# DonorsChoose

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

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

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

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

## About the DonorsChoose Data Set

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

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

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

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

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

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

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

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

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

## Notes on the Essay Data

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

- __project_essay_1:__ "Introduce us to your classroom"
- __project_essay_2:__ "Tell us more about your students"
- __project_essay_3:__ "Describe how your students will use the materials you're requesting"
- __project_essay_4:__ "Close by sharing why your project will make a difference"

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

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

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

In [1]:

```python
%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 [2]:

```python
project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')
project_data.head(2)
```

Out[2]:

| | 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 [3]:

```python
print("Number of data points in train data", project_data.shape)
print('-'*70)
print("The attributes of data :", project_data.columns.values)
```

```
('Number of data points in train data', (109248, 17))
----------------------------------------------------------------------
('The attributes of data :', array(['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'], dtype=object))
```

In [4]:

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

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

Out[4]:

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

## 1.2 Data Analysis

In [5]:

```
# 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()
```
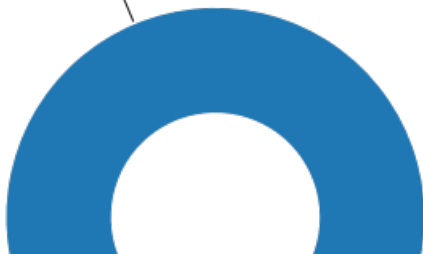
```
('Number of projects thar are approved for funding ', 92706, ', (', 0, '%)')
('Number of projects thar are not approved for funding ', 16542, ', (', 0, '%)')
```

Not Accepted

Out of 109248 applied projects, 84.85% are approved for funding.

### 1.2.1 Univariate Analysis: School State

In [6]:

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

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')
```

In [7]:

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

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

In [8]:

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

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

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

In [9]:

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

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

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

    if top:
        temp = temp[0:top]

    stack_plot(temp, xtick=col1, col2=col2, col3='total')
    print(temp.head(5))
```

```
    print("="*50)
    print(temp.tail(5))
```
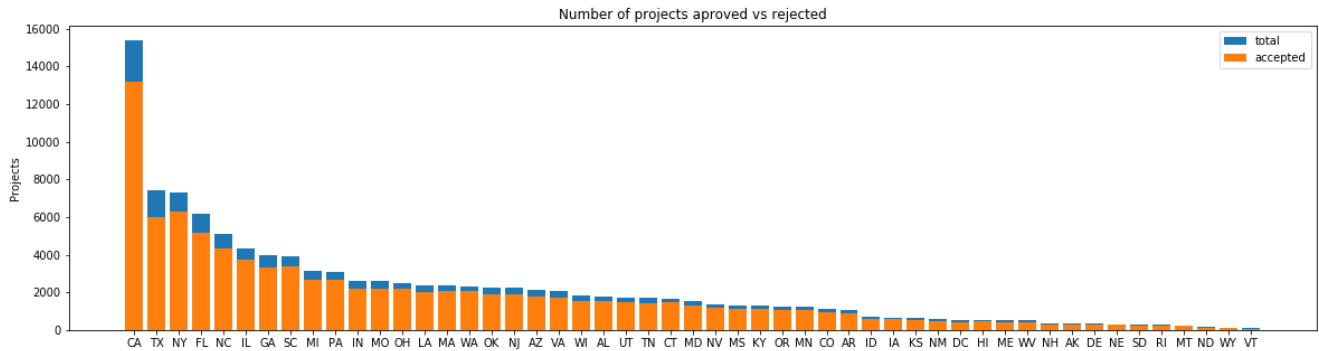
In [10]:

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


Number of projects aproved vs rejected

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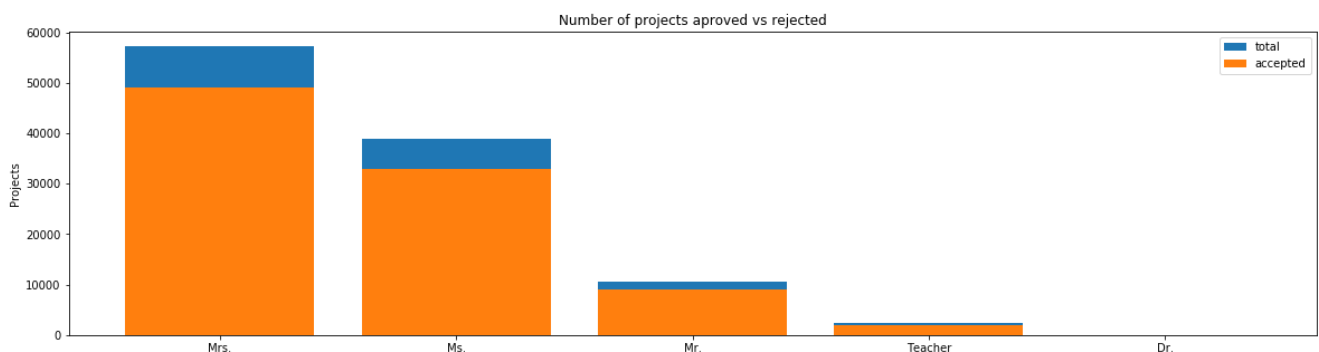
**SUMMARY: Every state has greater than 80% success rate in approval**

### 1.2.2 Univariate Analysis: teacher_prefix

In [11]:

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


Number of projects aproved vs rejected

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

**SUMMARY: The teachers with prefix 'Mrs.' have highest approval percentage.**

### 1.2.3 Univariate Analysis: project_grade_category

In [12]:

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



Number of projects aproved vs rejected

```
  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
```

**SUMMARY: The Project Grade Category 3-5 has highest project approval percentage. The Project Grade Category PreK-2 has highest number of applications**

### 1.2.4 Univariate Analysis: project_subject_categories

In [13]:

```
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 & Hunger"]
        if 'The' in j.split(): # this will split each of the catogory based on space "Math & Science"=> "Math","&", "Science"
            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
        j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
        temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_') # we are replacing the & value into
    cat_list.append(temp.strip())
```

In [14]:

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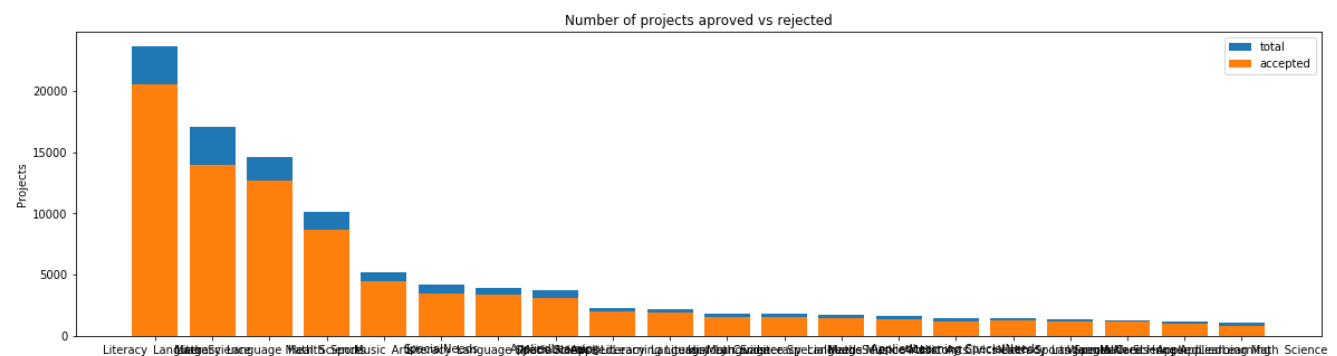
```
project_data.drop(['project_subject_categories'], axis=1, inplace=True)
project_data.head(2)
```

Out[14]:

| | 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 [15]:

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



```
              clean_categories  project_is_approved  total      Avg
24           Literacy_Language                20520  23655  0.867470
32                Math_Science                13991  17072  0.819529
28  Literacy_Language Math_Science            12725  14636  0.869432
8                Health_Sports                 8640  10177  0.848973
40                   Music_Arts                4429   5180  0.855019
==================================================
              clean_categories  project_is_approved  total      Avg
19  History_Civics Literacy_Language             1271   1421  0.894441
14        Health_Sports SpecialNeeds             1215   1391  0.873472
50            Warmth Care_Hunger                 1212   1309  0.925898
33        Math_Science AppliedLearning           1019   1220  0.835246
4         AppliedLearning Math_Science            855   1052  0.812738
```
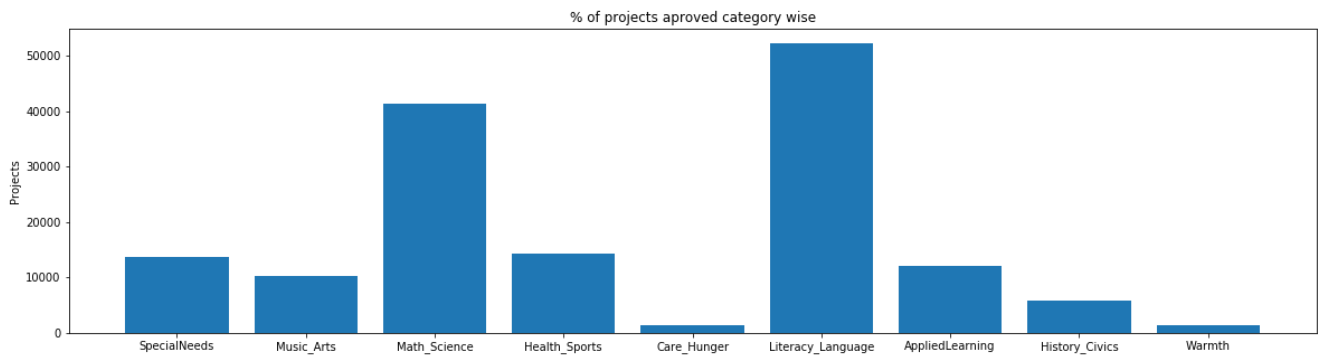
In [16]:

```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 [17]:

```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()
```

% of projects aproved category wise

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

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

**SUMMARY:The Project Grade Category Literacy_Language has highest project approvals**

### 1.2.5 Univariate Analysis: project_subject_subcategories

In [19]:

```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 [20]:

```python
project_data['clean_subcategories'] = sub_cat_list
project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
project_data.head(2)
```
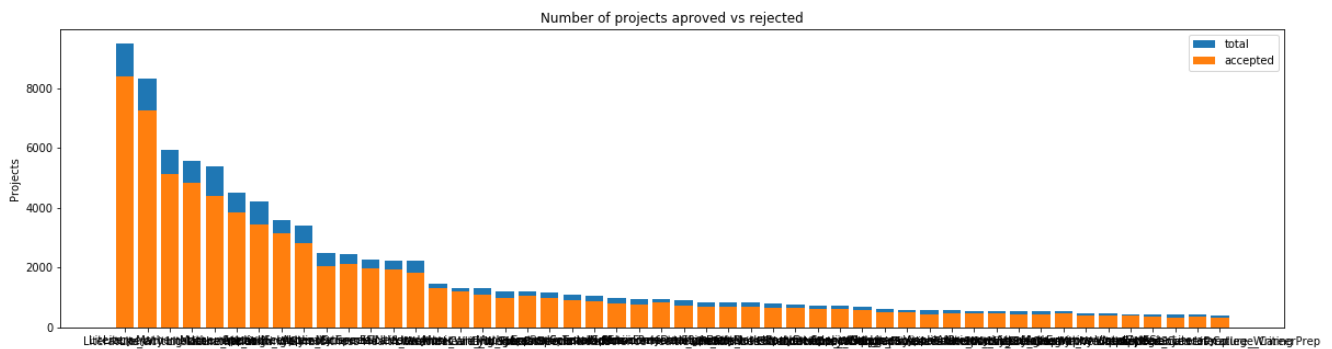
Out[20]:

| Unnamed:<br>0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_cate |

| | 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 [21]:

```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
```
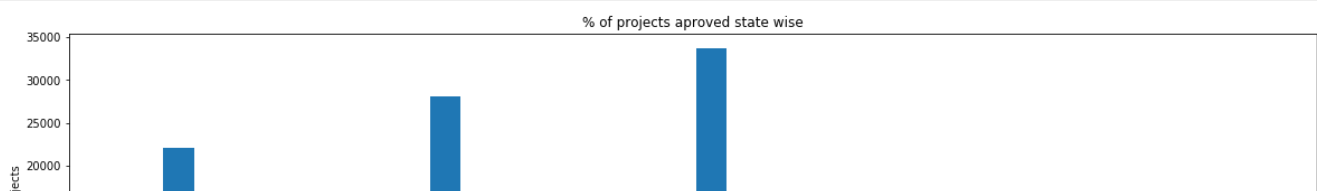
In [22]:

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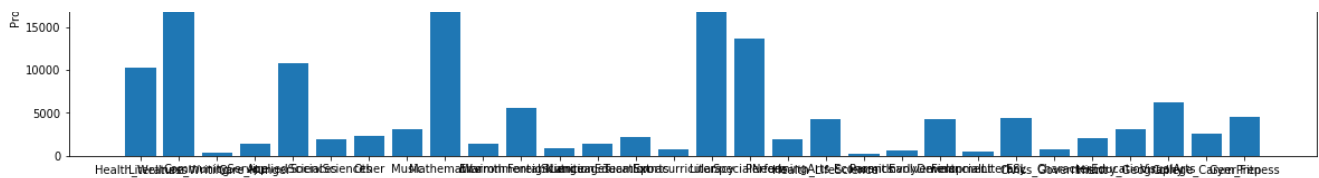
In [23]:

```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()
```



% of projects aproved state wise

In [24]:

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

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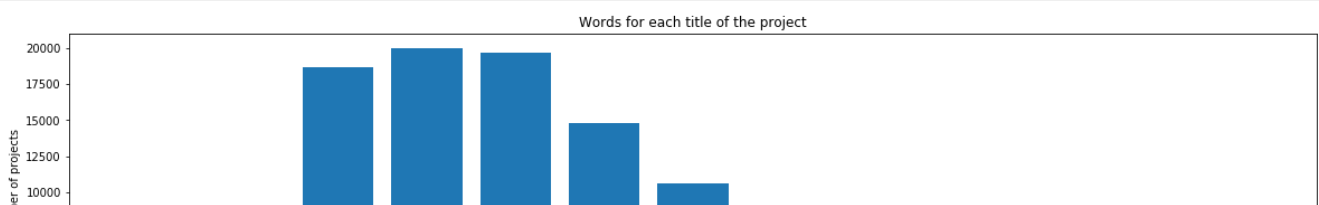
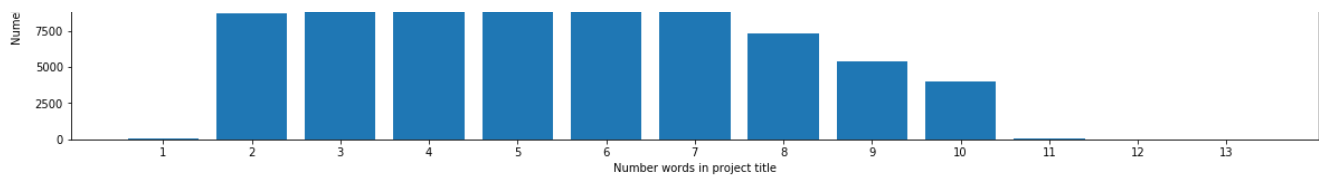**SUMMARY:The Project Grade Category Literacy has highest project approvals**

## 1.2.6 Univariate Analysis: Text features (Title)

In [25]:

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

**SUMMARY:Most of the projects contain Title with 4 words following title with 5 words.**
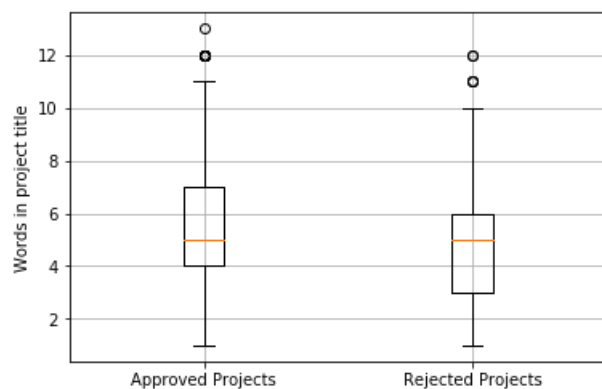
In [26]:

```
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
```
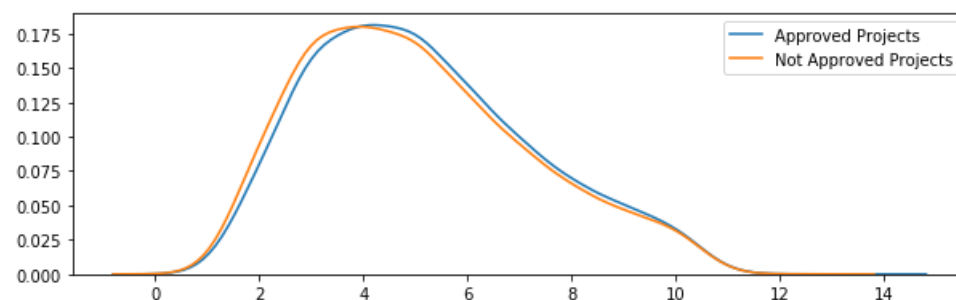
In [27]:

```
# 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 [28]:

```
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()
```



**SUMMARY:Project Titles with 6 or more words have higher chances of approval**

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

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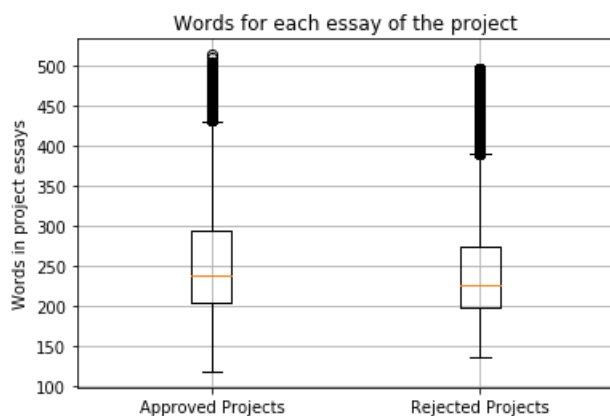
```
approved_word_count = project_data[project_data['project_is_approved']==1]['essay'].str.split().app
ly(len)
approved_word_count = approved_word_count.values

rejected_word_count = project_data[project_data['project_is_approved']==0]['essay'].str.split().app
ly(len)
rejected_word_count = rejected_word_count.values
```
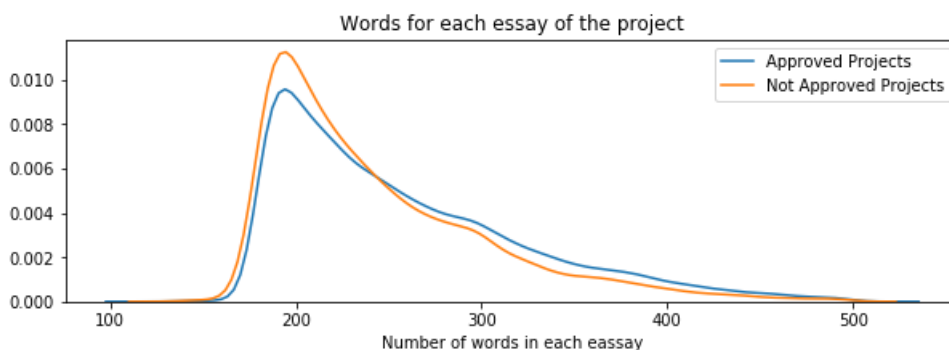
```
# https://glowingpython.blogspot.com/2012/09/boxplot-with-matplotlib.html
plt.boxplot([approved_word_count, rejected_word_count])
plt.title('Words for each essay of the project')
plt.xticks([1,2],('Approved Projects','Rejected Projects'))
plt.ylabel('Words in project essays')
plt.grid()
plt.show()
```

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



**SUMMARY:The Project Essays with more than 250 words have higher chances of approval.**

### 1.2.8 Univariate Analysis: Cost per project

In [33]:

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

Out[33]:

|   | 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 [34]:

```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[34]:

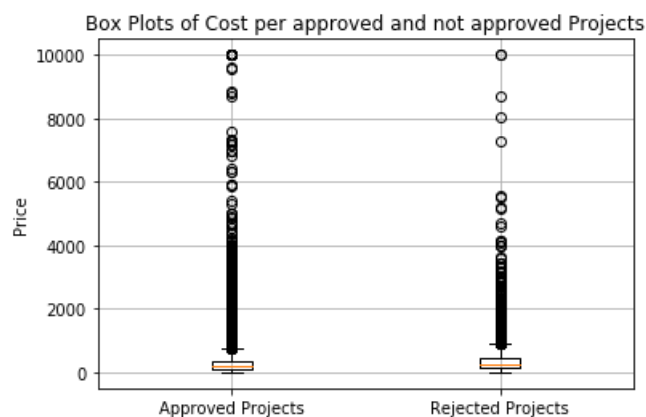|   | id | price | quantity |
|---|-----|-------|----------|
| **0** | p000001 | 459.56 | 7 |
| **1** | p000002 | 515.89 | 21 |

In [35]:

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

In [36]:

```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
```
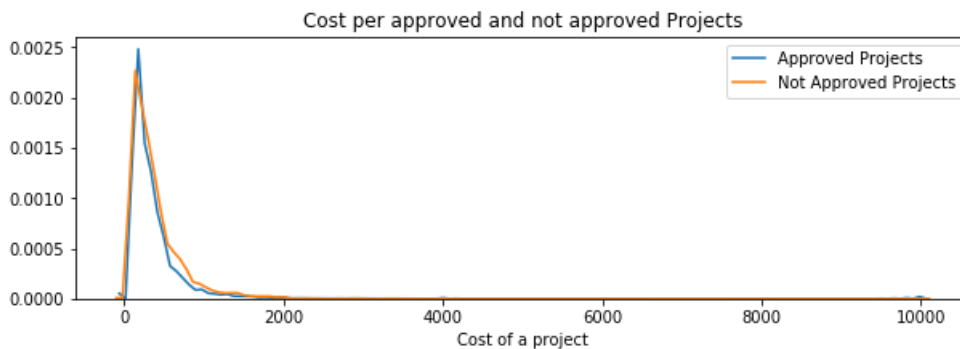
In [37]:

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



In [38]:

```
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        |
+------------+-------------------+-----------------------+
```

**SUMMARY: If the Cost per Project is lesser then chances of approval are greater.**

### 1.2.9 Univariate Analysis: teacher_number_of_previously_posted_projects
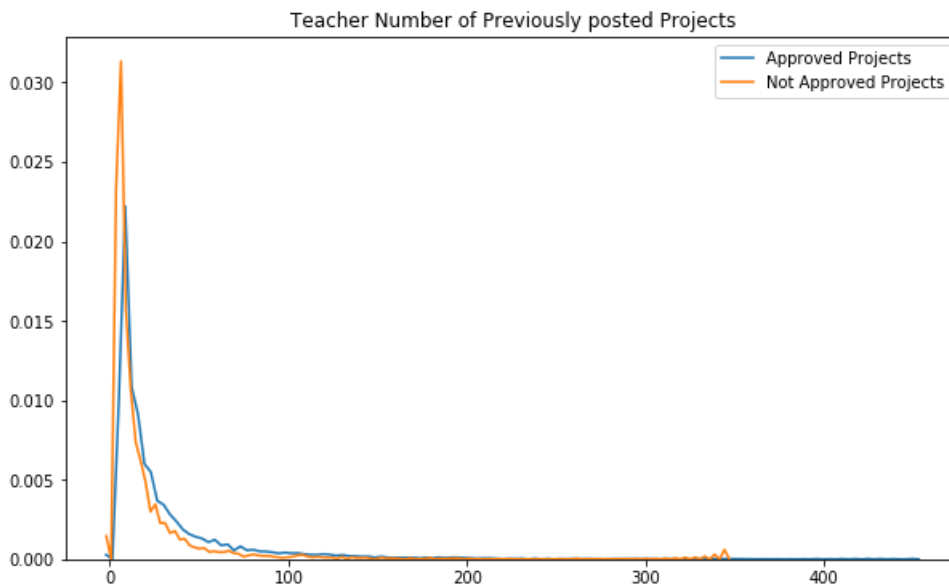
```
tn=pd.DataFrame(project_data,columns=['teacher_number_of_previously_posted_projects','project_is_ap
proved'])
```

```
approved = tn[tn['project_is_approved']==1]['teacher_number_of_previously_posted_projects'].values
not_approved=tn[tn['project_is_approved']==0]
['teacher_number_of_previously_posted_projects'].values
```

```
plt.figure(figsize=(10,6))
sns.distplot(approved, hist=False, label="Approved Projects")
sns.distplot(not_approved, hist=False, label="Not Approved Projects")
plt.title('Teacher Number of Previously posted Projects')
plt.legend()
plt.show()
```

```
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,i), 3), np.round(np.percentile(not_approved,i), 3)
])
print(x)
```

```
+------------+-------------------+-----------------------+
| Percentile | Approved Projects | Not Approved Projects |
+------------+-------------------+-----------------------+
|     0      |        0.0        |          0.0          |
|     5      |        0.0        |          0.0          |
|     10     |        0.0        |          0.0          |
|     15     |        0.0        |          0.0          |
|     20     |        0.0        |          0.0          |
|     25     |        0.0        |          0.0          |
|     30     |        1.0        |          0.0          |
|     35     |        1.0        |          1.0          |
|     40     |        1.0        |          1.0          |
|     45     |        2.0        |          1.0          |
|     50     |        2.0        |          2.0          |
|     55     |        3.0        |          2.0          |
|     60     |        4.0        |          3.0          |
|     65     |        5.0        |          3.0          |
|     70     |        7.0        |          4.0          |
|     75     |        9.0        |          6.0          |
|     80     |       13.0        |          8.0          |
|     85     |       19.0        |         11.0          |
|     90     |       30.0        |         17.0          |
|     95     |       57.0        |         31.0          |
|    100     |       451.0       |         345.0         |
```

```
+------------+-------------------+----------------------+
```

**Approved Projects is slightly ahead of non approved projects. So this means that teachers with more previous submissions have better chances of projects getting approved.**

```
t=tn.groupby(tn['teacher_number_of_previously_posted_projects'])
p=t.count()
state = tn.groupby(tn['teacher_number_of_previously_posted_projects']).sum()
k=state['project_is_approved']/p['project_is_approved']*100
print(k.head(10))
```

```
teacher_number_of_previously_posted_projects
0    82.135004
1    83.005356
2    84.106280
3    84.345992
4    84.542347
5    84.775833
6    85.517039
7    85.395764
8    86.218927
9    86.778969
Name: project_is_approved, dtype: float64
```

**So, the above percentages tell us that as the number of previously applied projects increases, the approval rate also increases.**

## 1.2.10 Univariate Analysis: project_resource_summary

```
prs=pd.DataFrame(project_data,columns=['project_resource_summary','project_is_approved'])
```

```
prs['project_resource_summary'] = prs['project_resource_summary'].astype(str)
# Extract digits from the project resource summary
prs['project_resource'] = prs['project_resource_summary'].str.extract('(\d+ )', expand=False).str.s
trip()
```

```
#Fill all NaN with zeros
prs['project_resource']=prs['project_resource'].fillna(0)
```

```
prt=prs.loc[prs['project_is_approved']==1]#positve
prt=prt.loc[prt['project_resource'].astype(int)>0]#non-zero value
size_of_non_zero_positives=prt.shape[0]
size_of_non_zero_positives=float(size_of_non_zero_positives)
print('Total number of non-zero project resource valued positives =
'+str(size_of_non_zero_positives))
```

```
Total number of non-zero project resource valued positives = 10802.0
```

```
non_zeros=prs.loc[prs['project_resource'].astype(int)>0]
size_of_non_zeros=non_zeros.shape[0]
size_of_non_zeros=float(size_of_non_zeros)
print('Total number of non-zero project resource valued positives =
'+str(size_of_non_zero_positives/size_of_non_zeros*100))
```

```
Total number of non-zero project resource valued positives = 90.1744719927
```

So, it can be observed here that **90 percent of rows which have a number in it are accepted. So if the project resource contains a number which can be describing a the quantity of a thing, then the chances of acceptance are high.**

## 1.3 Text preprocessing

### 1.3.1 Essay Text

In [50]:

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

Out[50]:

| | 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 [51]:

```
# printing some random essays.
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
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
==================================================
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

be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pic tures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade.  This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs lost of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project t o make our new school year a very successful one. Thank you!nannan
==================================================
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 don't want to sit and do worksheets. They want to learn to count by ju mping and playing. Physical engagement is the key to our success. The number toss and color and sh ape 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 [52]:

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

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

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

In [53]:

```python
sent = decontracted(project_data['essay'].values[20000])
print(sent)
print("="*100)
```

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 [54]:

```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 [55]:
```
#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 [56]:
```
# 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', 'e
ach', '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', "do
esn'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 [57]:
```
# Combining all the above statemennts
```

```
# combining all the above statements
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 [00:58<00:00, 1876.52it/s]
```

In [58]:

```
# after preprocesing
print(preprocessed_essays[2507])
print("="*100)
print(len(preprocessed_essays))
```

```
my students come diverse backgrounds most live poverty level our school considered smart school du
e changing neighborhood many students leave middle school year also get many new kids middle year
they no consistency life transient lives make difficult put focus school they deserve teacher give
s 100 i not not need i new teacher school i little it high poverty school parents not afford suppl
ies kids need school year i already spent much money supplies still need much i simply asking basi
c supplies items keep classroom organized i want students opportunities students schools if
students i basic needs covered keep focus learning nannan
====================================================================================================

109248
```

### 1.3.2 Project title Text

In [59]:

```
# similarly you can preprocess the titles also
```

In [60]:

```
project_data['project_title'].size
```

Out[60]:

```
109248
```

In [61]:

```
# printing some random titles.
print(project_data['project_title'].values[0])
print("="*50)
print(project_data['project_title'].values[150])
print("="*50)
print(project_data['project_title'].values[1000])
print("="*50)
print(project_data['project_title'].values[20000])
print("="*50)
print(project_data['project_title'].values[99999])
print("="*50)
```

```
Educational Support for English Learners at Home
==================================================
More Movement with Hokki Stools
==================================================
Sailing Into a Super 4th Grade Year
==================================================
We Need To Move It While We Input It!
==================================================
```

```
Inspiring Minds by Enhancing the Educational Experience
=======================================================
```

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

```
We Need To Move It While We Input It!
=======================================================
```

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

```
We Need To Move It While We Input It!
```

```python
title = re.sub('[^A-Za-z0-9]+', ' ', title)
print(title)
```

```
We Need To Move It While We Input It
```

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

```
100%|██████████| 109248/109248 [00:03<00:00, 32536.11it/s]
```

```python
print(preprocessed_titles[2867])
print("="*100)
print(len(preprocessed_titles))
```

```
a fresh start
====================================================================================================

109248
```

## 1. 4 Preparing data for models

```python
project_data.columns
```

```
Index([u'Unnamed: 0', u'id', u'teacher_id', u'teacher_prefix', u'school_state',
       u'project_submitted_datetime', u'project_grade_category',
```

```
      u'project_title', u'project_essay_1', u'project_essay_2',
      u'project_essay_3', u'project_essay_4', u'project_resource_summary',
      u'teacher_number_of_previously_posted_projects', u'project_is_approved',
      u'clean_categories', u'clean_subcategories', u'essay', u'price',
      u'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 [68]:

```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)
print(type(categories_one_hot))
```

```
['SpecialNeeds', 'Music_Arts', 'Math_Science', 'Health_Sports', 'Care_Hunger',
'Literacy_Language', 'AppliedLearning', 'History_Civics', 'Warmth']
('Shape of matrix after one hot encodig ', (109248, 9))
<class 'scipy.sparse.csr.csr_matrix'>
```

In [69]:

```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)
print(type(sub_categories_one_hot))
```

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

In [70]:

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

In [71]:

```
# we use count vectorizer to convert the values into one hot encoded features
my_counter = Counter()
for word in project_data['school_state'].values:
    my_counter.update(word.split())

state_dict = dict(my_counter)
sorted_state_dict = dict(sorted(state_dict.items(), key=lambda kv: kv[1]))

vectorizers = CountVectorizer(vocabulary=list(sorted_state_dict.keys()), lowercase=False, binary=T
rue)
vectorizers.fit(project_data['school_state'].values)
print(vectorizers.get_feature_names())

state_one_hot = vectorizers.transform(project_data['school_state'].values)
print("Shape of matrix after one hot encodig ",state_one_hot.shape)
print(type(state_one_hot))
```

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

◀                      ▶

In [72]:

```
project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('')
my_counter = Counter()
for word in project_data['teacher_prefix'].values:
    my_counter.update(word.split())

teacher_dict = dict(my_counter)
sorted_teacher_dict = dict(sorted(teacher_dict.items(), key=lambda kv: kv[1]))

vectorizers = CountVectorizer(vocabulary=list(sorted_teacher_dict.keys()), lowercase=False, binary
=True)
vectorizers.fit(project_data['teacher_prefix'].values)
print(vectorizers.get_feature_names())

teacher_one_hot = vectorizers.transform(project_data['teacher_prefix'].values)
print("Shape of matrix after one hot encodig ",teacher_one_hot.shape)
print(type(teacher_one_hot))
```

```
['Ms.', 'Mr.', 'Teacher', 'Mrs.', 'Dr.']
('Shape of matrix after one hot encodig ', (109248, 5))
<class 'scipy.sparse.csr.csr_matrix'>
```

In [73]:

```
#Combining all the above statements
preproc = []
# tqdm is for printing the status bar
for sent in project_data['project_grade_category']:
    sent = sent.replace('Grades ', 'Grade_')
    preproc.append(sent)
project_data['project_grade_category']=preproc
```

In [74]:

```
my_counter = Counter()
for word in project_data['project_grade_category'].values:
    my_counter.update(word.split())

grade_dict = dict(my_counter)
sorted_grade_dict = dict(sorted(grade_dict.items(), key=lambda kv: kv[1]))

vectorizer = CountVectorizer(vocabulary=list(sorted_grade_dict.keys()), lowercase=False, binary=Tr
```

```
ue)
vectorizer.fit(project_data['project_grade_category'].values)
print(vectorizer.get_feature_names())

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

```
['Grade_3-5', 'Grade_9-12', 'Grade_6-8', 'Grade_PreK-2']
('Shape of matrix after one hot encodig ', (109248, 4))
<class 'scipy.sparse.csr.csr_matrix'>
```

## 1.4.2 Vectorizing Text data

### 1.4.2.1 Bag of words

In [75]:

```
# 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)
print(type(text_bow))
```

```
('Shape of matrix after one hot encodig ', (109248, 16623))
<class 'scipy.sparse.csr.csr_matrix'>
```

### 1.4.2.2 Bag of Words on `project_title`

In [76]:

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

```
('Shape of matrix after one hot encodig ', (109248, 3329))
<class 'scipy.sparse.csr.csr_matrix'>
```

### 1.4.2.3 TFIDF vectorizer

In [77]:

```
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)
print(type(text_tfidf))
```

```
('Shape of matrix after one hot encodig ', (109248, 16623))
<class 'scipy.sparse.csr.csr_matrix'>
```

### 1.4.2.4 TFIDF Vectorizer on `project_title`

In [78]:

```
vectorizer = TfidfVectorizer(min_df=10)
title_tfidf = vectorizer.fit_transform(preprocessed_titles)
print("Shape of matrix after one hot encodig ",title_tfidf.shape)
print(type(title_tfidf))
```

```
('Shape of matrix after one hot encodig ', (109248, 3329))
<class 'scipy.sparse.csr.csr_matrix'>
```

**1.4.2.5 Using Pretrained Models: Avg W2V**

In [79]:

```python
# # Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
# def loadGloveModel(gloveFile):
#     print ("Loading Glove Model")
#     f = open(gloveFile,'r')
#     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')


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

# for i in preprocessed_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-
# save-and-load-variables-in-python/

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

In [80]:

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

In [81]:

```python
# average Word2Vec
# compute average word2vec for each essay.
avg_w2v_vectors_essay = []; # 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_essay.append(vector)

print(len(avg_w2v_vectors_essay))
print(len(avg_w2v_vectors_essay[0]))
print(type(avg_w2v_vectors_essay))
```

```
100%|████████| 109248/109248 [00:28<00:00, 3845.03it/s]
```

```
109248
300
<type 'list'>
```

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

In [82]:

```python
# average Word2Vec
# compute average word2vec for each title.
avg_w2v_vectors_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_titles): # 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]))
print(type(avg_w2v_vectors_title))
```

```
100%|████████| 109248/109248 [00:01<00:00, 63603.32it/s]
```

```
109248
300
<type 'list'>
```

### 1.4.2.7 Using Pretrained Models: TFIDF weighted W2V

In [83]:

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

In [84]:

```python
# average Word2Vec
# compute average word2vec for each essay.
tfidf_w2v_vectors_essay = []; # 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_essay.append(vector)

print(len(tfidf_w2v_vectors_essay))
print(len(tfidf_w2v_vectors_essay[0]))
print(type(tfidf_w2v_vectors_essay))
```

```
100%|██████████| 109248/109248 [03:40<00:00, 494.83it/s]
```

```
109248
300
<type 'list'>
```

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

In [85]:

```
tfidf_model = TfidfVectorizer()
tfidf_model.fit(preprocessed_titles)
# 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 [86]:

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

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

```
100%|██████████| 109248/109248 [00:03<00:00, 31285.81it/s]
```

```
109248
300
<type 'list'>
```

## 1.4.3 Vectorizing Numerical features

In [87]:

```
# 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
```

```python
# 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))
print(type(price_standardized))
```

```
<type 'numpy.ndarray'>
```

In [88]:

```python
price_standardized
```

Out[88]:

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

In [89]:

```python
# 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

# 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)

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

# Now standardize the data with above maen and variance.
project_standardized =
project_scalar.transform(project_data['teacher_number_of_previously_posted_projects'].values.reshap
e(-1, 1))
print(type(project_standardized))
```

```
<type 'numpy.ndarray'>
```

In [90]:

```python
project_standardized
```

Out[90]:

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

In [91]:

```
# 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

# 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)

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

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

```
<type 'numpy.ndarray'>
```

In [92]:

```
quantity_standardized
```

Out[92]:

```
array([[-0.24576291],
       [ 0.01842593],
       [-0.24576291],
       ...,
       [ 0.41470919],
       [-0.11366849],
       [-0.11366849]])
```

### 1.4.4 Merging all the above features

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

In [93]:

```
tfidf_w2v_vectors_title=np.asarray(tfidf_w2v_vectors_title)
avg_w2v_vectors_title=np.asarray(avg_w2v_vectors_title)
print(categories_one_hot.shape,type(categories_one_hot))
print(sub_categories_one_hot.shape,type(sub_categories_one_hot))
print(state_one_hot.shape,type(state_one_hot))
print(teacher_one_hot.shape,type(teacher_one_hot))
print(grade_one_hot.shape,type(grade_one_hot))

print(title_bow.shape,type(title_bow))
print(title_tfidf.shape,type(title_tfidf))
print(tfidf_w2v_vectors_title.shape,type(tfidf_w2v_vectors_title))
print(avg_w2v_vectors_title.shape,type(avg_w2v_vectors_title))
print(price_standardized.shape,type(price_standardized))
print(project_standardized.shape,type(project_standardized))
```

```
((109248, 9), <class 'scipy.sparse.csr.csr_matrix'>)
((109248, 30), <class 'scipy.sparse.csr.csr_matrix'>)
((109248, 51), <class 'scipy.sparse.csr.csr_matrix'>)
((109248, 5), <class 'scipy.sparse.csr.csr_matrix'>)
((109248, 4), <class 'scipy.sparse.csr.csr_matrix'>)
((109248, 3329), <class 'scipy.sparse.csr.csr_matrix'>)
((109248, 3329), <class 'scipy.sparse.csr.csr_matrix'>)
((109248, 300), <type 'numpy.ndarray'>)
((109248, 300), <type 'numpy.ndarray'>)
```

```
((109248, 1), <type 'numpy.ndarray'>)
((109248, 1), <type 'numpy.ndarray'>)
```

In [ ]:

In [94]:

```python
from scipy.sparse import hstack,isspmatrix
# 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[94]:

```
(109248, 16663)
```

In [95]:

```python
print(type(X))
```

```
<class 'scipy.sparse.coo.coo_matrix'>
```

# 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)
   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

## 2.1 TSNE with `BOW` encoding of `project_title` feature(5K Points)

In [96]:

```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

#Reference:: https://scikit-learn.org/stable/modules/generated/sklearn.manifold.TSNE.html
```

```
#Reference:: https://scikit-learn.org/stable/modules/generated/sklearn.manifold.TSNE.html
from sklearn.manifold import TSNE
X = hstack((categories_one_hot, sub_categories_one_hot, state_one_hot, grade_one_hot, teacher_one_h
ot, price_standardized,project_standardized, title_bow))
print(X.shape)
print(type(X))
print(isspmatrix(X))
X = X.tocsr() #convert sparse matrix in coordinate format to compressed sparse row matrix
print(type(X))
X_new = X[0:3500,:]
X_new = X_new.todense()
print(type(X_new))
print(isspmatrix(X_new))
print(X_new.shape)
X_new = StandardScaler().fit_transform(X_new)
```

```
(109248, 3430)
<class 'scipy.sparse.coo.coo_matrix'>
True
<class 'scipy.sparse.csr.csr_matrix'>
<class 'numpy.matrix'>
False
(3500, 3430)
```
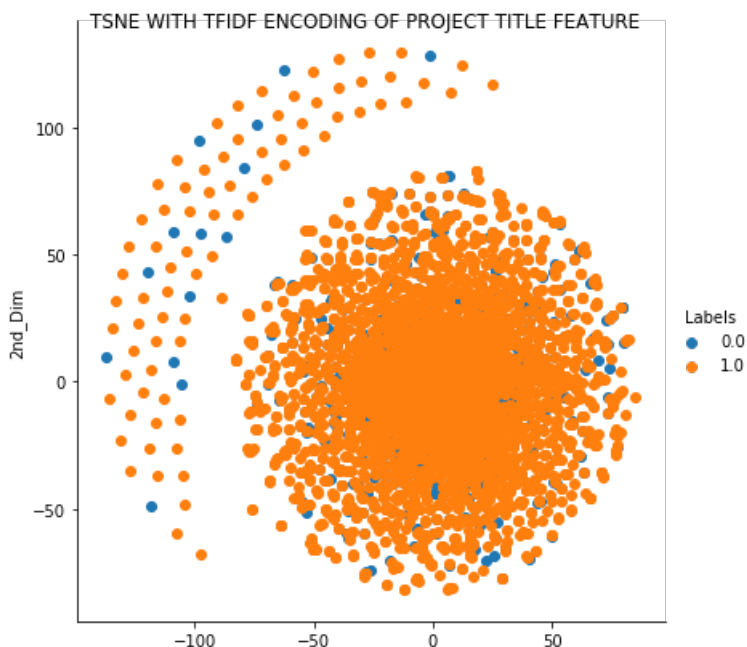
In [97]:

```
model = TSNE(n_components = 2, perplexity = 25.0, random_state = 0)
tsne_data_b = model.fit_transform(X_new)
labels = project_data["project_is_approved"]
labels_new = labels[0: 3500]
tsne_data_b = np.vstack((tsne_data_b.T, labels_new)).T
tsne_df_b = pd.DataFrame(tsne_data_b, columns = ("1st_Dim","2nd_Dim","Labels"))
print(tsne_df_b.shape)
```

```
(3500, 3)
```

In [98]:

```
sns.FacetGrid(tsne_df_b, hue = "Labels", height = 7).map(plt.scatter, "1st_Dim", "2nd_Dim").add_leg
end().fig.suptitle("TSNE WITH BOW ENCODING OF PROJECT TITLE FEATURE ")
plt.show()
```

## 2.2 TSNE with `TFIDF` encoding of `project_title` feature (5K Points)

In [99]:

```
# 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
X = hstack((categories_one_hot, sub_categories_one_hot, state_one_hot, grade_one_hot, teacher_one_h
ot, price_standardized,project_standardized, title_tfidf))
print(X.shape)
print(type(X))
print(isspmatrix(X))
X = X.tocsr()
X_new = X[0:3500,:]
X_new = X_new.todense()
print(type(X_new))
print(isspmatrix(X_new))
print(X_new.shape)
X_new = StandardScaler().fit_transform(X_new)
```

```
(109248, 3430)
<class 'scipy.sparse.coo.coo_matrix'>
True
<class 'numpy.matrix'>
False
(3500, 3430)
```

In [100]:

```
model = TSNE(n_components = 2, perplexity = 25.0, random_state = 0)
tsne_data_b = model.fit_transform(X_new)
labels = project_data["project_is_approved"]
labels_new = labels[0: 3500]
tsne_data_b = np.vstack((tsne_data_b.T, labels_new)).T
tsne_df_b = pd.DataFrame(tsne_data_b, columns = ("1st_Dim","2nd_Dim","Labels"))
print(tsne_df_b.shape)
```

```
(3500, 3)
```

In [101]:

```
sns.FacetGrid(tsne_df_b, hue = "Labels", height = 6).map(plt.scatter, "1st_Dim", "2nd_Dim").add_leg
end().fig.suptitle("TSNE WITH TFIDF ENCODING OF PROJECT TITLE FEATURE ")
plt.show()
```

## 2.3 TSNE with `AVG W2V` encoding of `project_title` feature (5K Points)

In [102]:

```
# 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
X = hstack((categories_one_hot, sub_categories_one_hot, state_one_hot, grade_one_hot, teacher_one_h
ot, price_standardized,project_standardized, avg_w2v_vectors_title))
print(X.shape)
print(type(X))
print(isspmatrix(X))
X = X.tocsr()
X_new = X[0:3500,:]
X_new = X_new.todense()
print(type(X_new))
print(isspmatrix(X_new))
print(X_new.shape)
X_new = StandardScaler().fit_transform(X_new)
```

```
(109248, 401)
<class 'scipy.sparse.coo.coo_matrix'>
True
<class 'numpy.matrix'>
False
(3500, 401)
```
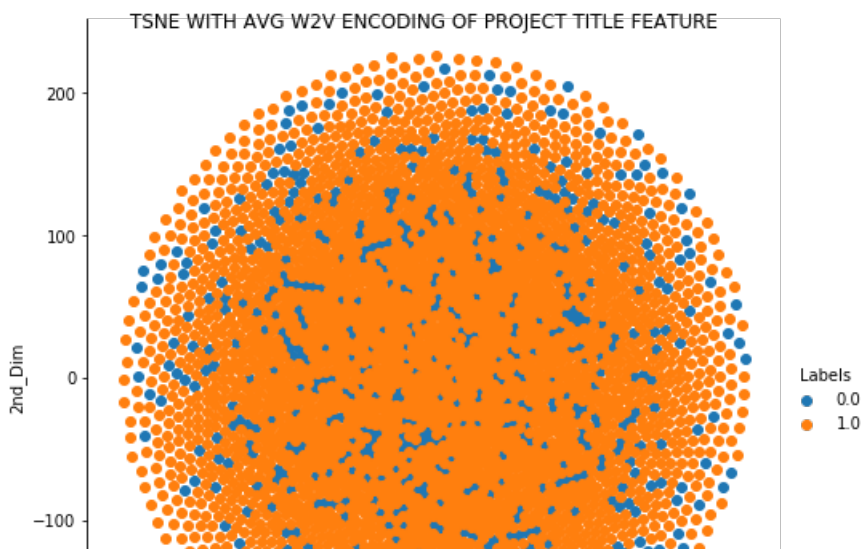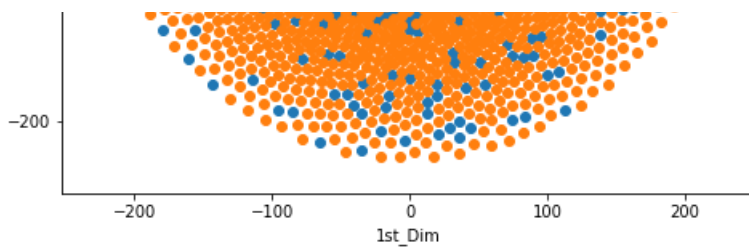
In [103]:

```
model = TSNE(n_components = 2, perplexity = 0.0, random_state = 0)
tsne_data_b = model.fit_transform(X_new)
labels = project_data["project_is_approved"]
labels_new = labels[0: 3500]
tsne_data_b = np.vstack((tsne_data_b.T, labels_new)).T
tsne_df_b = pd.DataFrame(tsne_data_b, columns = ("1st_Dim","2nd_Dim","Labels"))
print(tsne_df_b.shape)
```

```
(3500, 3)
```

In [104]:

```
sns.FacetGrid(tsne_df_b, hue = "Labels", height = 7).map(plt.scatter, "1st_Dim", "2nd_Dim").add_leg
end().fig.suptitle("TSNE WITH AVG W2V ENCODING OF PROJECT TITLE FEATURE ")
plt.show()
```

-200

-200    -100    0    100    200

1st_Dim

## 2.4 TSNE with `TFIDF Weighted W2V` encoding of `project_title` feature (5K Points)

In [105]:

```
# 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

X = hstack((categories_one_hot, sub_categories_one_hot, state_one_hot, grade_one_hot, teacher_one_hot, price_standardized,project_standardized, tfidf_w2v_vectors_title))
print(X.shape)
print(type(X))
print(isspmatrix(X))
X = X.tocsr()
X_new = X[0:3500,:]
X_new = X_new.todense()
print(type(X_new))
print(isspmatrix(X_new))
print(X_new.shape)
X_new = StandardScaler().fit_transform(X_new)
```

```
(109248, 401)
<class 'scipy.sparse.coo.coo_matrix'>
True
<class 'numpy.matrix'>
False
(3500, 401)
```
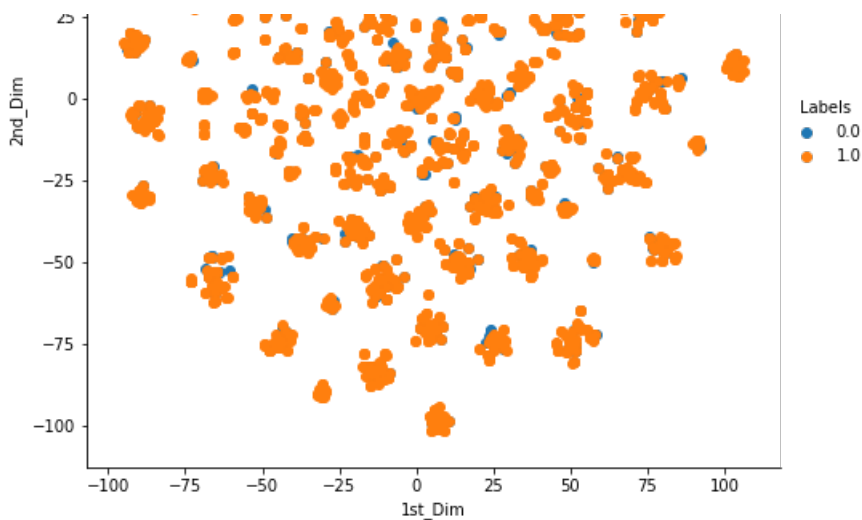
In [106]:

```
model = TSNE(n_components = 2, perplexity = 15.0, random_state = 0)
tsne_data_b = model.fit_transform(X_new)
labels = project_data["project_is_approved"]
labels_new = labels[0: 3500]
tsne_data_b = np.vstack((tsne_data_b.T, labels_new)).T
tsne_df_b = pd.DataFrame(tsne_data_b, columns = ("1st_Dim","2nd_Dim","Labels"))
print(tsne_df_b.shape)
```

```
(3500, 3)
```

In [107]:

```
sns.FacetGrid(tsne_df_b, hue = "Labels", height = 7).map(plt.scatter, "1st_Dim", "2nd_Dim").add_legend().fig.suptitle("TSNE WITH TFIDF WEIGHTED W2V ENCODING OF PROJECT TITLE FEATURE ")
plt.show()
```

## 2.5 TSNE with all features combined (5K Points)

In [108]:

```
# 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
X = hstack((categories_one_hot, sub_categories_one_hot, state_one_hot, grade_one_hot, teacher_one_h
ot, price_standardized,project_standardized,title_bow,title_tfidf,avg_w2v_vectors_title,
tfidf_w2v_vectors_title))
print(X.shape)
print(type(X))
print(isspmatrix(X))
X = X.tocsr()
X_new = X[0:3500,:]
X_new = X_new.todense()
print(type(X_new))
print(isspmatrix(X_new))
print(X_new.shape)
X_new = StandardScaler().fit_transform(X_new)
```

```
(109248, 7359)
<class 'scipy.sparse.coo.coo_matrix'>
True
<class 'numpy.matrix'>
False
(3500, 7359)
```
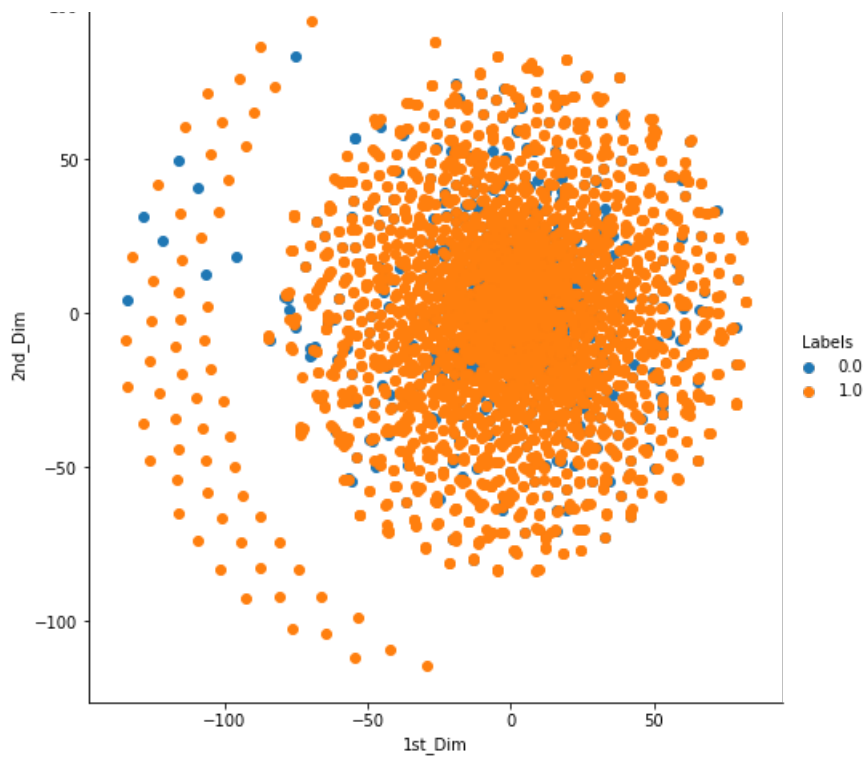
In [109]:

```
model = TSNE(n_components = 2, perplexity = 25.0, random_state = 0)
tsne_data_b = model.fit_transform(X_new)
labels = project_data["project_is_approved"]
labels_new = labels[0: 3500]
tsne_data_b = np.vstack((tsne_data_b.T, labels_new)).T
tsne_df_b = pd.DataFrame(tsne_data_b, columns = ("1st_Dim","2nd_Dim","Labels"))
print(tsne_df_b.shape)
```

```
(3500, 3)
```

In [110]:

```
sns.FacetGrid(tsne_df_b, hue = "Labels", height = 7).map(plt.scatter, "1st_Dim", "2nd_Dim").add_leg
end().fig.suptitle("TSNE WITH ALL FEATURES COMBINED")
plt.show()
```

TSNE WITH ALL FEATURES COMBINED

100

## 2.5 Summary

**The TSNE visualisation with BoW, TF-IDF, Avg Word2Vec, TF-IDF Weighted Word2Vec does not yield the expected result even after trying with various values of perplexity due to almost overlapping of points in all the cases.**

**Note:** Due to shortage of RAM, I have loaded glove.42B.300d.txt into glove_vector by running one kernel and then commenting out that snippet and then running another kernel. And I also have to perform tsne on 3500 points

```
In [ ]:
```