWeightOfEssay = 0;
    # Tf-Idf weighted word2vec vector of a particular essay
    tfIdfWeightedWord2VecEssayVector = np.zeros(300);
    # Splitting essay into list of words
    splittedEssay = essay.split();
    # Iterating over each word
    for word in splittedEssay:
        # Checking if word is in glove words and set of words used by t
        fIdf essay vectorizer
        if (word in gloveWords) and (word in tfIdfEssayWords):
            # Tf-Idf value of particular word in essay
            tfIdfValueWord = tfIdfEssayDictionary[word] * (essay.count(
```

```

word) / len(splittedEssay));
        # Making tf-idf weighted word2vec
        tfIdfWeightedWord2VecEssayVector += tfIdfValueWord * gloveModel[word];
        # Summing tf-idf weight of word to cumulative sum
        cumulativeSumTfIdfWeightOfEssay += tfIdfValueWord;
    if cumulativeSumTfIdfWeightOfEssay != 0:
        # Taking average of sum of vectors with tf-idf cumulative sum
        tfIdfWeightedWord2VecEssayVector = tfIdfWeightedWord2VecEssayVector / cumulativeSumTfIdfWeightOfEssay;
        # Appending the above calculated tf-idf weighted vector of particular essay to list of vectors of essays
        tfIdfWeightedWord2VecEssaysVectors.append(tfIdfWeightedWord2VecEssayVector);

```

In [168]:

```

print("Shape of Tf-Idf weighted Word2Vec vectorization matrix of project essays: {}, {}".format(len(tfIdfWeightedWord2VecEssaysVectors), len(tfIdfWeightedWord2VecEssaysVectors[0])));
equalsBorder(70);
print("Sample Essay: ");
equalsBorder(70);
print(preProcessedEssaysWithoutStopWords[0]);
equalsBorder(70);
print("Tf-Idf Weighted Word2Vec vector of sample essay: ");
equalsBorder(70);
print(tfIdfWeightedWord2VecEssaysVectors[0]);

```

Shape of Tf-Idf weighted Word2Vec vectorization matrix of project essays: 59200, 300  
=====

Sample Essay:

=====  
teacher speech language impaired working students low income urban school district service 70 students moderate severe speech language impairments students difficulty across academic social environments basic tasks take granted including saying good morning participating classroom discussions require practice coaching encouragement students thrive routine love coming school work hard every day overcome disabilities student

we love coming school work hard every day overcome disabilities student s work various speech language goals including receptive expressive language articulation voice fluency pragmatics achieve functional communication students require low tech materials learn others benefit multi sensory materials including computers web based applications art materials budgets tight teacher wish lists thing past new materials hard come students expand communication utilizing new picture cards verbs opposite associations etc ipads visual verbal prompting easel delivery carryover new skills thank consideration donating students nannan

=====

Tf-Idf Weighted Word2Vec vector of sample essay:

=====

```
[ 2.97570068e-02 -4.08024367e-02  6.15747140e-02 -5.06394831e-02
 -7.97055460e-02 -3.05888759e-02 -2.79663198e+00  1.54061875e-01
  8.74659724e-03  1.64910776e-01 -3.03789479e-02  5.51766850e-03
 -1.47102432e-02 -1.14807751e-01 -1.20403637e-01 -1.21715480e-01
  5.25914034e-02 -4.79782095e-03  8.16081569e-02 -2.94487849e-02
 -7.51081486e-02  6.21288181e-02  5.51783228e-02 -4.26891672e-02
 -2.76544808e-02 -1.076777871e-01  1.32945329e-01 -1.86255464e-01
 -9.41385788e-02 -1.52008668e-01 -2.65144143e-01 -9.39512106e-02
  2.16535664e-02  8.99035600e-02 -1.39614484e-01 -8.64725651e-02
  6.74551662e-02 -8.38368630e-02 -6.10851274e-02 -3.59473216e-02
 -3.37853130e-02  1.78586222e-01 -5.71712560e-02 -1.11205530e-01
  1.94448627e-02 -1.76651307e-02  5.00251588e-02  6.67337230e-02
 -7.41338384e-02 -7.72210580e-02 -6.44341665e-02 -4.01887903e-02
  2.24377389e-02  7.00131275e-02  7.19278337e-03 -8.06817087e-02
 -4.53586348e-02  4.97732402e-02 -1.49381260e-01 -8.82242931e-02
 -6.76033109e-02  6.52601200e-03  1.16473560e-01 -1.18711925e-01
 -7.87583324e-02  8.04330941e-02  1.13473218e-02 -1.00295737e-01
  6.77731818e-02  3.07324463e-02 -1.36586329e-01 -4.65168004e-02
  4.92202388e-02 -8.28169927e-02 -1.14693854e-02 -2.61930431e-01
 -7.53518957e-02 -2.68751618e-02  6.30228422e-02  3.66157786e-02
  2.68962672e-02 -4.45577609e-01  9.15501644e-02 -1.23542154e-01
 -7.22232179e-02 -2.37144869e-02  8.80670946e-02 -4.12562037e-02
  1.67544009e-01 -5.66263552e-03 -6.42550263e-02 -1.76757560e-03
  1.72872645e-02 -9.69649796e-02 -1.38905916e-01 -3.52474987e-02
 -2.29210254e+00  1.42588667e-01  3.31955532e-02  2.32817627e-01
 -1.07160971e-01 -6.62423664e-02  3.40573985e-02 -8.64771504e-02
  1.53011362e-01  5.63758795e-02  9.65061920e-02 -1.37557616e-01
                           1.22775577e-01  1.13386469e-01 -8.14579173e-02 -1.89502048e-02
```

$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$
7.80371086e-02	1.53675100e-01	1.18485991e-01	-1.28860960e-02
-1.74554728e-01	1.24357645e-02	6.23928678e-02	-5.49155377e-02
-9.20379966e-02	1.28871258e-02	-9.84488623e-03	-4.05970150e-02
5.59692862e-02	-3.74002687e-02	-5.34245750e-02	-1.54391483e-02
1.68481028e-02	1.11231987e-01	1.99240106e-01	-5.10983515e-02
-2.79988571e-02	-2.60127427e-02	-1.08896041e-01	-5.51923706e-02
1.42509829e-01	8.79783354e-03	1.30751802e-01	1.01105626e-01
8.63631303e-02	2.82712281e-02	5.75996575e-02	-3.25780507e-02
-7.57719427e-02	4.91623901e-02	1.46284559e-01	-8.28706818e-02
2.03526559e-01	-1.17279983e-01	-8.12247920e-03	-1.07132517e-01
-5.80078277e-02	6.20785098e-02	-3.80971708e-02	-3.00488045e-02
9.70286178e-02	-6.69969712e-02	-4.56021257e-02	3.83368336e-03
-1.62824173e-02	2.51130902e-02	6.32253488e-02	3.48912352e-02
-6.89540735e-02	-8.68971762e-02	-7.41942532e-02	2.08668304e-02
-5.55800730e-02	-7.45822268e-02	6.57329899e-02	7.20989362e-03
-2.98979976e-02	-1.24910837e-01	8.80316177e-02	1.24200391e-01
-6.58432059e-02	1.06253719e-02	-1.61665497e-01	4.14927295e-03
-8.83996881e-02	1.76050815e-01	8.37135217e-02	-4.15903453e-02
-1.45716185e-01	9.88005849e-02	1.98388592e-02	-7.61471899e-02
-7.53094079e-02	4.88763287e-02	3.91436491e-02	-3.85405714e-02
-1.05232899e-01	-1.90293910e-02	-9.91677574e-04	-7.44462645e-02
-8.65934708e-02	9.86971183e-02	6.78374014e-02	9.36204140e-02
2.14234242e-01	-7.70856785e-02	-1.54858488e-03	1.65236629e-02
-1.24011380e-01	-4.26764824e-02	9.55782621e-02	-1.15732797e-01
1.20066629e-01	-5.70072426e-02	3.63173955e-04	7.28944890e-02
4.62917856e-02	-7.18106582e-02	-2.12113660e-01	5.90467266e-02
-7.82132936e-02	-1.40361013e-01	-1.79293740e-02	-5.01351134e-02
-1.95352841e-01	-1.48006369e-01	-1.87516801e-01	5.30169276e-03
-1.56339191e+00	2.29586485e-01	-1.45287433e-01	6.77520495e-03
-9.04170866e-02	-1.11840808e-01	-3.68107689e-02	-6.01571902e-02
6.36776524e-02	-1.50558722e-01	-9.61919901e-03	-8.65455168e-02
4.71165558e-02	-1.38114441e-01	2.04658212e-02	1.45081801e-01
-1.06856291e-01	-1.29595610e-02	-1.56463386e-01	-1.17172174e-01
-3.56578145e-02	5.44812382e-02	5.20172721e-02	-1.08192964e-01
8.43958507e-03	-1.16972904e-01	-3.42512000e-02	1.91612272e-01
3.11983009e-02	-5.50160878e-02	4.44807702e-02	4.29839879e-02
-5.05977012e-02	-1.19824821e-01	2.48676405e-01	-1.83036446e-01
-6.76319983e-03	1.25206164e-01	8.62787670e-02	6.56376477e-02
2.14043216e-02	-1.44008965e-01	-8.51691663e-02	-1.46556630e-01

```
[1]: [1] 2.05179477e-02 3.12872307e-02 -4.75955006e-03 9.17673705e-02  
[2] -2.09507006e-01 1.60385239e-01 -2.79961998e-02 7.50958812e-02  
[3] 3.89373123e-02 -1.27010622e-02 -2.91518146e-02 9.40270106e-02  
[4] -1.02535489e-01 -5.54284809e-02 -1.65467495e-02 7.83054192e-02  
[5] -6.75082395e-02 9.84094417e-02 5.87090485e-03 4.50925329e-02  
[6] 8.84219840e-02 1.40194823e-01 6.36338546e-02 -3.88485344e-02  
[7] 4.87648368e-02 -1.12444048e-01 3.89317860e-02 -3.77539583e-02  
[8] -8.25666093e-02 2.35232142e-01 9.56531340e-02 4.66267198e-03]
```

## 2. Vectorizing project\_title

```
In [0]: # Initializing tfidf vectorizer  
tfIdfTitleTempVectorizer = TfidfVectorizer();  
# Vectorizing preprocessed titles using tfidf vectorizer initialized above  
tfIdfTitleTempVectorizer.fit(preProcessedProjectTitlesWithoutStopWords);  
# Saving dictionary in which each word is key and it's idf is value  
tfIdfTitleDictionary = dict(zip(tfIdfTitleTempVectorizer.get_feature_names(), list(tfIdfTitleTempVectorizer.idf_)));  
# Creating set of all unique words used by tfidf vectorizer  
tfIdfTitleWords = set(tfIdfTitleTempVectorizer.get_feature_names());
```

```
In [170]: # Creating list to save tf-idf weighted vectors of project titles  
tfIdfWeightedWord2VecTitlesVectors = [];  
# Iterating over each title  
for title in tqdm(preProcessedProjectTitlesWithoutStopWords):  
    # Sum of tf-idf values of all words in a particular project title  
    cumulativeSumTfIdfWeightOfTitle = 0;  
    # Tf-Idf weighted word2vec vector of a particular project title  
    tfIdfWeightedWord2VecTitleVector = np.zeros(300);  
    # Splitting title into list of words  
    splittedTitle = title.split();  
    # Iterating over each word  
    for word in splittedTitle:  
        # Checking if word is in glove words and set of words used by tfIdf title vectorizer
```

```

        if (word in gloveWords) and (word in tfIdfTitleWords):
            # Tf-Idf value of particular word in title
            tfIdfValueWord = tfIdfTitleDictionary[word] * (title.count(
word) / len(splittedTitle));
            # Making tf-idf weighted word2vec
            tfIdfWeightedWord2VecTitleVector += tfIdfValueWord * gloveM
odel[word];
            # Summing tf-idf weight of word to cumulative sum
            cumulativeSumTfIdfWeightOfTitle += tfIdfValueWord;
        if cumulativeSumTfIdfWeightOfTitle != 0:
            # Taking average of sum of vectors with tf-idf cumulative sum
            tfIdfWeightedWord2VecTitleVector = tfIdfWeightedWord2VecTitleVe
ctor / cumulativeSumTfIdfWeightOfTitle;
            # Appending the above calculated tf-idf weighted vector of particul
ar title to list of vectors of project titles
            tfIdfWeightedWord2VecTitlesVectors.append(tfIdfWeightedWord2VecTitl
eVector);

```

In [171]:

```

print("Shape of Tf-Idf weighted Word2Vec vectorization matrix of projec
t titles: {}, {}".format(len(tfIdfWeightedWord2VecTitlesVectors), len(t
fIdfWeightedWord2VecTitlesVectors[0])));
equalsBorder(70);
print("Sample Title: ");
equalsBorder(70);
print(preProcessedProjectTitlesWithoutStopWords[0]);
equalsBorder(70);
print("Tf-Idf Weighted Word2Vec vector of sample title: ");
equalsBorder(70);
print(tfIdfWeightedWord2VecTitlesVectors[0]);

```

Shape of Tf-Idf weighted Word2Vec vectorization matrix of project title
s: 59200, 300  
=====

Sample Title:

=====

```
help students achieve communication goals new supplies
```

=====

Tf-Idf Weighted Word2Vec vector of sample title:

```
=====
[ 3.43431605e-02 2.32133872e-01 5.46992729e-02 2.62744552e-02
-1.02905875e-01 -3.62331661e-02 -3.29133147e+00 1.93263154e-01
2.05205562e-01 3.12260111e-01 -2.48575596e-01 7.45463691e-02
-9.89481608e-02 -3.02114382e-01 -2.14039742e-01 6.50962125e-02
-1.57534553e-02 -1.95919490e-01 6.33712612e-02 -2.29365318e-01
5.37773428e-02 -2.62057852e-01 1.75713159e-01 -5.17775608e-02
1.44533605e-01 -1.03064103e-01 3.44732729e-01 4.37619247e-02
-2.26767693e-01 -1.69552806e-01 -3.34539369e-01 -1.54053661e-01
-8.76380555e-02 -1.66111807e-01 1.75586392e-01 -5.50606968e-02
-4.59361262e-02 -1.01533366e-01 -1.66117835e-01 -1.72373415e-01
-4.29131372e-02 1.03933927e-01 1.76387075e-01 -9.78060367e-02
-9.25308932e-03 -1.20003985e-02 -2.88839806e-02 1.52067898e-01
-1.66186577e-01 -1.60106287e-01 1.70730167e-01 -1.25009665e-01
-2.70889958e-02 4.38286360e-02 8.97206312e-02 7.53328204e-02
-1.10999289e-01 -1.35816127e-01 -7.13295685e-02 1.02674104e-01
-7.49755570e-02 -3.53524290e-02 -9.65811361e-02 -1.55034765e-01
4.66861678e-03 1.26715978e-01 1.99442911e-02 -9.99268369e-02
-7.71452368e-02 -3.92767383e-02 -2.93731836e-01 -9.49183744e-02
-6.55479645e-02 -2.67649224e-01 -2.67903676e-01 -3.16112607e-01
-2.05593938e-02 -5.17653134e-04 -2.92287443e-02 1.63047675e-02
1.65546048e-01 -4.94570140e-01 3.91362801e-03 1.68910858e-02
-1.05524140e-01 -9.87565637e-02 3.40490774e-02 -4.31717970e-02
2.79024437e-01 2.42811682e-03 1.00131725e-01 -1.62581748e-01
-2.96715419e-02 8.42788375e-02 1.30559455e-01 -5.16622652e-02
-2.72080472e+00 -5.12740202e-02 1.03698586e-01 1.24307525e-01
-1.14391400e-01 -6.62362852e-02 2.81989874e-01 5.83731910e-02
3.29072383e-01 6.26085097e-02 -1.36020996e-01 1.43862208e-01
1.04758535e-01 6.68401514e-02 -2.04593451e-01 -1.78882939e-01
-1.49301207e-02 2.61187798e-01 -5.32094863e-03 -2.79265842e-01
-4.31987088e-01 1.23995664e-01 7.81708685e-02 2.50929110e-02
-8.10892152e-02 -1.20513814e-01 -2.45190246e-01 -1.34722910e-01
-1.44049326e-02 1.29280641e-01 3.98609393e-02 4.05402368e-02
9.42418171e-02 -2.22405620e-03 8.98827718e-02 -7.39794505e-02
3.17377176e-02 3.60291896e-02 -1.59337566e-01 -7.45835038e-03
-7.73905194e-02 1.53733549e-01 4.81775521e-02 2.02321874e-01
-7.22821432e-03 -2.55089325e-03 7.15331491e-02 -1.57427691e-01
1.51967448e-01 5.68949725e-02 8.11004562e-02 9.53494554e-04
1.07815396e-01 9.80485604e-02 -4.67030861e-02 -3.82289044e-02
```

3.16251351e-01	7.08727677e-02	-2.22958028e-01	1.31704767e-01
2.03362226e-01	-1.01964442e-01	1.48931518e-01	-5.32437159e-01
1.86765082e-01	5.00701725e-02	-9.67291789e-02	-1.77104628e-01
-2.23163607e-01	8.52998532e-02	-1.57192270e-01	1.13217809e-01
2.45594762e-01	-2.84310485e-01	-6.48835795e-02	1.63509901e-01
3.31005668e-02	-2.74363606e-01	-3.70973104e-02	-7.25260907e-02
-1.34315823e-01	-4.68521034e-02	-2.76809435e-02	6.41367187e-02
7.61448287e-02	1.32235014e-01	-1.73900016e-01	9.62455800e-02
-1.27659168e-02	8.86739405e-03	4.44329135e-02	1.22480734e-01
-1.18461701e-01	5.97856743e-02	-1.26982952e-01	1.66958784e-01
-1.32512577e-01	-1.52382322e-01	1.41021869e-01	-2.16582302e-01
-1.93530253e-02	2.41685302e-01	-4.30957296e-02	2.64032314e-01
1.98536777e-01	-2.12559824e-01	-1.45326851e-01	1.60490138e-01
-2.03430262e-01	-1.80526380e-02	1.76620794e-01	-2.23528577e-02
1.95299795e-03	-6.80914872e-02	-6.91115693e-02	-7.65323554e-02
-7.37914227e-02	-2.67589797e-01	-2.53794373e-01	1.38650233e-01
1.60773366e-02	-2.07776087e-01	-5.39636049e-02	4.86431694e-02
-4.31646724e-01	-1.97162970e-01	-2.04391710e-01	-2.19731746e-01
-1.68846011e+00	3.44355129e-01	-1.62340797e-01	5.93206135e-02
-5.84586951e-02	-6.65633167e-02	1.16909780e-02	-2.18750646e-01
-1.00833986e-01	-4.00552226e-03	4.04272229e-02	3.02110780e-01
1.34562619e-01	-4.86281024e-02	1.34598833e-01	6.48676513e-02
3.28135384e-02	-4.28528971e-02	-3.39925533e-01	-1.76873381e-01
-1.35730167e-01	-2.01942499e-02	2.86311017e-02	-1.50596585e-01
-6.26448042e-02	-2.24008595e-01	-1.51068360e-02	5.19497898e-03
6.68559281e-02	-2.28623352e-01	1.15165999e-01	-7.39084467e-02
2.09552162e-01	6.94616487e-02	5.03225104e-02	-1.63883106e-01
2.09514881e-01	7.40297570e-02	8.44035585e-02	4.22629017e-02
5.16906868e-03	-2.82154613e-02	6.74318925e-02	-6.25753852e-02
1.67334715e-02	-3.63016698e-02	-1.59696780e-01	-8.23004772e-03
-6.68113971e-02	9.59830427e-02	1.10052726e-01	2.39849565e-01
-1.11705486e-01	-2.38074168e-01	5.32335655e-03	1.45109304e-01
-2.87235926e-01	-4.17143530e-02	1.93520321e-02	-4.89604101e-02
1.20851153e-01	2.39350883e-01	1.54110234e-02	7.75699070e-02
-6.10967942e-02	1.73882887e-02	1.38814382e-01	-1.44495206e-01
-5.51299442e-03	-9.36892958e-02	-2.53442572e-02	-1.01460385e-01
-3.43451670e-01	3.30964424e-01	3.02761112e-01	-1.37075513e-01]

## Vectorizing numerical features

### 1. Vectorizing price

```
In [0]: # Standardizing the price data using StandardScaler(Uses mean and std f  
or standardization)  
priceScaler = StandardScaler();  
priceScaler.fit(trainingData['price'].values.reshape(-1, 1));  
priceStandardized = priceScaler.transform(trainingData['price'].values.  
reshape(-1, 1));
```

```
In [173]: print("Shape of standardized matrix of prices: ", priceStandardized.sha  
pe);  
equalsBorder(70);  
print("Sample original prices: ");  
equalsBorder(70);  
print(trainingData['price'].values[0:5]);  
print("Sample standardized prices: ");  
equalsBorder(70);  
print(priceStandardized[0:5]);
```

Shape of standardized matrix of prices: (59200, 1)

=====

Sample original prices:

=====

[497.25 60.06 429. 12.52 535.07]

Sample standardized prices:

=====

[[ 0.44688774]  
[-0.69116967]  
[ 0.26922489]  
[-0.81492193]  
[ 0.54533768]]

### 2. Vectorizing quantity

```
In [0]: # Standardizing the quantity data using StandardScaler(Uses mean and std for standardization)
quantityScaler = StandardScaler();
quantityScaler.fit(trainingData['quantity'].values.reshape(-1, 1));
quantityStandardized = quantityScaler.transform(trainingData['quantity'].values.reshape(-1, 1));
```

```
In [175]: print("Shape of standardized matrix of quantities: ", quantityStandardized.shape);
equalsBorder(70);
print("Sample original quantities: ");
equalsBorder(70);
print(trainingData['quantity'].values[0:5]);
print("Sample standardized quantities: ");
equalsBorder(70);
print(quantityStandardized[0:5]);
```

Shape of standardized matrix of quantities: (59200, 1)

=====

Sample original quantities:

=====

[ 9 6 1 17 4]

Sample standardized quantities:

=====

```
[[ -0.3294401 ]
 [-0.4333539 ]
 [-0.60654357]
 [-0.05233663]
 [-0.50262976]]
```

### 3. Vectorizing teacher\_number\_of\_previously\_posted\_projects

```
In [0]: # Standardizing the teacher_number_of_previously_posted_projects data using StandardScaler(Uses mean and std for standardization)
previouslyPostedScaler = StandardScaler();
previouslyPostedScaler.fit(trainingData['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1));
previouslyPostedStandardized = previouslyPostedScaler.transform(trainin
```

```
gData['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1));
```

```
In [177]: print("Shape of standardized matrix of teacher_number_of_previously_pos  
ted_projects: ", previouslyPostedStandardized.shape);  
equalsBorder(70);  
print("Sample original quantities: ");  
equalsBorder(70);  
print(trainingData['teacher_number_of_previously_posted_projects'].valu  
es[0:5]);  
print("Sample standardized teacher_number_of_previously_posted_project  
s: ");  
equalsBorder(70);  
print(previouslyPostedStandardized[0:5]);
```

```
Shape of standardized matrix of teacher_number_of_previously_posted_pro  
jects: (59200, 1)  
=====  
Sample original quantities:  
=====  
[3 1 8 2 5]  
Sample standardized teacher_number_of_previously_posted_projects:  
=====  
[[ -0.26854127]  
 [-0.35097142]  
 [-0.0624659 ]  
 [-0.30975635]  
 [-0.18611112]]
```

```
In [0]: numberPoints = previouslyPostedStandardized.shape[0];  
# Categorical data  
categoriesVectorsSub = categoriesVectors[0:numberPoints];  
subCategoriesVectorsSub = subCategoriesVectors[0:numberPoints];  
teacherPrefixVectorsSub = teacherPrefixVectors[0:numberPoints];  
schoolStateVectorsSub = schoolStateVectors[0:numberPoints];  
projectGradeVectorsSub = projectGradeVectors[0:numberPoints];  
  
# Text data  
bowEssayModelSub = bowEssayModel[0:numberPoints];
```

```

bowTitleModelSub = bowTitleModel[0:numberOfPoints];
tfIdfEssayModelSub = tfIdfEssayModel[0:numberOfPoints];
tfIdfTitleModelSub = tfIdfTitleModel[0:numberOfPoints];
word2VecEssaysVectorsSub = word2VecEssaysVectors[0:numberOfPoints];
word2VecTitlesVectorsSub = word2VecTitlesVectors[0:numberOfPoints];
tfIdfWeightedWord2VecEssaysVectorsSub = tfIdfWeightedWord2VecEssaysVectors[0:numberOfPoints];
tfIdfWeightedWord2VecTitlesVectorsSub = tfIdfWeightedWord2VecTitlesVectors[0:numberOfPoints];

# Numerical data
priceStandardizedSub = priceStandardized[0:numberOfPoints];
quantityStandardizedSub = quantityStandardized[0:numberOfPoints];
previouslyPostedStandardizedSub = previouslyPostedStandardized[0:numberOfPoints];

```

```

In [0]: def getAvgTfIdfEssayVectors(arrayOfTexts):
    # Creating list to save tf-idf weighted vectors of essays
    tfIdfWeightedWord2VecEssaysVectors = [];
    # Iterating over each essay
    for essay in tqdm(arrayOfTexts):
        # Sum of tf-idf values of all words in a particular essay
        cumulativeSumTfIdfWeight0fEssay = 0;
        # Tf-Idf weighted word2vec vector of a particular essay
        tfIdfWeightedWord2VecEssayVector = np.zeros(300);
        # Splitting essay into list of words
        splittedEssay = essay.split();
        # Iterating over each word
        for word in splittedEssay:
            # Checking if word is in glove words and set of words used
            by tfIdf essay vectorizer
            if (word in gloveWords) and (word in tfIdfEssayWords):
                # Tf-Idf value of particular word in essay
                tfIdfValueWord = tfIdfEssayDictionary[word] * (essay.co
                unt(word) / len(splittedEssay));
                # Making tf-idf weighted word2vec
                tfIdfWeightedWord2VecEssayVector += tfIdfValueWord * gl
                oveModel[word];
                # Summing tf-idf weight of word to cumulative sum

```

```

        cumulativeSumTfIdfWeightOfEssay += tfIdfValueWord;
    if cumulativeSumTfIdfWeightOfEssay != 0:
        # Taking average of sum of vectors with tf-idf cumulative sum
        tfIdfWeightedWord2VecEssayVector = tfIdfWeightedWord2VecEssayVector / cumulativeSumTfIdfWeightOfEssay;
        # Appending the above calculated tf-idf weighted vector of particular essay to list of vectors of essays
        tfIdfWeightedWord2VecEssaysVectors.append(tfIdfWeightedWord2VecEssayVector);
    return tfIdfWeightedWord2VecEssaysVectors;

```

In [0]:

```

def getAvgTfIdfTitleVectors(arrayOfTexts):
    # Creating list to save tf-idf weighted vectors of project titles
    tfIdfWeightedWord2VecTitlesVectors = [];
    # Iterating over each title
    for title in tqdm(arrayOfTexts):
        # Sum of tf-idf values of all words in a particular project title
        cumulativeSumTfIdfWeightOfTitle = 0;
        # Tf-Idf weighted word2vec vector of a particular project title
        tfIdfWeightedWord2VecTitleVector = np.zeros(300);
        # Splitting title into list of words
        splittedTitle = title.split();
        # Iterating over each word
        for word in splittedTitle:
            # Checking if word is in glove words and set of words used by tfIdf title vectorizer
            if (word in gloveWords) and (word in tfIdfTitleWords):
                # Tf-Idf value of particular word in title
                tfIdfValueWord = tfIdfTitleDictionary[word] * (title.count(word) / len(splittedTitle));
                # Making tf-idf weighted word2vec
                tfIdfWeightedWord2VecTitleVector += tfIdfValueWord * gloveModel[word];
                # Summing tf-idf weight of word to cumulative sum
                cumulativeSumTfIdfWeightOfTitle += tfIdfValueWord;
            if cumulativeSumTfIdfWeightOfTitle != 0:
                # Taking average of sum of vectors with tf-idf cumulative sum

```

```

um
        tfIdfWeightedWord2VecTitleVector = tfIdfWeightedWord2VecTit
leVector / cumulativeSumTfIdfWeightOfTitle;
        # Appending the above calculated tf-idf weighted vector of part
        #icular title to list of vectors of project titles
        tfIdfWeightedWord2VecTitlesVectors.append(tfIdfWeightedWord2Vec
TitleVector);
    return tfIdfWeightedWord2VecTitlesVectors;

```

In [6]:

```
kFoldResultsDataFrame = pd.DataFrame(columns = ['Vectorizer', 'Model',
'Hyper Parameter - K', 'AUC']);
kFoldResultsDataFrame
```

Out[6]:

Vectorizer	Model	Hyper Parameter - K	AUC

## Preparing Test data for analysis

In [182]:

```

# Test data categorical features transformation
categoriesTransformedTestData = subjectsCategoriesVectorizer.transform(
testData['cleaned_categories']);
subCategoriesTransformedTestData = subjectsSubCategoriesVectorizer.tran
sform(testData['cleaned_sub_categories']);
teacherPrefixTransformedTestData = teacherPrefixVectorizer.transform(te
stData['teacher_prefix']);
schoolStateTransformedTestData = schoolStateVectorizer.transform(testDa
ta['school_state']);
projectGradeTransformedTestData = projectGradeVectorizer.transform(test
Data['project_grade_category']);

# Test data text features transformation
preProcessedEssaysTemp = preProcessingWithAndWithoutStopWords(testData[
'project_essay'])[1];
preProcessedTitlesTemp = preProcessingWithAndWithoutStopWords(testData[
'project_title'])[1];
bowEssayTransformedTestData = bowEssayVectorizer.transform(preProcessed
EssaysTemp);

```

```
bowTitleTransformedTestData = bowTitleVectorizer.transform(preProcessedTitlesTemp);
tfIdfEssayTransformedTestData = tfIdfEssayVectorizer.transform(preProcessedEssaysTemp);
tfIdfTitleTransformedTestData = tfIdfTitleVectorizer.transform(preProcessedTitlesTemp);
word2VecEssayTransformedTestData = getWord2VecVectors(preProcessedEssaysTemp);
word2VecTitleTransformedTestData = getWord2VecVectors(preProcessedTitlesTemp);
tfIdfWeightedEssayTransformedTestData = getAvgTfIdfEssayVectors(preProcessedEssaysTemp);
tfIdfWeightedTitleTransformedTestData = getAvgTfIdfTitleVectors(preProcessedTitlesTemp);

# Test data numerical features transformation
priceTransformedTestData = priceScaler.transform(testData['price'].values.reshape(-1, 1));
quantityTransformedTestData = quantityScaler.transform(testData['quantity'].values.reshape(-1, 1));
previouslyPostedTransformedTestData = previouslyPostedScaler.transform(testData['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1));
```

## Analysis of data containing BOW & Tf-Idf vectorized text features using K-NN

In [0]: testKValues = np.arange(1, 40, 2);

```

techniques = ['Bag of Words', 'Tf-Idf'];
for index, technique in enumerate(techniques):
    areaUnderRocValuesTrain = [];
    trainingMergedData = hstack((categoriesVectorsSub,\n                                subCategoriesVectorsSub,\n                                teacherPrefixVectorsSub,\n                                schoolStateVectorsSub,\n                                projectGradeVectorsSub,\n                                priceStandardizedSub,\n                                previouslyPostedStandardizedSub));
    testMergedData = hstack((categoriesTransformedTestData,\n                            subCategoriesTransformedTestD\n                            ata,\n                            teacherPrefixTransformedTestD\n                            ata,\n                            schoolStateTransformedTestDat\n                            a,\n                            projectGradeTransformedTestDa\n                            ta,\n                            priceTransformedTestData,\n                            previouslyPostedTransformedTe\n                            stData));

    if(index == 0):
        trainingMergedData = hstack((trainingMergedData,\n                                    bowTitleModelSub,\n                                    bowEssayModelSub));
        testMergedData = hstack((testMergedData,\n                                bowTitleTransformedTestData,\n                                bowEssayTransformedTestData));
    elif(index == 1):
        trainingMergedData = hstack((trainingMergedData,\n                                    tfIdfTitleModelSub,\n                                    tfIdfEssayModelSub));
        testMergedData = hstack((testMergedData,\n                                tfIdfTitleTransformedTestData,\n                                tfIdfEssayTransformedTestData));
    for testKValue in tqdm(testKValues):
        knnClassifier = KNeighborsClassifier(n_neighbors = testKValue,

```

```

algorithm = 'brute');
scores = cross_val_score(knnClassifier, trainingMergedData, classesTraining, cv = 5, scoring = 'roc_auc');
areaUnderRocValuesTrain.append(np.array(scores).mean());

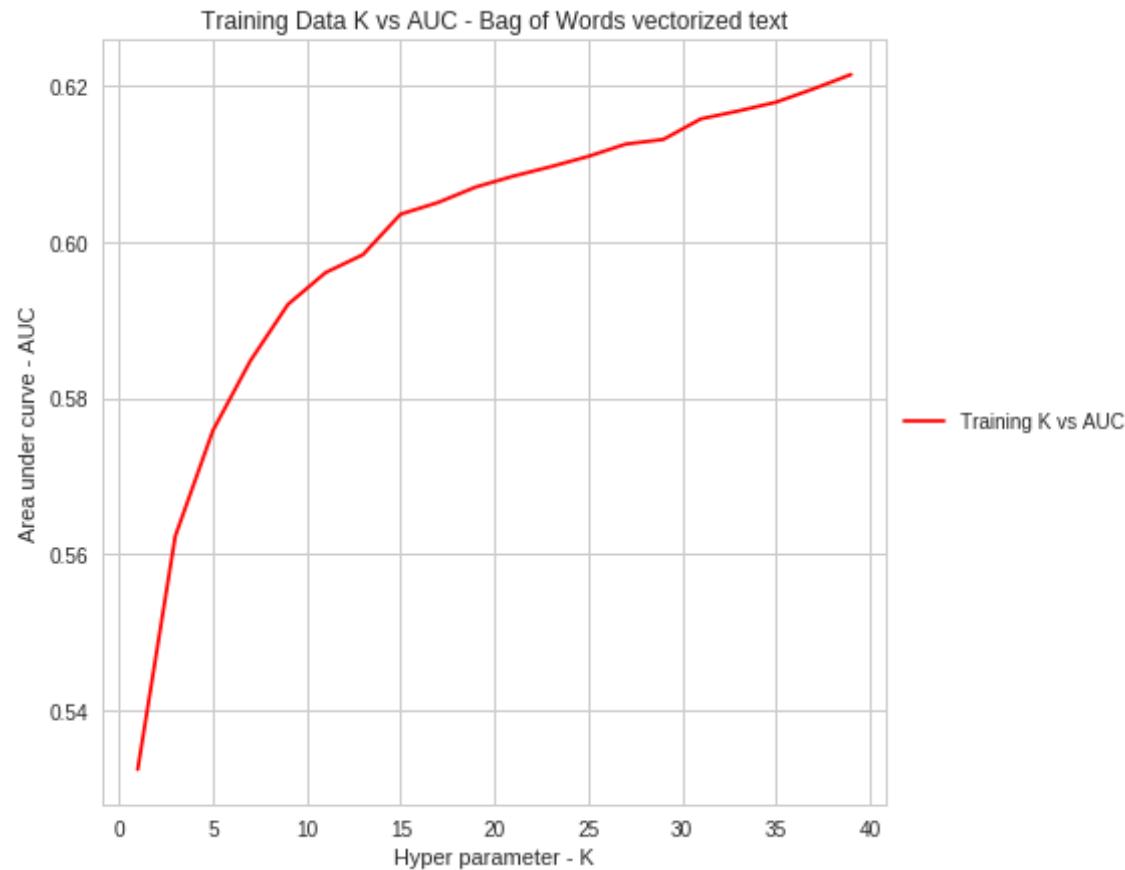
plt.plot(testKValues, areaUnderRocValuesTrain, 'r', label = "Training K vs AUC");
plt.title("Training Data K vs AUC - {} vectorized text".format(technique));
plt.xlabel("Hyper parameter - K");
plt.ylabel("Area under curve - AUC");
plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
plt.show();

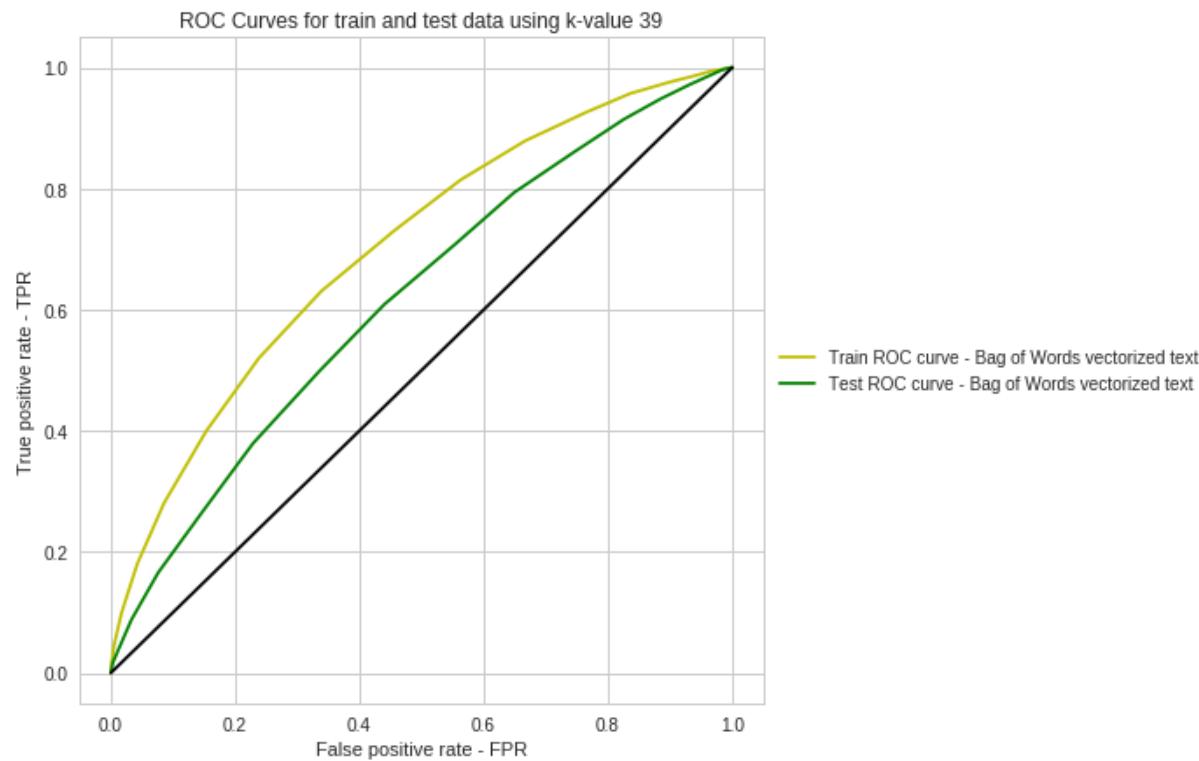
optimalKValue = testKValues[np.argmax(areaUnderRocValuesTrain)];
knnClassifier = KNeighborsClassifier(n_neighbors = optimalKValue, algorithm = 'brute');
knnClassifier.fit(trainingMergedData, classesTraining);
predProbScoresTraining = knnClassifier.predict_proba(trainingMergedData);
fprTrain, tprTrain, thresholdTrain = roc_curve(classesTraining, predProbScoresTraining[:, 1]);
predProbScoresTest = knnClassifier.predict_proba(testMergedData);
fprTest, tprTest, thresholdTest = roc_curve(classesTest, predProbScoresTest[:, 1]);
areaUnderRocValueTest = auc(fprTest, tprTest);
plt.plot(fprTrain, tprTrain, 'y', label="Train ROC curve - {} vectorized text".format(technique));
plt.plot(fprTest, tprTest, 'g', label="Test ROC curve - {} vectorized text".format(technique));
plt.plot([0, 1], [0, 1], 'k-');
plt.title("ROC Curves for train and test data using k-value {}".format(optimalKValue))
plt.xlabel('False positive rate - FPR');
plt.ylabel('True positive rate - TPR');
plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
plt.show();

equalsBorder(100);

```

```
        equalsBorder(100);
        print("Results of analysis using {} vectorized text features merged
with other features using K-NN brute force algorithm:".format(techniqu
e));
        equalsBorder(70);
        print("AUC values of train data: ");
        equalsBorder(40);
        print(areaUnderRocValuesTrain);
        equalsBorder(40);
        print("Optimal K-Value: ", optimalKValue);
        equalsBorder(40);
        print("AUC value of test data: ", areaUnderRocValueTest);
# Predicting classes of test data projects
predictionClassesTest = knnClassifier.predict(testMergedData);
equalsBorder(40);
# Printing confusion matrix
confusionMatrix = confusion_matrix(classesTest, predictionClassesTe
st);
# Creating dataframe for generated confusion matrix
confusionMatrixDataFrame = pd.DataFrame(data = confusionMatrix, ind
ex = ['Actual: NO', 'Actual: YES'], columns = ['Predicted: NO', 'Predic
ted: YES']);
print("Confusion Matrix : ");
equalsBorder(60);
sbrn.heatmap(confusionMatrixDataFrame, annot = True, fmt = 'd');
plt.show();
# Adding results to results dataframe
kFoldResultsDataFrame = kFoldResultsDataFrame.append({'Vectorizer':
technique, 'Model': 'Brute', 'Hyper Parameter - K': optimalKValue, 'AU
C': areaUnderRocValueTest}, ignore_index = True);
```





=====  
=====  
=====  
=====  
Results of analysis using Bag of Words vectorized text features merged  
with other features using K-NN brute force algorithm:  
=====

=====  
AUC values of train data:

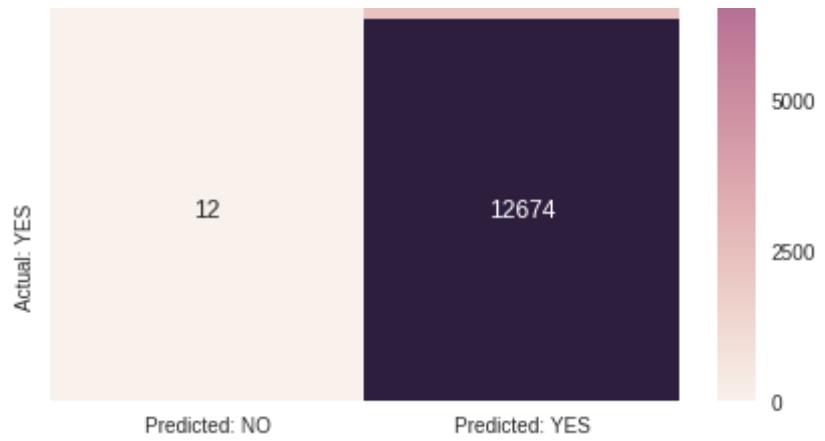
```
=====
[0.5324718468468468, 0.5623096846846847, 0.5757869588338339, 0.58473662
72522523, 0.5919456331331331, 0.5960126063563564, 0.5983142673923922,
0.6034868149399399, 0.6049877533783785, 0.6069565033783784, 0.608351695
4454453, 0.609563501001001, 0.6108810529279279, 0.6124558621121122, 0.6
130490646896897, 0.6156833708708709, 0.6166933965215217, 0.617812718968
9689, 0.6195215840840841, 0.6213565909659658]
=====
```

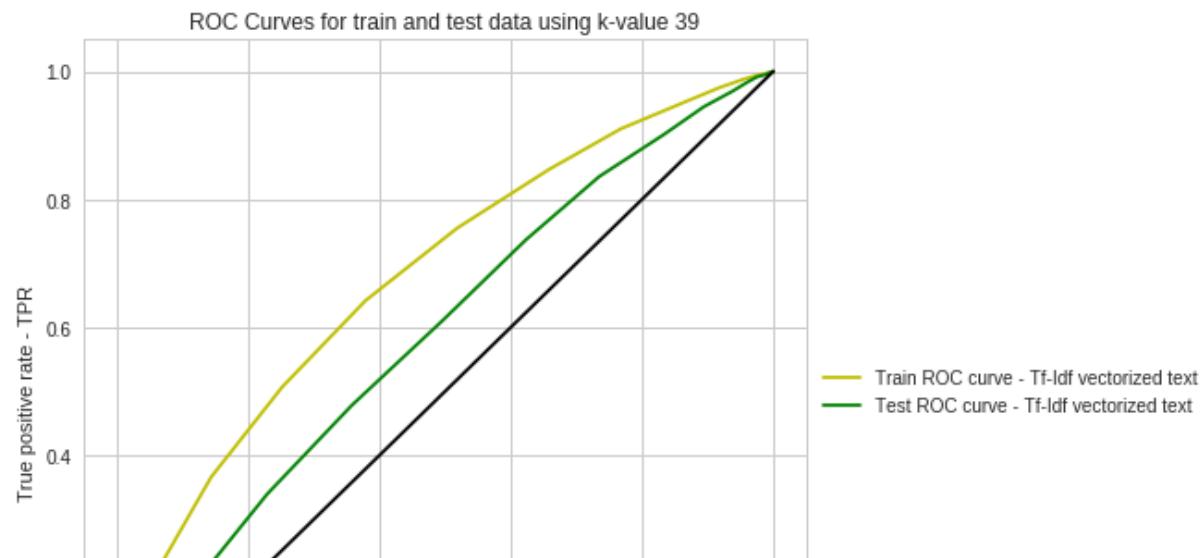
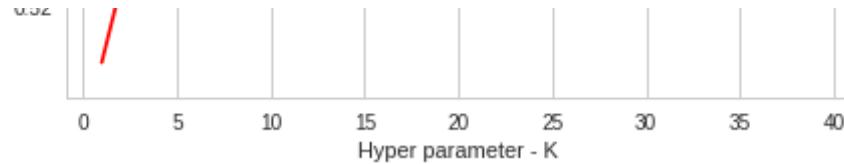
```
Optimal K-Value: 39
=====
```

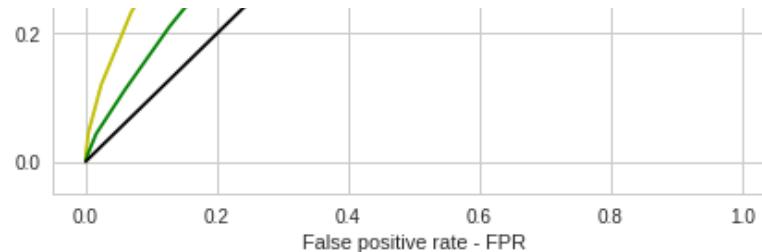
```
AUC value of test data: 0.6182960214071658
=====
```

```
Confusion Matrix :
=====
```









```
=====
=====
=====
=====
Results of analysis using Tf-Idf vectorized text features merged with other features using K-NN brute force algorithm:
```

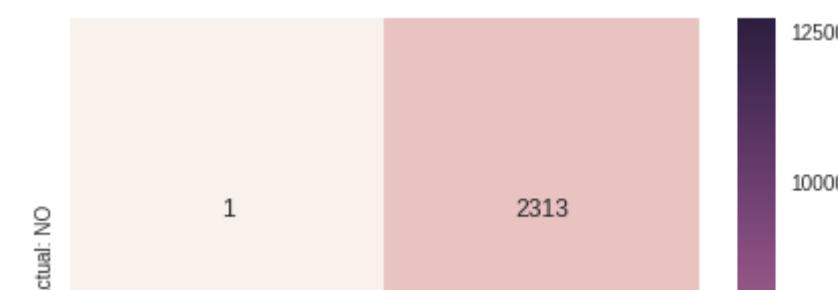
```
=====
AUC values of train data:
```

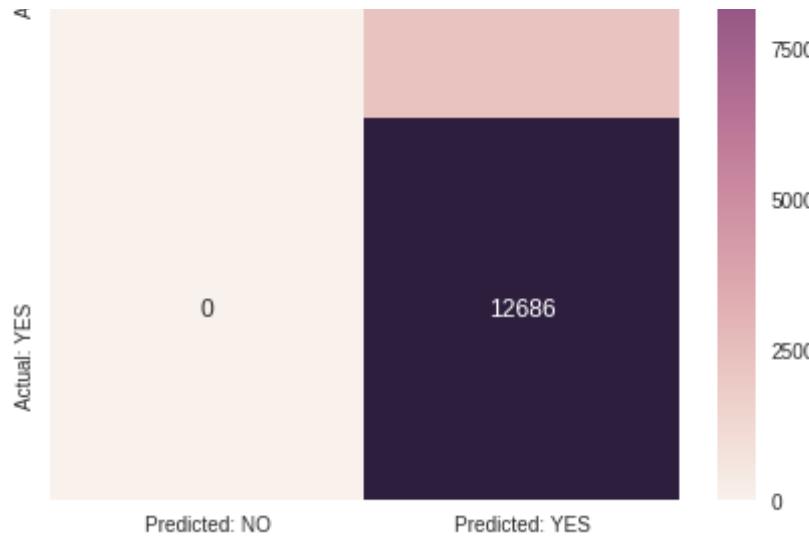
```
=====
[0.5139164164164164, 0.5296537162162162, 0.5383975225225225, 0.5484509196696696, 0.5542950137637638, 0.5595634071571572, 0.564910817067067, 0.5672973598598599, 0.5712998310810812, 0.5733333646146146, 0.5754397678928929, 0.5778523523523524, 0.5792881631631631, 0.5818879348098098, 0.5838870589339338, 0.584635354104104, 0.5834029341841842, 0.584898023023023, 0.5857724912412412, 0.5866925362862863]
```

```
=====
Optimal K-Value: 39
```

```
=====
AUC value of test data: 0.5887691411094189
```

```
=====
Confusion Matrix :
```





## Analysis of data containing average Word2Vec vectorized text features using K-NN

```
In [0]: testKValues = [1, 5, 10, 30, 50, 80]
techniques = ['Average Word2Vec', 'Tf-Idf Weighted Word2Vec'];
for index, technique in enumerate(techniques):
    areaUnderRocValuesTrain = []
    trainingMergedData = hstack((categoriesVectorsSub,\
                                 subCategoriesVectorsSub,\
                                 teacherPrefixVectorsSub,\
                                 schoolStateVectorsSub,\
                                 projectGradeVectorsSub,\
                                 priceStandardizedSub,\
                                 previouslyPostedStandardizedSub));
    testMergedData = hstack((categoriesTransformedTestData,\
                             subCategoriesTransformedTestD
ata,\\
ata,\\
teacherPrefixTransformedTestD
```

```

a,\                                         schoolStateTransformedTestDat
ta,\                                         projectGradeTransformedTestDa
stData));                                         priceTransformedTestData,\_
                                         previouslyPostedTransformedTe
    if(index == 0):
        trainingMergedData = hstack((trainingMergedData,\_
                                         word2VecTitlesVectorsSub,\_
                                         word2VecEssaysVectorsSub));
    testMergedData = hstack((testMergedData,\_
                                         word2VecTitleTransformedTestData,\_
                                         word2VecEssayTransformedTestData));
    elif(index == 1):
        trainingMergedData = hstack((trainingMergedData,\_
                                         tfIdfWeightedWord2VecTitlesVectors
Sub,\                                         tfIdfWeightedWord2VecEssaysVectors
Sub));
        testMergedData = hstack((testMergedData,\_
                                         tfIdfWeightedTitleTransformedTestData,
\                                         tfIdfWeightedEssayTransformedTestData
));
    for testKValue in tqdm(testKValues):
        knnClassifier = KNeighborsClassifier(n_neighbors = testKValue,
algorithm = 'brute');
        scores = cross_val_score(knnClassifier, trainingMergedData, cla
ssesTraining, cv = 5, scoring = 'roc_auc');
        areaUnderRocValuesTrain.append(np.array(scores).mean());

        plt.plot(testKValues, areaUnderRocValuesTrain, 'r', label = "Trains
ing K vs AUC");
        plt.title("Training Data K vs AUC - {} vectorized text".format(tech
nique));
        plt.xlabel("Hyper parameter - K");
        plt.ylabel("Area under curve - AUC");
        plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))

```

```

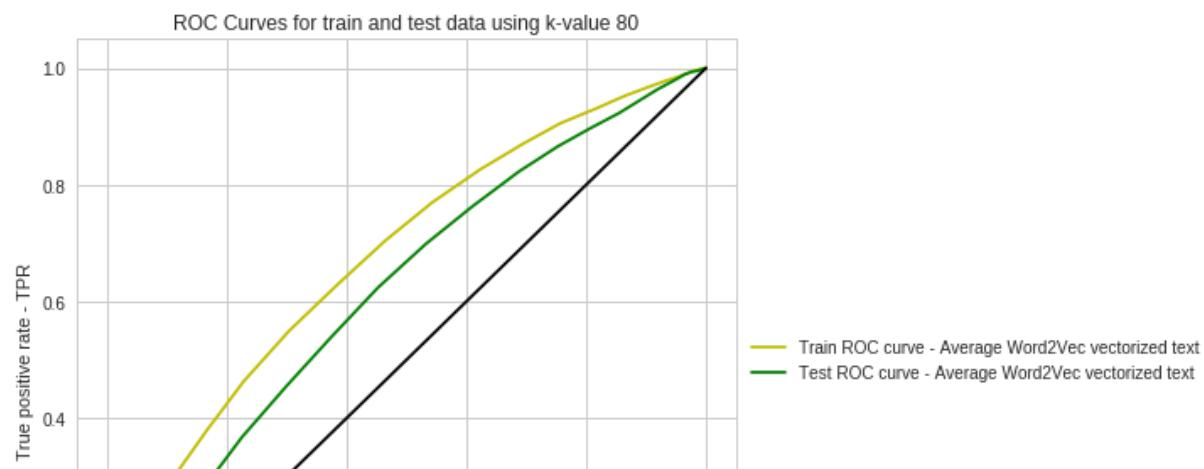
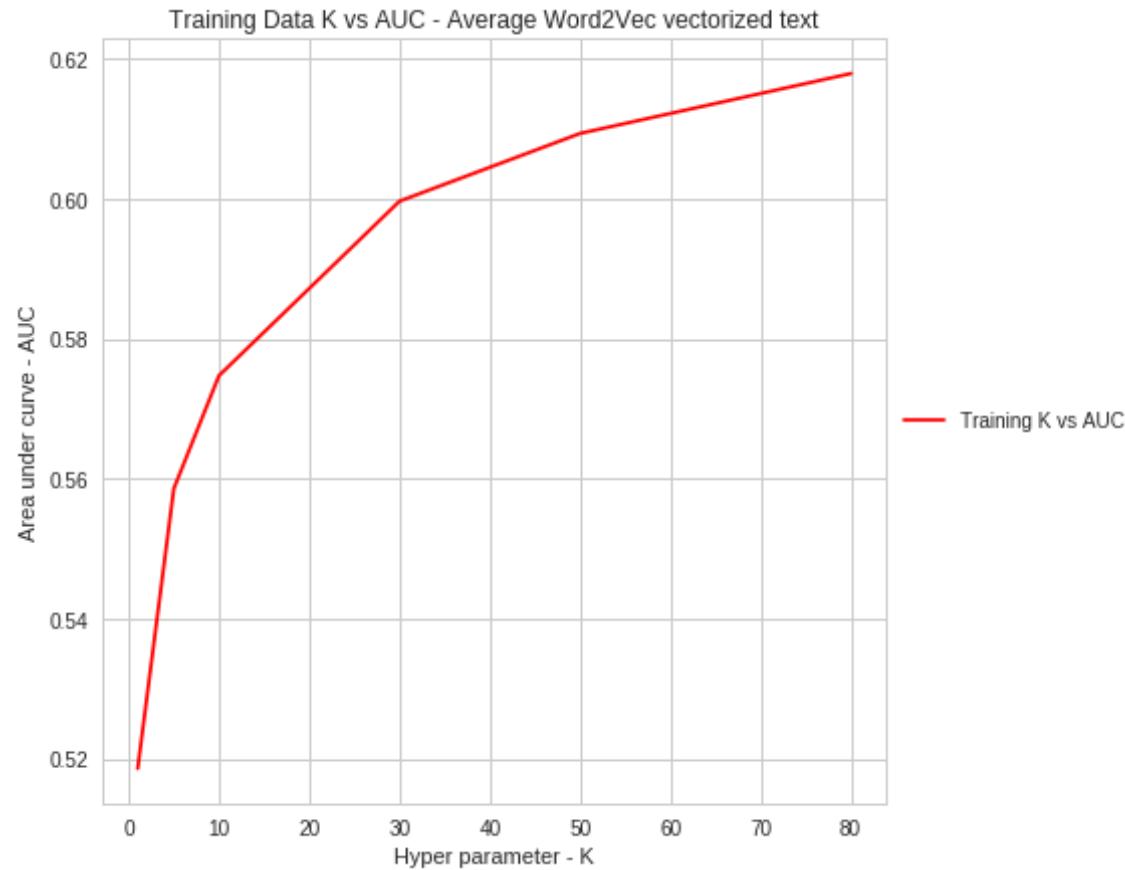
plt.show();

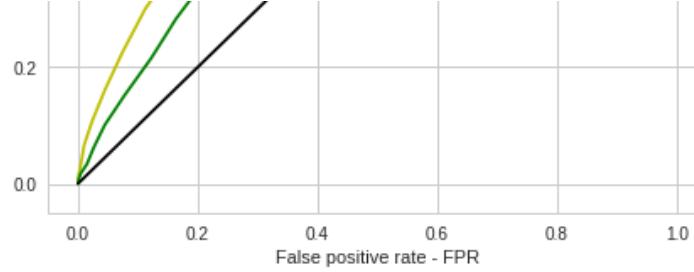
optimalKValue = testKValues[np.argmax(areaUnderRocValuesTrain)];
knnClassifier = KNeighborsClassifier(n_neighbors = optimalKValue, algorithm = 'brute');
knnClassifier.fit(trainingMergedData, classesTraining);
predProbScoresTraining = knnClassifier.predict_proba(trainingMergedData);
fprTrain, tprTrain, thresholdTrain = roc_curve(classesTraining, predProbScoresTraining[:, 1]);
predProbScoresTest = knnClassifier.predict_proba(testMergedData);
fprTest, tprTest, thresholdTest = roc_curve(classesTest, predProbScoresTest[:, 1]);
areaUnderRocValueTest = auc(fprTest, tprTest);
plt.plot(fprTrain, tprTrain, 'y', label="Train ROC curve - {} vectorized text".format(technique));
plt.plot(fprTest, tprTest, 'g', label="Test ROC curve - {} vectorized text".format(technique));
plt.plot([0, 1], [0, 1], 'k-');
plt.title("ROC Curves for train and test data using k-value {}".format(optimalKValue));
plt.xlabel('False positive rate - FPR');
plt.ylabel('True positive rate - TPR');
plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
plt.show();

equalsBorder(100);
equalsBorder(100);
print("Results of analysis using {} vectorized text features merged with other features using K-NN brute force algorithm:".format(technique));
equalsBorder(70);
print("AUC values of train data: ");
equalsBorder(40);
print(areaUnderRocValuesTrain);
equalsBorder(40);
print("Optimal K-Value: ", optimalKValue);
equalsBorder(40);
print("AUC value of test data: ", areaUnderRocValueTest);

```

```
# Predicting classes of test data projects
predictionClassesTest = knnClassifier.predict(testMergedData);
equalsBorder(40);
# Printing confusion matrix
confusionMatrix = confusion_matrix(classesTest, predictionClassesTe
st);
# Creating dataframe for generated confusion matrix
confusionMatrixDataFrame = pd.DataFrame(data = confusionMatrix, ind
ex = ['Actual: NO', 'Actual: YES'], columns = ['Predicted: NO', 'Predic
ted: YES']);
print("Confusion Matrix : ");
equalsBorder(60);
sbrn.heatmap(confusionMatrixDataFrame, annot = True, fmt = 'd');
plt.show();
# Adding results to results dataframe
kFoldResultsDataFrame = kFoldResultsDataFrame.append({'Vectorizer':
technique, 'Model': 'Brute', 'Hyper Parameter - K': optimalKValue, 'AU
C': areaUnderRocValueTest}, ignore_index = True);
```





=====

=====

=====

=====

Results of analysis using Average Word2Vec vectorized text features merged with other features using K-NN brute force algorithm:

=====

AUC values of train data:

=====

[0.5185848348348349, 0.5585838338338338, 0.5747256631631632, 0.5996525588088089, 0.6093320038788788, 0.6178955205205205]

=====

Optimal K-Value: 80

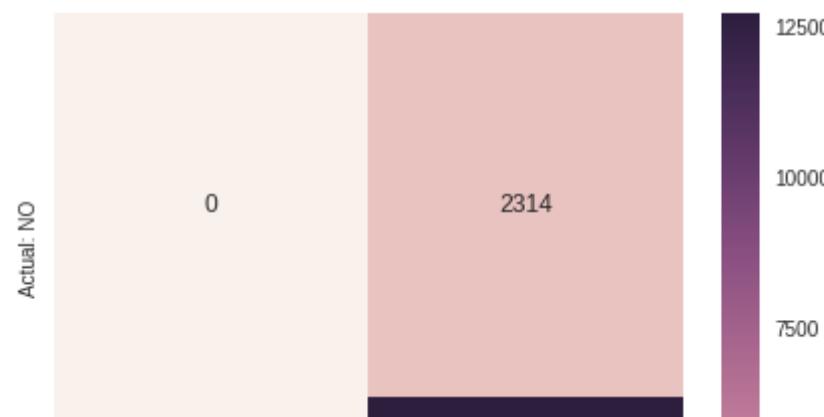
=====

AUC value of test data: 0.6158632325414428

=====

Confusion Matrix :

=====





## Analysis of data containing Tf-Idf weighted vectorized text features using K-NN

```
In [120]: testKValues = [1, 5, 10, 30, 50, 80]
techniques = ['Tf-Idf Weighted Word2Vec'];
for index, technique in enumerate(techniques):
    areaUnderRocValuesTrain = []
    trainingMergedData = hstack((categoriesVectorsSub,\n                                subCategoriesVectorsSub,\n                                teacherPrefixVectorsSub,\n                                schoolStateVectorsSub,\n                                projectGradeVectorsSub,\n                                priceStandardizedSub,\n                                previouslyPostedStandardizedSub));
    testMergedData = hstack((categoriesTransformedTestData,\n                            subCategoriesTransformedTestD\n                            ata,\n                            ata,\n                            a,\n                            ta,\n                            teacherPrefixTransformedTestData,\n                            schoolStateTransformedTestData,\n                            projectGradeTransformedTestData))
```

```

                priceTransformedTestData,\n
                previouslyPostedTransformedTe\n
stData));\n
    if(index == 0):\n
        trainingMergedData = hstack((trainingMergedData,\n
                                      tfIdfWeightedWord2VecTitlesVectors\n
Sub,\n
                                      tfIdfWeightedWord2VecEssaysVectors\n
Sub));\n
        testMergedData = hstack((testMergedData,\n
                                      tfIdfWeightedTitleTransformedTestData,\n
\\
                                      tfIdfWeightedEssayTransformedTestData\n
));\n
    for testKValue in tqdm(testKValues):\n
        knnClassifier = KNeighborsClassifier(n_neighbors = testKValue,\n
algorithm = 'brute');\n
        scores = cross_val_score(knnClassifier, trainingMergedData, cla\n
ssesTraining, cv = 5, scoring = 'roc_auc');\n
        areaUnderRocValuesTrain.append(np.array(scores).mean());\n
\n
        plt.plot(testKValues, areaUnderRocValuesTrain, 'r', label = "Trains\n
ing K vs AUC");\n
        plt.title("Training Data K vs AUC - {} vectorized text".format(tech\n
nique));\n
        plt.xlabel("Hyper parameter - K");\n
        plt.ylabel("Area under curve - AUC");\n
        plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))\n
        plt.show();\n
\n
    optimalKValue = testKValues[np.argmax(areaUnderRocValuesTrain)];\n
    knnClassifier = KNeighborsClassifier(n_neighbors = optimalKValue,\n
algorithm = 'brute');\n
    knnClassifier.fit(trainingMergedData, classesTraining);\n
    predProbScoresTraining = knnClassifier.predict_proba(trainingMerged\n
Data);\n
    fprTrain, tprTrain, thresholdTrain = roc_curve(classesTraining, pre\n
dProbScoresTraining[:, 1]);\n
    predProbScoresTest = knnClassifier.predict_proba(testMergedData);\n

```

```

        fprTest, tprTest, thresholdTest = roc_curve(classesTest, predProbScoresTest[:, 1]);
        areaUnderRocValueTest = auc(fprTest, tprTest);
        plt.plot(fprTrain, tprTrain, 'y', label="Train ROC curve - {} vectorized text".format(technique));
        plt.plot(fprTest, tprTest, 'g', label="Test ROC curve - {} vectorized text".format(technique));
        plt.plot([0, 1], [0, 1], 'k-');
        plt.title("ROC Curves for train and test data using k-value {}".format(optimalKValue))
        plt.xlabel('False positive rate - FPR');
        plt.ylabel('True positive rate - TPR');
        plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
        plt.show();

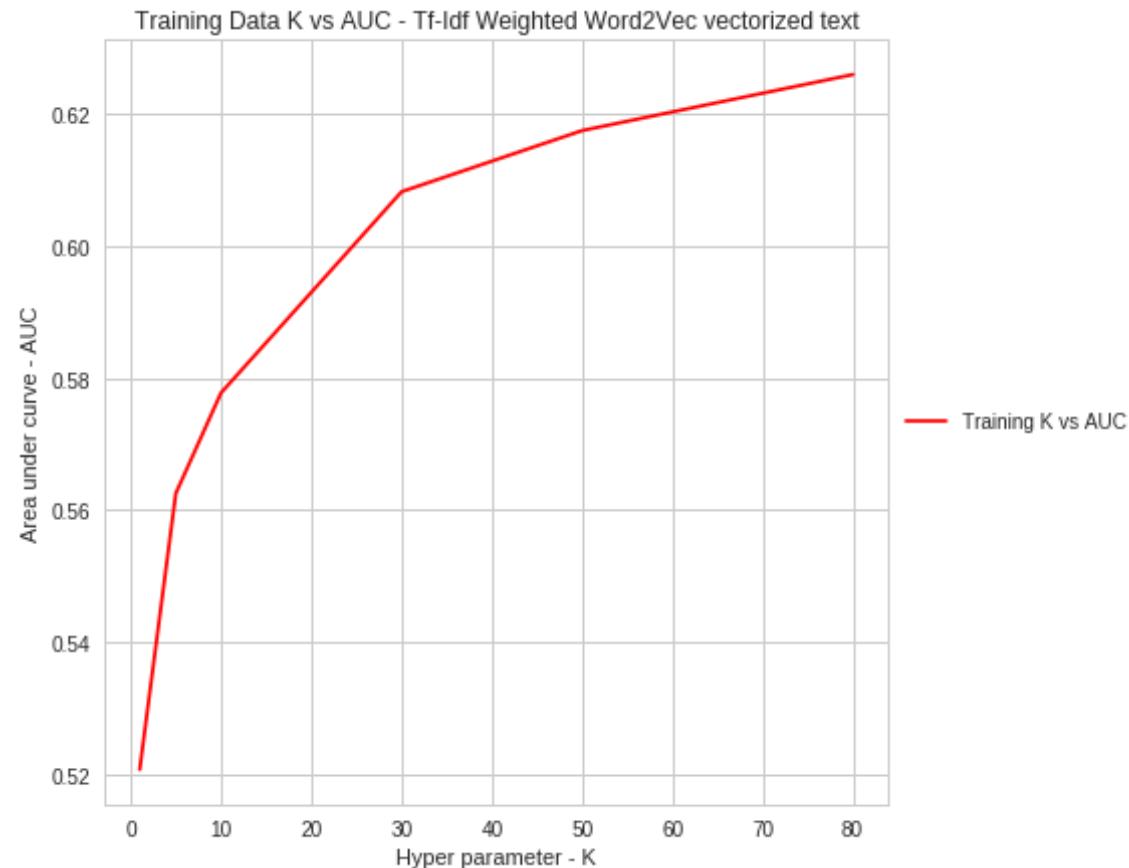
equalsBorder(100);
equalsBorder(100);
print("Results of analysis using {} vectorized text features merged with other features using K-NN brute force algorithm:".format(technique));
equalsBorder(70);
print("AUC values of train data: ");
equalsBorder(40);
print(areaUnderRocValuesTrain);
equalsBorder(40);
print("Optimal K-Value: ", optimalKValue);
equalsBorder(40);
print("AUC value of test data: ", areaUnderRocValueTest);
# Predicting classes of test data projects
predictionClassesTest = knnClassifier.predict(testMergedData);
equalsBorder(40);
# Printing confusion matrix
confusionMatrix = confusion_matrix(classesTest, predictionClassesTest);
# Creating dataframe for generated confusion matrix
confusionMatrixDataFrame = pd.DataFrame(data = confusionMatrix, index = ['Actual: NO', 'Actual: YES'], columns = ['Predicted: NO', 'Predicted: YES']);
print("Confusion Matrix : ");

```

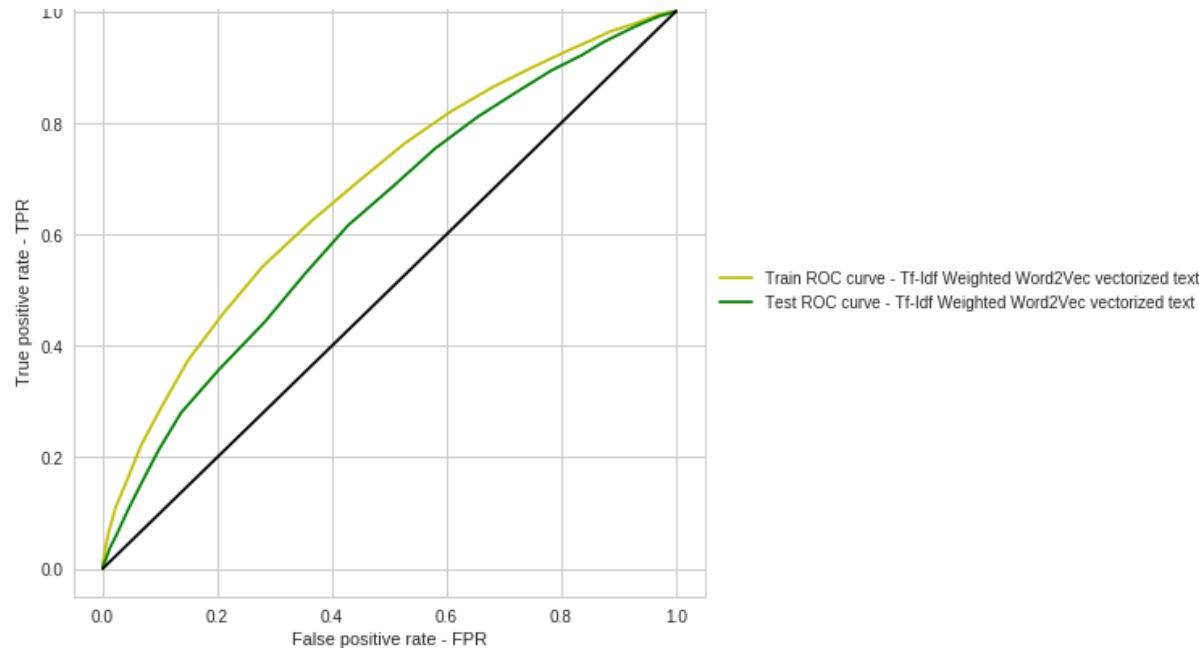
```

        equalsBorder(60);
        sbrn.heatmap(confusionMatrixDataFrame, annot = True, fmt = 'd');
        plt.show();
    # Adding results to results dataframe
    kFoldResultsDataFrame = kFoldResultsDataFrame.append({'Vectorizer': technique, 'Model': 'Brute', 'Hyper Parameter - K': optimalKValue, 'AU C': areaUnderRocValueTest}, ignore_index = True);

```



ROC Curves for train and test data using k-value 80



=====

=====

=====

=====

Results of analysis using Tf-Idf Weighted Word2Vec vectorized text features merged with other features using K-NN brute force algorithm:

=====

AUC values of train data:

=====

[0.5207207207207207, 0.5625490646896897, 0.5777487018268268, 0.608225412912913, 0.6174866898148148, 0.6259789633383384]

=====

Optimal K-Value: 80

=====

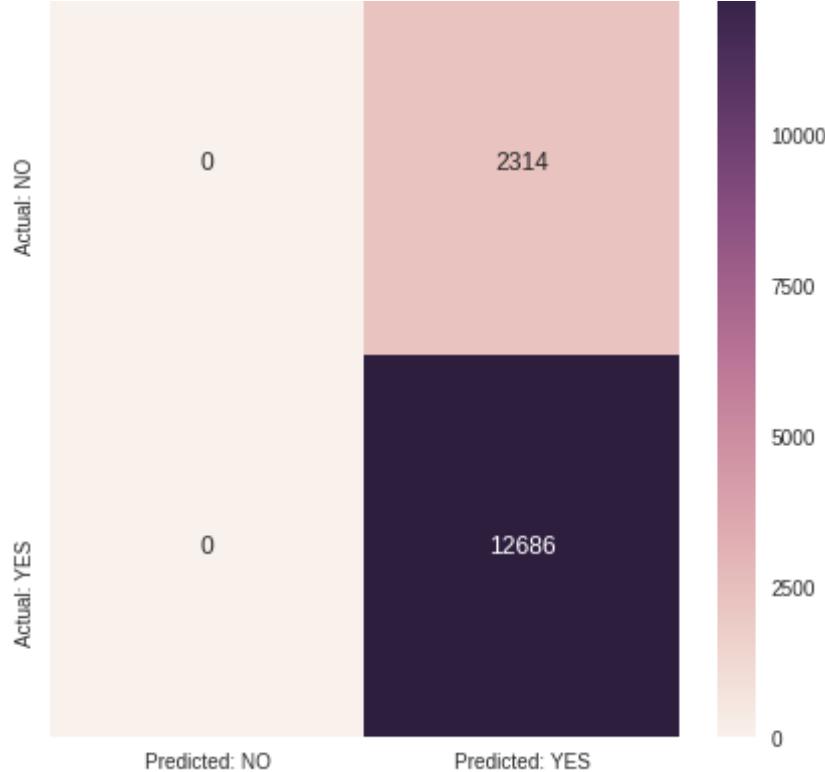
AUC value of test data: 0.6290808329532783

=====

Confusion Matrix :

=====





## Selecting top 2000 important features from data with Tf- Idf vectorized text

```
In [0]: trainingMergedData = hstack((categoriesVectorsSub,\n                                subCategoriesVectorsSub,\n                                teacherPrefixVectorsSub,\n                                schoolStateVectorsSub,\n                                projectGradeVectorsSub,\n                                priceStandardizedSub,\n                                previouslyPostedStandardizedSub));\ntestMergedData = hstack((categoriesTransformedTestData,\n                           subCategoriesTransformedTestData,
```

```
\teacherPrefixTransformedTestData,  
\schoolStateTransformedTestData,\  
projectGradeTransformedTestData,\  
priceTransformedTestData,\  
previouslyPostedTransformedTestData));  
trainingMergedData = hstack((trainingMergedData,\  
    tfIdfTitleModelSub,\  
    tfIdfEssayModelSub));  
testMergedData = hstack((testMergedData,\  
    tfIdfTitleTransformedTestData,\  
    tfIdfEssayTransformedTestData));
```

In [185]:

```
print("Training data shape: ", trainingMergedData.shape);  
print("Test data shape: ", testMergedData.shape);  
print("Classes Training shape: ", classesTraining.shape);
```

```
Training data shape: (59200, 15779)  
Test data shape: (15000, 15779)  
Classes Training shape: (59200,)
```

In [186]:

```
selectKBest = SelectKBest(f_classif, k = 2000);  
filteredFeaturesTrainingMergedData = selectKBest.fit_transform(training  
MergedData, classesTraining);  
filteredFeaturesTrainingMergedData.shape
```

Out[186]: (59200, 2000)

In [9]:

```
selectedFeaturesResultsDataFrame = pd.DataFrame(columns = ['Vectorizer',  
    'Model', 'Hyper Parameter - K', 'AUC']);  
selectedFeaturesResultsDataFrame
```

Out[9]:

Vectorizer	Model	Hyper Parameter - K	AUC

## Analysis on imbalanced data using top 2000 features of data & K-NN(k-fold cross validation)

```
In [126]: testKValues = np.arange(1, 40, 2);
areaUnderRocValuesTrain = [];
for testKValue in tqdm(testKValues):
    knnClassifier = KNeighborsClassifier(n_neighbors = testKValue, algorithm = 'brute');
    scores = cross_val_score(knnClassifier, filteredFeaturesTrainingMergedData, classesTraining, cv = 10, scoring = 'roc_auc');
    areaUnderRocValuesTrain.append(np.array(scores).mean());

plt.plot(testKValues, areaUnderRocValuesTrain, 'r', label = "Training K vs AUC");
plt.title("Training Data K vs AUC - {} vectorized text".format("Tf-Idf"));
plt.xlabel("Hyper parameter - K");
plt.ylabel("Area under curve - AUC");
plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
plt.show();

optimalKValue = testKValues[np.argmax(areaUnderRocValuesTrain)];
knnClassifier = KNeighborsClassifier(n_neighbors = optimalKValue, algorithm = 'brute');
knnClassifier.fit(filteredFeaturesTrainingMergedData, classesTraining);
predProbScoresTraining = knnClassifier.predict_proba(filteredFeaturesTrainingMergedData);
fprTrain, tprTrain, thresholdTrain = roc_curve(classesTraining, predProbScoresTraining[:, 1]);
filteredFeaturesTestMergedData = selectKBest.transform(testMergedData);
predProbScoresTest = knnClassifier.predict_proba(filteredFeaturesTestMergedData);
fprTest, tprTest, thresholdTest = roc_curve(classesTest, predProbScoresTest[:, 1]);
areaUnderRocValueTest = auc(fprTest, tprTest);
plt.plot(fprTrain, tprTrain, 'y', label="Train ROC curve - {} vectorized text".format("Tf-Idf"));
plt.plot(fprTest, tprTest, 'g', label="Test ROC curve - {} vectorized text".format("Tf-Idf"));
```

```

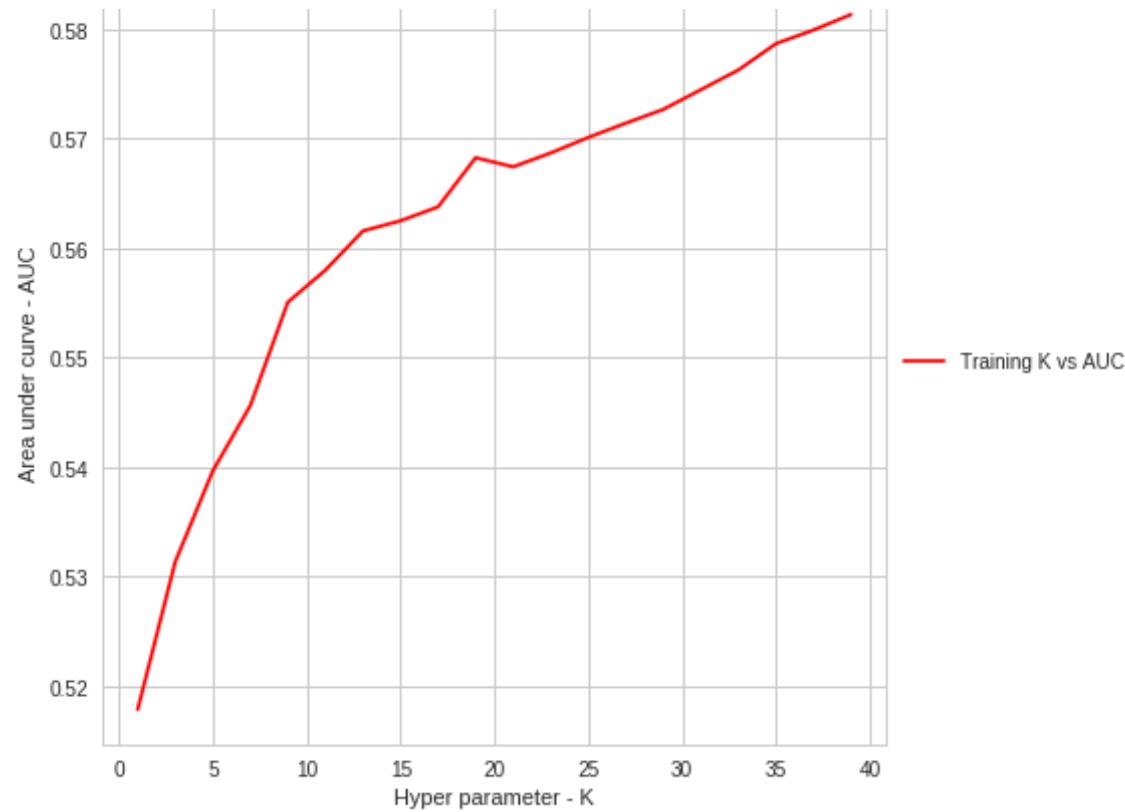
plt.plot([0, 1], [0, 1], 'k-');
plt.title("ROC Curves for train and test data using k-value {}".format(optimalKValue))
plt.xlabel('False positive rate - FPR');
plt.ylabel('True positive rate - TPR');
plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
plt.show();

print("Results of analysis using {} vectorized text features merged with other features using K-NN brute force algorithm:".format("Tf-Idf"));
equalsBorder(70);
print("AUC values of train data: ");
equalsBorder(40);
print(areaUnderRocValuesTrain);
equalsBorder(40);
print("Optimal K-Value: ", optimalKValue);
equalsBorder(40);
print("AUC value of test data: ", areaUnderRocValueTest);
# Predicting classes of test data projects
predictionClassesTest = knnClassifier.predict(filteredFeaturesTestMerge
dData);
equalsBorder(40);
# Adding results to results dataframe
selectedFeaturesResultsDataFrame = selectedFeaturesResultsDataFrame.app
end({'Vectorizer': "Tf-Idf", 'Model': 'Brute', 'Hyper Parameter - K': o
ptimalKValue, 'AUC': areaUnderRocValueTest}, ignore_index = True);
# Printing confusion matrix
confusionMatrix = confusion_matrix(classesTest, predictionClassesTest);
# Creating dataframe for generated confusion matrix
confusionMatrixDataFrame = pd.DataFrame(data = confusionMatrix, index =
['Actual: NO', 'Actual: YES'], columns = ['Predicted: NO', 'Predicted:
YES']);
print("Confusion Matrix : ");
equalsBorder(60);
confusionMatrixDataFrame

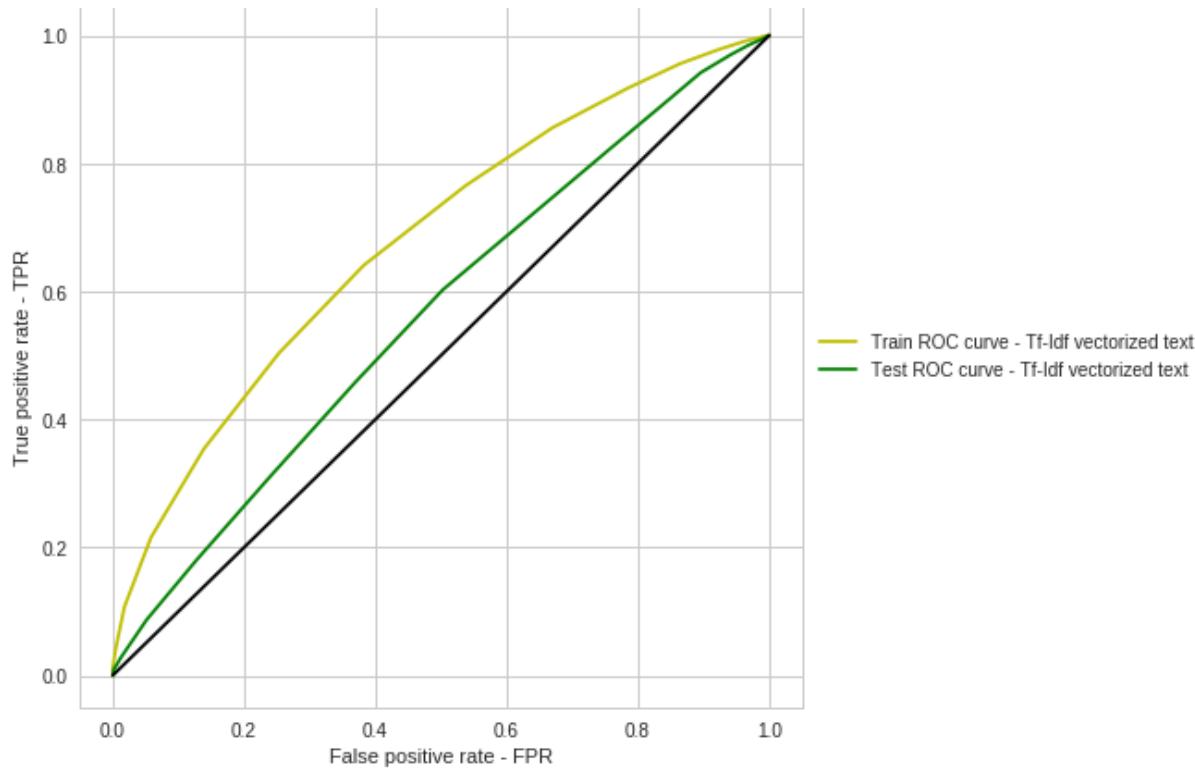
```

Training Data K vs AUC - Tf-Idf vectorized text





ROC Curves for train and test data using k-value 39



Results of analysis using Tf-Idf vectorized text features merged with other features using K-NN brute force algorithm:

=====

AUC values of train data:

=====

```
[0.517892892892893, 0.5313767830330332, 0.5396751438938939, 0.5456481481481481, 0.5550816754254254, 0.5580060372872873, 0.5615842092092092, 0.5625073511011011, 0.5637776526526527, 0.568262794044044, 0.5674354041541542, 0.5686971033533534, 0.5701153966466467, 0.571416478978979, 0.5726778966466467, 0.574476101101101, 0.5762777152152151, 0.5786796796796797, 0.5799234234234235, 0.5813418105605606]
```

=====

Optimal K-Value: 39

=====

AUC value of test data: 0.5656295685796047

=====

Confusion Matrix :

Out[126]:

	Predicted: NO	Predicted: YES
Actual: NO	0	2314
Actual: YES	0	12686

## Analysis on balanced data using top 2000 features of data & K-NN(k-fold cross validation)

In [187]:

```
testKValues = np.arange(1, 40, 2);
areaUnderRocValuesTrain = [];
for testKValue in tqdm(testKValues):
    knnClassifier = KNeighborsClassifier(n_neighbors = testKValue, algorithm = 'brute');
    scores = cross_val_score(knnClassifier, filteredFeaturesTrainingMergedData, classesTraining, cv = 10, scoring = 'roc_auc');
    areaUnderRocValuesTrain.append(np.array(scores).mean());

plt.plot(testKValues, areaUnderRocValuesTrain, 'r', label = "Training K vs AUC");
plt.title("Training Data K vs AUC - {} vectorized text".format("Tf-Idf"));
plt.xlabel("Hyper parameter - K");
plt.ylabel("Area under curve - AUC");
plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
plt.show();

optimalKValue = testKValues[np.argmax(areaUnderRocValuesTrain)];
knnClassifier = KNeighborsClassifier(n_neighbors = optimalKValue, algorithm = 'brute');
knnClassifier.fit(filteredFeaturesTrainingMergedData, classesTraining);
predProbScoresTraining = knnClassifier.predict_proba(filteredFeaturesTrainingMergedData);
fprTrain, tprTrain, thresholdTrain = roc_curve(classesTraining, predPro
```

```

bScoresTraining[:, 1]);
filteredFeaturesTestMergedData = selectKBest.transform(testMergedData);
predProbScoresTest = knnClassifier.predict_proba(filteredFeaturesTestMergedData);
fprTest, tprTest, thresholdTest = roc_curve(classesTest, predProbScoresTest[:, 1]);
areaUnderRocValueTest = auc(fprTest, tprTest);
plt.plot(fprTrain, tprTrain, 'y', label="Train ROC curve - {} vectorized text".format("Tf-Idf"));
plt.plot(fprTest, tprTest, 'g', label="Test ROC curve - {} vectorized text".format("Tf-Idf"));
plt.plot([0, 1], [0, 1], 'k-');
plt.title("ROC Curves for train and test data using k-value {}".format(optimalKValue));
plt.xlabel('False positive rate - FPR');
plt.ylabel('True positive rate - TPR');
plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
plt.show();

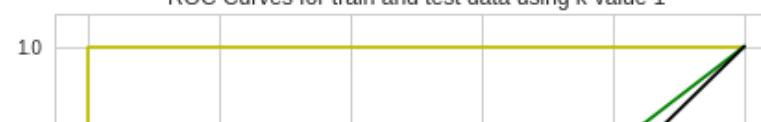
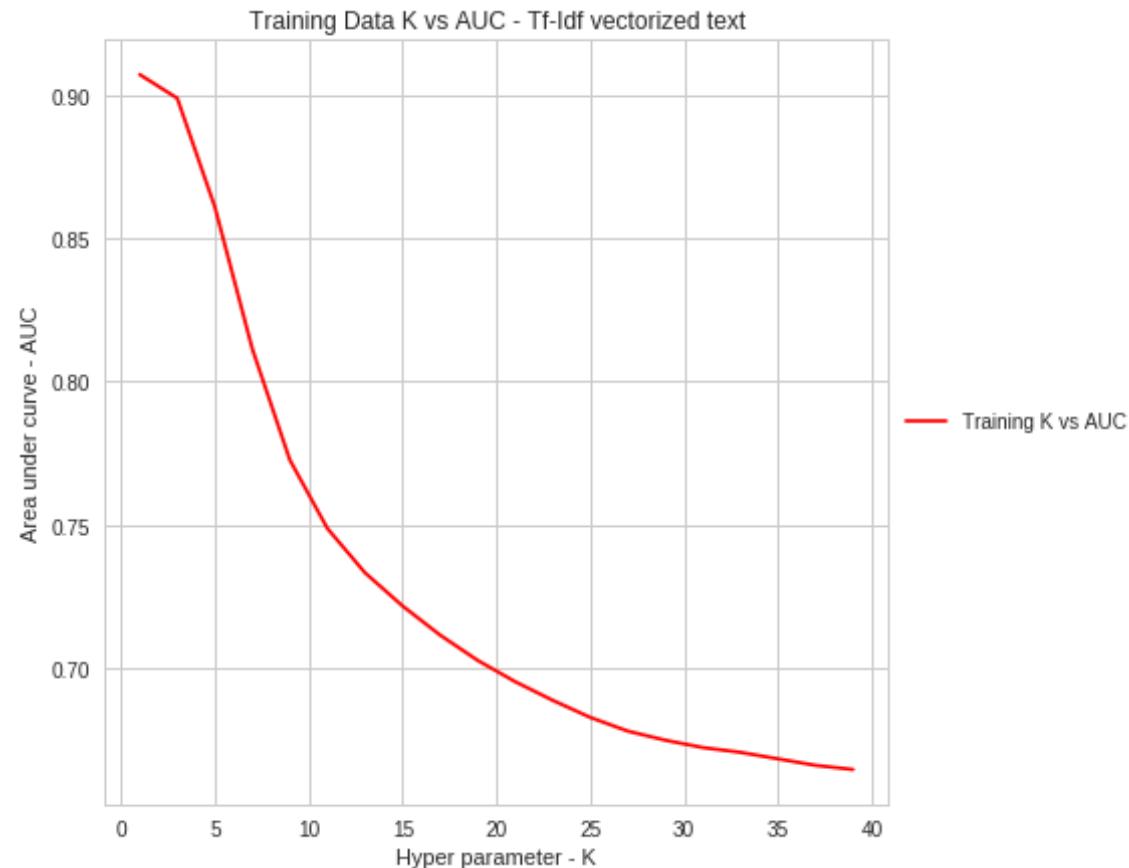
print("Results of analysis using {} vectorized text features merged with other features using K-NN brute force algorithm:".format("Tf-Idf"));
equalsBorder(70);
print("AUC values of train data: ");
equalsBorder(40);
print(areaUnderRocValuesTrain);
equalsBorder(40);
print("Optimal K-Value: ", optimalKValue);
equalsBorder(40);
print("AUC value of test data: ", areaUnderRocValueTest);
# Predicting classes of test data projects
predictionClassesTest = knnClassifier.predict(filteredFeaturesTestMergedData);
equalsBorder(40);
# Adding results to results dataframe
selectedFeaturesResultsDataFrame = selectedFeaturesResultsDataFrame.append({'Vectorizer': "Tf-Idf", 'Model': 'Brute', 'Hyper Parameter - K': optimalKValue, 'AUC': areaUnderRocValueTest}, ignore_index = True);
# Printing confusion matrix
confusionMatrix = confusion_matrix(classesTest, predictionClassesTest);

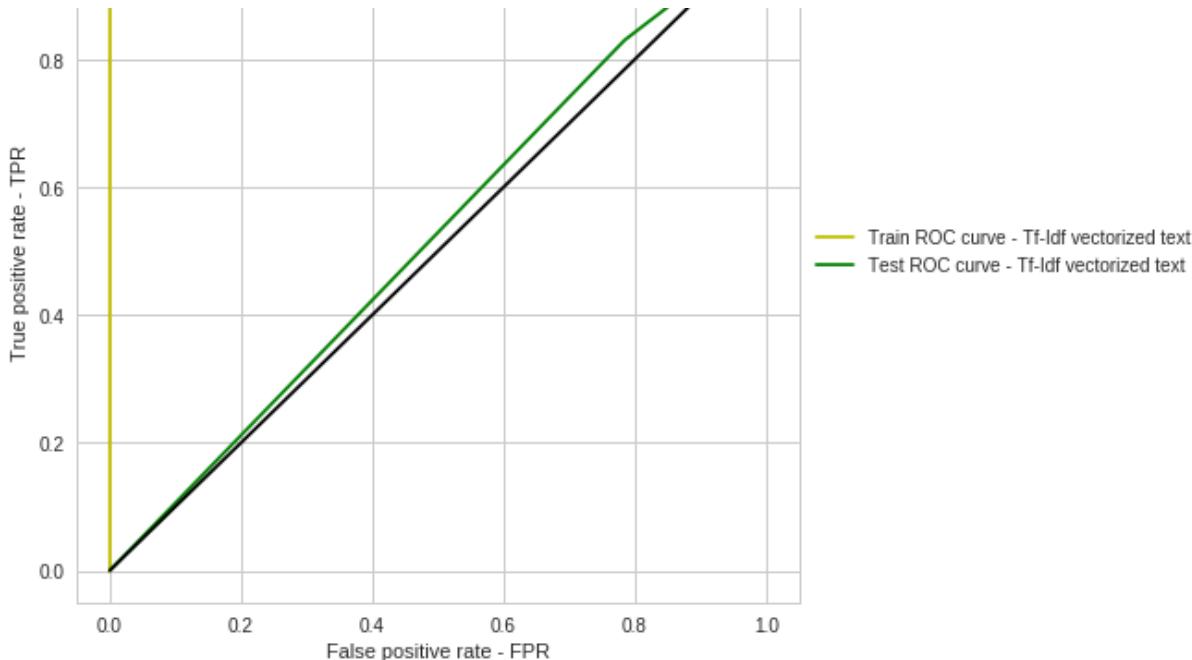
```

```

# Creating dataframe for generated confusion matrix
confusionMatrixDataFrame = pd.DataFrame(data = confusionMatrix, index =
[ 'Actual: NO', 'Actual: YES'], columns = [ 'Predicted: NO', 'Predicted:
YES']);
print("Confusion Matrix : ");
equalsBorder(60);
confusionMatrixDataFrame

```





Results of analysis using Tf-IDf vectorized text features merged with other features using K-NN brute force algorithm:

```
=====
AUC values of train data:
=====
```

```
[0.9074662162162163, 0.8991698262874361, 0.8611724342585829, 0.81140329
96256393, 0.7725532779401023, 0.7486323274287802, 0.7331517188641344,
0.7215411911066472, 0.7113543759130752, 0.7024989043097152, 0.695042172
6625274, 0.6885242079072316, 0.6825375330989774, 0.6777488187089116, 0.
6745635785701242, 0.6719822235664719, 0.6703189543005844, 0.66804206994
1563, 0.6657612650657414, 0.6643143090303141]
```

```
=====
Optimal K-Value: 1
=====
```

```
AUC value of test data: 0.5227675626606945
=====
```

```
Confusion Matrix :
=====
```

Out[18]:

	Predicted: NO	Predicted: YES
Actual: NO	499	1815
Actual: YES	2158	10528

## Summarizing results of above analysis using K-NN

Results of analysis on imbalanced data when K-NN(k-fold cross validation) is used

In [8]: kFoldResultsDataFrame

Out[8]:

	Vectorizer	Model	Hyper Parameter - K	AUC
0	Bag of words	Brute	39	0.618926
1	Tf-Idf	Brute	39	0.588769
2	Average Word2Vec	Brute	80	0.615863
3	Tf-Idf Weighted Word2Vec	Brute	80	0.629080

Results of analysis on imbalanced and balanced when K-NN(k-fold cross validation) with top 200 features

In [12]: selectedFeaturesResultsDataFrame

Out[12]:

	Vectorizer	Model	Hyper Parameter - K	AUC
0	Tf-Idf	Brute	39	0.565629

	Vectorizer	Model	Hyper Parameter - K	AUC
1	Tf-Idf	Brute	1	0.522767

## Conclusions of above analysis:

1. The best k-value by considering AUC values and difference between AUC value of cross-validate, test data will be 80 but the model trained with this k-value and imbalanced data is unable to predict the negative points and so it's like a dumb model.
2. While training with balanced data is able to predict negative and positive points considerably but with accuracy far less than training with imbalanced data.
3. It seems like K-NN cannot be used for solving the above problem if we want good prediction results.