# DonorsChoose

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

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

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

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

## About the DonorsChoose Data Set

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

| Feature | Description |
|---|---|
| `project_id` | A unique identifier for the proposed project. **Example:** `p036502` |
| `project_title` | Title of the project. **Examples:**<br><br>- `Art Will Make You Happy!`<br>- `First Grade Fun` |
| `project_grade_category` | Grade level of students for which the project is targeted. One of the following enumerated values:<br><br>- `Grades PreK-2`<br>- `Grades 3-5`<br>- `Grades 6-8`<br>- `Grades 9-12` |
| `project_subject_categories` | One or more (comma-separated) subject categories for the project from the following enumerated list of values:<br><br>- `Applied Learning`<br>- `Care & Hunger`<br>- `Health & Sports`<br>- `History & Civics`<br>- `Literacy & Language`<br>- `Math & Science`<br>- `Music & The Arts`<br>- `Special Needs`<br>- `Warmth`<br><br>**Examples:**<br><br>- `Music & The Arts`<br>- `Literacy & Language, Math & Science` |
| `school_state` | State where school is located ([Two-letter U.S. postal code](#)). **Example:** `WY` |
| `project_subject_subcategories` | One or more (comma-separated) subject subcategories for the project. **Examples:**<br><br>- `Literacy`<br>- `Literature & Writing, Social Sciences` |
| `project_resource_summary` | An explanation of the resources needed for the project. **Example:**<br><br>- `My students need hands on literacy materials to manage sensory needs!` |
| `project_essay_1` | First application essay[*] |
| `project_essay_2` | Second application essay[*] |
| `project_essay_3` | Third application essay[*] |

| Feature | Description |
| --- | --- |
| project_essay_4 | Fourth application essay |
| project_submitted_datetime | Datetime when project application was submitted. **Example:** `2016-04-28 12:43:56.245` |
| teacher_id | A unique identifier for the teacher of the proposed project. **Example:** `bdf8baa8fedef6bfeec7ae4ff1c15c56` |
| teacher_prefix | Teacher's title. One of the following enumerated values:<br><br>• `nan`<br>• `Dr.`<br>• `Mr.`<br>• `Mrs.`<br>• `Ms.`<br>• `Teacher.` |
| teacher_number_of_previously_posted_projects | Number of project applications previously submitted by the same teacher. **Example:** `2` |

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

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

| Feature | Description |
| --- | --- |
| id | A `project_id` value from the `train.csv` file. **Example:** `p036502` |
| description | Desciption of the resource. **Example:** `Tenor Saxophone Reeds, Box of 25` |
| quantity | Quantity of the resource required. **Example:** `3` |
| price | Price of the resource required. **Example:** `9.95` |

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

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

| Label | Description |
| --- | --- |
| project_is_approved | A binary flag indicating whether DonorsChoose approved the project. A value of `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.

In [1]:

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

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

```
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

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

## 1.1 Reading Data

In [2]:

```
project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')
```

In [3]:

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

```
Number of data points in train data (109248, 17)
--------------------------------------------------
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
 'project_submitted_datetime' 'project_grade_category'
 'project_subject_categories' 'project_subject_subcategories'
 'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
 'project_essay_4' 'project_resource_summary'
 'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

In [4]:

```
print("Number of data points in train data", resource_data.shape)
print(resource_data.columns.values)
resource_data.head(2)
```

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

Out[4]:

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

# 1.2 Data Analysis

```python
# PROVIDE CITATIONS TO YOUR CODE IF YOU TAKE IT FROM ANOTHER WEBSITE.
# https://matplotlib.org/gallery/pie_and_polar_charts/pie_and_donut_labels.html#sphx-glr-gallery-p
ie-and-polar-charts-pie-and-donut-labels-py


y_value_counts = project_data['project_is_approved'].value_counts()
print("Number of projects thar are approved for funding ", y_value_counts[1], ", (",
(y_value_counts[1]/(y_value_counts[1]+y_value_counts[0]))*100,"%)")
print("Number of projects thar are not approved for funding ", y_value_counts[0], ", (",
(y_value_counts[0]/(y_value_counts[1]+y_value_counts[0]))*100,"%)")

fig, ax = plt.subplots(figsize=(6, 6), subplot_kw=dict(aspect="equal"))
recipe = ["Accepted", "Not Accepted"]

data = [y_value_counts[1], y_value_counts[0]]

wedges, texts = ax.pie(data, wedgeprops=dict(width=0.5), startangle=-40)

bbox_props = dict(boxstyle="square,pad=0.3", fc="w", ec="k", lw=0.72)
kw = dict(xycoords='data', textcoords='data', arrowprops=dict(arrowstyle="-"),
          bbox=bbox_props, zorder=0, va="center")

for i, p in enumerate(wedges):
    ang = (p.theta2 - p.theta1)/2. + p.theta1
    y = np.sin(np.deg2rad(ang))
    x = np.cos(np.deg2rad(ang))
    horizontalalignment = {-1: "right", 1: "left"}[int(np.sign(x))]
    connectionstyle = "angle,angleA=0,angleB={}".format(ang)
    kw["arrowprops"].update({"connectionstyle": connectionstyle})
    ax.annotate(recipe[i], xy=(x, y), xytext=(1.35*np.sign(x), 1.4*y),
                horizontalalignment=horizontalalignment, **kw)

ax.set_title("Nmber of projects that are Accepted and not accepted")

plt.show()
```
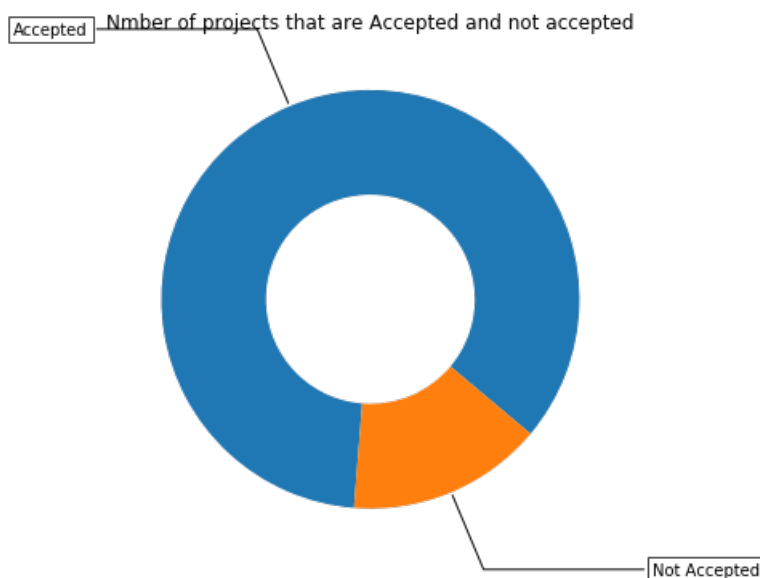
```
Number of projects thar are approved for funding  92706 , ( 84.85830404217927 %)
Number of projects thar are not approved for funding  16542 , ( 15.141695957820739 %)
```



### 1.2.1 Univariate Analysis: School State

```python
# Pandas dataframe groupby count, mean: https://stackoverflow.com/a/19385591/4084039
```

```python
temp = pd.DataFrame(project_data.groupby("school_state")
["project_is_approved"].apply(np.mean)).reset_index()
# if you have data which contain only 0 and 1, then the mean = percentage (think about it)
temp.columns = ['state_code', 'num_proposals']
```

In [7]:

```python
# https://www.csi.cuny.edu/sites/default/files/pdf/administration/ops/2letterstabbrev.pdf
temp.sort_values(by=['num_proposals'], inplace=True)
print("States with lowest % approvals")
print(temp.head(5))
print('='*50)
print("States with highest % approvals")
print(temp.tail(5))
```

```
States with lowest % approvals
   state_code  num_proposals
46         VT       0.800000
7          DC       0.802326
43         TX       0.813142
26         MT       0.816327
18         LA       0.831245
==================================================
States with highest % approvals
   state_code  num_proposals
30         NH       0.873563
35         OH       0.875152
47         WA       0.876178
28         ND       0.888112
8          DE       0.897959
```

In [8]:

```python
#stacked bar plots matplotlib:
https://matplotlib.org/gallery/lines_bars_and_markers/bar_stacked.html
def stack_plot(data, xtick, col2='project_is_approved', col3='total'):
    ind = np.arange(data.shape[0])

    plt.figure(figsize=(20,5))
    p1 = plt.bar(ind, data[col3].values)
    p2 = plt.bar(ind, data[col2].values)

    plt.ylabel('Projects')
    plt.title('Number of projects aproved vs rejected')
    plt.xticks(ind, list(data[xtick].values))
    plt.legend((p1[0], p2[0]), ('total', 'accepted'))
    plt.show()
```

In [9]:

```python
def univariate_barplots(data, col1, col2='project_is_approved', top=False):
    # Count number of zeros in dataframe python: https://stackoverflow.com/a/51540521/4084039
    temp = pd.DataFrame(project_data.groupby(col1)[col2].agg(lambda x: x.eq(1).sum())).reset_index(
)

    # Pandas dataframe grouby count: https://stackoverflow.com/a/19385591/4084039
    temp['total'] = pd.DataFrame(project_data.groupby(col1)
[col2].agg({'total':'count'})).reset_index()['total']
    temp['Avg'] = pd.DataFrame(project_data.groupby(col1)[col2].agg({'Avg':'mean'})).reset_index()[
'Avg']

    temp.sort_values(by=['total'],inplace=True, ascending=False)

    if top:
        temp = temp[0:top]

    stack_plot(temp, xtick=col1, col2=col2, col3='total')
    print(temp.head(5))
    print("="*50)
    print(temp.tail(5))
```
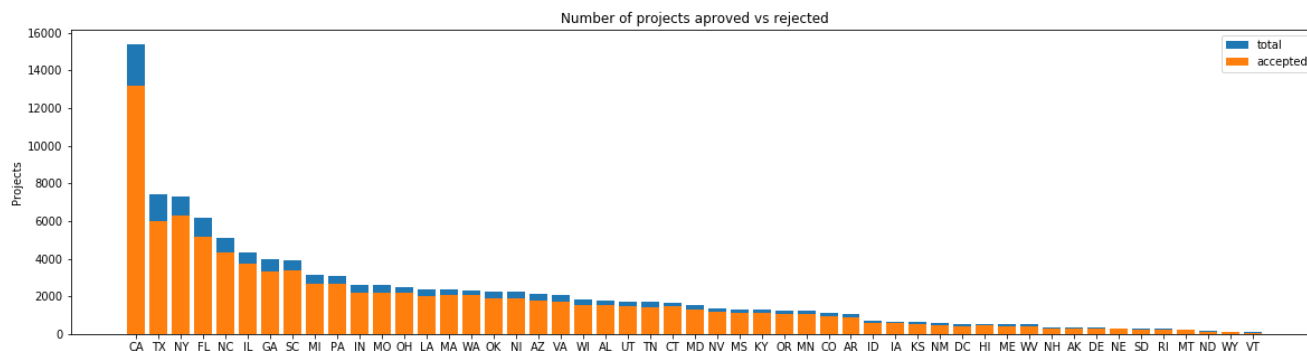
In [10]:

```
univariate_barplots(project_data, 'school_state', 'project_is_approved', False)
```



```
   school_state  project_is_approved   total        Avg
4            CA                13205    15388   0.858136
43           TX                 6014     7396   0.813142
34           NY                 6291     7318   0.859661
9            FL                 5144     6185   0.831690
27           NC                 4353     5091   0.855038
================================================
   school_state  project_is_approved   total        Avg
39           RI                  243      285   0.852632
26           MT                  200      245   0.816327
28           ND                  127      143   0.888112
50           WY                   82       98   0.836735
46           VT                   64       80   0.800000
```
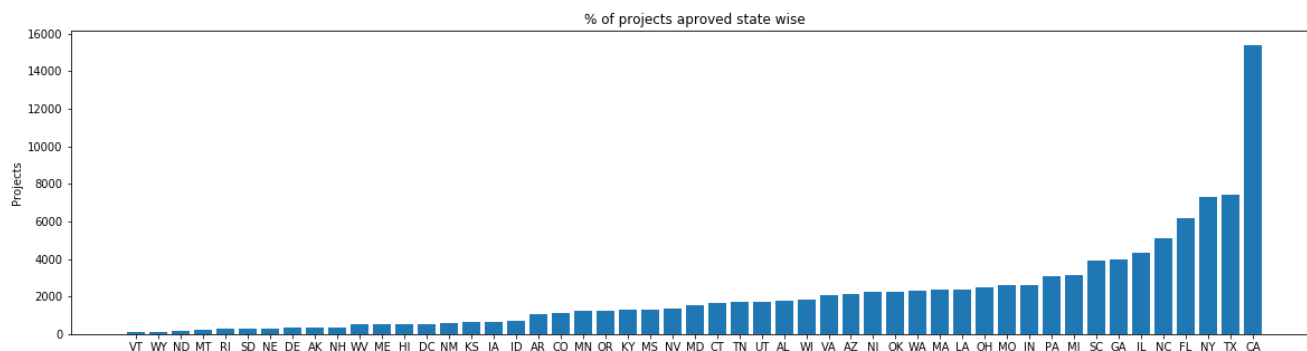
In [11]:

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

In [12]:

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


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

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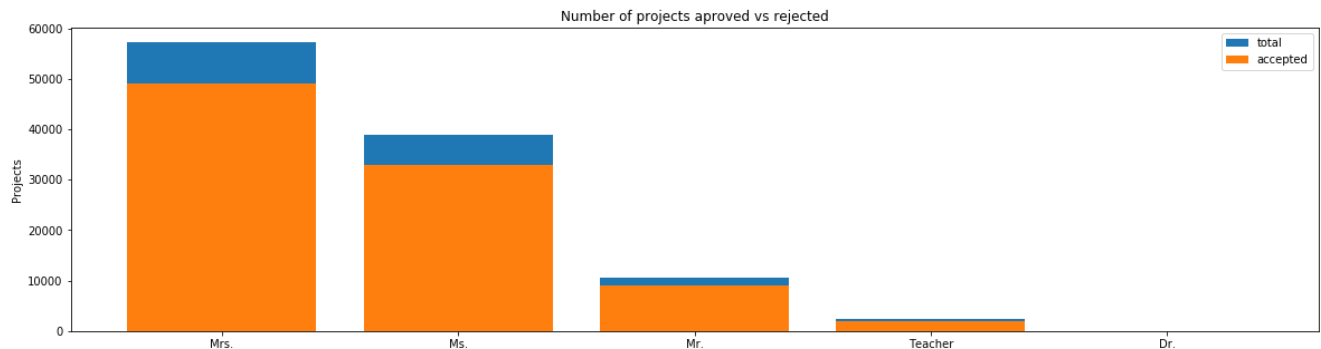


**SUMMARY:**

**1. Every state has greater than 80% success rate in approval**

**2. California has the highest number of approved projects**

---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

## 1.2.2 Univariate Analysis: teacher_prefix

In [13]:

```
univariate_barplots(project_data, 'teacher_prefix', 'project_is_approved' , top=False)
```



Number of projects aproved vs rejected

```
   teacher_prefix  project_is_approved   total        Avg
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
==================================================
   teacher_prefix  project_is_approved   total        Avg
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
```

In [14]:

```python
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()

#https://chrisalbon.com/python/data_wrangling/pandas_missing_data/
#project_data = project_data.dropna()

for word in project_data['teacher_prefix'].values:
    my_counter.update(str(word).split())
```
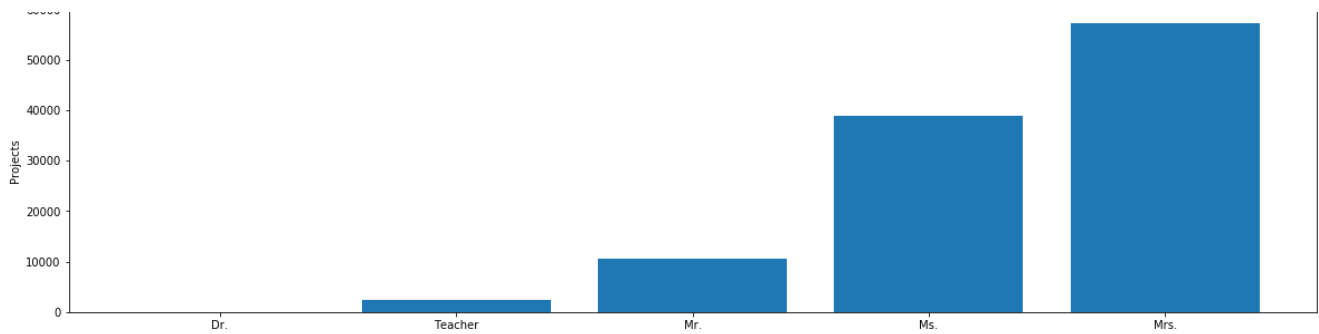
In [15]:

```python
# dict sort by value python: https://stackoverflow.com/a/613218/4084039
teacher_prefix_dict = dict(my_counter)
#https://thispointer.com/different-ways-to-remove-a-key-from-dictionary-in-python/
del teacher_prefix_dict['nan']
teacher_prefix_dict = dict(sorted(teacher_prefix_dict.items(), key=lambda kv: kv[1]))
ind = np.arange(len(teacher_prefix_dict))
plt.figure(figsize=(20,5))
p1 = plt.bar(ind, list(teacher_prefix_dict.values()))
plt.ylabel('Projects')
plt.title('number of projects aproved based on teacher prefix wise')
plt.xticks(ind, list(teacher_prefix_dict.keys()))
print(teacher_prefix_dict)
plt.show()
```

```
{'Dr.': 13, 'Teacher': 2360, 'Mr.': 10648, 'Ms.': 38955, 'Mrs.': 57269}
```
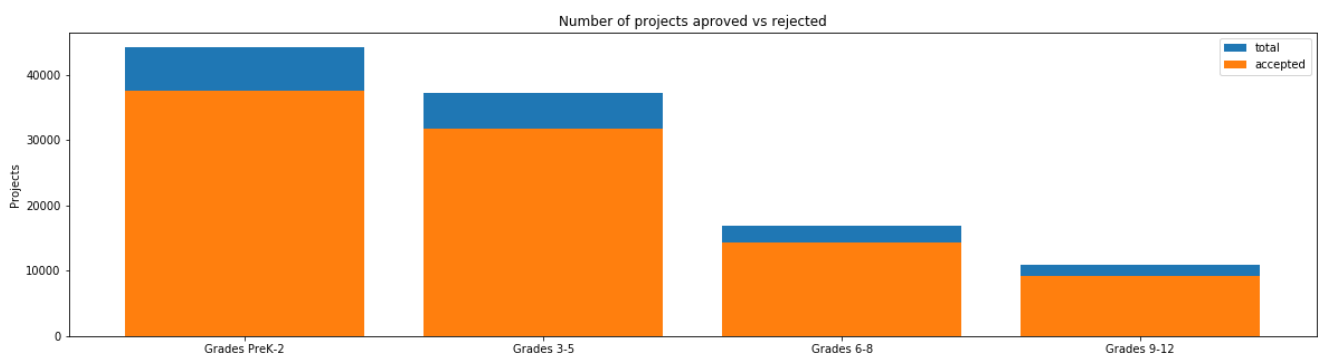
number of projects aproved based on teacher prefix wise

**Summary - If we do not take Dr. into consideration due to very less number of projects as compared to other prefixes, then teachers who have a proper prefix like Mr.,Miss,Mrs. get the approval rate of more than 84% as compared to when the prefix is just "teacher" which has 79% rate.**

--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

### 1.2.3 Univariate Analysis: project_grade_category

In [16]:

```
univariate_barplots(project_data, 'project_grade_category', 'project_is_approved', top=False)
```



```
  project_grade_category  project_is_approved  total       Avg
3           Grades PreK-2                37536  44225  0.848751
0             Grades 3-5                31729  37137  0.854377
1             Grades 6-8                14258  16923  0.842522
2            Grades 9-12                 9183  10963  0.837636
================================================
  project_grade_category  project_is_approved  total       Avg
3           Grades PreK-2                37536  44225  0.848751
0             Grades 3-5                31729  37137  0.854377
1             Grades 6-8                14258  16923  0.842522
2            Grades 9-12                 9183  10963  0.837636
```

In [17]:

```
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()
for word in project_data['project_grade_category'].values:
    my_counter.update(word.split("/"))
```
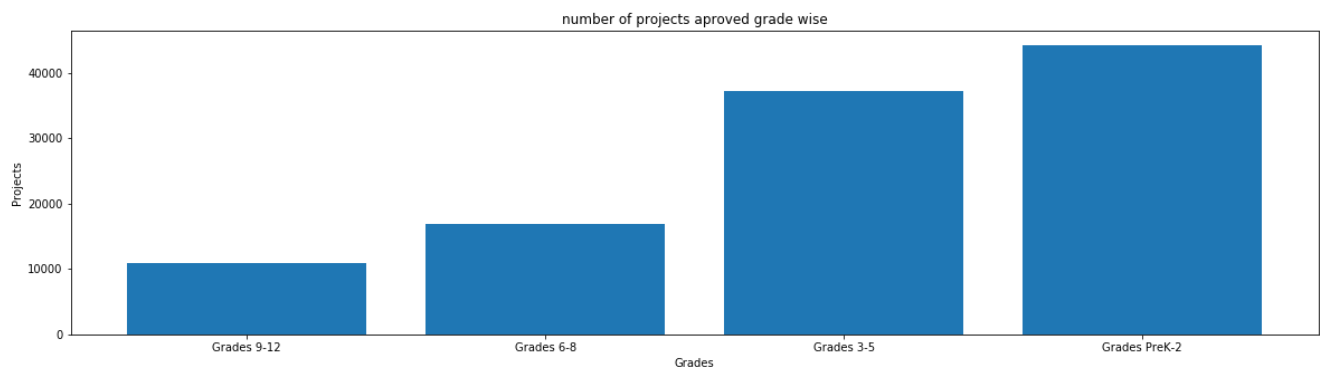
In [18]:

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

ind = np.arange(len(project_grade_category_dict))
```

```
plt.figure(figsize=(20,5))
p1 = plt.bar(ind, list(project_grade_category_dict.values()))

plt.ylabel('Projects')
plt.xlabel('Grades')
plt.title('number of projects aproved grade wise')
plt.xticks(ind, list(project_grade_category_dict.keys()))
plt.show()
```



**Summary - As the Grades increases the number of projects submitted decreases**

--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

### 1.2.4 Univariate Analysis: project_subject_categories

In [19]:

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

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

In [20]:

```
project_data['clean_categories'] = cat_list
project_data.drop(['project_subject_categories'], axis=1, inplace=True)
project_data.head(2)
```
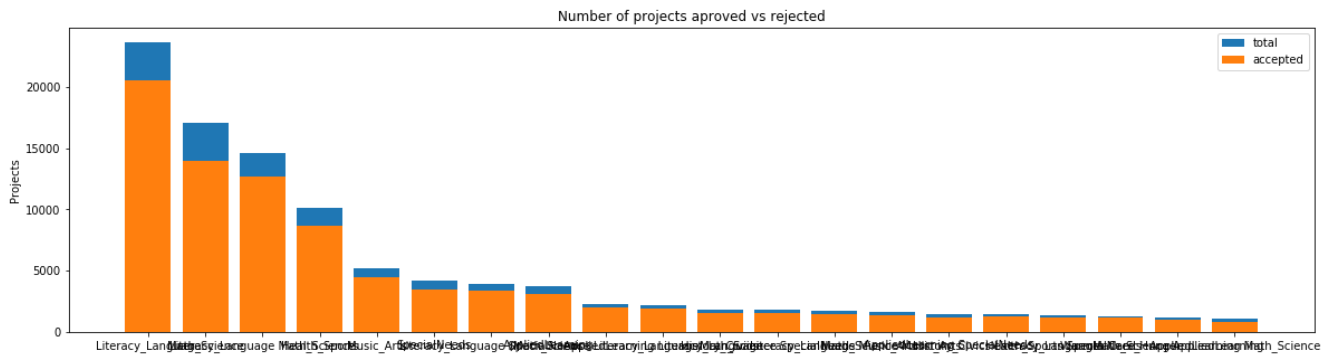
Out[20]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_cate |
|---|---|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | 2016-12-05 13:43:57 | Grades P |

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_cate |
|---|---|---|---|---|---|---|---|
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 2016-10-25 09:22:10 | Grade |



`In [21]:`

```python
univariate_barplots(project_data, 'clean_categories', 'project_is_approved', top=20)
```



| | clean_categories | project_is_approved | total | Avg |
|---|---|---|---|---|
| 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 |

==================================================

| | clean_categories | project_is_approved | total | Avg |
|---|---|---|---|---|
| 19 | History_Civics Literacy_Language | 1271 | 1421 | 0.894441 |
| 14 | Health_Sports SpecialNeeds | 1215 | 1391 | 0.873472 |
| 50 | Warmth Care_Hunger | 1212 | 1309 | 0.925898 |
| 33 | Math_Science AppliedLearning | 1019 | 1220 | 0.835246 |
| 4 | AppliedLearning Math_Science | 855 | 1052 | 0.812738 |

**Summary:**

**1. Projects for both Maths and Science when combined with Applied Learning has the least number of projects proposed as well approved.**

**2. Project approval rate of maths and science can be increased if they are combined with Literacy language**
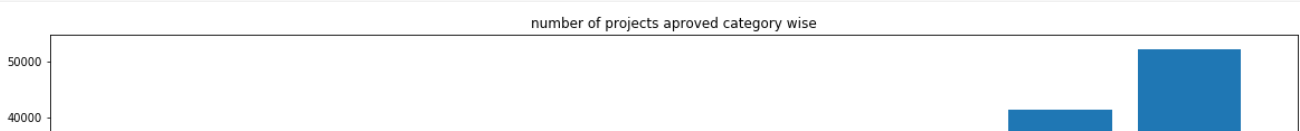
`In [22]:`

```python
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()
for word in project_data['clean_categories'].values:
    my_counter.update(word.split())
```
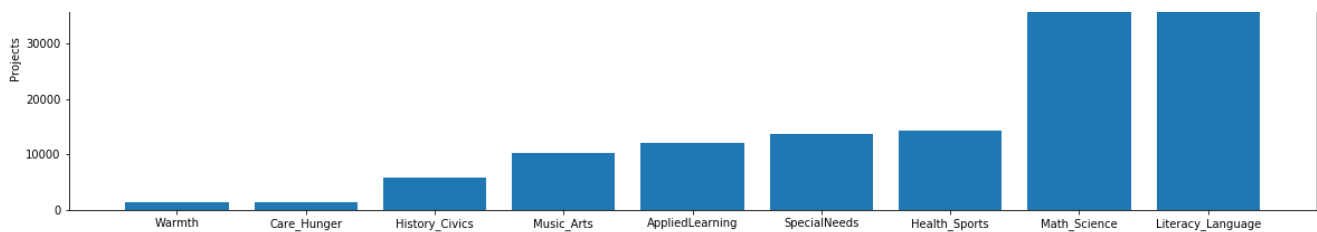
`In [23]:`

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


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

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

```
for i, j in sorted_cat_dict.items():
    print("{:20} :{:10}".format(i,j))
```

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

**Summary:**
**1. Literacy language has the highest number of approved projects and warmth has the lowest.**

## 1.2.5 Univariate Analysis: project_subject_subcategories

In [25]:

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

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

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

In [26]:

```
project_data['clean_subcategories'] = sub_cat_list
project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
project_data.head(2)
```
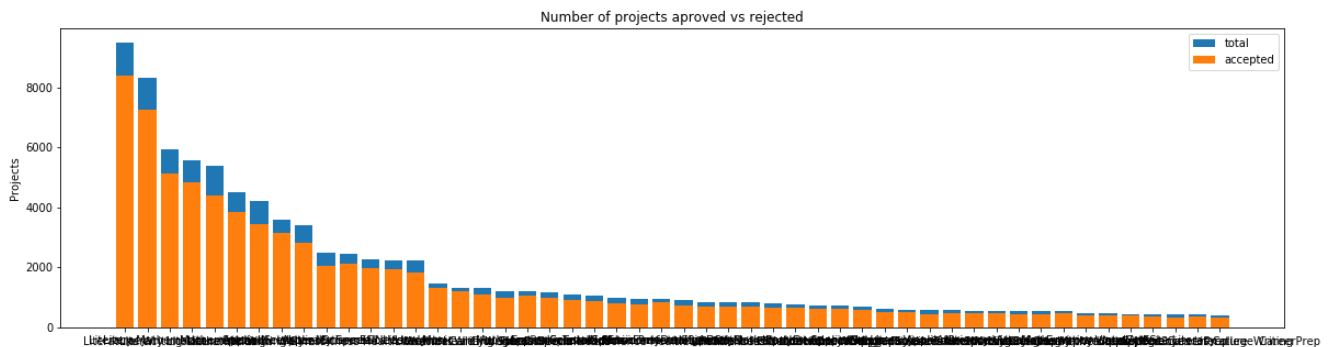
Out[26]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_cate |
|---|---|---|---|---|---|---|---|
| **0** | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | 2016-12-05 13:43:57 | Grades P |

| Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_cate |
|---|---|---|---|---|---|---|
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 2016-10-25 09:22:10 | Grade |

In [27]:

```
univariate_barplots(project_data, 'clean_subcategories', 'project_is_approved', top=50)
```



Number of projects aproved vs rejected

|  | clean_subcategories | project_is_approved | total | Avg |
|---|---|---|---|---|
| 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 |

==================================================

|  | clean_subcategories | project_is_approved | total | Avg |
|---|---|---|---|---|
| 196 | EnvironmentalScience Literacy | 389 | 444 | 0.876126 |
| 127 | ESL | 349 | 421 | 0.828979 |
| 79 | College_CareerPrep | 343 | 421 | 0.814727 |
| 17 | AppliedSciences Literature_Writing | 361 | 420 | 0.859524 |
| 3 | AppliedSciences College_CareerPrep | 330 | 405 | 0.814815 |

**Summary:**

**1. Literacy has the highest rate of Approved projects and mathematics has the lowest**
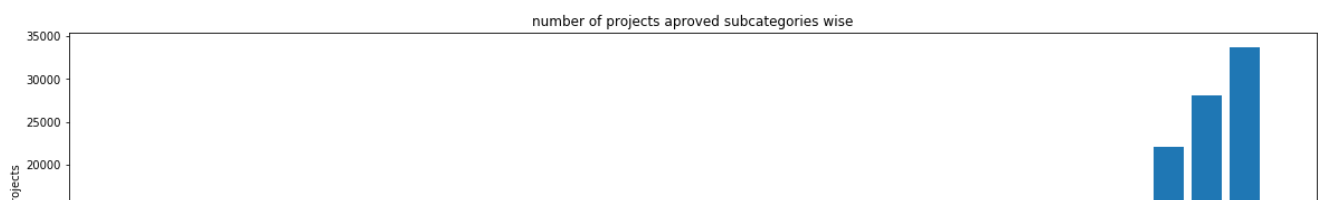
In [28]:

```
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()
for word in project_data['clean_subcategories'].values:
    my_counter.update(word.split())
```
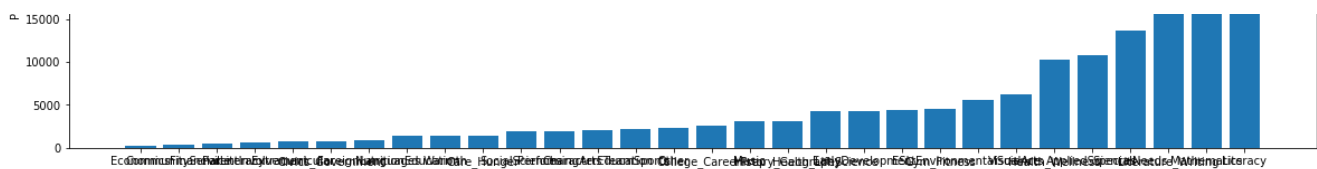
In [29]:

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


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

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



number of projects aproved subcategories wise

```
for i, j in sorted_sub_cat_dict.items():
    print("{:20} :{:10}".format(i,j))
```

```
Economics            :       269
CommunityService     :       441
FinancialLiteracy    :       568
ParentInvolvement    :       677
Extracurricular      :       810
Civics_Government    :       815
ForeignLanguages     :       890
NutritionEducation   :      1355
Warmth               :      1388
Care_Hunger          :      1388
SocialSciences       :      1920
PerformingArts       :      1961
CharacterEducation   :      2065
TeamSports           :      2192
Other                :      2372
College_CareerPrep   :      2568
Music                :      3145
History_Geography    :      3171
Health_LifeScience   :      4235
EarlyDevelopment     :      4254
ESL                  :      4367
Gym_Fitness          :      4509
EnvironmentalScience :      5591
VisualArts           :      6278
Health_Wellness      :     10234
AppliedSciences      :     10816
SpecialNeeds         :     13642
Literature_Writing   :     22179
Mathematics          :     28074
Literacy             :     33700
```

**Summary:**
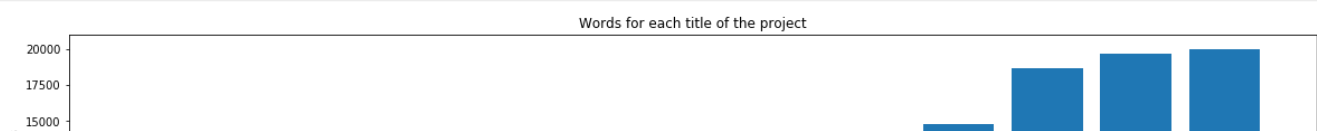**1. Here we can see Economics has the lowest number of approved projects.**

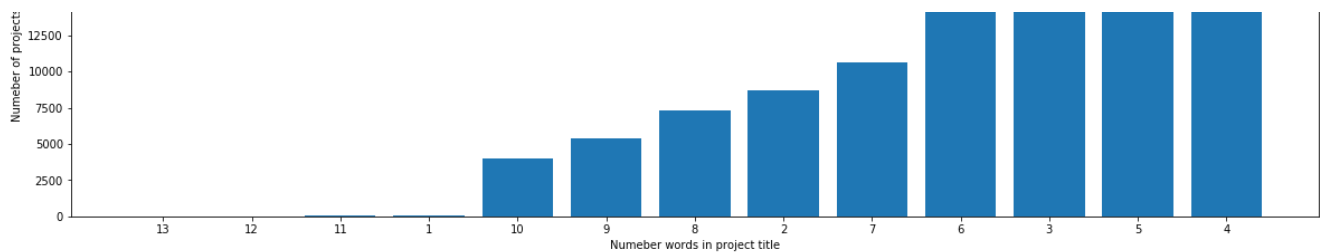### 1.2.6 Univariate Analysis: Text features (Title)

In [31]:

```
#How to calculate number of words in a string in DataFrame:
https://stackoverflow.com/a/37483537/4084039
word_count = project_data['project_title'].str.split().apply(len).value_counts()
word_dict = dict(word_count)
word_dict = dict(sorted(word_dict.items(), key=lambda kv: kv[1]))


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

plt.ylabel('Numeber of projects')
plt.xlabel('Numeber words in project title')
plt.title('Words for each title of the project')
plt.xticks(ind, list(word_dict.keys()))
plt.show()
```



Words for each title of the project

**Summary:**

**1. There are no projects with words greater than 11 in the title.**
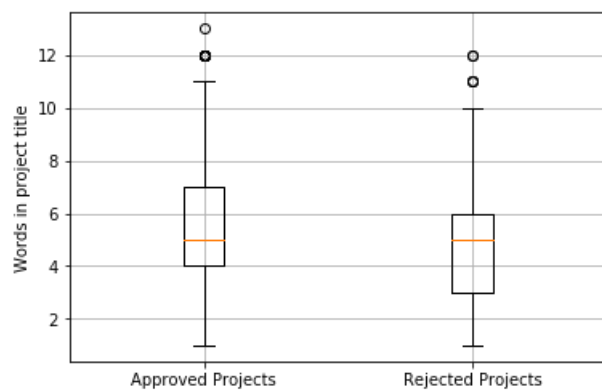
In [32]:

```
approved_title_word_count = project_data[project_data['project_is_approved']==1]['project_title'].
str.split().apply(len)
approved_title_word_count = approved_title_word_count.values

rejected_title_word_count = project_data[project_data['project_is_approved']==0]['project_title'].
str.split().apply(len)
rejected_title_word_count = rejected_title_word_count.values
```

In [33]:

```
# https://glowingpython.blogspot.com/2012/09/boxplot-with-matplotlib.html
plt.boxplot([approved_title_word_count, rejected_title_word_count])
plt.xticks([1,2],('Approved Projects','Rejected Projects'))
plt.ylabel('Words in project title')
plt.grid()
plt.show()
```
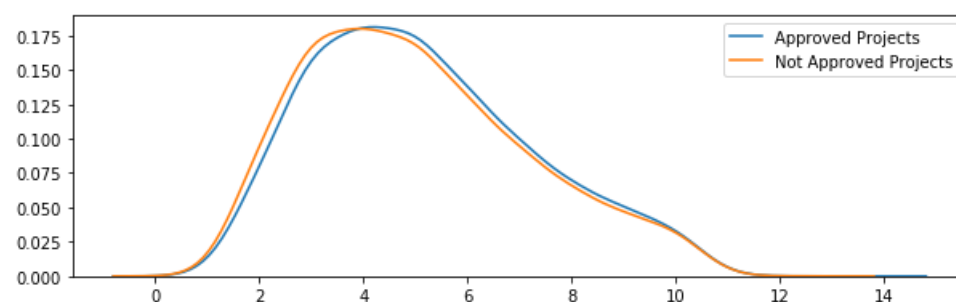


**Summary:**

**1. 75% of the approved projects have less than 7 words in the title.**

In [34]:

```
plt.figure(figsize=(10,3))
sns.kdeplot(approved_title_word_count,label="Approved Projects", bw=0.6)
sns.kdeplot(rejected_title_word_count,label="Not Approved Projects", bw=0.6)
plt.legend()
plt.show()
```

In [35]:

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

In [36]:

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

**Summary:**
1. Approved projects have slightly more number of words than those of rejected projects.

### 1.2.7 Univariate Analysis: Text features (Project Essay's)

In [37]:

```python
# merge two column text dataframe:
project_data["essay"] = project_data["project_essay_1"].map(str) +\
                        project_data["project_essay_2"].map(str) + \
                        project_data["project_essay_3"].map(str) + \
                        project_data["project_essay_4"].map(str)
```
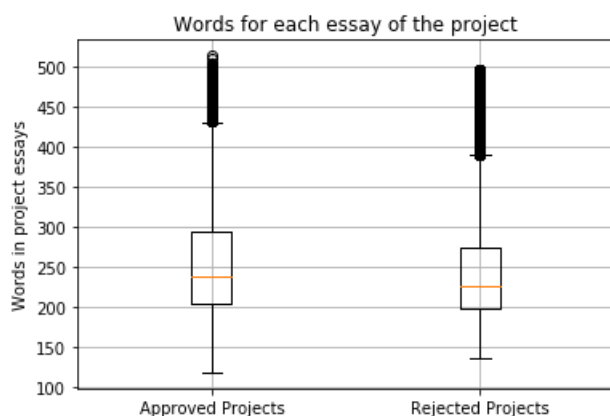
In [38]:

```python
approved_word_count = project_data[project_data['project_is_approved']==1]['essay'].str.split().app
ly(len)
approved_word_count = approved_word_count.values

rejected_word_count = project_data[project_data['project_is_approved']==0]['essay'].str.split().app
ly(len)
rejected_word_count = rejected_word_count.values
```

In [39]:

```python
# https://glowingpython.blogspot.com/2012/09/boxplot-with-matplotlib.html
plt.boxplot([approved_word_count, rejected_word_count])
plt.title('Words for each essay of the project')
plt.xticks([1,2],('Approved Projects','Rejected Projects'))
plt.ylabel('Words in project essays')
plt.grid()
plt.show()
```
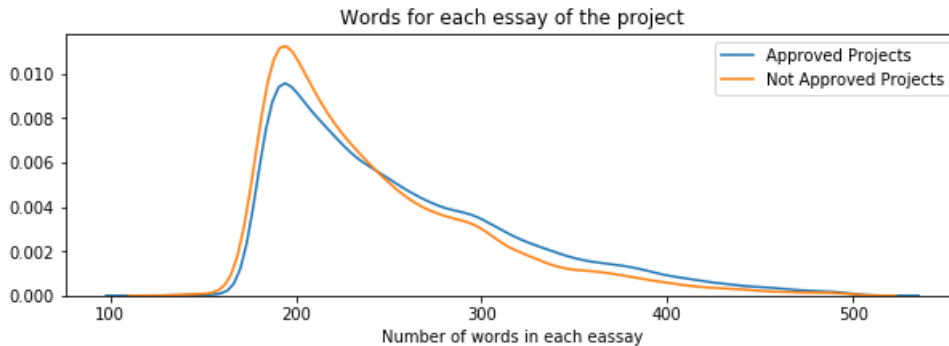
**Summary:**
1. Projects which have less than 120 words or more than 430 words in each essay are not approved.
2. Most of the approved project have less than 300 words each in essay

```
plt.figure(figsize=(10,3))
sns.distplot(approved_word_count, hist=False, label="Approved Projects")
sns.distplot(rejected_word_count, hist=False, label="Not Approved Projects")
plt.title('Words for each essay of the project')
plt.xlabel('Number of words in each eassay')
plt.legend()
plt.show()
```



**Summary:**
**1. Approved projects have more number of words in essay as compared to the rejected projects.**

## 1.2.8 Univariate Analysis: Cost per project

In [41]:

```
# we get the cost of the project using resource.csv file
resource_data.head(2)
```

Out[41]:

| | 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 |

In [42]:

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

project_data.columns
```

Out[42]:

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

In [43]:

```
# join two dataframes in python:
project_data = pd.merge(project_data, price_data, on='id', how='left')
```
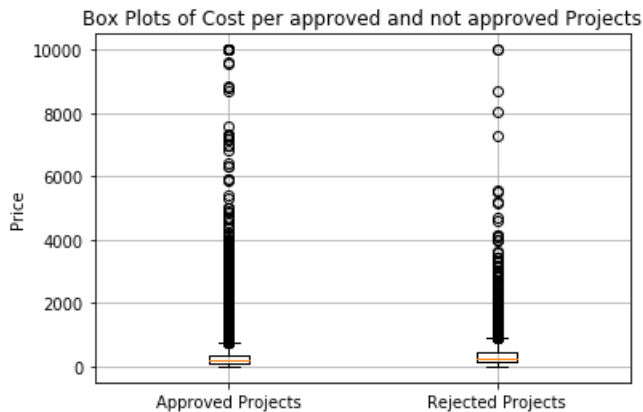
In [44]:

```python
approved_price = project_data[project_data['project_is_approved']==1]['price'].values

rejected_price = project_data[project_data['project_is_approved']==0]['price'].values
```
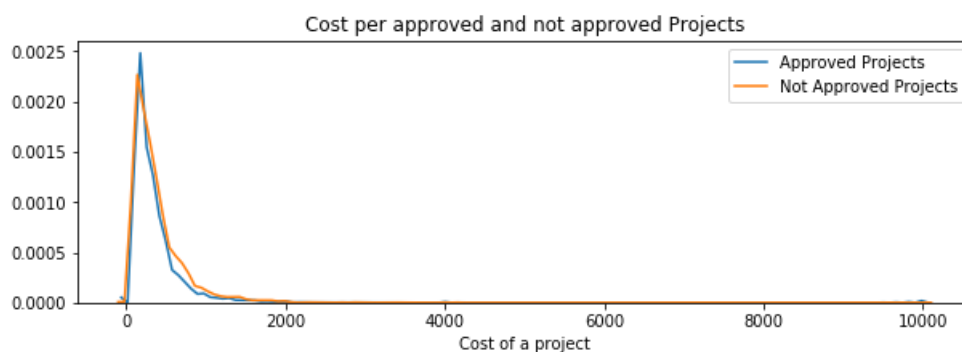
In [45]:

```python
# https://glowingpython.blogspot.com/2012/09/boxplot-with-matplotlib.html
plt.boxplot([approved_price, rejected_price])
plt.title('Box Plots of Cost per approved and not approved Projects')
plt.xticks([1,2],('Approved Projects','Rejected Projects'))
plt.ylabel('Price')
plt.grid()
plt.show()
```



In [46]:

```python
plt.figure(figsize=(10,3))
sns.distplot(approved_price, hist=False, label="Approved Projects")
sns.distplot(rejected_price, hist=False, label="Not Approved Projects")
plt.title('Cost per approved and not approved Projects')
plt.xlabel('Cost of a project')
plt.legend()
plt.show()
```



In [47]:

```python
# http://zetcode.com/python/prettytable/
from prettytable import PrettyTable

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

x = PrettyTable()
x.field_names = ["Percentile", "Approved Projects", "Not Approved Projects"]

for i in range(0,101,5):
    x.add_row([i,np.round(np.percentile(approved_price,i),3), np.round(np.percentile(rejected_price
,i), 3)])
print(x)
```

```
+------------+-------------------+-----------------------+
| Percentile | Approved Projects | Not Approved Projects |
+------------+-------------------+-----------------------+
```

```
|      0      |        0.66        |         1.97          |
|      5      |       13.59        |         41.9          |
|     10      |       33.88        |         73.67         |
|     15      |        58.0        |        99.109         |
|     20      |       77.38        |        118.56         |
|     25      |       99.95        |        140.892        |
|     30      |       116.68       |        162.23         |
|     35      |      137.232       |        184.014        |
|     40      |       157.0        |        208.632        |
|     45      |      178.265       |        235.106        |
|     50      |       198.99       |        263.145        |
|     55      |       223.99       |        292.61         |
|     60      |       255.63       |        325.144        |
|     65      |      285.412       |        362.39         |
|     70      |      321.225       |        399.99         |
|     75      |      366.075       |        449.945        |
|     80      |       411.67       |        519.282        |
|     85      |       479.0        |        618.276        |
|     90      |       593.11       |        739.356        |
|     95      |      801.598       |        992.486        |
|    100      |       9999.0       |        9999.0         |
+-------------+--------------------+-----------------------+
```

In [48]:

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

In [49]:

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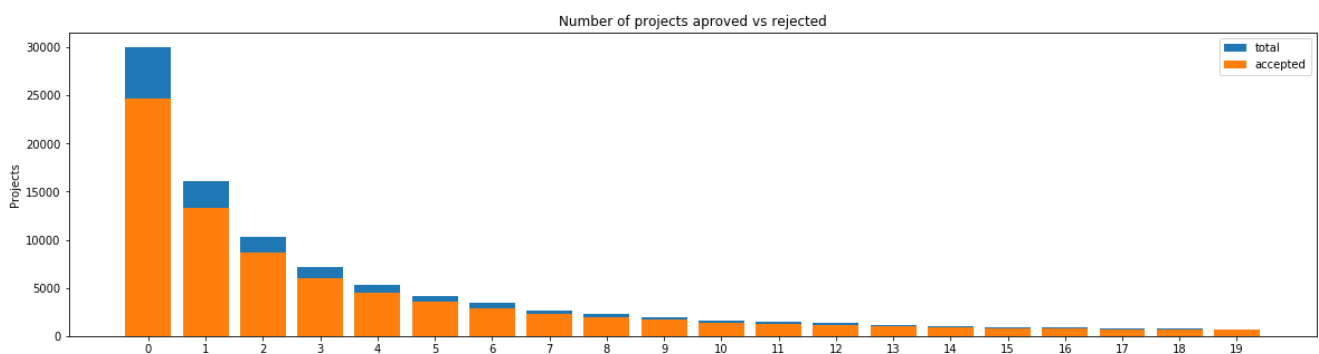
**Summary:**

**1. Projects which have more than 9999 budget will not be approved.**

**2. Approved project always have lower cost when compared to rejected project within same percentile range.**

## 1.2.9 Univariate Analysis: teacher_number_of_previously_posted_projects

In [50]:

```python
univariate_barplots(project_data, 'teacher_number_of_previously_posted_projects',
'project_is_approved', top=20)
```



```
   teacher_number_of_previously_posted_projects  project_is_approved  total  \
0                                             0                24652  30014
1                                             1                13329  16058
2                                             2                 8705  10350
3                                             3                 5997   7110
4                                             4                 4452   5266

      Avg
```

```
0  0.821350
1  0.830054
2  0.841063
3  0.843460
4  0.845423
==================================================
    teacher_number_of_previously_posted_projects  project_is_approved  total  \
15                                            15                   818    942
16                                            16                   769    894
17                                            17                   712    803
18                                            18                   666    772
19                                            19                   632    710

          Avg
15  0.868365
16  0.860179
17  0.886675
18  0.862694
19  0.890141
```

**Summary:**
**1. Rate of approval increase if the number of previously posted project increases.**
**2. Prior posting of project is no mandatory to get the approval the project. Nearly 82% proejcts are approved even if the teacher has not posted any project previously.**

In [51]:

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

In [52]:

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

### 1.2.10 Univariate Analysis: project_resource_summary

In [53]:

```python
#https://stackoverflow.com/questions/49773722/check-whether-a-dataframe-or-ndrray-contains-digits

project_data['isnumber'] = project_data['project_resource_summary'].str.findall('(\d+)').dropna().a
stype(bool)


count1=project_data['isnumber'].value_counts()
true=count1[1]

#https://stackoverflow.com/questions/27474921/compare-two-columns-using-pandas
project_data["effect"] = np.where((project_data["project_is_approved"]==1) & (project_data["isnumbe
r"]==True),'yes','no')


count=project_data["effect"].value_counts()
yes=count[1]

print("Total number of project posted - ",len(project_data['id']))
print("total number of project approved - ",len(project_data[project_data['project_is_approved']==
1]))
print("No. of times numerical value is used in resource summary -",true)
print("No. of times the project is approved when resource summary has numerical value -",yes)

approved_per = yes/true*100
print("Approved percentage when resource summary has numerical value -",approved_per)
```

```
approved_words = project_data[project_data['project_is_approved']==1]['project_resource_summary'].s
tr.split().apply(len)
approved_words = approved_words.values

rejected_words = project_data[project_data['project_is_approved']==0]['project_resource_summary'].s
tr.split().apply(len)
rejected_words = rejected_words.values

plt.boxplot([approved_words, rejected_words])
plt.title('Words for resurce summary of the project')
plt.xticks([1,2],('Approved Projects','Rejected Projects'))
plt.ylabel('Words in project resource summary')
plt.grid()
plt.show()

#Please do this on your own based on the data analysis that was done in the above cells

#Check if the `presence of the numerical digits` in the `project_resource_summary` effects the acc
eptance of the project or not. If you observe that `presence of the numerical digits` is helpful i
n the classification, please include it for further process or you can ignore it.
```
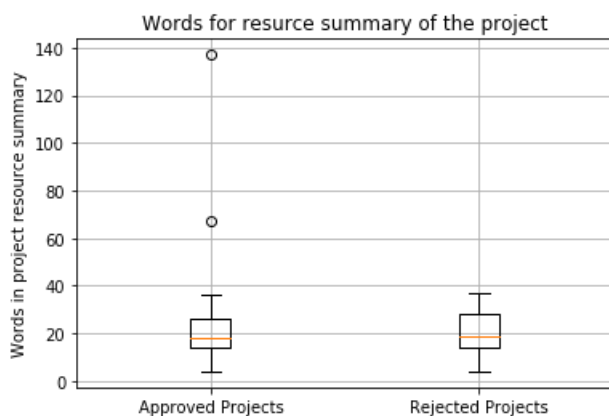
```
Total number of project posted -  109248
total number of project approved -  92706
No. of times numerical value is used in resource summary - 15756
No. of times the project is approved when resource summary has numerical value - 14090
Approved percentage when resource summary has numerical value - 89.42625031733942
```



**Summary:**
**1. Presence of Numerical value in the project resource summary has 89% of approval rate but still only 14% of posted project have numerical value in them.**

## 1.3 Text preprocessing

### 1.3.1 Essay Text

In [54]:

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

Out[54]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_cate |
|---|---|---|---|---|---|---|---|
| **0** | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | 2016-12-05 13:43:57 | Grades P |

| | Unnamed: 1 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_cate |
|---|---|---|---|---|---|---|---|
| 1 | 140946 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 2016-10-25 09:22:10 | Grade |

2 rows × 22 columns

In [55]:

```
# printing some random essays.
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)

#print(project_data['essay'].values[20000])
#print("="*50)
#print(project_data['essay'].values[99999])
#print("="*50)
```

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. \r\n\r\n We have over 24 languages represented in our English Learner program with students at every level of mastery.  We also have over 40 countries represented with the families within our school.  Each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures, beliefs, and respect.\"The limits of your language are the limits of your world.\"-Ludwig Wittgenstein  Our English learner's have a strong support system at home that begs for more resources.  Many times our parents are learning to read and speak English along side of their children.  Sometimes this creates barriers for parents to be able to help their child learn phonetics, letter recognition, and other reading skills.\r\n\r\nBy providing these dvd's and players, students are able to continue their mastery of the English language even if no one at home is able to assist.  All families with students within the Level 1 proficiency status, will be a offered to be a part of this program.  These educational videos will be specially chosen by the English Learner Teacher and will be sent home regularly to watch.  The 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 opportunity to check out a dvd player to use for the year.  The plan is to use these videos and educational dvd's for the years to come for other EL students.\r\nnannan
==================================================
The 51 fifth grade students that will cycle through my classroom this year all love learning, at least most of the time. At our school, 97.3% of the students receive free or reduced price lunch. Of the 560 students, 97.3% are minority students. \r\nThe school has a vibrant community that loves to get together and celebrate. Around Halloween there is a whole school parade to show off the beautiful costumes that students wear. On Cinco de Mayo we put on a big festival with crafts made by the students, dances, and games. At the end of the year the school hosts a carnival to celebrate the hard work put in during the school year, with a dunk tank being the most popular activity.My students will use these five brightly colored Hokki stools in place of regular, stationary, 4-legged chairs. As I will only have a total of ten in the classroom and not enough for each student to have an individual one, they will be used in a variety of ways. During independent reading time they will be used as special chairs students will each use on occasion. I will utilize them in place of chairs at my small group tables during math and reading times. The rest of the day they will be used by the students who need the highest amount of movement in their life in order to stay focused on school.\r\n\r\nWhenever asked what the classroom is missing, my students always say more Hokki Stools. They can't get their fill of the 5 stools we already have. When the students are sitting in group with me on the Hokki Stools, they are always moving, but at the same time doing their work. Anytime the students get to pick where they can sit, the Hokki Stools are the first to be taken. There are always students who head over to the kidney table to get one of the stools who are disappointed as there are not enough of them. \r\n\r\nWe ask a lot of students to sit for 7 hours a day. The Hokki stools will be a compromise that allow my students to do desk work and move at the same time. These stools will help students to meet their 60 minutes a day of movement by allowing them to activate their core muscles for balance while they sit. For many of my students, these chairs will take away the barrier that exists in schools for a child who can't sit still.nannan
==================================================
How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day.\r\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas.\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an \"open classroom\" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more.With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets. I will be able to help create the mood in our classroom setting to

g decor and the blue fish nets, I will be able to help create the mood in our classroom setting to
be one of a themed nautical environment. Creating a classroom environment is very important in the
success in each and every child's education. The nautical photo props will be used with each child
as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pic
tures of each child with them, have them developed, and then hung in our classroom ready for their
first day of 4th grade.  This kind gesture will set the tone before even the first day of school!
The nautical thank you cards will be used throughout the year by the students as they create thank
you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our
classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs lost of money out of
my own pocket on resources to get our classroom ready. Please consider helping with this project t
o make our new school year a very successful one. Thank you!nannan
==================================================

In [56]:

```python
# https://stackoverflow.com/a/47091490/4084039
import re

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

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

In [57]:

```python
sent = decontracted(project_data['essay'].values[1000])
print(sent)
print("="*50)
```

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of
desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work
hard to create a warm inviting themed room for my students look forward to coming to each day.\r\n
\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed races in
Arkansas.\r\nThey attend a Title I school, which means there is a high enough percentage of free a
nd reduced-price lunch to qualify. Our school is an \"open classroom\" concept, which is very uniq
ue as there are no walls separating the classrooms. These 9 and 10 year-old students are very eage
r learners; they are like sponges, absorbing all the information and experiences and keep on wanti
ng more.With these resources such as the comfy red throw pillows and the whimsical nautical hangin
g decor and the blue fish nets, I will be able to help create the mood in our classroom setting to
be one of a themed nautical environment. Creating a classroom environment is very important in the
success in each and every child is education. The nautical photo props will be used with each
child as they step foot into our classroom for the first time on Meet the Teacher evening. I will
take pictures of each child with them, have them developed, and then hung in our classroom ready f
or their first day of 4th grade.  This kind gesture will set the tone before even the first day of
school! The nautical thank you cards will be used throughout the year by the students as they crea
te thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make
our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs lost of money ou
t of my own pocket on resources to get our classroom ready. Please consider helping with this proj
ect to make our new school year a very successful one. Thank you!nannan
==================================================

In [58]:

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

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of
desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work
hard to create a warm inviting themed room for my students look forward to coming to each day.

My class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas. They attend a Title I school, which means there is a high enough percentage of free and reduced-price l unch to qualify. Our school is an open classroom concept, which is very unique as there are no w alls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more.With these re sources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed n autical environment. Creating a classroom environment is very important in the success in each and every child is education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I will take pictures of each ch ild with them, have them developed, and then hung in our classroom ready for their first day of 4t h grade. This kind gesture will set the tone before even the first day of school! The nautical th ank you cards will be used throughout the year by the students as they create thank you cards to t heir team groups. Your generous donations will help me to help make our classroom a fun, inviting, learning environment from day one. It costs lost of money out of my own pocket on res ources to get our classroom ready. Please consider helping with this project to make our new schoo l year a very successful one. Thank you!nannan

In [59]:

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

How do you remember your days of school Was it in a sterile environment with plain walls rows of d esks and a teacher in front of the room A typical day in our room is nothing like that I work hard to create a warm inviting themed room for my students look forward to coming to each day My class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas They attend a Title I school which means there is a high enough percentage of free and reduced price lunch to qualify Our school is an open classroom concept which is very unique as there are no walls separating the classrooms These 9 and 10 year old students are very eager learners they are like sponges absorbin g all the information and experiences and keep on wanting more With these resources such as the co mfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets I will be ab le to help create the mood in our classroom setting to be one of a themed nautical environment Cre ating a classroom environment is very important in the success in each and every child is education The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening I will take pictures of each child with t hem have them developed and then hung in our classroom ready for their first day of 4th grade This kind gesture will set the tone before even the first day of school The nautical thank you cards wi ll be used throughout the year by the students as they create thank you cards to their team groups Your generous donations will help me to help make our classroom a fun inviting learning environmen t from day one It costs lost of money out of my own pocket on resources to get our classroom ready Please consider helping with this project to make our new school year a very successful one Thank you nannan

In [60]:

```python
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', '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', "isn't", 'ma', 'mightn', "mightn't", 'mustn',\
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn'
```

```
          musun't , needn , needn't , shall , shall't , shouldn , shouldn't , wasn ,
"wasn't", 'weren', "weren't", \
          'won', "won't", 'wouldn', "wouldn't"]
```

In [61]:

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

```
100%|████████████████████████████████████████████████████| 109248/109248
[01:37<00:00, 1120.26it/s]
```

In [62]:

```python
# after preprocesing
preprocessed_essays[1000]
```

Out[62]:

```
'how remember days school was sterile environment plain walls rows desks teacher front room a typi
cal day room nothing like i work hard create warm inviting themed room students look forward comin
g day my class made 28 wonderfully unique boys girls mixed races arkansas they attend title i
school means high enough percentage free reduced price lunch qualify our school open classroom
concept unique no walls separating classrooms these 9 10 year old students eager learners like spo
nges absorbing information experiences keep wanting with resources comfy red throw pillows
whimsical nautical hanging decor blue fish nets i able help create mood classroom setting one them
ed nautical environment creating classroom environment important success every child education the
nautical photo props used child step foot classroom first time meet teacher evening i take
pictures child developed hung classroom ready first day 4th grade this kind gesture set tone even
first day school the nautical thank cards used throughout year students create thank cards team gr
oups your generous donations help help make classroom fun inviting learning environment day one it
costs lost money pocket resources get classroom ready please consider helping project make new sch
ool year successful one thank nannan'
```

### 1.3.2 Project title Text

In [63]:

```python
project_data['project_title'].head(5)
```

Out[63]:

```
0       Educational Support for English Learners at Home
1                 Wanted: Projector for Hungry Learners
2       Soccer Equipment for AWESOME Middle School Stu...
3                              Techie Kindergarteners
4                              Interactive Math Tools
Name: project_title, dtype: object
```

In [64]:

```python
from tqdm import tqdm
preprocessed_title = []
# tqdm is for printing the status bar
for sentance in tqdm(project_data['project_title'].values):
    sent = decontracted(sentance)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
```

```
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    preprocessed_title.append(sent.lower().strip())
```

In [65]:

```
preprocessed_title[2222]
```

Out[65]:

```
'seating for tablet station'
```

In [66]:

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

## 1. 4 Preparing data for models

In [67]:

```
project_data.columns
```

Out[67]:

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

we are going to consider

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

    - project_title : text data
    - text : text data
    - project_resource_summary: text data

    - quantity : numerical
    - teacher_number_of_previously_posted_projects : numerical
    - price : numerical
```

### 1.4.1 Vectorizing Categorical data

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

In [68]:

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


categories_one_hot = vectorizer.transform(project_data['clean_categories'].values)
print("Shape of matrix after one hot encodig ",categories_one_hot.shape)
```

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

In [69]:

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


sub_categories_one_hot = vectorizer.transform(project_data['clean_subcategories'].values)
print("Shape of matrix after one hot encodig ",sub_categories_one_hot.shape)
```

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular',
'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger',
'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other',
'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL
', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences',
'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
Shape of matrix after one hot encodig  (109248, 30)
```

In [70]:

```
vectorizer = CountVectorizer(vocabulary=list(school_state_dict.keys()), lowercase=False, binary=Tr
ue)
vectorizer.fit(project_data['school_state'].values)
print(vectorizer.get_feature_names())


school_state_one_hot = vectorizer.transform(project_data['school_state'].values)
print("Shape of matrix after one hot encodig ",school_state_one_hot.shape)



# Please do the similar feature encoding with state, teacher_prefix and project_grade_category als
o
```

```
['VT', 'WY', 'ND', 'MT', 'RI', 'SD', 'NE', 'DE', 'AK', 'NH', 'WV', 'ME', 'HI', 'DC', 'NM', 'KS', 'I
A', 'ID', 'AR', 'CO', 'MN', 'OR', 'KY', 'MS', 'NV', 'MD', 'CT', 'TN', 'UT', 'AL', 'WI', 'VA', 'AZ',
'NJ', 'OK', 'WA', 'MA', 'LA', 'OH', 'MO', 'IN', 'PA', 'MI', 'SC', 'GA', 'IL', 'NC', 'FL', 'NY', 'TX
', 'CA']
Shape of matrix after one hot encodig  (109248, 51)
```

In [71]:

```
vectorizer = CountVectorizer(vocabulary=list(project_grade_category_dict.keys()), lowercase=False,
binary=True)
vectorizer.fit(project_data['project_grade_category'].values)
print(vectorizer.get_feature_names())
```

```
project_grade_category_one_hot =
vectorizer.transform(project_data['project_grade_category'].values)
print("Shape of matrix after one hot encodig ",project_grade_category_one_hot.shape)
```

```
['Grades 9-12', 'Grades 6-8', 'Grades 3-5', 'Grades PreK-2']
Shape of matrix after one hot encodig  (109248, 4)
```

In [72]:

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

#https://stackoverflow.com/questions/39303912/tfidfvectorizer-in-scikit-learn-valueerror-np-nan-is
-an-invalid-document
vectorizer.fit(project_data['teacher_prefix'].values.astype(str))
print(vectorizer.get_feature_names())


teacher_prefix_one_hot = vectorizer.transform(project_data['teacher_prefix'].values.astype(str))
print("Shape of matrix after one hot encodig ",teacher_prefix_one_hot.shape)
```

```
['Dr.', 'Teacher', 'Mr.', 'Ms.', 'Mrs.']
Shape of matrix after one hot encodig  (109248, 5)
```

## 1.4.2 Vectorizing Text data

### 1.4.2.2 Bag of Words on `Text`

In [73]:

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

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

### 1.4.2.2 Bag of Words on `project_title`

In [74]:

```
vectorizer = CountVectorizer(min_df=10)
title_bow = vectorizer.fit_transform(preprocessed_title)
print("Shape of matrix after one hot encodig ",title_bow.shape)
```

```
Shape of matrix after one hot encodig  (109248, 3329)
```

### 1.4.2.3 Bag of Words on `project_resource_summary`

In [75]:

```
vectorizer = CountVectorizer(min_df=10)
project_resource_summary_bow = vectorizer.fit_transform(preprocessed_project_resource_summary)
print("Shape of matrix after one hot encodig ",project_resource_summary_bow.shape)
```

```
Shape of matrix after one hot encodig  (109248, 5797)
```

### 1.4.2.3 TFIDF vectorizer

In [76]:

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

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

### 1.4.2.4 TFIDF Vectorizer on `project_title`

In [77]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
title_tfidf = vectorizer.fit_transform(preprocessed_title)
print("Shape of matrix after one hot encodig ",title_tfidf.shape)


# Similarly you can vectorize for title also
```

Shape of matrix after one hot encodig  (109248, 3329)

### 1.4.2.5 Using Pretrained Models: Avg W2V

In [78]:

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

# ===========================
Output:

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

# ===========================

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

for i in preproced_titles:
    words.extend(i.split(' '))
print("all the words in the coupus", len(words))
words = set(words)
print("the unique words in the coupus", len(words))

inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove vectors and our coupus", \
      len(inter_words),"(",np.round(len(inter_words)/len(words)*100,3),"%)")

words_courpus = {}
words_glove = set(model.keys())
for i in words:
    if i in words_glove:
        words_courpus[i] = model[i]
print("word 2 vec length", len(words_courpus))
```

```python
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/

import pickle
with open('glove_vectors', 'wb') as f:
    pickle.dump(words_courpus, f)


'''
```

Out[78]:

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef
loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\'r\',
encoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\n
word = splitLine[0]\n        embedding = np.array([float(val) for val in splitLine[1:]])\n        m
odel[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
loadGloveModel(\'glove.42B.300d.txt\')\n\n# ===============================\nOutput:\n    \nLoading G
love Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495  words loaded!\n\n#
===========================\n\nwords = []\nfor i in preproced_texts:\n    words.extend(i.split(\'
\'))\n\nfor i in preproced_titles:\n    words.extend(i.split(\' \'))\nprint("all the words in the
coupus", len(words))\nwords = set(words)\nprint("the unique words in the coupus",
len(words))\n\ninter_words = set(model.keys()).intersection(words)\nprint("The number of words tha
t are present in both glove vectors and our coupus",        len(inter_words),"
(",np.round(len(inter_words)/len(words)*100,3),"%)")\n\nwords_courpus = {}\nwords_glove =
set(model.keys())\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\n
print("word 2 vec length", len(words_courpus))\n\n\n# stronging variables into pickle files python
: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport pic
kle\nwith open(\'glove_vectors\', \'wb\') as f:\n    pickle.dump(words_courpus, f)\n\n\n'
```

In [79]:

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

In [80]:

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

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

```
100%|████████████████████████████████████████████████████| 109248/109248
[00:55<00:00, 1978.31it/s]
```

```
109248
300
```

### 1.4.2.6 Using Pretrained Models: AVG W2V on `project_title`

In [81]:

```python
avg_w2v_vectors_title = []; # the avg-w2v for each title is stored in this list
for sentence in tqdm(preprocessed_title): # for each title
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words =0; # num of words with a valid vector in the title
```

```
    cnt_words =0; # num of words with a valid vector in the title
    for word in sentence.split(): # for each word in a title
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_title.append(vector)

print(len(avg_w2v_vectors_title))
print(len(avg_w2v_vectors_title[0]))


# Similarly you can vectorize for title also
```

```
100%|████████████████████████████████████████████████████████| 109248/109248
[00:02<00:00, 39184.35it/s]
```

```
109248
300
```

### 1.4.2.7 Using Pretrained Models: TFIDF weighted W2V

In [82]:

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

In [83]:

```python
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors.append(vector)

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

```
100%|████████████████████████████████████████████████████████| 109248/109248
[07:02<00:00, 258.44it/s]
```

```
109248
300
```

### 1.4.2.9 Using Pretrained Models: TFIDF weighted W2V on `project_title`

In [84]:

```python
#Similarly you can vectorize for title also

tfidf_model = TfidfVectorizer()
tfidf_model.fit(preprocessed_title)
```

```
tfidf_model.fit(preprocessed_title)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())
```

In [85]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_title): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_title.append(vector)

print(len(tfidf_w2v_vectors_title))
print(len(tfidf_w2v_vectors_title[0]))
```

```
100%|████████████████████████████████████████████████████████████████| 109248/109248
[00:06<00:00, 17302.50it/s]
```

```
109248
300
```

### 1.4.3 Vectorizing Numerical features

In [86]:

```
# check this one: https://www.youtube.com/watch?v=0HOqOcln3Z4&t=530s
# standardization sklearn: https://scikit-
learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html
from sklearn.preprocessing import StandardScaler

# price_standardized = standardScalar.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329.   ... 399.   287.
73   5.5 ].
# Reshape your data either using array.reshape(-1, 1)

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

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

```
Mean : 298.1193425966608, Standard deviation : 367.49634838483496
```

In [87]:

```
price_standardized
```

Out[87]:

```
array([[-0.3905327 ],
       [ 0.00239637],
       [ 0.59519138],
       ...,
       [-0.15825829],
```

```
         [-0.61243967],
         [-0.51216657]])
```

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

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

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

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

Mean : 16.965610354422964, Standard deviation : 26.182821919093175

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

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

```
quantity_standardized
```

```
array([[ 0.23047132],
       [-0.60977424],
       [ 0.19227834],
       ...,
       [-0.4951953 ],
       [-0.03687954],
       [-0.45700232]])
```

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

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

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

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

Mean : 11.153165275336848, Standard deviation : 27.77702641477403

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

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

```
teacher_number_of_previously_posted_projects_standardized
```

```
Out[91]:

array([[-0.40152481],
       [-0.14951799],
       [-0.36552384],
       ...,
       [-0.29352189],
       [-0.40152481],
       [-0.40152481]])
```

### 1.4.4 Merging all the above features

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

```
In [92]:

print(school_state_one_hot.shape)
print(project_grade_category_one_hot.shape)
print(teacher_prefix_one_hot.shape)
print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(text_bow.shape)
print(title_bow.shape)
print(project_resource_summary_bow.shape)
print(price_standardized.shape)
print(quantity_standardized.shape)
print(teacher_number_of_previously_posted_projects_standardized.shape)
```

```
(109248, 51)
(109248, 4)
(109248, 5)
(109248, 9)
(109248, 30)
(109248, 16623)
(109248, 3329)
(109248, 5797)
(109248, 1)
(109248, 1)
(109248, 1)
```

```
In [93]:

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

```
Out[93]:

(109248, 16663)
```

# Assignment 2: Apply TSNE

If you are using any code snippet from the internet, you have to provide the reference/citations, as we did in the above cells. Otherwise, it will be treated as plagiarism without citations.

1. In the above cells we have plotted and analyzed many features. Please observe the plots and write the observations in markdown cells below every plot.
2. EDA: Please complete the analysis of the feature: teacher_number_of_previously_posted_projects
3. Build the data matrix using these features
   - school_state : categorical data (one hot encoding)
   - clean_categories : categorical data (one hot encoding)
   - clean_subcategories : categorical data (one hot encoding)
   - teacher_prefix : categorical data (one hot encoding)
   - project_grade_category : categorical data (one hot encoding)

- project_grade_category : categorical data (one hot encoding)
- project_title : text data (BOW, TFIDF, AVG W2V, TFIDF W2V)
- price : numerical
- teacher_number_of_previously_posted_projects : numerical

4. Now, plot FOUR t-SNE plots with each of these feature sets.
   A. categorical, numerical features + project_title(BOW)
   B. categorical, numerical features + project_title(TFIDF)
   C. categorical, numerical features + project_title(AVG W2V)
   D. categorical, numerical features + project_title(TFIDF W2V)
5. Concatenate all the features and Apply TNSE on the final data matrix
6. Note 1: The TSNE accepts only dense matrices
7. Note 2: Consider only 5k to 6k data points to avoid memory issues. If you run into memory error issues, reduce the number of data points but clearly state the number of datat-poins you are using

In [94]:

```python
# this is the example code for TSNE
import numpy as np
from sklearn.manifold import TSNE
from sklearn import datasets
import pandas as pd
import matplotlib.pyplot as plt

iris = datasets.load_iris()
x = iris['data']
y = iris['target']

tsne = TSNE(n_components=2, perplexity=30, learning_rate=200)

X_embedding = tsne.fit_transform(x)
# if x is a sparse matrix you need to pass it as X_embedding = tsne.fit_transform(x.toarray()) , .
toarray() will convert the sparse matrix into dense matrix

for_tsne = np.hstack((X_embedding, y.reshape(-1,1)))

for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x','Dimension_y','Score'])
colors = {0:'red', 1:'blue', 2:'green'}
plt.scatter(for_tsne_df['Dimension_x'], for_tsne_df['Dimension_y'], c=for_tsne_df['Score'].apply(lambda x: colors[x]))
plt.show()
```



## 2.1 TSNE with `BOW` encoding of `project_title` feature (5000 datapoints)

In [95]:

```python
import numpy as np
from sklearn.manifold import TSNE
from sklearn import datasets
import pandas as pd
import matplotlib.pyplot as plt

X1 = hstack((school_state_one_hot, categories_one_hot, sub_categories_one_hot,
project_grade_category_one_hot
          , teacher_prefix_one_hot, title_bow, price_standardized,
teacher_number_of_previously_posted_projects_standardized, quantity_standardized))
```

```
X1 = X1.toarray()
X1 = X1[0:5000,:]

Y1 = project_data['project_is_approved']
Y1 = Y1[0:5000]
```
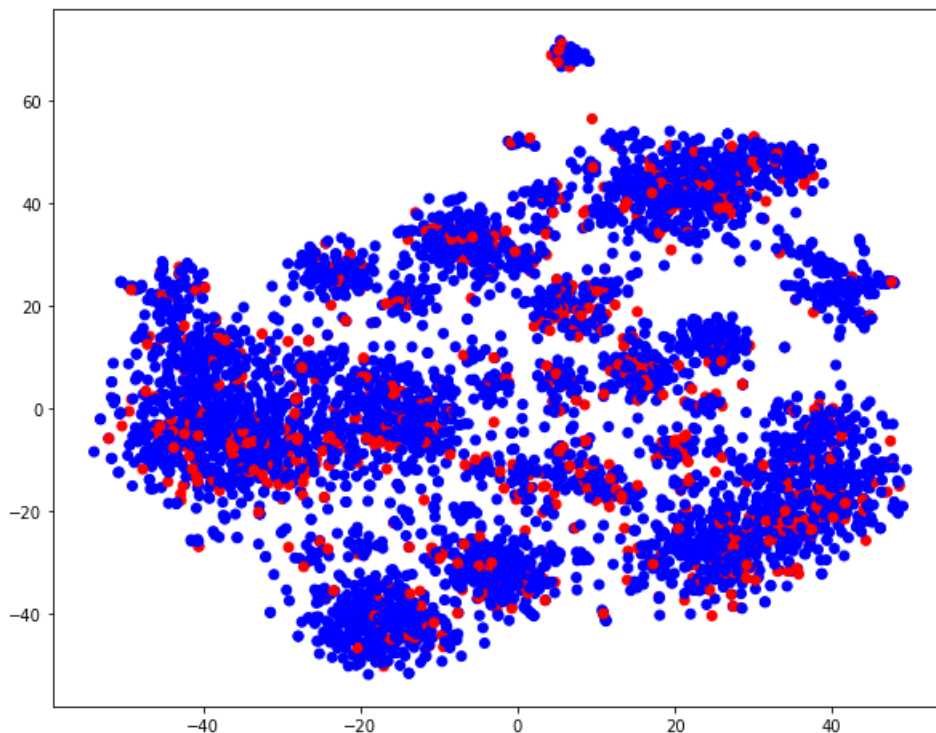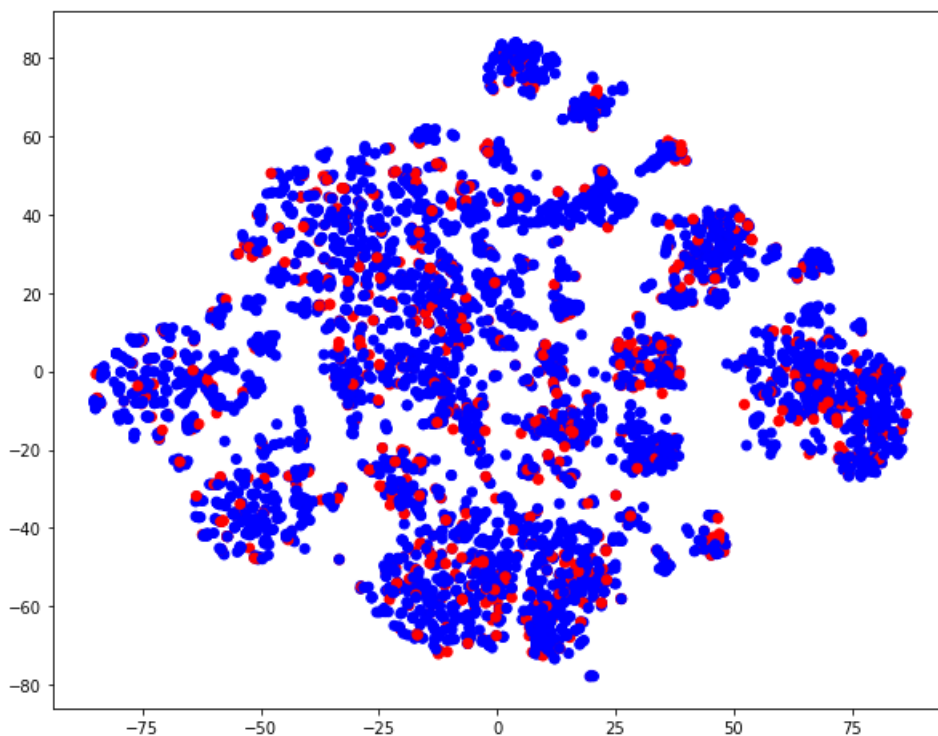
```
model = TSNE(n_components=2, perplexity=30, learning_rate=200)
tsne_data = model.fit_transform(X1)
```

```
tsne_data = np.vstack((tsne_data.T, Y1)).T
tsne_df = pd.DataFrame(tsne_data, columns = ("1st_Dim","2nd_Dim","Labels"))
```

```
plt.figure(figsize=(10, 8))
plt.scatter(tsne_df['1st_Dim'], tsne_df['2nd_Dim'], c=tsne_df['Labels'].apply(lambda x: colors[x]))
plt.show()
```



## 2.2 TSNE with `TFIDF` encoding of `project_title` feature (5000 Data points)

```
X2 = hstack((school_state_one_hot, categories_one_hot, sub_categories_one_hot,
project_grade_category_one_hot
         , teacher_prefix_one_hot, title_tfidf, price_standardized,
teacher_number_of_previously_posted_projects_standardized, quantity_standardized))

X2 = X2.toarray()
X2 = X2[0:5000,:]

Y1 = project_data['project_is_approved']
Y1 = Y1[0:5000]
```
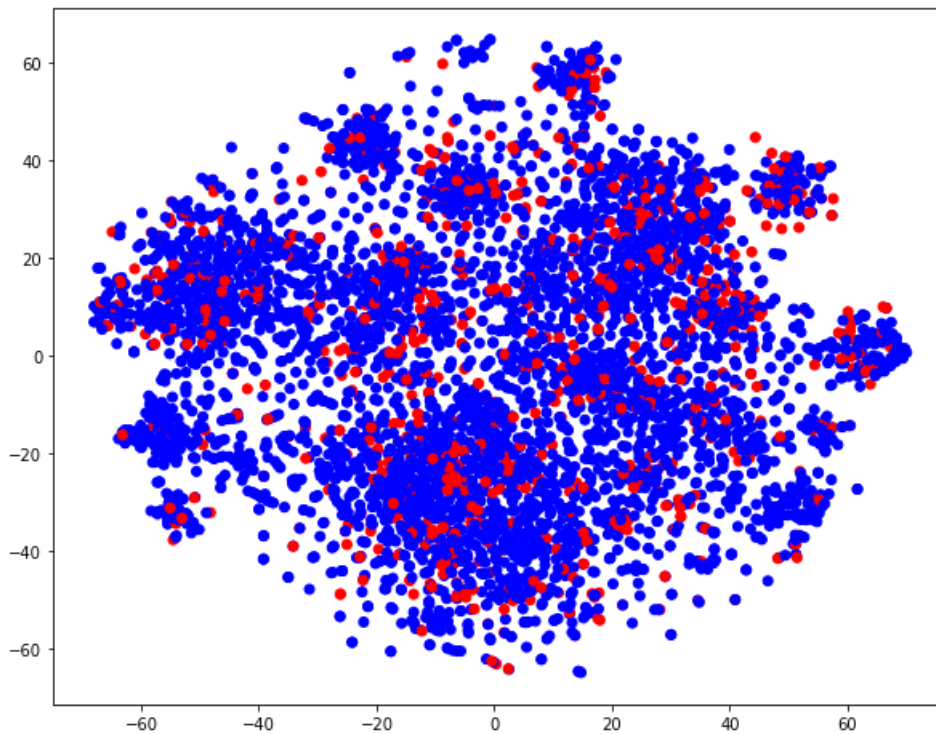
```
model = TSNE(n_components=2, perplexity=30, learning_rate=200)
tsne_data2 = model.fit_transform(X2)
```

In [102]:

```
tsne_data2 = np.vstack((tsne_data2.T, Y1)).T
tsne_df2 = pd.DataFrame(tsne_data2, columns = ("1st_Dim","2nd_Dim","Labels"))
plt.figure(figsize=(10, 8))
plt.scatter(tsne_df2['1st_Dim'], tsne_df2['2nd_Dim'], c=tsne_df2['Labels'].apply(lambda x: colors[x
]))
plt.show()
```



## 2.3 TSNE with `AVG W2V` encoding of `project_title` feature (5000 Data points)

In [103]:

```
X3 = hstack((school_state_one_hot, categories_one_hot, sub_categories_one_hot,
project_grade_category_one_hot
            , teacher_prefix_one_hot, avg_w2v_vectors_title, price_standardized,
teacher_number_of_previously_posted_projects_standardized, quantity_standardized))

X3 = X3.toarray()
X3 = X3[0:5000,:]

Y1 = project_data['project_is_approved']
Y1 = Y1[0:5000]
```
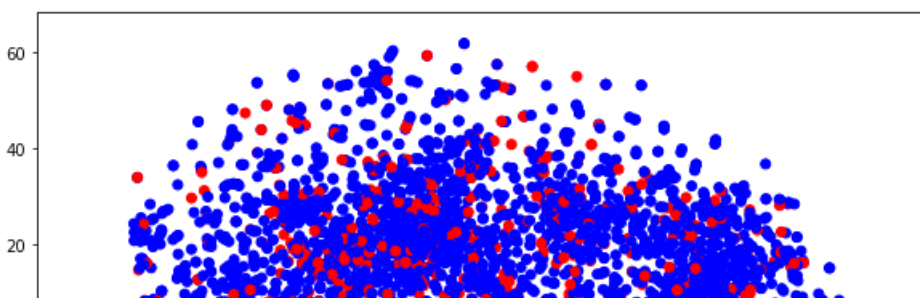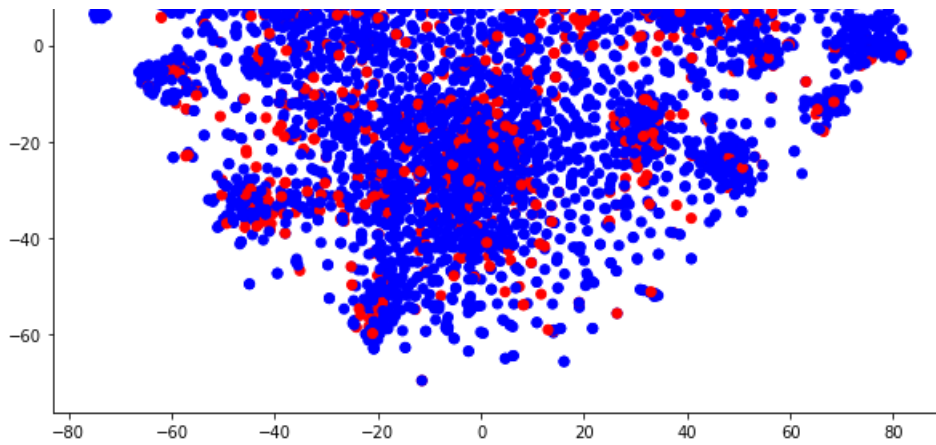
In [104]:

```
model = TSNE(n_components=2, perplexity=30, learning_rate=200)
tsne_data3 = model.fit_transform(X3)
```

In [105]:

```
tsne_data3 = np.vstack((tsne_data3.T, Y1)).T
tsne_df3 = pd.DataFrame(tsne_data3, columns = ("1st_Dim","2nd_Dim","Labels"))
plt.figure(figsize=(10, 8))
plt.scatter(tsne_df3['1st_Dim'], tsne_df3['2nd_Dim'], c=tsne_df3['Labels'].apply(lambda x: colors[x
]))
plt.show()
```

## 2.4 TSNE with `TFIDF Weighted W2V` encoding of `project_title` feature (5000 Data points)

```
X4 = hstack((school_state_one_hot, categories_one_hot, sub_categories_one_hot,
project_grade_category_one_hot
            , teacher_prefix_one_hot, tfidf_w2v_vectors_title, price_standardized,
teacher_number_of_previously_posted_projects_standardized, quantity_standardized))

X4 = X4.toarray()
X4 = X4[0:5000,:]

Y1 = project_data['project_is_approved']
Y1 = Y1[0:5000]
```

```
model = TSNE(n_components=2, perplexity=30, learning_rate=200)
tsne_data4 = model.fit_transform(X4)
```

```
tsne_data4 = np.vstack((tsne_data4.T, Y1)).T
tsne_df4 = pd.DataFrame(tsne_data4, columns = ("1st_Dim","2nd_Dim","Labels"))
plt.figure(figsize=(10, 8))
plt.scatter(tsne_df4['1st_Dim'], tsne_df4['2nd_Dim'], c=tsne_df4['Labels'].apply(lambda x: colors[x
]))
plt.show()
```

## 2.5 TSNE with `BOW`, `TFIDF`, `AVG W2V`, `TFIDF Weighted W2V` encoding of `project_title` feature (5000 Data points)

In [109]:

```
X5 = hstack((school_state_one_hot, categories_one_hot, sub_categories_one_hot,
project_grade_category_one_hot
            , teacher_prefix_one_hot, title_bow, title_tfidf, avg_w2v_vectors_title,
tfidf_w2v_vectors_title,
            price_standardized, teacher_number_of_previously_posted_projects_standardized, quantit
y_standardized))

X5 = X5.toarray()
X5 = X5[0:5000,:]

Y1 = project_data['project_is_approved']
Y1 = Y1[0:5000]
```
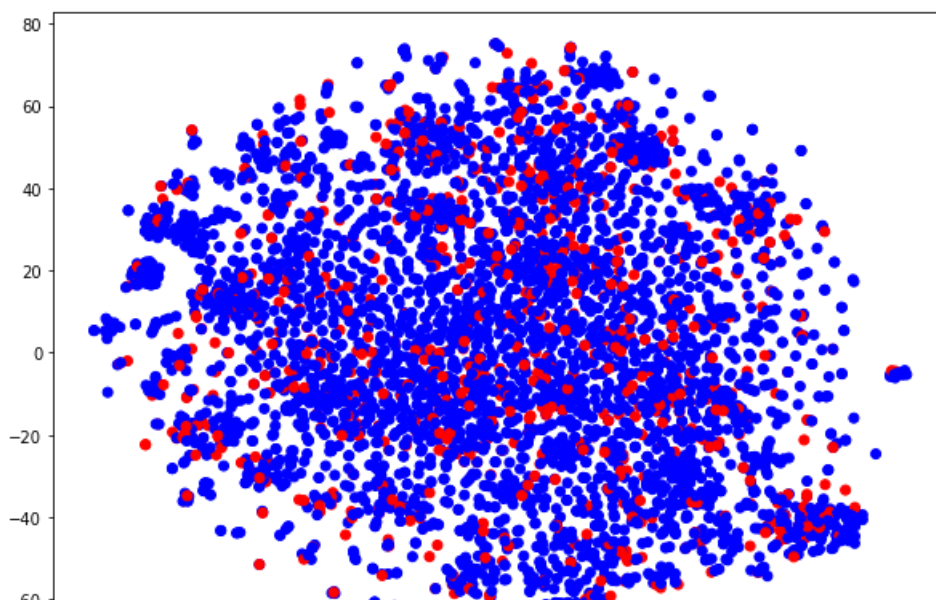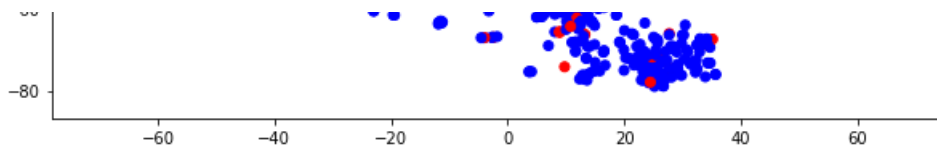
In [110]:

```
model = TSNE(n_components=2, perplexity=30, learning_rate=200)
tsne_data5 = model.fit_transform(X5)
```

In [111]:

```
tsne_data5 = np.vstack((tsne_data5.T, Y1)).T
tsne_df5 = pd.DataFrame(tsne_data5, columns = ("1st_Dim","2nd_Dim","Labels"))
plt.figure(figsize=(10, 8))
plt.scatter(tsne_df5['1st_Dim'], tsne_df5['2nd_Dim'], c=tsne_df5['Labels'].apply(lambda x: colors[x
]))
plt.show()
```

## 2.6 Summary

All these above visualization of TSNE with Bag of Words, TF-IDF, Avg Word2Vec, TF-IDF Weighted Word2Vec doesn't seem to show any clear picture of clustering of similar points

## 2.7 Conclusion

1. Every state has greater than 80% success rate in approval
2. California has the highest number of approved projects
3. If we do not take Dr. into consideration due to very less number of projects as compared to other prefixes, then teachers who have a proper prefix like Mr.,Miss,Mrs. get the approval rate of more than 84% as compared to when the prefix is just "teacher" which has 79% rate.
4. As the Grades increases the number of projects submitted decreases
5. Project approval rate of maths and science can be increased if they are combined with Literacy language
6. There are no projects with words greater than 11 in the title
7. 75% of the approved projects have less than 7 words in the title.
8. Projects which have less than 120 words or more than 430 words in each essay are not approved.
9. Projects which have more than 9999 budget will not be approved.
10. Rate of approval increase if the number of previously posted project increases.
11. It is not mandatory for a teacher to have previously posted projects to get the approval their new project.
12. Presence of Numerical value in the project resource summary has 89% of approval rate but still only 14% of posted project have numerical value in them.
13. All these above visualization of TSNE with Bag of Words, TF-IDF, Avg Word2Vec, TF-IDF Weighted Word2Vec doesn't seem to show any clear picture of clustering of similar points.