

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature	Description
<code>project_id</code>	A unique identifier for the proposed project. Example
<code>project_title</code>	Title of the project. Examples: <ul style="list-style-type: none"> • Art Will Make You Happy! • First Grade Fun
<code>project_grade_category</code>	Grade level of students for which the project is targeted. Enumerated values: <ul style="list-style-type: none"> • Grades PreK-2 • Grades 3-5 • Grades 6-8 • Grades 9-12
<code>project_subject_categories</code>	One or more (comma-separated) subject categories from the following enumerated list of values: <ul style="list-style-type: none"> • Applied Learning • Care & Hunger • Health & Sports • History & Civics • Literacy & Language • Math & Science • Music & The Arts • Special Needs • Warmth Examples: <ul style="list-style-type: none"> • Music & The Arts • Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (<u>Two-letter U.S. postal code</u> (https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations)). Example: WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories. Examples: <ul style="list-style-type: none"> • Literacy • Literature & Writing, Social Sciences
<code>project_resource_summary</code>	An explanation of the resources needed for the project. <ul style="list-style-type: none"> • My students need hands on literacy materials to address sensory needs!

Feature	Description
project_essay_1	First application essay*
project_essay_2	Second application essay*
project_essay_3	Third application essay*
project_essay_4	Fourth application essay*
project_submitted_datetime	Datetime when project application was submitted. Example: 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> • nan • Dr. • Mr. • Mrs. • Ms. • Teacher.
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the teacher. Example: 2

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

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

Feature	Description
id	A project_id value from the <code>train.csv</code> file. Example: p036502
description	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
quantity	Quantity of the resource required. Example: 3
price	Price of the resource required. Example: 9.95

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

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

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



Notes on the Essay Data

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

- __project_essay_1:__ "Introduce us to your classroom"
- __project_essay_2:__ "Tell us more about your students"
- __project_essay_3:__ "Describe how your students will use the materials you're requesting"
- __project_essay_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]:

```
%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 chart_studio.plotly import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()

from collections import Counter
from prettytable import PrettyTable

from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import roc_auc_score
import math
```

In [2]:

```
df1 = pd.read_csv('train_data.csv',nrows=80000)
df2 =pd.read_csv('resources.csv', nrows =80000)
```

In [3]:

```
print("Number of data points in train data", df1.shape)
print('-'*50)
print("The attributes of data :", df1.columns.values)
```

Number of data points in train data (80000, 17)

```
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix'
'school_state'
'project_submitted_datetime' 'project_grade_category'
'project_subject_categories' 'project_subject_subcategories'
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

In [4]:

```
print("Number of data points in train data", df2.shape)
print(df2.columns.values)
print(df2.head(2))
```

Number of data points in train data (80000, 4)

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

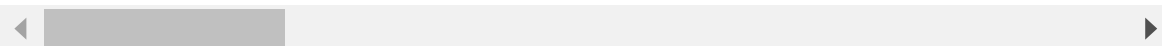
	id	description	quantity	\
0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	
1	p069063	Bouncy Bands for Desks (Blue support pipes)	3	
	price			
0	149.00			
1	14.95			

In [5]:

df1.head()

Out[5]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_s
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



In [6]:

```
# merge two column text dataframe:
df1["essay"] = df1["project_essay_1"].map(str) + \
                df1["project_essay_2"].map(str) + \
                df1["project_essay_3"].map(str) + \
                df1["project_essay_4"].map(str)
```

preprocessing of project subject categories

In [7]:

```
categories = list(df1['project_subject_categories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039

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

df1['clean_categories'] = cat_list
df1.drop(['project_subject_categories'], axis=1, inplace=True)

from collections import Counter
my_counter = Counter()
for word in df1['clean_categories'].values:
    my_counter.update(word.split())

cat_dict = dict(my_counter)
sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
sorted_cat_dict
```

Out[7]:

```
{'Warmth': 1029,
 'Care_Hunger': 1029,
 'History_Civics': 4319,
 'Music_Arts': 7540,
 'AppliedLearning': 8917,
 'SpecialNeeds': 9974,
 'Health_Sports': 10445,
 'Math_Science': 30285,
 'Literacy_Language': 38266}
```

preprocessing project subject subcategories

In [8]:

```
sub_categories = list(df1['project_subject_subcategories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039

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

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

df1['clean_subcategories'] = sub_cat_list
df1.drop(['project_subject_subcategories'], axis=1, inplace=True)

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in df1['clean_subcategories'].values:
    my_counter.update(word.split())

sub_cat_dict = dict(my_counter)
sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
sorted_sub_cat_dict
```

Out[8]:

```
{'Economics': 198,
 'CommunityService': 332,
 'FinancialLiteracy': 417,
 'ParentInvolvement': 496,
 'Civics_Government': 591,
 'Extracurricular': 607,
 'ForeignLanguages': 623,
 'NutritionEducation': 995,
 'Warmth': 1029,
 'Care_Hunger': 1029,
 'SocialSciences': 1381,
 'PerformingArts': 1439,
 'CharacterEducation': 1499,
 'TeamSports': 1596,
 'Other': 1759,
 'College_CareerPrep': 1865,
 'History_Geography': 2310,
 'Music': 2321,
 'Health_LifeScience': 3059,
 'EarlyDevelopment': 3130,
 'ESL': 3169,
 'Gym_Fitness': 3337,
 'EnvironmentalScience': 4060,
 'VisualArts': 4575,
 'Health_Wellness': 7509,
 'AppliedSciences': 7904,
 'SpecialNeeds': 9974,
 'Literature_Writing': 16239,
 'Mathematics': 20552,
 'Literacy': 24757}
```

In [9]:

```
df1['teacher_prefix'] = df1['teacher_prefix'].fillna(df1['teacher_prefix'].mode().iloc[0])

prefix = list(df1['teacher_prefix'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039

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

prefix_list = []
for i in prefix:
    temp = ""
    if "." in i:
        i = i.replace('.', '')
        temp += i.strip() + " #" abc ".strip() will return "abc", remove the trailing spaces
    prefix_list.append(temp.strip())

df1['teachers_prefix'] = prefix_list
df1.drop(['teacher_prefix'], inplace = True, axis = 1)
```

In [10]:

```
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in df1['teachers_prefix'].values:
    my_counter.update(word.split())

prefix_dict = dict(my_counter)
sorted_prefix_dict = dict(sorted(prefix_dict.items(), key=lambda kv: kv[1]))
sorted_prefix_dict
```

Out[10]:

```
{'Dr': 8, 'Mr': 7770, 'Ms': 28741, 'Mrs': 41776}
```

preprocessing project grade category

In [11]:

```
grade_categories = list(df1['project_grade_category'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039

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

grade_cat_list = []
for i in grade_categories:
    temp = ""
    if "-" in i:
        i = i.replace('-', '_')
        i = i.replace(' ', '_')
        temp += i.strip() + ' #" abc ".strip() will return "abc", remove the trailing spaces
    grade_cat_list.append(temp.strip())

df1['grade_category'] = grade_cat_list
df1.drop(['project_grade_category'], axis=1, inplace=True)
```

In [12]:

```
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in df1['grade_category'].values:
    my_counter.update(word.split())

grade_dict = dict(my_counter)
sorted_grade_dict = dict(sorted(grade_dict.items(), key=lambda kv: kv[1]))

sorted_grade_dict
```

Out[12]:

```
{'Grades_9_12': 8021,
 'Grades_6_8': 12383,
 'Grades_3_5': 27244,
 'Grades_PreK_2': 32352}
```

train test splitting

In [13]:

```
y = df1['project_is_approved'].values
X = df1.drop(['Unnamed: 0', 'teacher_id', 'project_submitted_datetime'], axis=1)
X.head(1)
```

Out[13]:

	id	school_state	project_title	project_essay_1	project_essay_2	project_ess
0	p253737	IN	Educational Support for English Learners at Home	My students are English learners that are work...	"The limits of your language are the limits o...	NaN

In [14]:

```
from sklearn.model_selection import train_test_split

# split the data into test and train by maintaining same distribution of output variable 'y' [stratify=y]
X_train, X_test, y_train, y_test = train_test_split(X, y, stratify=y, test_size=0.2)
# split the train data into train and cross validation by maintaining same distribution of output variable 'y_train' [stratify=y_train]
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_train, test_size=0.2)
```

In [15]:

```
# 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 [16]:

```
sent = decontracted(df1['essay'].values[20000])
print(sent)
print("="*50)
```

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

=====

In [17]:

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

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

In [18]:

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

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

In [19]:

```
# 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'r
e", "you've",\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him',
'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 't
hey', 'them', 'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "th
at'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'ha
d', 'having', 'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as'
, 'until', 'while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through'
, 'during', 'before', 'after',\
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'ov
er', 'under', 'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'an
y', 'both', 'each', 'few', 'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too'
, 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'no
w', 'd', 'll', 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't",
'doesn', "doesn't", 'hadn',\
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'migh
tn', "mightn't", 'mustn',\
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'w
asn', "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [20]:

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

```
100%|████████████████████████████████████████████████████████████████████████████████| 51200/51200 [00:55<00:00, 928.25it/s]
```

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

In [22]:

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

In [23]:

```
# Combining all the above students
from tqdm import tqdm
preprocessed_train_title = []
# tqdm is for printing the status bar
for sentence in tqdm(X_train['project_title'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\r', ' ')
    sent = sent.replace('\n', ' ')
    sent = sent.replace('\t', ' ')
    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_train_title.append(sent.lower().strip())
```

16/65


```
preprocessed_cv_title = []
# tqdm is for printing the status bar
for sentence in tqdm(X_cv['project_title'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
    preprocessed_cv_title.append(sent.lower().strip())
```

In [25]:

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

17/65

In [26]:

```
# it returns a dict, keys as aprvd labels and values as the number of data points in th
at aprvd
train_aprvd_distribution = X_train['project_is_approved'].value_counts().sort_index()
test_aprvd_distribution = X_test['project_is_approved'].value_counts().sort_index()
cv_aprvd_distribution = X_cv['project_is_approved'].value_counts().sort_index()

my_colors = 'rgbkymc'
train_aprvd_distribution.plot(kind='bar')
plt.xlabel('approved')
plt.ylabel('Data points per approved')
plt.title('Distribution of yi in train data')
plt.grid()
plt.show()

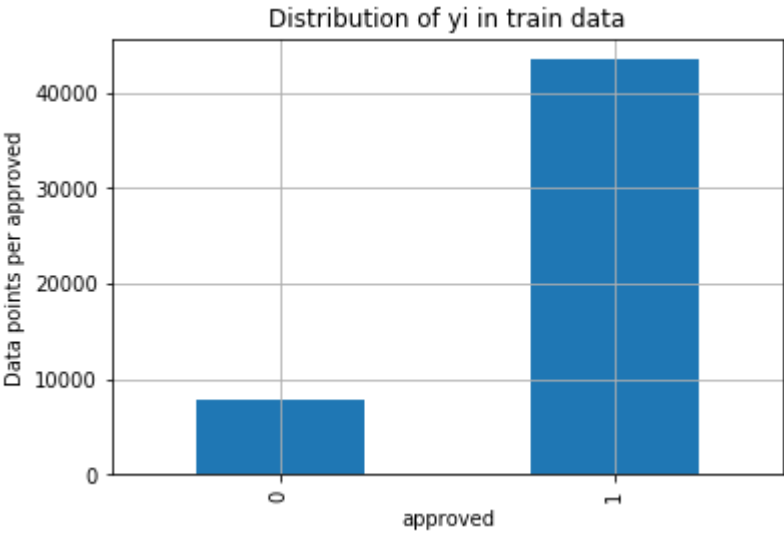
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_aprvd_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_aprvd_distribution.values)
for i in sorted_yi:
    print('Number of data points in approved', i+1, ':', train_aprvd_distribution.values
[i], '(', np.round((train_aprvd_distribution.values[i]/X_train.shape[0]*100), 3), '%)')

print('-'*80)
my_colors = 'rgbkymc'
test_aprvd_distribution.plot(kind='bar')
plt.xlabel('approved')
plt.ylabel('Data points per approved')
plt.title('Distribution of yi in test data')
plt.grid()
plt.show()

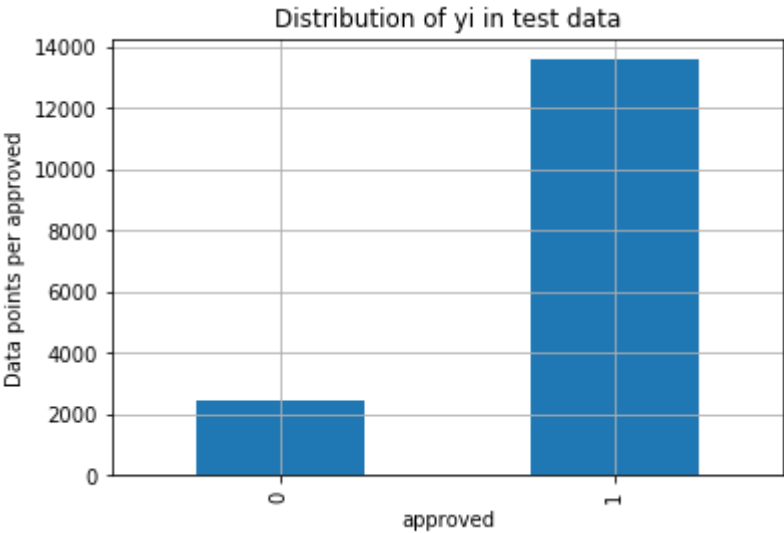
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_aprvd_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-test_aprvd_distribution.values)
for i in sorted_yi:
    print('Number of data points in approved', i+1, ':', test_aprvd_distribution.values[
i], '(', np.round((test_aprvd_distribution.values[i]/X_test.shape[0]*100), 3), '%)')

print('-'*80)
my_colors = 'rgbkymc'
cv_aprvd_distribution.plot(kind='bar')
plt.xlabel('approved')
plt.ylabel('Data points per approved')
plt.title('Distribution of yi in cross validation data')
plt.grid()
plt.show()

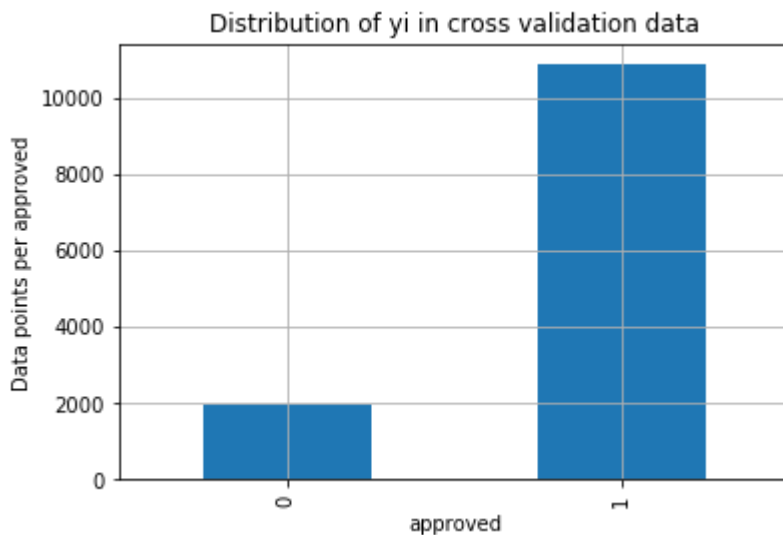
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_aprvd_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_aprvd_distribution.values)
for i in sorted_yi:
    print('Number of data points in approved', i+1, ':', cv_aprvd_distribution.values[i
], '(', np.round((cv_aprvd_distribution.values[i]/X_cv.shape[0]*100), 3), '%)')
```



Number of data points in approved 2 : 43439 (84.842 %)
Number of data points in approved 1 : 7761 (15.158 %)



Number of data points in approved 2 : 13575 (84.844 %)
Number of data points in approved 1 : 2425 (15.156 %)



Number of data points in approved 2 : 10860 (84.844 %)

Number of data points in approved 1 : 1940 (15.156 %)

Response encoding

for school state

In [27]:

```
def get_project_fea_dict(alpha, feature, df):

    value_count = X_train[feature].value_counts()
    project_dict = dict()

    # denominator will contain the number of time that particular feature occurred in whole data
    for i, denominator in value_count.items():
        vec = []
        for k in range(0,2):
            aprvd_cnt = X_train.loc[(X_train['project_is_approved']==k) & (X_train[feature]==i)]

            vec.append((aprvd_cnt.shape[0] + alpha*10)/ (denominator + 90*alpha))

        project_dict[i]=vec
    return project_dict

def get_project_feature(alpha, feature, df):
    project_dict = get_project_fea_dict(alpha, feature, df)
    value_count = X_train[feature].value_counts()
    project_fea = []
    # for every feature values in the given data frame we will check if it is there in the train data then we will add the feature to project_fea
    # if not we will add [1/2,1/2] to project_fea
    for index, row in df.iterrows():
        if row[feature] in dict(value_count).keys():
            project_fea.append(project_dict[row[feature]])
        else:
            project_fea.append([1/2,1/2])
    return project_fea
```

In [28]:

```
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()
for word in df1['school_state'].values:
    my_counter.update(word.split())

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

Out[28]:

```
{'VT': 58,  
'WY': 79,  
'ND': 106,  
'MT': 168,  
'RI': 206,  
'SD': 221,  
'NE': 236,  
'NH': 237,  
'DE': 250,  
'AK': 256,  
'WV': 354,  
'HI': 369,  
'ME': 369,  
'DC': 382,  
'NM': 398,  
'KS': 460,  
'IA': 486,  
'ID': 501,  
'AR': 734,  
'CO': 858,  
'MN': 870,  
'OR': 904,  
'KY': 955,  
'MS': 955,  
'NV': 1016,  
'MD': 1087,  
'TN': 1202,  
'CT': 1235,  
'UT': 1270,  
'AL': 1273,  
'WI': 1331,  
'VA': 1513,  
'AZ': 1561,  
'NJ': 1625,  
'OK': 1710,  
'WA': 1715,  
'LA': 1764,  
'MA': 1765,  
'OH': 1819,  
'MO': 1896,  
'IN': 1897,  
'PA': 2237,  
'MI': 2341,  
'SC': 2881,  
'GA': 2908,  
'IL': 3178,  
'NC': 3737,  
'FL': 4568,  
'NY': 5391,  
'TX': 5406,  
'CA': 11262}
```

In [29]:

```
#response-coding of the feature
# alpha is used for laplace smoothing
alpha = 1
# train feature
train_state_responseCoding = np.array(get_project_feature(alpha, "school_state", X_train))
# test feature
test_state_responseCoding = np.array(get_project_feature(alpha, "school_state", X_test))
# cross validation feature
cv_state_responseCoding = np.array(get_project_feature(alpha, "school_state", X_cv))
```

In [30]:

```
print("Shape of X_train after response encodig ",train_state_responseCoding.shape, y_train.shape)
print("Shape of X_cv after response encodig ",cv_state_responseCoding.shape, y_cv.shape)
print("Shape of X_test after response encodig ",test_state_responseCoding.shape, y_test.shape)
```

```
Shape of X_train after response encodig (51200, 2) (51200,)
Shape of X_cv after response encodig (12800, 2) (12800,)
Shape of X_test after response encodig (16000, 2) (16000,)
```

for teacher_prefix

In [31]:

```
#response-coding of the feature
# alpha is used for laplace smoothing
alpha = 1
# train feature
train_teacher_responseCoding = np.array(get_project_feature(alpha, "teachers_prefix", X_train))
# test feature
test_teacher_responseCoding = np.array(get_project_feature(alpha, "teachers_prefix", X_test))
# cross validation feature
cv_teacher_responseCoding = np.array(get_project_feature(alpha, "teachers_prefix", X_cv))
```

In [32]:

```
print("Shape of X_train after response encodig ",train_teacher_responseCoding.shape, y_train.shape)
print("Shape of X_cv after response encodig ",cv_teacher_responseCoding.shape, y_cv.shape)
print("Shape of X_test after response encodig ",test_teacher_responseCoding.shape, y_test.shape)
```

```
Shape of X_train after response encodig (51200, 2) (51200,)
Shape of X_cv after response encodig (12800, 2) (12800,)
Shape of X_test after response encodig (16000, 2) (16000,)
```

for project grade category

In [33]:

```
#response-coding of the feature  
# alpha is used for laplace smoothing  
alpha = 1  
# train feature  
train_grade_responseCoding = np.array(get_project_feature(alpha, "grade_category", X_train))  
# test feature  
test_grade_responseCoding = np.array(get_project_feature(alpha, "grade_category", X_test))  
# cross validation feature  
cv_grade_responseCoding = np.array(get_project_feature(alpha, "grade_category", X_cv))
```

In [34]:

```
print("Shape of X_train after response encoding ",train_grade_responseCoding.shape, y_train.shape)  
print("Shape of X_cv after response encoding ",cv_grade_responseCoding.shape, y_cv.shape)  
print("Shape of X_test after response encoding ",test_grade_responseCoding.shape, y_test.shape)
```

```
Shape of X_train after response encoding (51200, 2) (51200,)  
Shape of X_cv after response encoding (12800, 2) (12800,)  
Shape of X_test after response encoding (16000, 2) (16000,)
```

for project subject categories

In [35]:

```
#response-coding of the feature  
# alpha is used for laplace smoothing  
alpha = 1  
# train feature  
train_category_responseCoding = np.array(get_project_feature(alpha, "clean_categories", X_train))  
# test feature  
test_category_responseCoding = np.array(get_project_feature(alpha, "clean_categories", X_test))  
# cross validation feature  
cv_category_responseCoding = np.array(get_project_feature(alpha, "clean_categories", X_cv))
```


In [36]:

```
print("Shape of X_train after response encoding ",train_category_responseCoding.shape, y_train.shape)
print("Shape of X_cv after response encoding ",cv_category_responseCoding.shape, y_cv.shape)
print("Shape of X_test after response encoding ",test_category_responseCoding.shape, y_test.shape)
```

```
Shape of X_train after response encoding (51200, 2) (51200,)
Shape of X_cv after response encoding (12800, 2) (12800,)
Shape of X_test after response encoding (16000, 2) (16000,)
```

for project categories

In [37]:

```
#response-coding of the feature
# alpha is used for laplace smoothing
alpha = 1
# train feature
train_subcategory_responseCoding = np.array(get_project_feature(alpha, "clean_subcategories", X_train))
# test feature
test_subcategory_responseCoding = np.array(get_project_feature(alpha, "clean_subcategories", X_test))
# cross validation feature
cv_subcategory_responseCoding = np.array(get_project_feature(alpha, "clean_subcategories", X_cv))
```

In [38]:

```
print("Shape of X_train after response encoding ",train_subcategory_responseCoding.shape, y_train.shape)
print("Shape of X_cv after response encoding ",cv_subcategory_responseCoding.shape, y_cv.shape)
print("Shape of X_test after response encoding ",test_subcategory_responseCoding.shape, y_test.shape)
```

```
Shape of X_train after response encoding (51200, 2) (51200,)
Shape of X_cv after response encoding (12800, 2) (12800,)
Shape of X_test after response encoding (16000, 2) (16000,)
```

bag of words on essay

In [39]:

```
vec1 = CountVectorizer(min_df=10)
# fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vec1.fit_transform(preprocessed_train_essay)
X_cv_essay_bow = vec1.transform(preprocessed_cv_essay)
X_test_essay_bow = vec1.transform(preprocessed_test_essay)

print("After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
print(X_cv_essay_bow.shape, y_cv.shape)
print(X_test_essay_bow.shape, y_test.shape)
```

```
After vectorizations
(51200, 12247) (51200,)
(12800, 12247) (12800,)
(16000, 12247) (16000,)
```

bag of words on title

In [40]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vec2 = CountVectorizer(min_df=10)

# fit has to happen only on train data
# we use the fitted CountVectorizer to convert the text to vector
X_train_title_bow = vec2.fit_transform(preprocessed_train_title)
X_cv_title_bow = vec2.transform(preprocessed_cv_title)
X_test_title_bow = vec2.transform(preprocessed_test_title)

print("After vectorizations")
print(X_train_title_bow.shape, y_train.shape)
print(X_cv_title_bow.shape, y_cv.shape)
print(X_test_title_bow.shape, y_test.shape)
```

```
After vectorizations
(51200, 2053) (51200,)
(12800, 2053) (12800,)
(16000, 2053) (16000,)
```

tfidf on essay

In [41]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vec3 = TfidfVectorizer(min_df=10, max_features =5000)

# fit has to happen only on train data
X_train_ess_tfidf = vec3.fit_transform(preprocessed_train_essay)

# we use the fitted Tfidf Vectorizer to convert the text to vector
X_cv_ess_tfidf = vec3.transform(preprocessed_cv_essay)
X_test_ess_tfidf = vec3.transform(preprocessed_test_essay)

print("After vectorizations")
print(X_train_ess_tfidf.shape, y_train.shape)
print(X_cv_ess_tfidf.shape, y_cv.shape)
print(X_test_ess_tfidf.shape, y_test.shape)
```

After vectorizations
(51200, 5000) (51200,)
(12800, 5000) (12800,)
(16000, 5000) (16000,)

tfidf on title

In [42]:

```
vec4 = TfidfVectorizer(min_df=10)
X_train_title_tfidf = vec4.fit_transform(preprocessed_train_title)

# we use the fitted Tfidf Vectorizer to convert the text to vector
X_cv_title_tfidf = vec4.transform(preprocessed_cv_title)
X_test_title_tfidf = vec4.transform(preprocessed_test_title)

print("After vectorizations")
print(X_train_title_tfidf.shape, y_train.shape)
print(X_cv_title_tfidf.shape, y_cv.shape)
print(X_test_title_tfidf.shape, y_test.shape)
```

After vectorizations
(51200, 2053) (51200,)
(12800, 2053) (12800,)
(16000, 2053) (16000,)

Numerical features

In [43]:

```
# https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indexes-for-all-groups-in-one-step
price_data = df2.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
price_data.head(2)
```

Out[43]:

	id	price	quantity
0	p000027	782.13	15
1	p000031	357.98	2

In [44]:

```
# join two dataframes in python:
df1 = pd.merge(df1, price_data, on='id', how='left')
```

In [45]:

```
# https://stackoverflow.com/questions/32617811/imputation-of-missing-values-for-categories-in-pandas
#replacing nan with most frequently occurring element
df1['price'] = df1['price'].fillna(df1['price'].mode().iloc[0])
```

In [46]:

```
X_train = pd.merge(X_train, price_data, how = 'left', on = 'id')
X_cv = pd.merge(X_cv, price_data, how = 'left', on = 'id')
X_test = pd.merge(X_test, price_data, how = 'left', on = 'id')
```

In [47]:

```
# https://stackoverflow.com/questions/32617811/imputation-of-missing-values-for-categories-in-pandas
#replacing nan with most frequently occurring element
X_train['price'] = X_train['price'].fillna(X_train['price'].mode().iloc[0])
X_cv['price'] = X_cv['price'].fillna(X_train['price'].mode().iloc[0])
X_test['price'] = X_test['price'].fillna(X_test['price'].mode().iloc[0])
print(X_train['price'].isnull().sum())
print(X_cv['price'].isnull().sum())
print(X_test['price'].isnull().sum())
```

0

0

0

In [48]:

```
# price
# check this one: https://www.youtube.com/watch?v=0H0q0cLn3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html
from sklearn.preprocessing import StandardScaler

# price_standardized = standardScaler.fit(df1['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ...
# 399. 287.73 5.5 ].
# Reshape your data either using array.reshape(-1, 1)

price_scalar = StandardScaler()

X_train_price_stndrd = price_scalar.fit_transform(X_train['price'].values.reshape(-1,1))
X_cv_price_stndrd = price_scalar.transform(X_cv['price'].values.reshape(-1,1))
X_test_price_stndrd = price_scalar.transform(X_test['price'].values.reshape(-1,1))
```

In [49]:

```
X_train_price_stndrd
```

Out[49]:

```
array([[0.09102914],
       [0.09102914],
       [0.09102914],
       ...,
       [0.09102914],
       [0.09102914],
       [0.09102914]])
```

In [50]:

```
X_cv_price_stndrd
```

Out[50]:

```
array([[ 0.09102914],
       [ 0.09102914],
       [-4.39125154],
       ...,
       [ 0.09102914],
       [ 0.09102914],
       [ 0.09102914]])
```

In [51]:

```
X_test_price_stdndr
```

Out[51]:

```
array([[ -1.69162597],
       [ -1.69162597],
       [ -1.69162597],
       ...,
       [ -1.69162597],
       [ -1.69162597],
       [ -1.69162597]])
```

In [52]:

```
df1['quantity'] = df1['quantity'].fillna(df1['quantity'].mode().iloc[0])
```

In [53]:

```
# https://stackoverflow.com/questions/32617811/imputation-of-missing-values-for-categories-in-pandas
#replacing nan with most frequently occurring element

X_train['quantity'] = X_train['quantity'].fillna(X_train['quantity'].mode().iloc[0])
X_cv['quantity'] = X_cv['quantity'].fillna(X_cv['quantity'].mode().iloc[0])
X_test['quantity'] = X_test['quantity'].fillna(X_test['quantity'].mode().iloc[0])
```

In [54]:

```
# quantity
quantity_scalar = StandardScaler()
X_train_quantity_stdndr = quantity_scalar.fit_transform(X_train['quantity'].values.reshape(-1,1))

X_cv_quantity_stdndr = quantity_scalar.transform(X_cv['quantity'].values.reshape(-1,1))
X_test_quantity_stdndr = quantity_scalar.transform(X_test['quantity'].values.reshape(-1,1))
```

In [55]:

```
X_train_quantity_stdndr
```

Out[55]:

```
array([[ -0.11135748],
       [ -0.11135748],
       [ -0.11135748],
       ...,
       [ -0.11135748],
       [ -0.11135748],
       [ -0.11135748]])
```

In [56]:

```
X_cv_quantity_stndrd
```

Out[56]:

```
array([[ -0.11135748],
       [ -0.11135748],
       [  4.10490497],
       ...,
       [ -0.11135748],
       [ -0.11135748],
       [ -0.11135748]])
```

In [57]:

```
X_test_quantity_stndrd
```

Out[57]:

```
array([[ -0.11135748],
       [ -0.11135748],
       [ -0.11135748],
       ...,
       [ -0.11135748],
       [ -0.11135748],
       [ -0.11135748]])
```

In [58]:

```
# previous_year_projects
# finding the mean and standard deviation of this data
price_scalar.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
X_train_prev_proj_stndrd = price_scalar.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

X_test_prev_proj_stndrd = price_scalar.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
X_cv_prev_proj_stndrd = price_scalar.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
```

In [59]:

```
X_train_prev_proj_stndrd
```

Out[59]:

```
array([[ 6.53997261],
       [-0.36636304],
       [ 0.24513543],
       ...,
       [-0.36636304],
       [-0.00665805],
       [-0.36636304]])
```

In [60]:

```
X_cv_prev_proj_stndrd
```

Out[60]:

```
array([[ 0.96454539],
       [-0.36636304],
       [-0.33039254],
       ...,
       [ 6.6119136 ],
       [ 0.13722394],
       [-0.25845154]])
```

In [61]:

```
X_test_prev_proj_stndrd
```

Out[61]:

```
array([[ 1.10842739],
       [-0.40233353],
       [ 0.49692892],
       ...,
       [-0.40233353],
       [ 1.79186685],
       [-0.40233353]])
```

Concatinating all

In [62]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatinating a sparse matrix and a dense matirx
X_train_bow = hstack((X_train_title_bow,X_train_essay_bow,train_category_responseCoding
,train_subcategory_responseCoding,
                    train_grade_responseCoding,train_teacher_responseCoding,train_sta
te_responseCoding,
                    X_train_price_stndrd, X_train_quantity_stndrd, X_train_prev_proj_
stndrd))

print(X_train_bow.shape, y_train.shape)

(51200, 14313) (51200,)
```

In [63]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatinating a sparse matrix and a dense matirx
X_cv_bow = hstack((X_cv_title_bow,X_cv_essay_bow,cv_category_responseCoding,cv_subcateg
ory_responseCoding,
                  cv_grade_responseCoding,cv_teacher_responseCoding,cv_state_respon
seCoding,
                  X_cv_price_stndrd, X_cv_quantity_stndrd, X_cv_prev_proj_stndrd))

print(X_cv_bow.shape, y_cv.shape)

(12800, 14313) (12800,)
```


In [64]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix
X_test_bow = hstack((X_test_title_bow,X_test_essay_bow,test_category_responseCoding,test_subcategory_responseCoding,
                    test_grade_responseCoding,test_teacher_responseCoding,test_state_responseCoding,
                    X_test_price_stdndrd, X_test_quantity_stdndrd, X_test_prev_proj_stdndrd))

print(X_test_bow.shape, y_test.shape)

(16000, 14313) (16000,)
```

In [65]:

```
X_train_bow = X_train_bow.tocsr()
X_cv_bow = X_cv_bow.tocsr()
X_test_bow = X_test_bow.tocsr()
```

In [66]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix
X_train_tfidf = hstack((X_train_title_tfidf,X_train_ess_tfidf,train_category_responseCoding,train_subcategory_responseCoding,
                    train_grade_responseCoding,train_teacher_responseCoding,train_state_responseCoding,
                    X_train_price_stdndrd,X_train_quantity_stdndrd, X_train_prev_proj_stdndrd))

print(X_train_tfidf.shape, y_train.shape)

(51200, 7066) (51200,)
```

In [67]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix
X_cv_tfidf = hstack((X_cv_title_tfidf,X_cv_ess_tfidf,cv_category_responseCoding,cv_subcategory_responseCoding,
                    cv_grade_responseCoding,cv_teacher_responseCoding,cv_state_responseCoding,
                    X_cv_price_stdndrd,X_cv_quantity_stdndrd, X_cv_prev_proj_stdndrd))

print(X_cv_tfidf.shape, y_cv.shape)

(12800, 7066) (12800,)
```

In [68]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix
X_test_tfidf = hstack((X_test_title_tfidf,X_test_ess_tfidf,test_category_responseCoding
,test_subcategory_responseCoding,
                        test_grade_responseCoding,test_teacher_responseCoding,test_state_responseCoding,
                        X_test_price_stdndrd,X_test_quantity_stdndrd, X_test_prev_proj_stdndrd))

print(X_test_tfidf.shape, y_test.shape)
```

(16000, 7066) (16000,)

In [69]:

```
X_train_tfidf = X_train_tfidf.tocsr()
X_cv_tfidf = X_cv_tfidf.tocsr()
X_test_tfidf = X_test_tfidf.tocsr()
```

Pretrained avg_w2v model

In [70]:

```
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039

def loadGloveModel(gloveFile):

    print ("Loading Glove Model")

    f = open(gloveFile,'r', encoding = 'utf8')

    model = {}

    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding

    print ("Done.",len(model)," words loaded!")

    return model
```

In [71]:

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

0it [00:00, ?it/s]

Loading Glove Model

1917494it [09:38, 3311.79it/s]

Done. 1917494 words loaded!

```
glove_words = set(model.keys())
```

```
# compute average word2vec for each review.
def func(wordlist):

    train_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in tqdm(wordlist): # for each review/sentence
        vector = np.zeros(300) # as word vectors are of zero length # we are taking the 300 dimensions very large
        cnt_words = 0; # num of words with a valid vector in the sentence/review
        for word in sentence.split(): # for each word in a review/sentence
            if word in glove_words:
                vector += model[word]
                cnt_words += 1
        if cnt_words != 0:
            vector /= cnt_words
        train_avg_w2v_vectors.append(vector)

    print(len(train_avg_w2v_vectors))
    print(len(train_avg_w2v_vectors[0]))
    return train_avg_w2v_vectors
```

```
# FOR ESSAYS
X_train_avg_w2v_ess=func(preprocessed_train_essay)
X_cv_avg_w2v_ess=func(preprocessed_cv_essay)
X_test_avg_w2v_ess=func(preprocessed_test_essay)
```

```
100%|██████████|  
██████████| 51200/51200 [00:26<00:00, 1942.62it/s]  
4%|█████|  
| 449/12800 [00:00<00:05, 2132.54it/s]  
  
51200  
300  
  
100%|██████████|  
██████████| 12800/12800 [00:06<00:00, 1987.16it/s]  
1%|████|  
| 198/16000 [00:00<00:08, 1910.98it/s]  
  
12800  
300  
  
100%|██████████|  
██████████| 16000/16000 [00:08<00:00, 1933.17it/s]  
  
16000  
300
```

```
X_train_avg_w2v_ess=np.array(X_train_avg_w2v_ess)
X_cv_avg_w2v_ess=np.array(X_cv_avg_w2v_ess)
X_test_avg_w2v_ess =np.array(X_test_avg_w2v_ess)
```

```
# FOR TITLES
X_train_avg_w2v_title=func(preprocessed_train_title)
X_cv_avg_w2v_title=func(preprocessed_cv_title)
X_test_avg_w2v_title=func(preprocessed_test_title)
```

```
X_train_avg_w2v_title=np.asarray(X_train_avg_w2v_title)
X_cv_avg_w2v_title=np.asarray(X_cv_avg_w2v_title)
X_test_avg_w2v_title =np.asarray(X_test_avg_w2v_title)
```

```
X_train_avg_w2v = np.hstack((X_train_avg_w2v_ess,X_train_avg_w2v_title,train_category_r
responseCoding,train_subcategory_responseCoding,
                                train_grade_responseCoding,train_teacher_responseCoding,tra
in_state_responseCoding,
                                X_train_price_stdndr, X_train_quantity_stdndr, X_train_prev
_proj_stdndr ))
```

```
X_cv_avg_w2v = np.hstack((X_cv_avg_w2v_ess,X_cv_avg_w2v_title,cv_category_responseCoding,
cv_subcategory_responseCoding,
cv_grade_responseCoding,cv_teacher_responseCoding,cv_state_
responseCoding,
X_cv_price_stndrd, X_cv_quantity_stndrd, X_cv_prev_proj_stn
drd ))
```

In [80]:

```
X_test_avg_w2v = np.hstack((X_test_avg_w2v_ess,X_test_avg_w2v_title,test_category_respo
nseCoding,test_subcategory_responseCoding,
                             test_grade_responseCoding,test_teacher_responseCoding,test_
state_responseCoding,
                             X_test_price_stndrd, X_test_quantity_stndrd, X_test_prev_pr
oj_stndrd ))
```

In [81]:

```
print(X_train_avg_w2v.shape, y_train.shape)
print(X_cv_avg_w2v.shape, y_cv.shape)
print(X_test_avg_w2v.shape, y_test.shape)
```

```
(51200, 613) (51200,)
(12800, 613) (12800,)
(16000, 613) (16000,)
```

In [82]:

```
X_train_avg_w2v = X_train_avg_w2v[0:20000]
X_cv_avg_w2v = X_cv_avg_w2v[0:20000]
X_test_avg_w2v = X_test_avg_w2v[0:20000]
```

In [83]:

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

In [84]:

```
def func(wordlist):
    # average Word2Vec
    # compute average word2vec for each review.
    tfidf_w2v_ess = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in tqdm(wordlist): # for each review/sentence
        vector = np.zeros(300) # as word vectors are of zero length
        tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
        for word in sentence.split(): # for each word in a review/sentence
            if (word in glove_words) and (word in tfidf_words):
                vec = model[word] # getting the vector for each word
                # here we are multiplying idf value(dictionary[word]) and the tf value
                ((sentence.count(word)/len(sentence.split()))
                 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())))
            # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
        if tf_idf_weight != 0:
            vector /= tf_idf_weight
        tfidf_w2v_ess.append(vector)

    print(len(tfidf_w2v_ess))
    print(len(tfidf_w2v_ess[0]))
    return tfidf_w2v_ess
```

```
#For essays
X_train_tfidf_w2v_ess =func(preprocessed_train_essay)
X_test_tfidf_w2v_ess =func(preprocessed_test_essay)
X_cv_tfidf_w2v_ess =func(preprocessed_cv_essay)
```

```
X_train_tfidf_w2v_ess = np.array(X_train_tfidf_w2v_ess)
X_test_tfidf_w2v_ess = np.array(X_test_tfidf_w2v_ess)
X_cv_tfidf_w2v_ess = np.array(X_cv_tfidf_w2v_ess)
```

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


In [93]:

```
print(X_train_tfidf_w2v.shape, y_train.shape)
print(X_cv_tfidf_w2v.shape, y_cv.shape)
print(X_test_tfidf_w2v.shape, y_test.shape)
```

```
(51200, 613) (51200,)
(12800, 613) (12800,)
(16000, 613) (16000,)
```

In [94]:

```
X_train_tfidf_w2v = X_train_tfidf_w2v[0:20000]
X_cv_tfidf_w2v = X_cv_tfidf_w2v[0:20000]
X_test_tfidf_w2v = X_test_tfidf_w2v[0:20000]
```

Random Forest Classifier on bag of words

In [116]:

```
from sklearn.metrics import f1_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import f1_score
from sklearn.model_selection import RandomizedSearchCV
from scipy.stats import randint as sp_randint
from xgboost.sklearn import XGBClassifier
from scipy.stats import uniform

param_dist = {"n_estimators":[10, 50, 100, 150, 200, 300, 500, 1000],
              "max_depth":[2, 3, 4, 5, 6, 7, 8, 9, 10],
              "min_samples_split": sp_randint(110,190),
              "min_samples_leaf": sp_randint(20,75)}

clf = RandomForestClassifier(random_state=42, class_weight = 'balanced')

rf = RandomizedSearchCV(clf, param_distributions=param_dist,return_train_score = True,
                        n_iter=5,cv=10,scoring='f1',random_state=42)

rf.fit(X_train_bow,y_train)

train_auc= rf.cv_results_['mean_train_score']
train_auc_std= rf.cv_results_['std_train_score']
cv_auc = rf.cv_results_['mean_test_score']
cv_auc_std= rf.cv_results_['std_test_score']
```

In [117]:

```
train_auc
```

Out[117]:

```
array([0.75948755, 0.76957798, 0.76308804, 0.7551134 , 0.77102644])
```


In [118]:

```
train_auc_std
```

Out[118]:

```
array([0.00270147, 0.001997 , 0.00331074, 0.00155028, 0.00351632])
```

In [119]:

```
cv_auc
```

Out[119]:

```
array([0.75095856, 0.75822382, 0.75256765, 0.7472037 , 0.76074461])
```

In [120]:

```
cv_auc_std
```

Out[120]:

```
array([0.00759748, 0.00825959, 0.00610561, 0.00704392, 0.00689054])
```

In [121]:

```
rf.best_estimator_
```

Out[121]:

```
RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight='balanced',
                        criterion='gini', max_depth=7, max_features='auto',
                        max_leaf_nodes=None, max_samples=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=21, min_samples_split=173,
                        min_weight_fraction_leaf=0.0, n_estimators=150,
                        n_jobs=None, oob_score=False, random_state=42, verbose=0,
                        warm_start=False)
```

In [122]:

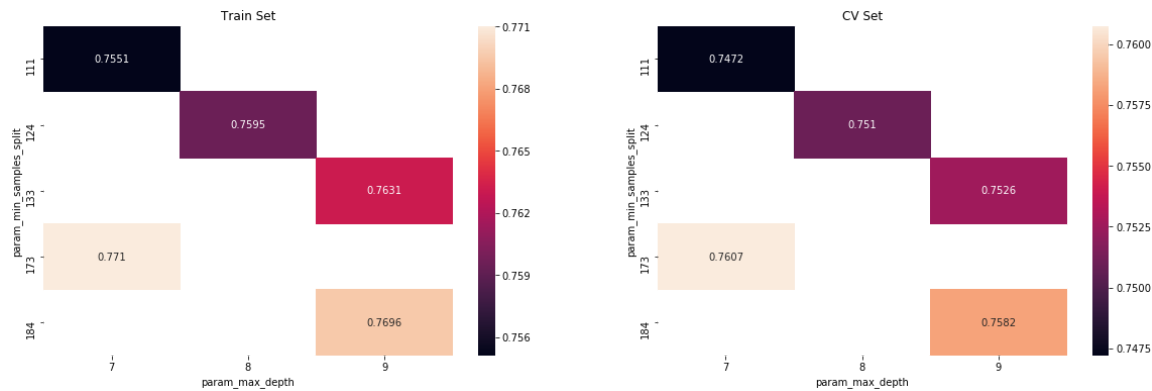
```
max_scores = pd.DataFrame(rf.cv_results_).groupby(['param_min_samples_split', 'param_max_depth']).max().unstack()[['mean_test_score', 'mean_train_score']]

fig, ax = plt.subplots(1,2, figsize=(20,6))

sns.heatmap(max_scores.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sns.heatmap(max_scores.mean_test_score, annot = True, fmt='.4g', ax=ax[1])

ax[0].set_title('Train Set')
ax[1].set_title('CV Set')

plt.show()
```



Hyperparameter tuning

In [124]:

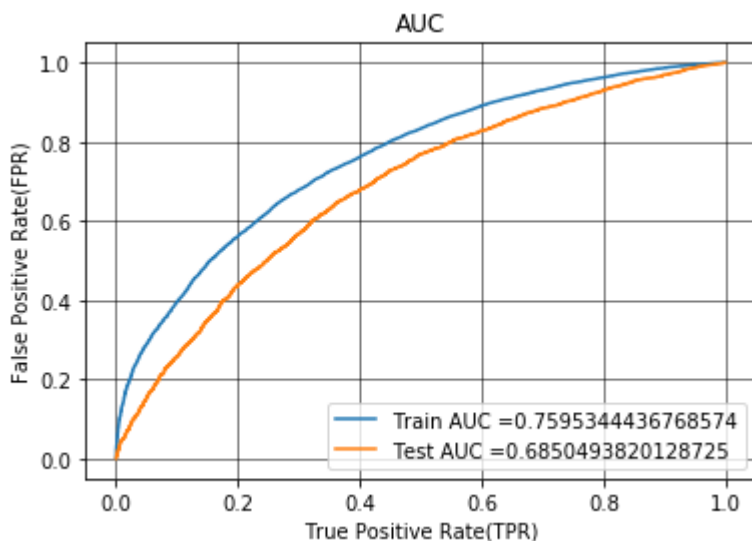
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
model = RandomForestClassifier(max_depth = 7, n_estimators = 150, class_weight= 'balanced')

model.fit(X_train_bow,y_train)

y_train_pred = model.predict_proba(X_train_bow)[:,-1]
y_test_pred = model.predict_proba(X_test_bow)[:,-1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()
```



best threshold

In [125]:

```
# we are writing our own function for predict, with defined threshould
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.rou
nd(t,3))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

confusion matrix

In [126]:

```
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
print('='*100)
```

the maximum value of tpr*(1-fpr) 0.47610066261693607 for threshold 0.5

Train confusion matrix

```
[[ 5538  2223]
 [14456 28983]]
```

Test confusion matrix

```
[[1529  896]
 [4704 8871]]
```

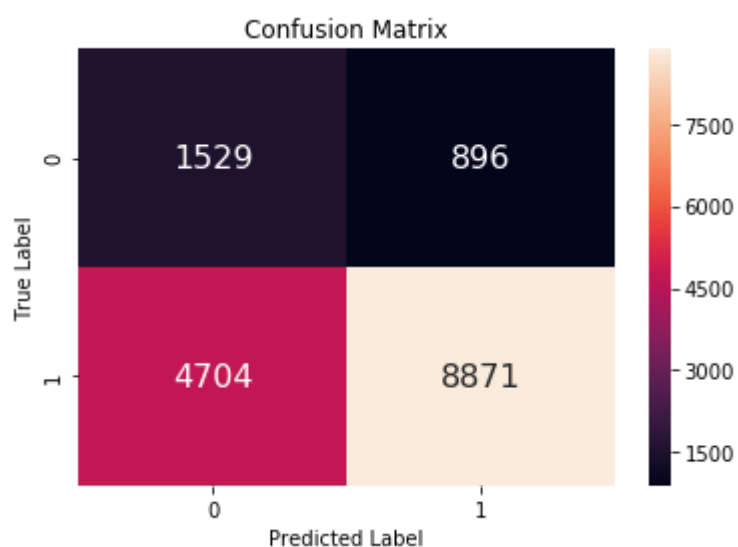
```
=====
=====
```

In [127]:

```
#stackoverflow.com/questions/54018742/valueerror-classification-metrics-cant-handle-a-matrix-of-unknown-and-binary-targets  
matrix = confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))  
sns.heatmap(matrix, annot=True, annot_kws={'size':16}, fmt='g')  
plt.ylabel('True Label')  
plt.xlabel('Predicted Label')  
plt.title('Confusion Matrix')
```

Out[127]:

Text(0.5, 1, 'Confusion Matrix')

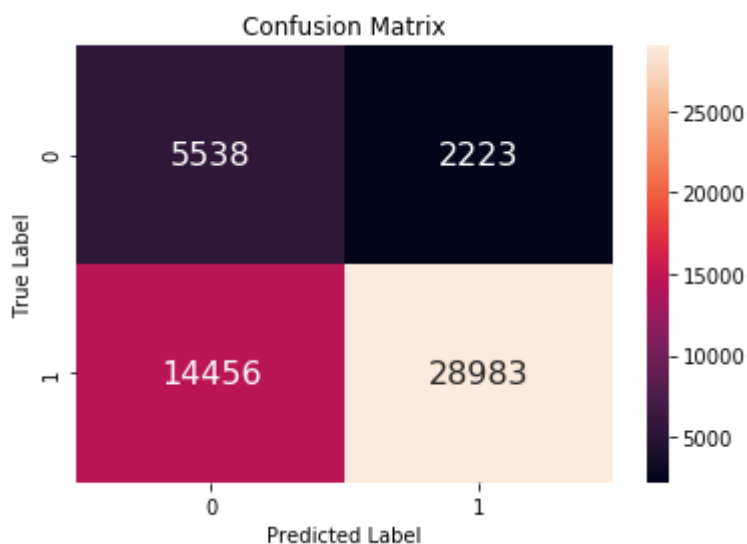


In [128]:

```
#stackoverflow.com/questions/54018742/valueerror-classification-metrics-cant-handle-a-matrix-of-unknown-and-binary-targets  
matrix = confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))  
sns.heatmap(matrix, annot=True, annot_kws={'size':16}, fmt='g')  
plt.ylabel('True Label')  
plt.xlabel('Predicted Label')  
plt.title('Confusion Matrix')
```

Out[128]:

Text(0.5, 1, 'Confusion Matrix')



Random forest classifier on tfidf

In [129]:

```
param_dist = {"n_estimators":[10, 50, 100, 150, 200, 300, 500, 1000],
              "max_depth": [2, 3, 4, 5, 6, 7, 8, 9, 10],
              "min_samples_split": sp_randint(110,190),
              "min_samples_leaf": sp_randint(20,75)}

clf = RandomForestClassifier(random_state=42, class_weight = 'balanced')

rf = RandomizedSearchCV(clf, param_distributions=param_dist,return_train_score = True,
                        n_iter=5,cv=10,scoring='f1',random_state=42)

rf.fit(X_train_tfidf,y_train)

train_auc= rf.cv_results_['mean_train_score']
train_auc_std= rf.cv_results_['std_train_score']
cv_auc = rf.cv_results_['mean_test_score']
cv_auc_std= rf.cv_results_['std_test_score']
```

In [130]:

```
rf.best_estimator_
```

Out[130]:

```
RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight='balanced',
                        criterion='gini', max_depth=9, max_features='auto',
                        max_leaf_nodes=None, max_samples=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=40, min_samples_split=184,
                        min_weight_fraction_leaf=0.0, n_estimators=100,
                        n_jobs=None, oob_score=False, random_state=42, verbose=0,
                        warm_start=False)
```

In [131]:

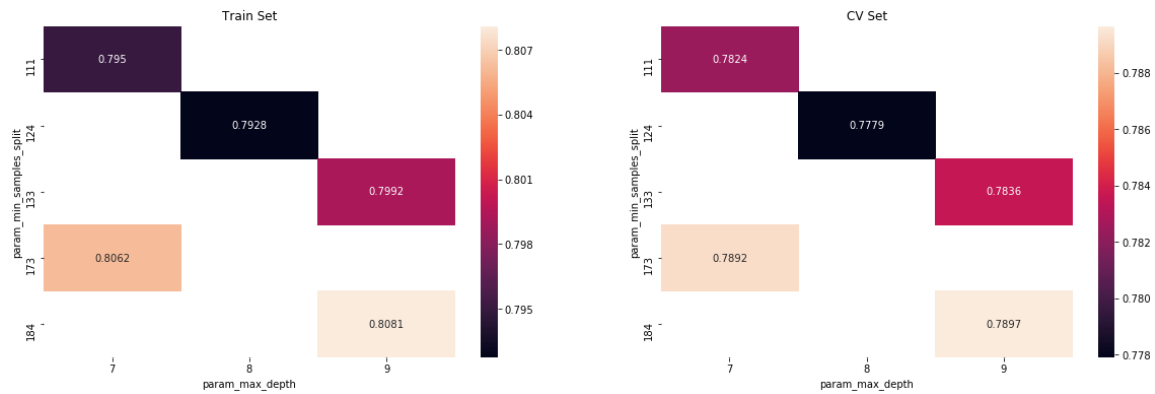
```
max_scores = pd.DataFrame(rf.cv_results_).groupby(['param_min_samples_split', 'param_max_depth']).max().unstack()[['mean_test_score', 'mean_train_score']]

fig, ax = plt.subplots(1,2, figsize=(20,6))

sns.heatmap(max_scores.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sns.heatmap(max_scores.mean_test_score, annot = True, fmt='.4g', ax=ax[1])

ax[0].set_title('Train Set')
ax[1].set_title('CV Set')

plt.show()
```



Hyperparameter tuning

In [132]:

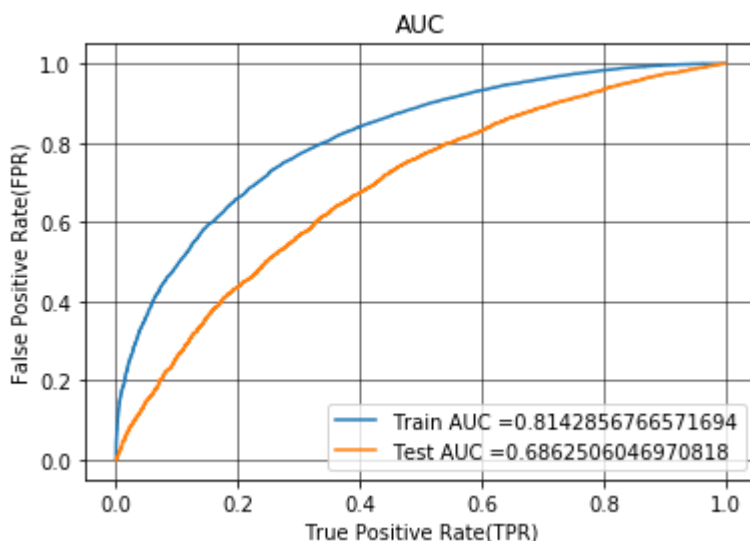
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
model = RandomForestClassifier(max_depth = 9, n_estimators = 100, class_weight= 'balanced')

model.fit(X_train_tfidf,y_train)

y_train_pred = model.predict_proba(X_train_tfidf)[:,-1]
y_test_pred = model.predict_proba(X_test_tfidf)[:,-1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid(color='black', linestyle='--', linewidth=0.5)
plt.show()
```



best threshold

In [133]:

```
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.rou
nd(t,3))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

Confusion matrix

In [134]:

```
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
print('='*100)
```

the maximum value of tpr*(1-fpr) 0.5440109365026429 for threshold 0.505

Train confusion matrix

```
[[ 5775  1986]
 [11681 31758]]
```

Test confusion matrix

```
[[1400 1025]
 [4141 9434]]
```

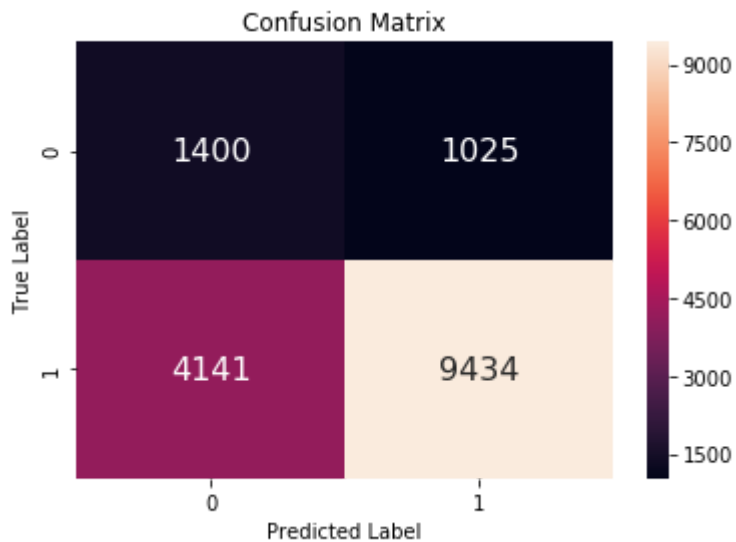
```
=====
=====
```

In [135]:

```
#stackoverflow.com/questions/54018742/valueerror-classification-metrics-cant-handle-a-matrix-of-unknown-and-binary-targets  
matrix = confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))  
sns.heatmap(matrix, annot=True, annot_kws={'size':16}, fmt='g')  
plt.ylabel('True Label')  
plt.xlabel('Predicted Label')  
plt.title('Confusion Matrix')
```

Out[135]:

Text(0.5, 1, 'Confusion Matrix')

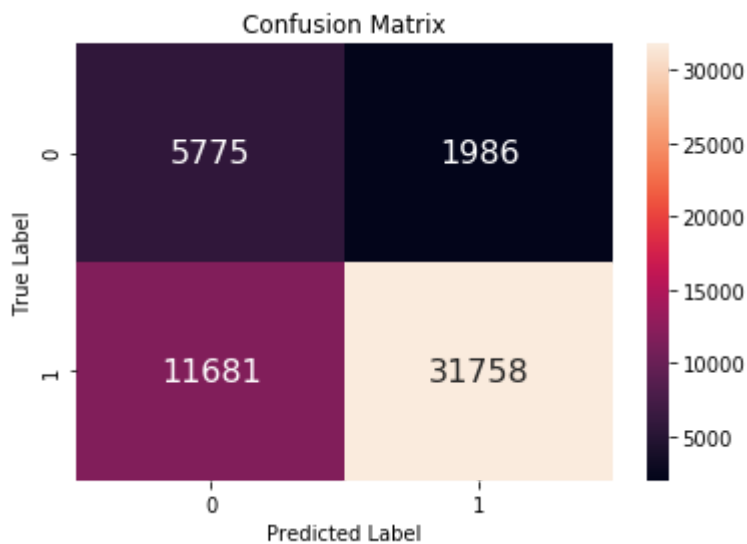


In [136]:

```
#stackoverflow.com/questions/54018742/valueerror-classification-metrics-cant-handle-a-matrix-of-unknown-and-binary-targets
matrix = confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
sns.heatmap(matrix, annot=True, annot_kws={'size':16}, fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Confusion Matrix')
```

Out[136]:

Text(0.5, 1, 'Confusion Matrix')



Random forest classifier on avg_w2v

In [138]:

```
y_train1 = y_train[0:20000]
y_cv1 = y_cv[0:20000]
y_test1 = y_test[0:20000]
```

In [139]:

```
param_dist = {"n_estimators":[10, 50, 100, 150, 200, 300, 500, 1000],
              "max_depth": [2, 3, 4, 5, 6, 7, 8, 9, 10],
              "min_samples_split": sp_randint(110,190),
              "min_samples_leaf": sp_randint(20,75)}

clf = RandomForestClassifier(random_state=42, class_weight = 'balanced')

rf = RandomizedSearchCV(clf, param_distributions=param_dist,return_train_score = True,
                        n_iter=5,cv=10,scoring='f1',random_state=42)

rf.fit(X_train_avg_w2v,y_train1)

train_auc= rf.cv_results_['mean_train_score']
train_auc_std= rf.cv_results_['std_train_score']
cv_auc = rf.cv_results_['mean_test_score']
cv_auc_std= rf.cv_results_['std_test_score']
```

In [140]:

```
rf.best_estimator_
```

Out[140]:

```
RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight='balanced',
                        criterion='gini', max_depth=9, max_features='auto',
                        max_leaf_nodes=None, max_samples=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=72, min_samples_split=133,
                        min_weight_fraction_leaf=0.0, n_estimators=100,
                        n_jobs=None, oob_score=False, random_state=42, verbose=0,
                        warm_start=False)
```

In [141]:

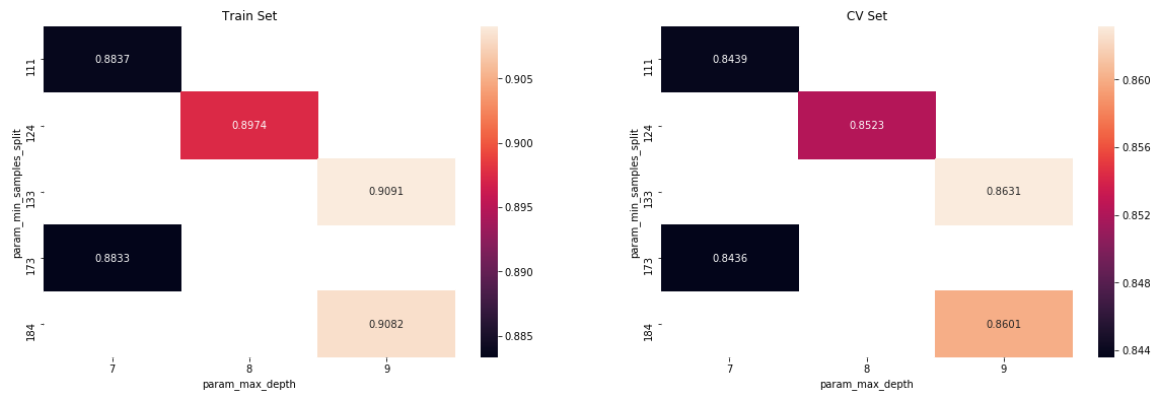
```
max_scores = pd.DataFrame(rf.cv_results_).groupby(['param_min_samples_split', 'param_max_depth']).max().unstack()[['mean_test_score', 'mean_train_score']]

fig, ax = plt.subplots(1,2, figsize=(20,6))

sns.heatmap(max_scores.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sns.heatmap(max_scores.mean_test_score, annot = True, fmt='.4g', ax=ax[1])

ax[0].set_title('Train Set')
ax[1].set_title('CV Set')

plt.show()
```



Hyperparameter tuning

In [142]:

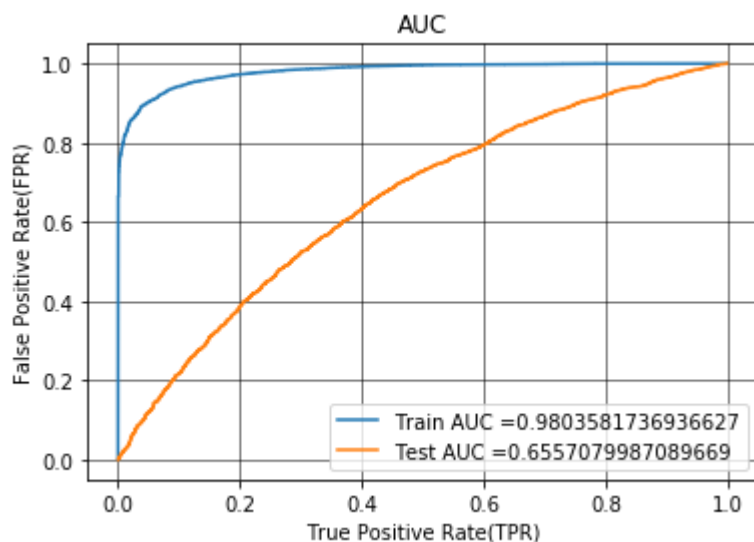
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
model = RandomForestClassifier(max_depth = 9 , n_estimators = 100, class_weight= 'balanced')

model.fit(X_train_avg_w2v, y_train1)

y_train_pred = model.predict_proba(X_train_avg_w2v)[: , 1]
y_test_pred = model.predict_proba(X_test_avg_w2v)[: , 1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train1, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test1, y_test_pred)

plt.plot(train_fpr, train_tpr, label="Train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()
```



best threshold

In [143]:

```
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.rou
nd(t,3))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

confusion matrix

In [144]:

```
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train1, predict_with_best_t(y_train_pred, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test1, predict_with_best_t(y_test_pred, best_t)))
print('='*100)
```

the maximum value of tpr*(1-fpr) 0.8583300568527676 for threshold 0.544

Train confusion matrix

```
[[ 2855   137]
 [ 1709 15299]]
```

Test confusion matrix

```
[[   863   1562]
 [ 2278 11297]]
```

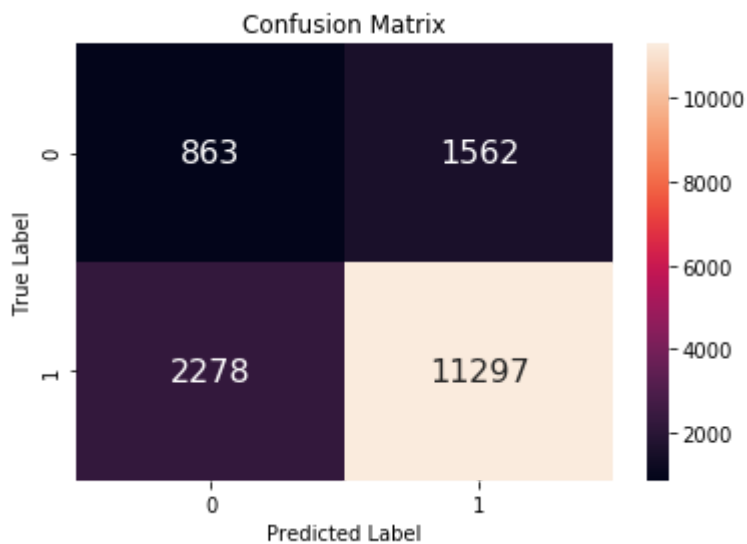
```
=====
=====
```


In [145]:

```
#stackoverflow.com/questions/54018742/valueerror-classification-metrics-cant-handle-a-matrix-of-unknown-and-binary-targets  
matrix = confusion_matrix(y_test1, predict_with_best_t(y_test_pred, best_t))  
sns.heatmap(matrix, annot=True, annot_kws={'size':16}, fmt='g')  
plt.ylabel('True Label')  
plt.xlabel('Predicted Label')  
plt.title('Confusion Matrix')
```

Out[145]:

Text(0.5, 1, 'Confusion Matrix')

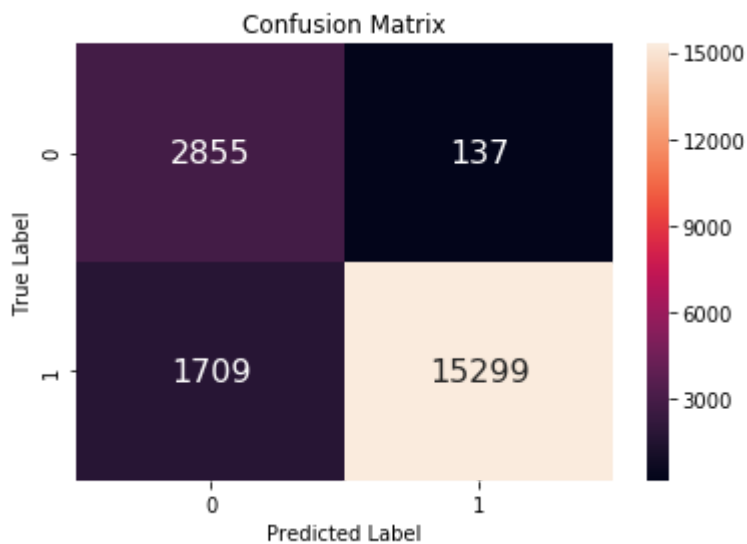


In [146]:

```
#stackoverflow.com/questions/54018742/valueerror-classification-metrics-cant-handle-a-matrix-of-unknown-and-binary-targets  
matrix = confusion_matrix(y_train1, predict_with_best_t(y_train_pred, best_t))  
sns.heatmap(matrix, annot=True, annot_kws={'size':16}, fmt='g')  
plt.ylabel('True Label')  
plt.xlabel('Predicted Label')  
plt.title('Confusion Matrix')
```

Out[146]:

Text(0.5, 1, 'Confusion Matrix')



Random forest classifier on tfidf_w2v

In [147]:

```
param_dist = {"n_estimators":[10, 50, 100, 150, 200, 300, 500, 1000],
              "max_depth": [2, 3, 4, 5, 6, 7, 8, 9, 10],
              "min_samples_split": sp_randint(110,190),
              "min_samples_leaf": sp_randint(20,75)}

clf = RandomForestClassifier(random_state=42, class_weight = 'balanced')

rf = RandomizedSearchCV(clf, param_distributions=param_dist,return_train_score = True,
                        n_iter=5,cv=10,scoring='f1',random_state=42)

rf.fit(X_train_tfidf_w2v,y_train1)

train_auc= rf.cv_results_['mean_train_score']
train_auc_std= rf.cv_results_['std_train_score']
cv_auc = rf.cv_results_['mean_test_score']
cv_auc_std= rf.cv_results_['std_test_score']
```

In [149]:

```
rf.best_estimator_
```

Out[149]:

```
RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight='balanced',
                        criterion='gini', max_depth=9, max_features='auto',
                        max_leaf_nodes=None, max_samples=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=72, min_samples_split=133,
                        min_weight_fraction_leaf=0.0, n_estimators=100,
                        n_jobs=None, oob_score=False, random_state=42, verbose=0,
                        warm_start=False)
```

In [150]:

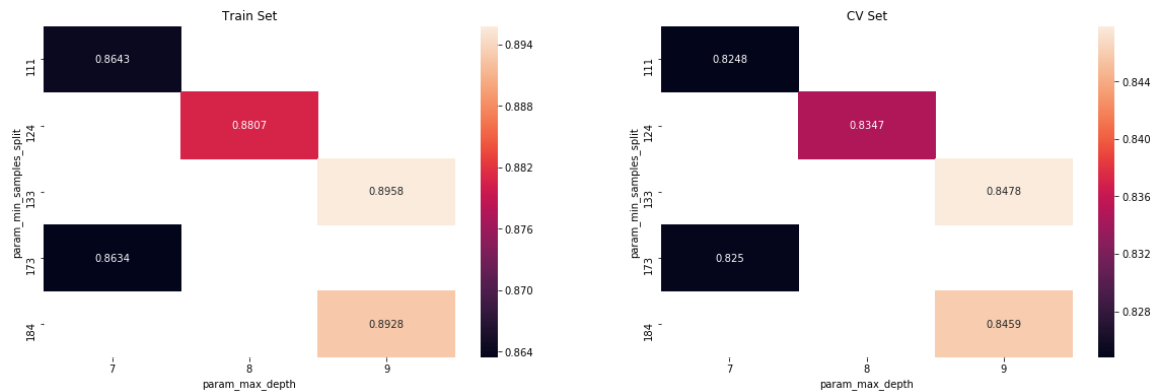
```
max_scores = pd.DataFrame(rf.cv_results_).groupby(['param_min_samples_split', 'param_max_depth']).max().unstack()[['mean_test_score', 'mean_train_score']]

fig, ax = plt.subplots(1,2, figsize=(20,6))

sns.heatmap(max_scores.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sns.heatmap(max_scores.mean_test_score, annot = True, fmt='.4g', ax=ax[1])

ax[0].set_title('Train Set')
ax[1].set_title('CV Set')

plt.show()
```



Hyperparameter tuning

In [151]:

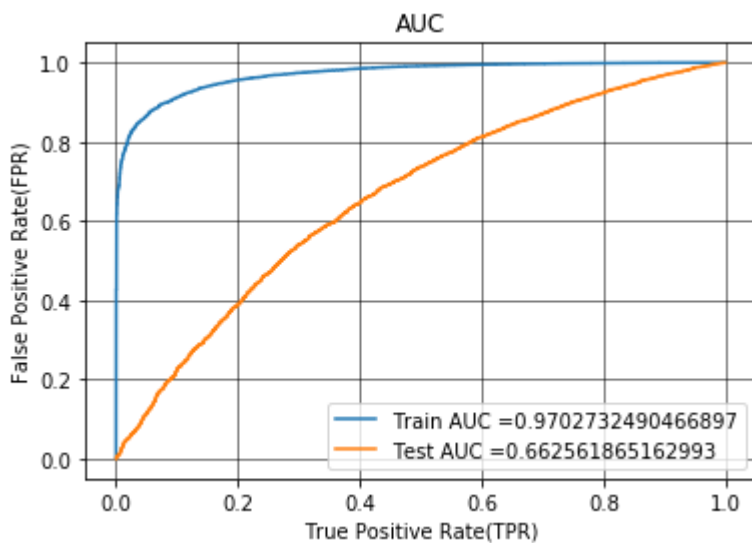
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
model = RandomForestClassifier(max_depth = 9, n_estimators = 100, class_weight= 'balanced')

model.fit(X_train_tfidf_w2v,y_train1)

y_train_pred = model.predict_proba(X_train_tfidf_w2v)[:,-1]
y_test_pred = model.predict_proba(X_test_tfidf_w2v)[:,-1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train1, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test1, y_test_pred)

plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid(color='black', linestyle='--', linewidth=0.5)
plt.show()
```



best threshold

In [152]:

```
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.rou
nd(t,3))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

confusion matrix

In [153]:

```
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train1, predict_with_best_t(y_train_pred, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test1, predict_with_best_t(y_test_pred, best_t)))
print('='*100)
```

the maximum value of $tpr*(1-fpr)$ 0.8296152549790976 for threshold 0.521

Train confusion matrix

```
[[ 2776   216]
```

```
 [ 1800 15208]]
```

Test confusion matrix

```
[[   868  1557]
```

```
 [ 2177 11398]]
```

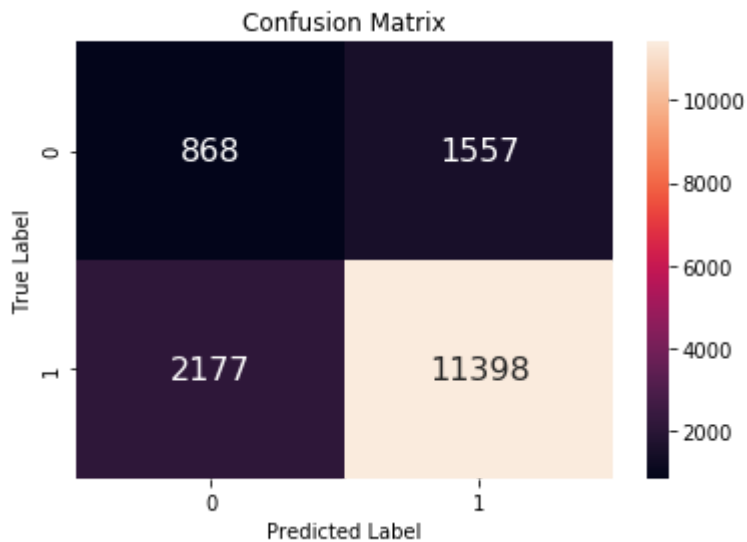
```
=====
=====
```

In [154]:

```
#stackoverflow.com/questions/54018742/valueerror-classification-metrics-cant-handle-a-matrix-of-unknown-and-binary-targets  
matrix = confusion_matrix(y_test1, predict_with_best_t(y_test_pred, best_t))  
sns.heatmap(matrix, annot=True, annot_kws={'size':16}, fmt='g')  
plt.ylabel('True Label')  
plt.xlabel('Predicted Label')  
plt.title('Confusion Matrix')
```

Out[154]:

Text(0.5, 1, 'Confusion Matrix')

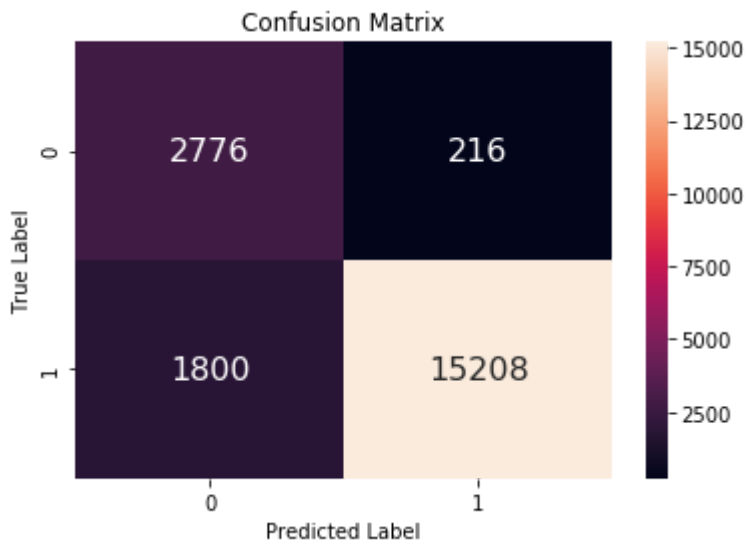


In [156]:

```
#stackoverflow.com/questions/54018742/valueerror-classification-metrics-cant-handle-a-matrix-of-unknown-and-binary-targets  
matrix = confusion_matrix(y_train1, predict_with_best_t(y_train_pred, best_t))  
sns.heatmap(matrix, annot=True, annot_kws={'size':16}, fmt='g')  
plt.ylabel('True Label')  
plt.xlabel('Predicted Label')  
plt.title('Confusion Matrix')
```

Out[156]:

Text(0.5, 1, 'Confusion Matrix')



Summary

In [158]:

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

ptable = PrettyTable()
ptable.title = 'Classification Report'

ptable.field_names = ["Vectorization", "Model", "max_depth", "n_estimator", "AUC"]

ptable.add_row(["BOW", "RF", 7, 150, 68.50])
ptable.add_row(["tf-idf", "RF", 9, 100, 68.62])
ptable.add_row(["avg-w2v", "RF", 9, 100, 65.57])
ptable.add_row(["tf-idf-w2v", "RF", 9, 100, 66.25])

print(ptable)
```

Vectorization	Model	max_depth	n_estimator	AUC
BOW	RF	7	150	68.5
tf-idf	RF	9	100	68.62
avg-w2v	RF	9	100	65.57
tf-idf-w2v	RF	9	100	66.25

Initially after loading the dataset if any null values exists replace them with most occurring element then split the data into training, validation, testing data and preprocessed the data to avoid the leakage.

Applied bag of words, tfidf, avg_w2v, tfidf_w2v featurising on the data. After concatenating all the features applied random forest on each.

More accuracy is obtained for tfidf.