# 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:** `Art Will Make You Happy!` `First Grade Fun` |
| **project_grade_category** | Grade level of students for which the project is targeted. One of the following enumerated values: `Grades PreK-2` `Grades 3-5` `Grades 6-8` `Grades 9-12` |
| **project_subject_categories** | One or more (comma-separated) subject categories for the project from the following enumerated list of values: `Applied Learning` `Care & Hunger` `Health & Sports` `History & Civics` `Literacy & Language` `Math & Science` `Music & The Arts` `Special Needs` `Warmth` **Examples:** `Music & The Arts` `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:** `Literacy` `Literature & Writing, Social Sciences` |
| **project_resource_summary** | An explanation of the resources needed for the project.**Example:** `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 [4]:

```
%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
```

```python
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 [5]:

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

In [6]:

```python
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 [7]:

```python
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[7]:

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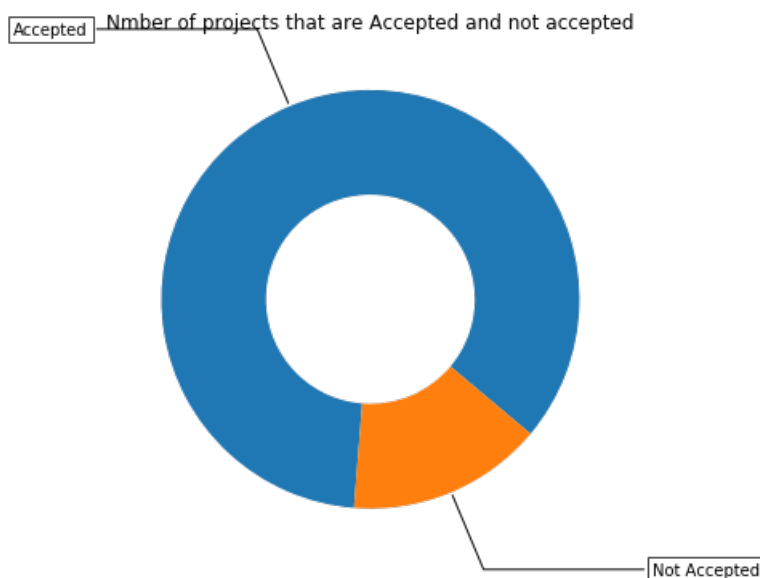
```
Number of projects thar are approved for funding  92706 , ( 84.85830404217927 %)
Number of projects thar are not approved for funding  16542 , ( 15.141695957820739 %)
```



### 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']

'''# How to plot US state heatmap: https://datascience.stackexchange.com/a/9620

scl = [[0.0, 'rgb(242,240,247)'],[0.2, 'rgb(218,218,235)'],[0.4, 'rgb(188,189,220)'],\
            [0.6, 'rgb(158,154,200)'],[0.8, 'rgb(117,107,177)'],[1.0, 'rgb(84,39,143)']]

data = [ dict(
        type='choropleth',
        colorscale = scl,
        autocolorscale = False,
        locations = temp['state_code'],
        z = temp['num_proposals'].astype(float),
        locationmode = 'USA-states',
        text = temp['state_code'],
        marker = dict(line = dict (color = 'rgb(255,255,255)',width = 2)),
        colorbar = dict(title = "% of pro")
    ) ]

layout = dict(
        title = 'Project Proposals % of Acceptance Rate by US States',
        geo = dict(
            scope='usa',
            projection=dict( type='albers usa' ),
            showlakes = True,
            lakecolor = 'rgb(255, 255, 255)',
        ),
    )

fig = go.Figure(data=data, layout=layout)
offline.iplot(fig, filename='us-map-heat-map')
'''
```

Out[9]:

```
'# How to plot US state heatmap: https://datascience.stackexchange.com/a/9620\n\nscl = [[0.0, \'rg
b(242,240,247)\'],[0.2, \'rgb(218,218,235)\'],[0.4, \'rgb(188,189,220)\'],          [0.6, \'rgb(1
58,154,200)\'],[0.8, \'rgb(117,107,177)\'],[1.0, \'rgb(84,39,143)\']]\n\ndata = [ dict(\n        ty
pe=\'choropleth\',\n        colorscale = scl,\n        autocolorscale = False,\n        locations =
temp[\'state_code\'],\n        z = temp[\'num_proposals\'].astype(float),\n        locationmode = \
'USA-states\',\n        text = temp[\'state_code\'],\n        marker = dict(line = dict (color = \'
rgb(255,255,255)\',width = 2)),\n        colorbar = dict(title = "% of pro")\n    ) ]\n\nlayout = d
ict(\n        title = \'Project Proposals % of Acceptance Rate by US States\',\n        geo = dict(
\n            scope=\'usa\',\n            projection=dict( type=\'albers usa\' ),\n            show
akes = True,\n            lakecolor = \'rgb(255, 255, 255)\',\n        ),\n    )\n\nfig =
go.Figure(data=data, layout=layout)\noffline.iplot(fig, filename=\'us-map-heat-map\')\n'
```

In [10]:

```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
```

## Observation:

1. Every state has more than 80% project approval rate.

In [11]:

```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 [12]:

```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 [13]:

```python
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
```

**SUMMARY: Every state has greater than 80% success rate in approval**

## 1.2.2 Univariate Analysis: teacher_prefix

In [14]:

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



```
   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
```

## 1.2.3 Univariate Analysis: project_grade_category

In [15]:

```
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
```

# Observation:

1. Every Project grade category has more than 83% approval rate

### 1.2.4 Univariate Analysis: project_subject_categories

In [16]:

```python
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 [17]:

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

Out[17]:

| | 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 |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 2016-10-25 09:22:10 | Grade |

In [18]:

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



Number of projects aproved vs rejected

```
      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
```

# Observation:

1. 92% of projects are approved with Warth and Care Hunger as categories

In [19]:

```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 [20]:

```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('% of projects aproved category wise')
plt.xticks(ind, list(sorted_cat_dict.keys()))
plt.show()
```



In [21]:

```python
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
```

# Observation:

1. Literacy_Language is the highly used category for the Approved projects.

### 1.2.5 Univariate Analysis: project_subject_subcategories

In [22]:

```python
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 [23]:

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

Out[23]:

| | 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 |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 2016-10-25 09:22:10 | Grade |

In [24]:

```python
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 |

In [25]:

```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_subcategories'].values:
    my_counter.update(word.split())
```

In [26]:

```python
# 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('% of projects aproved state wise')
plt.xticks(ind, list(sorted_sub_cat_dict.keys()))
plt.show()
```



In [27]:

```python
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
```

# Observation:

1. Literacy is highly used sub category for Approved prjects

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

In [28]:

```python
#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()
```



# Observation:

1. Large number of projects has 4 words title.
2. Very few projects have title words more than 10 words.

```
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 [30]:

```
# 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()
```



In [31]:

```
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()
```



# Observation:

1. Both for Approved and rejected projects the mean is almost same with reference to Project title.
2. Approval rate is slightly higher for the projects with more words in title. However this feature can't separate Approval and rejected projects well enough.

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

In [32]:

```
# 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 [33]:

```
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 [34]:

```
# 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()
```



In [35]:

```
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()
```



## Observation:

1. Both for Approved and rejected projects the mean is almost same.
2. Approval rate is slightly higher for the projects with more words in essays.

### 1.2.8 Univariate Analysis: Cost per project

In [36]:

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

Out[36]:

|   | 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 [37]:

```python
# 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)
```

Out[37]:

|   | id | price | quantity |
|---|----|-----|-----|
| 0 | p000001 | 459.56 | 7 |
| 1 | p000002 | 515.89 | 21 |

In [38]:

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

In [39]:

```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 [40]:

```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()
```

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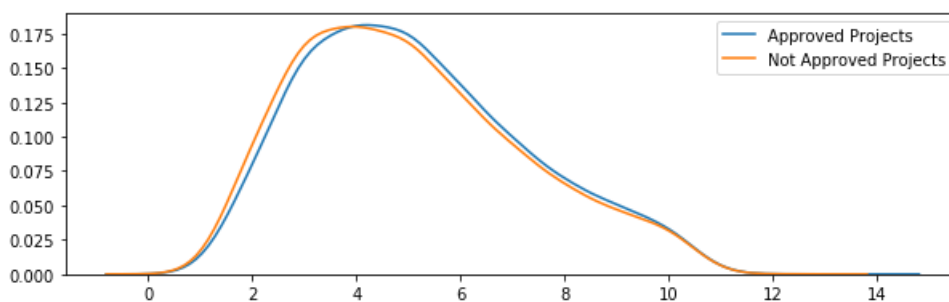
```
# http://zetcode.com/python/prettytable/
from prettytable import PrettyTable

#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_pric
e,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        |
+------------+-------------------+-----------------------+
```

# Observation:

1. Price of approved projects is cheaper compare to rejected projects

## 1.2.9 Univariate Analysis: teacher_number_of_previously_posted_projects

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

In [43]:

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



Number of projects aproved vs rejected

|   | 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 |
|---|---|---|---|
| 242 | 242 | 1 | 1 |
| 268 | 270 | 1 | 1 |
| 234 | 234 | 1 | 1 |
| 335 | 347 | 1 | 1 |
| 373 | 451 | 1 | 1 |

|   | Avg |
|---|---|
| 242 | 1.0 |
| 268 | 1.0 |
| 234 | 1.0 |
| 335 | 1.0 |
| 373 | 1.0 |

# Observation :

1. Project Approval rate ~ 100% for the Teachers who has greater number of previously posted projects.
2. However large number of teachers haven't projects previously

# Conclusion for EDA:

1. Very good data points are collected for each feature.
2. However following analysis can't help in efficiently separating Project Approval and Rejection.
3. Teachers previously posted projects, Project title, Project category and sub category are some useful categorial features for classification.
4. Text data analysis might help in further classification

## 1.2.10 Univariate Analysis: project_resource_summary

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 acceptance of the project or not. If you observe that `presence of the numerical digits` is helpful in the classification, please include it for further process or you can ignore it.

# 1.3 Text preprocessing

## 1.3.1 Essay Text

In [44]:

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

Out[44]:

| | 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 |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 2016-10-25 09:22:10 | Grade |

In [45]:

```
# 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 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 pictures 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 to make our new school year a very successful one. Thank you!nannan

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

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch.  Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore.Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. The want to be able to move as they learn or so they say.Wobble chairs are the answer and I love then because they develop their core, which enhances gross motor and in Turn fine motor skills. \r\nThey also want to learn through games, my kids don't want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves.nannan

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

The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The great teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is makeup is 97.6% African-American, making up the largest segment of the student body. A typical school in Dallas is made up of 23.2% African-American students. Most of the students are on free or reduced lunch. We aren't receiving doctors, lawyers, or engineers children from rich backgrounds or neighborhoods. As an educator I am inspiring minds of young children and we focus not only on academics but one smart, effective, efficient, and disciplined students with good character.In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time.\r\nThe cart will allow me to have more room for storage of things that are needed for the day and has an extra part to it I can use.  The table top chart has all of the letter, words and pictures for students to learn about different letters and it is more accessible.nannan

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

In [46]:

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

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

```python
        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 [47]:

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

My kindergarten students have varied disabilities ranging from speech and language delays, cogniti
ve delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work th
eir hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out
for my students. I teach in a Title I school where most of the students receive free or reduced pr
ice lunch.  Despite their disabilities and limitations, my students love coming to school and come
eager to learn and explore.Have you ever felt like you had ants in your pants and you needed to gr
oove and move as you were in a meeting? This is how my kids feel all the time. The want to be able
to move as they learn or so they say.Wobble chairs are the answer and I love then because they dev
elop their core, which enhances gross motor and in Turn fine motor skills. \r\nThey also want to l
earn through games, my kids do not want to sit and do worksheets. They want to learn to count by j
umping and playing. Physical engagement is the key to our success. The number toss and color and s
hape mats can make that happen. My students will forget they are doing work and just have the fun
a 6 year old deserves.nannan
==================================================

In [48]:

```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)
```

My kindergarten students have varied disabilities ranging from speech and language delays, cogniti
ve delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work th
eir hardest working past their limitations.     The materials we have are the ones I seek out for
my students. I teach in a Title I school where most of the students receive free or reduced price
lunch.  Despite their disabilities and limitations, my students love coming to school and come eag
er to learn and explore.Have you ever felt like you had ants in your pants and you needed to groov
e and move as you were in a meeting? This is how my kids feel all the time. The want to be able to
move as they learn or so they say.Wobble chairs are the answer and I love then because they develo
p their core, which enhances gross motor and in Turn fine motor skills.   They also want to learn t
hrough games, my kids do not want to sit and do worksheets. They want to learn to count by jumping
and playing. Physical engagement is the key to our success. The number toss and color and shape ma
ts can make that happen. My students will forget they are doing work and just have the fun a 6 yea
r old deserves.nannan

In [49]:

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

My kindergarten students have varied disabilities ranging from speech and language delays cognitiv
e delays gross fine motor delays to autism They are eager beavers and always strive to work their
hardest working past their limitations The materials we have are the ones I seek out for my studen
ts I teach in a Title I school where most of the students receive free or reduced price lunch
Despite their disabilities and limitations my students love coming to school and come eager to lea
rn and explore Have you ever felt like you had ants in your pants and you needed to groove and mov
e as you were in a meeting This is how my kids feel all the time The want to be able to move as th
ey learn or so they say Wobble chairs are the answer and I love then because they develop their co
re which enhances gross motor and in Turn fine motor skills They also want to learn through games
my kids do not want to sit and do worksheets They want to learn to count by jumping and playing Ph

ysical engagement is the key to our success The number toss and color and shape mats can make that happen My students will forget they are doing work and just have the fun a 6 year old deserves nan nan

In [50]:

```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', "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [51]:

```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
[01:04<00:00, 1699.54it/s]
```

In [113]:

```python
# after preprocesing
preprocessed_essays[20000]
```

Out[113]:

'my kindergarten students varied disabilities ranging speech language delays cognitive delays gross fine motor delays autism they eager beavers always strive work hardest working past limitations the materials ones i seek students i teach title i school students receive free reduced price lunch despite disabilities limitations students love coming school come eager learn explore have ever felt like ants pants needed groove move meeting this kids feel time the want able move learn say wobble chairs answer i love develop core enhances gross motor turn fine motor skills they also want learn games kids not want sit worksheets they want learn count jumping playing physical engagement key success the number toss color shape mats make happen my students forget work fun 6 year old deserves nannan'

```python
from tqdm import tqdm
Sample_preprocessed_essay = []

for sentance in tqdm(final_approvals['essay'].values):
    samplesentEssay = decontracted(sentance)
    samplesentEssay = samplesentEssay.replace('\\r', ' ')
    samplesentEssay = samplesentEssay.replace('\\"', ' ')
    samplesentEssay = samplesentEssay.replace('\\n', ' ')
    samplesentEssay = re.sub('[^A-Za-z0-9]+', ' ', samplesentEssay)
    # https://gist.github.com/sebleier/554280
    samplesentEssay = ' '.join(e for e in samplesentEssay.split() if e not in stopwords)
    Sample_preprocessed_essay.append(samplesentEssay.lower().strip())
```

```
100%|████████████████████████████████████████████████████| 8000/8000
[00:09<00:00, 863.77it/s]
```

### 1.3.2 Project title Text

```python
# similarly you can preprocess the titles also
# Using above lines of code for preprocessing Title text

from tqdm import tqdm
preprocessed_title = []

for sentance in tqdm(project_data['project_title'].values):
    sentTitle = decontracted(sentance)
    sentTitle = sentTitle.replace('\\r', ' ')
    sentTitle = sentTitle.replace('\\"', ' ')
    sentTitle = sentTitle.replace('\\n', ' ')
    sentTitle = re.sub('[^A-Za-z0-9]+', ' ', sentTitle)
    # https://gist.github.com/sebleier/554280
    sentTitle = ' '.join(e for e in sentTitle.split() if e not in stopwords)
    preprocessed_title.append(sentTitle.lower().strip())
```

```
100%|████████████████████████████████████████████████████| 109248/109248
[00:05<00:00, 21673.34it/s]
```

## 1. 4 Preparing data for models

```python
project_data.columns
```

```
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'],
      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 [55]:

```python
# 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 [175]:

```python
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True
)
vectorizer.fit(final_approvals['clean_categories'].values)
print(vectorizer.get_feature_names())


sample_categories_one_hot = vectorizer.transform(final_approvals['clean_categories'].values)
print("Shape of matrix after one hot encodig ",sample_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  (8000, 9)
```

In [56]:

```python
# 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 [176]:

```python
# 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=
```

```
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=
True)
vectorizer.fit(final_approvals['clean_subcategories'].values)
print(vectorizer.get_feature_names())


sample_sub_categories_one_hot = vectorizer.transform(final_approvals['clean_subcategories'].values
)
print("Shape of matrix after one hot encodig ",sample_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  (8000, 30)
```

In [178]:

```
# Please do the similar feature encoding with state, teacher_prefix and project_grade_category als
o
#Using above lines of code

# project_grade_category

from collections import Counter
my_counter = Counter()
for word in project_data['project_grade_category'].values:
    my_counter.update(word.split())
grade_cat_dict = dict(my_counter)
sorted_grade_cat_dict = dict(sorted(grade_cat_dict.items(), key=lambda kv: kv[1]))
vectorizer_state = CountVectorizer(vocabulary=list(sorted_grade_cat_dict.keys()), lowercase=False,
binary=True)
vectorizer_state.fit(project_data['project_grade_category'].values)
print(vectorizer_state.get_feature_names())
categories_one_hot_grade =
vectorizer_state.transform(project_data['project_grade_category'].values)
print("Shape of matrix after one hot encodig ",categories_one_hot_grade.shape)
```

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

In [179]:

```
from collections import Counter
my_counter = Counter()
for word in final_approvals['project_grade_category'].values:
    my_counter.update(word.split())
grade_cat_dict = dict(my_counter)
sorted_grade_cat_dict = dict(sorted(grade_cat_dict.items(), key=lambda kv: kv[1]))
vectorizer_state = CountVectorizer(vocabulary=list(sorted_grade_cat_dict.keys()), lowercase=False,
binary=True)
vectorizer_state.fit(final_approvals['project_grade_category'].values)
print(vectorizer_state.get_feature_names())
sample_categories_one_hot_grade =
vectorizer_state.transform(final_approvals['project_grade_category'].values)
print("Shape of matrix after one hot encodig ",sample_categories_one_hot_grade.shape)
```

```
['9-12', '6-8', '3-5', 'PreK-2', 'Grades']
Shape of matrix after one hot encodig  (8000, 5)
```

In [55]:

```
# Subject State

from collections import Counter
my_counter = Counter()
for word in project_data['school_state'].values:
    my_counter.update(word.split())
state_cat_dict = dict(my_counter)
sorted_state_cat_dict = dict(sorted(state_cat_dict.items(), key=lambda kv: kv[1]))
vectorizer state = CountVectorizer(vocabulary=list(sorted state cat dict keys()), lowercase=False
```

```
vectorizer_state = CountVectorizer(vocabulary=list(sorted_state_cat_dict.keys()), lowercase=False,
binary=True)
vectorizer_state.fit(project_data['school_state'].values)
print(vectorizer_state.get_feature_names())
categories_one_hot_state = vectorizer_state.transform(project_data['school_state'].values)
print("Shape of matrix after one hot encodig ",categories_one_hot_state.shape)
```

```
['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 [180]:

```python
from collections import Counter
my_counter = Counter()
for word in final_approvals['school_state'].values:
    my_counter.update(word.split())
state_cat_dict = dict(my_counter)
sorted_state_cat_dict = dict(sorted(state_cat_dict.items(), key=lambda kv: kv[1]))
vectorizer_state = CountVectorizer(vocabulary=list(sorted_state_cat_dict.keys()), lowercase=False,
binary=True)
vectorizer_state.fit(final_approvals['school_state'].values)
print(vectorizer_state.get_feature_names())
sample_categories_one_hot_state =
vectorizer_state.transform(final_approvals['school_state'].values)
print("Shape of matrix after one hot encodig ",sample_categories_one_hot_state.shape)
```

```
['ND', 'VT', 'WY', 'RI', 'MT', 'AK', 'NE', 'NH', 'HI', 'DE', 'SD', 'WV', 'ME', 'NM', 'KS', 'DC', 'I
D', 'IA', 'AR', 'MN', 'MS', 'CO', 'KY', 'OR', 'CT', 'NV', 'MD', 'TN', 'UT', 'WI', 'MA', 'AL', 'NJ',
'VA', 'AZ', 'LA', 'WA', 'OH', 'OK', 'IN', 'MO', 'PA', 'MI', 'SC', 'GA', 'IL', 'NC', 'FL', 'NY', 'TX
', 'CA']
Shape of matrix after one hot encodig  (8000, 51)
◄ ▶
```

## 1.4.2 Vectorizing Text data

### 1.4.2.1 Bag of words

In [156]:

```python
# 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 [116]:

```python
# you can vectorize the title also
# before you vectorize the title make sure you preprocess it
# After preprocessing the Project Title

preprocessed_title[20000]
```

Out[116]:

```
'we need to move it while we input it'
```

In [88]:

```python
# Similarly you can vectorize for title also
# Using above lines of code

vectorizerTitle = CountVectorizer(min_df=10)
```

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

```
Shape of matrix after one hot encodig  (109248, 132)
```

In [89]:

```
data_approved = project_data[project_data['project_is_approved'] == 1].sample(n = 4000)
print('Shape of projects approved', data_approved.shape)

data_rejected = project_data[project_data['project_is_approved'] == 0].sample(n = 4000)
print('Shape of projects approved', data_rejected.shape)

final_approvals = pd.concat([data_approved, data_rejected])
print('Shape of final approvals', final_approvals.shape)

sampleTitle_8000 = final_approvals['project_title']
```

```
Shape of projects approved (4000, 20)
Shape of projects approved (4000, 20)
Shape of final approvals (8000, 20)
```

In [123]:

```
# similarly you can preprocess the titles also
# Using above lines of code for preprocessing Title text

from tqdm import tqdm
Sample_preprocessed_title = []

for sentance in tqdm(final_approvals['project_title'].values):
    samplesentTitle = decontracted(sentance)
    samplesentTitle = samplesentTitle.replace('\\r', ' ')
    samplesentTitle = samplesentTitle.replace('\\"', ' ')
    samplesentTitle = samplesentTitle.replace('\\n', ' ')
    samplesentTitle = re.sub('[^A-Za-z0-9]+', ' ', samplesentTitle)
    # https://gist.github.com/sebleier/554280
    samplesentTitle = ' '.join(e for e in samplesentTitle.split() if e not in stopwords)
    Sample_preprocessed_title.append(samplesentTitle.lower().strip())
```

```
100%|████████████████████████████████████████████████████████████████████| 8000/8000
[00:00<00:00, 32071.20it/s]
```

In [125]:

```
Sample_preprocessed_title[100]
```

Out[125]:

```
'stem fairy tales'
```

In [126]:

```
count_vect = CountVectorizer(min_df=10)
from sklearn.preprocessing import StandardScaler

std_scaler = StandardScaler(with_mean=False)

sampleTitle_8000 = count_vect.fit_transform(Sample_preprocessed_title)
sampleTitle_8000 = std_scaler.fit_transform(sampleTitle_8000)
sampleTitle_8000 = sampleTitle_8000.todense()
print(sampleTitle_8000.shape)
```

```
C:\Users\psudheer\AppData\Local\Continuum\anaconda3\lib\site-
packages\sklearn\utils\validation.py:595: DataConversionWarning:

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

C:\Users\psudheer\AppData\Local\Continuum\anaconda3\lib\site-
packages\sklearn\utils\validation.py:595: DataConversionWarning:
```

```
(8000, 558)
```

**1.4.2.3 TFIDF vectorizer**

In [143]:

```python
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 [92]:

```python
# Similarly you can vectorize for title also
# Using above lines of code

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

```
Shape of matrix after one hot encodig  (109248, 132)
```

In [127]:

```python
# Similarly you can vectorize for title also
# Using above lines of code

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.preprocessing import StandardScaler
std_scaler = StandardScaler(with_mean=False)
samplevectorizerTitle = TfidfVectorizer(min_df=10)
sample_tfidf_title = samplevectorizerTitle.fit_transform(Sample_preprocessed_title)
sample_tfidf_title = std_scaler.fit_transform(sample_tfidf_title)
sample_tfidf_title = sample_tfidf_title.todense()
print("Shape of matrix after one hot encodig ",sample_tfidf_title.shape)
```

```
Shape of matrix after one hot encodig  (8000, 558)
```

**1.4.2.5 Using Pretrained Models: Avg W2V**

In [203]:

```python
'''
# 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))


# 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[203]:

'\n\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\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 word
s in the coupus", len(words))\n\ninter_words = set(model.keys()).intersection(words)\nprint("The n
umber of words that 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'

In [141]:

```
# 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 [142]:

```
# average Word2Vec
# compute average word2vec for each review.
```

```python
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
[01:02<00:00, 1741.29it/s]

109248
300

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

In [206]:

```python
# Similarly you can vectorize for title also
# Using above lines of code
# average Word2Vec for Project Title

avg_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
    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_title.append(vector)

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

100%|████████████████████████████████████████████████████████████████| 109248/109248
[00:30<00:00, 3563.00it/s]

109248
300

In [103]:

```python
sample_avg_w2v_vectors_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(Sample_preprocessed_title): # 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
    sample_avg_w2v_vectors_title.append(vector)

print(len(sample_avg_w2v_vectors_title))
print(len(sample_avg_w2v_vectors_title[0]))
```

100%|████████████████████████████████████████████████████████████████| 8000/8000
[00:03<00:00, 2218.78it/s]

```
8000
300
```

## 1.4.2.7 Using Pretrained Models: TFIDF weighted W2V

In [144]:

```
# 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 [145]:

```
# 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
[09:05<00:00, 200.23it/s]
```

```
109248
300
```

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

In [146]:

```
# Similarly you can vectorize for title also

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

In [147]:

```
# Using above lines of code
# 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_title):
            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_title[word]*(sentence.count(word)/len(sentence.split())) # getting
the tfidf value for each word
        vector += (vec * tf_idf) # calculating tfidf weighted w2v
        tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_title.append(vector)

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

```
109248
300
```

In [148]:

```
sample_tfidf_model_title = TfidfVectorizer()
sample_tfidf_model_title.fit(Sample_preprocessed_title)
# we are converting a dictionary with word as a key, and the idf as a value
sample_dictionary_title = dict(zip(sample_tfidf_model_title.get_feature_names(),
list(sample_tfidf_model_title.idf_)))
sample_tfidf_words_title = set(sample_tfidf_model_title.get_feature_names())
```

In [149]:

```
sample_tfidf_w2v_vectors_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(Sample_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 sample_tfidf_words_title):
            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 = sample_dictionary_title[word]*(sentence.count(word)/len(sentence.split())) # g
etting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    sample_tfidf_w2v_vectors_title.append(vector)

print(len(sample_tfidf_w2v_vectors_title))
print(len(sample_tfidf_w2v_vectors_title[0]))
```

```
8000
300
```

### 1.4.3 Vectorizing Numerical features

In [158]:

```
# 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 maen and variance.
price_standardized = price_scalar.transform(project_data['price'].values.reshape(-1, 1))
```

Mean : 298.1193425966608, Standard deviation : 367.49634838483496

```
price_standardized
```

```
array([[-0.3905327 ],
       [ 0.00239637],
       [ 0.59519138],
       ...,
       [-0.15825829],
       [-0.61243967],
       [-0.51216657]])
```

```
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(final_approvals['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 maen and variance.
sample_price_standardized = price_scalar.transform(final_approvals['price'].values.reshape(-1, 1))
```

Mean : 321.14584625, Standard deviation : 340.44757528543784

```
sample_price_standardized
```

```
array([[-0.71290226],
       [ 0.93011135],
       [-0.62883645],
       ...,
       [ 0.34958144],
       [-0.92832456],
       [-0.62428362]])
```

### 1.4.4 Merging all the above features

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

```
print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(text_bow.shape)
print(price_standardized.shape)
```

```
(109248, 9)
(109248, 30)
(109248, 16623)
(109248, 1)
```

In [167]:

```python
# 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[167]:

```
(109248, 16663)
```

In [186]:

```python
print(sample_categories_one_hot.shape)
print(sample_sub_categories_one_hot.shape)
print(sample_categories_one_hot_grade.shape)
print(sample_categories_one_hot_state.shape)
print(sampleTitle_8000.shape)
print(sample_price_standardized.shape)
```

```
(8000, 9)
(8000, 30)
(8000, 5)
(8000, 51)
(8000, 558)
(8000, 1)
```

In [187]:

```python
# 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_all = hstack((sample_categories_one_hot, sample_sub_categories_one_hot,
sample_categories_one_hot_grade,
sample_categories_one_hot_state,sampleTitle_8000,sample_price_standardized))
X_all.shape
```

Out[187]:

```
(8000, 654)
```

# 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_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)

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 [1]:

```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()
```
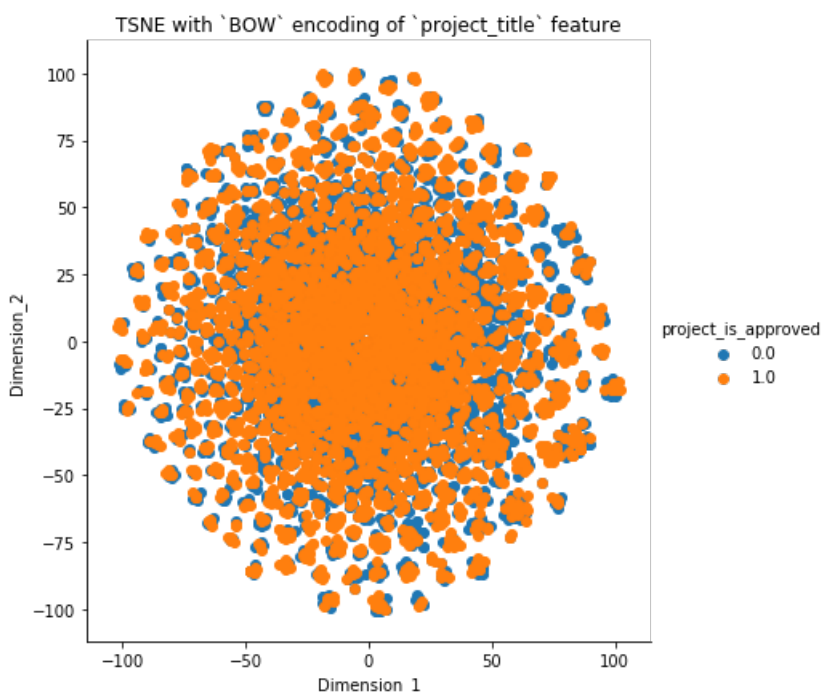


## 2.1 TSNE with `BOW` encoding of `project_title` feature

In [128]:

```python
# please write all of the code with proper documentation and proper titles for each subsection
# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label
# Have referred some of the Kernels from https://www.kaggle.com/snap/amazon-fine-food-reviews

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

```
# storing 8k data points in approve lists
approve = final_approvals['project_is_approved']

# 8k Sample Datapoints are taken into consideration
# Shape of Sample BOW project title is (8000, 558)
# Iteration for 8000 Datapoints with a perplexity of 50
model = TSNE(n_components=2, perplexity=150, n_iter = 2000, random_state=0)

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

for_tsne = model.fit_transform(sampleTitle_8000)
for_tsne = np.vstack((for_tsne.T, approve)).T

# X- axis is Dimension_ 1 and Y-axis is Dimension_2

for_tsne_df = pd.DataFrame(data=for_tsne,
columns=['Dimension_1','Dimension_2','project_is_approved'])
sns.FacetGrid(for_tsne_df, hue="project_is_approved", size=6).map(plt.scatter, 'Dimension_1', 'Dime
nsion_2').add_legend()
plt.title("TSNE with `BOW` encoding of `project_title` feature")
plt.show()
```
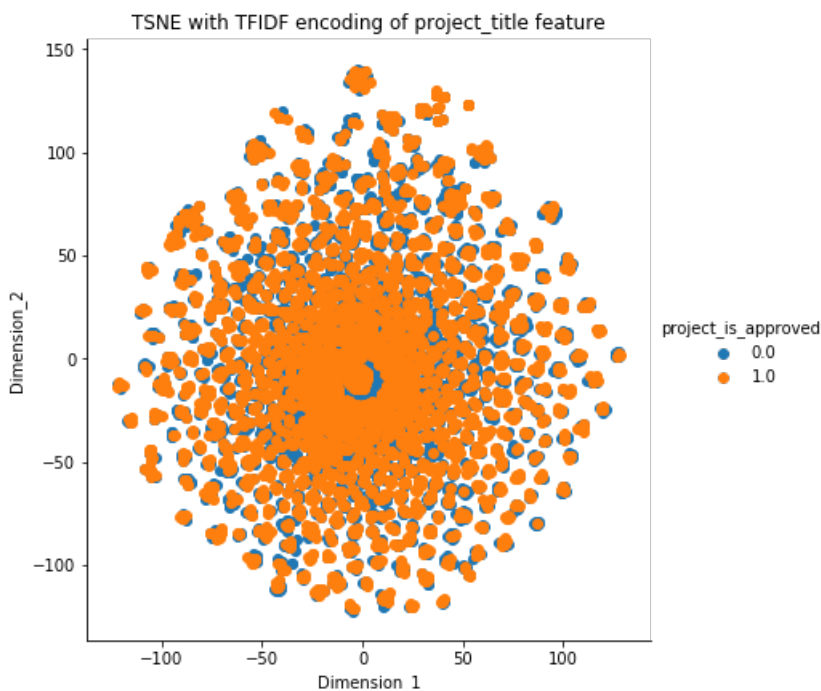
```
C:\Users\psudheer\AppData\Local\Continuum\anaconda3\lib\site-packages\seaborn\axisgrid.py:230: Use
rWarning:

The `size` paramter has been renamed to `height`; please update your code.
```



## Observation:

1. Points are not very well separated for Project approved or rejected.
2. Lot of overlapping is seen for projects approved vs rejected

### 2.2 TSNE with `TFIDF` encoding of `project_title` feature

In [131]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label
```

```python
# Have referred some of the Kernels from https://www.kaggle.com/snap/amazon-fine-food-reviews

from sklearn.manifold import TSNE
from sklearn import datasets
import pandas as pd
import matplotlib.pyplot as plt

approve = final_approvals['project_is_approved']


# 8k Sample Datapoints are taken into consideration
# Shape of Sample BOW project title is (8000, 558)
# Iteration for 8000 Datapoints with a perplexity of 50
model = TSNE(n_components=2, perplexity=50, n_iter = 2000, random_state=0)

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

for_tsne_tfidf = model.fit_transform(sample_tfidf_title)
for_tsne_tfidf = np.vstack((for_tsne_tfidf.T, approve)).T

for_tsne_tfidf_df = pd.DataFrame(data=for_tsne_tfidf,
columns=['Dimension_1','Dimension_2','project_is_approved'])
sns.FacetGrid(for_tsne_tfidf_df, hue="project_is_approved", size=6).map(plt.scatter, 'Dimension_1',
'Dimension_2').add_legend()
plt.title("TSNE with TFIDF encoding of project_title feature")
plt.show()
```

```
C:\Users\psudheer\AppData\Local\Continuum\anaconda3\lib\site-packages\seaborn\axisgrid.py:230: Use
rWarning:

The `size` paramter has been renamed to `height`; please update your code.
```



TSNE with TFIDF encoding of project_title feature

## Observation:

1. Points are not very well separated for Project approved or rejected.
2. Lot of overlapping is seen for projects approved vs rejected


## 2.3 TSNE with `AVG W2V` encoding of `project_title` feature

In [132]:

```python
# please write all the code with proper documentation, and proper titles for each subsection
# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label

# Have referred some of the Kernels from https://www.kaggle.com/snap/amazon-fine-food-reviews

from sklearn.manifold import TSNE
from sklearn import datasets
import pandas as pd
import matplotlib.pyplot as plt


approve = final_approvals['project_is_approved']

# 8k Sample Datapoints are taken into consideration
# Shape of Sample BOW project title is (8000, 300)
# Iteration for 8000 Datapoints with a perplexity of 50
model = TSNE(n_components=2, perplexity=50, n_iter = 2000, random_state=0)

for_tsne_avgW2V = model.fit_transform(sample_avg_w2v_vectors_title)
for_tsne_avgW2V = np.vstack((for_tsne_avgW2V.T, approve)).T

for_tsne_avgW2V_df = pd.DataFrame(data=for_tsne_avgW2V,
columns=['Dimension_1','Dimension_2','project_is_approved'])
sns.FacetGrid(for_tsne_avgW2V_df, hue="project_is_approved", size=6).map(plt.scatter, 'Dimension_1'
, 'Dimension_2').add_legend()
plt.title("TSNE with AvgW2V encoding of project_title feature")
plt.show()
```
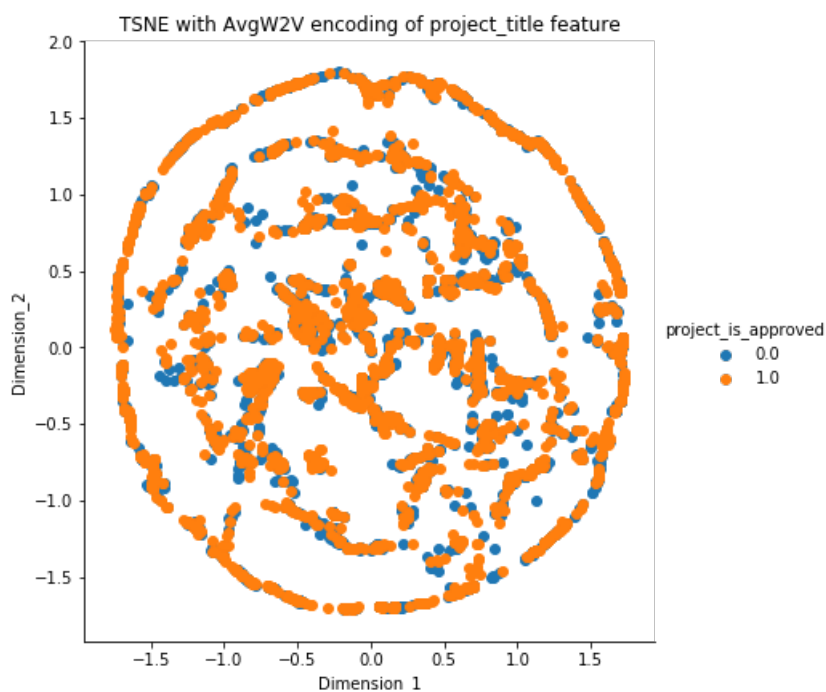
```
C:\Users\psudheer\AppData\Local\Continuum\anaconda3\lib\site-packages\seaborn\axisgrid.py:230: Use
rWarning:

The `size` paramter has been renamed to `height`; please update your code.
```



## Observation:

1. Points are not very well separated for Project approved or rejected.
2. Lot of overlapping is seen for projects approved vs rejected

## 2.4 TSNE with `TFIDF Weighted W2V` encoding of `project_title` feature

```python
# please write all the code with proper documentation, and proper titles for each subsection
# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label
# Have referred some of the Kernels from https://www.kaggle.com/snap/amazon-fine-food-reviews

from sklearn.manifold import TSNE
from sklearn import datasets
import pandas as pd
import matplotlib.pyplot as plt


approve = final_approvals['project_is_approved']

# 8k Sample Datapoints are taken into consideration
# Shape of Sample BOW project title is (8000, 300)
# Iteration for 8000 Datapoints with a perplexity of 50
model = TSNE(n_components=2, perplexity=50, n_iter = 2000, random_state=0)

for_tsne_tfidfW2V = model.fit_transform(sample_tfidf_w2v_vectors_title)
for_tsne_tfidfW2V = np.vstack((for_tsne_tfidfW2V.T, approve)).T

for_tsne_tfidfW2V_df = pd.DataFrame(data=for_tsne_tfidfW2V, columns=['Dimension_1','Dimension_2','
project_is_approved'])
sns.FacetGrid(for_tsne_tfidfW2V_df, hue="project_is_approved", size=6).map(plt.scatter, 'Dimension_
1', 'Dimension_2').add_legend()
plt.title("TSNE with Weighted W2V encoding of project_title feature")
plt.show()
```
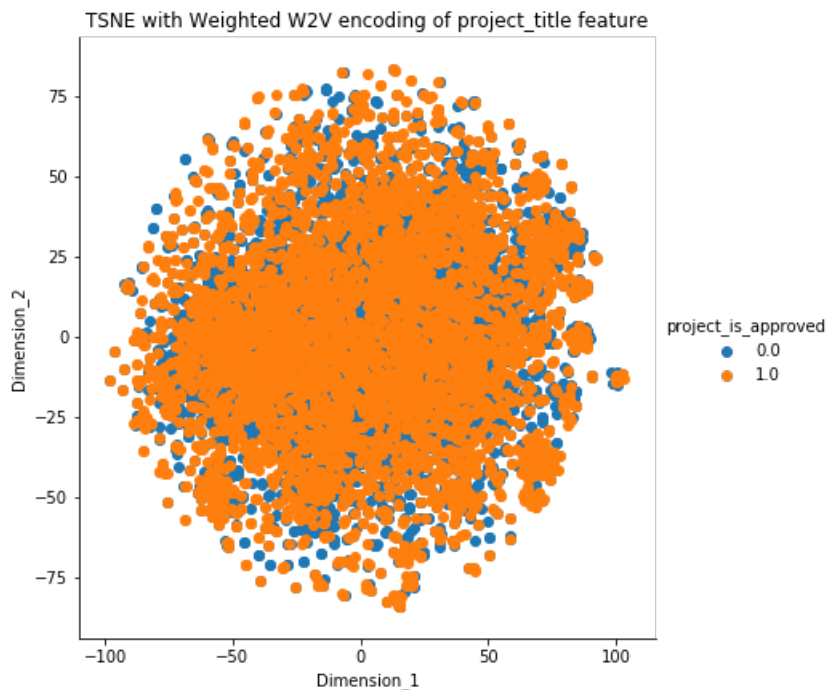
```
C:\Users\psudheer\AppData\Local\Continuum\anaconda3\lib\site-packages\seaborn\axisgrid.py:230: Use
rWarning:

The `size` paramter has been renamed to `height`; please update your code.
```



TSNE with Weighted W2V encoding of project_title feature

## Observation:

1. Points are not very well separated for Project approved or rejected.
2. Lot of overlapping is seen for projects approved vs rejected

```python
# Have referred some of the Kernels from https://www.kaggle.com/snap/amazon-fine-food-reviews

from sklearn.manifold import TSNE
from sklearn import datasets
import pandas as pd
import matplotlib.pyplot as plt

approve = final_approvals['project_is_approved']


# 8k Sample Datapoints are taken into consideration
# Shape of Sample BOW project title is (8000, 654)
# Iteration for 8000 Datapoints with a perplexity of 50
model = TSNE(n_components=2, perplexity=50, n_iter = 2000, random_state=0)

# X_all is a horizontial vertical stacking of all the features
for_all = model.fit_transform(X_all.toarray())
for_all = np.vstack((for_all.T, approve)).T

for_all_df = pd.DataFrame(data=for_all, columns=['Dimension_1','Dimension_2','project_is_approved'
])
sns.FacetGrid(for_all_df, hue="project_is_approved", size=6).map(plt.scatter, 'Dimension_1', 'Dimen
sion_2').add_legend()
plt.title("TSNE with all features")
plt.show()
```
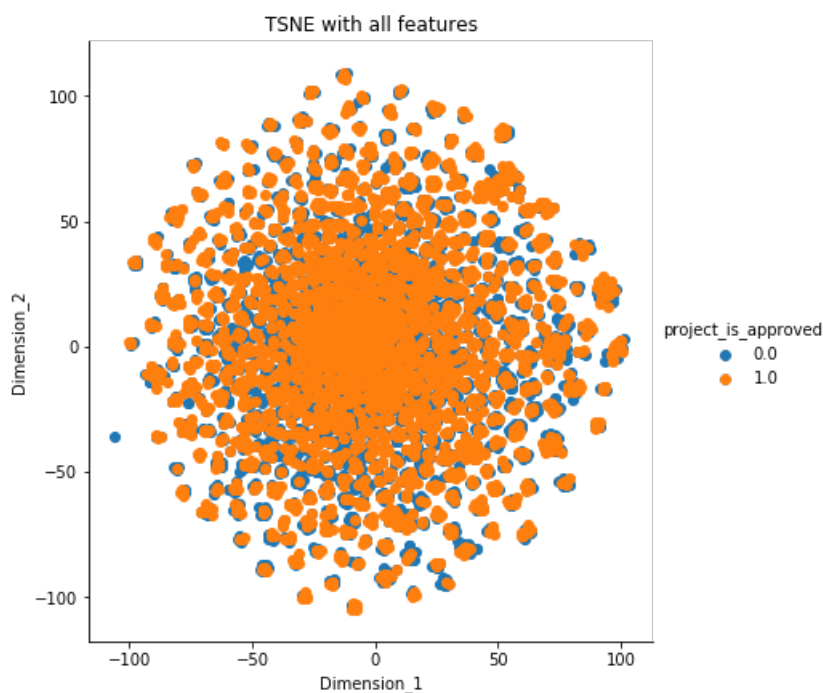
```
C:\Users\psudheer\AppData\Local\Continuum\anaconda3\lib\site-packages\seaborn\axisgrid.py:230: Use
rWarning:

The `size` paramter has been renamed to `height`; please update your code.
```


TSNE with all features

## Observation:

1. Even on merging all categorial, Text and Numerical data the points are still not very well separated for Project approved or rejected.
2. Lot of overlapping is seen for projects approved vs rejected

## 2.5 Summary

## Conclusion

# Conclusion

1. None of the TSNE representation gave us a well separation between projects Approved vs Rejected.
2. We cannot draw a plane or write a condition to separate projects Approved vs Rejected using this analysis.
3. We have to follow some other alternative method for a well separation between projects Approved vs Rejected. .