

DonorsChoose

DonorsChoose.org receives hundreds of thousands of project proposals each year for classroom projects in need of funding. Right now, a large number of volunteers is needed to manually screen each submission before it's approved to be posted on the DonorsChoose.org website.

Next year, DonorsChoose.org expects to receive close to 500,000 project proposals. As a result, there are three main problems they need to solve:

- How to scale current manual processes and resources to screen 500,000 projects so that they can be posted as quickly and as efficiently as possible
- How to increase the consistency of project vetting across different volunteers to improve the experience for teachers
- How to focus volunteer time on the applications that need the most assistance

The goal of the competition is to predict whether or not a DonorsChoose.org project proposal submitted by a teacher will be approved, using the text of project descriptions as well as additional metadata about the project, teacher, and school. DonorsChoose.org can then use this information to identify projects most likely to need further review before approval.

About the DonorsChoose Data Set

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

Feature		Desc
project_id		A unique identifier for the proposed project. Example: p0
		Title of the project. Example:
project_title	• •	Art Will Make You H First Grad
		Grade level of students for which the project is targeted. One of the following enumerated values
project_grade_category	• • • •	Grades P
		Grade
		Grade
		Grades
project_subject_categories	• • • • • • • • • •	One or more (comma-separated) subject categories for the project from the following enumerated list of values
		Applied Learning
		Care & Health
		Health & Safety
		History & Culture
		Literacy & Language
		Math & Science
		Music & The Arts
		Special Education
		World Languages
	• •	Example:
		Music & The Arts, Literacy & Language, Math & Science

Feature	Desc
<code>school_state</code>	State where school is located (Two-letter U.S. postal code) (https://en.wikipedia.org/wiki/List of U.S. state abbreviations#Postal codes) Example: CA
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. Example: Literature & Writing, Social Sciences <ul style="list-style-type: none"> Literature & Writing Literature & Writing, Social Sciences
<code>project_resource_summary</code>	An explanation of the resources needed for the project. Example: My students need hands on literacy materials to meet sensory needs!<
<code>project_essay_1</code>	First application
<code>project_essay_2</code>	Second application
<code>project_essay_3</code>	Third application
<code>project_essay_4</code>	Fourth application
<code>project_submitted_datetime</code>	Datetime when project application was submitted. Example: 2016-01-01 12:43:50
<code>teacher_id</code>	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c
<code>teacher_prefix</code>	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> Teacher Assistant Teacher Paraprofessional Volunteer Other
<code>teacher_number_of_previously_posted_projects</code>	Number of project applications previously submitted by the same teacher. Example: 1

* 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
<code>id</code>	A <code>project_id</code> value from the <code>train.csv</code> file. Example: p036502
<code>description</code>	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
<code>quantity</code>	Quantity of the resource required. Example: 3
<code>price</code>	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
<code>project_is_approved</code>	A binary flag indicating whether DonorsChoose approved the project. A value of <code>0</code> indicates the project was not approved, and a value of <code>1</code> 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]:

```
1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4
5 # Files:
6 import os
7
8 # Data:
9 import sqlite3
10 import pandas as pd
11 import numpy as np
12 from collections import Counter
13
14 # Visuals:
15 import matplotlib.pyplot as plt
16 import seaborn as sns
17 from plotly import plotly
18 import plotly.offline as offline
19 import plotly.graph_objs as go
20 offline.init_notebook_mode()
21 from prettytable import PrettyTable
22
23 # Text:
24 import re
25 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
26 from nltk.corpus import stopwords
27 from nltk.stem.wordnet import WordNetLemmatizer
28 import nltk
29 from nltk.stem.porter import PorterStemmer
30 import string
31 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
32 from gensim.models import Word2Vec
33 from gensim.models import KeyedVectors
34 #from sklearn.feature_extraction.text import TfidfTransformer
35
36 # Metrics:
37 from sklearn import metrics
38 from sklearn.metrics import confusion_matrix, roc_curve, auc
39
40 # Misc:
41 import pickle
42 from tqdm import tqdm
```

c:\users\byron\applications\pythonmaster\lib\site-packages\gensim\utils.py:1
212: UserWarning:

detected Windows; aliasing chunkize to chunkize_serial

1.1 Reading Data

In [2]:

```
1 project_data = pd.read_csv('data/train_data.csv')
2 resource_data = pd.read_csv('data/resources.csv')
```

In [3]:

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

Number of data points in train data (109248, 17)

```
-----
The attributes of data : ['index' '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]:

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

Number of data points in train data (1541272, 4)

['id' 'description' 'quantity' 'price']

Out[4]:

	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 [5]:

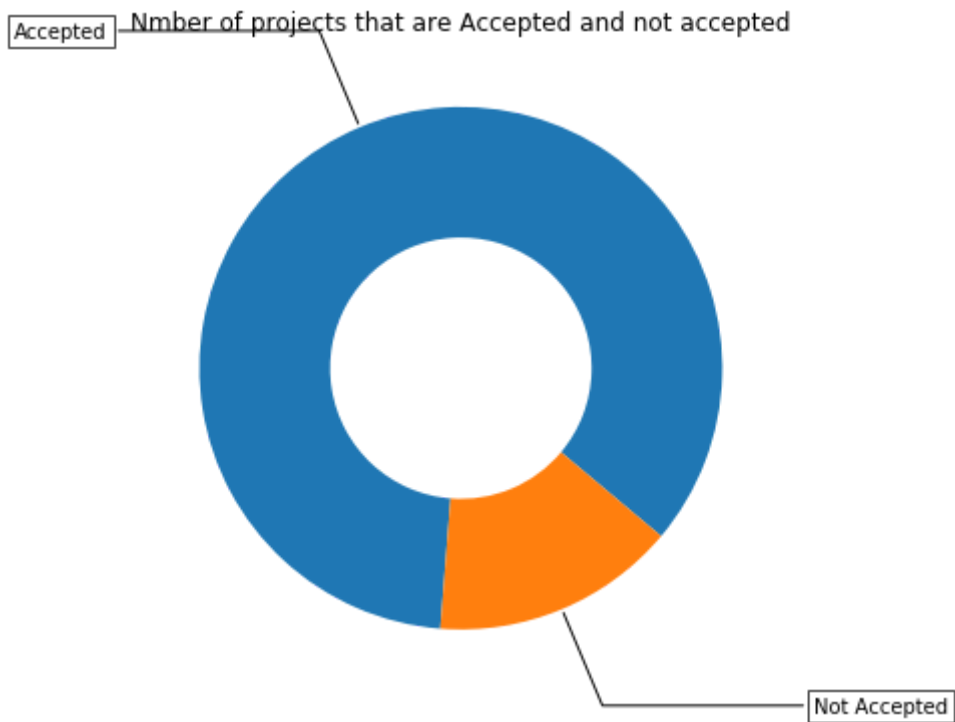
```

1 # this code is taken from
2 # https://matplotlib.org/gallery/pie_and_polar_charts/pie_and_donut_labels.html#sphx-g
3
4
5 y_value_counts = project_data['project_is_approved'].value_counts()
6 print("Number of projects thar are approved for funding ", y_value_counts[1], ", (", (
7 print("Number of projects thar are not approved for funding ", y_value_counts[0], ", (
8
9 fig, ax = plt.subplots(figsize=(6, 6), subplot_kw=dict(aspect="equal"))
10 recipe = ["Accepted", "Not Accepted"]
11
12 data = [y_value_counts[1], y_value_counts[0]]
13
14 wedges, texts = ax.pie(data, wedgeprops=dict(width=0.5), startangle=-40)
15
16 bbox_props = dict(boxstyle="square,pad=0.3", fc="w", ec="k", lw=0.72)
17 kw = dict(xycoords='data', textcoords='data', arrowprops=dict(arrowstyle="-"),
18           bbox=bbox_props, zorder=0, va="center")
19
20 for i, p in enumerate(wedges):
21     ang = (p.theta2 - p.theta1)/2. + p.theta1
22     y = np.sin(np.deg2rad(ang))
23     x = np.cos(np.deg2rad(ang))
24     horizontalalignment = {-1: "right", 1: "left"}[int(np.sign(x))]
25     connectionstyle = "angle,angleA=0,angleB={}".format(ang)
26     kw["arrowprops"].update({"connectionstyle": connectionstyle})
27     ax.annotate(recipe[i], xy=(x, y), xytext=(1.35*np.sign(x), 1.4*y),
28               horizontalalignment=horizontalalignment, **kw)
29
30 ax.set_title("Nmber of projects that are Accepted and not accepted")
31
32 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 %)



1.2.1 Univariate Analysis: School State

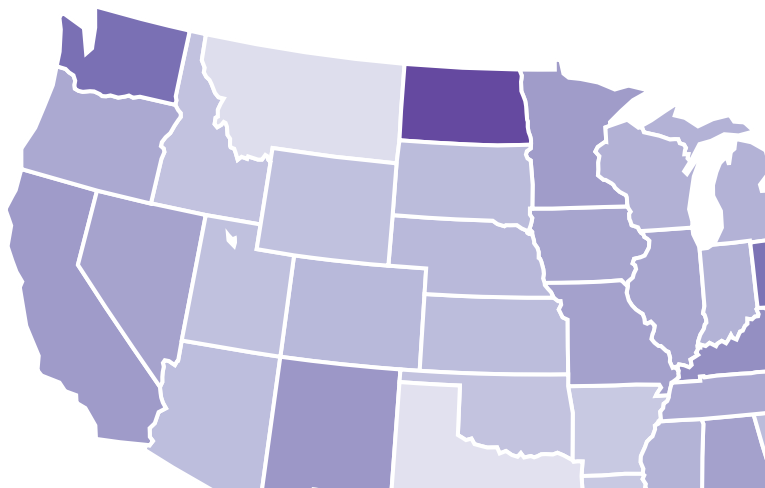
In [6]:

```

1 # Pandas dataframe grouby count, mean: https://stackoverflow.com/a/19385591/4084039
2
3 temp = pd.DataFrame(project_data.groupby("school_state")["project_is_approved"].apply(
4 # if you have data which contain only 0 and 1, then the mean = percentage (think about
5 temp.columns = ['state_code', 'num_proposals']
6
7 # How to plot US state heatmap: https://datascience.stackexchange.com/a/9620
8
9 scl = [[0.0, 'rgb(242,240,247)'],[0.2, 'rgb(218,218,235)'],[0.4, 'rgb(188,189,220)'],\
10         [0.6, 'rgb(158,154,200)'],[0.8, 'rgb(117,107,177)'],[1.0, 'rgb(84,39,143)']
11
12 data = [ dict(
13     type='choropleth',
14     colorscale = scl,
15     autocolorscale = False,
16     locations = temp['state_code'],
17     z = temp['num_proposals'].astype(float),
18     locationmode = 'USA-states',
19     text = temp['state_code'],
20     marker = dict(line = dict (color = 'rgb(255,255,255)',width = 2)),
21     colorbar = dict(title = "% of pro")
22 ) ]
23
24 layout = dict(
25     title = 'Project Proposals % of Acceptance Rate by US States',
26     geo = dict(
27         scope='usa',
28         projection=dict( type='albers usa' ),
29         showlakes = True,
30         lakecolor = 'rgb(255, 255, 255)',
31     ),
32 )
33
34 fig = go.Figure(data=data, layout=layout)
35 offline.ipplot(fig, filename='us-map-heat-map')

```

Project Proposals % of Acceptance Rate





In [7]:

```

1 # https://www.csi.cuny.edu/sites/default/files/pdf/administration/ops/2Letterstabbrev.pdf
2 temp.sort_values(by=['num_proposals'], inplace=True)
3 print("States with lowest % approvals")
4 print(temp.head(5))
5 print('='*50)
6 print("States with highest % approvals")
7 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 [8]:

```

1 #stacked bar plots matplotlib: https://matplotlib.org/gallery/lines_bars_and_markers/stacked_bar.html
2 def stack_plot(data, xtick, col2='project_is_approved', col3='total'):
3     ind = np.arange(data.shape[0])
4
5     plt.figure(figsize=(20,5))
6     p1 = plt.bar(ind, data[col3].values)
7     p2 = plt.bar(ind, data[col2].values)
8
9     plt.ylabel('Projects')
10    plt.title('% of projects aproved state wise')
11    plt.xticks(ind, list(data[xtick].values))
12    plt.legend((p1[0], p2[0]), ('total', 'accepted'))
13    plt.show()

```

In [9]:

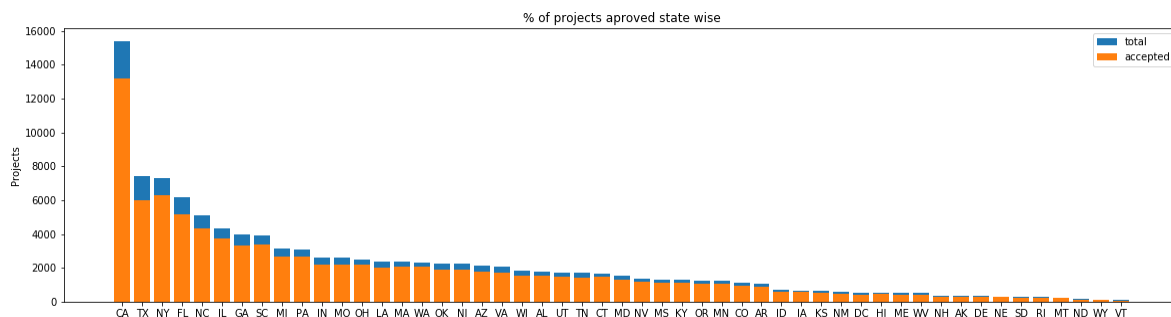
```

1 def univariate_barplots(data, col1, col2='project_is_approved', top=False):
2     # Count number of zeros in dataframe python: https://stackoverflow.com/a/51540521/4
3     temp = pd.DataFrame(project_data.groupby(col1)[col2].agg(lambda x: x.eq(1).sum()))
4
5     # Pandas dataframe grouby count: https://stackoverflow.com/a/19385591/4084039
6     temp['total'] = pd.DataFrame(project_data.groupby(col1)[col2].agg({'total': 'count'}))
7     temp['Avg'] = pd.DataFrame(project_data.groupby(col1)[col2].agg({'Avg': 'mean'})).reset_index()
8
9     temp.sort_values(by=['total'], inplace=True, ascending=False)
10
11     if top:
12         temp = temp[0:top]
13
14     stack_plot(temp, xtick=col1, col2=col2, col3='total')
15     print(temp.head(5))
16     print("="*50)
17     print(temp.tail(5))

```

In [10]:

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



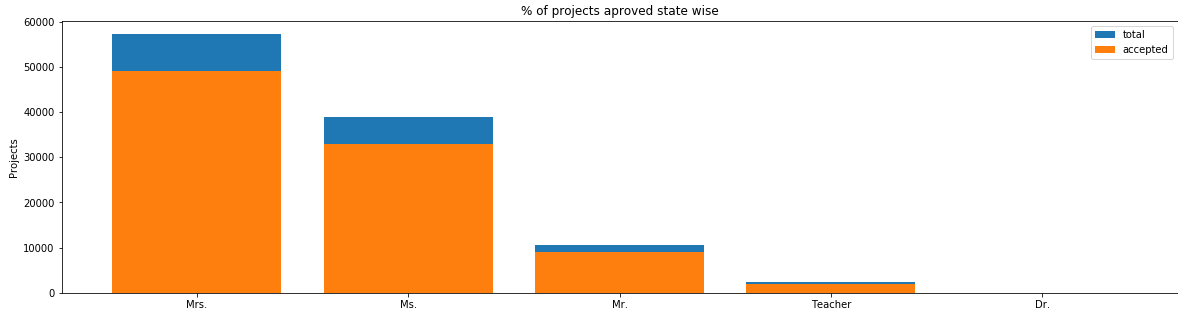
	school_state	project_is_approved	total	Avg
4	CA	13205	15388	0.858136
43	TX	6014	7396	0.813142
34	NY	6291	7318	0.859661
9	FL	5144	6185	0.831690
27	NC	4353	5091	0.855038
=====				
	school_state	project_is_approved	total	Avg
39	RI	243	285	0.852632
26	MT	200	245	0.816327
28	ND	127	143	0.888112
50	WY	82	98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

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

1.2.2 Univariate Analysis: teacher_prefix

In [11]:

```
1 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

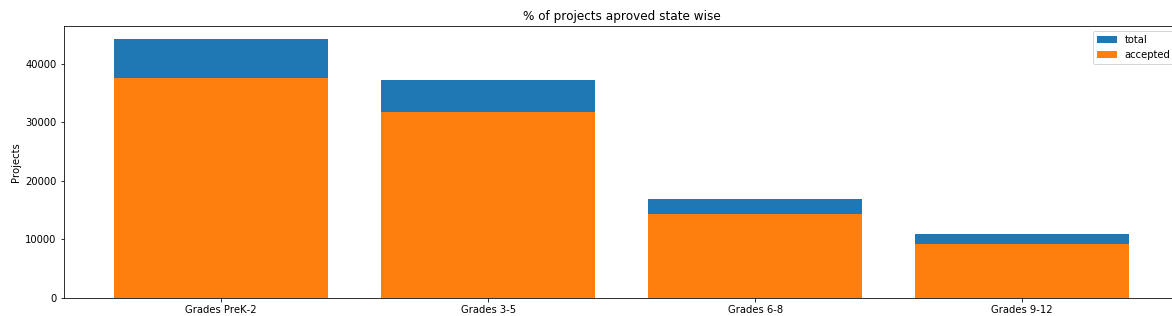
=====

	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 [12]:

```
1 univariate_barplots(project_data, 'project_grade_category', 'project_is_approved', top:
```



	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

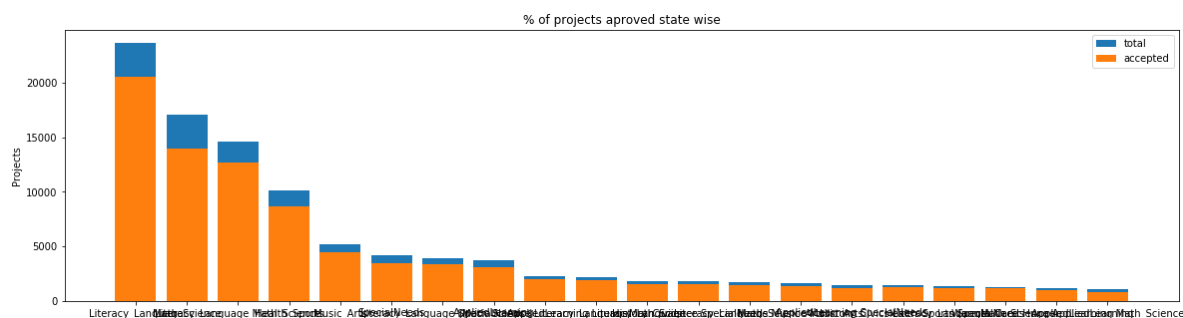
In [13]:

```
1 catogories = list(project_data['project_subject_categories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/4
3
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-str
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-pyt
7 cat_list = []
8 for i in catogories:
9     temp = ""
10    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
11    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth
12        if 'The' in j.split(): # this will split each of the catogory based on space "
13            j=j.replace('The','') # if we have the words "The" we are going to replace
14            j = j.replace(' ', '') # we are placeing all the ' '(space) with ''(empty) ex:"
15            temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing sp
16            temp = temp.replace('&','_') # we are replacing the & value into
17    cat_list.append(temp.strip())
```

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

	index	id	teacher_id	teacher_prefix	school_state	project_su
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	20
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL	20

```
1 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 [16]:

```

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

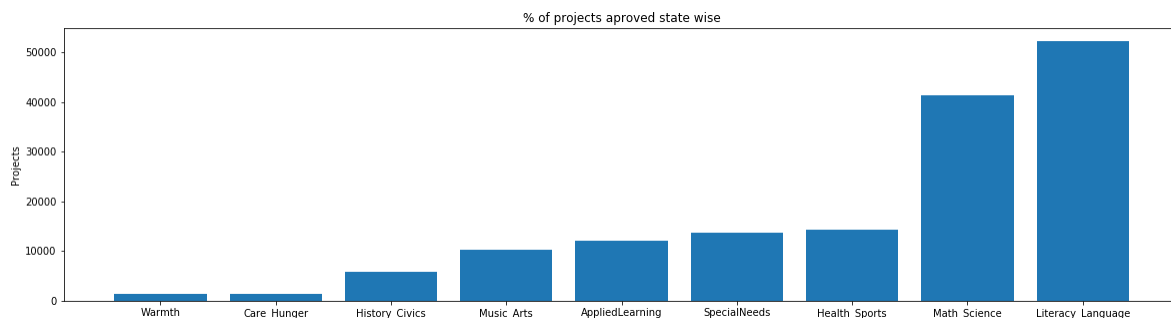
```

In [17]:

```

1 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
2 cat_dict = dict(my_counter)
3 sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
4
5
6 ind = np.arange(len(sorted_cat_dict))
7 plt.figure(figsize=(20,5))
8 p1 = plt.bar(ind, list(sorted_cat_dict.values()))
9
10 plt.ylabel('Projects')
11 plt.title('% of projects aproved state wise')
12 plt.xticks(ind, list(sorted_cat_dict.keys()))
13 plt.show()

```



In [18]:

```

1 for i, j in sorted_cat_dict.items():
2     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 [19]:

```

1 sub_categories = list(project_data['project_subject_subcategories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/4
3
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-str
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-pyt
7
8 sub_cat_list = []
9 for i in sub_categories:
10     temp = ""
11     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
12     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth
13         if 'The' in j.split(): # this will split each of the category based on space "
14             j=j.replace('The','') # if we have the words "The" we are going to replace
15             j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"
16             temp +=j.strip()+" #" "abc ".strip() will return "abc", remove the trailing sp
17             temp = temp.replace('&','_')
18     sub_cat_list.append(temp.strip())

```

In [20]:

```

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

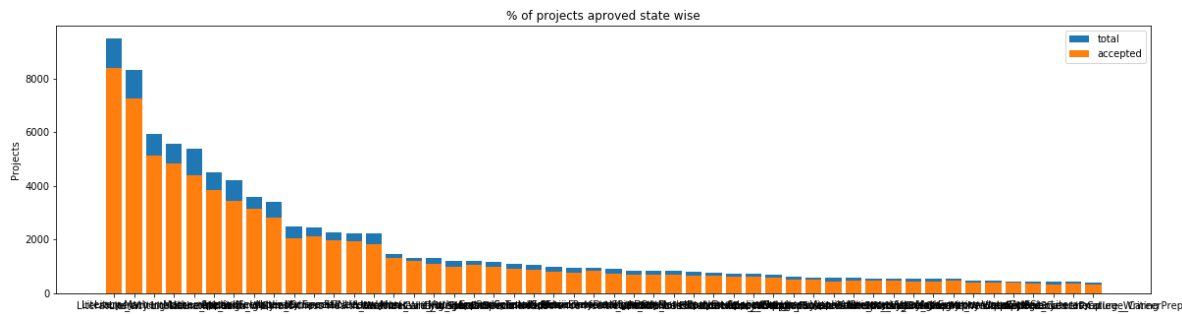
```

Out[20]:

	index	id	teacher_id	teacher_prefix	school_state	project_su
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	20
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	20

In [21]:

```
1 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

=====

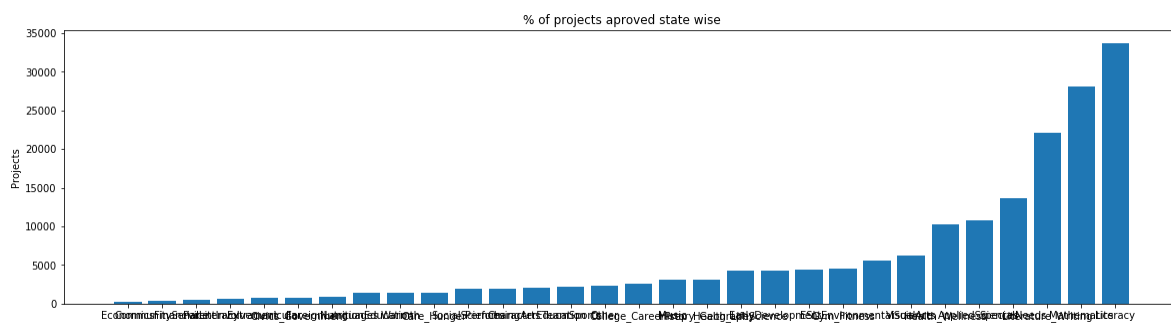
	clean_subcategories	project_is_approved	total	Av
g				
196	EnvironmentalScience Literacy	389	444	0.87612
6				
127	ESL	349	421	0.82897
9				
79	College_CareerPrep	343	421	0.81472
7				
17	AppliedSciences Literature_Writing	361	420	0.85952
4				
3	AppliedSciences College_CareerPrep	330	405	0.81481
5				

In [22]:

```
1 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084035
2 from collections import Counter
3 my_counter = Counter()
4 for word in project_data['clean_subcategories'].values:
5     my_counter.update(word.split())
```


In [23]:

```
1 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
2 sub_cat_dict = dict(my_counter)
3 sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
4
5
6 ind = np.arange(len(sorted_sub_cat_dict))
7 plt.figure(figsize=(20,5))
8 p1 = plt.bar(ind, list(sorted_sub_cat_dict.values()))
9
10 plt.ylabel('Projects')
11 plt.title('% of projects aproved state wise')
12 plt.xticks(ind, list(sorted_sub_cat_dict.keys()))
13 plt.show()
```



In [24]:

```
1 for i, j in sorted_sub_cat_dict.items():  
2     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

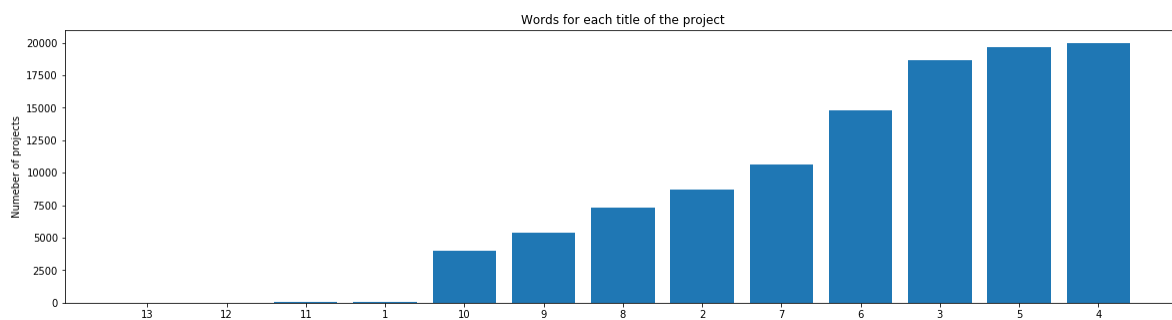
1.2.6 Univariate Analysis: Text features (Title)

In [25]:

```

1 #How to calculate number of words in a string in DataFrame: https://stackoverflow.com/
2 word_count = project_data['project_title'].str.split().apply(len).value_counts()
3 word_dict = dict(word_count)
4 word_dict = dict(sorted(word_dict.items(), key=lambda kv: kv[1]))
5
6
7 ind = np.arange(len(word_dict))
8 plt.figure(figsize=(20,5))
9 p1 = plt.bar(ind, list(word_dict.values()))
10
11 plt.ylabel('Numeber of projects')
12 plt.title('Words for each title of the project')
13 plt.xticks(ind, list(word_dict.keys()))
14 plt.show()

```



In [26]:

```

1 approved_word_count = project_data[project_data['project_is_approved']==1]['project_title']
2 approved_word_count = approved_word_count.values
3
4 rejected_word_count = project_data[project_data['project_is_approved']==0]['project_title']
5 rejected_word_count = rejected_word_count.values

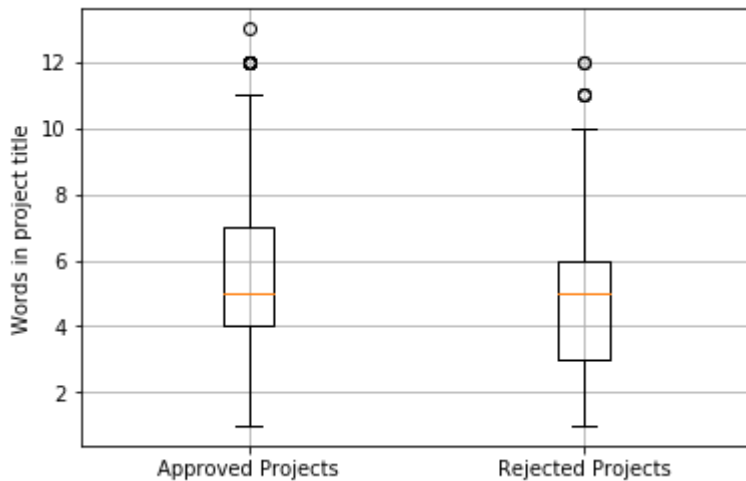
```

In [27]:

```

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

```

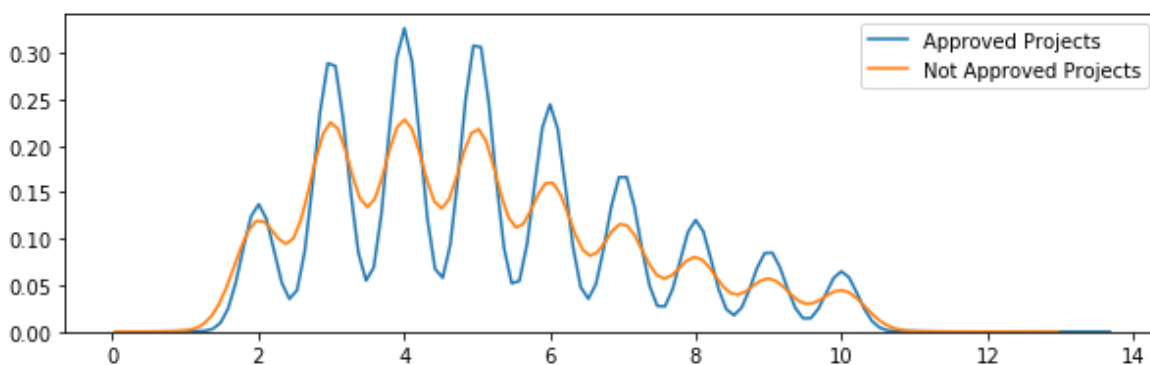


In [28]:

```

1 plt.figure(figsize=(10,3))
2 sns.distplot(approved_word_count, hist=False, label="Approved Projects")
3 sns.distplot(rejected_word_count, hist=False, label="Not Approved Projects")
4 plt.legend()
5 plt.show()

```



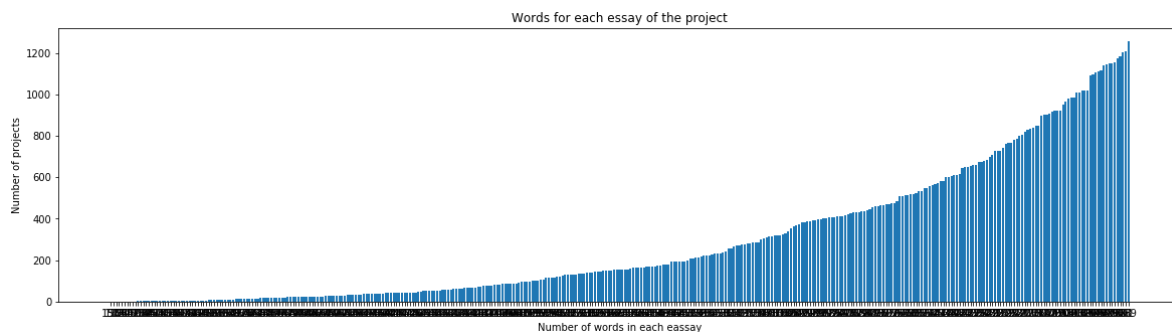
1.2.7 Univariate Analysis: Text features (Project Essay's)

In [29]:

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

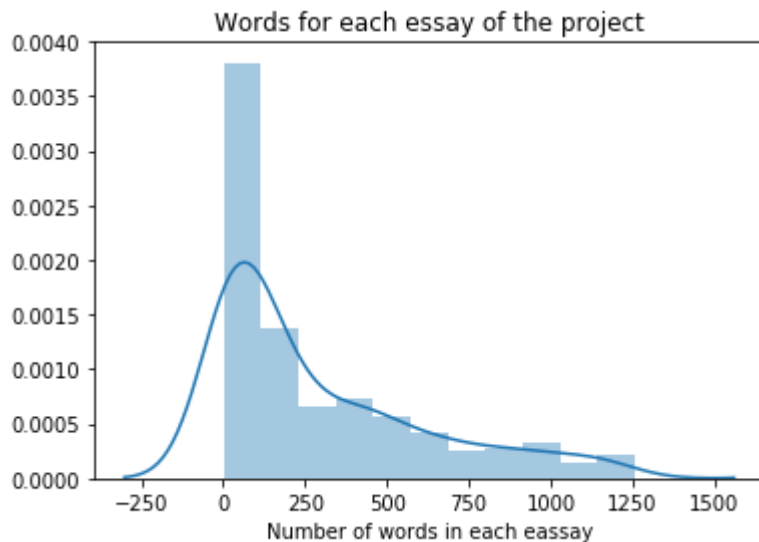
In [30]:

```
1 #How to calculate number of words in a string in DataFrame: https://stackoverflow.com/
2 word_count = project_data['essay'].str.split().apply(len).value_counts()
3 word_dict = dict(word_count)
4 word_dict = dict(sorted(word_dict.items(), key=lambda kv: kv[1]))
5
6
7 ind = np.arange(len(word_dict))
8 plt.figure(figsize=(20,5))
9 p1 = plt.bar(ind, list(word_dict.values()))
10
11 plt.ylabel('Number of projects')
12 plt.xlabel('Number of words in each eassay')
13 plt.title('Words for each essay of the project')
14 plt.xticks(ind, list(word_dict.keys()))
15 plt.show()
```



In [31]:

```
1 sns.distplot(word_count.values)
2 plt.title('Words for each essay of the project')
3 plt.xlabel('Number of words in each eassay')
4 plt.show()
5
```



In [32]:

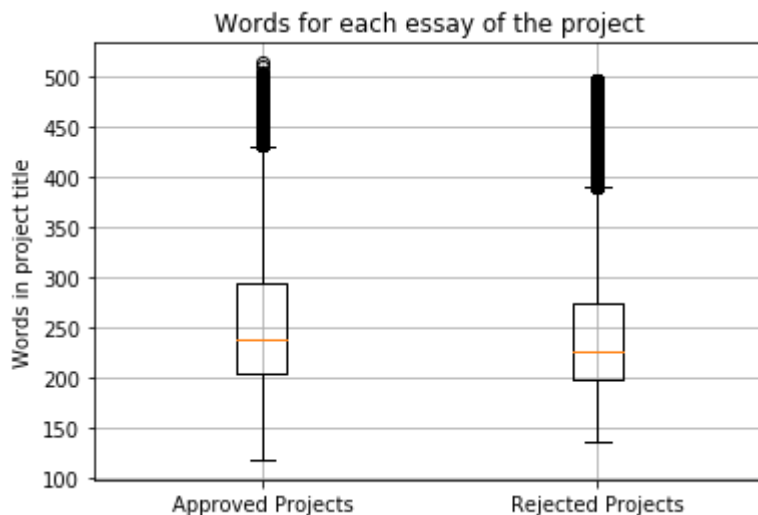
```
1 approved_word_count = project_data[project_data['project_is_approved']==1]['essay'].st
2 approved_word_count = approved_word_count.values
3
4 rejected_word_count = project_data[project_data['project_is_approved']==0]['essay'].st
5 rejected_word_count = rejected_word_count.values
```

In [33]:

```

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

```

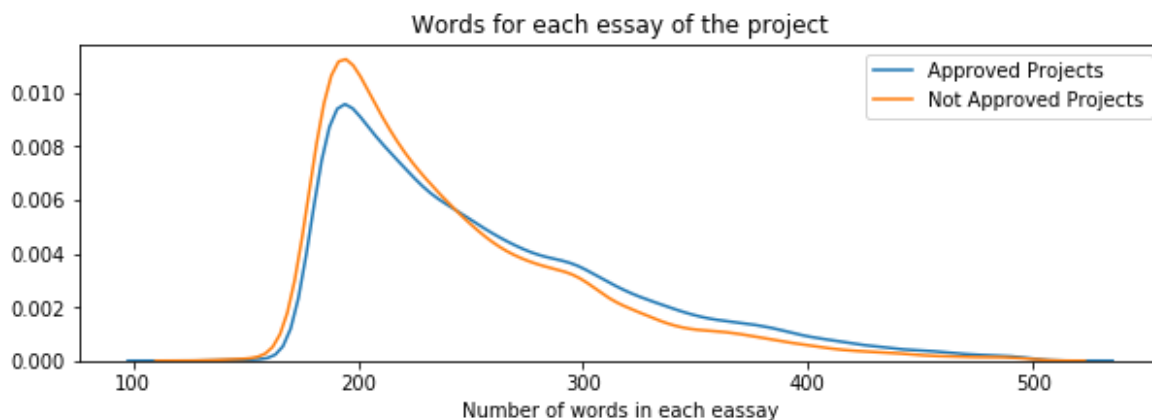


In [34]:

```

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

```



1.2.8 Univariate Analysis: Cost per project

In [35]:

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

Out[35]:

	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 [36]:

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

Out[36]:

	id	price	quantity
0	p000001	459.56	7
1	p000002	515.89	21

In [37]:

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

In [38]:

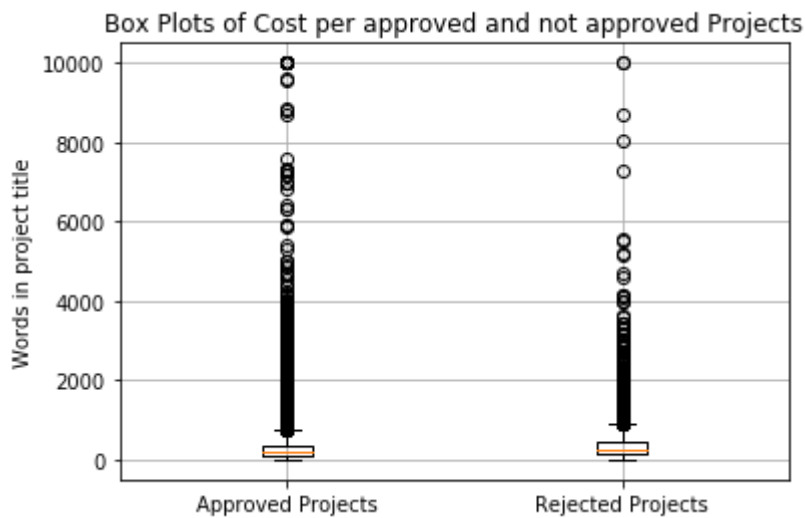
```
1 approved_price = project_data[project_data['project_is_approved']==1]['price'].values
2
3 rejected_price = project_data[project_data['project_is_approved']==0]['price'].values
```


In [39]:

```

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

```

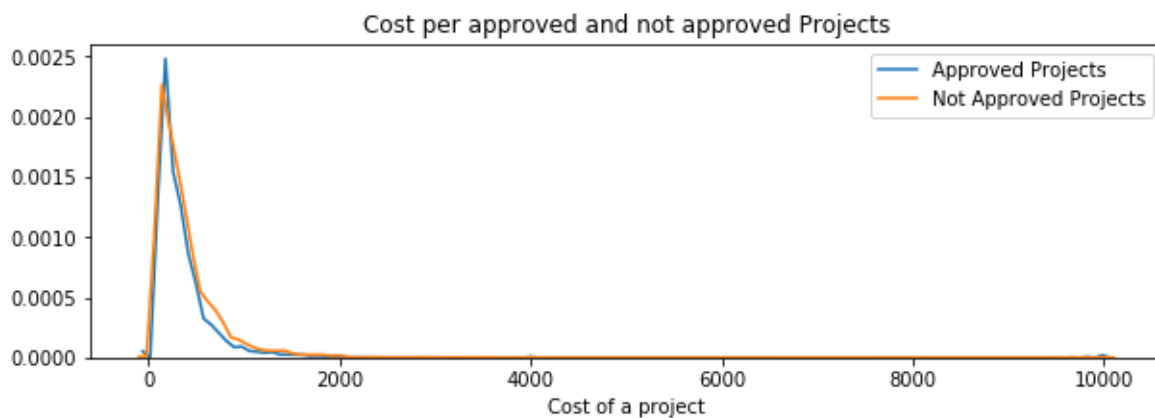


In [40]:

```

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

```



In [41]:

```

1 # http://zetcode.com/python/prettytable/
2 table = PrettyTable()
3 table.field_names = ["Percentile", "Approved Projects", "Not Approved Projects"]
4
5 for i in range(0,101,5):
6     table.add_row([i,np.round(np.percentile(approved_price,i), 3), np.round(np.percentile(not_approved_price,i), 3)])
7 print(table)

```

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 [42]:

```

1 teacher_projects = project_data.loc[:, ['teacher_id', 'teacher_number_of_previously_posted_projects']]
2 teacher_projects
3 plt.boxplot([teacher_projects.query('project_is_approved == 1')['teacher_number_of_previously_posted_projects'],
4               teacher_projects.query('project_is_approved == 0')['teacher_number_of_previously_posted_projects']])
5 plt.title('Box Plots of previously posted projects')
6 plt.xticks([1,2], ('Approved Projects', 'Rejected Projects'))
7 plt.grid()
8 plt.show()

```



1.2.10 Univariate Analysis: project_resource_summary

In [43]:

```

1 def check_numeric(x):
2     for i in x.split():
3         if i.isnumeric() == True:
4             return True
5         else:
6             continue
7     return False
8
9 project_data['project_resource_summary_num'] = project_data['project_resource_summary'].apply(check_numeric)

```

In [44]:

```
1 projects_resource_summary_num = project_data.loc[:,['project_resource_summary_num','project_is_approved']]
2 projects_resource_summary_num
```

Out[44]:

	project_resource_summary_num	project_is_approved
0	False	82562
1	True	10144

In [45]:

```
1 plt.bar(projects_resource_summary_num['project_resource_summary_num'],projects_resource_summary_num['project_is_approved'])
2 plt.ylabel('Number of projects')
3 plt.xlabel('number in resource summary')
4 plt.title('Split on resource summary')
5 plt.xticks(projects_resource_summary_num['project_resource_summary_num'])
6 plt.show()
```

