# Donors Choose

**By Aziz Presswala**

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 Dataset

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

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

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

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

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

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

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

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

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

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

## Notes on the Essay Data

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

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

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

In [1]:

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

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

# 1.1 Reading Data

In [47]:

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

In [48]:

```python
print("Number of data points in train data", project_data.shape)
print('-'*80)
print("The attributes of data :", project_data.columns.values)
print('-'*80)
print("The number of attributes in dataset :", len(project_data.columns.values))
```

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

In [49]:

```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[49]:

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

# 1.2 Data Analysis

In [8]:

```python
# this code is taken from
# https://matplotlib.org/gallery/pie_and_polar_charts/pie_and_donut_labels.html#sphx-glr-ga

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

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

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

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

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

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

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

plt.show()
```

Number of projects thar are approved for funding  92706 , ( 84.8583040421792
7 %)
Number of projects thar are not approved for funding  16542 , ( 15.141695957
820739 %)

Nmber of projects that are Accepted and not accepted

## 1.2.1 Univariate Analysis: School State

In [9]:

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

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

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

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

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

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

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

Project Proposals % of Acceptance Rate

In [10]:

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

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

In [11]:

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

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

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

In [12]:

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

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

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

    if top:
        temp = temp[0:top]

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

In [13]:

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



% of projects aproved state wise

```
   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
```

**Every state is having more than 80% success rate in approval**

## 1.2.2 Univariate Analysis: teacher_prefix

In [14]:

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



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

## 1.2.3 Univariate Analysis: project_grade_category

In [15]:

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



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

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

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

## 1.2.4 Univariate Analysis: project_subject_categories

In [16]:

```
catogories = list(project_data['project_subject_categories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/473019
# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
cat_list = []
for i in catogories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "
        if 'The' in j.split(): # this will split each of the catogory based on space "Math
            j=j.replace('The','') # if we have the words "The" we are going to replace it w
        j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math
        temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_') # we are replacing the & value into
    cat_list.append(temp.strip())
```

In [17]:

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

Out[17]:

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

In [18]:

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



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

In [19]:

```
# 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 [20]:

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


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

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



In [21]:

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

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

## 1.2.5 Univariate Analysis: project_subject_subcategories

In [22]:

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

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

sub_cat_list = []
for i in sub_catogories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "
        if 'The' in j.split(): # this will split each of the catogory based on space "Math
            j=j.replace('The','') # if we have the words "The" we are going to replace it w
        j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math
        temp +=j.strip()+" "+"#" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_')
    sub_cat_list.append(temp.strip())
```
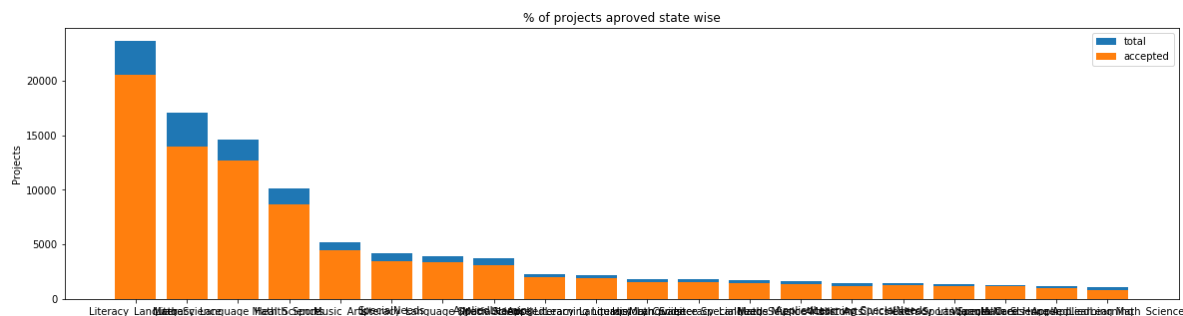
In [23]:

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

Out[23]:

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

In [24]:

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



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

```
==================================================
```

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

In [25]:

```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())
```
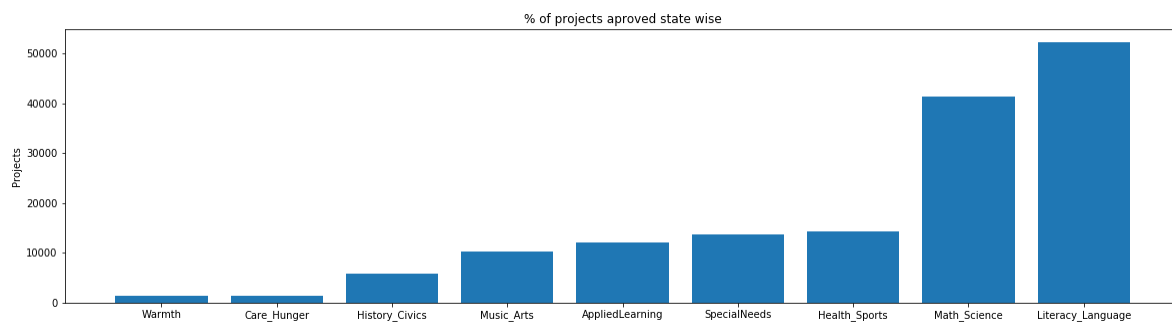
In [26]:

```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 [27]:

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

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

## 1.2.6 Univariate Analysis: Text features (Title)

In [28]:

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


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

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



In [29]:

```python
approved_word_count = project_data[project_data['project_is_approved']==1]['project_title']
approved_word_count = approved_word_count.values

rejected_word_count = project_data[project_data['project_is_approved']==0]['project_title']
rejected_word_count = rejected_word_count.values
```

In [30]:

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

In [31]:

```python
plt.figure(figsize=(10,3))
sns.distplot(approved_word_count, hist=False, label="Approved Projects")
sns.distplot(rejected_word_count, hist=False, label="Not Approved Projects")
plt.legend()
plt.show()
```



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

In [32]:

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

In [33]:

```python
#How to calculate number of words in a string in DataFrame: https://stackoverflow.com/a/374
word_count = project_data['essay'].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('Number of projects')
plt.xlabel('Number of words in each eassay')
plt.title('Words for each essay of the project')
plt.xticks(ind, list(word_dict.keys()))
plt.show()
```



In [34]:

```python
sns.distplot(word_count.values)
plt.title('Words for each essay of the project')
plt.xlabel('Number of words in each eassay')
plt.show()
```



In [35]:

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

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

In [36]:

```
# 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 title')
plt.grid()
plt.show()
```



In [37]:

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



## 1.2.8 Univariate Analysis: Cost per project

In [38]:

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

Out[38]:

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

In [39]:

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

Out[39]:

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

In [40]:

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

In [41]:

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

rejected_price = project_data[project_data['project_is_approved']==0]['price'].values
```
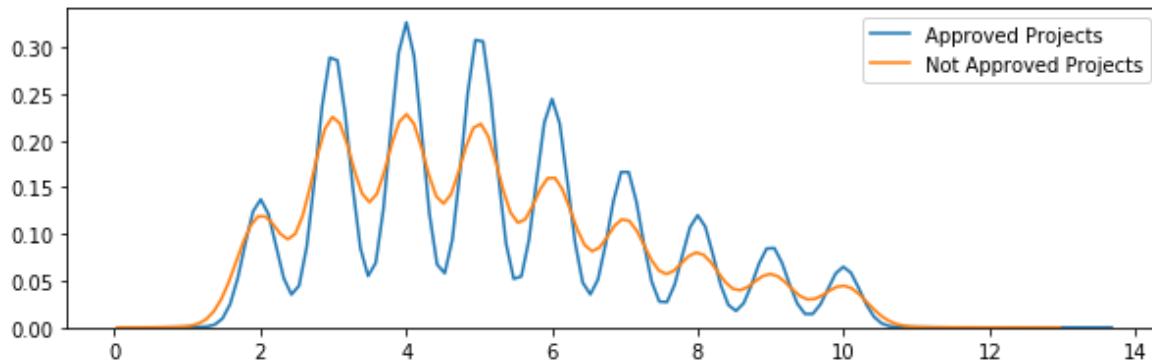
In [42]:

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



In [44]:

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

In [45]:

```python
# http://zetcode.com/python/prettytable/
from prettytable import 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(rejec
print(x)
```

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

## 1.2.9 Univariate Analysis: teacher_number_of_previously_posted_projects

In [62]:

```python
count = project_data['teacher_number_of_previously_posted_projects'].value_counts()
count_dict = dict(count)
count_dict = dict(sorted(count_dict.items(), key=lambda kv: kv[1]))


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

plt.ylabel('Number of occurences')
plt.xlabel('Number of previously posted projects')
plt.xticks(ind, list(count_dict.keys()))
plt.show()
```
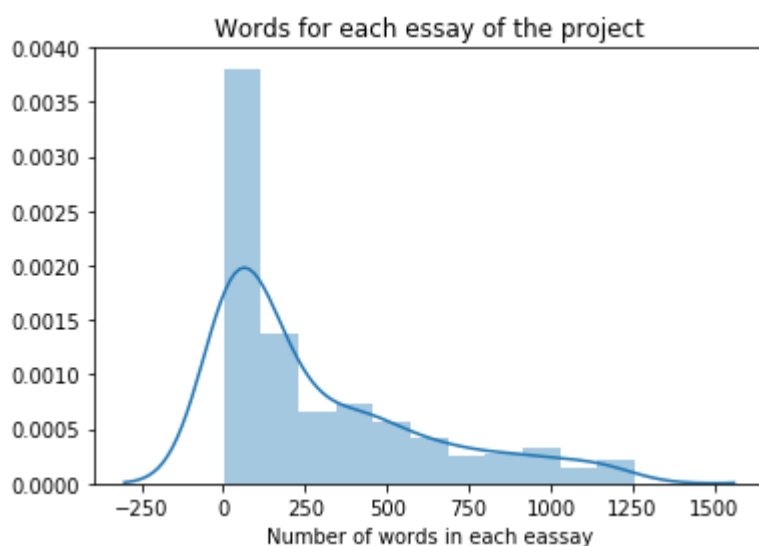


**Observation - Not interpretable**

In [56]:

```python
approved_number = project_data[project_data['project_is_approved']==1]['teacher_number_of_p

rejected_number = project_data[project_data['project_is_approved']==0]['teacher_number_of_p
```

In [57]:

```python
# https://glowingpython.blogspot.com/2012/09/boxplot-with-matplotlib.html
plt.boxplot([approved_number, rejected_number])
plt.xticks([1,2],('Approved Projects','Rejected Projects'))
plt.ylabel('Number of projects previously posted by a teacher')
plt.grid()
plt.show()
```

**Observation - Teachers who previously posted more projects have a higher probability of project approval**

In [58]:

```python
plt.figure(figsize=(10,3))
sns.distplot(approved_number, hist=False, label="Approved Projects")
sns.distplot(rejected_number, hist=False, label="Not Approved Projects")
plt.title('Number of projects posted previously per approved and not approved Projects')
plt.xlabel('Number of projects posted previously')
plt.legend()
plt.show()
```

In [63]:

```python
# http://zetcode.com/python/prettytable/
from prettytable import 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_number,i), 3), np.round(np.percentile(reje
print(x)
```

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

**Observation - As observed in box plots, Teachers who have submitted more projects previously, their projects are more likely to be accepted.**

## 1.2.10 Univariate Analysis: project_resource_summary

In [65]:

```python
#How to calculate number of words in a string in DataFrame: https://stackoverflow.com/a/374
word_count = project_data['project_resource_summary'].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('Number of projects')
plt.xlabel('Number of words in each project summary')
plt.title('Words for each summary of the project')
plt.xticks(ind, list(word_dict.keys()))
plt.show()
```



**Observation - Most project summary have words in range 11-20 ie.project_summary is generally short.**

In [68]:

```python
approved_word_count = project_data[project_data['project_is_approved']==1]['project_resourc
approved_word_count = approved_word_count.values

rejected_word_count = project_data[project_data['project_is_approved']==0]['project_resourc
rejected_word_count = rejected_word_count.values
```

In [69]:

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

Words for each essay of the project

**Observation - It is difficult to differentiate whether a project will be approved or not based on the project_summary feature because their distributions are almost similar**

## 1.2.11 Univariate Analysis: project_summary_numerical

It is a numerical feature that indiactes denoted a the count of numerical values present in the project_summary text

In [91]:

```python
project_summary_numerical = []
for sentance in tqdm(project_data['project_resource_summary']):
    sent1 = decontracted(sentance)
    sent1 = ' '.join(e for e in sent1.split() if e.isdigit())
    k=len(sent1)
    project_summary_numerical.append(k)
project_data['project_summary_num'] = project_summary_numerical
```

```
100%|████████████████████████████████████████████████████████████████|
| 109248/109248 [00:03<00:00, 35185.63it/s]
```

In [92]:

```
project_data['project_summary_num'].value_counts().plot.bar()
```

Out[92]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x1c6256984a8>
```



**Observation - most of the project_resource_summary do not contain any numerical values**

# 2. Preprocessing Categorical Features: project_grade_category

In [50]:

```
project_data['project_grade_category'].value_counts()
```

Out[50]:

```
Grades PreK-2    44225
Grades 3-5       37137
Grades 6-8       16923
Grades 9-12      10963
Name: project_grade_category, dtype: int64
```

In [51]:

```python
# https://stackoverflow.com/questions/36383821/pandas-dataframe-apply-function-to-column-st
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace
project_data['project_grade_category'] = project_data['project_grade_category'].str.lower()
project_data['project_grade_category'].value_counts()
```

Out[51]:

```
grades_prek_2      44225
grades_3_5         37137
grades_6_8         16923
grades_9_12        10963
Name: project_grade_category, dtype: int64
```

# 3. Preprocessing Categorical Features: project_subject_categories

In [52]:

```python
project_data['project_subject_categories'].value_counts()
```

Out[52]:

```
Literacy & Language                         23655
Math & Science                              17072
Literacy & Language, Math & Science         14636
Health & Sports                             10177
Music & The Arts                             5180
Special Needs                                4226
Literacy & Language, Special Needs           3961
Applied Learning                             3771
Math & Science, Literacy & Language          2289
Applied Learning, Literacy & Language        2191
History & Civics                             1851
Math & Science, Special Needs                1840
Literacy & Language, Music & The Arts        1757
Math & Science, Music & The Arts             1642
Applied Learning, Special Needs              1467
History & Civics, Literacy & Language        1421
Health & Sports, Special Needs               1391
Warmth, Care & Hunger                        1309
```

In [53]:

```
project_data['project_subject_categories'] = project_data['project_subject_categories'].str
project_data['project_subject_categories'] = project_data['project_subject_categories'].str
project_data['project_subject_categories'] = project_data['project_subject_categories'].str
project_data['project_subject_categories'] = project_data['project_subject_categories'].str
project_data['project_subject_categories'] = project_data['project_subject_categories'].str
project_data['project_subject_categories'].value_counts()
```

Out[53]:

```
literacy_language                  23655
math_science                       17072
literacy_language_math_science     14636
health_sports                      10177
music_arts                          5180
specialneeds                        4226
literacy_language_specialneeds      3961
appliedlearning                     3771
math_science_literacy_language      2289
appliedlearning_literacy_language   2191
history_civics                      1851
math_science_specialneeds           1840
literacy_language_music_arts        1757
math_science_music_arts             1642
appliedlearning_specialneeds        1467
history_civics_literacy_language    1421
health_sports_specialneeds          1391
warmth care hunger                  1309
```

# 4. Preprocessing Categorical Features: teacher_prefix

In [54]:

```
project_data['teacher_prefix'].value_counts()
```

Out[54]:

```
Mrs.       57269
Ms.        38955
Mr.        10648
Teacher     2360
Dr.           13
Name: teacher_prefix, dtype: int64
```

In [55]:

```
# check if we have any nan values are there
print(project_data['teacher_prefix'].isnull().values.any())
print("number of nan values",project_data['teacher_prefix'].isnull().values.sum())
```

```
True
number of nan values 3
```

> numebr of missing values are very less in number, we can replace it with Mrs. as most of the projects are submitted by Mrs.

In [56]:

```
project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('Mrs.')
```

In [57]:

```
project_data['teacher_prefix'].value_counts()
```

Out[57]:

```
Mrs.        57272
Ms.         38955
Mr.         10648
Teacher      2360
Dr.            13
Name: teacher_prefix, dtype: int64
```

In [58]:

```
project_data['teacher_prefix'] = project_data['teacher_prefix'].str.replace('.','')
project_data['teacher_prefix'] = project_data['teacher_prefix'].str.lower()
project_data['teacher_prefix'].value_counts()
```

Out[58]:

```
mrs        57272
ms         38955
mr         10648
teacher     2360
dr            13
Name: teacher_prefix, dtype: int64
```

# 5. Preprocessing Categorical Features: project_subject_subcategories

In [59]:

```
project_data['project_subject_subcategories'].value_counts()
```

Out[59]:

```
Literacy                              9486
Literacy, Mathematics                 8325
Literature & Writing, Mathematics     5923
Literacy, Literature & Writing        5571
Mathematics                           5379
Literature & Writing                  4501
Special Needs                         4226
Health & Wellness                     3583
Applied Sciences, Mathematics         3399
Applied Sciences                      2492
Literacy, Special Needs               2440
Gym & Fitness, Health & Wellness      2264
ESL, Literacy                         2234
Visual Arts                           2217
Music                                 1472
Warmth, Care & Hunger                 1309
Literature & Writing, Special Needs   1306
Gym & Fitness                         1195
```

In [60]:

```
project_data['project_subject_subcategories'] = project_data['project_subject_subcategories
project_data['project_subject_subcategories'] = project_data['project_subject_subcategories
project_data['project_subject_subcategories'] = project_data['project_subject_subcategories
project_data['project_subject_subcategories'] = project_data['project_subject_subcategories
project_data['project_subject_subcategories'] = project_data['project_subject_subcategories
project_data['project_subject_subcategories'].value_counts()
```

Out[60]:

```
literacy                              9486
literacy_mathematics                  8325
literature_writing_mathematics        5923
literacy_literature_writing           5571
mathematics                           5379
literature_writing                    4501
specialneeds                          4226
health_wellness                       3583
appliedsciences_mathematics           3399
appliedsciences                       2492
literacy_specialneeds                 2440
gym_fitness_health_wellness           2264
esl_literacy                          2234
visualarts                            2217
music                                 1472
warmth_care_hunger                    1309
literature_writing_specialneeds       1306
gym_fitness                           1195
```

# 6. Preprocessing Categorical Features: school_state

In [61]:

```
project_data['school_state'].value_counts()
```

Out[61]:

```
CA    15388
TX     7396
NY     7318
FL     6185
NC     5091
IL     4350
GA     3963
SC     3936
MI     3161
PA     3109
IN     2620
MO     2576
OH     2467
LA     2394
MA     2389
WA     2334
OK     2276
NJ     2237
```

convert all of them into small letters

In [62]:

```
project_data['school_state'] = project_data['school_state'].str.lower()
project_data['school_state'].value_counts()
```

Out[62]:

```
ca    15388
tx     7396
ny     7318
fl     6185
nc     5091
il     4350
ga     3963
sc     3936
mi     3161
pa     3109
in     2620
mo     2576
oh     2467
la     2394
ma     2389
wa     2334
ok     2276
nj     2237
```

# 7. Preprocessing Categorical Features: project_title

In [63]:

```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 [64]:

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

In [65]:

```python
project_data['project_title'].head(5)
```

Out[65]:

```
0        Educational Support for English Learners at Home
1                   Wanted: Projector for Hungry Learners
2      Soccer Equipment for AWESOME Middle School Stu...
3                                  Techie Kindergarteners
4                                   Interactive Math Tools
Name: project_title, dtype: object
```

In [66]:

```python
print("printing some random reviews")
print(9, project_data['project_title'].values[9])
print(34, project_data['project_title'].values[34])
print(147, project_data['project_title'].values[147])
```

```
printing some random reviews
9 Just For the Love of Reading--\r\nPure Pleasure
34 \"Have A Ball!!!\"
147 Who needs a Chromebook?\r\nWE DO!!
```

In [67]:

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

In [68]:

```python
preprocessed_titles = preprocess_text(project_data['project_title'].values)
```

```
100%|████████████████████████████████████████████████████████████████|
| 109248/109248 [00:04<00:00, 24157.00it/s]
```

In [69]:

```python
print("printing some random reviews")
print(9, preprocessed_titles[9])
print(34, preprocessed_titles[34])
print(147, preprocessed_titles[147])
```

```
printing some random reviews
9 love reading pure pleasure
34 ball
147 needs chromebook
```

In [84]:

```python
project_data['project_title'] = preprocessed_titles
```

# 8. Preprocessing Categorical Features: essay

In [70]:

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

In [71]:

```python
print("printing some random essay")
print(19, project_data['essay'].values[19])
print('-'*50)
print(134, project_data['essay'].values[134])
print('-'*50)
print(247, project_data['essay'].values[247])
```

printing some random essay
9 Over 95% of my students are on free or reduced lunch.  I have a few who
are homeless, but despite that, they come to school with an eagerness to l
earn.  My students are inquisitive eager learners who  embrace the challen
ge of not having great books and other resources  every day.  Many of them
are not afforded the opportunity to engage with these big colorful pages o
f a book on a regular basis at home and they don't travel to the public li
brary.  \r\nIt is my duty as a teacher to do all I can to provide each stu
dent an opportunity to succeed in every aspect of life. \r\nReading is Fun
damental! My students will read these books over and over again while boos
ting their comprehension skills. These books will be used for read alouds,
partner reading and for Independent reading. \r\nThey will engage in readi
ng to build their \"Love for Reading\" by reading for pure enjoyment. They
will be introduced to some new authors as well as some old favorites. I wa
nt my students to be ready for the 21st Century and know the pleasure of h
olding a good hard back book in hand. There's nothing like a good book to
read!  \r\nMy students will soar in Reading, and more because of your cons
ideration and generous funding contribution. This will help build stamina
and prepare for 3rd grade. Thank you so much for reading our proposal!nann
an

In [72]:

```python
preprocessed_essays = preprocess_text(project_data['essay'].values)
```

```
100%|████████████████████████████████████████████████████████|
| 109248/109248 [01:42<00:00, 1071.05it/s]
```

In [73]:

```python
print("printing some random essay")
print(19, preprocessed_essays[19])
print('-'*50)
print(134, preprocessed_essays[134])
print('-'*50)
print(247, preprocessed_essays[247])
```

```
printing some random essay
19 apart urban district many students come financially disadvantaged homes
4th 5th grade students face many challenges classroom class work build com
munity learners encourages hard work perseverance many challenges students
face also important child feel safe positive learning environment strive f
ind creative ways create atmosphere children teach many students struggle
sit still extended periods time students asked read work independent activ
ties several time day students need requested opportunity move completing
activities hokki stools provide students quiet option continue move compet
ing necessary school work classroom use stools seating option read aloud m
ini lesson time well 7 stools classroom allow 30 students opportunity use
stools throughout day keep active build core muscles continue great learni
ng engage every day nannan
--------------------------------------------------
134 students unique wonderful come belong vibrant diverse urban community
whose population close 300 000 students ages 5 12 members self contained s
pecial education class perspective schools students identified one 14 disa
bilities named individuals disabilities education act idea majority autism
intellectual disabilities multiple disabilities integral part school distr
```

In [85]:

```python
project_data['essay'] = preprocessed_essays
```

# 9. Preprocessing Numerical Values: price

In [74]:

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

Out[74]:

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

In [75]:

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

In [76]:

```python
project_data['price'].head()
```

Out[76]:

```
0    154.60
1    299.00
2    516.85
3    232.90
4     67.98
Name: price, dtype: float64
```

# 9.1 applying StandardScaler

In [77]:

```python
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaler.fit(project_data['price'].values.reshape(-1, 1))
project_data['std_price']=scaler.transform(project_data['price'].values.reshape(-1, 1) )
```

In [78]:

```python
project_data['std_price'].head()
```

Out[78]:

```
0   -0.390533
1    0.002396
2    0.595191
3   -0.177469
4   -0.626236
Name: std_price, dtype: float64
```

# 9.2 applying MinMaxScaler

In [79]:

```python
from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()
scaler.fit(project_data['price'].values.reshape(-1, 1))
project_data['nrm_price']=scaler.transform(project_data['price'].values.reshape(-1, 1))
```

In [80]:

```python
project_data['nrm_price'].head()
```

Out[80]:

```
0    0.015397
1    0.029839
2    0.051628
3    0.023228
4    0.006733
Name: nrm_price, dtype: float64
```

# 10. Preprocessing Numerical Values: project_summary_numerical

In [81]:

```python
# we will extract numerical digits from the project resource summary
project_summary_numerical = []
for sentance in tqdm(project_data['project_resource_summary']):
    sent1 = decontracted(sentance)
    sent1 = ' '.join(e for e in sent1.split() if e.isdigit())
    k=len(sent1)
    project_summary_numerical.append(k)

project_data["project_summary_numerical"] = project_summary_numerical
```

```
100%|████████████████████████████████████████████████|
| 109248/109248 [00:02<00:00, 44160.95it/s]
```

In [82]:

```python
project_data['project_summary_numerical'].value_counts()
```

Out[82]:

```
0     97975
1      6004
2      2834
3      1187
4       449
5       398
7       124
6       115
8        50
9        43
11       20
14       16
10       13
12        6
13        5
15        3
19        2
16        2
```

## Final Data after Preprocessing

In [86]:

```
project_data.head()
```

Out[86]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_ |
|---|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | mrs | in | |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | mr | fl | |
| 2 | 21895 | p182444 | 3465aaf82da834c0582ebd0ef8040ca0 | ms | az | |
| 3 | 45 | p246581 | f3cb9bffbba169bef1a77b243e620b60 | mrs | ky | |
| 4 | 172407 | p104768 | be1f7507a41f8479dc06f047086a39ec | mrs | tx | |

5 rows × 23 columns

In [87]:

```
project_data.shape
```

Out[87]:

```
(109248, 23)
```

In [88]:

```
# saving the preprocessed dataset
project_data.to_csv('preprocessed_data.csv')
```