

```

%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

!pip install chart_studio
import chart_studio.plotly as py

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

```

```

Requirement already satisfied: chart_studio in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from chart_studio)
Requirement already satisfied: retrying>=1.3.3 in /usr/local/lib/python3.7/dist-packages (from chart_studio)
Requirement already satisfied: plotly in /usr/local/lib/python3.7/dist-packages (from chart_studio)
Requirement already satisfied: requests in /usr/local/lib/python3.7/dist-packages (from chart_studio)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/dist-packages (from chart_studio)
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-packages (from chart_studio)
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/dist-packages (from chart_studio)
Requirement already satisfied: urllib3!=1.25.0,!1.25.1,<1.26,>=1.21.1 in /usr/local,

```

```
from google.colab import files
files= files.upload()
```

Choose Files

No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving train_data.csv to train_data (1).csv

```
project_data = pd.read_csv('train_data (1).csv')
```

```
#Load the Resources File
from google.colab import files
files= files.upload()
```

Choose Files

No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving resources.csv to resources.csv

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

```
project_data.head(5)
```

	Unnamed: 0		id	teacher_id	teacher_prefix	school_state
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

```
resource_data
```

	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
2	p069063	Cory Stories: A Kid's Book About Living With Adhd	1	8.45
3	p069063	Dixon Ticonderoga Wood-Cased #2 HB Pencils, Bo...	2	13.59
4	p069063	EDUCATIONAL INSIGHTS FLUORESCENT LIGHT FILTERS...	3	24.95
...
1541267	p031981	AmazonBasics 9 Volt Everyday Alkaline Batterie...	1	9.99
1541268	p031981	AmazonBasics AAA Performance Alkaline Batterie...	1	6.99
1541269	p031981	Black Electrical Tape (GIANT 3 PACK) Each Roll...	6	8.99
1541270	n031981	Flormoon DC Motor Mini Electric Motor 0.5-3V 1...	2	8.14

```
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_sta
'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']

```
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']

	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

```
# this code is taken from
# https://matplotlib.org/gallery/pie\_and\_polar\_charts/pie\_and\_donut\_labels.html#sphx-glr-g
```

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

```

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)

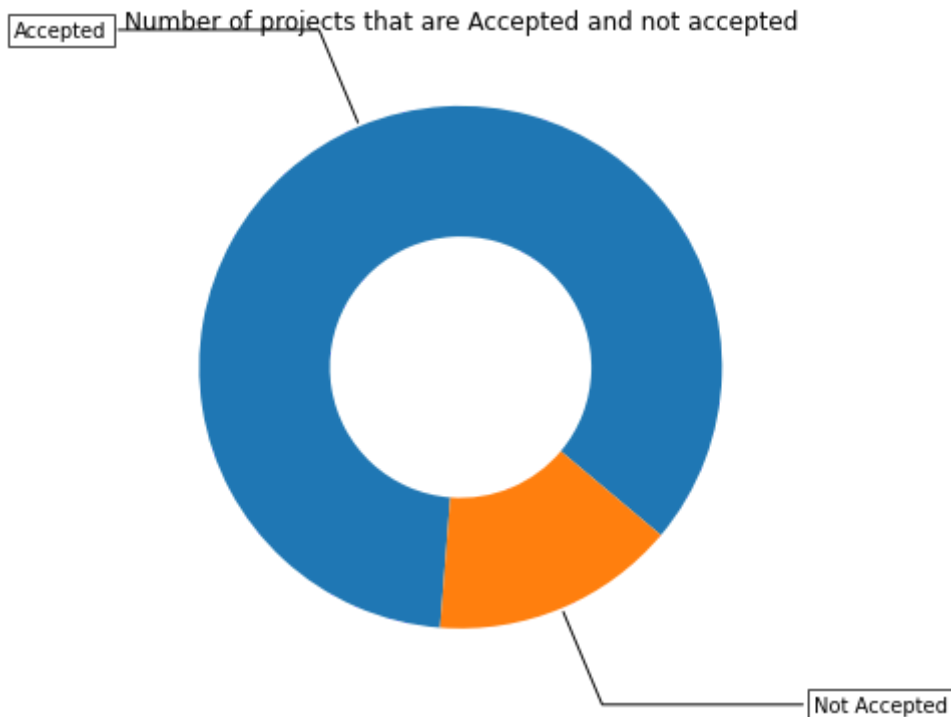
ax.set_title("Number of projects that are Accepted and not accepted")

plt.show()

```



Number of projects that are approved for funding 92706 , (84.85830404217927 %)
 Number of projects that are not approved for funding 16542 , (15.141695957820739 %)



▼ 1.2.1 Univariate Analysis: School State

Pandas dataframe grouby count, mean: <https://stackoverflow.com/a/19385591/4084039>

```
temp = pd.DataFrame(project_data.groupby("school_state")["project_is_approved"].apply(np.m
# 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')
```

```
# 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

#stacked bar plots matplotlib: https://matplotlib.org/gallery/lines_bars_and_markers/bar_s

```
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]), ('rejected', 'accepted'))
```

```
    plt.show()
```

```
def univariate_barplots(data, col1, col2='project_is_approved', top=False):
```

```
    # Count number of zeros in dataframe python: https://stackoverflow.com/a/51540521/4084
```

```
    temp = pd.DataFrame(project_data.groupby(col1)[col2].agg(lambda x: x.eq(1).sum())).res
```

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

```
    temp['total'] = pd.DataFrame(project_data.groupby(col1)[col2].agg([('total', 'count')]))
```

```
    temp['Avg'] = pd.DataFrame(project_data.groupby(col1)[col2].agg([('Avg', 'mean')])).res
```

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

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

```
if top:
```

```
    temp = temp[0:top]
```

```
stack_plot(temp, xtick=col1, col2=col2, col3='total')
```

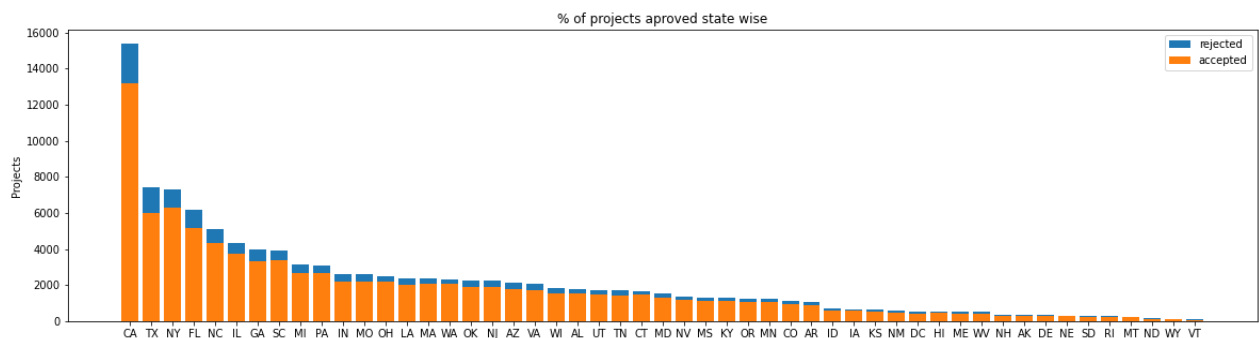
```
print(temp.head(5))
```

```
print("="*50)
```

```
print(temp.tail(5))
```

▼ 1.2.2 Univariate Analysis: teacher_prefix

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



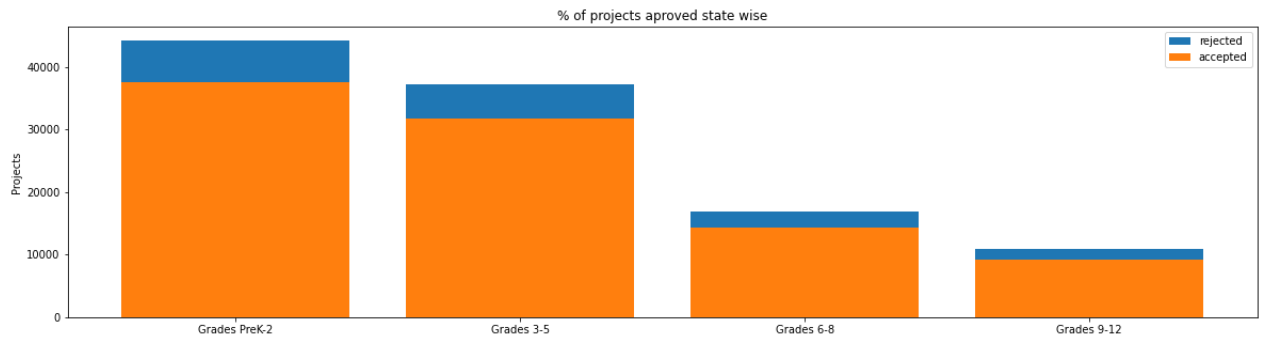
	school_state	project_is_approved	total	Avg
4	CA	13205	15388	0.858136
43	TX	6014	7396	0.813142
34	NY	6291	7318	0.859661
9	FL	5144	6185	0.831690
27	NC	4353	5091	0.855038

```
=====
```

	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

▼ 1.2.3 Univariate Analysis: project_grade_category

```
univariate_barplots(project_data, 'project_grade_category', 'project_is_approved', top=Fa1
```



	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

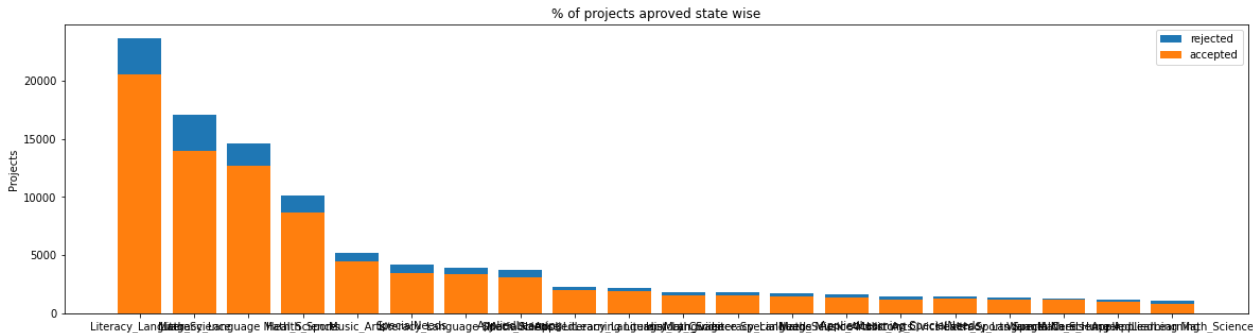
```
catogories = list(project_data['project_subject_categories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301

# 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
        j = j.replace(' ','') # we are placing 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())

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


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

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



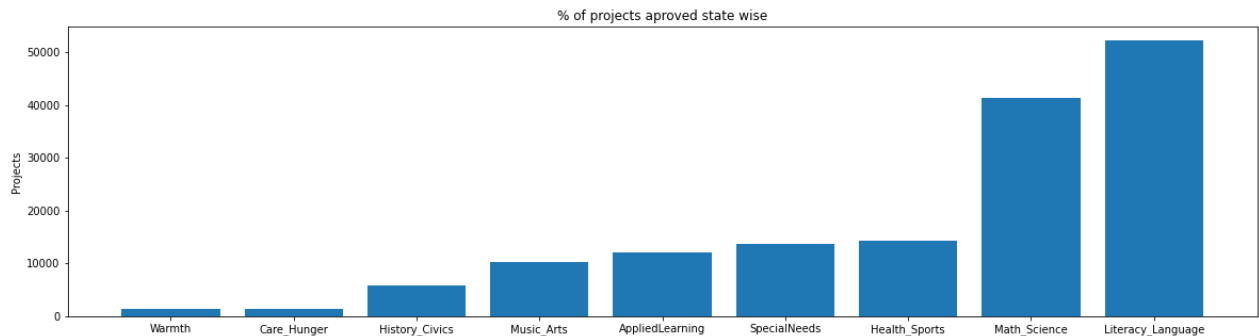
	clean_categories	project_is_approved	total	Avg
24	Literacy_Language	20520	23655	0.867470
32	Math_Science	13991	17072	0.819529
28	Literacy_Language Math_Science	12725	14636	0.869432
8	Health_Sports	8640	10177	0.848973
40	Music_Arts	4429	5180	0.855019
=====				
	clean_categories	project_is_approved	total	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

```
# 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())

# 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()
```



```
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

```
sub_catogories = list(project_data['project_subject_subcategories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301
```

```
# 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
```

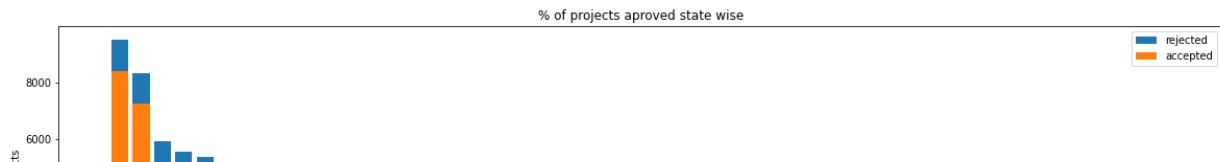
```
sub_cat_list = []
for i in sub_catogories:
    temp = ""
    # consider we have text like this "Math & Science Health Care & Hunger"
```

```
# 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 category based on space "Math
        j=j.replace('The','') # if we have the words "The" we are going to replace it
    j = j.replace(' ','') # we are placing 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())
```

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

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

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



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())
```

```
my_counter.most_common(10)
```

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))
```

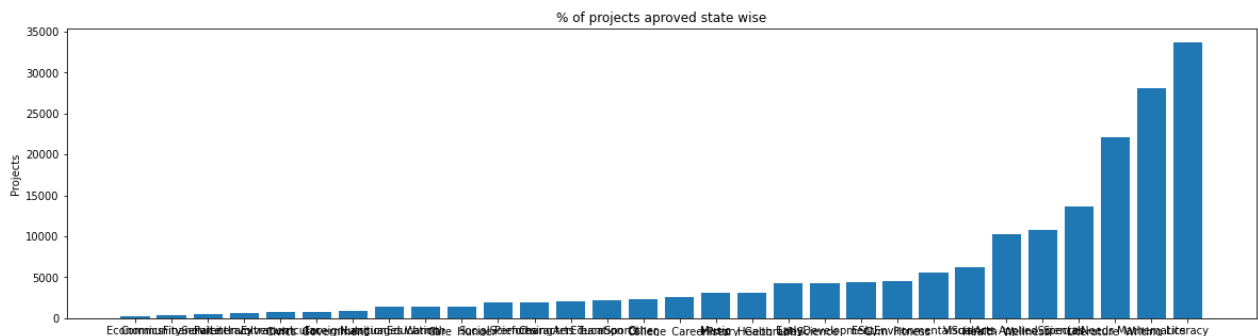
```
p1 = plt.bar(ind, list(sorted_sub_cat_dict.values()))
```

```
plt.ylabel('Projects')
```

```
plt.title('% of projects aproved state wise')
```

```
plt.xticks(ind, list(sorted_sub_cat_dict.keys()))
```

```
plt.show()
```



```
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)

#How to calculate number of words in a string in DataFrame: <https://stackoverflow.com/a/37>

```
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))
```

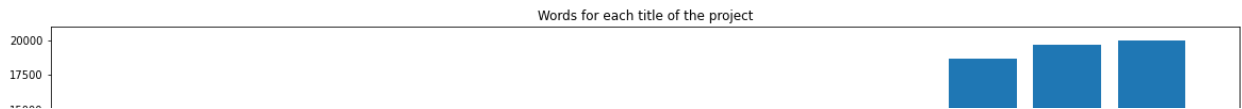
```
p1 = plt.bar(ind, list(word_dict.values()))
```

```
plt.ylabel('Numeber of projects')
```

```
plt.title('Words for each title of the project')
```

```
plt.xticks(ind, list(word_dict.keys()))
```

```
plt.show()
```

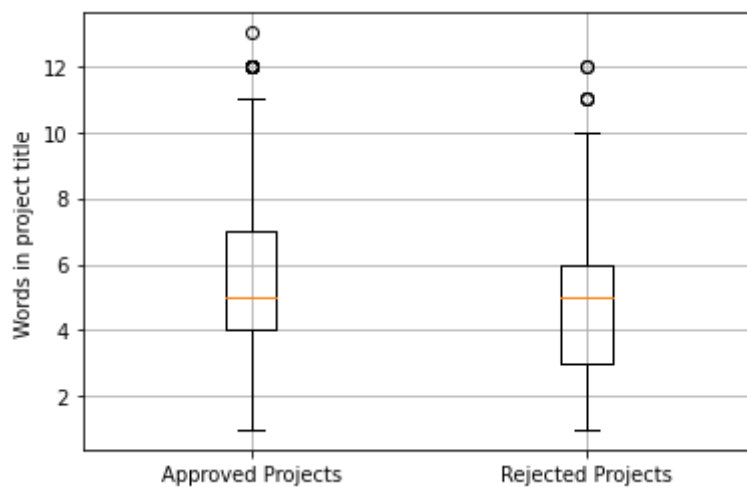


```
approved_word_count = project_data[project_data['project_is_approved']==1]['project_title']
approved_word_count = approved_word_count.values
```

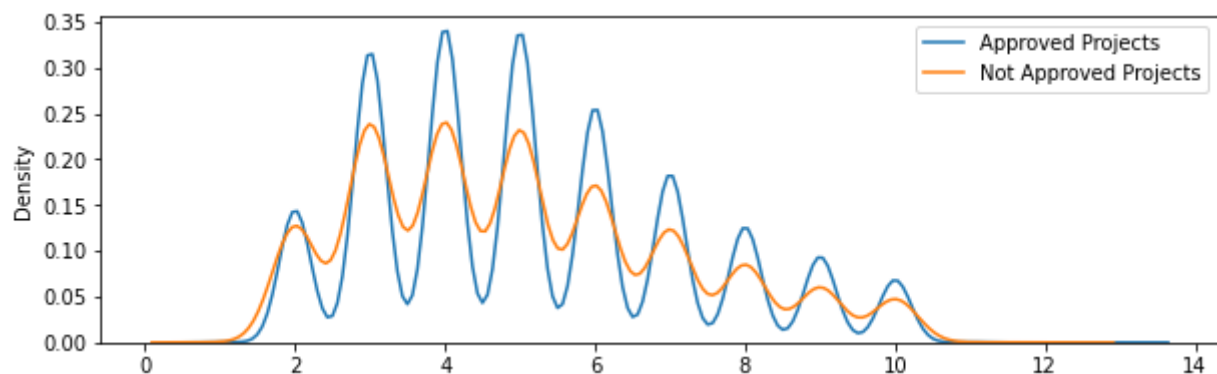
```
rejected_word_count = project_data[project_data['project_is_approved']==0]['project_title']
rejected_word_count = rejected_word_count.values
```

```
# 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()
```



```
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)

```
# merge two column text dataframe:
```

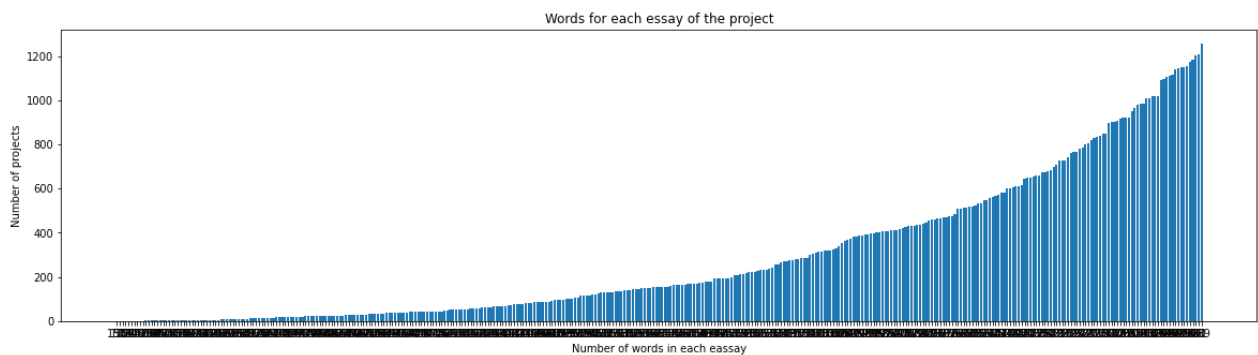
```
project_data["essay"] = project_data["project_essay_1"].map(str) + \
    project_data["project_essay_2"].map(str) + \
    project_data["project_essay_3"].map(str) + \
    project_data["project_essay_4"].map(str)
```

#How to calculate number of words in a string in DataFrame: <https://stackoverflow.com/a/37>

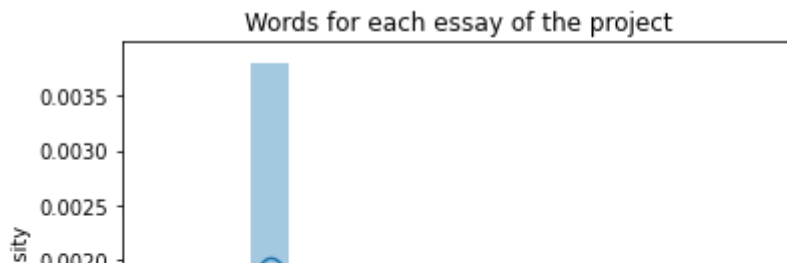
```
word_count = project_data['essay'].str.split().apply(len).value_counts()
word_dict = dict(word_count)
word_dict = dict(sorted(word_dict.items(), key=lambda kv: kv[1]))
```

```
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()
```



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

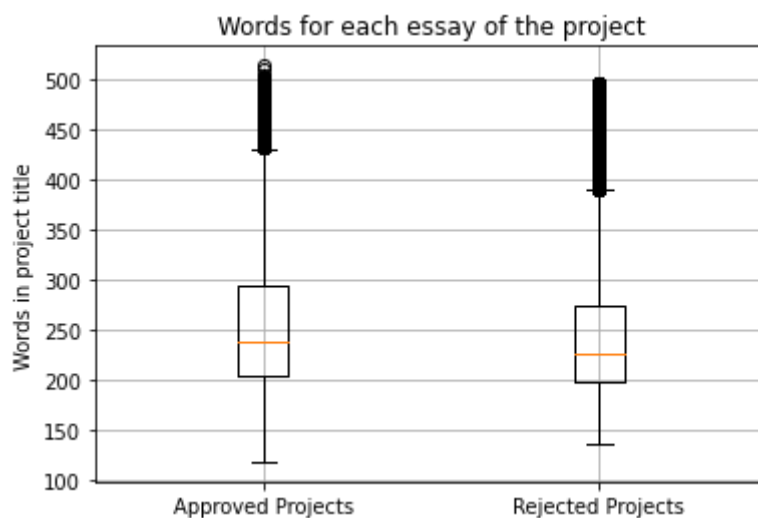


```
approved_word_count = project_data[project_data['project_is_approved']==1]['essay'].str.sp
approved_word_count = approved_word_count.values
```

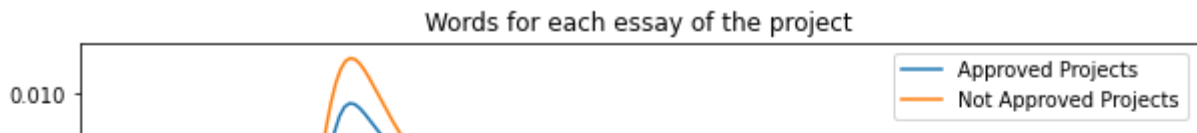
```
rejected_word_count = project_data[project_data['project_is_approved']==0]['essay'].str.sp
rejected_word_count = rejected_word_count.values
```

```
-250    0    250    500    750    1000    1250    1500
```

```
# 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()
```



```
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



```
# 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-g
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')
```

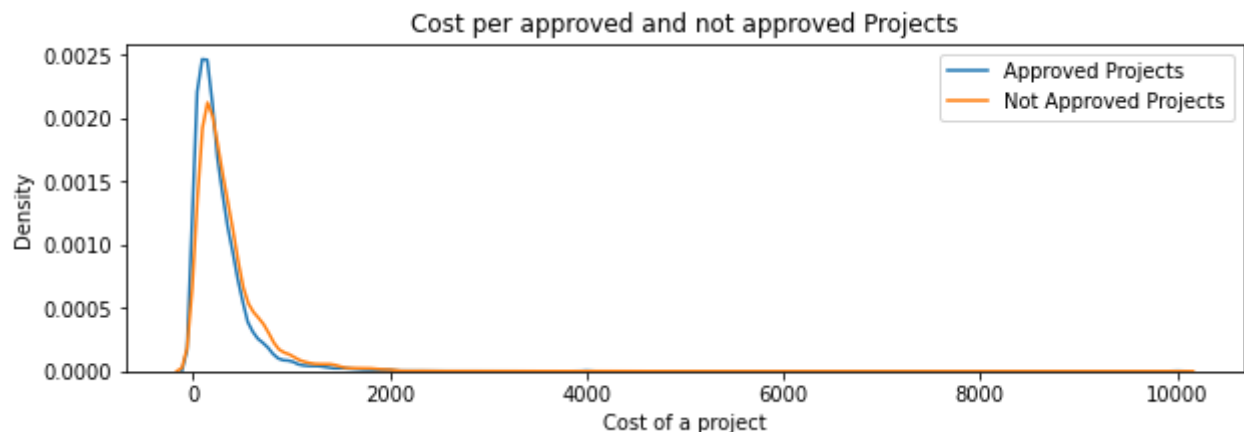
```
approved_price = project_data[project_data['project_is_approved']==1]['price'].values
```

```
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('Words in project title')
plt.grid()
plt.show()
```



```
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()
```



```
# 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(rejected_price,i), 3)])
print(x)
```

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

	85		479.0		618.276	
	90		593.11		739.356	
	95		801.598		992.486	
	100		9999.0		9999.0	
+-----+-----+-----+-----+						

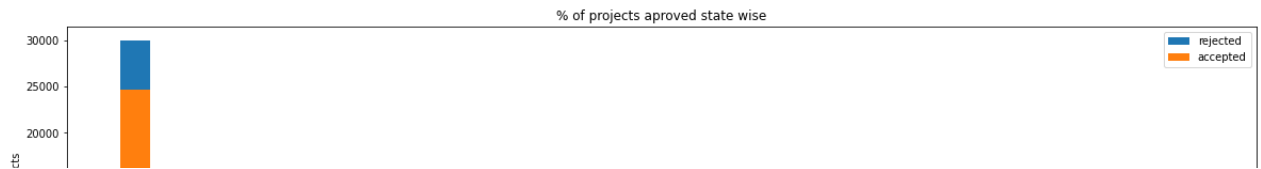
1.2.9 Univariate Analysis:

teacher_number_of_previously_posted_projects

```
overview= project_data["teacher_number_of_previously_posted_projects"].describe()
overview
```

```
count      109248.000000
mean         11.153165
std          27.777154
min           0.000000
25%           0.000000
50%           2.000000
75%           9.000000
max          451.000000
Name: teacher_number_of_previously_posted_projects, dtype: float64
```

```
univariate_barplots(project_data, 'teacher_number_of_previously_posted_projects', 'project
```



https://mode.com/example-gallery/python_histogram/

```
from matplotlib.ticker import StrMethodFormatter
```

```
ax = project_data.hist(column='teacher_number_of_previously_posted_projects', bins=25, gri
```

```
ax = ax[0]
```

```
for x in ax:
```

```
    # Despine
```

```
    x.spines['right'].set_visible(False)
```

```
    x.spines['top'].set_visible(False)
```

```
    x.spines['left'].set_visible(False)
```

```
    # Switch off ticks
```

```
    x.tick_params(axis="both", which="both", bottom="off", top="off", labelbottom="on", le
```

```
    # Draw horizontal axis lines
```

```
    vals = x.get_yticks()
```

```
    for tick in vals:
```

```
        x.axhline(y=tick, linestyle='dashed', alpha=0.4, color='#eeeeee', zorder=1)
```

```
    # Remove title
```

```
    x.set_title("")
```

```
    # Set x-axis label
```

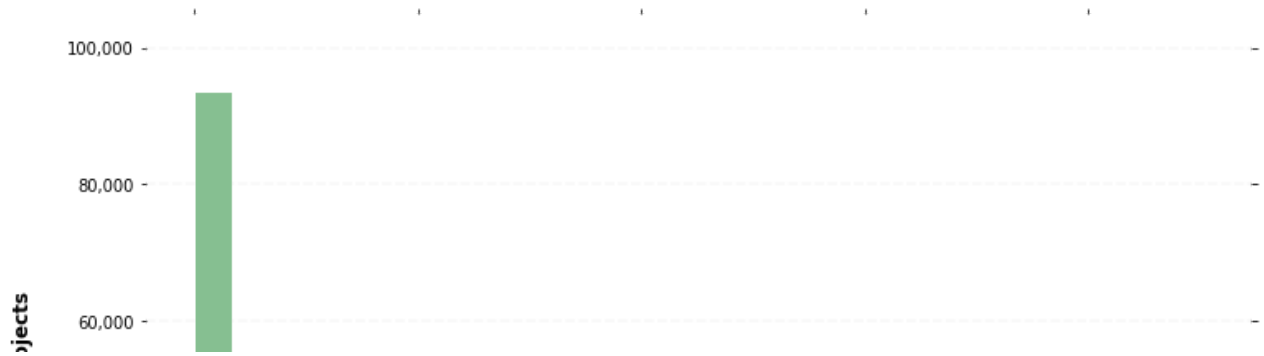
```
    x.set_xlabel("No. of Teachers", labelpad=20, weight='bold', size=12)
```

```
    # Set y-axis label
```

```
    x.set_ylabel("No. of Projects", labelpad=20, weight='bold', size=12)
```

```
    # Format y-axis label
```

```
    x.yaxis.set_major_formatter(StrMethodFormatter('{x:,g}'))
```



```
previous_projects_approved = project_data[project_data['project_is_approved']==1]['teacher_id']
```

```
previous_projects_rejected = project_data[project_data['project_is_approved']==0]['teacher_id']
```

```
# 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(previous_projects_approved,i), 3), np.round(np.percentile(previous_projects_rejected,i), 3)])
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

▼ 1.2.10 Univariate Analysis: project_resource_summary

1. Please do this by yourself

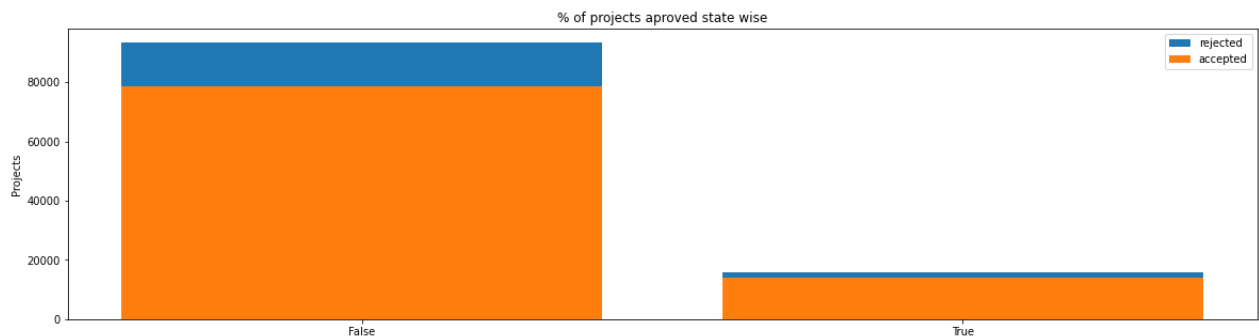
2. Check the presence of the numerical digits in the project_resource_summary effects the acceptance of the project
3. If you feel like it will helpfull in the classification, please include in the further process or you can ignore it.

```
#checking number in a string:
import re
def hasNumbers(inputString):
    return bool(re.search(r'\d', inputString)) #https://stackoverflow.com/questions/19859
```

```
summary = project_data['project_resource_summary'].values
digits_in_summary = []
for sent in summary:
    numb = hasNumbers(sent)
    digits_in_summary.append(numb)
```

```
#adding the digits_in_summary column in our dataframe
se = pd.Series(digits_in_summary)
project_data['digits_in_summary'] = se.values
```

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



	digits_in_summary	project_is_approved	total	Avg
0	False	78616	93492	0.840885
1	True	14090	15756	0.894263

=====

	digits_in_summary	project_is_approved	total	Avg
0	False	78616	93492	0.840885
1	True	14090	15756	0.894263

False = Summary without Digits, True = Summary with Digits

