# Applying tsne on Donors Choose

## 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 (https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes)](https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes)). **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[*] |
| `project_essay_4` | Fourth application essay[*] |
| `project_submitted_datetime` | Datetime when project application was submitted. **Example:** `2016-04-28 12:43:56.245` |
| `teacher_id` | A unique identifier for the teacher of the proposed project. **Example:** `bdf8baa8fedef6bfeec7ae4ff1c15c56` |

| | |
|---|---|
| **teacher_prefix** | Teacher's title. One of the following enumerated values:<br><br>• nan<br>• Dr.<br>• Mr.<br>• Mrs.<br>• Ms.<br>• Teacher. |
| **teacher_number_of_previously_posted_projects** | Number of project applications previously submitted by the same teacher. **Example:** 2 |

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

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

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

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

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

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

## Notes on the Essay Data

Prior to May 17, 2016, the prompts for the essays were as follows:
- __project_essay_1:__ "Introduce us to your classroom"
- __project_essay_2:__ "Tell us more about your students"
- __project_essay_3:__ "Describe how your students will use the materials you're requesting"
- __project_essay_4:__ "Close by sharing why your project will make a difference"

Starting on May 17, 2016, the number of essays was reduced from 4 to 2, and the prompts for the first 2 essays were changed to the following:
- __project_essay_1:__ "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- __project_essay_2:__ "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

## Importing libraries

```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
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


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

## Reading Data

In [3]:

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

In [4]:

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

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

In [5]:

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

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

Out[5]:

| | 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 |

## Data Analysis

In [6]:

```python
# PROVIDE CITATIONS TO YOUR CODE IF YOU TAKE IT FROM ANOTHER WEBSITE.
# https://matplotlib.org/gallery/pie_and_polar_charts/pie_and_donut_labels.html#sphx-glr-gallery-pie-and-polar-ch
arts-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_va
lue_counts[1]+y_value_counts[0]))*100,"%)")

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

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

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

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

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

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

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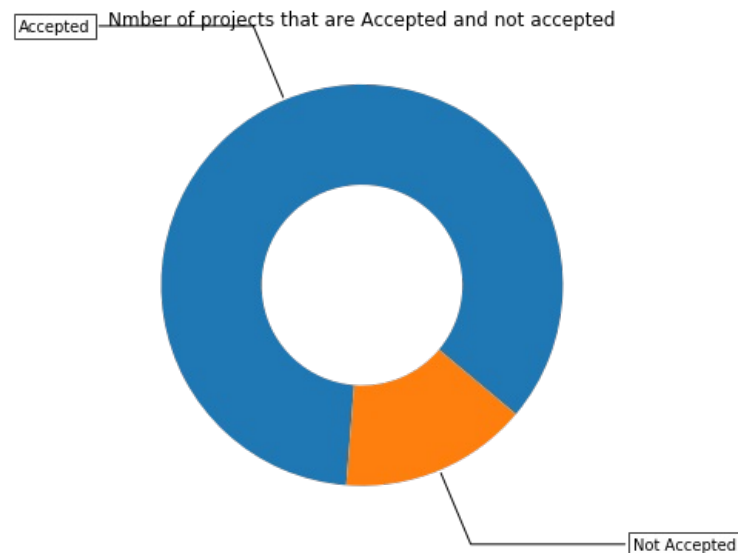
```
Number of projects thar are approved for funding  92706 , ( 84.85830404217927 %)
Number of projects thar are not approved for funding  16542 , ( 15.141695957820739 %)
```



## Univariate Analysis: School State

```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')
'''
```

Out[7]:

```
'# How to plot US state heatmap: https://datascience.stackexchange.com/a/9620\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, \'rgb(117,107,177)\'],[1.0, \'rgb(84,39,143)\']]\n\ndata = [ dict(\n        type=\
'choropleth\',\n          colorscale = scl,\n          autocolorscale = False,\n          locations = temp
[\'state_code\'],\n          z = temp[\'num_proposals\'].astype(float),\n          locationmode = \'USA-
states\',\n          text = temp[\'state_code\'],\n          marker = dict(line = dict (color = \'rgb(25
5,255,255)\',width = 2)),\n          colorbar = dict(title = "% of pro")\n     ) ]\n\nlayout = dict(\n
title = \'Project Proposals % of Acceptance Rate by US States\',\n          geo = dict(\n               s
cope=\'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)\
nofflineiplot(fig, filename=\'us-map-heat-map\')\n'
```

In [8]:

```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
```

```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()
```

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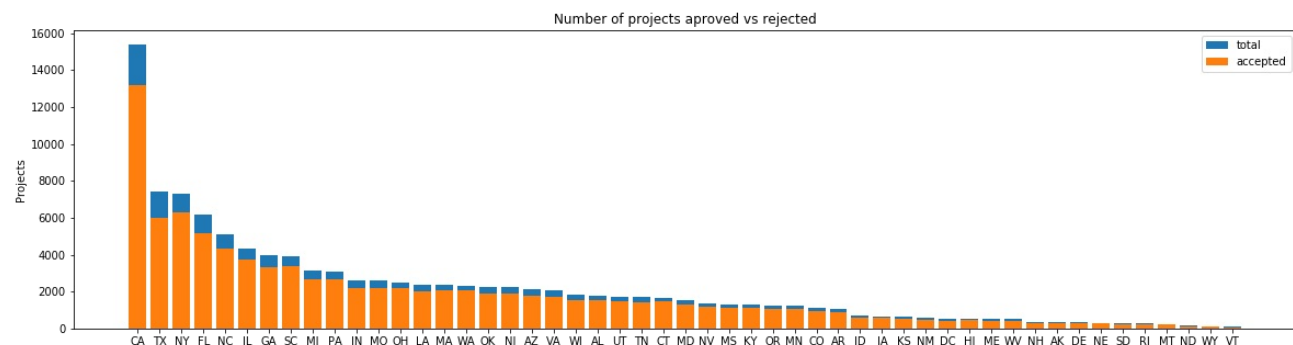
```python
univariate_barplots(project_data, 'school_state', 'project_is_approved', False)
```



```
   school_state  project_is_approved  total       Avg
4            CA                13205  15388  0.858136
43           TX                 6014   7396  0.813142
34           NY                 6291   7318  0.859661
9            FL                 5144   6185  0.831690
27           NC                 4353   5091  0.855038
==================================================
   school_state  project_is_approved  total       Avg
39           RI                  243    285  0.852632
26           MT                  200    245  0.816327
28           ND                  127    143  0.888112
50           WY                   82     98  0.836735
46           VT                   64     80  0.800000
```
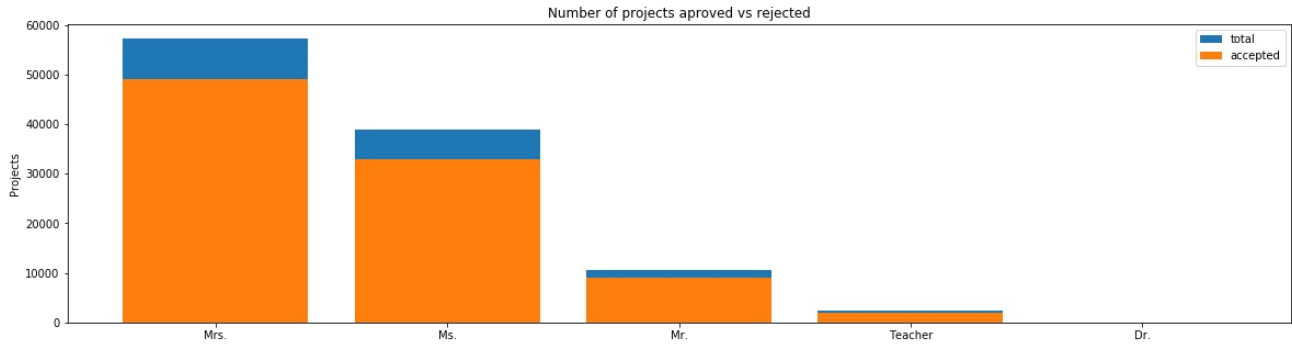
**SUMMARY: Every state has greater than 80% success rate in approval**

## Univariate Analysis: teacher_prefix

```
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
```
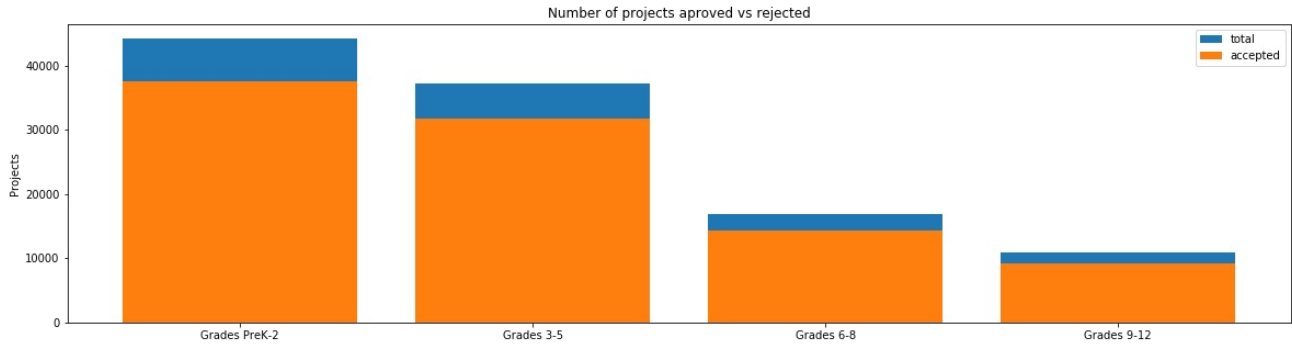
## Univariate Analysis: project_grade_category

In [13]:

```
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
```

## Univariate Analysis: project_subject_categories

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

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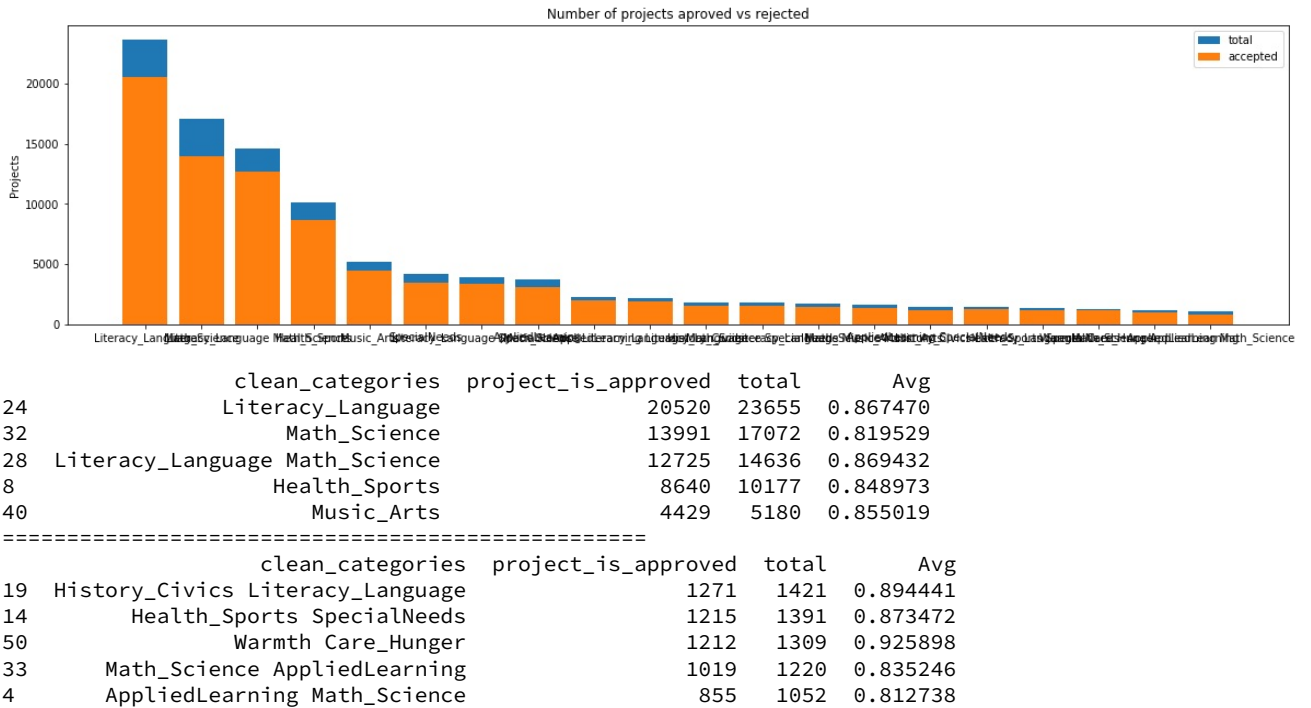```

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

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_ca |
|---|---|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | 2016-12-05 13:43:57 | Grades |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 2016-10-25 09:22:10 | Gra |

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



```
        clean_categories  project_is_approved  total       Avg
24      Literacy_Language                20520  23655  0.867470
32           Math_Science                13991  17072  0.819529
28  Literacy_Language Math_Science       12725  14636  0.869432
8           Health_Sports                 8640  10177  0.848973
40             Music_Arts                 4429   5180  0.855019
=================================================
        clean_categories  project_is_approved  total       Avg
19  History_Civics Literacy_Language         1271   1421  0.894441
14      Health_Sports SpecialNeeds           1215   1391  0.873472
50          Warmth Care_Hunger               1212   1309  0.925898
33      Math_Science AppliedLearning         1019   1220  0.835246
4       AppliedLearning Math_Science          855   1052  0.812738
```
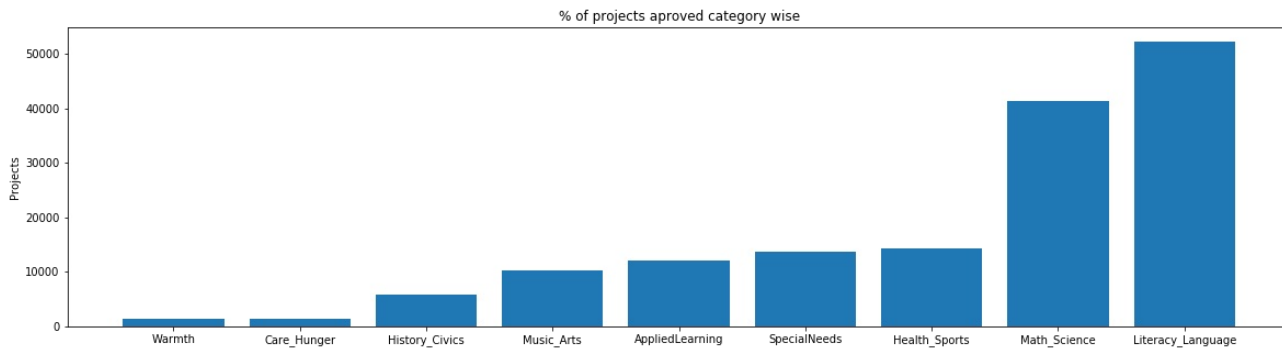
```python
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()
for word in project_data['clean_categories'].values:
    my_counter.update(word.split())
```

```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()
```

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

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

## Univariate Analysis: project_subject_subcategories

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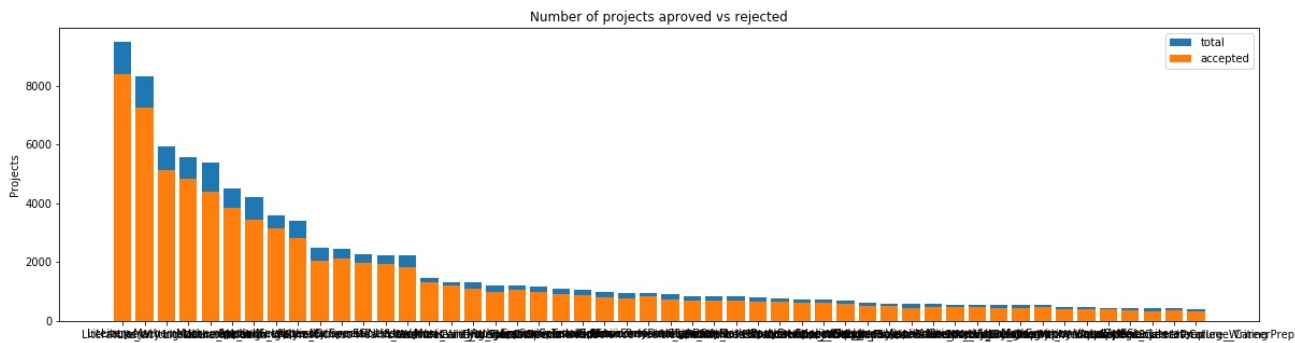
```python
project_data['clean_subcategories'] = sub_cat_list
project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
project_data.head(2)
```

Out[21]:

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

In [22]:

```python
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  0.876126
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
```
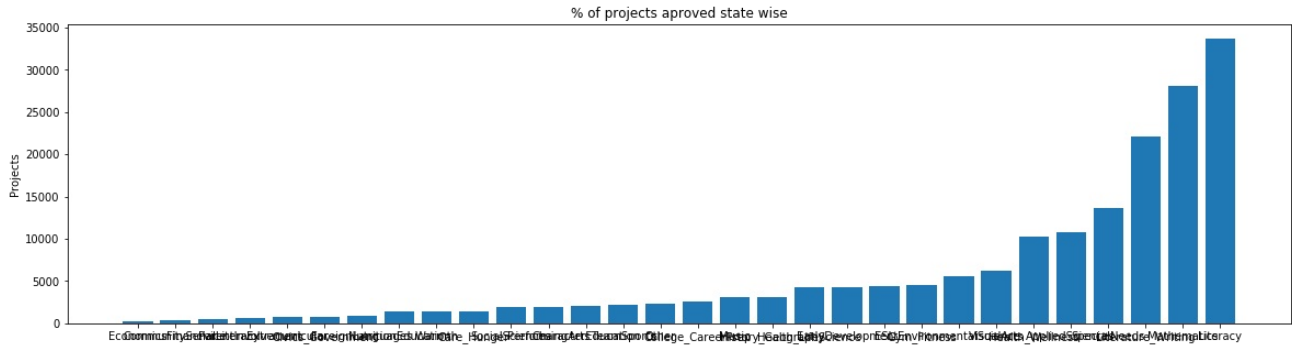
In [23]:

```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())
```

```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()
```

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

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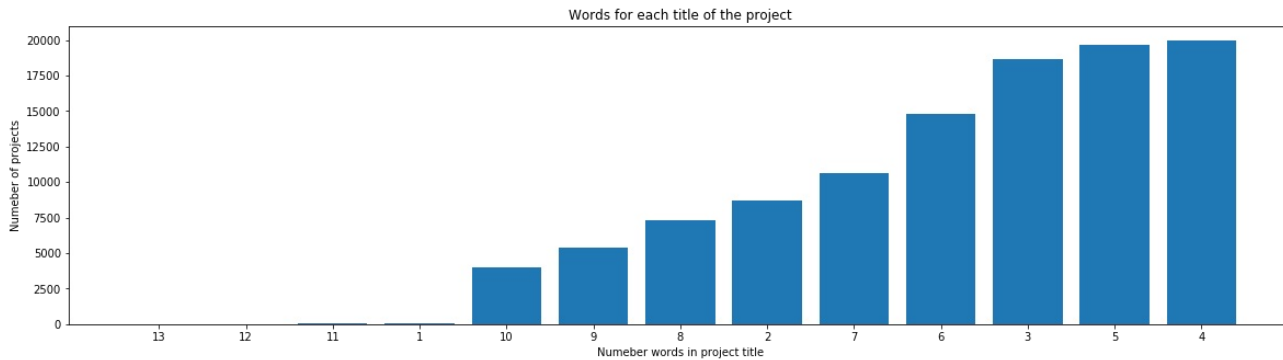
## Univariate Analysis: Text features (Title)

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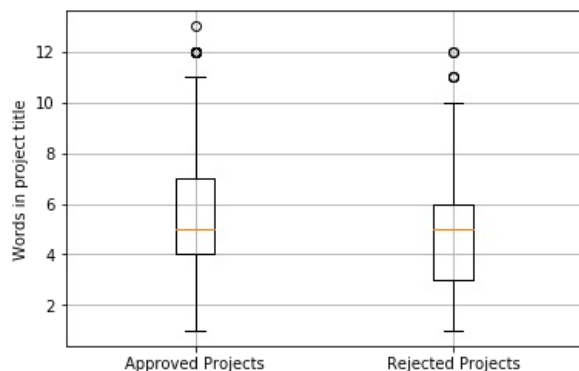
```
approved_title_word_count = project_data[project_data['project_is_approved']==1]['project_title'].str.split().app
ly(len)
approved_title_word_count = approved_title_word_count.values

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

```
# 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()
```



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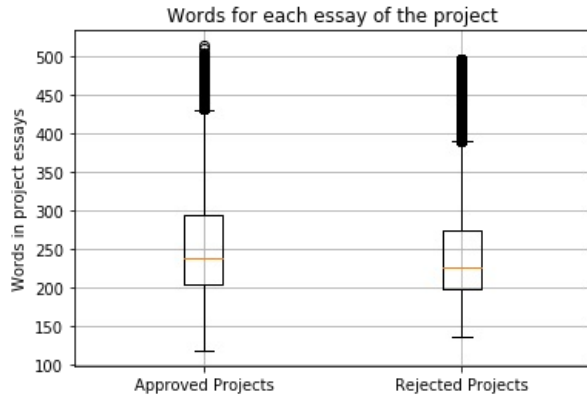
```
approved_word_count = project_data[project_data['project_is_approved']==1]['essay'].str.split().apply(len)
approved_word_count = approved_word_count.values

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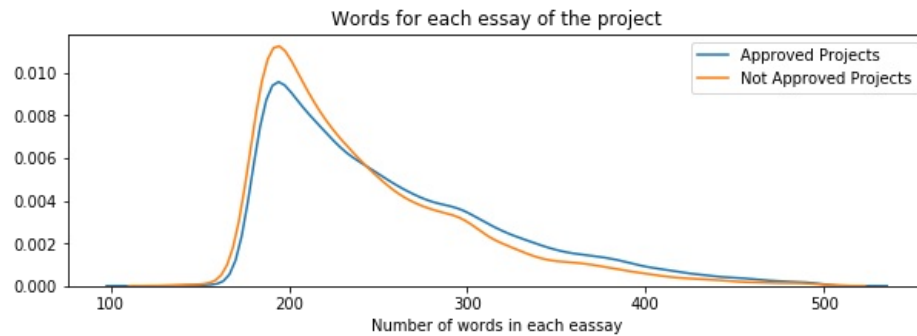
```
# https://glowingpython.blogspot.com/2012/09/boxplot-with-matplotlib.html
plt.boxplot([approved_word_count, rejected_word_count])
plt.title('Words for each essay of the project')
plt.xticks([1,2],('Approved Projects','Rejected Projects'))
plt.ylabel('Words in project essays')
plt.grid()
plt.show()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

```
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()
```



## Univariate Analysis: Cost per project

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

|   | 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 |

```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]:

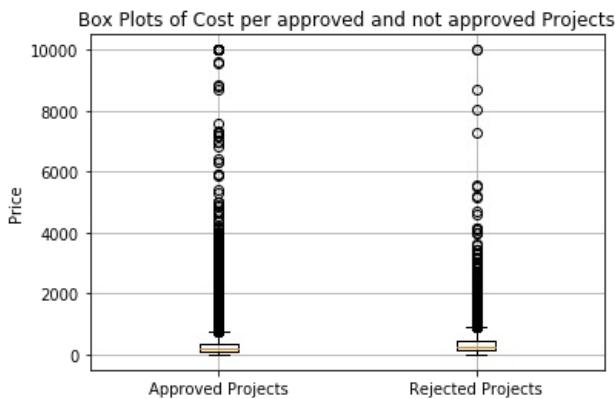|   | id | price | quantity |
|---|---|---|---|
| **0** | p000001 | 459.56 | 7 |
| **1** | p000002 | 515.89 | 21 |

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
```
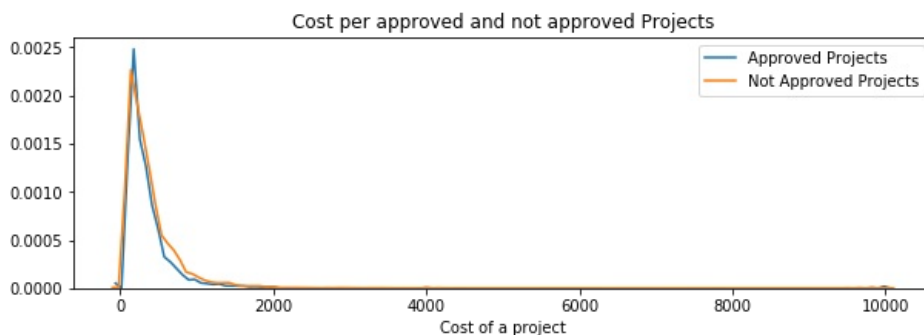
In [37]:

```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]:

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

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

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

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

for i in range(0,101,5):
    x.add_row([i,np.round(np.percentile(approved_price,i), 3), np.round(np.percentile(rejected_price,i), 3)])
print(x)
```
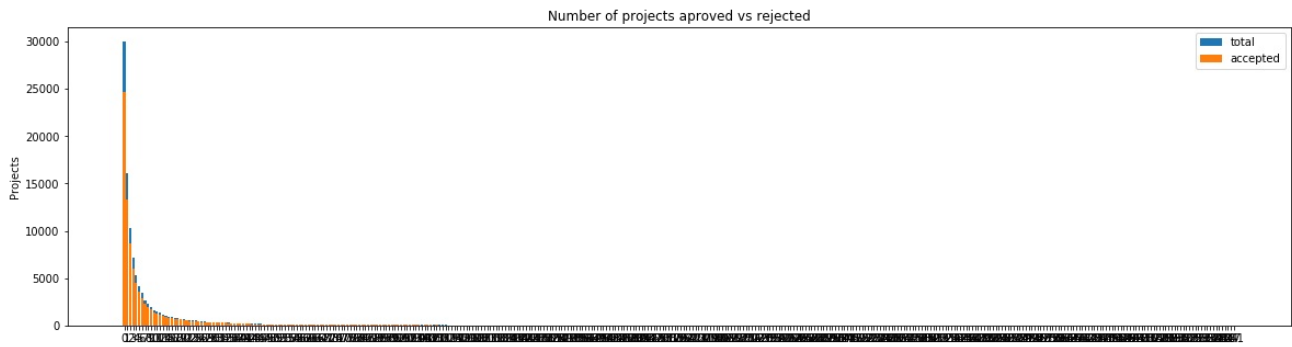
```
+------------+-------------------+-----------------------+
| Percentile | Approved Projects | Not Approved Projects |
+------------+-------------------+-----------------------+
|     0      |        0.66       |          1.97         |
|     5      |       13.59       |          41.9         |
|     10     |       33.88       |         73.67         |
|     15     |        58.0       |         99.109        |
|     20     |       77.38       |         118.56        |
|     25     |       99.95       |        140.892        |
|     30     |       116.68      |         162.23        |
|     35     |      137.232      |        184.014        |
|     40     |       157.0       |        208.632        |
|     45     |      178.265      |        235.106        |
|     50     |       198.99      |        263.145        |
|     55     |       223.99      |         292.61        |
|     60     |       255.63      |        325.144        |
|     65     |      285.412      |         362.39        |
|     70     |      321.225      |         399.99        |
|     75     |      366.075      |        449.945        |
|     80     |       411.67      |        519.282        |
|     85     |       479.0       |        618.276        |
|     90     |       593.11      |        739.356        |
|     95     |      801.598      |        992.486        |
|    100     |       9999.0      |         9999.0        |
+------------+-------------------+-----------------------+
```

**Univariate Analysis: teacher_number_of_previously_posted_projects**

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



Number of projects aproved vs rejected

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

        Avg
0  0.821350
1  0.830054
2  0.841063
3  0.843460
4  0.845423
=================================================
     teacher_number_of_previously_posted_projects  project_is_approved  total  \
242                                           242                    1      1
268                                           270                    1      1
234                                           234                    1      1
335                                           347                    1      1
373                                           451                    1      1

      Avg
242  1.0
268  1.0
234  1.0
335  1.0
373  1.0
```

## Summary

1.We observe that it is not mandatory for a teacher to have proposed any project prior. Maximum number of teachers, nearly 82% of the approved projects have been submitted by teachers with no prior project proposals. 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 given the teacher has proposed atleast 19 different projects. 3.There is alot of variability in the number of projects previously proposed by the teacher varying from 0 to more than 20.

## project_resource_summary

In [41]:
```
## Let us separate the data and carry out our work only on the required Project Resource Summaries.

summaries = []

for a in project_data["project_resource_summary"] :
    summaries.append(a)

summaries[0:10]
```

Out[41]:
```
['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 iPa
ds!',
 '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 inter
est in learning.',
 "My students need three devices and three management licenses for small group's easy access to newl
y-implemented online programs--Go Noodle Plus, for increased in-class physical activity and Light Sa
il, an interactive reading program.",
 'My students need great books to use during Independent Reading, Read Alouds, Partner Reading and A
uthor Studies.']
```

In [42]:
```
## 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[42]:
```
109248
```

In [43]:
```
## Identifying the numbers from the project summaries and storing the values as a key value pair in a dictionary
to
## avoid missing the position of the value within the huge ocean of summary data.

numeric_summary_values = {}

for x in tqdm(range(len(summaries))):
    for s in summaries[x].split():
        if s.isdigit() :
            numeric_summary_values[x] = int(s)
```

```
100%|████████████████████████████████| 109248/109248 [00:01<00:00, 79046.11it/s]
```

In [44]:
```
numeric_summary_values[14]
```

Out[44]:
```
5
```

In [45]:
```
## We only have the key value pairs for Summaries containing Numeric values, so in this step

numeric_digits = {}

for c in range(len(summaries)) :
    if c in numeric_summary_values.keys() :
        numeric_digits[c] = numeric_summary_values[c]
    else :
        numeric_digits[c] = 0
```

```
In [46]:
```

```
for i in range (20) :
    print(numeric_digits[i])
```

```
0
0
0
0
0
0
0
0
0
0
0
0
0
0
5
0
2
0
0
7
```

```
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(15)
```

Out[50]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_ |
|---|---|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | 2016-12-05 13:43:57 | Grad |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 2016-10-25 09:22:10 | G |
| 2 | 21895 | p182444 | 3465aaf82da834c0582ebd0ef8040ca0 | Ms. | AZ | 2016-08-31 12:03:56 | G |
| 3 | 45 | p246581 | f3cb9bffbba169bef1a77b243e620b60 | Mrs. | KY | 2016-10-06 21:16:17 | Grad |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 4 | 172407 | p104768 | be1f7507a41f8479dc06f047086a39ec | Mrs. | TX | 2016-07-11 01:10:09 | Grad |
| 5 | 141660 | p154343 | a50a390e8327a95b77b9e495b58b9a6e | Mrs. | FL | 2017-04-08 22:40:43 | G |
| 6 | 21147 | p099819 | 9b40170bfa65e399981717ee8731efc3 | Mrs. | CT | 2017-02-17 19:58:56 | G |
| 7 | 94142 | p092424 | 5bfd3d12fae3d2fe88684bbac570c9d2 | Ms. | GA | 2016-09-01 00:02:15 | G |
| 8 | 112489 | p045029 | 487448f5226005d08d36bdd75f095b31 | Mrs. | SC | 2016-09-25 17:00:26 | Grad |
| 9 | 158561 | p001713 | 140eeac1885c820ad5592a409a3a8994 | Ms. | NC | 2016-11-17 18:18:56 | Grad |
| 10 | 43184 | p040307 | 363788b51d40d978fe276bcb1f8a2b35 | Mrs. | CA | 2017-01-04 16:40:30 | G |
| 11 | 127083 | p251806 | 4ba7c721133ef651ca54a03551746708 | Ms. | CA | 2016-11-14 22:57:28 | Grad |
| 12 | 19090 | p051126 | 5e52c92b7e3c472aad247a239d345543 | Mrs. | NY | 2016-05-23 15:46:02 | G |
| 13 | 15126 | p003874 | 178f6ae765cd4e0fb143a77c47fd65e2 | Mrs. | OK | 2016-10-17 09:49:27 | Grad |
| 14 | 62232 | p233127 | 424819801de22a60bba7d0f4354d0258 | Ms. | MA | 2017-02-14 16:29:10 | Grad |

15 rows x 21 columns

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

Number of projects aproved vs rejected



```
     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
```

## SUMMARY

1.The project summaries containing numeric values have a very high acceptance rate of 90%. Well, proper numbered requirements suggest clarity in the proposals and hence Alot of people tend to donate for a better cause, that is to help children

2.It is obvious from the graph that majority of the projects do not have numeric values stating the requirement of certain products

# Text preprocessing

## Project_title

In [52]:

```python
# printing some random essays.
print(project_data['project_title'].values[9])
print("="*50)
print(project_data['project_title'].values[34])
print("="*50)
print(project_data['project_title'].values[79])
print("="*50)
print(project_data['project_title'].values[101])
print("="*50)
print(project_data['project_title'].values[1111])
print("="*50)
```

```
Just For the Love of Reading--\r\nPure Pleasure
=================================================
\"Have A Ball!!!\"
=================================================
Make Music, Make Our Year!
=================================================
Fun & Physically Fit
=================================================
\"Chrome Bookworms\" Readers and Writers of the Future
=================================================
```

```
In [53]:
# 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 [54]:
sent = decontracted(project_data['project_title'].values[34])
print(sent)
print("="*50)
```

```
\"Have A Ball!!!\"
==================================================
```

```
In [55]:
# \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)
```

```
 Have A Ball!!!
```

```
In [56]:
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

```
 Have A Ball
```

```
In [57]:
# 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', 'aft\
er',\
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'fu\
rther',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few',\
'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', '\
re', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn\
',\
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn',\
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "\
weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [58]:

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

```
100%|████████████████████████████| 109248/109248 [00:07<00:00, 14988.72it/s]
```

In [59]:

```python
print(preprocessed_title[34])
print('-'*50)
print(preprocessed_title[1111])
```

```
have a ball
--------------------------------------------------
chrome bookworms readers writers future
```

## preprocessing of school_state

In [60]:

```python
school_state = list(project_data['school_state'].values)

school_state_list = []

for state in school_state:
    school_state_list.append(state.lower())

# Now replace the "school_state" column by the cleaned one.
project_data['clean_school_state'] = school_state_list
project_data.drop(['school_state'], axis=1, inplace=True)
```

In [61]:

```python
#final result
clean_school_state = list(project_data['clean_school_state'].values)
print(list(set(clean_school_state)))
```

```
['tn', 'nc', 'ak', 'de', 'pa', 'fl', 'id', 'nd', 'nv', 'me', 'ms', 'oh', 'ca', 'wy', 'mt', 'tx', 'ga
', 'ok', 'wa', 'ct', 'la', 'nj', 'al', 'ia', 'or', 'ar', 'az', 'mn', 'nm', 'ny', 'vt', 'wi', 'ut', '
ne', 'nh', 'sd', 'mo', 'ma', 'il', 'ri', 'sc', 'co', 'in', 'ky', 'wv', 'ks', 'hi', 'md', 'mi', 'dc',
'va']
```

## preprocessing of teacher_prefix

In [62]:

```python
##https://stackoverflow.com/questions/39303912/tfidfvectorizer-in-scikit-learn-valueerror-np-nan-is-an-invalid-do
cument
project_data['teacher_prefix'] = project_data['teacher_prefix'].apply(lambda x: np.str_(x))
```

In [63]:

```python
teacher_prefix = list(project_data['teacher_prefix'].values)

teacher_prefix_list = []

for prefix in teacher_prefix:
    prefix = prefix.replace('.','')
    teacher_prefix_list.append(prefix.lower())

# Now replace the "teacher_prefix" column by the cleaned one.
project_data['clean_teacher_prefix'] = teacher_prefix_list
project_data.drop(['teacher_prefix'], axis=1, inplace=True)
```

```
##final result
clean_teacher_prefix = list(project_data['clean_teacher_prefix'].values)
print(list(set(clean_teacher_prefix)))
```

```
['mr', 'mrs', 'nan', 'teacher', 'ms', 'dr']
```

## preprocessing of project_grade_category

```
project_grade_category = list(project_data['project_grade_category'].values)

project_grade_category_list = []

for pgc in project_grade_category:
    pgc = pgc.lower()
    pgc_replace = pgc.replace(' ', '_')
    pgc_final_replace = pgc_replace.replace('-', '_')

    project_grade_category_list.append(pgc_final_replace)

# Now replace the "school_state" column by the cleaned one.
project_data['clean_project_grade_category'] = project_grade_category_list
project_data.drop(['project_grade_category'], axis=1, inplace=True)
```

```
##final result
clean_project_grade_category = list(project_data['clean_project_grade_category'].values)
print(list(set(clean_project_grade_category)))
```

```
['grades_prek_2', 'grades_9_12', 'grades_6_8', 'grades_3_5']
```

# Preparing data for models

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

# Vectorizing Categorical data

```
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 encoding ",categories_one_hot.shape)
```

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

```
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 encoding ",sub_categories_one_hot.shape)
```

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular', 'Civi
cs_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger', 'SocialSciences',
'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other', 'College_CareerPrep', 'Music', 'Histo
ry_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL', 'Gym_Fitness', 'EnvironmentalScience
', 'VisualArts', 'Health_Wellness', 'AppliedSciences', 'SpecialNeeds', 'Literature_Writing', 'Mathem
atics', 'Literacy']
Shape of matrix after one hot encoding  (109248, 30)
```

```
#One Hot Encode – School States
my_counter = Counter()
for state in project_data['clean_school_state'].values:
    my_counter.update(state.split())
```

```
school_state_cat_dict = dict(my_counter)
sorted_school_state_cat_dict = dict(sorted(school_state_cat_dict.items(), key=lambda kv: kv[1]))
```

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

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

```
['vt', 'wy', 'nd', 'mt', 'ri', 'sd', 'ne', 'de', 'ak', 'nh', 'wv', 'me', 'hi', 'dc', 'nm', 'ks', 'ia
', 'id', 'ar', 'co', 'mn', 'or', 'ky', 'ms', 'nv', 'md', 'ct', 'tn', 'ut', 'al', 'wi', 'va', 'az', '
nj', 'ok', 'wa', 'ma', 'la', 'oh', 'mo', 'in', 'pa', 'mi', 'sc', 'ga', 'il', 'nc', 'fl', 'ny', 'tx',
'ca']
Shape of matrix after one hot encoding  (109248, 51)
```

```
#One Hot Encode – Project Grade Category
my_counter = Counter()
for project_grade in project_data['clean_project_grade_category'].values:
    my_counter.update(project_grade.split())
```

```
project_grade_cat_dict = dict(my_counter)
sorted_project_grade_cat_dict = dict(sorted(project_grade_cat_dict.items(), key=lambda kv: kv[1]))
```

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

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

```
['grades_9_12', 'grades_6_8', 'grades_3_5', 'grades_prek_2']
Shape of matrix after one hot encoding  (109248, 4)
```

```
#one hot encode teacher prefix
my_counter = Counter()
for teacher_prefix in project_data['clean_teacher_prefix'].values:
    teacher_prefix = str(teacher_prefix)
    my_counter.update(teacher_prefix.split())
```

```
teacher_prefix_cat_dict = dict(my_counter)
sorted_teacher_prefix_cat_dict = dict(sorted(teacher_prefix_cat_dict.items(), key=lambda kv: kv[1]))
```

```
## https://stackoverflow.com/questions/39303912/tfidfvectorizer-in-scikit-learn-valueerror-np-nan-is-an-invalid-d
ocument

vectorizer = CountVectorizer(vocabulary=list(sorted_teacher_prefix_cat_dict.keys()), lowercase=False, binary=True
)
vectorizer.fit(project_data['clean_teacher_prefix'].values.astype("U"))
print(vectorizer.get_feature_names())

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

```
['nan', 'dr', 'teacher', 'mr', 'ms', 'mrs']
Shape of matrix after one hot encoding  (109248, 6)
```

# Vectorizing Text data

## Bag of words on project title

In [85]:

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

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

## TFIDF Vectorizer on project_title

In [86]:

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

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

## Using Pretrained Models: Avg W2V on preprocessed title

```
'''
# 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-var
iables-in-python/

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


'''
```

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef loadGloveMod
el(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\'r\', encoding="utf8")\n
model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\n        word = splitLine[0]\
n        embedding = np.array([float(val) for val in splitLine[1:]])\n        model[word] = embeddin
g\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel = loadGloveModel(\'glove
.42B.300d.txt\')\n\n# ============================\nOutput:\n    \nLoading Glove Model\n1917495it [0
6:32, 4879.69it/s]\nDone. 1917495  words loaded!\n\n# ============================\n\nwords = []\nfo
r i in preproced_texts:\n    words.extend(i.split(\' \'))\n\nfor i in preproced_titles:\n    words.e
xtend(i.split(\' \'))\nprint("all the words in the coupus", len(words))\nwords = set(words)\nprint("
the unique words in the coupus", len(words))\n\ninter_words = set(model.keys()).intersection(words)\
nprint("The number of words that are present in both glove vectors and our coupus",       len(inter_
words),"(",np.round(len(inter_words)/len(words)*100,3),"%)")\n\nwords_courpus = {}\nwords_glove = se
t(model.keys())\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\npri
nt("word 2 vec length", len(words_courpus))\n\n\n# stronging variables into pickle files python: htt
p://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport pickle\nwi
th open(\'glove_vectors\', \'wb\') as f:\n    pickle.dump(words_courpus, f)\n\n\n'
```

```python
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-var
iables-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 [89]:

```python
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_title): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_title.append(vector)

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

```
100%|████████████████████████████| 109248/109248 [00:05<00:00, 18603.84it/s]

109248
300
```

## Using Pretrained Models: TFIDF weighted W2V on project_title

In [90]:

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

In [91]:

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

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

```
100%|████████████████████████████| 109248/109248 [00:17<00:00, 6321.42it/s]

109248
300
```

## Vectorizing Numerical features

```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
from sklearn.preprocessing import StandardScaler

# price_standardized = standardScalar.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329.   ... 399.   287.73   5.5 ].
# Reshape your data either using array.reshape(-1, 1)

price_scalar = StandardScaler()
price_scalar.fit(project_data['price'].values.reshape(-1,1)) # finding the mean and standard deviation of this da
ta
print("Mean : {}".format(price_scalar.mean_[0]))
print("Standard deviation : {}".format(np.sqrt(price_scalar.var_[0])))
# Now standardize the data with above maen and variance.
price_standardized = price_scalar.transform(project_data['price'].values.reshape(-1, 1))
```

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

```python
price_standardized
```

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

```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
from sklearn.preprocessing import StandardScaler

# price_standardized = standardScalar.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329.   ... 399.   287.73   5.5 ].
# Reshape your data either using array.reshape(-1, 1)

quantity_scalar = StandardScaler()
quantity_scalar.fit(project_data['quantity'].values.reshape(-1,1)) # finding the mean and standard deviation of t
his data
print(f"Mean : {quantity_scalar.mean_[0]}, Standard deviation : {np.sqrt(quantity_scalar.var_[0])}")

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

```
Mean : 16.965610354422964, Standard deviation : 26.182821919093175
```

```python
quantity_standardized .shape
```

```
(109248, 1)
```

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

# price_standardized = standardScalar.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329.   ... 399.   287.73   5.5 ].
# Reshape your data either using array.reshape(-1, 1)

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

# Now standardize the data with above maen and variance.
teacher_proje_standardized = teacher_proje_scalar.transform(project_data['teacher_number_of_previously_posted_pro
jects'].values.reshape(-1, 1))
```

Mean : 11.153165275336848, Standard deviation : 27.77702641477403

In [97]:

```
teacher_proje_standardized.shape
```

Out[97]:

(109248, 1)

# Merging all the above features

**merging categorical data**

In [98]:

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

Out[98]:

(109248, 100)

**merging numerical data**

In [99]:

```
##https://stackoverflow.com/questions/55756294/scipy-sparse-hstack-valueerror-blocks-must-be-2-d
import scipy as sp
numerical=sp.hstack((quantity_standardized,teacher_proje_standardized,price_standardized))
numerical.shape
```

Out[99]:

(109248, 3)

## 2.1 TSNE with BOW encoding of project_title feature (considering 5000 data points)

In [100]:

```
data=hstack((school_state_categories_one_hot,project_grade_categories_one_hot,teacher_prefix_categories_one_hot,c
ategories_one_hot,sub_categories_one_hot,
          quantity_standardized,teacher_proje_standardized,price_standardized,text_bow_title))
data.shape
```

Out[100]:

(109248, 3432)

```
#https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.csr_matrix.html
x=data.tocsr()
x_new=x[0:5000,:]
x_new.shape
```

Out[104]:

```
(5000, 3433)
```

In [105]:

```
y=project_data['project_is_approved']
y_new=y[0:5000]
y_new.shape
```

Out[105]:

```
(5000,)
```

In [107]:

```python
import seaborn as sn
from sklearn.manifold import TSNE
model = TSNE(n_components=2, perplexity = 30.0,random_state=0)


tsne_data = model.fit_transform(x_new.toarray())


# creating a new data frame which help us in ploting the result data
tsne_data_bow = np.vstack((tsne_data.T, y_new)).T
tsne_df_bow = pd.DataFrame(data=tsne_data_bow, columns=("Dim_1", "Dim_2", "label"))

# Ploting the result of tsne
sn.FacetGrid(tsne_df_bow, hue="label", size=10).map(plt.scatter, 'Dim_1', 'Dim_2').add_legend()
plt.title('TSNE with BOW encoding of project_title feature')
plt.show()
```



TSNE with BOW encoding of project_title feature

```
tsne_df_bow.shape
```

Out[108]:

(5000, 3)

**Summary**

1.We observe alot of overlapping in the datapoints and the points are well scattered, unable to draw any proper conclusion

## 2.2 TSNE with TFIDF encoding of project_title feature (considering 5000 data points)

In [101]:

```
data=hstack((categorical,numerical,text_tfidf_title))
data.shape
```

Out[101]:

(109248, 3432)

In [110]:

```
#https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.csr_matrix.html
x=data.tocsr()
x_new=x[0:5000,:]
```

In [111]:

```
model = TSNE(n_components = 2, perplexity = 30.0, random_state = 0)
tsne_data_tfidf = model.fit_transform(x_new.toarray())

tsne_data_tfidf = np.vstack((tsne_data_tfidf.T, y_new)).T
tsne_df_tfidf= pd.DataFrame(tsne_data_tfidf, columns = ("1st_Dim","2nd_Dim","labels"))

sn.FacetGrid(tsne_df_tfidf, hue="labels", size=10).map(plt.scatter, "1st_Dim","2nd_Dim").add_legend()
plt.title(' TSNE with TFIDF encoding of project_title feature')
plt.show()
```

TSNE with TFIDF encoding of project_title feature

In [112]:

```
tsne_df_tfidf.shape
```

Out[112]:

```
(5000, 3)
```

## 2.3 TSNE with AVG W2V encoding of project_title feature

In [102]:

```
data=hstack((categorical,numerical,avg_w2v_vectors_title))
data.shape
```
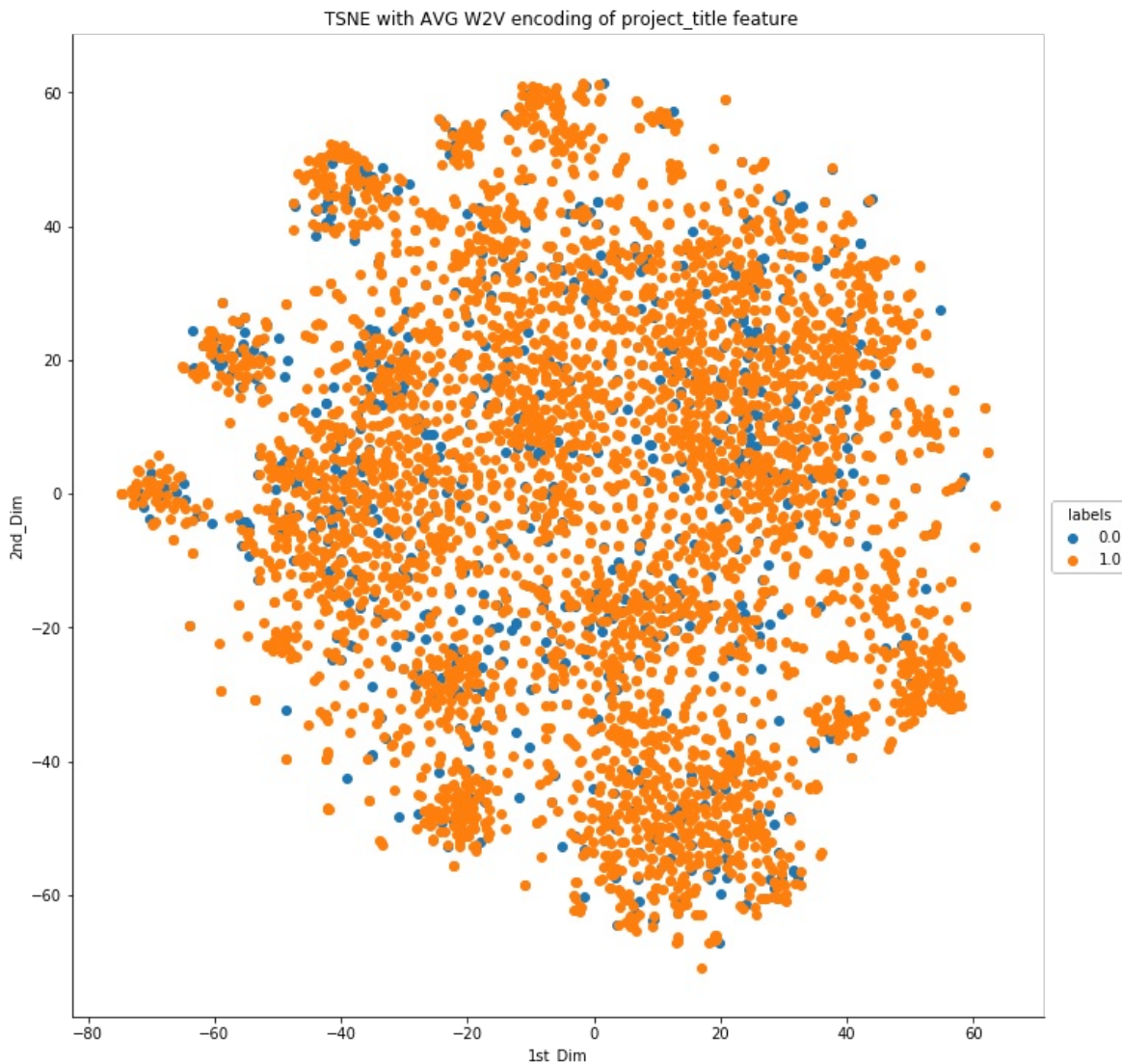
Out[102]:

```
(109248, 403)
```

In [114]:

```
x=data.tocsr()
x_new=x[0:5000,:]
```

```
model = TSNE(n_components = 2, perplexity = 30.0, random_state = 0)
tsne_data_avg_w2v = model.fit_transform(x_new.toarray())

tsne_data_avg_w2v = np.vstack((tsne_data_avg_w2v.T, y_new)).T
tsne_df_avg_w2v= pd.DataFrame(tsne_data_avg_w2v, columns = ("1st_Dim","2nd_Dim","labels"))

sn.FacetGrid(tsne_df_avg_w2v, hue="labels", size=10).map(plt.scatter, "1st_Dim","2nd_Dim").add_legend()
plt.title(' TSNE with AVG W2V encoding of project_title feature')
plt.show()
```



In [116]:

```
tsne_df_avg_w2v.shape
```

Out[116]:

```
(5000, 3)
```

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

In [117]:

```
data=hstack((categorical,numerical,tfidf_w2v_vectors_title))
data.shape
```

Out[117]:

```
(109248, 404)
```

In [118]:

```
x=data.tocsr()
x_new=x[0:5000,:]
```

```
model = TSNE(n_components = 2, perplexity = 30.0, random_state = 0)
tsne_data_wigh_w2v = model.fit_transform(x_new.toarray())

tsne_data_wigh_w2v = np.vstack((tsne_data_wigh_w2v.T, y_new)).T
tsne_df_wigh_w2v= pd.DataFrame(tsne_data_wigh_w2v, columns = ("1st_Dim","2nd_Dim","labels"))

sn.FacetGrid(tsne_df_wigh_w2v, hue="labels", size=10).map(plt.scatter, "1st_Dim","2nd_Dim").add_legend()
plt.title(' TSNE with TFIDF Weighted W2V encoding of project_title feature ')
plt.show()
```


TSNE with TFIDF Weighted W2V encoding of project_title feature

In [120]:

```
tsne_data_wigh_w2v .shape
```

Out[120]:

```
(5000, 3)
```

**Summary**

This visualisation of TSNE with 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

# TSNE with BOW, TFIDF, AVG W2V, TFIDF Weighted W2V encoding of project_title feature (considering 5000 data points)

```
data=hstack((categorical,numerical,text_bow_title,text_tfidf_title,avg_w2v_vectors_title,tfidf_w2v_vectors_title)
)
data.shape
```

Out[121]:

(109248, 7362)

In [122]:

```
x=data.tocsr()
x_new=x[0:5000,:]
x_new.shape
```
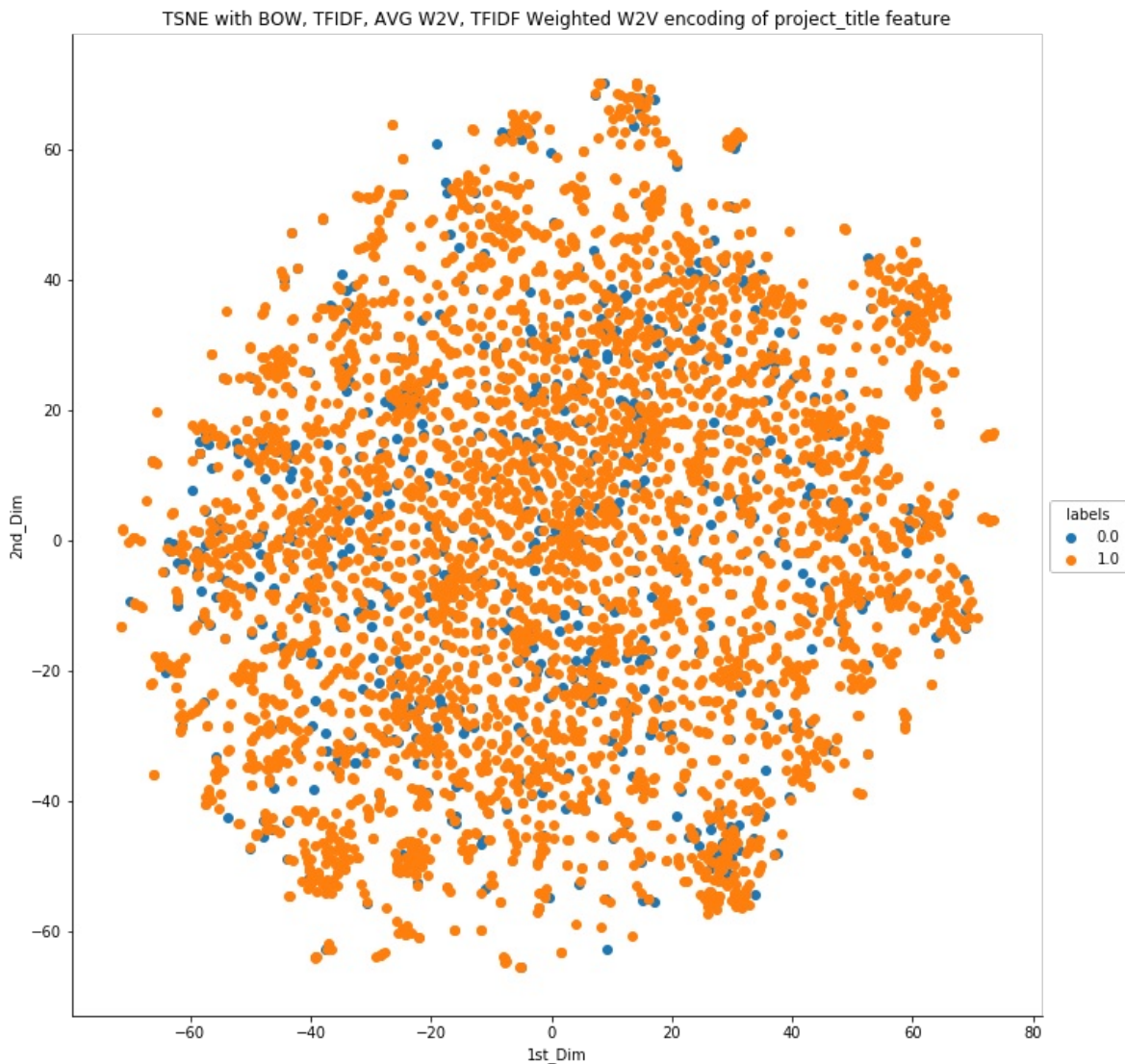
Out[122]:

(5000, 7362)

In [123]:

```
model = TSNE(n_components = 2, perplexity = 30.0, random_state = 0)
tsne_data_complete = model.fit_transform(x_new.toarray())

tsne_data_complete = np.vstack((tsne_data_complete.T, y_new)).T
tsne_df_complete= pd.DataFrame(tsne_data_complete, columns = ("1st_Dim","2nd_Dim","labels"))

sn.FacetGrid(tsne_df_complete, hue="labels", size=10).map(plt.scatter, "1st_Dim","2nd_Dim").add_legend()
plt.title('TSNE with BOW, TFIDF, AVG W2V, TFIDF Weighted W2V encoding of project_title feature')
plt.show()
```



TSNE with BOW, TFIDF, AVG W2V, TFIDF Weighted W2V encoding of project_title feature

In [124]:

```
tsne_df_complete.shape
```

Out[124]:

(5000, 3)

**conclusion**

<>.female Teachers have the maximum number of projects proposed and accepted compared to the male teachers.

<>.There are alot of projects proposed for the students between Pre Kindergarden and 2nd Grade while for the rest it keeps decreasing as the Grades increase.

<>.We also notice that Students between the 9th Grade and 12th Grade have the lowest number of projects proposed as well as accepted.

<>.Projects belonging to the Literacy and Language categories have the highest number of projects proposed under. The maximum number of accepted projects also belong to this category, having an acceptance rate of nearly 87%.

<>.Projects belonging to both Maths and Science have acceptance rate of nearly 82% while introducing the concept of Literacy and Language to this can increase its accpetance rate to nearly 87%

<>.Projects belonging to both Maths and Science when combined with Applied Learning has the least number of projects proposed as well approved.¶

<>.There is also Variability in Acceptance rate, projects under the category Warmth, Care and Hunger have an accpetance rate of 93.5%

<>.The highest number of projects are registered under Literacy and Langauage with 52,239 projects, followed by Maths and Science having 41,421 projects.

<>.The sub-Category Literacy has the highest number of projects approved with 8371 projects. Also the accpetance rate is 88%.

<>. From figure 'TSE with BOW and TFIDF' we can see multiples of small cluster of datapoint but most of the approved project and rejected project datapoint overlapped hence we cannot draw decision line to separate both classes

<>.Even if perplexity value varies TSNE output remains same, we can see that approved project and rejected project datapoint overlapped we can conclude that most of the words used in project title are same for both approved project and rejected project.

In [ ]: