# 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><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><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><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><br>- `Music & The Arts`<br>- `Literacy & Language, Math & Science` |
| `school_state` | State where school is located ([Two-letter U.S. postal code](#)). **Example:** `WY` |
| `project_subject_subcategories` | One or more (comma-separated) subject subcategories for the project.**Examples:**<br><br>- `Literacy`<br>- `Literature & Writing, Social Sciences` |
| `project_resource_summary` | An explanation of the resources needed for the project.**Example:**<br><br>- `My students need hands on literacy materials to manage sensory needs!` |
| `project_essay_1` | First application essay[*] |
| `project_essay_2` | Second application essay[*] |
| `project_essay_3` | Third application essay[*] |

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

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

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

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

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

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

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

## Notes on the Essay Data

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

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

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

In [1]:

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
```

```python
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os

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

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

## 1.1 Reading Data

In [2]:

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

In [3]:

```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 [4]:

```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[4]:

| id | description | quantity | price |
|---|---|---|---|
| LC652 - Lakeshore Double-Space Mobile Drying |

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

# 1.2 Data Analysis

In [5]:

```python
# PROVIDE CITATIONS TO YOUR CODE IF YOU TAKE IT FROM ANOTHER WEBSITE.
# https://matplotlib.org/gallery/pie_and_polar_charts/pie_and_donut_labels.html#sphx-glr-gallery-p
ie-and-polar-charts-pie-and-donut-labels-py


y_value_counts = project_data['project_is_approved'].value_counts()
print("Number of projects thar are approved for funding ", y_value_counts[1], ", (",
(y_value_counts[1]/(y_value_counts[1]+y_value_counts[0]))*100,"%)")
print("Number of projects thar are not approved for funding ", y_value_counts[0], ", (",
(y_value_counts[0]/(y_value_counts[1]+y_value_counts[0]))*100,"%)")

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

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

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

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

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

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

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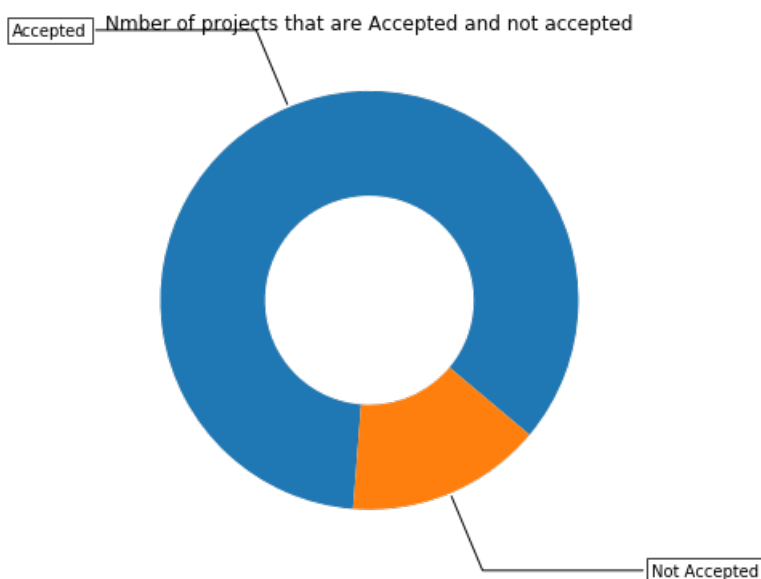
Number of projects thar are approved for funding  92706 , ( 84.85830404217927 %)
Number of projects thar are not approved for funding  16542 , ( 15.141695957820739 %)

**Summary:**

1. 85% of the projects got approved
2. 15% of projects got rejected

### 1.2.1 Univariate Analysis: School State

In [6]:

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

temp = pd.DataFrame(project_data.groupby("school_state")
["project_is_approved"].apply(np.mean)).reset_index()
# if you have data which contain only 0 and 1, then the mean = percentage (think about it)
temp.columns = ['state_code', 'num_proposals']

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

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

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

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

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

Out[6]:

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

In [7]:

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

```
States with lowest % approvals
   state_code  num_proposals
46         VT       0.800000
```

```
4.          V1          0.800000
7           DC          0.802326
43          TX          0.813142
26          MT          0.816327
18          LA          0.831245
==================================================
States with highest % approvals
    state_code  num_proposals
30          NH          0.873563
35          OH          0.875152
47          WA          0.876178
28          ND          0.888112
8           DE          0.897959
```

In [8]:

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

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

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

In [9]:

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

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

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

    if top:
        temp = temp[0:top]

    stack_plot(temp, xtick=col1, col2=col2, col3='total')
    print(temp.head(5))
    print("="*50)
    print(temp.tail(5))
```
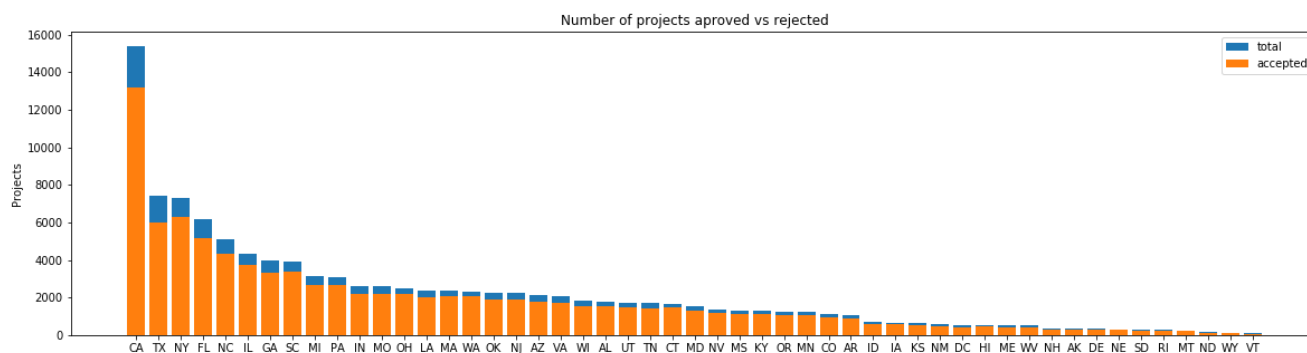
In [10]:

```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
```
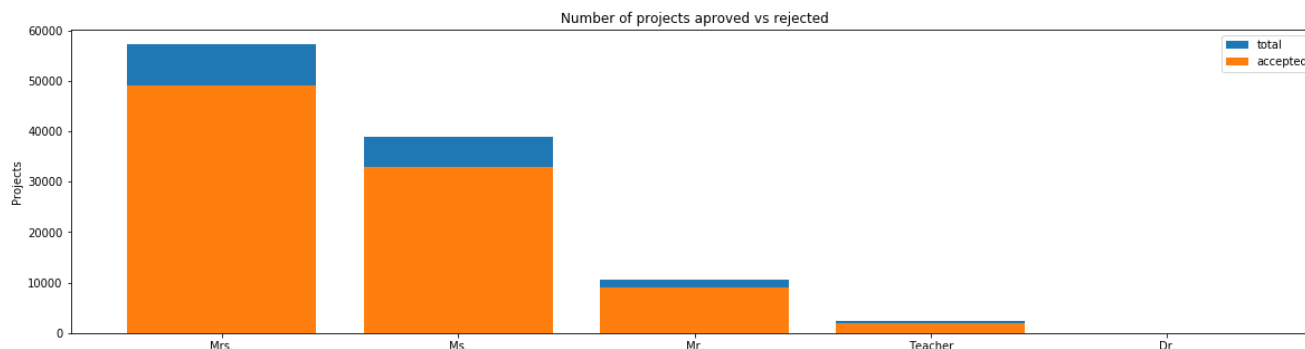
# summary:

1. Total number of projects submitted and number of projects approved by different states of U.S is plotted.
2. blue colour indicates total number of projects
3. orange colour indicates number of projects approved.
4. we can observe that the states which have submitted more projects in total have highest percentage of approval. The states which have submited lesser projects have lesser chance of approval greater than 80%
5. Highest number of projects approved is by state ND with 89% of approval, then comes NY and CA
6. The least project submitted and approved is by State VT with 80% approval rate then comes MT with 82% od approval rate.
7. overall every state is having more than 80% approval rate.

### 1.2.2 Univariate Analysis: teacher_prefix

In [11]:

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



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

# Summary

1. The plot is for project approval based on Teacher_prefix
2. Women holding prefixes Mrs. and Ms. have submitted more projects and approval rate is also more for them
3. Among all prefixes, Mrs. is the prefix having highest approval rate of 86%, then comes Ms. with 84% approval rate.
4. Dr. prefix people have submitted lesser projects(13) among 9 are selected which leads to 69% approval rate which is the lowest.
5. The women contribute more in submitting projects in that married women have highest success rate of approval.

### 1.2.3 Univariate Analysis: project_grade_category

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



Number of projects aproved vs rejected

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

# summary:

1. The above plot ia about project approval based on Grades.
2. The highest approval rate is for Grades 3-5 with 85% approval rate.
3. more number of projects submitted by Grades Prek-2 with 84% aprroval rate
4. Less number of project submitted and approved for Grades 9-12 with approval rate 84%
5. As the Grade increases the total number of projects submitted is decreased.

### 1.2.4 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 & H
unger"]
        if 'The' in j.split(): # this will split each of the catogory based on space "Math & Scienc
e"=> "Math","&", "Science"
            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i
.e removing 'The')
        j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
Science"=>"Math&Science"
        temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_') # we are replacing the & value into
    cat_list.append(temp.strip())
```
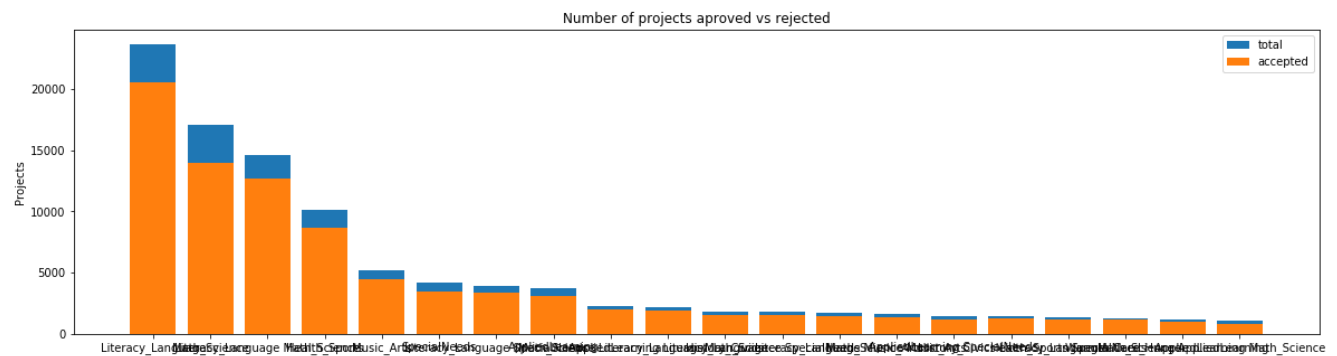
In [14]:

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

Out[14]:

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

In [15]:

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



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

# Summary:

1. plot is about project approval rate based on project subject categories.
2. Though the total projects are less in number but Warmth care_hunger subject category have highest approval rating of 93%.
3. The joint subject categories like have higher approval rating, like history_civics literacy_language(89%),Literacy_language math_science (87%)
4. AppliedLearning Math_science category has submitted least projects among all and approval rating is less i.e. 81%
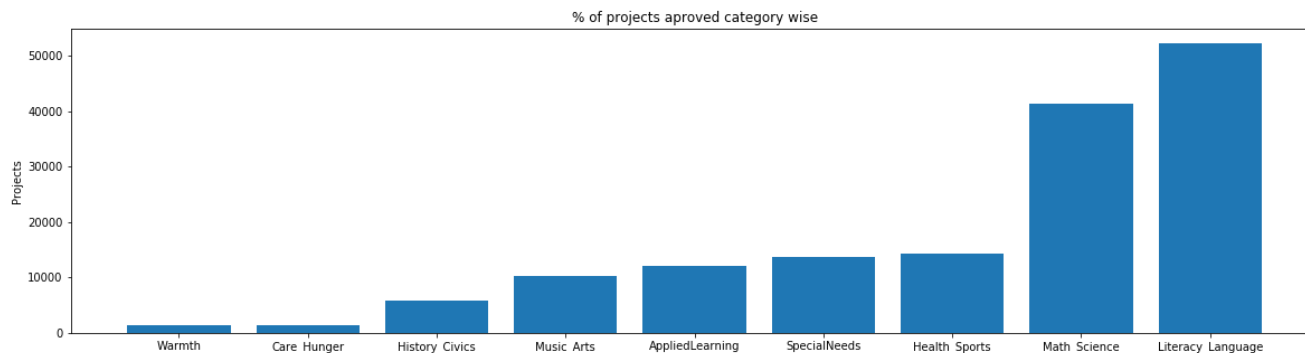
In [16]:

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

In [17]:

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


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

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

```
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
```

# Summary:

1. Plot is about total number ofprojects submitted based on individual subject categories.
2. Higest number of projects i.e. 52239 are submitted by Literacy and Language, then comes Math and Science, it's about 41421.
3. Least number i.e 1388 are submitted by Warmth, Care and Hunger, then comes History and Civics, it's about 5914.

### 1.2.5 Univariate Analysis: project_subject_subcategories

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

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python

sub_cat_list = []
for i in sub_catogories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & H
unger"]
        if 'The' in j.split(): # this will split each of the catogory based on space "Math & Scienc
e"=> "Math","&", "Science"
```

```
           j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i
.e removing 'The')
         j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
Science"=>"Math&Science"
         temp +=j.strip()+" "#" abc ".strip() will return "abc", remove the trailing spaces
         temp = temp.replace('&','_')
     sub_cat_list.append(temp.strip())
```
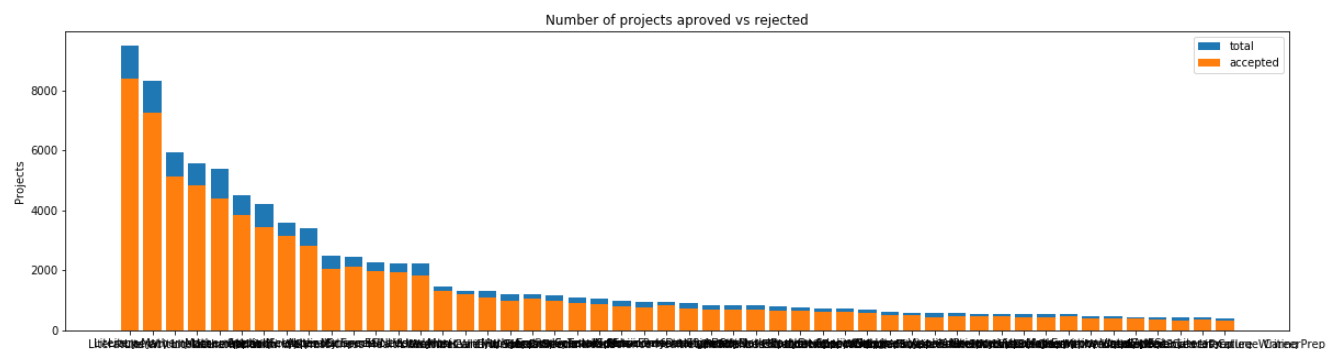
In [20]:

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

Out[20]:

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

In [21]:

```
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
```

# Summary:

1. The plot is about project approval rating based on project subject subcategories.
2. the highest project approval rate is for Literacy of 88% for 8325 and this category has submitted highest number of projects about 9486. next comes literacy mathematics with approval rate of 87% with 8325 number of projects.
3. Lowest project approval rate is for college careerPrep(421) and AppliedSciences College_CareerPrep(405) of 81%, then comes Mathematics(5379) with 82% approval rate.
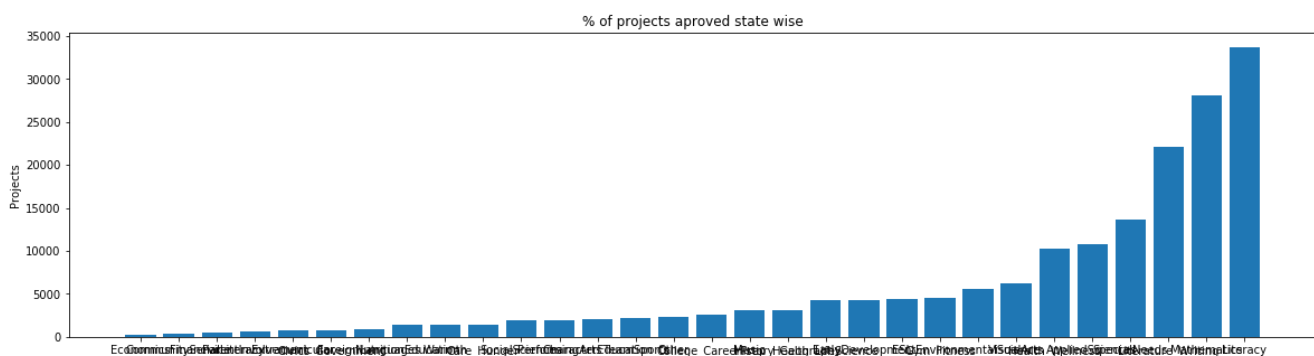
In [22]:

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

In [23]:

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


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

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



In [24]:

```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
```

# Summary :

1. The plot is about total number of projects submitted based on individual project subject subcategories.
2. The Highest number of projects are submitted by Literacy with 33700 projects. Then comes mathematics with 28074 projects
3. Least number of projects are submitted by Economics with 269 projects, next least one is Community Service with 441 projects.
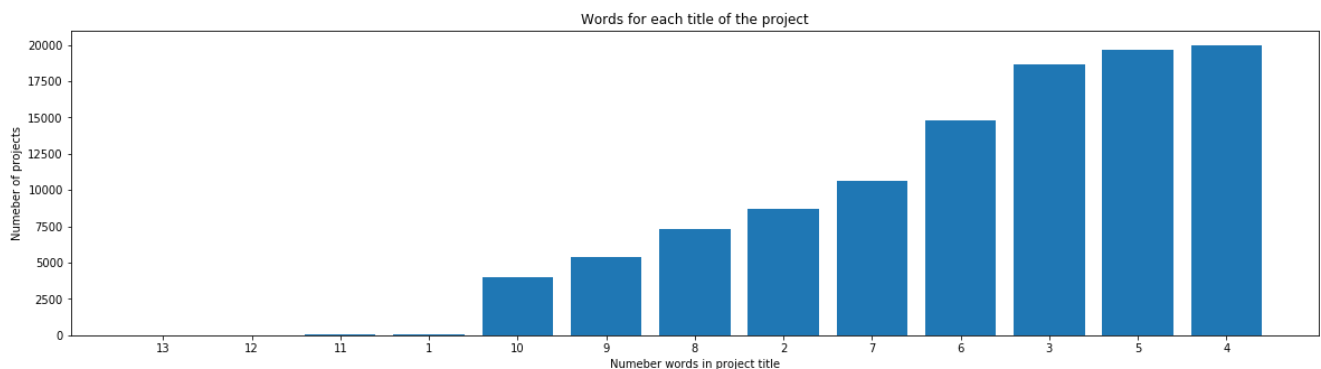
### 1.2.6 Univariate Analysis: Text features (Title)

In [25]:

```python
#How to calculate number of words in a string in DataFrame:
https://stackoverflow.com/a/37483537/4084039
word_count = project_data['project_title'].str.split().apply(len).value_counts()
word_dict = dict(word_count)
word_dict = dict(sorted(word_dict.items(), key=lambda kv: kv[1]))


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

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

Words for each title of the project



# Summary :

1. The plot is about to calculate the number of words in Project Title.
2. Most of the Projects(Around 20000) have 4 word in title, next comes 5 words for around 18000 projects.
3. Around 4000 projects have 10 words in title.
4. Least number of projects are found with 1 word, and 11 and more number of words.
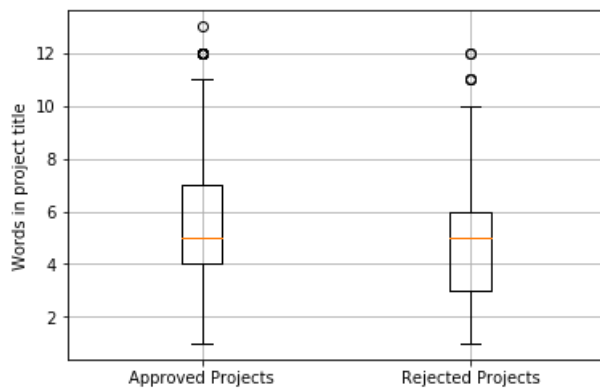
In [26]:

```python
approved_title_word_count = project_data[project_data['project_is_approved']==1]['project_title'].
str.split().apply(len)
approved_title_word_count = approved_title_word_count.values

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

In [27]:

```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()
```
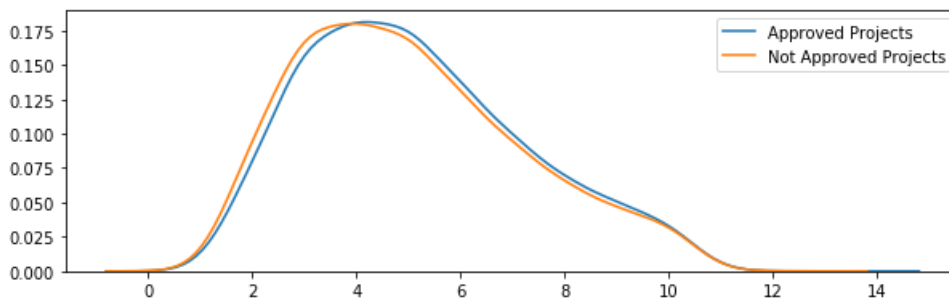
# Summary :

1. The BoxPlot is about number of words in project for approved projects and Rejected Projects.
2. 50th Percentile is very similar to both approved and Rejected Projects and it shows number words at 5oth percentile are 5.
3. Approved projects Boxplot is slightly higher than Rejected projects plot which means number of words in project title of approved projects are slightly more than rejected projects.
4. 25th and 75th percentile for approved projects are 4 and 7 words respectively.
5. 25th and 75th percentile for rejected projects are 3 and 6 words respectively.

In [28]:

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



# Summary :

1. This is the PDF plot of number of words in Project Title.
2. Blue line indicates approved projects and Orange line indicates Not approved projects.
3. Approved projects line is slightly ahead than the not approved projects indicating number of words are more in approved projects than not approved projects.

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

In [29]:

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

In [30]:

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

rejected_word_count = project_data[project_data['project_is_approved']==0]['essay'].str.split().app
ly(len)
rejected_word_count = rejected_word_count.values
```
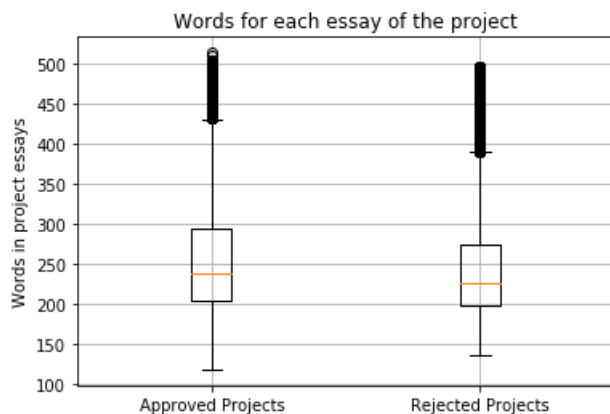
In [31]:

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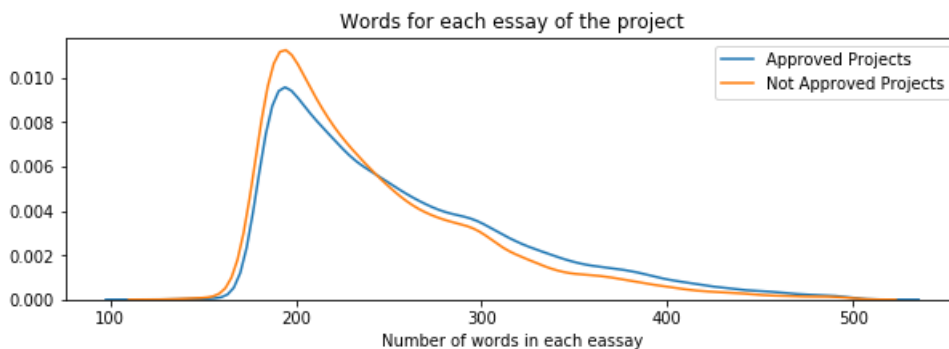


## Summary :

1. The BoxPlot is for number of words in project essays.
2. 50th percentile values are similar but Approved projects are slightly higher than Rejected Projects. Therefore Approved Project essays have more number of words than Rejected Projects.

In [32]:

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



## Summary :

1. PDF plot is for number of words varying in approved and not approved projects.
2. Blue line indicates Approved projects which is slightly ahead than not approved projects indicated by orange line.

3. If the words are between 190 to 250 they have highest approval rate and it approval rate range is for 190 to 500 words.

### 1.2.8 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 |

```
# 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)
```

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

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

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

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

```
# 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()
```

# Summary :

1. BoxPlot is for approved and not approved projects based on cost of the projects.
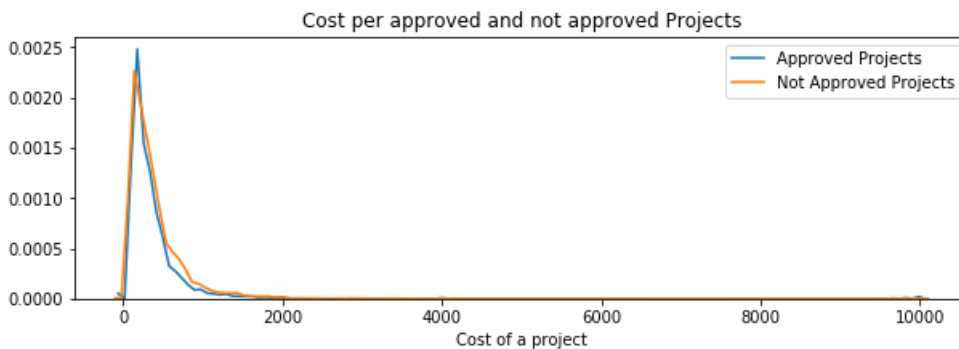2. 50th, 25th and 75th percentile are almost similar for approved and not approved projects.
3. The 25th and 75th percentile of not approved projects are slightly higher than approved projects which indicates that higher cost projects are not approved compared to approved projects.

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()
```



# Summary :

1. PDF plot is shown for approved and not approved projects based on cost of projects.
2. Both the plots are almost similar and if project cost is exeeding 500 or more then not approved project indicated by orange line is ahead of approved projects indicated by blue line.
3. projects of higher cost are not approved compared to approved projects.

In [39]:

```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               |

```
|    50    |    198.99   |    205.145   |
|    55    |    223.99   |    292.61    |
|    60    |    255.63   |    325.144   |
|    65    |   285.412   |    362.39    |
|    70    |   321.225   |    399.99    |
|    75    |   366.075   |    449.945   |
|    80    |    411.67   |    519.282   |
|    85    |    479.0    |    618.276   |
|    90    |    593.11   |    739.356   |
|    95    |   801.598   |    992.486   |
|   100    |    9999.0   |    9999.0    |
+----------+-------------+--------------+
```

# Summary :

1. Since Boxplot and PDF were not able analyse much the percentiles are plotted.
2. For almost all the percentiles cost of Approved projects are less than not approved projects.

   ```
    a. For 25th percentile cost of approved projects is 100 where as for not approved proje
   cts it is 141.
    b. For 50th percentile cost of approved projects is 199 where as for not approved proje
   cts it is 263.
    c. For 75th percentile cost of approved projects is 366 where as for not approved proje
   cts it is 450.
   ```

3. Overall maximum cost of approved projects should be 10000.

## 1.2.9 Univariate Analysis: teacher_number_of_previously_posted_projects

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

In [40]:

```python
# Count the total number of projects posted by Teachers

univariate_barplots(project_data,
'teacher_number_of_previously_posted_projects','project_is_approved', top=40)
```



```
   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  \
35                                            35                   243    267
36                                            36                   231    261
38                                            38                   224    246
```

```
37                                          37          196   226
40                                          40          202   221

        Avg
35  0.910112
36  0.885057
38  0.910569
37  0.867257
40  0.914027
```

# Summary:

1. This is plot for project approved based on number of projects submitted by same teacher.
2. The Teachers who have submitted more applications are less in number for example:

    a. 40 projects submitted adds upto total 221 and 202 are selected with 91% approval rate.

    b. 35 and 38 projects submitted have approval rate of 91%

    The approval rate is high for same teachers submitting projects above 36 as project aproval rate is greater than 87%.

1. teachers who submitted projects for the first time are more in number 30014, among which 24652 are excepted, leading to project approval rate of 82%.

### 1.2.10 Univariate Analysis: project_resource_summary

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

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

In [41]:

```python
res_summary= project_data['project_resource_summary'].value_counts()


s=project_data['project_resource_summary']
res_summary_list=[]

for l in s:
    res_summary_list.append(l)

print("number of summaries are",len(res_summary_list))

res_summary_list[0:20]
```

```
number of summaries are 109248
```

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 materia
ls for the upcoming Soccer season.',
 'My students need to engage in Reading and Math in a way that will inspire them with these Mini i
Pads!',
 'My students need hands on practice in mathematics. Having fun and personalized journals and char
ts 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 al
l different types of needs, flexible seating would assist not only these students with special nee
ds, 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 int
erest in learning.',
 "My students need three devices and three management licenses for small group's easy access to ne
wly-implemented online programs--Go Noodle Plus, for increased in-class physical activity and Ligh
t Sail, an interactive reading program.",
 'My students need great books to use during Independent Reading, Read Alouds, Partner Reading and
Author Studies.'
```

Author Studies.',
 'My students need books by their favorite authors like Chris Grabenstein, Raina Telgemeier and Ja
mes Patterson to keep building their vocabulary.',
 'My students need paper, three chromebooks, and a keyboard cover to enhance our learning center a
nd introduce parents to the technology and programs being used at our site.',
 'My students need 3D and 4D life science activity kits so they can get the full effect of the les
sons and work in their groups and store their models in sturdy bins.',
 'My students need access to technology that will allow them to communicate with others.',
 'My students need 5 tablets for our classroom technology center PLEASE!!!',
 'My students need activities to play during recess, which takes place on a large green open
space. There is no equipment, so there is strong need for kick balls, soccer balls, hula hoops, pa
rachute, and frisbees.',
 'My students need 2 LeapPad that will engage them in critical thinking skills, vocabulary and com
munication skills that they need to be successful.',
 'My students need Chromebooks to publish written work, continue working n their typing skills, ma
nipulate reading passages by highlighting and analyzing a text, and so much more!',
 'My students need privacy partitions to use while testing or while working on independent work.',
 'My students need 7 Hokki stools to encourage and allow for movement during times of stationary w
ork in the classroom.']

In [42]:

```python
num_res_summary={}

for n in tqdm(range(len(res_summary_list))):
    for r in res_summary_list[n].split():
        if r.isdigit():
            num_res_summary[n]=int(r)

print(num_res_summary[16])
print(num_res_summary[19])
```

```
100%|██████████████████████████████████████████████████| 109248/109248
[00:00<00:00, 190561.94it/s]
```

2
7

## 1.3 Text preprocessing

### 1.3.1 Essay Text

In [43]:

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

Out[43]:

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

In [44]:

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

```python
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
print("="*50)
print(project_data['essay'].values[99999])
print("="*50)
```

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. \r\n\r\n We have over 24 languages represented in our English Learner program with students at every level of mastery.  We also have over 40 countries represented with the families within our school.  Each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures, beliefs, and respect.\"The limits of your language are the limits of your world.\"-Ludwig Wittgenstein  Our English learner's have a strong support system at home that begs for more resources.  Many times our parents are learning to read and speak English along side of their children.  Sometimes this creates barriers for parents to be able to help their child learn phonetics, letter recognition, and other reading skills.\r\n\r\nBy providing these dvd's and players, students are able to continue their mastery of the English language even if no one at home is able to assist.  All families with students within the Level 1 proficiency status, will be a offered to be a part of this program.  These educational videos will be specially chosen by the English Learner Teacher and will be sent home regularly to watch.  The videos are to help the child develop early reading skills.\r\n\r\nParents that do not have access to a dvd player will have the opportunity to check out a dvd player to use for the year.  The plan is to use these videos and educational dvd's for the years to come for other EL students.\r\nnannan
==================================================
The 51 fifth grade students that will cycle through my classroom this year all love learning, at least most of the time. At our school, 97.3% of the students receive free or reduced price lunch. Of the 560 students, 97.3% are minority students. \r\nThe school has a vibrant community that loves to get together and celebrate. Around Halloween there is a whole school parade to show off the beautiful costumes that students wear. On Cinco de Mayo we put on a big festival with crafts made by the students, dances, and games. At the end of the year the school hosts a carnival to celebrate the hard work put in during the school year, with a dunk tank being the most popular activity.My students will use these five brightly colored Hokki stools in place of regular, stationary, 4-legged chairs. As I will only have a total of ten in the classroom and not enough for each student to have an individual one, they will be used in a variety of ways. During independent reading time they will be used as special chairs students will each use on occasion. I will utilize them in place of chairs at my small group tables during math and reading times. The rest of the day they will be used by the students who need the highest amount of movement in their life in order to stay focused on school.\r\n\r\nWhenever asked what the classroom is missing, my students always say more Hokki Stools. They can't get their fill of the 5 stools we already have. When the students are sitting in group with me on the Hokki Stools, they are always moving, but at the same time doing their work. Anytime the students get to pick where they can sit, the Hokki Stools are the first to be taken. There are always students who head over to the kidney table to get one of the stools who are disappointed as there are not enough of them. \r\n\r\nWe ask a lot of students to sit for 7 hours a day. The Hokki stools will be a compromise that allow my students to do desk work and move at the same time. These stools will help students to meet their 60 minutes a day of movement by allowing them to activate their core muscles for balance while they sit. For many of my students, these chairs will take away the barrier that exists in schools for a child who can't sit still.nannan
==================================================
How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day.\r\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas.\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an \"open classroom\" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more.With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade.  This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs lost of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you!nannan
==================================================
My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out

for my students. I teach in a Title I school where most of the students receive free or reduced pr
ice lunch.  Despite their disabilities and limitations, my students love coming to school and come
eager to learn and explore.Have you ever felt like you had ants in your pants and you needed to gr
oove and move as you were in a meeting? This is how my kids feel all the time. The want to be able
to move as they learn or so they say.Wobble chairs are the answer and I love then because they dev
elop their core, which enhances gross motor and in Turn fine motor skills. \r\nThey also want to l
earn through games, my kids don't want to sit and do worksheets. They want to learn to count by ju
mping and playing. Physical engagement is the key to our success. The number toss and color and sh
ape mats can make that happen. My students will forget they are doing work and just have the fun a
6 year old deserves.nannan
==================================================
The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The grea
t teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is makeup is 97.6% Af
rican-American, making up the largest segment of the student body. A typical school in Dallas is m
ade up of 23.2% African-American students. Most of the students are on free or reduced lunch. We a
ren'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 smar
t, effective, efficient, and disciplined students with good character.In our classroom we can util
ize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the so
und enough to receive the message. Due to the volume of my speaker my students can't hear videos o
r books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my s
tudents will be able to hear and I can stop, pause and replay it at any time.\r\nThe cart will all
ow me to have more room for storage of things that are needed for the day and has an extra part to
it I can use.  The table top chart has all of the letter, words and pictures for students to learn
about different letters and it is more accessible.nannan
==================================================

In [45]:

```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 [46]:

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

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

In [47]:

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

```python
sent = sent.replace('\\n', ' ')
print(sent)
```

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

In [48]:

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

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

In [49]:

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

In [50]:

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

```
100%|████████████████████████████████████████████████████████████████| 109248/109248
[05:40<00:00, 320.59it/s]
```

In [51]:

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

Out[51]:

'my kindergarten students varied disabilities ranging speech language delays cognitive delays gros
s 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 lunc
h 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 w
obble 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 de
serves nannan'

## 1.3.2 Project title Text

In [52]:

```
# similarly you can preprocess the titles also
# project_title
projtitle=list(project_data['project_title'].values)
print(len(projtitle))
```

```
109248
```

In [53]:

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

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

In [54]:

```
sent1 = decontracted(project_data['project_title'].values[15000])
print(sent1)
print("="*50)
```

\r\nThe \"i\" Classroom
==================================================

In [55]:

```
# \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
sent1 = sent1.replace('\\r', ' ')
sent1 = sent1.replace('\\"', ' ')
sent1 = sent1.replace('\\n', ' ')
print(sent1)
```

  The  i  Classroom

In [56]:

```
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent1 = re.sub('[^A-Za-z0-9]+', ' ', sent1)
print(sent1)
```

 The i Classroom

In [57]:

```
# Combining all the above statemennts
from tqdm import tqdm
preprocessed_projtitles = []
# tqdm is for printing the status bar
for sentance in tqdm(project_data['project_title'].values):
    sent1 = decontracted(sentance)
    sent1 = sent1.replace('\\r', ' ')
    sent1 = sent1.replace('\\"', ' ')
    sent1 = sent1.replace('\\n', ' ')
    sent1 = re.sub('[^A-Za-z0-9]+', ' ', sent1)
    # https://gist.github.com/sebleier/554280
    sent1 = ' '.join(e for e in sent1.split() if e not in stopwords)
    preprocessed_projtitles.append(sent1.lower().strip())
```

```
100%|████████████████████████████████████████████████| 109248/109248
[00:16<00:00, 6655.87it/s]
```

In [58]:

```
# after preprocesing
print(preprocessed_projtitles[10000])
print("="*50)
print(preprocessed_projtitles[20000])
print("="*50)
print(preprocessed_projtitles[30000])
print("="*50)
print(preprocessed_projtitles[40000])
print("="*50)
print(preprocessed_projtitles[50000])
```

family book clubs
==================================================
we need to move it while we input it
==================================================
weaving together
==================================================
let them read let them learn let them thrive
==================================================
help bridgeport students improve their listening skills

# 1. 4 Preparing data for models

```
project_data.columns
```

Out[59]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
       'project_submitted_datetime', 'project_grade_category', 'project_title',
       'project_essay_1', 'project_essay_2', 'project_essay_3',
       'project_essay_4', 'project_resource_summary',
       'teacher_number_of_previously_posted_projects', 'project_is_approved',
       'clean_categories', 'clean_subcategories', 'essay', 'price',
       'quantity'],
      dtype='object')
```

we are going to consider

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

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

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

## 1.4.1 Vectorizing Categorical data

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

In [60]:

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


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

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

In [61]:

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


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

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

In [62]:

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

In [63]:

```python
# school_state
# count all the words
from collections import Counter
school_state=Counter()
for state in project_data['school_state'].values:
    school_state.update(state.split())
```

In [64]:

```python
# dict sort by value
state_dict=dict(school_state)
sorted_school_dict=dict(sorted(state_dict.items(), key=lambda kv: kv[1]))
```

In [65]:

```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_school_dict.keys()), lowercase=False, binary=T
rue)
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', 'HI', 'DC', 'NM', 'KS', 'I
A', 'ID', 'AR', 'CO', 'MN', 'OR', 'KY', 'MS', 'NV', 'MD', 'CT', 'TN', 'UT', 'AL', 'WI', 'VA', 'AZ',
'NJ', 'OK', 'WA', 'MA', 'LA', 'OH', 'MO', 'IN', 'PA', 'MI', 'SC', 'GA', 'IL', 'NC', 'FL', 'NY', 'TX
', 'CA']
Shape of matrix after one hot encodig  (109248, 51)
```

In [66]:

```python
# teacher_prefix
# count all the words
from collections import Counter
teacher_prefix=Counter()
for prefix in project_data['teacher_prefix'].values:
    prefix=str(prefix)
    teacher_prefix.update(prefix.split())
```

In [67]:

```python
# dict sort by value
prefix_dict=dict(teacher_prefix)
sorted_teacher_prefix_dict=dict(sorted(prefix_dict.items(), key=lambda kv: kv[1]))
```

In [68]:

```python
# we use count vectorizer to convert the values into one hot encoded features
# teacher_prefix
```

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary=list(sorted_teacher_prefix_dict.keys()), lowercase=False,
binary=True)
vectorizer.fit(project_data['teacher_prefix'].values.astype("U"))
print(vectorizer.get_feature_names())


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

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

In [69]:

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

```
Grades PreK-2
==================================================
Grades 3-5
==================================================
Grades 3-5
==================================================
Grades PreK-2
==================================================
Grades PreK-2
==================================================
```

In [70]:

```
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 [71]:

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

```
Grades PreK-2
==================================================
```

In [72]:

```
Grade = Grade.replace('Grades', ' ')
```

```
print(Grade)
```

```
  PreK-2
```

In [73]:

```python
from tqdm import tqdm
preprocessed_Grades = []
# tqdm is for printing the status bar
for grades in tqdm(project_data['project_grade_category'].values):
    Grade = decontracted(grades)
    Grade=Grade.replace('Grades', ' ')
    #sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    #Grade = ' '.join(e for e in Grade.split() if e not in stopwords)
    preprocessed_Grades.append(Grade.lower().strip())
```

```
100%|████████████████████████████████████████████████████████| 109248/109248
[00:07<00:00, 14773.20it/s]
```

In [74]:

```python
preprocessed_Grades[20000]
```

Out[74]:

```
'prek-2'
```

In [75]:

```python
# project_grade_category
from collections import Counter
project_grade_category=Counter()
for grade in project_data['project_grade_category'].values:
    #grade=str(grade)
    grade=grade.replace('Grades', ' ')
    project_grade_category.update(grade.split())
```

In [76]:

```python
# dict sort by value
grade_dict=dict(project_grade_category)
sorted_project_grade_cat_dict=dict(sorted(grade_dict.items(), key=lambda kv: kv[1]))
```

In [77]:

```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_project_grade_cat_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(project_data['project_grade_category'].values)
print(vectorizer.get_feature_names())


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

```
['9-12', '6-8', '3-5', 'PreK-2']
Shape of matrix after one hot encodig  (109248, 4)
```

### 1.4.2 Vectorizing Text data

#### 1.4.2.1 Bag of words

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

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

### 1.4.2.2 Bag of Words on `project_title`

```
# you can vectorize the title also
# before you vectorize the title make sure you preprocess it
# Similarly you can vectorize for title also
vectorizer = CountVectorizer(min_df=10)
text_bow = vectorizer.fit_transform(preprocessed_projtitles)
print("Shape of matrix after one hot encodig ",text_bow.shape)
```

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

### 1.4.2.3 TFIDF vectorizer

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

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

### 1.4.2.4 TFIDF Vectorizer on `project_title`

```
# Similarly you can vectorize for title also
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
text_tfidf = vectorizer.fit_transform(preprocessed_projtitles)
print("Shape of matrix after one hot encodig ",text_tfidf.shape)
```

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

### 1.4.2.5 Using Pretrained Models: Avg W2V

```
# 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')
```

Loading Glove Model

```
1917495it [23:23, 1366.08it/s]
```

Done. 1917495  words loaded!

In [83]:

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

for i in preprocessed_projtitles:
    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-sav
e-and-load-variables-in-python/

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

```
all the words in the coupus 17014413
the unique words in the coupus 58968
The number of words that are present in both glove vectors and our coupus 51503 ( 87.341 %)
word 2 vec length 51503
```

In [84]:

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

In [85]:

```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(preprocessed_essays,disable=True): # 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]))
```

```
109248
300
```

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

In [86]:

```python
# Similarly you can vectorize for title also
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_projtitles, disable=True): # 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]))
```

```
109248
300
```

### 1.4.2.7 Using Pretrained Models: TFIDF weighted W2V

In [87]:

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

In [88]:

```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(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors.append(vector)

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

```
100%|███████████████████████████████████████████████████████| 109248/109248
[21:35<00:00, 84.32it/s]
```

```
109248
300
```

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

```
# Similarly you can vectorize for title also
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(preprocessed_projtitles)
# 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())
```

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_projtitles): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors.append(vector)

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

```
100%|████████████████████████████████████████████████████| 109248/109248
[00:18<00:00, 5774.41it/s]
```

```
109248
300
```

### 1.4.3 Vectorizing Numerical features

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

# price_standardized = standardScalar.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329.   ... 399.   287.
73   5.5 ].
# Reshape your data either using array.reshape(-1, 1)

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

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

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

```
price_standardized
```

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

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

# Now standardize the data with above maen and variance.
teacher_num_standardized =
teacher_num_scalar.transform(project_data['teacher_number_of_previously_posted_projects'].values.r
eshape(-1, 1))

import warnings
warnings.filterwarnings('ignore')
```

```
C:\Users\Santosh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:
DataConversionWarning:

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

```
Mean : 11.153165275336848, Standard deviation : 367.49634838483496
```

```
C:\Users\Santosh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:
DataConversionWarning:

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

```
teacher_num_standardized
```

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

### 1.4.4 Merging all the above features

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

```
print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(school_state_categories_one_hot.shape)
print(teacher_prefix_one_hot.shape)
print(project_grade_categories_one_hot.shape)
print(text_bow.shape)
print(price_standardized.shape)
print(teacher_num_standardized.shape)
```

```
(109248, 9)
(109248, 30)
(109248, 51)
(109248, 6)
(109248, 4)
(109248, 3329)
(109248, 1)
(109248, 1)
```

In [96]:

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

Out[96]:

```
(109248, 3431)
```

# Assignment 2: Apply TSNE

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

1. In the above cells we have plotted and analyzed many features. Please observe the plots and write the observations in markdown cells below every plot.
2. EDA: Please complete the analysis of the feature: teacher_number_of_previously_posted_projects
3.      Build the data matrix using these features
   - school_state : categorical data (one hot encoding)
   - clean_categories : categorical data (one hot encoding)
   - clean_subcategories : categorical data (one hot encoding)
   - teacher_prefix : categorical data (one hot encoding)
   - project_grade_category : categorical data (one hot encoding)
   - project_title : text data (BOW, TFIDF, AVG W2V, TFIDF W2V)
   - price : numerical
   - teacher_number_of_previously_posted_projects : numerical
4. Now, plot FOUR t-SNE plots with each of these feature sets.
   A. categorical, numerical features + project_title(BOW)
   B. categorical, numerical features + project_title(TFIDF)
   C. categorical, numerical features + project_title(AVG W2V)
   D. categorical, numerical features + project_title(TFIDF W2V)
5. Concatenate all the features and Apply TNSE on the final data matrix
6. Note 1: The TSNE accepts only dense matrices
7. Note 2: Consider only 5k to 6k data points to avoid memory issues. If you run into memory error issues, reduce the number of data points but clearly state the number of datat-poins you are using

In [97]:

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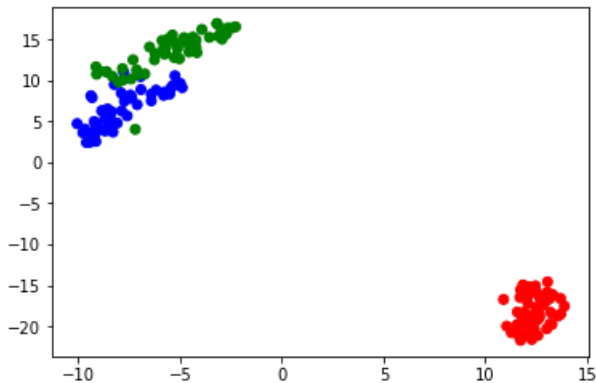
```
# iſ X is a sparse matrix you need to pass it as X_embedding = tsne.fit_transform(X.toarray()) , .
toarray() will convert the sparse matrix into dense matrix

for_tsne = np.hstack((X_embedding, y.reshape(-1,1)))
for_tsne_df = pd.DataFrame(data=for_tsne, columns=['Dimension_x','Dimension_y','Score'])
colors = {0:'red', 1:'blue', 2:'green'}
plt.scatter(for_tsne_df['Dimension_x'], for_tsne_df['Dimension_y'], c=for_tsne_df['Score'].apply(la
mbda x: colors[x]))
plt.show()
```



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

**please write all of the code with proper documentation and proper titles for each subsection**

**when you plot any graph make sure you use**

**a. Title, that describes your plot, this will be very helpful to the reader**

**b. Legends if needed**

**c. X-axis label**

**d. Y-axis label**

In [98]:

```
print("Shape of matrix after one hot encodig ",text_bow.shape)
```

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

In [99]:

```
print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(school_state_categories_one_hot.shape)
print(teacher_prefix_one_hot.shape)
print(project_grade_categories_one_hot.shape)
print(text_bow.shape)
print(price_standardized.shape)
print(teacher_num_standardized.shape)

# 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_n = hstack((categories_one_hot,
```

```
           sub_categories_one_hot,school_state_categories_one_hot,teacher_prefix_one_hot,
                   project_grade_categories_one_hot, text_bow,
           price_standardized,teacher_num_standardized))

x_n.shape
```

```
(109248, 9)
(109248, 30)
(109248, 51)
(109248, 6)
(109248, 4)
(109248, 3329)
(109248, 1)
(109248, 1)
```

```
(109248, 3431)
```

In [100]:

```python
from sklearn.manifold import TSNE

print(x_n.shape)


x_n=x_n.tocsr() # https://stackoverflow.com/questions/30163830/accessing-elements-in-coo-matrix
x_n=x_n[0:1000, :]

# TypeError: A sparse matrix was passed, but dense data is required for method="barnes_hut".
#Use X.toarray() to convert to a dense numpy array if the array is small enough for it to fit in m
emory.
#Otherwise consider dimensionality reduction techniques (e.g. TruncatedSVD)


x_nt=x_n.toarray()

data_5k=x_nt[0:1000, :]
print('The Shape of data', data_5k.shape)
```

```
(109248, 3431)
The Shape of data (1000, 3431)
```

In [101]:

```python
model = TSNE(n_components=2, perplexity=50, random_state=0)
# configuring the parameteres
# the number of components = 2
# default perplexity = 30
# default learning rate = 200
# default Maximum number of iterations for the optimization = 1000

tsne_data_bow = model.fit_transform(data_5k)
```

In [102]:

```python
tsne_data_bow.shape
```

Out[102]:

```
(1000, 2)
```

In [103]:

```python
label=project_data['project_is_approved']
label_5k=label[0:1000]
label_5k.shape
```

Out[103]:

```
(1000,)
```

```
# creating a new data frame which help us in ploting the result data
tsne_data_bow = np.vstack((tsne_data_bow.T, label_5k)).T
tsne_df_bow = pd.DataFrame(data=tsne_data_bow, columns=("Dim_1", "Dim_2", "label"))
```
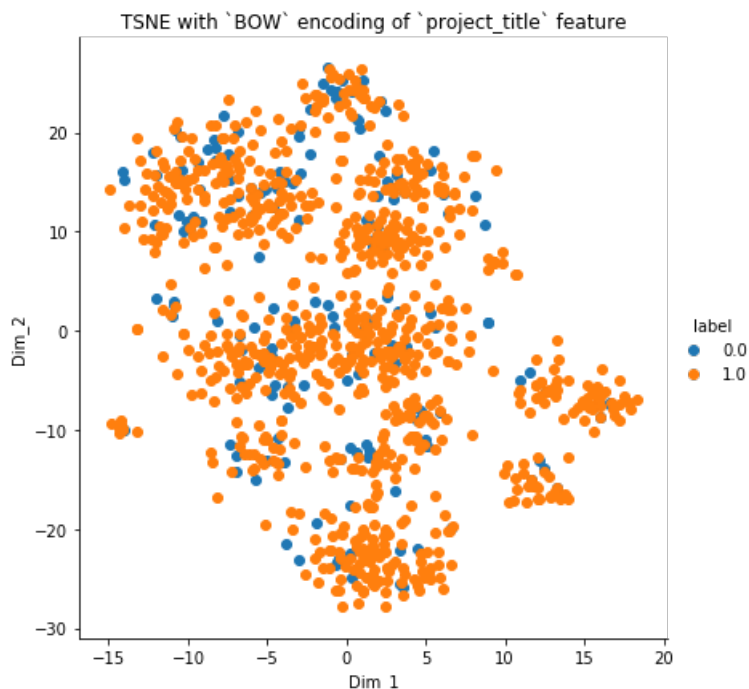
In [105]:

```
tsne_df_bow.shape
```

Out[105]:

```
(1000, 3)
```

In [106]:

```
# Ploting the result of tsne
sns.FacetGrid(tsne_df_bow, hue="label", size=6).map(plt.scatter, 'Dim_1', 'Dim_2').add_legend()
plt.title("TSNE with `BOW` encoding of `project_title` feature")
plt.show()
```



## Summary

1.The above plot shows TSNE with BOW encoding of project title

1. orange dots(1) indicate approved
2. blue dots(0) indicate not approved
3. Since blue points are suppressed under orange points we can not conclude anything about it

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

## please write all the code with proper documentation, and proper titles for each subsection

## when you plot any graph make sure you use

#a. Title, that describes your plot, this will be very helpful to the reader

```
#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
```

In [107]:

```
print("Shape of matrix after one hot encodig ",text_tfidf.shape)
```

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

In [108]:

```
x_n = hstack((categories_one_hot,
sub_categories_one_hot,school_state_categories_one_hot,teacher_prefix_one_hot,
              project_grade_categories_one_hot, text_tfidf,
price_standardized,teacher_num_standardized))

x_n.shape
```

Out[108]:

(109248, 3431)

In [109]:

```
#Tfidf
project_title_tfidf=x_n
project_title_tfidf=project_title_tfidf.tocsr() #
https://stackoverflow.com/questions/30163830/accessing-elements-in-coo-matrix

project_title_tfidf=project_title_tfidf.toarray()

data_tfidf_1k=project_title_tfidf[0:1000, :]
print('The Shape of data', data_tfidf_1k.shape)
```

The Shape of data (1000, 3431)

In [110]:

```
model = TSNE(n_components=2, perplexity=80, random_state=0)

tsne_data_tfidf = model.fit_transform(data_tfidf_1k)
```

In [111]:

```
tsne_data_tfidf.shape
```

Out[111]:

(1000, 2)

In [112]:

```
# creating a new data frame which help us in ploting the result data
tsne_data_tfidf = np.vstack((tsne_data_tfidf.T, label_5k)).T
tsne_df_tfidf = pd.DataFrame(data=tsne_data_tfidf, columns=("Dim_1", "Dim_2", "label"))
```
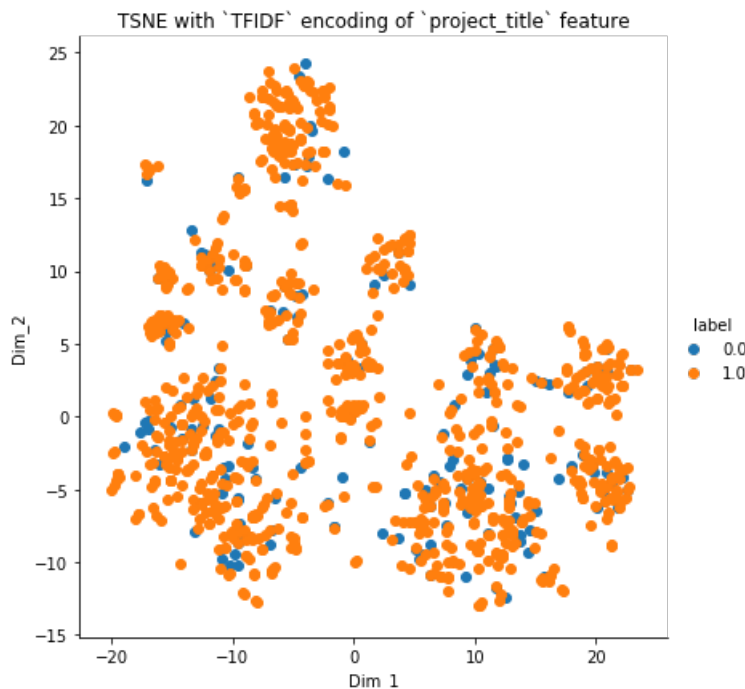
In [113]:

```
tsne_df_tfidf.shape
```

Out[113]:

(1000, 3)

In [114]:

```
# Ploting the result of tsne
sns.FacetGrid(tsne_df_tfidf, hue="label", size=6).map(plt.scatter, 'Dim_1', 'Dim_2').add_legend()
plt.title('TSNE with `TFIDF` encoding of `project_title` feature')
plt.show()
```


TSNE with `TFIDF` encoding of `project_title` feature

## Summary

1.The above plot shows TSNE with `TFIDF` encoding of `project_title` feature

1. orange dots(1) indicate approved
2. blue dots(0) indicate not approved
3. Since blue points are suppressed under orange points we can not conclude anything about it

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

### please write all the code with proper documentation, and proper titles for each subsection

### when you plot any graph make sure you use

```
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

In [115]:

```
print(len(avg_w2v_vectors))
```

```
109248
```

In [116]:

```
x_n = hstack((categories_one_hot,
sub_categories_one_hot,school_state_categories_one_hot,teacher_prefix_one_hot,
            project_grade_categories_one_hot, avg_w2v_vectors,
price_standardized,teacher_num_standardized))
```

```
price_standardized,teacher_num_standardized,))
x_n.shape
```

Out[116]:

```
(109248, 402)
```

In [117]:

```
#avg_w2v
project_title_avgw2k=x_n

project_title_avgw2k=project_title_avgw2k.tocsr() #
https://stackoverflow.com/questions/30163830/accessing-elements-in-coo-matrix

project_title_avgw2k=project_title_avgw2k.toarray()

data_avgw2k_1k=project_title_avgw2k[0:1000]

print('The length of data', len(data_avgw2k_1k))
data_avgw2k_1k.shape
```

The length of data 1000

Out[117]:

```
(1000, 402)
```

In [118]:

```
model = TSNE(n_components=2, perplexity=80, random_state=0)

tsne_data_avgw2k = model.fit_transform(data_avgw2k_1k)
```

In [119]:

```
tsne_data_avgw2k.shape
```

Out[119]:

```
(1000, 2)
```

In [120]:

```
# creating a new data frame which help us in ploting the result data
tsne_data_avgw2k = np.vstack((tsne_data_avgw2k.T, label_5k)).T
tsne_df_avgw2k = pd.DataFrame(data=tsne_data_avgw2k, columns=("Dim_1", "Dim_2", "label"))
```
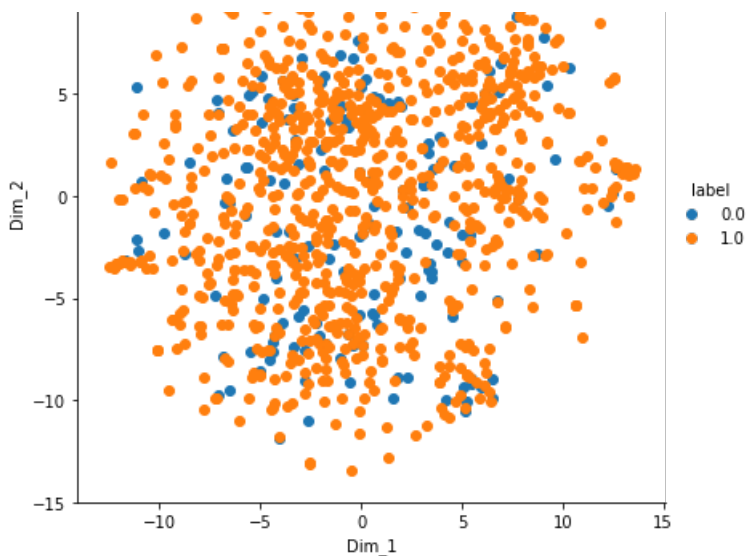
In [121]:

```
tsne_df_avgw2k.shape
```

Out[121]:

```
(1000, 3)
```

In [122]:

```
# Ploting the result of tsne
sns.FacetGrid(tsne_df_avgw2k, hue="label", size=6).map(plt.scatter, 'Dim_1', 'Dim_2').add_legend()
plt.title('TSNE with `AVG W2V` encoding of `project_title` feature')
plt.show()
```

# Summary

1. Plot shows Projects Approved or not approved using TSNE with `AVG W2V` encoding of `project_title` feature
2. orange dots(1) indicate approved
3. blue dots(0) indicate not approved
4. Since blue points are suppressed under orange points we can not conclude anything about it

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

## please write all the code with proper documentation, and proper titles for each subsection

## when you plot any graph make sure you use

```
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

In [123]:

```
print(len(tfidf_w2v_vectors))
```

109248

In [124]:

```
x_n = hstack((categories_one_hot,
sub_categories_one_hot,school_state_categories_one_hot,teacher_prefix_one_hot,
            project_grade_categories_one_hot, tfidf_w2v_vectors, price_standardized,teacher_num_s
tandardized))

x_n.shape
```

Out[124]:

(109248, 402)

In [125]:

```
#tfidf_w2v_vectors
project_title_tfidf_w2v=x_n
```

```
project_title_tfidf_w2v=project_title_tfidf_w2v.tocsr() #
https://stackoverflow.com/questions/30163830/accessing-elements-in-coo-matrix

project_title_tfidf_w2v=project_title_tfidf_w2v.toarray()

data_tfidf_w2v_1k=project_title_tfidf_w2v[0:1000]

data_tfidf_w2v_1k.shape
```

Out[125]:

```
(1000, 402)
```

In [126]:

```
model = TSNE(n_components=2, perplexity=80, random_state=0)

tsne_data_tfidf_w2v = model.fit_transform(data_tfidf_w2v_1k)
```

In [127]:

```
tsne_data_tfidf_w2v.shape
```

Out[127]:

```
(1000, 2)
```

In [128]:

```
# creating a new data frame which help us in ploting the result data
tsne_data_tfidf_w2v = np.vstack((tsne_data_tfidf_w2v.T, label_5k)).T
tsne_df_tfidf_w2v = pd.DataFrame(data=tsne_data_tfidf_w2v, columns=("Dim_1", "Dim_2", "label"))
```
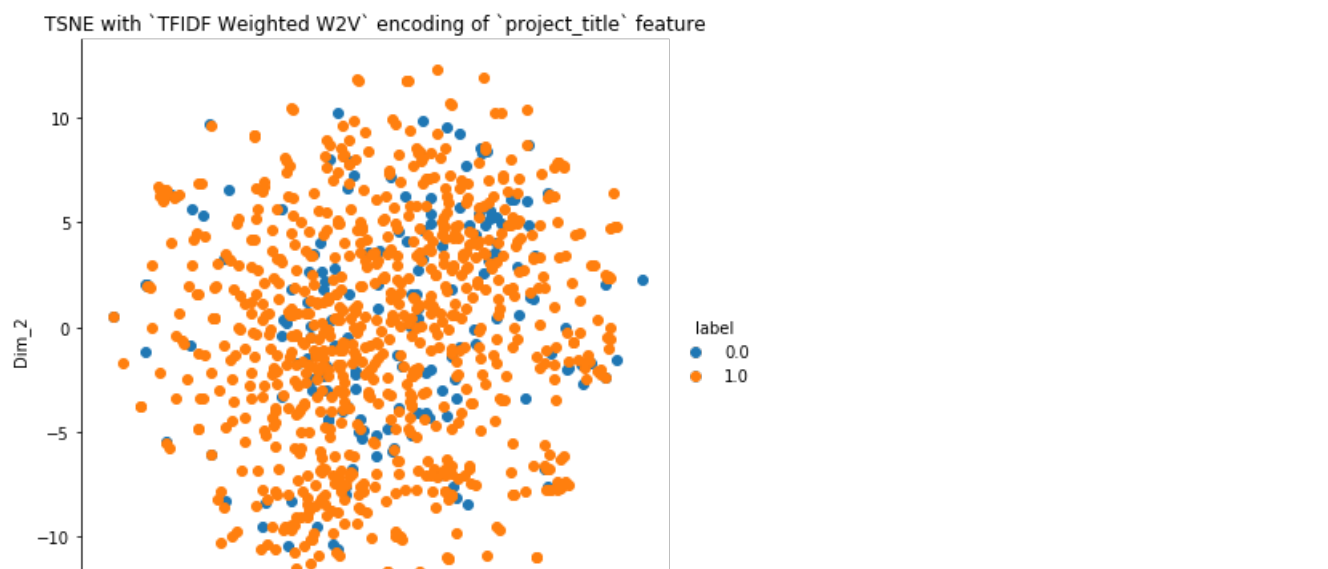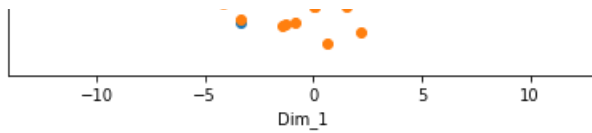
In [129]:

```
tsne_df_tfidf_w2v.shape
```

Out[129]:

```
(1000, 3)
```

In [130]:

```
# Ploting the result of tsne
sns.FacetGrid(tsne_df_tfidf_w2v, hue="label", size=6).map(plt.scatter, 'Dim_1', 'Dim_2').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

# Summary

1. Plot shows Projects Approved or not approved using TSNE with `TFIDF Weighted W2V` encoding of `project_title` feature
2. orange dots(1) indicate approved
3. blue dots(0) indicate not approved
4. Since blue points are suppressed under orange points we can not conclude anything about it

### 2.5 Summary

# Write few sentences about the results that you obtained and the observations you made.

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

1. The Objective of this Analysis is to predict if the project proposal submitted by a teacher to DonorsChoose.org will be approved or not, using text of project descriptions as well data's like teacher, school, state and resources.
2. Overall 85% projects are accepted and 15% are Rejected
3. Considering project approval by state wise, states which have submitted more projects have higher chances of approval and The states which have submited lesser projects have lesser chance of approval overall every state is having more than 80% approval rate
4. Considering project approval by prefixes The women contribute more in submitting projects in that married women have highest success rate of approval as Mrs. is the prefix having highest approval rate of 86%
5. Considering project approval by Grades, As the Grade increases the total number of projects submitted is decreased The highest approval rate is for Grades 3-5 with 85% approval rate. Less number of project submitted and approved for Grades 9-12 with approval rate 84%
6. Considering project approval rate based on project subject categories The joint subject categories like have higher approval rating, like history_civics literacy_language(89%),Literacy_language math_science (87%), Though the total projects are less in number but Warmth care_hunger subject category have highest approval rating of 93%, AppliedLearning Math_science category has submitted least projects among all and approval rating is less i.e 81%
7. Considering project approval rate based on individual subject categories Higest number of projects i.e. 52239 are submitted by Literacy and Language, then comes Math and Science, it's about 41421. Least number i.e 1388 are submitted by Warmth, Care and Hunger, then comes History and Civics, it's about 5914.
8. Considering project approval rate based on the number of words in Project Title. Most of the Projects(Around 20000) have 4 word in title, next comes 5 words for around 18000 projects, Least number of projects are found with 1 word, and 11 and more number of words
9. Considering project approval rate based on cost of projects, projects of higher cost are not approved compared to approved projects, and maximum cost of approved projects are within 10000
10. Considering project approval rate based on number of projects submitted by same teacher,The approval rate is high for same teachers submitting projects above 36 as project aproval rate is greater than 87%. Teachers who submitted projects for the first time are more in number 30014, among which 24652 are accepted, leading to project approval rate of 82%.
11. Considering project approval rate based on TSNE analysis with BOW, TFIDF, AVG Word2Vec, TFIDF weighted Word2Vec is much of help as analysis can't be done due to approved and not approved project points are overlapping.

In [ ]: