# 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 | Desc |
|---|---|
| project_id | A unique identifier for the proposed project. **Example:** p0 |
| project_title | Title of the project. **Exa**<br>Art Will Make You H<br>First Grad |
| project_grade_category | Grade level of students for which the project is targeted. One of the fo<br>enumerated v<br>Grades P<br>Grade<br>Grade<br>Grades |

| Feature | Desc |
|---|---|
| | One or more (comma-separated) subject categories for the project fr |
| | following enumerated list of v |
| | • Applied Lea |
| | • Care & H |
| | • Health & S |
| | • History & C |
| | • Literacy & Lan |
| project_subject_categories | • Math & Sc |
| | • Music & The |
| | • Special |
| | • W |
| | Exan |
| | • Music & The |
| | • Literacy & Language, Math & Sc |
| school_state | State where school is located ([Two-letter U.S. post](https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_c)  (https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_c Exampl |
| | One or more (comma-separated) subject subcategories for the |
| | Exan |
| project_subject_subcategories | • Lit |
| | • Literature & Writing, Social Sci |
| | An explanation of the resources needed for the project. **Exa** |
| project_resource_summary | • My students need hands on literacy materials to m: |
| | sensory needs!< |
| project_essay_1 | First application |
| project_essay_2 | Second application |
| project_essay_3 | Third application |
| project_essay_4 | Fourth application |
| project_submitted_datetime | Datetime when project application was submitted. **Example:** 2016-( 12:43:5 |
| teacher_id | A unique identifier for the teacher of the proposed project. **Ex:** bdf8baa8fedef6bfeec7ae4ff1c |
| | Teacher's title. One of the following enumerated v |
| | • |
| | • |
| teacher_prefix | • |
| | • |
| | • |
| | Tea |
| teacher_number_of_previously_posted_projects | Number of project applications previously submitted by the same te |
| | Examp |

[*] 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 |

| Feature | Description |
|---|---|
| **quantity** | Quantity of the resource required. **Example:** 3 |
| **price** | Price of the resource required. **Example:** 9.95 |

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

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

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

## Notes on the Essay Data

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

- **project_essay_1:** "Introduce us to your classroom"
- **project_essay_2:** "Tell us more about your students"
- **project_essay_3:** "Describe how your students will use the materials you're requesting"
- **project_essay_4:** "Close by sharing why your project will make a difference"

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

- **project_essay_1:** "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- **project_essay_2:** "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

In [197]:

```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

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

# 1.1 Reading Data

In [198]:

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

In [199]:

```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' 's
chool_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 [200]:

```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[200]:

|   | 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 [201]:

```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-ga


y_value_counts = project_data['project_is_approved'].value_counts()
print("Number of projects thar are approved for funding ", y_value_counts[1], ", (", (y_val
print("Number of projects thar are not approved for funding ", y_value_counts[0], ", (", (y

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 , ( 84.8583040421792
7 %)
Number of projects thar are not approved for funding  16542 , ( 15.141695957
820739 %)

Accepted — Nmber of projects that are Accepted and not accepted



Not Accepted

# OBSERVATION

1. 85% of the projects are approved whereas the rest of the projects are not approved.

## 1.2.1 Univariate Analysis: School State

In [202]:

```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.me
# 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[202]:

```
'# How to plot US state heatmap: https://datascience.stackexchange.com/a/962
0\n\nscl (https://datascience.stackexchange.com/a/9620\n\nscl) = [[0.0, \'rg
b(242,240,247)\'],[0.2, \'rgb(218,218,235)\'],[0.4, \'rgb(188,189,220)\'],
[0.6, \'rgb(158,154,200)\'],[0.8, \'rgb(117,107,177)\'],[1.0, \'rgb(84,39,14
3)\']]\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        locationmo
de = \'USA-states\',\n        text = temp[\'state_code\'],\n        marker =
dict(line = dict (color = \'rgb(255,255,255)\',width = 2)),\n        colorba
r = dict(title = "% of pro")\n    ) ]\n\nlayout = dict(\n        title = \'P
roject 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 [203]:

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

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

# Observation

1. DE in USA is the state with highest number of approval percentage and second highest being the ND with 88% and the WA with third highest with 87% .
2. VT,DC and TX are the lowest approval rates having 80% and 80% and 81% respectively.

In [204]:

```python
#stacked bar plots matplotlib: https://matplotlib.org/gallery/lines_bars_and_markers/bar_st
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 [205]:

```python
def univariate_barplots(data, col1, col2='project_is_approved', top=False):
    # Count number of zeros in dataframe python: https://stackoverflow.com/a/51540521/40840
    temp = pd.DataFrame(project_data.groupby(col1)[col2].agg(lambda x: x.eq(1).sum())).rese

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

    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))
```
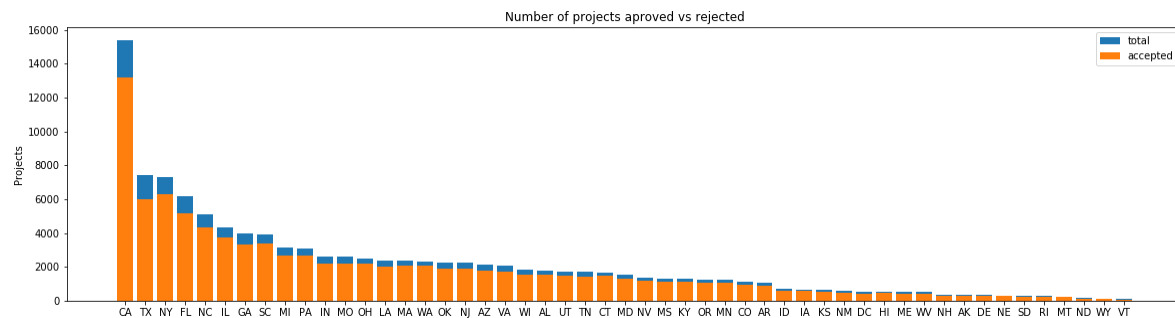
In [206]:

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



Number of projects aproved vs rejected

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

# Summary

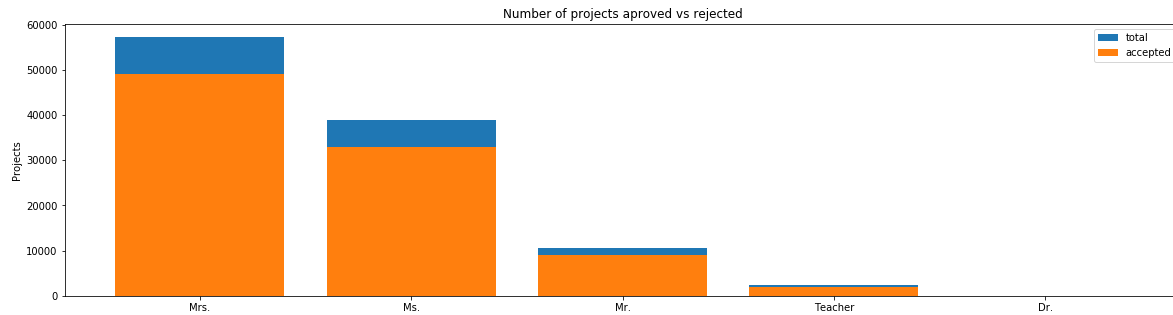1.all states approved percentage is more than 80%.

2.The number of projects submitted varies across the states.

3.CA has the highest approval rate compared to all other states with 85% approval rate having 13205 projects approved out of 15388.

4.VT has the lowest acceptance with lowest proposal but however only 16 projects seems to be rejected which is better than any other rejection counts in other states

## 1.2.2 Univariate Analysis: teacher_prefix

In [207]:

```
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
```
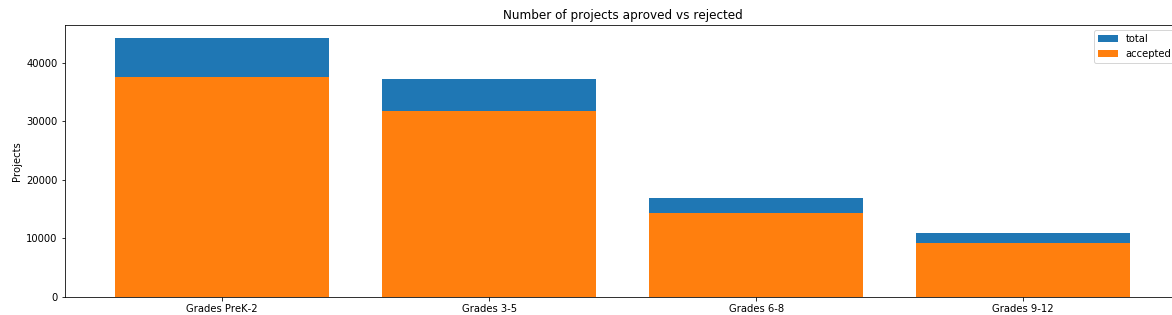
# Observation

1.it seems that the lady teachers who are experienced(just assumption since they are married) have most projects proposed and got accepted.

2.Womens projects has been propsed and approved more than mens projects.

3.Teachers with prefix "Dr" has propsed less(only 13 projects) and mostly(9 projects) has been approved.

## 1.2.3 Univariate Analysis: project_grade_category

In [208]:

```
univariate_barplots(project_data, 'project_grade_category', 'project_is_approved', top=Fals
```



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

# Observation

1.the highest number of projects were propsed and accepted for kids in pre kindergarden grade to second grade which is totally 44225 out of which 37536 projects were approved.

2.Mostly the acceptance percentage is 85% .

3.9th to 12th grade students projects were the lowest number of projects proposed and accepted.

### 1.2.4 Univariate Analysis: project_subject_categories

In [209]:

```
catogories = list(project_data['project_subject_categories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/473019

# 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", "
        if 'The' in j.split(): # this will split each of the catogory based on space "Math
            j=j.replace('The','') # if we have the words "The" we are going to replace it w
        j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math
        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())
```
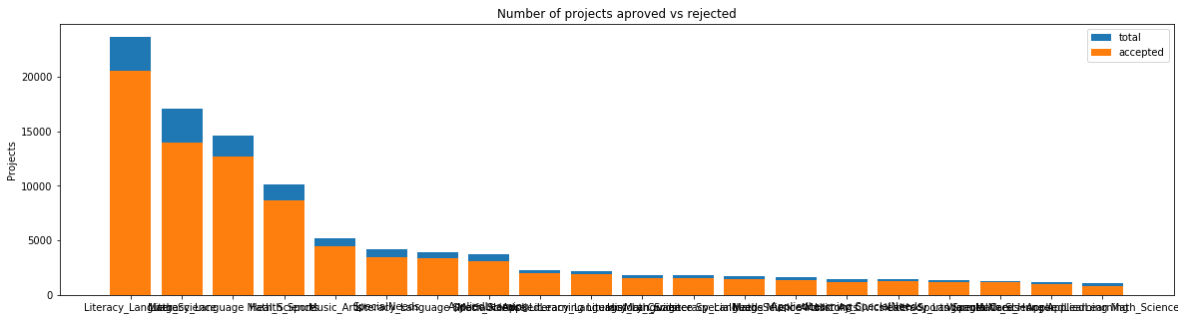
In [210]:

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

Out[210]:

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

In [211]:

```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 |

# observation

1.Literacy_Language category projects are proposed a lot and also accepted a lot with an acceptance percentage nearly 87. 2.Maths and science category projects have 82% acceptance while when literacy_language is combined its acceptance percentange is 86 percentage.
3.Number of projects proposed in each category varies widely.

4.AppliedLearning combined with Math Science have lower number of projects propsed and Lower numer of projects accepted.

5.Warmth Care Hunger category projects are accepted a lot with acceptance percentage of 92%.
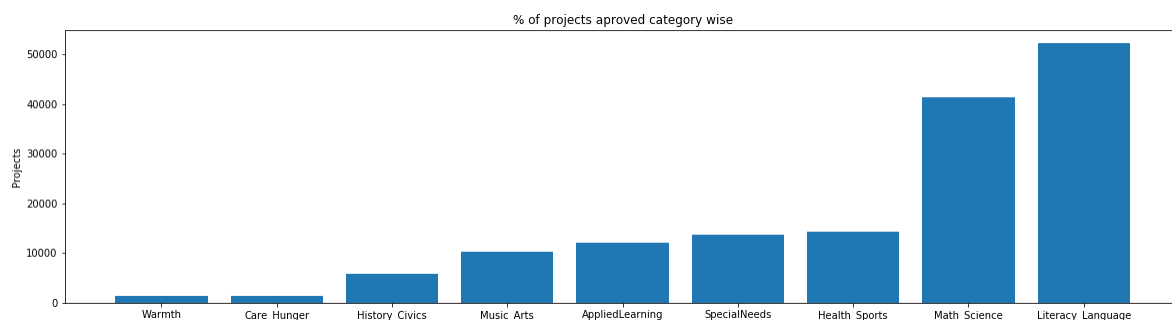
In [212]:

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

In [213]:

```python
# dict sort by value python: https://stackoverflow.com/a/613218/4084039
cat_dict = dict(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 [214]:

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

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

# Observation

1.Most number of projects are proposed in Literacy and Language category with 52239 projects and Second mostly proposed projects are in the category Math and Science with 41421 projects

2.Most accepted projects category Warmth and Care Hunger have only 1388 projects proposed.

## 1.2.5 Univariate Analysis: project_subject_subcategories

In [215]:

```python
sub_catogories = list(project_data['project_subject_subcategories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/473019

# 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", "
        if 'The' in j.split(): # this will split each of the catogory based on space "Math
            j=j.replace('The','') # if we have the words "The" we are going to replace it w
        j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math
        temp +=j.strip()+" "+" # abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_')
    sub_cat_list.append(temp.strip())
```

In [216]:

```python
project_data['clean_subcategories'] = sub_cat_list
project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
project_data.head(2)
```
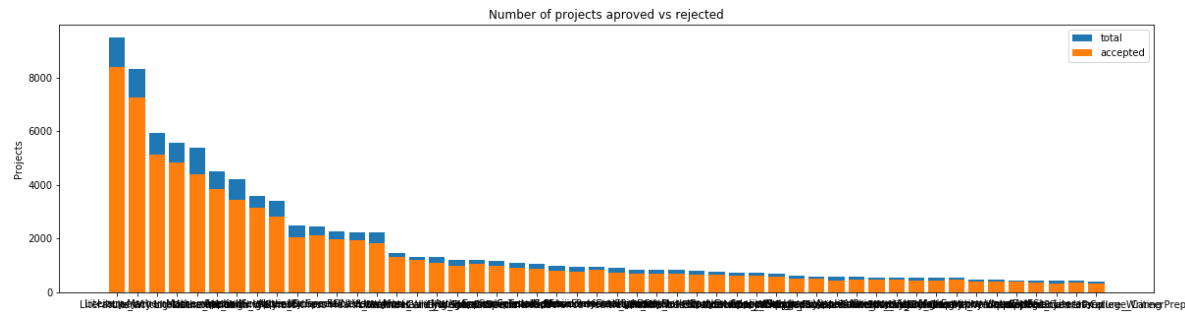
Out[216]:

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

In [217]:

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



Number of projects aproved vs rejected

|  | clean_subcategories | project_is_approved | total | Avg |
|---|---|---|---|---|
| 317 | Literacy | 8371 | 9486 | 0.882458 |
| 319 | Literacy Mathematics | 7260 | 8325 | 0.872072 |
| 331 | Literature_Writing Mathematics | 5140 | 5923 | 0.867803 |
| 318 | Literacy Literature_Writing | 4823 | 5571 | 0.865733 |
| 342 | Mathematics | 4385 | 5379 | 0.815207 |

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

|  | clean_subcategories | project_is_approved | total | Avg |
|---|---|---|---|---|
| 196 | EnvironmentalScience Literacy | 389 | 444 | 0.87612 |
| 127 | ESL | 349 | 421 | 0.82897 |
| 79 | College_CareerPrep | 343 | 421 | 0.81472 |
| 17 | AppliedSciences Literature_Writing | 361 | 420 | 0.85952 |
| 3 | AppliedSciences College_CareerPrep | 330 | 405 | 0.81481 |

# Observation

1.Literacy sub category projects are the most number of projects proposed (9486) out of which 8371 projects are approved with acceptance of 88 percentage. 2.AppliedSciences College and CareerPrep subcategory projects are the least number of projects proposed(405 projects) and accepted(330 projects).

In [218]:

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

In [219]:

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


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

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

In [220]:

```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
```

## 1.2.6 Univariate Analysis: Text features (Title)

In [221]:

```python
#How to calculate number of words in a string in DataFrame: https://stackoverflow.com/a/374
word_count = project_data['project_title'].str.split().apply(len).value_counts()
word_dict = dict(word_count)
word_dict = dict(sorted(word_dict.items(), key=lambda kv: kv[1]))


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

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



## Observation

1.projects have mostly 3 or 4 or 5 words and no project title seems to have more than ten or less than two word.

In [222]:

```python
approved_title_word_count = project_data[project_data['project_is_approved']==1]['project_t
approved_title_word_count = approved_title_word_count.values

rejected_title_word_count = project_data[project_data['project_is_approved']==0]['project_t
rejected_title_word_count = rejected_title_word_count.values
```

In [223]:

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

In [224]:

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

# Observation

1.It seems that a project is accepted if it have more number of words in the title

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

In [225]:

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

In [226]:

```python
approved_word_count = project_data[project_data['project_is_approved']==1]['essay'].str.spl
approved_word_count = approved_word_count.values

rejected_word_count = project_data[project_data['project_is_approved']==0]['essay'].str.spl
rejected_word_count = rejected_word_count.values
```
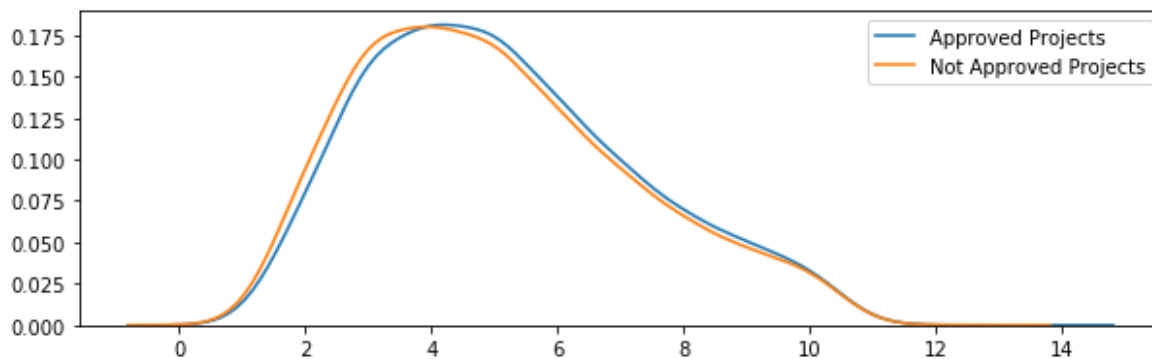
In [227]:

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



# Observation

1.If the number of words in essay increases then the chance for accepting the project may increase. this can be noted after 50.0 percentile in approved projects.

In [228]:

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



# observation

1.The projects which have number of words in essay more than 250 words have higher chance of accepting
.This can be observed by the increasing density of approved projects curve from 250.

### 1.2.8 Univariate Analysis: Cost per project

In [229]:

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

Out[229]:

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

In [230]:

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

Out[230]:

|   | 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 [231]:

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

In [232]:

```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 [233]:

```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()
```
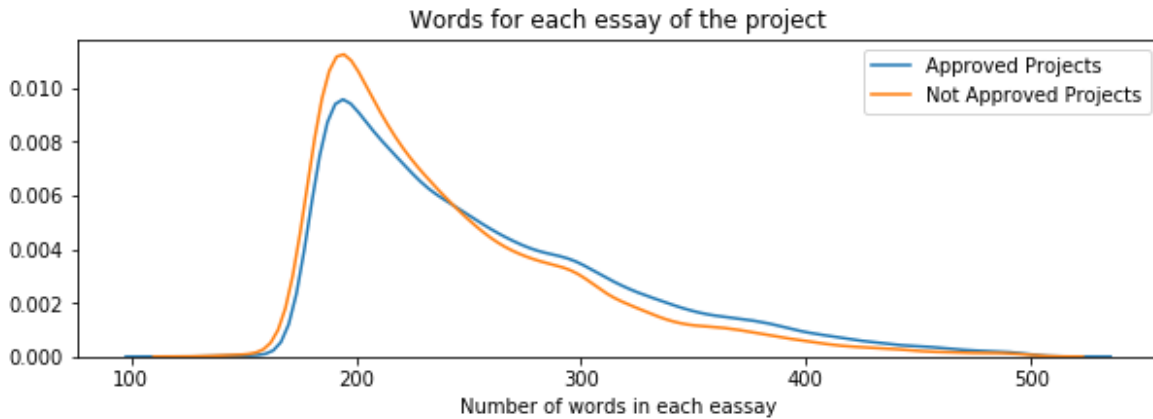


In [234]:

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

In [235]:

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

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

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(reje
print(x)
```

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

# Observation
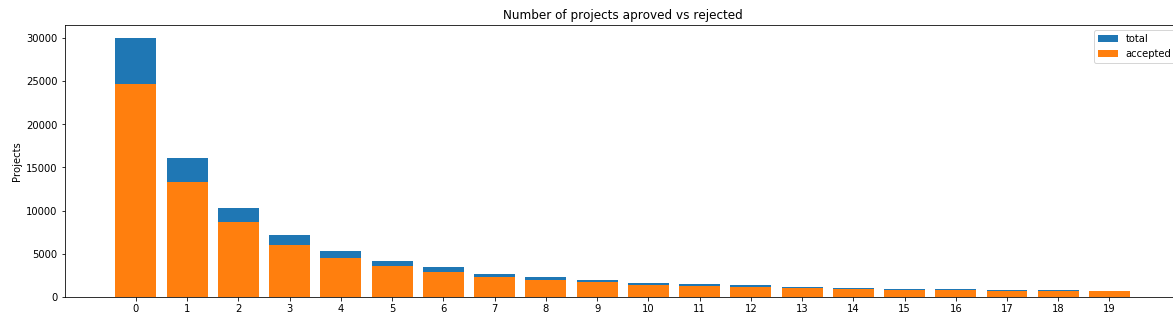
1.The Less cost projects are approved more than that of projects with high cost .

2.This can be seen from the the table in 50 th percentile having a project approved with a cost of 199 dollar whereas the rejected project having a cost of 263 which is more than that of the approved project.

## 1.2.9 Univariate Analysis: teacher_number_of_previously_posted_projects

In [236]:

```
univariate_barplots(project_data, 'teacher_number_of_previously_posted_projects','project_i
```



Number of projects aproved vs rejected

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

|   | Avg |
|---|---|
| 0 | 0.821350 |
| 1 | 0.830054 |
| 2 | 0.841063 |
| 3 | 0.843460 |
| 4 | 0.845423 |

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

|   | teacher_number_of_previously_posted_projects | project_is_approved | total |
|---|---|---|---|
| \ | | | |
| 15 | 15 | 818 | 942 |
| 16 | 16 | 769 | 894 |
| 17 | 17 | 712 | 803 |
| 18 | 18 | 666 | 772 |
| 19 | 19 | 632 | 710 |

|   | Avg |
|---|---|
| 15 | 0.868365 |
| 16 | 0.860179 |
| 17 | 0.886675 |
| 18 | 0.862694 |
| 19 | 0.890141 |

# Observation

1. It is observed that if a teacher have proposed more number of projects previously than the rate of approval also seems to be increasing .Even the teachers with no projects proposed also have a good rate of approval(82%).
2. But if the number of previously submitted projects increases then approval rate is maximum(89%).

### 1.2.10 Univariate Analysis: project_resource_summary

In [237]:

```python
summarys = []
for a in project_data["project_resource_summary"] :
    summarys.append(a)
summarys[0:10]
```

Out[237]:

```
['My students need opportunities to practice beginning reading skills in Eng
lish at home.',
 'My students need a projector to help with viewing educational programs',
 'My students need shine guards, athletic socks, Soccer Balls, goalie glove
s, 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 personal
ized journals and charts will help them be more involved in our daily Math r
outines.',
 '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 assi
st not only these students with special needs, but all students.',
 'My students need some dependable laptops for daily classroom use for readi
ng and math.',
 'My students need ipads to help them access a world of online resources tha
t will spark their interest in learning.',
 "My students need three devices and three management licenses for small gro
up's easy access to newly-implemented online programs--Go Noodle Plus, for i
ncreased in-class physical activity and Light Sail, an interactive reading p
rogram.",
 'My students need great books to use during Independent Reading, Read Aloud
s, Partner Reading and Author Studies.']
```

In [238]:

```python
NumericSummaryValues = {}
for x in range(len(summarys)):
    for s in summarys[x].split():
        if s.isdigit() :
            NumericSummaryValues[x] = int(s)
```

In [239]:

```python
NumericDigits = {}
for c in range(len(summarys)) :
    if c in NumericSummaryValues.keys() :
        NumericDigits[c] = NumericSummaryValues[c]
    else :
        NumericDigits[c] = 0
```

In [240]:

```python
DigitInSummary = []

for a in NumericDigits.values() :
    if a > 0 :
        DigitInSummary.append(1)
    else :
        DigitInSummary.append(0)
```

In [241]:

```python
DigitInSummary[0:10]
```

Out[241]:

```
[0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
```

In [242]:

```
project_data['digit_in_summary'] = DigitInSummary
project_data.head(10)
```

Out[242]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | projec |
|---|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | 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 | |

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | projec |
|---|---|---|---|---|---|---|
| **9** | 158561 | p001713 | 140eeac1885c820ad5592a409a3a8994 | Ms. | NC | |

10 rows × 21 columns

In [243]:

```
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
```

# Observation

1.Most of the approved projects don't have digits but the projects with digits have great acceptance rate of more than 90%(approximately).

## 1.3 Text preprocessing

### 1.3.1 Essay Text

In [244]:

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

Out[244]:

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

2 rows × 21 columns

In [245]:

```python
# 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 hav
e 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 a
nd experiences to us that open our eyes to new cultures, beliefs, and respec
t.\"The limits of your language are the limits of your world.\"-Ludwig Wittg
enstein  Our English learner's have a strong support system at home that beg
s for more resources.  Many times our parents are learning to read and speak
English along side of their children.  Sometimes this creates barriers for p
arents to be able to help their child learn phonetics, letter recognition, a
nd other reading skills.\r\n\r\nBy providing these dvd's and players, studen
ts 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 pr
oficiency status, will be a offered to be a part of this program.  These edu
cational 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 deve
lop early reading skills.\r\n\r\nParents that do not have access to a dvd pl
ayer will have the opportunity to check out a dvd player to use for the yea
r.  The plan is to use these videos and educational dvd's for the years to c
ome for other EL students.\r\nnannan
==================================================
The 51 fifth grade students that will cycle through my classroom this year a
ll love learning, at least most of the time. At our school, 97.3% of the stu
dents receive free or reduced price lunch. Of the 560 students, 97.3% are mi
nority students. \r\nThe school has a vibrant community that loves to get to
gether and celebrate. Around Halloween there is a whole school parade to sho
w 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 e
nd 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 regul
ar, stationary, 4-legged chairs. As I will only have a total of ten in the c
lassroom and not enough for each student to have an individual one, they wil
l 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 th
em in place of chairs at my small group tables during math and reading time
s. The rest of the day they will be used by the students who need the highes
t amount of movement in their life in order to stay focused on school.\r\n\r
\nWhenever asked what the classroom is missing, my students always say more
Hokki Stools. They can't get their fill of the 5 stools we already have. Whe
n the students are sitting in group with me on the Hokki Stools, they are al
ways 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. Ther
e are always students who head over to the kidney table to get one of the st

ools who are disappointed as there are not enough of them. \r\n\r\nWe ask a lot of students to sit for 7 hours a day. The Hokki stools will be a compromise that allow my students to do desk work and move at the same time. These stools will help students to meet their 60 minutes a day of movement by allowing them to activate their core muscles for balance while they sit. For many of my students, these chairs will take away the barrier that exists in schools for a child who can't sit still.nannan

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

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

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

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

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

The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The great teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is makeup is 97.6% African-American, making up the largest segment of the student body. A typical school in Dallas is made up of 23.2% African-American students. Most of the students are on free or reduced lunch. We aren't receiving doctors, lawyers, or engineers children from rich backgrounds or neighborhoods. As an educator I am inspiring minds of young children and we focus not only on academics but one smart, effective, efficie

nt, and disciplined students with good character.In our classroom we can uti
lize 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 o
f my speaker my students can't hear videos or books clearly and it isn't mak
ing the lessons as meaningful. But with the bluetooth speaker my students wi
ll be able to hear and I can stop, pause and replay it at any time.\r\nThe c
art will allow me to have more room for storage of things that are needed fo
r the day and has an extra part to it I can use.  The table top chart has al
l of the letter, words and pictures for students to learn about different le
tters and it is more accessible.nannan
==================================================

In [246]:

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

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

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

In [247]:

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

My kindergarten students have varied disabilities ranging from speech and la
nguage delays, cognitive delays, gross/fine motor delays, to autism. They ar
e eager beavers and always strive to work their hardest working past their l
imitations. \r\n\r\nThe materials we have are the ones I seek out for my stu
dents. I teach in a Title I school where most of the students receive free o
r reduced price lunch.  Despite their disabilities and limitations, my stude
nts 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 yo
u were in a meeting? This is how my kids feel all the time. The want to be a
ble to move as they learn or so they say.Wobble chairs are the answer and I
love then because they develop their core, which enhances gross motor and in
Turn fine motor skills. \r\nThey also want to learn through games, my kids d
o 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 [248]:

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

My kindergarten students have varied disabilities ranging from speech and la
nguage delays, cognitive delays, gross/fine motor delays, to autism. They ar
e eager beavers and always strive to work their hardest working past their l
imitations.       The materials we have are the ones I seek out for my student
s. I teach in a Title I school where most of the students receive free or re
duced 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 wer
e in a meeting? This is how my kids feel all the time. The want to be able t
o move as they learn or so they say.Wobble chairs are the answer and I love
then because they develop their core, which enhances gross motor and in Turn
fine motor skills.   They also want to learn through games, my kids do not w
ant to sit and do worksheets. They want to learn to count by jumping and pla
ying. Physical engagement is the key to our success. The number toss and col
or and shape mats can make that happen. My students will forget they are doi
ng work and just have the fun a 6 year old deserves.nannan

In [249]:

```
#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 la
nguage delays cognitive delays gross fine motor delays to autism They are ea
ger beavers and always strive to work their hardest working past their limit
ations 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 a
nts 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 le
arn or so they say Wobble chairs are the answer and I love then because they
develop their core which enhances gross motor and in Turn fine motor skills
They also want to learn through games my kids do not want to sit and do work
sheets They want to learn to count by jumping and playing Physical engagemen
t is the key to our success The number toss and color and shape mats can mak
e that happen My students will forget they are doing work and just have the
fun a 6 year old deserves nannan

In [250]:

```python
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they'
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'l
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had',
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'u
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'd
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over',
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any',
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'v
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now',
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'dc
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn',
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn'
            'won', "won't", 'wouldn', "wouldn't"]
```

In [251]:

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

```
100%|████████████████████████████████████████████████████████████████
██| 109248/109248 [02:06<00:00, 862.24it/s]
```

In [252]:

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

Out[252]:

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

## 1.3.2 Project title Text

In [253]:

```python
#printing random texts
print(project_data['project_title'].values[0])
print("="*50)
print(project_data['project_title'].values[100])
print("="*50)
print(project_data['project_title'].values[1000])
print("="*50)
print(project_data['project_title'].values[10000])
print("="*50)
print(project_data['project_title'].values[99999])
print("="*50)
```

```
Educational Support for English Learners at Home
==================================================
21st Century learners, 21st century technology!
==================================================
Sailing Into a Super 4th Grade Year
==================================================
Family Book Clubs
==================================================
Inspiring Minds by Enhancing the Educational Experience
==================================================
```

In [254]:

```python
PreprocessedTitles = []
for titles in tqdm(project_data["project_title"]):
    title = decontracted(titles)
    title = title.replace('\\r', ' ')
    title = title.replace('\\"', ' ')
    title = title.replace('\\n', ' ')
    title = re.sub('[^A-Za-z0-9]+', ' ', title)
    title = ' '.join(f for f in title.split() if f not in stopwords)
    PreprocessedTitles.append(title.lower().strip())
```

```
100%|████████████████████████████████████████████|
| 109248/109248 [00:05<00:00, 20125.78it/s]
```

In [255]:

```python
print(PreprocessedTitles[0])
print("="*50)
print(PreprocessedTitles[50])
print("="*50)
print(PreprocessedTitles[500])
print("="*50)
print(PreprocessedTitles[5000])
print("="*50)
print(PreprocessedTitles[10000])
print("="*50)
```

```
educational support english learners home
==================================================
be active be energized
==================================================
classroom chromebooks college bound seniors
==================================================
bouncing our wiggles worries away
==================================================
family book clubs
==================================================
```

# 1. 4 Preparing data for models

In [256]:

```python
project_data.columns
```

Out[256]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
       'project_submitted_datetime', 'project_grade_category', 'project_titl
e',
       'project_essay_1', 'project_essay_2', 'project_essay_3',
       'project_essay_4', 'project_resource_summary',
       'teacher_number_of_previously_posted_projects', 'project_is_approve
d',
       'clean_categories', 'clean_subcategories', 'essay', 'price', 'quantit
y',
       'digit_in_summary'],
      dtype='object')
```

we are going to consider

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

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

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

## 1.4.1 Vectorizing Categorical data

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

In [257]:

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


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

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

In [258]:

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


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

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

In [259]:

```python
# encodig with state
counter = Counter()
for state in project_data['school_state'].values:
    counter.update(state.split())

school_state_cat_dict = dict(counter)
sorted_school_state_cat_dict = dict(sorted(school_state_cat_dict.items(), key=lambda kv: kv


vectorizer = CountVectorizer(vocabulary=list(sorted_school_state_cat_dict.keys()), lowercas
vectorizer.fit(project_data['school_state'].values)
print(vectorizer.get_feature_names())
school_state_categories_one_hot = vectorizer.transform(project_data['school_state'].values)
print("Shape of matrix after one hot encodig ",school_state_categories_one_hot.shape)
```

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

In [260]:

```python
#encodig with project grade category
counter = Counter()
for project_grade in project_data['project_grade_category'].values:
    counter.update(project_grade.split())

project_grade_cat_dict = dict(counter)
sorted_project_grade_cat_dict = dict(sorted(project_grade_cat_dict.items(), key=lambda kv:
vectorizer = CountVectorizer(vocabulary=list(sorted_project_grade_cat_dict.keys()), lowerca
vectorizer.fit(project_data['project_grade_category'].values)
print(vectorizer.get_feature_names())
project_grade_categories_one_hot = vectorizer.transform(project_data['project_grade_categor
print("Shape of matrix after one hot encodig ",project_grade_categories_one_hot.shape)
```

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

In [261]:

```python
#encodig with Teacher Prefix
#Refered below link
# https://stackoverflow.com/questions/39303912/tfidfvectorizer-in-scikit-learn-valueerror-n

counter = Counter()
for teacher_prefix in project_data['teacher_prefix'].values:
    teacher_prefix = str(teacher_prefix)
    counter.update(teacher_prefix.split())

teacher_prefix_cat_dict = dict(counter)
sorted_teacher_prefix_cat_dict = dict(sorted(teacher_prefix_cat_dict.items(), key=lambda kv

vectorizer = CountVectorizer(vocabulary=list(sorted_teacher_prefix_cat_dict.keys()), lowerc
vectorizer.fit(project_data['teacher_prefix'].values.astype("U"))
print(vectorizer.get_feature_names())
teacher_prefix_categories_one_hot =vectorizer.transform(project_data['teacher_prefix'].valu
print("Shape of matrix after one hot encodig ",teacher_prefix_categories_one_hot.shape)
```

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

## 1.4.2 Vectorizing Text data

### 1.4.2.1 Bag of words

In [262]:

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

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

### 1.4.2.2 Bag of Words on `project_title`

In [263]:

```
# you can vectorize the title also
vectorizer = CountVectorizer(min_df=10)
title_bow = vectorizer.fit_transform(PreprocessedTitles)
print("Shape of matrix after one hot encodig ",title_bow.shape)
```

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

### 1.4.2.3 TFIDF vectorizer

In [264]:

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

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

### 1.4.2.4 TFIDF Vectorizer on `project_title`

In [265]:

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

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

### 1.4.2.5 Using Pretrained Models: Avg W2V

In [266]:

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

In [267]:

```python
model = loadGloveModel('glove.42B.300d.txt')
```

Loading Glove Model

1917495it [24:05, 1326.65it/s]

Done. 1917495  words loaded!

In [268]:

```python
words = []
for i in PreprocessedEssays:
    words.extend(i.split(' '))

for i in PreprocessedTitles:
    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))
```

all the words in the coupus 17014413
the unique words in the coupus 58968

In [269]:

```python
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))
```

The number of words that are present in both glove vectors and our coupus 51
503 ( 87.341 %)
word 2 vec length 51503

In [270]:

```python
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickl
import pickle
with open('glove_vectors', 'wb') as f:
    pickle.dump(words_courpus, f)
```

In [271]:

```python
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickl
# 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 [272]:

```python
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(PreprocessedEssays): # 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 [07:58<00:00, 228.24it/s]

109248
300
```

```html
<h4><font color='red'> 1.4.2.6 Using Pretrained Models: AVG W2V on `project_title`</font>
</h4>
```

In [273]:

```python
avg_w2v_vectors_titles = [];
for sentence in tqdm(PreprocessedTitles):
    vector = np.zeros(300)
    count_words =0;
    for word in sentence.split():
        if word in glove_words:
            vector += model[word]
            count_words += 1
    if count_words != 0:
        vector /= count_words
    avg_w2v_vectors_titles.append(vector)
print(len(avg_w2v_vectors_titles))
print(len(avg_w2v_vectors_titles[0]))
```

```
100%|████████████████████████████████████████████████████████████████████
█| 109248/109248 [00:26<00:00, 4070.31it/s]

109248
300
```

#### 1.4.2.7 Using Pretrained Models: TFIDF weighted W2V

In [274]:

```python
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(PreprocessedEssays)
# 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 [275]:

```python
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(PreprocessedEssays): # 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((sentenc
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # gettir
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors.append(vector)

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

```
100%|████████████████████████████████████████████████████████
██| 109248/109248 [05:17<00:00, 359.71it/s]

109248
300
```

```html
<h4><font color='red'> 1.4.2.9 Using Pretrained Models: TFIDF weighted W2V on
`project_title`</font></h4>
```

In [276]:

```python
# Similarly you can vectorize for title also
tfidf_model = TfidfVectorizer()
tfidf_model.fit(PreprocessedTitles)
# 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 [277]:

```python
tfidf_w2v_vectors_title = []; # the avg-w2v for every sentences is stored in this list
for sentence in tqdm(PreprocessedTitles):
    vector = np.zeros(300)
    tf_idf_weight =0;
    for word in sentence.split():
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word]
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
            vector += (vec * tf_idf)
            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 [01:21<00:00, 1335.34it/s]

109248
300
```

### 1.4.3 Vectorizing Numerical features

In [278]:

```python
# check this one: https://www.youtube.com/watch?v=0HOqOcln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.prepro
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.
# 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 standar
print(f"Mean : {price_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])}

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

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

In [279]:

```python
price_standardized
```

Out[279]:

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

# Observation
1. The projects costs only within 700 dollars since mean cost is 298 dollars and standard deviation is 367 dollars

### Vectorizing Quantity

In [280]:

```python
import warnings
warnings.filterwarnings("ignore")
quantity_scalar = StandardScaler()
quantity_scalar.fit(project_data['quantity'].values.reshape(-1,1))
print("Mean : {}".format(quantity_scalar.mean_[0]))
print("Standard deviation : {}".format(np.sqrt(quantity_scalar.var_[0])))
quantity_standardized = quantity_scalar.transform(project_data['quantity'].values.reshape(-
```

Mean : 16.965610354422964
Standard deviation : 26.182821919093175

In [281]:

```python
quantity_standardized
```

Out[281]:

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

# observation
The project will require an mean of 17 items and the standard deviation is 26.Donors will donate money for any of those projects depending on the items.

### Vectorizing  previous posted projects

In [282]:

```python
prev_projects_scalar = StandardScaler()
prev_projects_scalar.fit(project_data['teacher_number_of_previously_posted_projects'].value
print("Mean : {}".format(prev_projects_scalar.mean_[0]))
print("Standard deviation : {}".format(np.sqrt(prev_projects_scalar.var_[0])))
prev_projects_standardized =prev_projects_scalar.transform(project_data['teacher_number_of_
```

Mean : 11.153165275336848
Standard deviation : 27.77702641477403

In [283]:

```
prev_projects_standardized
```

Out[283]:

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

# Observation

1.Teacher have proposed atleast an average of 11 different projects.Seems to be more active in helping their students.

### 1.4.4 Merging all the above features

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

In [284]:

```
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 [285]:

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

Out[285]:

```
(109248, 16663)
```

# Assignment 2: Apply TSNE

<font color=#F4274F>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.</font>

```
<ol>
    <li> 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.</li>
```

```
    <li> EDA: Please complete the analysis of the feature:
teacher_number_of_previously_posted_projects</li>
    <li>
        <ul>Build the data matrix using these features
            <li>school_state : categorical data (one hot encoding)</li>
            <li>clean_categories : categorical data (one hot encoding)</li>
            <li>clean_subcategories : categorical data (one hot encoding)</li>
            <li>teacher_prefix : categorical data (one hot encoding)</li>
            <li>project_grade_category : categorical data (one hot encoding)</li>
            <li>project_title : text data (BOW, TFIDF, AVG W2V, TFIDF W2V)</li>
            <li>price : numerical</li>
            <li>teacher_number_of_previously_posted_projects : numerical</li>
        </ul>
    </li>
    <li> Now, plot FOUR t-SNE plots with each of these feature sets.
        <ol>
            <li>categorical, numerical features + project_title(BOW)</li>
            <li>categorical, numerical features + project_title(TFIDF)</li>
            <li>categorical, numerical features + project_title(AVG W2V)</li>
            <li>categorical, numerical features + project_title(TFIDF W2V)</li>
        </ol>
    </li>
    <li> Concatenate all the features and Apply TNSE on the final data matrix </li>
    <li> <font color='blue'>Note 1: The TSNE accepts only dense matrices</font></li>
    <li> <font color='blue'>Note 2: Consider only 5k to 6k data points to avoid memory
issues. If you run into memory error issues, reduce the number of data points but clearly
state the number of datat-poins you are using</font></li>
</ol>
```
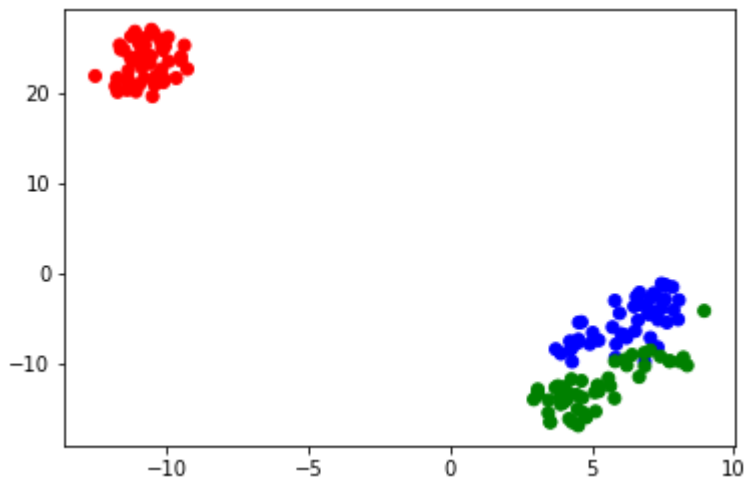
In [286]:

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

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

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

X_embedding = tsne.fit_transform(x)
# if x is a sparse matrix you need to pass it as X_embedding = tsne.fit_transform(x.toarray

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'].
plt.show()
```



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

In [287]:

```
# please write all of the code with proper documentation and proper titles for each subsect
# 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

print("BOW : {}".format(title_bow.shape))
print("TFIDF : {}".format(title_tfidf.shape))
print("AVG W2V : ({}, {})".format(len(avg_w2v_vectors_titles), len(avg_w2v_vectors_titles[0
print("TFIDF W2V : ({}, {})".format(len(tfidf_w2v_vectors_title),
len(tfidf_w2v_vectors_title[0])))
```

```
BOW : (109248, 3329)
TFIDF : (109248, 3329)
AVG W2V : (109248, 300)
TFIDF W2V : (109248, 300)
```

In [288]:

```
x = hstack((categories_one_hot, sub_categories_one_hot, school_state_categories_one_hot,pro
x.shape
```

Out[288]:

```
(109248, 3433)
```

In [289]:

```
from sklearn.manifold import TSNE
TSNE_model = TSNE(n_components = 2, perplexity=30, learning_rate=200) #creating model
x = x.tocsr()  #https://stackoverflow.com/questions/30163830/accessing-elements-in-coo-matr
x_5k = x[0:5000,:]
```

In [290]:

```
tsne_data = TSNE_model.fit_transform(x_5k.toarray())
#Don't touch
```

In [291]:

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

Out[291]:

```
5000
```

In [292]:

```
# https://github.com/Tejas163/Data-Science/blob/master/Project-1-Amazon%20Fine%20Food%20Rev

tsne_data = np.vstack((tsne_data.T, labels_new)).T
tsne_df = pd.DataFrame(tsne_data, columns = ("1stDim","2ndDim","Labels"))
```

In [293]:

```
tsne_df.shape
```

Out[293]:

```
(5000, 3)
```

In [294]:

```
# https://github.com/Tejas163/Data-Science/blob/master/Project-1-Amazon%20Fine%20Food%20Rev

sns.FacetGrid(tsne_df, hue = "Labels", size = 10).map(plt.scatter, "1stDim", "2ndDim").add_
plt.title("TSNE WITH BOW ENCODING OF PROJECT TITLE FEATURE ")
plt.show()
```



TSNE WITH BOW ENCODING OF PROJECT TITLE FEATURE

# Observation

1. There can't be observed anything since the data points are scattered

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

In [295]:

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

x = hstack((categories_one_hot, sub_categories_one_hot, school_state_categories_one_hot,prc
x.shape
```

Out[295]:

```
(109248, 3433)
```

In [296]:

```
TSNE_model = TSNE(n_components = 2, perplexity=30, learning_rate=200)
x = x.tocsr()
x_5k = x[0:5000,:]
```

In [297]:

```
tsne_data_tfidf = TSNE_model.fit_transform(x_5k.toarray())
```
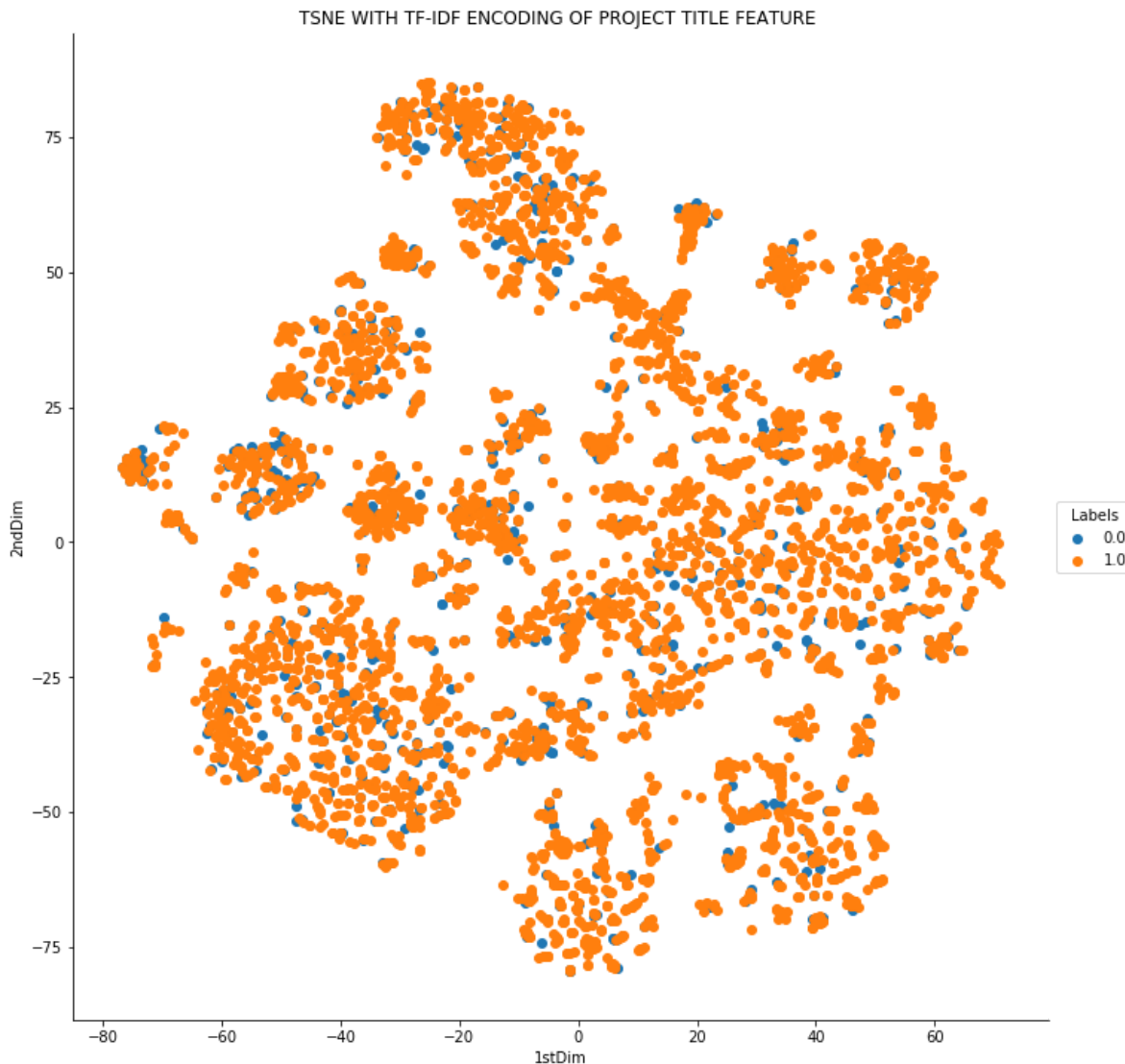
In [298]:

```
tsne_data_tfidf = np.vstack((tsne_data_tfidf.T, labels_new)).T
tsne_df_tfidf = pd.DataFrame(tsne_data_tfidf, columns = ("1stDim","2ndDim","Labels"))
tsne_df_tfidf.shape
```

Out[298]:

```
(5000, 3)
```

In [299]:

```python
sns.FacetGrid(tsne_df_tfidf, hue = "Labels", size = 10).map(plt.scatter, "1stDim", "2ndDim"
plt.title("TSNE WITH TF-IDF ENCODING OF PROJECT TITLE FEATURE ")
plt.show()
```



TSNE WITH TF-IDF ENCODING OF PROJECT TITLE FEATURE

# Observation

1. Nothing can be observed since the datapoints are still together.

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

In [300]:

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

x = hstack((categories_one_hot, sub_categories_one_hot, school_state_categories_one_hot,pro
quantity_standardized, prev_projects_standardized, avg_w2v_vectors_titles))
x.shape
```

Out[300]:

```
(109248, 404)
```

In [301]:

```python
TSNE_model = TSNE(n_components = 2, perplexity=30, learning_rate=200)
x = x.tocsr()
x_5k = x[0:5000,:]
```

In [302]:

```python
tsne_data_avg_w2v = TSNE_model.fit_transform(x_5k.toarray())
```

In [303]:

```python
tsne_data_avg_w2v = np.vstack((tsne_data_avg_w2v.T, labels_new)).T
tsne_df_avg_w2v = pd.DataFrame(tsne_data_avg_w2v, columns = ("1stDim","2ndDim","Labels"))
```

In [304]:

```python
tsne_df_avg_w2v.shape
```

Out[304]:

```
(5000, 3)
```

In [305]:

```
sns.FacetGrid(tsne_df_avg_w2v, hue = "Labels", size = 10).map(plt.scatter, "1stDim", "2ndDi
plt.title("TSNE WITH AVG W2V ENCODING OF PROJECT TITLE FEATURE ")
plt.show()
```



TSNE WITH AVG W2V ENCODING OF PROJECT TITLE FEATURE

# Observation

Since there is no clusters can be observed,drawing a conclusion in this case seems impossible

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

In [306]:

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

Out[306]:

(109248, 404)

In [307]:

```python
TSNE_model = TSNE(n_components = 2, perplexity=30, learning_rate=200)
x = x.tocsr()
x_5k = x[0:5000,:]
```
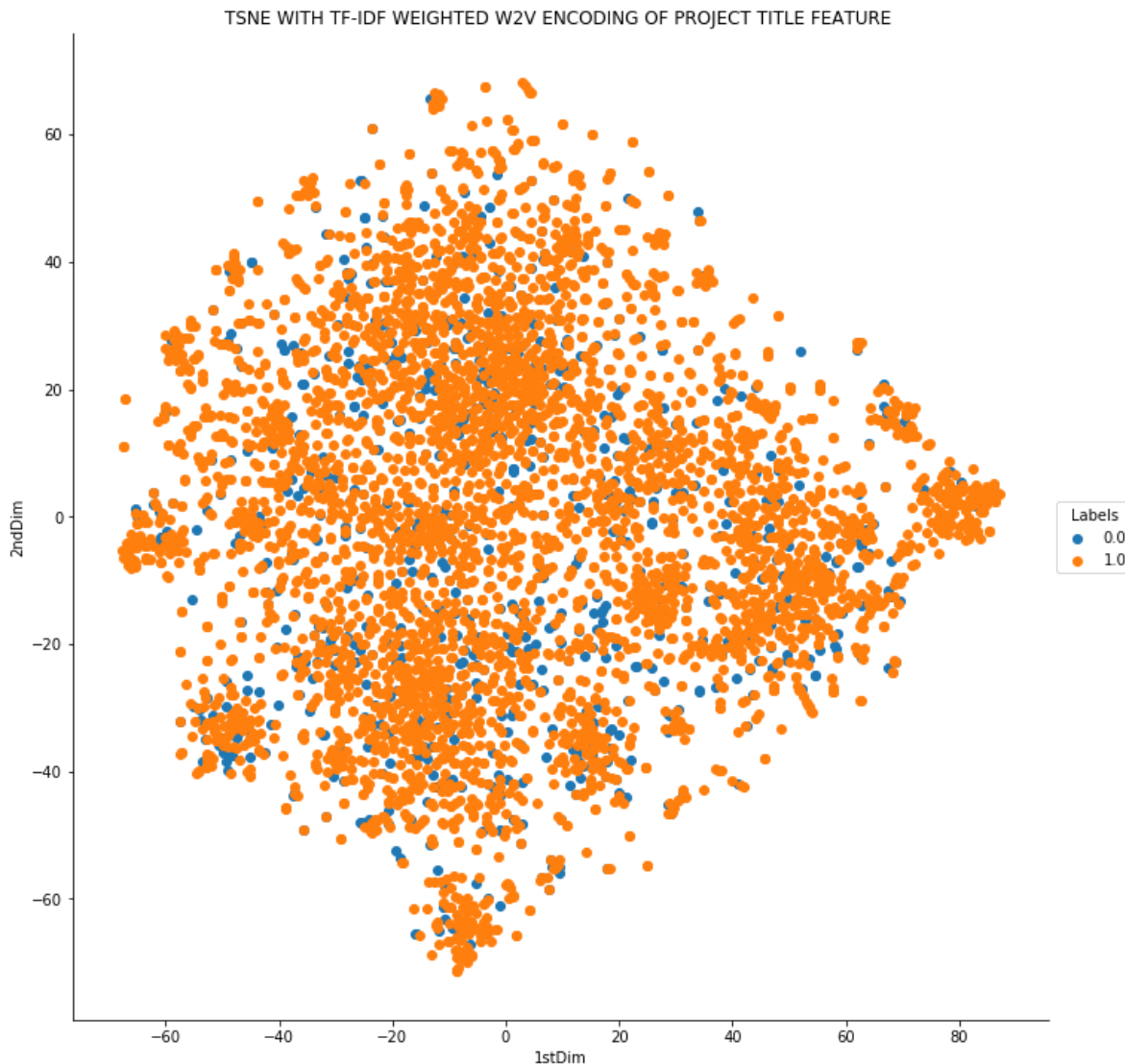
In [308]:

```python
tsne_data_tfidf_weigh_w2v = TSNE_model.fit_transform(x_5k.toarray())
```

In [309]:

```python
tsne_data_tfidf_weigh_w2v = np.vstack((tsne_data_tfidf_weigh_w2v.T, labels_new)).T
tsne_df_tfidf_weigh_w2v = pd.DataFrame(tsne_data_tfidf_weigh_w2v, columns = ("1stDim","2ndD
```

In [310]:

```
sns.FacetGrid(tsne_df_tfidf_weigh_w2v, hue = "Labels", size = 10).map(plt.scatter, "1stDim"
plt.title("TSNE WITH TF-IDF WEIGHTED W2V ENCODING OF PROJECT TITLE FEATURE ")
plt.show()
```

TSNE WITH TF-IDF WEIGHTED W2V ENCODING OF PROJECT TITLE FEATURE



# observation

1. Even here also the data points don't form any observable clusters and coincides with each other.No result can be drawn.

## 2.5 Summary

Write few sentences about the results that you obtained and the observations you made.
1.DE(90%) is the state with most projects accepted followed by ND(89%) and WA(88%).

2.VT(80%) is the state with lowest number of projects accepted and DC(80%) being second lowest and TX(81%) being the third lowest state.

3.For Prek to second grade students, the most projects were propsed and the number of project decreases as grades increases.

4.Female teachers projects were approved more than male teachers.

5.For 9th to 12th grade students, the projects proposed are less and those projects that were accepted are also less.

6.Most of the accepted projects come under Literacy and language (87%)

7.hunger warmth and care category projects are most welcomed and donated (93.5%)

8.Most of the registered projects are in Literacy and Language category and followed by Maths and Science

9.The highest acceptance rate 88% in subcategory is Literature.

10.Health and wellness subcategory has the lowest number of projects proposed.

11.Most projects have 3 to 5 words in the title and no one has exceeded 10 words.

12.Approved projects have more number of words in essays when compared to that of rejected projects

13.Lower cost projects are approved a lot than that of high cost projects .

14.New teacher's proposals are also welcomed a lot and 82% projects are approved even the proposed teacher have no prior proposals

15.when teachers propose different projects(19) the acceptance of the project seems to be high.

16.Numerical value in project proposals also seems to increase the acceptance rate for the project.

17.Average cost of the project is 298 dollars and On an average requires 17 items to be buyed for a project

18.TSNE with Bag of Words, TF-IDF, Avg Word2Vec, TF-IDF Weighted Word2Vec doesn't give any useful observation since datapoints are so close together.

```
In [ ]:

```