

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. Example: 123456789
<code>project_title</code>		Title of the project. Example: Art Will Make You
<code>project_grade_category</code>		Grade level of students for which the project is targeted. One of the following enumerated list of categories: Example: Grades K-2
<code>project_subject_categories</code>		One or more (comma-separated) subject categories for the project from the following enumerated list of categories: Example: Applied & Technical, Care & Health, History & Social Studies, Literacy & Language, Math & Science, Music & Arts, Special Education
<code>project_subject_subcategories</code>		One or more (comma-separated) subject subcategories for the project from the following enumerated list of categories: Example: Music & Arts, Literacy & Language, Math & Science
<code>school_state</code>		State where school is located (Two-letter U.S. postal abbreviations). Example: CA
<code>project_resource_summary</code>		An explanation of the resources needed for the project. Example: My students need hands on literacy materials to support sensory
<code>project_essay_1</code>		First application essay
<code>project_essay_2</code>		Second application essay
<code>project_essay_3</code>		Third application essay
<code>project_essay_4</code>		Fourth application essay
<code>project_submitted_datetime</code>		Datetime when project application was submitted. Example: 2018-01-12T12:43:00
<code>teacher_id</code>		A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4f1

Feature		Description
		Teacher's title. One of the following enumerated values: <code>teacher_prefix</code>
		• • • • • •
<code>teacher_number_of_previously_posted_projects</code>		Number of project applications previously submitted by the same teacher. Example: 1

* 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
<code>id</code>	A <code>project_id</code> value from the <code>train.csv</code> file. Example: <code>p036502</code>
<code>description</code>	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
<code>quantity</code>	Quantity of the resource required. Example: 3
<code>price</code>	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
<code>project_is_approved</code>	A binary flag indicating whether DonorsChoose approved the project. A value of <code>0</code> indicates the project was not approved.

Notes on the Essay Data

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

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

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

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

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

In [1]:

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
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 chart_studio import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
```

```
C:\Users\harsh\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

1.1 Reading Data

In [2]:

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

In [3]:

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

Number of data points in train data (109248, 17)

The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix'
'school_state'
'project_submitted_datetime' 'project_grade_category'
'project_subject_categories' 'project_subject_subcategories'
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']

In [4]:

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

Number of data points in train data (1541272, 4)

['id' 'description' 'quantity' 'price']

Out[4]:

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

1.2 Data Analysis

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-gl
r-gallery-pie-and-polar-charts-pie-and-donut-labels-py

y_value_counts = project_data['project_is_approved'].value_counts()
print("Number of projects that 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 that 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)

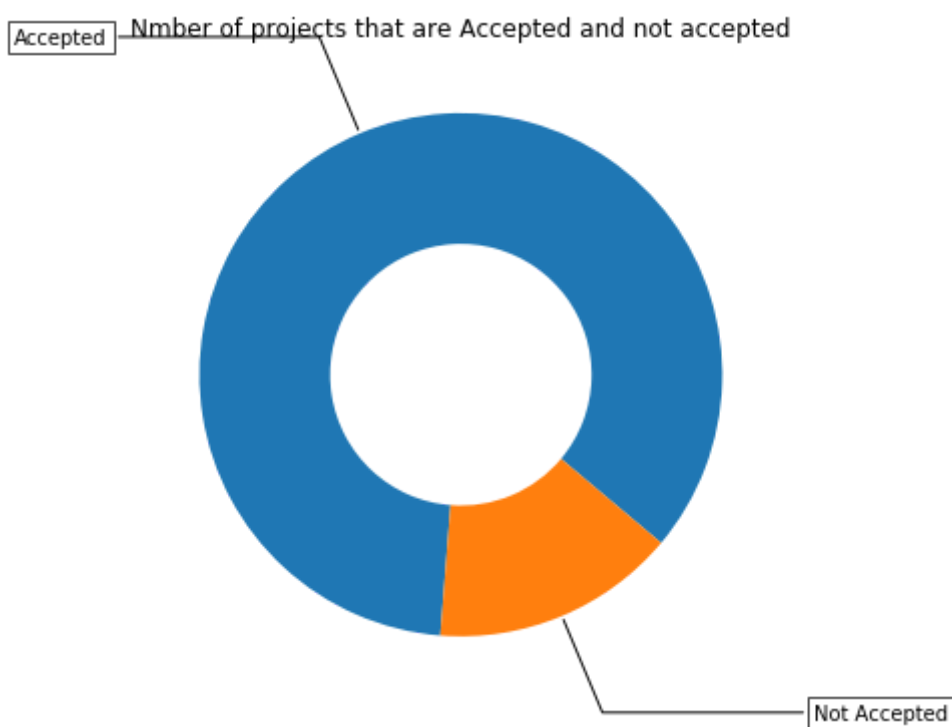
ax.set_title("Nmber of projects that are Accepted and not accepted")

plt.show()

```

Number of projects that are approved for funding 92706 , (84.85830404217 927 %)

Number of projects that are not approved for funding 16542 , (15.1416959 57820739 %)



1.2.1 Univariate Analysis: School State

In [6]:

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

temp = pd.DataFrame(project_data.groupby("school_state")["project_is_approved"].apply(n
p.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(84,39,14
3)']]

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[6]:

```
'''# How to plot US state heatmap: https://datascience.stackexchange.com/a/9
620\n\nscl = [[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, \'r
gb(117,107,177)\'],[1.0, \'rgb(84,39,143)\']] \n\ndata = [ dict(\n      t
ype=\'choropleth\',\n      colorscale = scl,\n      autocolorscale = F
alse,\n      locations = temp[\'state_code\'],\n      z = temp[\'num_p
roposals\'].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 p
ro")\n    ) ] \n\nlayout = dict(\n      title = \'Project Proposals % of
Acceptance Rate by US States\',\n      geo = dict(\n          scope=
\'usa\',\n          projection=dict( type=\'albers usa\' ),\n
showlakes = 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 [7]:

```
# https://www.csi.cuny.edu/sites/default/files/pdf/administration/ops/2letterstabbrev.p
df
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]:

```
#stacked bar plots matplotlib: https://matplotlib.org/gallery/lines\_bars\_and\_markers/ba
r_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]:

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

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

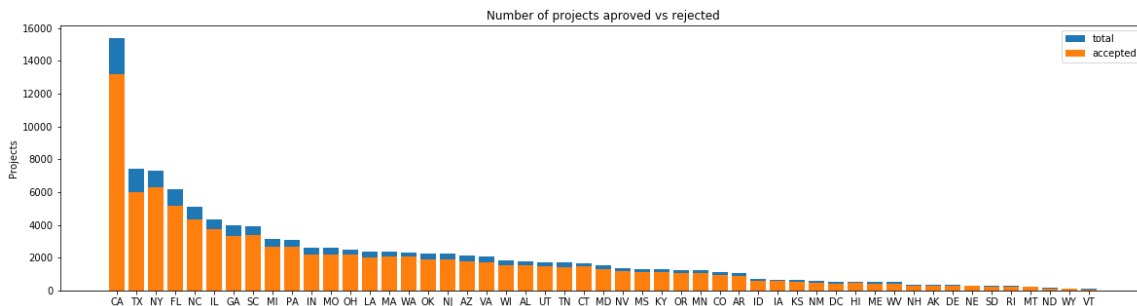
    temp.sort_values(by=['total'], inplace=True, ascending=False)

    if top:
        temp = temp[0:top]

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

In [10]:

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



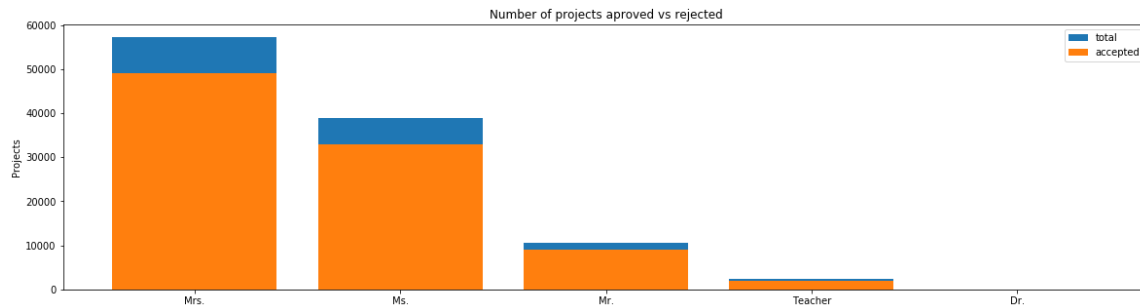
	school_state	project_is_approved	total	Avg
4	CA	13205	15388	0.858136
43	TX	6014	7396	0.813142
34	NY	6291	7318	0.859661
9	FL	5144	6185	0.831690
27	NC	4353	5091	0.855038
=====				
	school_state	project_is_approved	total	Avg
39	RI	243	285	0.852632
26	MT	200	245	0.816327
28	ND	127	143	0.888112
50	WY	82	98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

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)
```



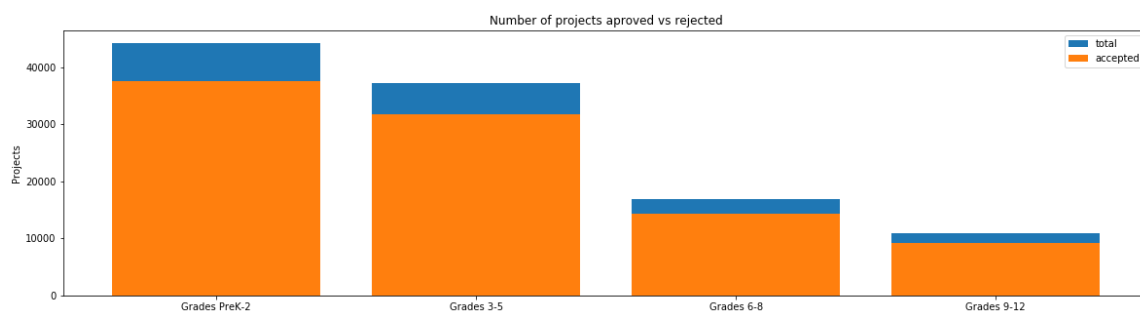
	teacher_prefix	project_is_approved	total	Avg
2	Mrs.	48997	57269	0.855559
3	Ms.	32860	38955	0.843537
1	Mr.	8960	10648	0.841473
4	Teacher	1877	2360	0.795339
0	Dr.	9	13	0.692308

	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 [12]:

```
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

SUMMARY: LOWER CLASSES ARE MORE LIKELY TO GET PROJECTS FUNDED.

1.2.4 Univariate Analysis: project subject categories

In [13]:

```
categories = 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 categories:
    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 category 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 replacing 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]:

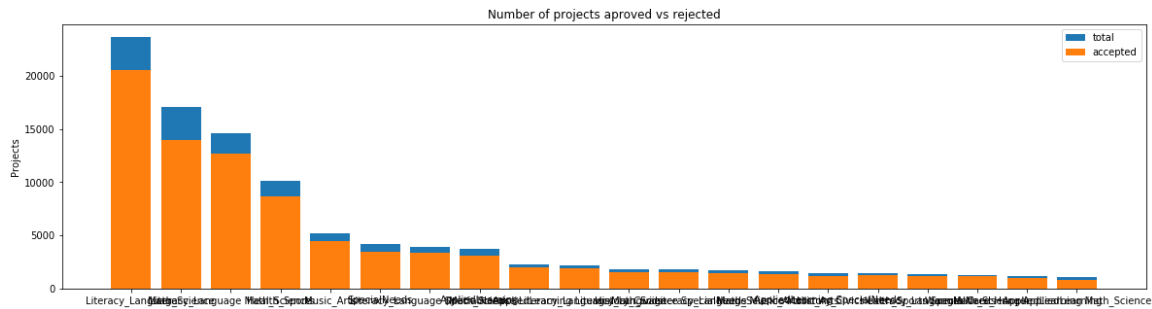
```
project_data['clean_categories'] = cat_list
project_data.drop(['project_subject_categories'], axis=1, inplace=True)
project_data.head(2)
```

Out[14]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	proj
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL

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

SUMMARY: here we can conclude that basic projects like literacy , language,health and hunger related subjects are more likely to be get funded.this is because people feel more obliged to hep those in need.

In [16]:

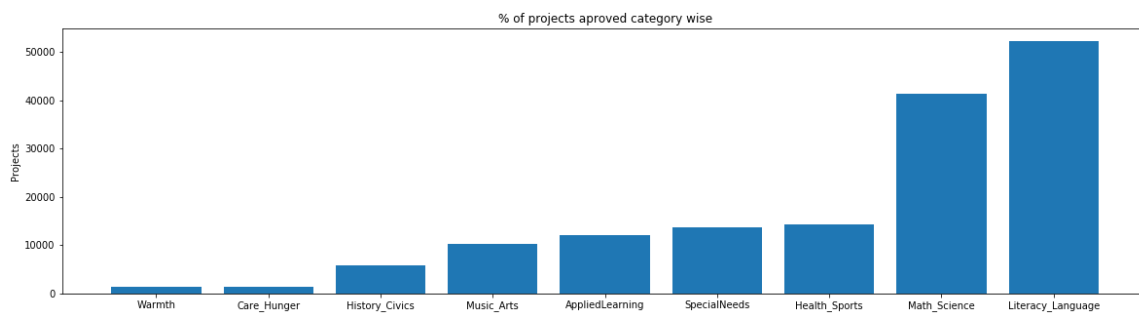
```
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()
for word in project_data['clean_categories'].values:
    my_counter.update(word.split())
```

In [17]:

```
# 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 [18]:

```
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
```

1.2.5 Univariate Analysis: project_subject_subcategories

In [19]:

```

sub_categories = 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_categories:
    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 category 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 replacing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " #"
    temp = temp.replace('&', '_')
    sub_cat_list.append(temp.strip())

```

In [20]:

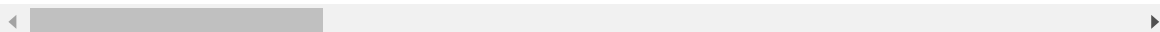
```

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

```

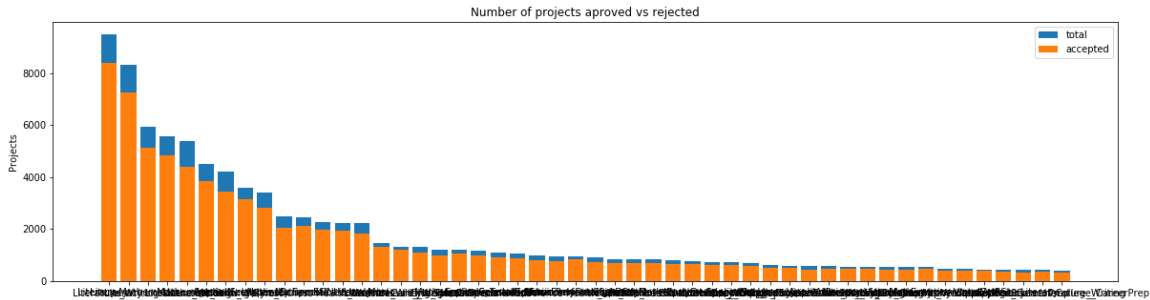
Out[20]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	proj
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL



In [21]:

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



	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	0.876
126				
127	ESL	349	421	0.828
979				
79	College_CareerPrep	343	421	0.814
727				
17	AppliedSciences Literature_Writing	361	420	0.859
524				
3	AppliedSciences College_CareerPrep	330	405	0.814
815				

In [22]:

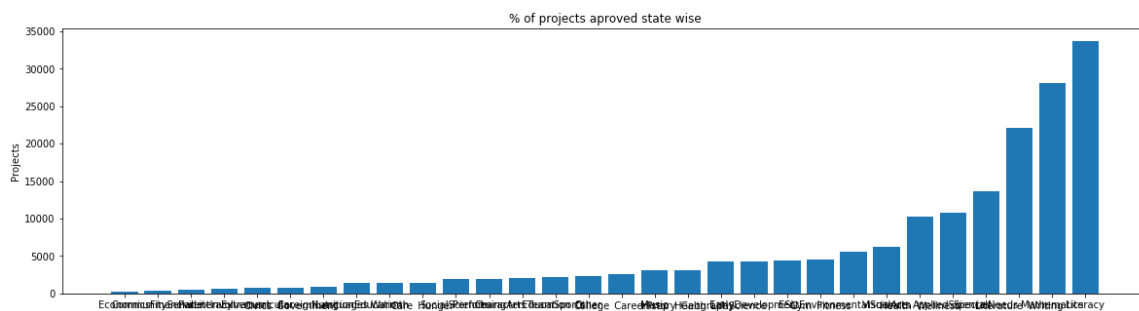
```
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()
for word in project_data['clean_subcategories'].values:
    my_counter.update(word.split())
```



```
# 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 approved state wise')
plt.xticks(ind, list(sorted_sub_cat_dict.keys()))
plt.show()
```



In [24]:

```
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

1.2.6 Univariate Analysis: Text features (Title)

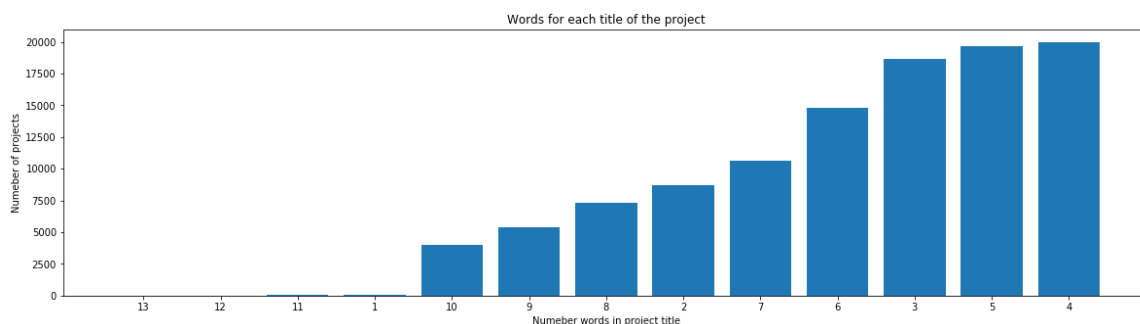
In [25]:

#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()
```



SUMMARY :

1. Most of the projects have 4 words in the title.
2. There are hardly any project titles containing more than 10 words.
3. Roughly most of the projects have 3, 4 or 5 words in the title.

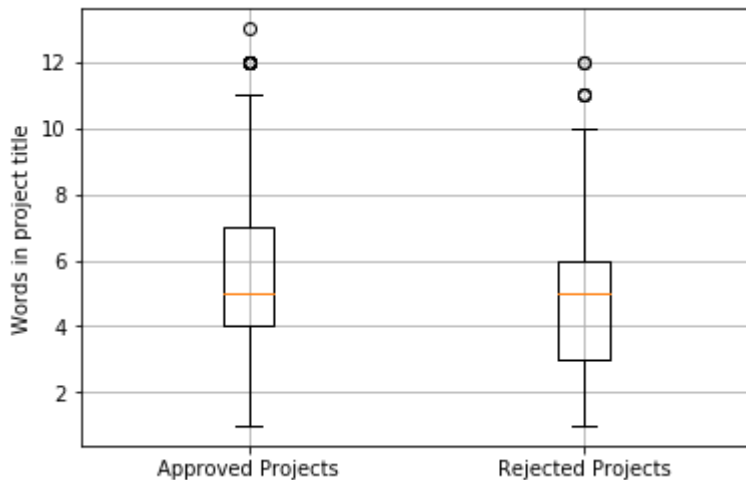
In [26]:

```
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]:

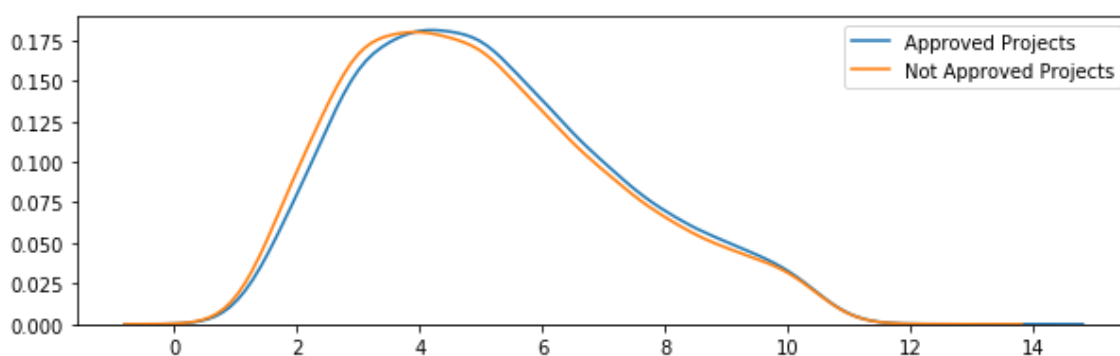
```
# https://glowingpython.blogspot.com/2012/09/boxplot-with-matplotlib.html
plt.boxplot([approved_title_word_count, rejected_title_word_count])
plt.xticks([1,2],('Approved Projects','Rejected Projects'))
plt.ylabel('Words in project title')
plt.grid()
plt.show()
```



SUMMARY : Approved projects have a slightly more number of words in the project essays when compared to the projects that have not been approved. This difference can be noticed in the plot

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()
```



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

In [29]:

```
# 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)
```

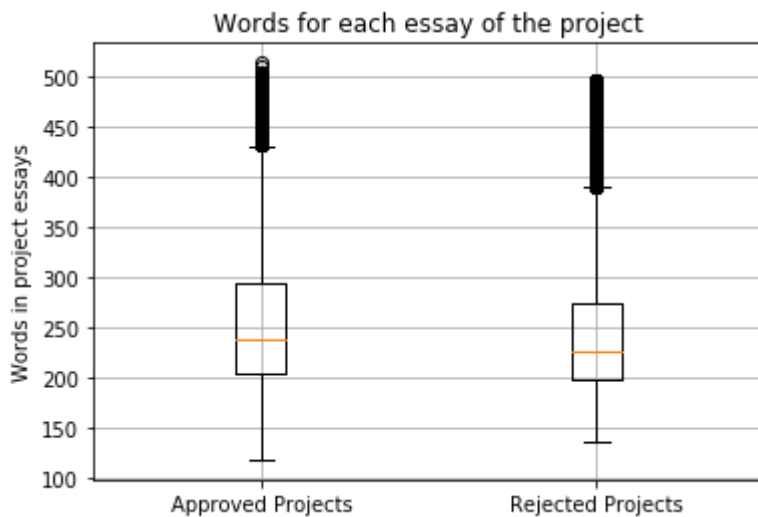
In [30]:

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

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

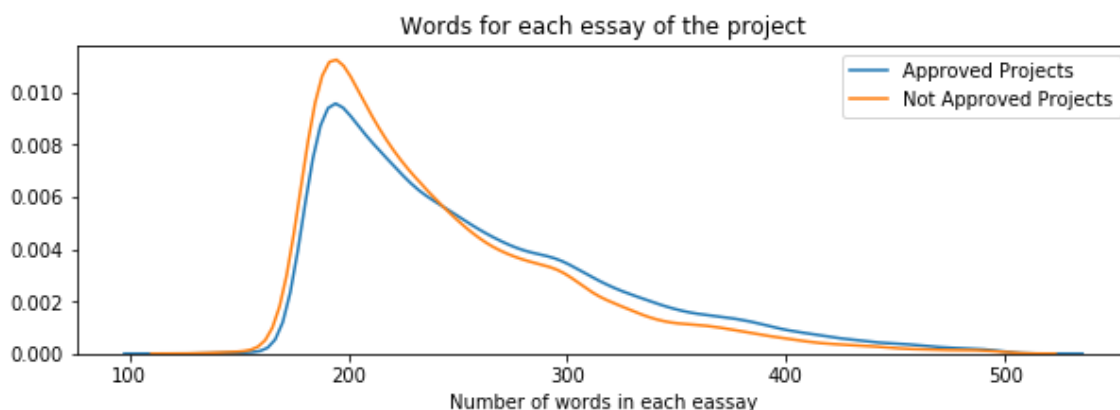
In [31]:

```
# 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 [32]:

```
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()
```



1.2.8 Univariate Analysis: Cost per project

In [33]:

```
# 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]:

```
# 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]:

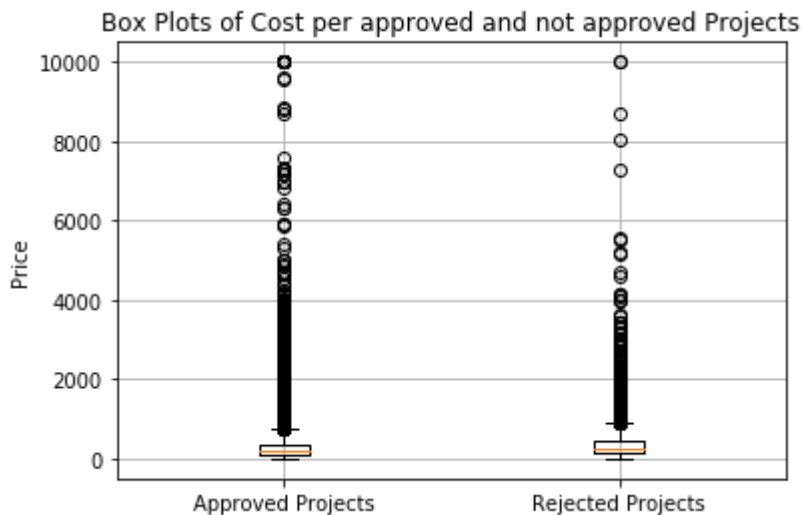
```
# join two dataframes in python:
project_data = pd.merge(project_data, price_data, on='id', how='left')
```

In [36]:

```
approved_price = project_data[project_data['project_is_approved']==1]['price'].values
rejected_price = project_data[project_data['project_is_approved']==0]['price'].values
```

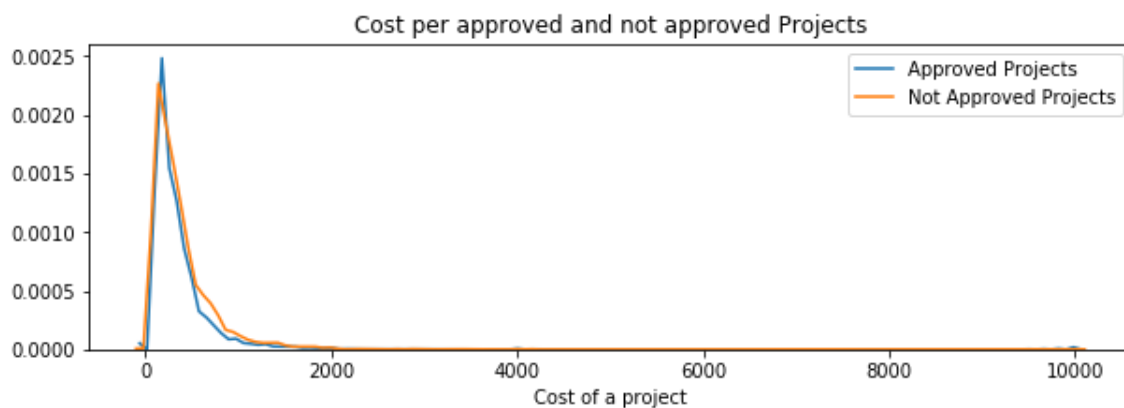
In [37]:

```
# 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()
```



In [39]:

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

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

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(r
ejected_price,i), 3)])
print(x)
```

Percentile	Approved Projects	Not Approved Projects
0	0.66	1.97
5	13.59	41.9
10	33.88	73.67
15	58.0	99.109
20	77.38	118.56
25	99.95	140.892
30	116.68	162.23
35	137.232	184.014
40	157.0	208.632
45	178.265	235.106
50	198.99	263.145
55	223.99	292.61
60	255.63	325.144
65	285.412	362.39
70	321.225	399.99
75	366.075	449.945
80	411.67	519.282
85	479.0	618.276
90	593.11	739.356
95	801.598	992.486
100	9999.0	9999.0

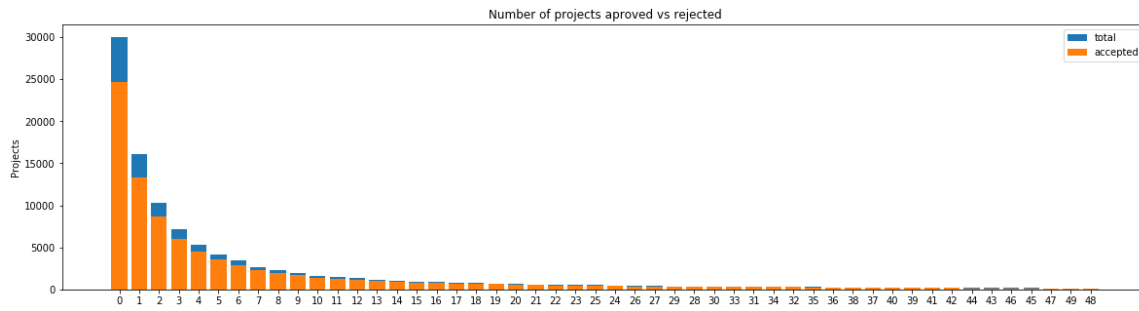
1.2.9 Univariate Analysis: teacher_number_of_previously_posted_projects

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

In [40]:

```
#print(project_data.head(2))
univariate_barplots(project_data, 'teacher_number_of_previously_posted_projects', 'project_is_approved', top=50)

#teachers who have submitted more projects before are more highly likely to be approved the current project.
```



teacher_number_of_previously_posted_projects	project_is_approved	tota
0	0	24652
1	1	13329
2	2	8705
3	3	5997
4	4	4452

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	tot
46	149	1
45	141	1
47	129	1
49	128	1
48	135	1

Avg

46	0.908537
45	0.921569
47	0.895833
49	0.895105
48	0.964286

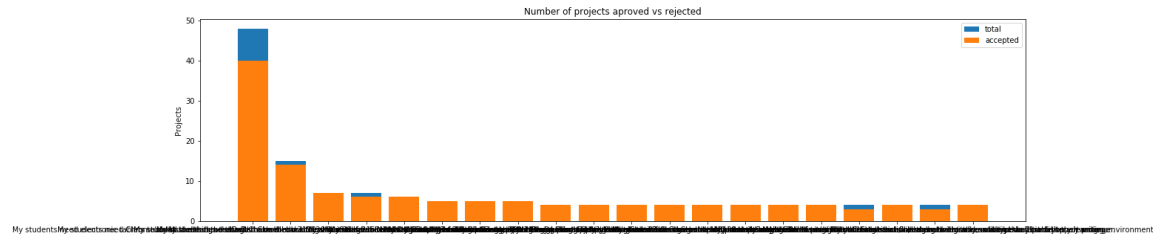
SUMMARY : 1.We observe that it is not mandatory for a teacher to have proposed any project prior. Maximum number of teachers, New talent and efforts are well appreciated. 2.Very few teachers who have proposed more than 20 projects have got approval. But the rate of approval is higher. 3.there is a wide range of previous submissions.

1.2.10 Univariate Analysis: project_resource_summary

In [41]:

```
univariate_barplots(project_data, 'project_resource_summary', 'project_is_approved', to
p=20)
```

#normal bivariate plot without considering the numbers in the text data



```

                                project_resource_summary  project_is_appro
ved \
56539  My students need electronic tablets to do all ...
40
10193  My students need Chromebooks to do all the thi...
14
18828  My students need a Dell Chromebook 3120 and a ...
7
51417  My students need chromebooks to do all the thi...
6
18819  My students need a Dell Chromebook 3120 11 6 C...
6
```

	total	Avg
56539	48	0.833333
10193	15	0.933333
18828	7	1.000000
51417	7	0.857143
18819	6	1.000000

```

=====
                                project_resource_summary  project_is_appro
ved \
3452   My students need 3 Chromebook computers to eng...
4
10510  My students need Chromebooks to open the doors...
3
8058   My students need 8 Hokki Stools to allow them ...
4
91082  My students need tablets to do all the things ...
3
10007  My students need Chromebooks to access the Lex...
4
```

	total	Avg
3452	4	1.00
10510	4	0.75
8058	4	1.00
91082	4	0.75
10007	4	1.00

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.

In [42]:

```
## Let us separate the data and carry out our work only on the required Project Resource Summaries.  
# code snippets taken from https://github.com/shashimanyam  
  
summaries = []  
for z in project_data["project_resource_summary"] :  
    summaries.append(z)  
  
summaries[0:10]
```

Out[42]:

```
['My students need opportunities to practice beginning reading skills in English at home.',  
 'My students need a projector to help with viewing educational programs',  
 'My students need shine guards, athletic socks, Soccer Balls, goalie gloves, and training materials for the upcoming Soccer season.',  
 'My students need to engage in Reading and Math in a way that will inspire them with these Mini iPads!',  
 'My students need hands on practice in mathematics. Having fun and personalized journals and charts will help them be more involved in our daily Math routines.',  
 'My students need movement to be successful. Being that I have a variety of students that have all different types of needs, flexible seating would assist not only these students with special needs, but all students.',  
 'My students need some dependable laptops for daily classroom use for reading and math.',  
 'My students need ipads to help them access a world of online resources that will spark their interest in learning.',  
 "My students need three devices and three management licenses for small group's easy access to newly-implemented online programs--Go Noodle Plus, for increased in-class physical activity and Light Sail, an interactive reading program.",  
 'My students need great books to use during Independent Reading, Read Alouds, Partner Reading and Author Studies.']
```

In [43]:

```
## The Length of the obtained List of Project summaries should match the total number of project summaries in  
## the project data. Just to ensure  
len(summaries)
```

Out[43]:

109248

In [47]:

```
len(numeric_digits)
```

Out[47]:

109248

In [48]:

```
## Converting the key value pairs to 1 or 0 based on presence of Numeric Values.  
digit_in_summary = []  
for a in numeric_digits.values() :  
    if a > 0 :  
        digit_in_summary.append(1)  
    else :  
        digit_in_summary.append(0)
```

In [49]:

```
digit_in_summary[0:20]
```

Out[49]:

[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1]

In [50]:

```
project_data['digit_in_summary'] = digit_in_summary  
project_data.head(20)
```

Out[50]:

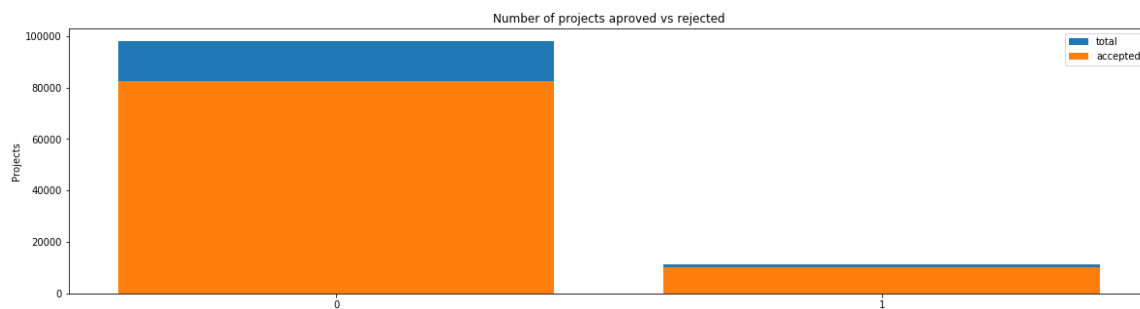
	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	pr
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ	
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY	
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX	
5	141660	p154343	a50a390e8327a95b77b9e495b58b9a6e	Mrs.	FL	
6	21147	p099819	9b40170bfa65e399981717ee8731efc3	Mrs.	CT	
7	94142	p092424	5bfd3d12fae3d2fe88684bbac570c9d2	Ms.	GA	
8	112489	p045029	487448f5226005d08d36bdd75f095b31	Mrs.	SC	
9	158561	p001713	140eeac1885c820ad5592a409a3a8994	Ms.	NC	

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	pr
10	43184	p040307	363788b51d40d978fe276bcb1f8a2b35	Mrs.	CA	
11	127083	p251806	4ba7c721133ef651ca54a03551746708	Ms.	CA	
12	19090	p051126	5e52c92b7e3c472aad247a239d345543	Mrs.	NY	
13	15126	p003874	178f6ae765cd4e0fb143a77c47fd65e2	Mrs.	OK	
14	62232	p233127	424819801de22a60bba7d0f4354d0258	Ms.	MA	
15	67303	p132832	bb6d6d054824fa01576ab38dfa2be160	Ms.	TX	
16	127215	p174627	4ad7e280fddff889e1355cc9f29c3b89	Mrs.	FL	
17	157771	p152491	e39abda057354c979c5b075cffbe5f88	Ms.	NV	
18	122186	p196421	fcd9b003fc1891383f340a89da02a1a6	Mrs.	GA	
19	146331	p058343	8e07a98deb1bc74c75b97521e05b1691	Ms.	OH	

20 rows × 21 columns

In [51]:

```
univariate_barplots(project_data, 'digit_in_summary', 'project_is_approved', top=2)
```



digit_in_summary	project_is_approved	total	Avg	
0	0	82563	98012	0.842376
1	1	10143	11236	0.902723

=====

digit_in_summary	project_is_approved	total	Avg	
0	0	82563	98012	0.842376
1	1	10143	11236	0.902723

INFERENCE: This plot clearly shows that the number of approval for those who clearly mention the quantity of the numbers are significantly high. This is because people usually choose to donate to projects which seem to be authentic and mentioning the numbers give the donor a clear idea of the context and scope of the project.

1.3 Text preprocessing

1.3.1 Essay Text

In [52]:

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

Out[52]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	proj
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL	

2 rows × 21 columns



In [53]:

```
# 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.\r\n\r\nThe limits of your language are the limits of your world.\r\n\r\n-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\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 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\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\n\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\n\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\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\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\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed r

aces 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 them because they develop their core, which enhances gross motor and in turn fine motor skills. \r\n\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\n\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 [54]:

```
# 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 [55]:

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

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. They want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them 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 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 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 [56]:

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

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. 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 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. They want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, 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. 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

In [57]:

```
#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 cognitive 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 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 They want to be able to move as they learn or so they say Wobble chairs are the answer and I love them because they develop their core 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 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

In [58]:

```
# 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'r
e", "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', 't
hey', 'them', 'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "th
at'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'ha
d', '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', 'ov
er', 'under', 'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'an
y', '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", 'no
w', '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', 'migh
tn', "mighntn't", 'mustn',\
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'w
asn', "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [59]:

```
# Combining all the above statements
from tqdm import tqdm
preprocessed_essays = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\n', ' ')
    sent = sent.replace('\\t', ' ')
    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 [02:16<00:00, 799.02it/s]
```

In [60]:

```
# after preprocessing  
preprocessed_essays[20000]
```

Out[60]:

```
'my kindergarten students varied disabilities ranging speech language dela  
ys cognitive delays gross fine motor delays autism they eager beavers alwa  
ys 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 le  
arn 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 de  
velop core enhances gross motor turn fine motor skills they also want lear  
n games kids not want sit worksheets they want learn count jumping playing  
physical engagement key success the number toss color shape mats make happ  
en my students forget work fun 6 year old deserves nannan'
```

1.3.2 Project title Text

In [61]:

```
# similarly you can preprocess the titles also
print(project_data.head(2))
```

```

    Unnamed: 0      id      teacher_id teacher_prefix \
0      160221  p253737  c90749f5d961ff158d4b4d1e7dc665fc  Mrs.
1      140945  p258326  897464ce9ddc600bced1151f324dd63a    Mr.

    school_state  project_submitted_datetime  project_grade_category \
0             IN      2016-12-05 13:43:57      Grades PreK-2
1             FL      2016-10-25 09:22:10      Grades 6-8

                                project_title \
0  Educational Support for English Learners at Home
1                Wanted: Projector for Hungry Learners

                                project_essay_1 \
0  My students are English learners that are work...
1  Our students arrive to our school eager to lea...

                                project_essay_2      ... \
0  \"The limits of your language are the limits o...      ...
1  The projector we need for our school is very c...      ...

    project_essay_4      project_resource_summary \
0      NaN  My students need opportunities to practice beg...
1      NaN  My students need a projector to help with view...

    teacher_number_of_previously_posted_projects  project_is_approved \
0                                             0                      0
1                                             7                      1

                                clean_categories      clean_subcategories \
0      Literacy_Language      ESL Literacy
1  History_Civics Health_Sports  Civics_Government TeamSports

                                essay  price  quantity \
0  My students are English learners that are work...  154.6      23
1  Our students arrive to our school eager to lea...  299.0      1

    digit_in_summary
0      0
1      0

[2 rows x 21 columns]
```

In [62]:

```

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
=====

```

In [63]:

```

sent = decontracted(project_data['project_title'].values[10000])
print(sent)
print("="*50)

```

```

Family Book Clubs
=====

```

In [64]:

```

sent = sent.replace('\\r', ' ')
sent = sent.replace('\\n', ' ')
sent = sent.replace('\\t', ' ')
print(sent)

```

```

Family Book Clubs

```

In [65]:

```

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

```

```

Family Book Clubs

```

In [66]:

```
# 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'r
e", "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', 't
hey', 'them', 'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "th
at'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'ha
d', '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', 'ov
er', 'under', 'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'an
y', '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", 'no
w', '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', 'migh
tn', "mighntn't", 'mustn',\
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'w
asn', "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [67]:

```
# Combining all the above statemennts
from tqdm import tqdm
preprocessed_title = []
# 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_title.append(sent.lower().strip())
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
109248/109248 [00:05<00:00, 19460.18it/s]
```

In [68]:

```
# after preprocessing  
preprocessed_title[40000]
```

Out[68]:

'let them read let them learn let them thrive'

1. 4 Preparing data for models

In [69]:

```
project_data.columns
```

Out[69]:

```
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',  
      'digit_in_summary'],  
      dtype='object')
```

In [70]:

```
#print(project_data.head(2))
```

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

In [71]:

```
# 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', 'AppliedLearnin
g', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encodig (109248, 9)
```

In [72]:

```
# 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'].value
s)
print("Shape of matrix after one hot encodig ",sub_categories_one_hot.shape)
```

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvemen
t', 'Extracurricular', 'Civics_Government', 'ForeignLanguages', 'Nutrition
Education', 'Warmth', 'Care_Hunger', 'SocialSciences', 'PerformingArts',
'CharacterEducation', 'TeamSports', 'Other', 'College_CareerPrep', 'Musi
c', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL',
'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'A
ppliedSciences', 'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Lit
eracy']
Shape of matrix after one hot encodig (109248, 30)
```

In [73]:

```
# Please do the similar feature encoding with state, teacher_prefix and project_grade_category also
# we use count vectorizer to convert the values into one hot encoded features
vectorizer = CountVectorizer( lowercase=False, binary=True)
vectorizer.fit(project_data['teacher_prefix'].values.astype('U'))
print(vectorizer.get_feature_names())

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

```
['Dr', 'Mr', 'Mrs', 'Ms', 'Teacher', 'nan']
Shape of matrix after one hot encoding (109248, 6)
```

In [74]:

```
# Please do the similar feature encoding with state, teacher_prefix and project_grade_category also
# we use count vectorizer to convert the values into one hot encoded features
vectorizer = CountVectorizer( lowercase=False, binary=True)
vectorizer.fit(project_data['school_state'].values)
print(vectorizer.get_feature_names())

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

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


In [133]:

```
# Please do the similar feature encoding with state, teacher_prefix and project_grade_category also
# we use count vectorizer to convert the values into one hot encoded features
my_counter = Counter()
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('\W', '_')

for project_grade in project_data['project_grade_category'].values:
    my_counter.update(project_grade.split())

print(my_counter)
project_grade_cat_dict = dict(my_counter)
#del project_grade_cat_dict['']
#sorted_project_grade_cat_dict = dict(sorted(project_grade_cat_dict.items(), key=lambda kv: kv[1]))

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

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

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

```
Counter({'GradesPreK2': 44225, 'Grades35': 37137, 'Grades68': 16923, 'Grades912': 10963})
['GradesPreK2', 'Grades68', 'Grades35', 'Grades912']
Shape of matrix after one hot encoding (109248, 4)
```

RESPECTIVE CATEGORICAL DATA COLUMNS HAVE BEEN SUCCESSFULLY VECTORIZED

1.4.2 Vectorizing Text data

1.4.2.1 Bag of words

In [91]:

```
# 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 encoding ", text_bow.shape)
```

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

1.4.2.2 Bag of Words on `project_title`

In [92]:

```
# you can vectorize the title also  
# before you vectorize the title make sure you preprocess it  
# We are considering only the words which appeared in at least 10 documents(rows or projects).  
vectorizer = CountVectorizer(min_df=10)  
title_bow = vectorizer.fit_transform(preprocessed_title)  
print("Shape of matrix after one hot encoding ",title_bow.shape)
```

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

1.4.2.3 TFIDF vectorizer

In [93]:

```
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 encoding ",text_tfidf.shape)
```

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

1.4.2.4 TFIDF Vectorizer on `project_title`

In [94]:

```
# Similarly you can vectorize for title also  
vectorizer = TfidfVectorizer(min_df=10)  
title_tfidf = vectorizer.fit_transform(preprocessed_title)  
print("Shape of matrix after one hot encoding ",title_tfidf.shape)
```

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

1.4.2.5 Using Pretrained Models: Avg W2V

In [95]:

```

'''
# 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-save-and-load-variables-in-python/

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

'''

```

Out[95]:

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/3823034
9/4084039\ndef loadGloveModel(gloveFile):\n    print ("Loading Glove Mode
l")\n    f = open(gloveFile,\'r\', encoding="utf8")\n    model = {}\n    f
or line in tqdm(f):\n        splitLine = line.split()\n        word = spli
tLine[0]\n        embedding = np.array([float(val) for val in splitLine
[1:]])\n        model[word] = embedding\n    print ("Done.",len(model)," w
ords loaded!")\n    return model\nmodel = loadGloveModel(\'glove.42B.300d.
txt\')\n\n# =====\nOutput:\n    \nLoading Glove Mod
el\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n# ===
=====
\n\nwords = []\nfor i in preproced_texts:\n    wor
ds.extend(i.split(\' \'))\n\nfor i in preproced_titles:\n    words.extend
(i.split(\' \'))\n\nprint("all the words in the coupus", len(words))\nwords
= set(words)\n\nprint("the unique words in the coupus", len(words))\n\ninter
_words = set(model.keys()).intersection(words)\n\nprint("The number of words
that are present in both glove vectors and our coupus", len(inter_wo
rds), "(" ,np.round(len(inter_words)/len(words)*100,3), "%")\n\nwords_courpu
s = {}\nwords_glove = set(model.keys())\nfor i in words:\n    if i in word
s_glove:\n        words_courpus[i] = model[i]\n\nprint("word 2 vec length",
len(words_courpus))\n\n\n# stronging variables into pickle files python: h
ttp://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-
python/\n\nimport pickle\nwith open(\'glove_vectors\', \'wb\') as f:\n
pickle.dump(words_courpus, f)\n\n\n'
```

In [96]:

```
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-p
ickle-to-save-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 [97]:

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

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

```
100%|████████████████████████████████████████████████████████████████████████████████| 109248/109248 [01:01<00:00, 1778.43it/s]
```

```
109248
300
```

1.4.2.6 Using Pretrained Models: AVG W2V on `project_title`

In [98]:

```
# Similarly you can vectorize for title also

# average Word2Vec
# compute average word2vec for each title.
avg_w2v_vectors = []; # the avg-w2v for each sentence is stored in this list
for sentence in tqdm(preprocessed_title): # for each 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
    for word in sentence.split(): # for each word in a 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]))
#print(avg_w2v_vectors[0])
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
109248/109248 [00:03<00:00, 28868.37it/s]

109248
300
```

1.4.2.7 Using Pretrained Models: TFIDF weighted W2V

In [99]:

```
# 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 [104]:

price_standardized

Out[104]:

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

In [105]:

```
prev_projects_scalar = StandardScaler()
prev_projects_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 : { prev_projects_scalar.mean_[0]}, Standard deviation : {np.sqrt( prev_projects_scalar.var_[0])}")

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

C:\Users\harsh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:59

5: DataConversionWarning:

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

Mean : 11.153165275336848, Standard deviation : 27.77702641477403

C:\Users\harsh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:59

5: DataConversionWarning:

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

In [106]:

prev_projects_standardized

Out[106]:

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

1.4.4 Merging all the above features

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

In [107]:

```
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 [108]:

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

Out[108]:

```
(109248, 16663)
```

Assignment 2: Apply TSNE

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

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

In [109]:

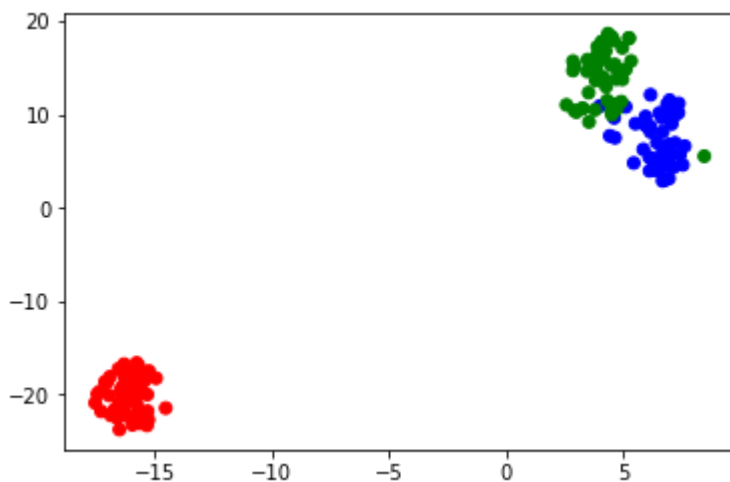
```
# this is the example code for TSNE
import numpy as np
from sklearn.manifold import TSNE
from sklearn import datasets
import pandas as pd
import matplotlib.pyplot as plt

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 [110]:

```
# please write all of the code with proper documentation and proper titles for each sub section  
# 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  
    #tsne plots for our data  
#some code snippets taken from https://github.com/shashimanyam/  
  
data = hstack((categories_one_hot, sub_categories_one_hot, school_state_one_hot, project_grade_category_one_hot, teacher_prefix_one_hot, price_standardized, prev_projects_standardized, title_bow))  
data.shape
```

Out[110]:

(109248, 3431)

In [112]:

```
from sklearn.manifold import TSNE  
data = data.tocsr()  
data_new = data[0:5000,:]  
#taking 5k data points
```

In [113]:

```
data_new = data_new.toarray()  
model = TSNE(n_components = 2, perplexity = 100.0, random_state = 0)  
tsne_data_b = model.fit_transform(data_new)
```

In [114]:

```
labels = project_data["project_is_approved"]  
labels_new = labels[0: 5000]  
len(labels_new)
```

Out[114]:

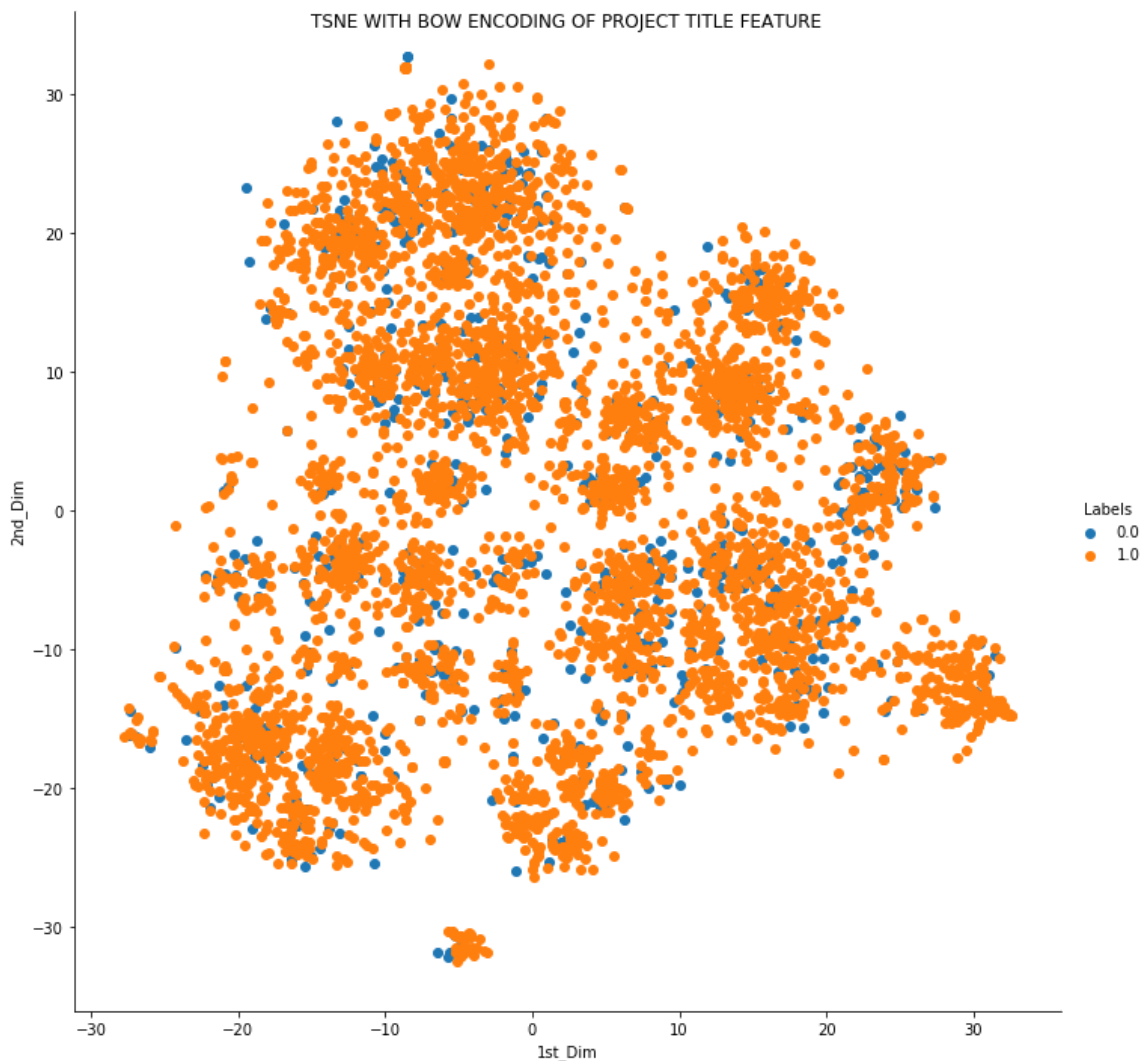
5000

In [115]:

```
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"))
```

In [116]:

```
sns.FacetGrid(tsne_df_b, hue = "Labels", height = 10).map(plt.scatter, "1st_Dim", "2nd_Dim").add_legend().fig.suptitle("TSNE WITH BOW ENCODING OF PROJECT TITLE FEATURE")  
plt.show()
```



In [117]:

```
tsne_df_b.shape
```

Out[117]:

(5000, 3)

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

In [118]:

```
# 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  
data = hstack((categories_one_hot, sub_categories_one_hot, school_state_one_hot, project_grade_category_one_hot, teacher_prefix_one_hot, price_standardized, prev_projects_standardized, title_tfidf))  
data.shape
```

Out[118]:

(109248, 3431)

In [105]:

```
from sklearn.manifold import TSNE  
data = data.tocsr()  
data_new = data[0:5000,:]
```

In [106]:

```
data_new = data_new.toarray()  
model = TSNE(n_components = 2, perplexity = 100.0, random_state = 0)  
tsne_data_tfidf = model.fit_transform(data_new)
```

In [107]:

```
labels = project_data["project_is_approved"]  
labels_new = labels[0: 5000]  
len(labels_new)
```

Out[107]:

5000

In [108]:

```
tsne_data_tfidf = np.vstack((tsne_data_tfidf.T, labels_new)).T  
tsne_df_tfidf = pd.DataFrame(tsne_data_tfidf, columns = ("1st_Dim", "2nd_Dim", "Labels"))
```

In [109]:

```
sns.FacetGrid(tsne_df_tfidf, hue = "Labels", height = 10).map(plt.scatter, "1st_Dim",  
"2nd_Dim").add_legend().fig.suptitle("TSNE WITH TFIDF ENCODING OF PROJECT TITLE FEATURE  
")  
plt.show()
```



In [110]:

```
tsne_df_tfidf.shape
```

Out[110]:

(5000, 3)

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

In [111]:

```
data = hstack((categories_one_hot, sub_categories_one_hot, school_state_one_hot, project_grade_category_one_hot, teacher_prefix_one_hot, price_standardized, prev_projects_standardized, avg_w2v_vectors))
data.shape
```

Out[111]:

(109248, 401)

In [112]:

```
# 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
from sklearn.manifold import TSNE
data = data.tocsr()
data_new = data[0:5000,:]
```

In [113]:

```
data_new = data_new.toarray()
model = TSNE(n_components = 2, perplexity = 100.0, random_state = 0)
tsne_data_avg_w2v_vectors = model.fit_transform(data_new)
```

In [114]:

```
labels = project_data["project_is_approved"]
labels_new = labels[0: 5000]
len(labels_new)
```

Out[114]:

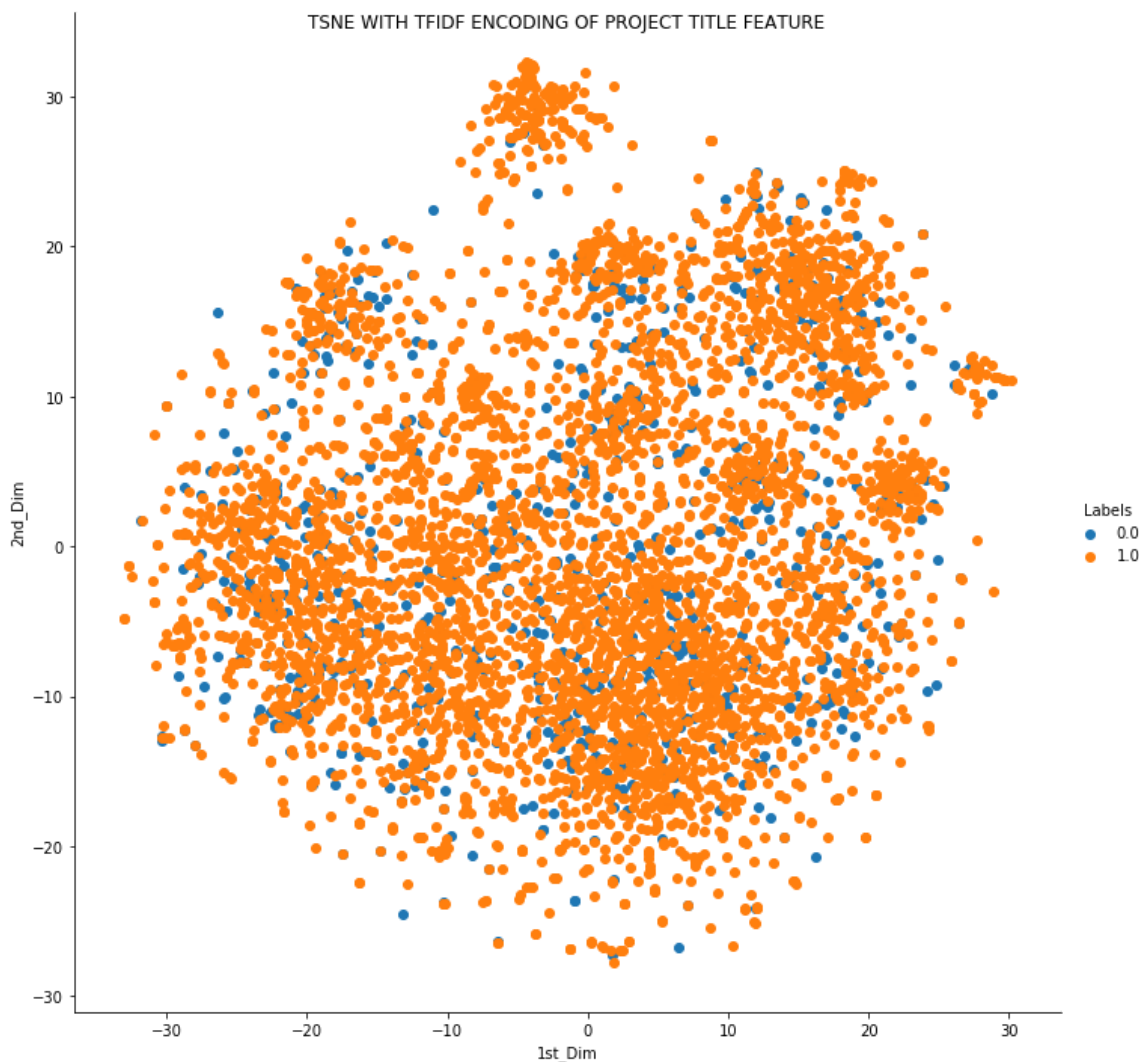
5000

In [115]:

```
tsne_data_avg_w2v_vectors = np.vstack((tsne_data_avg_w2v_vectors.T, labels_new)).T
tsne_df_avg_w2v_vectors = pd.DataFrame(tsne_data_avg_w2v_vectors, columns = ("1st_Dim", "2nd_Dim", "Labels"))
```

In [116]:

```
sns.FacetGrid(tsne_df_avg_w2v_vectors, hue = "Labels", height = 10).map(plt.scatter, "1st_Dim", "2nd_Dim").add_legend().fig.suptitle("TSNE WITH TFIDF ENCODING OF PROJECT TITLE FEATURE ")  
plt.show()
```



In [117]:

```
tsne_df_avg_w2v_vectors.shape
```

Out[117]:

(5000, 3)

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

In [118]:

```
# 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  
data = hstack((categories_one_hot, sub_categories_one_hot, school_state_one_hot, project_grade_category_one_hot, teacher_prefix_one_hot, price_standardized, prev_projects_standardized, tfidf_w2v_vectors))  
data.shape
```

Out[118]:

(109248, 401)

In [119]:

```
from sklearn.manifold import TSNE  
data = data.tocsr()  
data_new = data[0:5000,:]
```

In [120]:

```
data_new = data_new.toarray()  
model = TSNE(n_components = 2, perplexity = 100.0, random_state = 0)  
tsne_data_tfidf_w2v_vectors = model.fit_transform(data_new)
```

In [121]:

```
labels = project_data["project_is_approved"]  
labels_new = labels[0: 5000]  
len(labels_new)
```

Out[121]:

5000

In [122]:

```
tsne_data_tfidf_w2v_vectors = np.vstack((tsne_data_tfidf_w2v_vectors.T, labels_new)).T  
tsne_df_tfidf_w2v_vectors = pd.DataFrame(tsne_data_tfidf_w2v_vectors, columns = ("1st_Dim", "2nd_Dim", "Labels"))
```

In [123]:

```
sns.FacetGrid(tsne_df_tfidf_w2v_vectors, hue = "Labels", height = 10).map(plt.scatter,  
"1st_Dim", "2nd_Dim").add_legend().fig.suptitle("TSNE WITH TFIDF ENCODING OF PROJECT TI  
TLE FEATURE ")  
plt.show()
```



In [124]:

```
tsne_df_avg_w2v_vectors.shape
```

Out[124]:

(5000, 3)

2.5 Summary

few sentences about the results that obtained and the observations made.

1. Some states like Delaware, Washington and North Dakota have very high acceptance rate while states like Vermont and Texas have the lowest approvals. 2. A lot of projects are accepted for smaller classes. 3. Projects from females are accepted at a higher rate than men. 4. Projects belonging to the Literacy and Language categories have the highest number of projects proposed under. 5. The sub-Category Literacy has the highest number of projects approved. 6. Roughly most of the projects have 3, 4 or 5 words in the title. 7. The number of words in the Project Essays of Approved Projects are slightly more than the number of words in the Project Essays of the Rejected Projects. 8. We observe that it is not mandatory for a teacher to have proposed any project prior. New talent and efforts are well appreciated, but older submitters are respected as well. 9. The project summaries containing numeric values have a very high acceptance rate. 10. Visualisation of TSNE with Bag of Words, TF-IDF, Avg Word2Vec, TF-IDF Weighted Word2Vec does not seem to yield the expected result of clustering similar data points. Hence we would have to try any other method.

In []: