

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: p036502
<code>project_title</code>	<ul style="list-style-type: none">••	Title of the project. Examples: Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: Grades PreK-2 Grades 3-5 Grades 6-8 Grades 9-12

Feature	Description	
	One or more (comma-separated) subject categories for the project from the following enumerated list of values:	
project_subject_categories	<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 	
school_state	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: WY	
	One or more (comma-separated) subject subcategories for the project. Examples:	
project_subject_subcategories	<ul style="list-style-type: none"> • Literacy • Literature & Writing, Social Sciences 	
	An explanation of the resources needed for the project. Example:	
project_resource_summary	<ul style="list-style-type: none"> • My students need hands on literacy materials to manage sensory needs!</code 	
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: 2016-04-28 12:43:56.245	
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56	
	Teacher's title. One of the following enumerated values:	
teacher_prefix	<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 same teacher. Example: 2	

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

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

Feature	Description
<code>id</code>	A <code>project_id</code> value from the <code>train.csv</code> file. Example: <code>p036502</code>
<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 [0]:

```
1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4 import sqlite3
5 import pandas as pd
6 import numpy as np
7 import nltk
8 import string
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11 from sklearn.feature_extraction.text import TfidfTransformer
12 from sklearn.feature_extraction.text import TfidfVectorizer
13 from sklearn.feature_extraction.text import CountVectorizer
14 from sklearn.metrics import confusion_matrix
15 from sklearn import metrics
16 from sklearn.metrics import roc_curve, auc
17 from nltk.stem.porter import PorterStemmer
18 import re
19 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
20 import string
21 from nltk.corpus import stopwords
22 from nltk.stem import PorterStemmer
23 from nltk.stem.wordnet import WordNetLemmatizer
24 from gensim.models import Word2Vec
25 from gensim.models import KeyedVectors
26 import pickle
27 from tqdm import tqdm
28 import os
29 import chart_studio.plotly
30 # from plotly import plotly
31 import plotly.offline as offline
32 import plotly.graph_objs as go
33 offline.init_notebook_mode()
34 from collections import Counter
35 from scipy.sparse import hstack, vstack
36 from sklearn.model_selection import train_test_split
37 from sklearn.neighbors import KNeighborsClassifier
38 from sklearn.metrics import accuracy_score
39 from sklearn.model_selection import cross_val_score
40 from sklearn import model_selection
41 from sklearn.preprocessing import StandardScaler
42 from sklearn.model_selection import RandomizedSearchCV
43 #from sklearn.impute import SimpleImputer
44 from sklearn.datasets import load_digits
45 #from sklearn.feature_selection import SelectKBest, chi2
```

```

46 from sklearn.model_selection import GridSearchCV
47 from sklearn.feature_selection import SelectKBest,f_classif
48 from prettytable import PrettyTable
49 import pdb

```

1.1 Reading Data

```

In [0]: 1 Project_data = pd.read_csv('train_data25K.csv')
        2 Resource_data = pd.read_csv('resources.csv')
        3 print(Project_data.shape)
        4 print(Resource_data.shape)

```

```

(25000, 17)
(1541272, 4)

```

```

In [0]: 1 # how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
        2 cols = ['Date' if x=='project_submitted_datetime' else x for x in list(Project_data.columns)]
        3 #sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
        4 Project_data['Date'] = pd.to_datetime(Project_data['project_submitted_datetime'])
        5 Project_data.drop('project_submitted_datetime', axis=1, inplace=True)
        6 Project_data.sort_values(by=['Date'], inplace=True)
        7 # how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
        8 Project_data = Project_data[cols]
        9 Project_data.head(2)

```

Out[77]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_subject_categories	project
3287	159755	p147002	6ada7036aeb258d3653589d1f2a5b815	Mrs.	CA	2016-01-05 02:02:00	Grades 3-5	Literacy & Language, Special Needs	
19437	146532	p024903	55f60249d65840ee198285acdc455838	Mrs.	CA	2016-01-05 02:57:00	Grades 3-5	Math & Science, Literacy & Language	Health

1.2 preprocessing of project_subject_categories

```
In [0]: 1 categories = list(Project_data['project_subject_categories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
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-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7 cat_list = []
8 for i in categories:
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", "Care & Hunger"]
12        if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
13            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
14            j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
15            temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
16            temp = temp.replace('&','_') # we are replacing the & value into
17    cat_list.append(temp.strip())
18
19 Project_data['clean_categories'] = cat_list
20 Project_data.drop(['project_subject_categories'], axis=1, inplace=True)
21
22 from collections import Counter
23 my_counter = Counter()
24 for word in Project_data['clean_categories'].values:
25     my_counter.update(word.split())
26
27 cat_dict = dict(my_counter)
28 sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
29
```

1.3 preprocessing of project_subject_subcategories

```

In [0]: 1 sub_categories = list(Project_data['project_subject_subcategories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
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-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
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", "Care & Hunger"]
13         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
14             j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
15             j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
16             temp +=j.strip()+" #" "abc ".strip() will return "abc", remove the trailing spaces
17             temp = temp.replace('&','_')
18     sub_cat_list.append(temp.strip())
19
20 Project_data['clean_subcategories'] = sub_cat_list
21 Project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
22
23 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
24 my_counter = Counter()
25 for word in Project_data['clean_subcategories'].values:
26     my_counter.update(word.split())
27
28 sub_cat_dict = dict(my_counter)
29 sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))

```

```

In [0]: 1 y = Project_data['project_is_approved'].values
2 Project_data.drop(['project_is_approved'], axis=1, inplace=True)
3 X = Project_data

```

```
In [0]: 1 #Splitting the Dataset into three Train,CV and Test
2 X1, X_Test, Y1, Y_Test = train_test_split(X, y, test_size=0.33, random_state=0)
3 X_Train, X_CV, Y_Train, Y_CV = train_test_split(X1, Y1, test_size=0.33, random_state=0)
4 print('Shape of the X_Train data is {0} and Y_Train data is: {1}'.format(X_Train.shape,Y_Train.shape[0]))
5 print('Shape of the X_CV data is {0} and Y_CV data is : {1}'.format(X_CV.shape,Y_CV.shape[0]))
6 print('Shape of the X_Test data is {0} and Y_Test data is : {1}'.format(X_Test.shape,Y_Test.shape[0]))
```

Shape of the X_Train data is (11222, 16) and Y_Train data is: 11222

Shape of the X_CV data is (5528, 16) and Y_CV data is : 5528

Shape of the X_Test data is (8250, 16) and Y_Test data is : 8250

1.3 Text preprocessing

```
In [0]: 1 # merge two column text dataframe:
2 X_Train["essay"] = X_Train["project_essay_1"].map(str) + \
3                 X_Train["project_essay_2"].map(str) + \
4                 X_Train["project_essay_3"].map(str) + \
5                 X_Train["project_essay_4"].map(str)
6
7 X_CV["essay"] = X_CV["project_essay_1"].map(str) + \
8                X_CV["project_essay_2"].map(str) + \
9                X_CV["project_essay_3"].map(str) + \
10               X_CV["project_essay_4"].map(str)
11
12 X_Test["essay"] = X_Test["project_essay_1"].map(str) + \
13                  X_Test["project_essay_2"].map(str) + \
14                  X_Test["project_essay_3"].map(str) + \
15                  X_Test["project_essay_4"].map(str)
```


In [0]:

```
1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8     # general
9     phrase = re.sub(r"n't", " not", phrase)
10    phrase = re.sub(r"\ 're", " are", phrase)
11    phrase = re.sub(r"\ 's", " is", phrase)
12    phrase = re.sub(r"\ 'd", " would", phrase)
13    phrase = re.sub(r"\ 'll", " will", phrase)
14    phrase = re.sub(r"\ 't", " not", phrase)
15    phrase = re.sub(r"\ 've", " have", phrase)
16    phrase = re.sub(r"\ 'm", " am", phrase)
17    return phrase
```

In [0]:

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

In [0]:

```
1 # Combining all the above stundents
2 # tqdm is for printing the status bar
3
4 #-----PreProcessing of Essays in Train data set-----
5 preprocessed_essays_Train = []
6 for sentence in tqdm(X_Train['essay'].values):
7     sent = decontracted(sentence)
8     sent = sent.replace('\\r', ' ')
9     sent = sent.replace('\\\"', ' ')
10    sent = sent.replace('\\n', ' ')
11    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12    # https://gist.github.com/sebleier/554280
13    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
14    preprocessed_essays_Train.append(sent.lower().strip())
15    # pdb.set_trace()
16
17 #-----PreProcessing of Essays in CV data set-----
18 preprocessed_essays_CV = []
19 for sentence in tqdm(X_CV['essay'].values):
20     sent = decontracted(sentence)
21     sent = sent.replace('\\r', ' ')
22     sent = sent.replace('\\\"', ' ')
23     sent = sent.replace('\\n', ' ')
24     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
25     # https://gist.github.com/sebleier/554280
26     sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
27     preprocessed_essays_CV.append(sent.lower().strip())
28     # pdb.set_trace()
29
30 #-----PreProcessing of Essays in Test data set-----
31 preprocessed_essays_Test = []
32 for sentence in tqdm(X_Test['essay'].values):
33     sent = decontracted(sentence)
34     sent = sent.replace('\\r', ' ')
35     sent = sent.replace('\\\"', ' ')
36     sent = sent.replace('\\n', ' ')
37     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
38     # https://gist.github.com/sebleier/554280
39     sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
40     preprocessed_essays_Test.append(sent.lower().strip())
41    # pdb.set_trace()
```

100%|██████████| 11222/11222 [00:05<00:00, 2018.50it/s]

100%|██████████| 5528/5528 [00:02<00:00, 2027.93it/s]

1.4 Preprocessing of project_title

In [0]:

```
1 # Combining all the above stundents
2 # tqdm is for printing the status bar
3
4 #-----PreProcessing of Project Title in Train data set-----
5 preprocessed_titles_Train = []
6 for sentence in tqdm(X_Train['project_title'].values):
7     sent = decontracted(sentence)
8     sent = sent.replace('\\r', ' ')
9     sent = sent.replace('\\\"', ' ')
10    sent = sent.replace('\\n', ' ')
11    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12    # https://gist.github.com/sebleier/554280
13    sent = ' '.join(e for e in sent.split() if e not in stopwords)
14    preprocessed_titles_Train.append(sent.lower().strip())
15    # pdb.set_trace()
16
17 #-----PreProcessing of Project Title in CV data set-----
18 preprocessed_titles_CV = []
19 for sentence in tqdm(X_CV['project_title'].values):
20     sent = decontracted(sentence)
21     sent = sent.replace('\\r', ' ')
22     sent = sent.replace('\\\"', ' ')
23     sent = sent.replace('\\n', ' ')
24     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
25     # https://gist.github.com/sebleier/554280
26     sent = ' '.join(e for e in sent.split() if e not in stopwords)
27     preprocessed_titles_CV.append(sent.lower().strip())
28     # pdb.set_trace()
29
30 #-----PreProcessing of Project Title in Test data set-----
31 preprocessed_titles_Test = []
32 for sentence in tqdm(X_Test['project_title'].values):
33     sent = decontracted(sentence)
34     sent = sent.replace('\\r', ' ')
35     sent = sent.replace('\\\"', ' ')
36     sent = sent.replace('\\n', ' ')
37     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
38     # https://gist.github.com/sebleier/554280
39     sent = ' '.join(e for e in sent.split() if e not in stopwords)
40     preprocessed_titles_Test.append(sent.lower().strip())
41    # pdb.set_trace()
```

100%|██████████| 11222/11222 [00:00<00:00, 44587.37it/s]

100%|██████████| 5528/5528 [00:00<00:00, 43345.50it/s]

100%|██████████| 8250/8250 [00:00<00:00, 43852.51it/s]

1.5 Preparing data for models

1.5.1 Vectorizing Categorical data

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

In [0]:

```
1 #-----Vectorizing categorical data for Train,CV and Test-----
2
3 # we use count vectorizer to convert the values into one hot encoding
4 vectorizer_cat = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
5 categories_one_hot_Train = vectorizer_cat.fit_transform(X_Train['clean_categories'].values)
6 categories_one_hot_CV = vectorizer_cat.transform(X_CV['clean_categories'].values)
7 categories_one_hot_Test = vectorizer_cat.transform(X_Test['clean_categories'].values)
8 print(vectorizer_cat.get_feature_names())
9 print("-"*120)
10 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(categories_one_hot_Train.shape))
11 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(categories_one_hot_CV.shape))
12 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(categories_one_hot_Test.shape))
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
```

```
Shape of Train dataset matrix after one hot encoding is: (11222, 9)
```

```
Shape of CV dataset matrix after one hot encoding is: (5528, 9)
```

```
Shape of Test dataset matrix after one hot encoding is: (8250, 9)
```

```
In [0]: 1  #-----Vectorizing categorical data for Train,CV and Test-----
2
3  # we use count vectorizer to convert the values into one hot encoding
4  vectorizer_sub_cat = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
5  sub_categories_one_hot_Train = vectorizer_sub_cat.fit_transform(X_Train['clean_subcategories'].values)
6  sub_categories_one_hot_CV = vectorizer_sub_cat.transform(X_CV['clean_subcategories'].values)
7  sub_categories_one_hot_Test = vectorizer_sub_cat.transform(X_Test['clean_subcategories'].values)
8  print(vectorizer_sub_cat.get_feature_names())
9  print("-"*120)
10 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_Train.shape))
11 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_CV.shape))
12 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_Test.shape))
```

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular', 'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger', 'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other', 'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'ESL', 'EarlyDevelopment', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences', 'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
```

```
-----
Shape of Train dataset matrix after one hot encoding is: (11222, 30)
Shape of CV dataset matrix after one hot encoding is: (5528, 30)
Shape of Test dataset matrix after one hot encoding is: (8250, 30)
```

School State

```

In [0]: 1  #-----Vectorizing categorical data of School state for Train dataset-----
2
3  school_catogories_Train = list(X_Train['school_state'].values)
4  school_list_Train = []
5  for sent in school_catogories_Train:
6      school_list_Train.append(sent.lower().strip())
7  X_Train['school_categories'] = school_list_Train
8  X_Train.drop(['school_state'], axis=1, inplace=True)
9
10 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
11 my_counter_school_Train = Counter()
12 for word in X_Train['school_categories'].values:
13     my_counter_school_Train.update(word.split())
14
15 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
16 school_dict_Train = dict(my_counter_school_Train)
17 sorted_school_dict_Train = dict(sorted(school_dict_Train.items(), key=lambda kv: kv[1]))
18
19 vectorizer_school = CountVectorizer(vocabulary=list(sorted_school_dict_Train.keys()), lowercase=False, binary=True)
20 vectorizer_school.fit(X_Train['school_categories'].values)
21 #print(vectorizer.get_feature_names())
22
23 school_one_hot_Train = vectorizer_school.transform(X_Train['school_categories'].values)
24
25 #-----Vectorizing categorical data of School state for CV dataset-----
26
27 school_catogories_CV = list(X_CV['school_state'].values)
28 school_list_CV = []
29 for sent in school_catogories_CV:
30     school_list_CV.append(sent.lower().strip())
31 X_CV['school_categories'] = school_list_CV
32 X_CV.drop(['school_state'], axis=1, inplace=True)
33
34 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
35 my_counter_school_CV = Counter()
36 for word in X_CV['school_categories'].values:
37     my_counter_school_CV.update(word.split())
38
39 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
40 school_dict_CV = dict(my_counter_school_CV)
41 sorted_school_dict_CV = dict(sorted(school_dict_CV.items(), key=lambda kv: kv[1]))
42 school_one_hot_CV = vectorizer_school.transform(X_CV['school_categories'].values)
43
44 #-----Vectorizing categorical data of School state for Test dataset-----
45

```

```

46 school_catogories_Test = list(X_Test['school_state'].values)
47 school_list_Test = []
48 for sent in school_catogories_Test:
49     school_list_Test.append(sent.lower().strip())
50 X_Test['school_categories'] = school_list_Test
51 X_Test.drop(['school_state'], axis=1, inplace=True)
52
53 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
54 my_counter_school_Test = Counter()
55 for word in X_Test['school_categories'].values:
56     my_counter_school_Test.update(word.split())
57
58 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
59 school_dict_Test = dict(my_counter_school_Test)
60 sorted_school_dict_Test = dict(sorted(school_dict_Test.items(), key=lambda kv: kv[1]))
61 school_one_hot_Test = vectorizer_school.transform(X_Test['school_categories'].values)
62 print("-"*120)
63 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(school_one_hot_Train.shape))
64 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(school_one_hot_CV.shape))
65 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(school_one_hot_Test.shape))

```

```

Shape of Train dataset matrix after one hot encoding is: (11222, 51)
Shape of CV dataset matrix after one hot encoding is: (5528, 51)
Shape of Test dataset matrix after one hot encoding is: (8250, 51)

```

Prefix

In [0]:

```
1  #-----Vectorizing categorical data of Teacher Prefix for Train dataset-----
2
3  # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
4  # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5  # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6  # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7  prefix_catogories_Train = list(X_Train['teacher_prefix'].values)
8  prefix_list_Train = []
9  for sent in prefix_catogories_Train:
10     sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
11     # https://gist.github.com/sebleier/554280
12     sent = ' '.join(e for e in sent.split())
13     prefix_list_Train.append(sent.lower().strip())
14  X_Train['prefix_catogories'] = prefix_list_Train
15  X_Train.drop(['teacher_prefix'], axis=1, inplace=True)
16
17  # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
18  my_counter_prefix_Train = Counter()
19  for word in X_Train['prefix_catogories'].values:
20     my_counter_prefix_Train.update(word.split())
21
22  # dict sort by value python: https://stackoverflow.com/a/613218/4084039
23  prefix_dict_Train = dict(my_counter_prefix_Train)
24  sorted_prefix_dict_Train = dict(sorted(prefix_dict_Train.items(), key=lambda kv: kv[1]))
25
26
27  vectorizer_prefix = CountVectorizer(vocabulary=list(sorted_prefix_dict_Train.keys()), lowercase=False, binary=True)
28  vectorizer_prefix.fit(X_Train['prefix_catogories'].values)
29  #print(vectorizer.get_feature_names())
30
31  prefix_one_hot_Train = vectorizer_prefix.transform(X_Train['prefix_catogories'].values)
32  #print("Shape of matrix after one hot encodig ",prefix_one_hot.shape)
33
34  #-----Vectorizing categorical data of Teacher Prefix for CV dataset-----
35
36  prefix_catogories_CV = list(X_CV['teacher_prefix'].values)
37  prefix_list_CV = []
38  for sent in prefix_catogories_CV:
39     sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
40     # https://gist.github.com/sebleier/554280
41     sent = ' '.join(e for e in sent.split())
42     prefix_list_CV.append(sent.lower().strip())
43  X_CV['prefix_catogories'] = prefix_list_CV
44  X_CV.drop(['teacher_prefix'], axis=1, inplace=True)
45
```

```

46 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
47 my_counter_prefix_CV = Counter()
48 for word in X_CV['prefix_catogories'].values:
49     my_counter_prefix_CV.update(word.split())
50
51 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
52 prefix_dict_CV = dict(my_counter_prefix_CV)
53 sorted_prefix_dict_CV = dict(sorted(prefix_dict_CV.items(), key=lambda kv: kv[1]))
54 prefix_one_hot_CV = vectorizer_prefix.transform(X_CV['prefix_catogories'].values)
55
56 #-----Vectorizing categorical data of Teacher Prefix for Test dataset-----
57
58 prefix_catogories_Test = list(X_Test['teacher_prefix'].values)
59 prefix_list_Test = []
60 for sent in prefix_catogories_Test:
61     sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
62     # https://gist.github.com/sebleier/554280
63     sent = ' '.join(e for e in sent.split())
64     prefix_list_Test.append(sent.lower().strip())
65 X_Test['prefix_catogories'] = prefix_list_Test
66 X_Test.drop(['teacher_prefix'], axis=1, inplace=True)
67
68 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
69 my_counter_prefix_Test = Counter()
70 for word in X_Test['prefix_catogories'].values:
71     my_counter_prefix_Test.update(word.split())
72
73 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
74 prefix_dict_Test = dict(my_counter_prefix_Test)
75 sorted_prefix_dict_Test = dict(sorted(prefix_dict_Test.items(), key=lambda kv: kv[1]))
76 prefix_one_hot_Test = vectorizer_prefix.transform(X_Test['prefix_catogories'].values)
77 print("-"*120)
78 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_Train.shape))
79 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_CV.shape))
80 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_Test.shape))

```

```

-----
Shape of Train dataset matrix after one hot encoding is: (11222, 4)
Shape of CV dataset matrix after one hot encoding is: (5528, 4)
Shape of Test dataset matrix after one hot encoding is: (8250, 4)

```

project_grade_category

In [0]:

```
1  #-----Vectorizing categorical data of Project Grade for Train dataset-----
2
3  # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
4  # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5  # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6  # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7  grade_catogories_Train = list(X_Train['project_grade_category'].values)
8  grade_list_Train = []
9  for sent in grade_catogories_Train:
10     sent = sent.replace('-', '_')
11     sent = sent.replace(' ', '_')
12     # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
13     # https://gist.github.com/sebleier/554280
14     sent = ' '.join(e for e in sent.split())
15     grade_list_Train.append(sent.lower().strip())
16
17 # temp = temp.replace('-', '_')
18 X_Train['new_grade_category'] = grade_list_Train
19 X_Train.drop(['project_grade_category'], axis=1, inplace=True)
20
21 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
22 my_counter_grade_Train = Counter()
23 for word in X_Train['new_grade_category'].values:
24     my_counter_grade_Train.update(word.split())
25
26 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
27 grade_dict_Train = dict(my_counter_grade_Train)
28 sorted_grade_dict_Train = dict(sorted(grade_dict_Train.items(), key=lambda kv: kv[1]))
29
30 vectorizer_grade = CountVectorizer(vocabulary=list(sorted_grade_dict_Train.keys()), lowercase=False, binary=True)
31 vectorizer_grade.fit(X_Train['new_grade_category'].values)
32 #print(vectorizer_grade.get_feature_names())
33
34 grade_one_hot_Train = vectorizer_grade.transform(X_Train['new_grade_category'].values)
35
36 #-----Vectorizing categorical data of Project Grade for CV dataset-----
37
38 grade_catogories_CV = list(X_CV['project_grade_category'].values)
39 grade_list_CV = []
40 for sent in grade_catogories_CV:
41     sent = sent.replace('-', '_')
42     sent = sent.replace(' ', '_')
43     # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
44     # https://gist.github.com/sebleier/554280
45     sent = ' '.join(e for e in sent.split())
```

```

46     grade_list_CV.append(sent.lower().strip())
47
48     # temp = temp.replace('-', '_')
49     X_CV['new_grade_category'] = grade_list_CV
50     X_CV.drop(['project_grade_category'], axis=1, inplace=True)
51
52     # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
53     my_counter_grade_CV = Counter()
54     for word in X_CV['new_grade_category'].values:
55         my_counter_grade_CV.update(word.split())
56
57     # dict sort by value python: https://stackoverflow.com/a/613218/4084039
58     grade_dict_CV = dict(my_counter_grade_CV)
59     sorted_grade_dict_CV = dict(sorted(grade_dict_CV.items(), key=lambda kv: kv[1]))
60
61     grade_one_hot_CV = vectorizer_grade.transform(X_CV['new_grade_category'].values)
62
63     #-----Vectorizing categorical data of Project Grade for Train dataset-----
64
65     grade_categories_Test = list(X_Test['project_grade_category'].values)
66     grade_list_Test = []
67     for sent in grade_categories_Test:
68         sent = sent.replace('-', '_')
69         sent = sent.replace(' ', '_')
70         # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
71         # https://gist.github.com/sebleier/554280
72         sent = ' '.join(e for e in sent.split())
73         grade_list_Test.append(sent.lower().strip())
74
75     # temp = temp.replace('-', '_')
76     X_Test['new_grade_category'] = grade_list_Test
77     X_Test.drop(['project_grade_category'], axis=1, inplace=True)
78
79     # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
80     my_counter_grade_Test = Counter()
81     for word in X_Test['new_grade_category'].values:
82         my_counter_grade_Test.update(word.split())
83
84     # dict sort by value python: https://stackoverflow.com/a/613218/4084039
85     grade_dict_Test = dict(my_counter_grade_Test)
86     sorted_grade_dict_Test = dict(sorted(grade_dict_Test.items(), key=lambda kv: kv[1]))
87
88     grade_one_hot_Test = vectorizer_grade.transform(X_Test['new_grade_category'].values)
89     print("-"*120)
90     print('Shape of Train dataset matrix after one hot encoding is: {}'.format(grade_one_hot_Train.shape))
91     print('Shape of CV dataset matrix after one hot encoding is: {}'.format(grade_one_hot_CV.shape))

```

```
92 print('Shape of Testdataset matrix after one hot encoding is: {0}'.format(grade_one_hot_Test.shape))
```

```
-----  
Shape of Train dataset matrix after one hot encoding is: (11222, 4)  
Shape of CV dataset matrix after one hot encoding is: (5528, 4)  
Shape of Test dataset matrix after one hot encoding is: (8250, 4)
```

1.5.2 Vectorizing Numerical features

```
In [0]: 1 price_data = Resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()  
2 X_Train = pd.merge(X_Train, price_data, on='id', how='left')  
3 X_CV = pd.merge(X_CV, price_data, on='id', how='left')  
4 X_Test = pd.merge(X_Test, price_data, on='id', how='left')
```

```
In [0]: 1 # check this one: https://www.youtube.com/watch?v=0H0qOcln3Z4&t=530s  
2 # standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html  
3 # price_standardized = StandardScaler.fit(project_data['price'].values)  
4 # this will rise the error  
5 # ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.73 5.5 ].  
6 # Reshape your data either using array.reshape(-1, 1)  
7  
8 price_scalar = StandardScaler()  
9 median_price = Resource_data['price'].median()  
10 X_Train['price'] = X_Train['price'].fillna(median_price)  
11 X_CV['price'] = X_CV['price'].fillna(median_price)  
12 X_Test['price'] = X_Test['price'].fillna(median_price)  
13 price_scalar.fit(X_Train['price'].values.reshape(-1,1)) # finding the mean and standard deviation of this data  
14  
15  
16 # Now standardize the data with above maen and variance.  
17 price_standardized_Train = price_scalar.transform(X_Train['price'].values.reshape(-1, 1))  
18 price_standardized_CV = price_scalar.transform(X_CV['price'].values.reshape(-1, 1))  
19 price_standardized_Test = price_scalar.transform(X_Test['price'].values.reshape(-1, 1))
```

1.5.3 Vectorizing Text data

1.5.3.1 Bag of words

```
In [0]: 1 # We are considering only the words which appeared in at least 10 documents(rows or projects).
2 vectorizer_essays_bow = CountVectorizer(min_df=10)
3 text_bow_Train = vectorizer_essays_bow.fit_transform(preprocessed_essays_Train)
4 text_bow_CV = vectorizer_essays_bow.transform(preprocessed_essays_CV)
5 text_bow_Test = vectorizer_essays_bow.transform(preprocessed_essays_Test)
6 print("-"*120)
7 print("Applying Bag Of Words for Text Data")
8 print("-"*120)
9 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(text_bow_Train.shape))
10 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(text_bow_CV.shape))
11 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(text_bow_Test.shape))
```

Applying Bag Of Words for Text Data

Shape of Train dataset matrix after one hot encoding is: (11222, 6478)
Shape of CV dataset matrix after one hot encoding is: (5528, 6478)
Shape of Test dataset matrix after one hot encoding is: (8250, 6478)

Bag of Words for Project Title

```
In [0]: 1 # you can vectorize the title also
2 # before you vectorize the title make sure you preprocess it
3 vectorizer_titles_bow = CountVectorizer(min_df=10)
4 title_bow_Train = vectorizer_titles_bow.fit_transform(preprocessed_titles_Train)
5 title_bow_CV = vectorizer_titles_bow.transform(preprocessed_titles_CV)
6 title_bow_Test = vectorizer_titles_bow.transform(preprocessed_titles_Test)
7 print("-"*120)
8 print("Applying Bag Of Words for Project Title Data")
9 print("-"*120)
10 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(title_bow_Train.shape))
11 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(title_bow_CV.shape))
12 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(title_bow_Test.shape))
13
```

Applying Bag Of Words for Project Title Data

Shape of Train dataset matrix after one hot encoding is: (11222, 731)
Shape of CV dataset matrix after one hot encoding is: (5528, 731)
Shape of Test dataset matrix after one hot encoding is: (8250, 731)

1.5.2.2 TFIDF vectorizer

```
In [0]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 vectorizer_essays_tfidf = TfidfVectorizer(min_df=10)
3 text_tfidf_Train = vectorizer_essays_tfidf.fit_transform(preprocessed_essays_Train)
4 text_tfidf_CV = vectorizer_essays_tfidf.transform(preprocessed_essays_CV)
5 text_tfidf_Test = vectorizer_essays_tfidf.transform(preprocessed_essays_Test)
6 print("-"*120)
7 print("Applying TFIDF for Text Data")
8 print("-"*120)
9 print('Shape of Train dataset matrix after one hot encoding is: {0}'.format(text_tfidf_Train.shape))
10 print('Shape of CV dataset matrix after one hot encoding is: {0}'.format(text_tfidf_CV.shape))
11 print('Shape of Test dataset matrix after one hot encoding is: {0}'.format(text_tfidf_Test.shape))
```

Applying TFIDF for Text Data

Shape of Train dataset matrix after one hot encoding is: (11222, 6478)
Shape of CV dataset matrix after one hot encoding is: (5528, 6478)
Shape of Test dataset matrix after one hot encoding is: (8250, 6478)

TFIDF vectorizer for Project Title

```
In [0]: 1 vectorizer_titles_tfidf = TfidfVectorizer(min_df=10)
2 title_tfidf_Train = vectorizer_titles_tfidf.fit_transform(preprocessed_titles_Train)
3 title_tfidf_CV = vectorizer_titles_tfidf.transform(preprocessed_titles_CV)
4 title_tfidf_Test = vectorizer_titles_tfidf.transform(preprocessed_titles_Test)
5 print("-"*120)
6 print("Applying TFIDF for Project Title")
7 print("-"*120)
8 print('Shape of Train dataset matrix after one hot encoding is: {0}'.format(title_tfidf_Train.shape))
9 print('Shape of CV dataset matrix after one hot encoding is: {0}'.format(title_tfidf_CV.shape))
10 print('Shape of Test dataset matrix after one hot encoding is: {0}'.format(title_tfidf_Test.shape))
```

Applying TFIDF for Project Title

Shape of Train dataset matrix after one hot encoding is: (11222, 731)
Shape of CV dataset matrix after one hot encoding is: (5528, 731)
Shape of Test dataset matrix after one hot encoding is: (8250, 731)

1.5.2.3 Using Pretrained Models: Avg W2V

```
In [0]: 1 # stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-pyth
2 # make sure you have the glove_vectors file
3 with open('glove_vectors', 'rb') as f:
4     model = pickle.load(f)
5     glove_words = set(model.keys())
```


In [0]:

```
1 # average Word2Vec
2 # compute average word2vec for each review.
3 avg_w2v_vectors_Train = []; # the avg-w2v for each sentence/review is stored in this list
4 for sentence in tqdm(preprocessed_essays_Train): # for each review/sentence
5     vector = np.zeros(300) # as word vectors are of zero length
6     cnt_words = 0; # num of words with a valid vector in the sentence/review
7     for word in sentence.split(): # for each word in a review/sentence
8         if word in glove_words:
9             vector += model[word]
10            cnt_words += 1
11    if cnt_words != 0:
12        vector /= cnt_words
13    avg_w2v_vectors_Train.append(vector)
14 #-----
15
16 avg_w2v_vectors_CV = []; # the avg-w2v for each sentence/review is stored in this list
17 for sentence in tqdm(preprocessed_essays_CV): # for each review/sentence
18     vector = np.zeros(300) # as word vectors are of zero length
19     cnt_words = 0; # num of words with a valid vector in the sentence/review
20     for word in sentence.split(): # for each word in a review/sentence
21         if word in glove_words:
22             vector += model[word]
23             cnt_words += 1
24     if cnt_words != 0:
25         vector /= cnt_words
26     avg_w2v_vectors_CV.append(vector)
27 #-----
28
29 avg_w2v_vectors_Test = []; # the avg-w2v for each sentence/review is stored in this list
30 for sentence in tqdm(preprocessed_essays_Test): # for each review/sentence
31     vector = np.zeros(300) # as word vectors are of zero length
32     cnt_words = 0; # num of words with a valid vector in the sentence/review
33     for word in sentence.split(): # for each word in a review/sentence
34         if word in glove_words:
35             vector += model[word]
36             cnt_words += 1
37     if cnt_words != 0:
38         vector /= cnt_words
39     avg_w2v_vectors_Test.append(vector)
40
41 print(len(avg_w2v_vectors_Test))
42 print(len(avg_w2v_vectors_Test[1]))
```

100%|██████████| 11222/11222 [00:03<00:00, 3135.83it/s]

100%|██████████| 5528/5528 [00:01<00:00, 3163.21it/s]

100%|██████████| 8250/8250 [00:02<00:00, 3260.32it/s]

8250
300

AVG W2V on project_title

In [0]:

```
1  # Similarly you can vectorize for title also
2  # compute average word2vec for each title.
3  avg_w2v_vectors_title_Train = []; # the avg-w2v for each sentence/review is stored in this list
4  for sentence in tqdm(preprocessed_titles_Train): # for each review/sentence
5      vector_title = np.zeros(300) # as word vectors are of zero length
6      cnt_title_words = 0; # num of words with a valid vector in the sentence/review
7      for word in sentence.split(): # for each word in a review/sentence
8          if word in glove_words:
9              vector_title += model[word]
10             cnt_title_words += 1
11     if cnt_title_words != 0:
12         vector_title /= cnt_title_words
13     avg_w2v_vectors_title_Train.append(vector_title)
14
15
16  #-----
17  avg_w2v_vectors_title_CV = []; # the avg-w2v for each sentence/review is stored in this list
18  for sentence in tqdm(preprocessed_titles_CV): # for each review/sentence
19      vector_title = np.zeros(300) # as word vectors are of zero length
20      cnt_title_words = 0; # num of words with a valid vector in the sentence/review
21      for word in sentence.split(): # for each word in a review/sentence
22          if word in glove_words:
23              vector_title += model[word]
24              cnt_title_words += 1
25      if cnt_title_words != 0:
26          vector_title /= cnt_title_words
27      avg_w2v_vectors_title_CV.append(vector_title)
28
29  #-----
30  avg_w2v_vectors_title_Test = []; # the avg-w2v for each sentence/review is stored in this list
31  for sentence in tqdm(preprocessed_titles_Test): # for each review/sentence
32      vector_title = np.zeros(300) # as word vectors are of zero length
33      cnt_title_words = 0; # num of words with a valid vector in the sentence/review
34      for word in sentence.split(): # for each word in a review/sentence
35          if word in glove_words:
36              vector_title += model[word]
37              cnt_title_words += 1
38      if cnt_title_words != 0:
39          vector_title /= cnt_title_words
40      avg_w2v_vectors_title_Test.append(vector_title)
41
42  print(len(avg_w2v_vectors_title_Test))
43  print(len(avg_w2v_vectors_title_Test[0]))
```

100%|██████████| 11222/11222 [00:00<00:00, 57909.34it/s]

100%|██████████| 5528/5528 [00:00<00:00, 59084.94it/s]

100%|██████████| 8250/8250 [00:00<00:00, 58992.42it/s]

8250

300

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [0]:

```
1 tfidf_model_essays = TfidfVectorizer()
2 tfidf_model_essays.fit(preprocessed_essays_Train)
3 # we are converting a dictionary with word as a key, and the idf as a value
4 dictionary = dict(zip(tfidf_model_essays.get_feature_names(), list(tfidf_model_essays.idf_)))
5 tfidf_words_essays = set(tfidf_model_essays.get_feature_names())
```

In [0]:

```
1 # average Word2Vec
2 # compute average word2vec for each review.
3 tfidf_w2v_vectors_Train = []; # the avg-w2v for each sentence/review is stored in this list
4 for sentence in tqdm(preprocessed_essays_Train): # for each review/sentence
5     vector = np.zeros(300) # as word vectors are of zero length
6     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
7     for word in sentence.split(): # for each word in a review/sentence
8         if (word in glove_words) and (word in tfidf_words_essays):
9             vec = model[word] # getting the vector for each word
10             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))
11             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
12             vector += (vec * tf_idf) # calculating tfidf weighted w2v
13             tf_idf_weight += tf_idf
14     if tf_idf_weight != 0:
15         vector /= tf_idf_weight
16     tfidf_w2v_vectors_Train.append(vector)
17
18 #-----
19 tfidf_w2v_vectors_CV = []; # the avg-w2v for each sentence/review is stored in this list
20 for sentence in tqdm(preprocessed_essays_CV): # for each review/sentence
21     vector = np.zeros(300) # as word vectors are of zero length
22     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
23     for word in sentence.split(): # for each word in a review/sentence
24         if (word in glove_words) and (word in tfidf_words_essays):
25             vec = model[word] # getting the vector for each word
26             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))
27             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
28             vector += (vec * tf_idf) # calculating tfidf weighted w2v
29             tf_idf_weight += tf_idf
30     if tf_idf_weight != 0:
31         vector /= tf_idf_weight
32     tfidf_w2v_vectors_CV.append(vector)
33 #-----
34 tfidf_w2v_vectors_Test = []; # the avg-w2v for each sentence/review is stored in this list
35 for sentence in tqdm(preprocessed_essays_Test): # for each review/sentence
36     vector = np.zeros(300) # as word vectors are of zero length
37     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
38     for word in sentence.split(): # for each word in a review/sentence
39         if (word in glove_words) and (word in tfidf_words_essays):
40             vec = model[word] # getting the vector for each word
41             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))
42             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
43             vector += (vec * tf_idf) # calculating tfidf weighted w2v
44             tf_idf_weight += tf_idf
45     if tf_idf_weight != 0:
```

```
46         vector /= tf_idf_weight
47         tfidf_w2v_vectors_Test.append(vector)
48
49     print(len(tfidf_w2v_vectors_Test))
50     print(len(tfidf_w2v_vectors_Test[0]))
51
```

```
100%|██████████| 11222/11222 [00:20<00:00, 551.58it/s]
100%|██████████| 5528/5528 [00:10<00:00, 548.87it/s]
100%|██████████| 8250/8250 [00:15<00:00, 536.89it/s]
```

8250

300

Using Pretrained Models: TFIDF weighted W2V on project_title

In [0]:

```
1 # Similarly you can vectorize for title also
2 tfidf_model_title = TfidfVectorizer()
3 tfidf_model_title.fit(preprocessed_titles_Train)
4 # we are converting a dictionary with word as a key, and the idf as a value
5 dictionary = dict(zip(tfidf_model_title.get_feature_names(), list(tfidf_model_title.idf_)))
6 tfidf_words_title = set(tfidf_model_title.get_feature_names())
7
8 # compute tfidf word2vec for each title.
9 tfidf_w2v_vectors_title_Train = []; # the avg-w2v for each sentence/review is stored in this list
10 for sentence in tqdm(preprocessed_titles_Train): # for each review/sentence
11     vector_title = np.zeros(300) # as word vectors are of zero length
12     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
13     for word in sentence.split(): # for each word in a review/sentence
14         if (word in glove_words) and (word in tfidf_words_title):
15             vec = model[word] # getting the vector for each word
16             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))
17             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
18             vector_title += (vector_title * tf_idf) # calculating tfidf weighted w2v
19             tf_idf_weight += tf_idf
20     if tf_idf_weight != 0:
21         vector_title /= tf_idf_weight
22     tfidf_w2v_vectors_title_Train.append(vector_title)
23 #-----
24
25 tfidf_w2v_vectors_title_CV = []; # the avg-w2v for each sentence/review is stored in this list
26 for sentence in tqdm(preprocessed_titles_CV): # for each review/sentence
27     vector_title = np.zeros(300) # as word vectors are of zero length
28     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
29     for word in sentence.split(): # for each word in a review/sentence
30         if (word in glove_words) and (word in tfidf_words_title):
31             vec = model[word] # getting the vector for each word
32             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))
33             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
34             vector_title += (vector_title * tf_idf) # calculating tfidf weighted w2v
35             tf_idf_weight += tf_idf
36     if tf_idf_weight != 0:
37         vector_title /= tf_idf_weight
38     tfidf_w2v_vectors_title_CV.append(vector_title)
39 #-----
40
41
42 tfidf_w2v_vectors_title_Test = []; # the avg-w2v for each sentence/review is stored in this list
43 for sentence in tqdm(preprocessed_titles_Test): # for each review/sentence
44     vector_title = np.zeros(300) # as word vectors are of zero length
45     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
```

```

46     for word in sentence.split(): # for each word in a review/sentence
47         if (word in glove_words) and (word in tfidf_words_title):
48             vec = model[word] # getting the vector for each word
49             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))
50             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
51             vector_title += (vector_title * tf_idf) # calculating tfidf weighted w2v
52             tf_idf_weight += tf_idf
53         if tf_idf_weight != 0:
54             vector_title /= tf_idf_weight
55         tfidf_w2v_vectors_title_Test.append(vector_title)
56
57 print(len(tfidf_w2v_vectors_title_Test))
58 print(len(tfidf_w2v_vectors_title_Test[0]))
59
60

```

```

100%|██████████| 11222/11222 [00:00<00:00, 28350.87it/s]
100%|██████████| 5528/5528 [00:00<00:00, 35369.22it/s]
100%|██████████| 8250/8250 [00:00<00:00, 35216.34it/s]

8250
300

```

Assignment 3: Apply KNN

1. [Task-1] Apply KNN(brute force version) on these feature sets

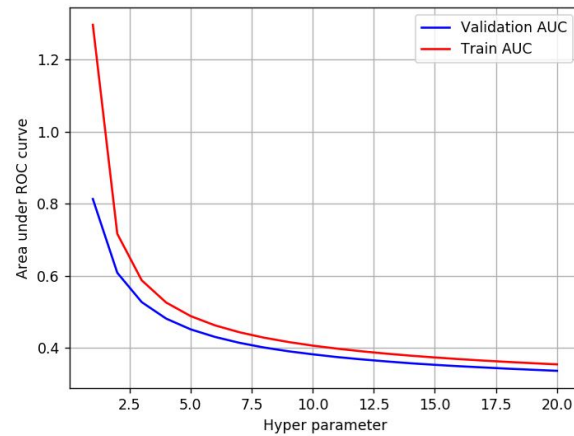
- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_essay (BOW)
- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_essay (TFIDF)
- **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_essay (AVG W2V)
- **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

2. Hyper paramter tuning to find best K

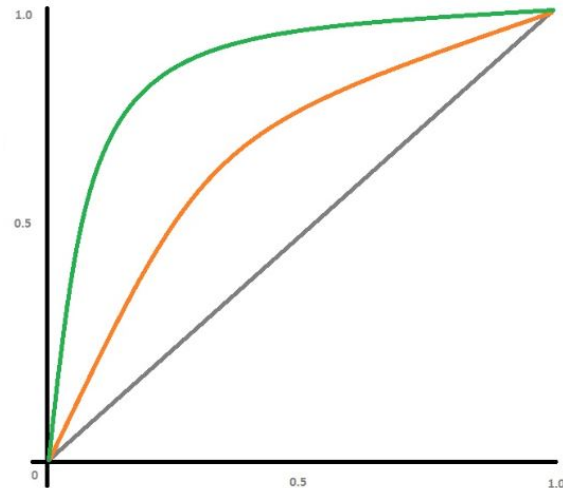
- Find the best hyper parameter which results in the maximum [AUC \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) value
- Find the best hyper paramter using k-fold cross validation (or) simple cross validation data
- Use gridsearch-cv or randomsearch-cv or write your own for loops to do this task

3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, as shown in the figure



- Once you find the best hyper parameter, you need to train your model-M using the best hyper-param. Now, find the AUC on test data and plot the ROC curve on both train and test using model-M.



- Along with plotting ROC curve, you need to print the [confusion matrix](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/>) with predicted and original labels of test data points

	Predicted: NO	Predicted: YES
Actual: NO	TN = ??	FP = ??
Actual: YES	FN = ??	TP = ??

- Select top 2000 features from feature **Set 2** using `SelectKBest` (https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.SelectKBest.html) and then apply KNN on top of these features

- ```
from sklearn.datasets import load_digits
from sklearn.feature_selection import SelectKBest, chi2
X, y = load_digits(return_X_y=True)
X.shape
X_new = SelectKBest(chi2, k=20).fit_transform(X, y)
X_new.shape
=====
output:
(1797, 64)
(1797, 20)
```

- Repeat the steps 2 and 3 on the data matrix after feature selection

## 5. Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library [link](http://zetcode.com/python/prettytable/) (<http://zetcode.com/python/prettytable/>).

| Vectorizer | Model | Hyper parameter | AUC  |
|------------|-------|-----------------|------|
| BOW        | Brute | 7               | 0.78 |
| TFIDF      | Brute | 12              | 0.79 |
| W2V        | Brute | 10              | 0.78 |
| TFIDFW2V   | Brute | 6               | 0.78 |

### Note: Data Leakage

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
2. To avoid the issue of data-leakag, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on you train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf). (<https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf>)

## **2. K Nearest Neighbor**

### **2.1 Splitting data into Train and cross validation(or test): Stratified Sampling**

### **2.2 Make Data Model Ready: encoding numerical, categorical features**

### **2.3 Make Data Model Ready: encoding eassay, and project\_title**

### **2.4 Appling KNN on different kind of featurization as mentioned in the instructions**

Apply KNN on different kind of featurization as mentioned in the instructions

For Every model that you work on make sure you do the step 2 and step 3 of instructions

#### **1.5.4 Merging all the above features**

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

```
In [0]: 1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2 BOW_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_one_
3 print(BOW_Train.shape)
4 TFIDF_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_one
5 print(TFIDF_Train.shape)
6 AVG_W2V_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_
7 print(AVG_W2V_Train.shape)
8 TFIDF_W2V_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefi
9 print(TFIDF_W2V_Train.shape)
10 TFIDF_TrainKB = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_
11 print(TFIDF_TrainKB.shape)
```

```
(11222, 7308)
(11222, 7308)
(11222, 699)
(11222, 699)
(11222, 7308)
```

```
In [0]: 1 BOW_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,price_st
2 print(BOW_CV.shape)
3 TFIDF_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,price_
4 print(TFIDF_CV.shape)
5 AVG_W2V_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,pric
6 print(AVG_W2V_CV.shape)
7 TFIDF_W2V_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,pr
8 print(TFIDF_W2V_CV.shape)
9 TFIDF_CVKB = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,pric
10 print(TFIDF_CVKB.shape)
```

```
(5528, 7308)
(5528, 7308)
(5528, 699)
(5528, 699)
(5528, 7308)
```

```
In [0]: 1 BOW_Test = hstack((categories_one_hot_Test,sub_categories_one_hot_Test,school_one_hot_Test,grade_one_hot_Test,prefix_one_hot_Test))
2 print(BOW_Test.shape)
3 TFIDF_Test = hstack((categories_one_hot_Test,sub_categories_one_hot_Test,school_one_hot_Test,grade_one_hot_Test,prefix_one_hot_Test))
4 print(TFIDF_Test.shape)
5 AVG_W2V_Test = hstack((categories_one_hot_Test,sub_categories_one_hot_Test,school_one_hot_Test,grade_one_hot_Test,prefix_one_hot_Test))
6 print(AVG_W2V_Test.shape)
7 TFIDF_W2V_Test = hstack((categories_one_hot_Test,sub_categories_one_hot_Test,school_one_hot_Test,grade_one_hot_Test,prefix_one_hot_Test))
8 print(TFIDF_W2V_Test.shape)
9 TFIDF_Test_KB = hstack((categories_one_hot_Test,sub_categories_one_hot_Test,school_one_hot_Test,grade_one_hot_Test,prefix_one_hot_Test))
10 print(TFIDF_Test_KB.shape)
```

(8250, 7308)

(8250, 7308)

(8250, 699)

(8250, 699)

(8250, 7308)

## Loading Test Pickle files

In [0]:

```
1 pck = open('BOW_Test', 'wb')
2 pickle.dump(BOW_Test, pck)
3 pck = open('BOW_Test', 'rb')
4 BOW_Test = pickle.load(pck)
5 pck.close()
6
7 pck = open('TFIDF_Test', 'wb')
8 pickle.dump(TFIDF_Test, pck)
9 pck = open('TFIDF_Test', 'rb')
10 TFIDF_Test = pickle.load(pck)
11 pck.close()
12
13 pck = open('AVG_W2V_Test', 'wb')
14 pickle.dump(AVG_W2V_Test, pck)
15 pck = open('AVG_W2V_Test', 'rb')
16 AVG_W2V_Test = pickle.load(pck)
17 pck.close()
18
19 pck = open('TFIDF_W2V_Test', 'wb')
20 pickle.dump(TFIDF_W2V_Test, pck)
21 pck = open('TFIDF_W2V_Test', 'rb')
22 TFIDF_W2V_Test = pickle.load(pck)
23 pck.close()
24
25 pck = open('TFIDF_Test_KB', 'wb')
26 pickle.dump(TFIDF_Test_KB, pck)
27 pck = open('TFIDF_Test_KB', 'rb')
28 TFIDF_Test_KB = pickle.load(pck)
29 pck.close()
```

```
In [0]: 1 BOW_TCV = vstack((BOW_Train,BOW_CV))
2 print(BOW_TCV.shape)
3 TFIDF_TCV = vstack((TFIDF_Train,TFIDF_CV))
4 print(TFIDF_TCV.shape)
5 AVG_W2V_TCV = vstack((AVG_W2V_Train,AVG_W2V_CV))
6 print(AVG_W2V_TCV.shape)
7 TFIDF_W2V_TCV = vstack((TFIDF_W2V_Train,TFIDF_W2V_CV))
8 print(TFIDF_W2V_TCV.shape)
9 TFIDF_KB_TCV = vstack((TFIDF_TrainKB,TFIDF_CVKB))
10 print(TFIDF_KB_TCV.shape)
```

(16750, 7308)

(16750, 7308)

(16750, 699)

(16750, 699)

(16750, 7308)

## Loading Train Pickle files

```
In [0]: 1 pck = open('BOW_TCV', 'wb')
2 pickle.dump(BOW_TCV, pck)
3 pck = open('BOW_TCV', 'rb')
4 BOW_TCV = pickle.load(pck)
5 pck.close()
6
7 pck = open('TFIDF_TCV', 'wb')
8 pickle.dump(TFIDF_TCV, pck)
9 pck = open('TFIDF_TCV', 'rb')
10 TFIDF_TCV = pickle.load(pck)
11 pck.close()
12
13 pck = open('AVG_W2V_TCV', 'wb')
14 pickle.dump(AVG_W2V_TCV, pck)
15 pck = open('AVG_W2V_TCV', 'rb')
16 AVG_W2V_TCV = pickle.load(pck)
17 pck.close()
18
19 pck = open('TFIDF_W2V_TCV', 'wb')
20 pickle.dump(TFIDF_W2V_TCV, pck)
21 pck = open('TFIDF_W2V_TCV', 'rb')
22 TFIDF_W2V_TCV = pickle.load(pck)
23 pck.close()
24
25 pck = open('TFIDF_KB_TCV', 'wb')
26 pickle.dump(TFIDF_KB_TCV, pck)
27 pck = open('TFIDF_KB_TCV', 'rb')
28 TFIDF_KB_TCV = pickle.load(pck)
29 pck.close()
```

```
In [0]: 1 def batch_predict(clf, data):
2 y_data_pred = []
3 tr_loop = data.shape[0] - data.shape[0]%1000
4 for i in range(0, tr_loop, 1000):
5 y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
6 y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
7 return y_data_pred
```

### 2.4.1 Applying KNN brute force on BOW, SET 1

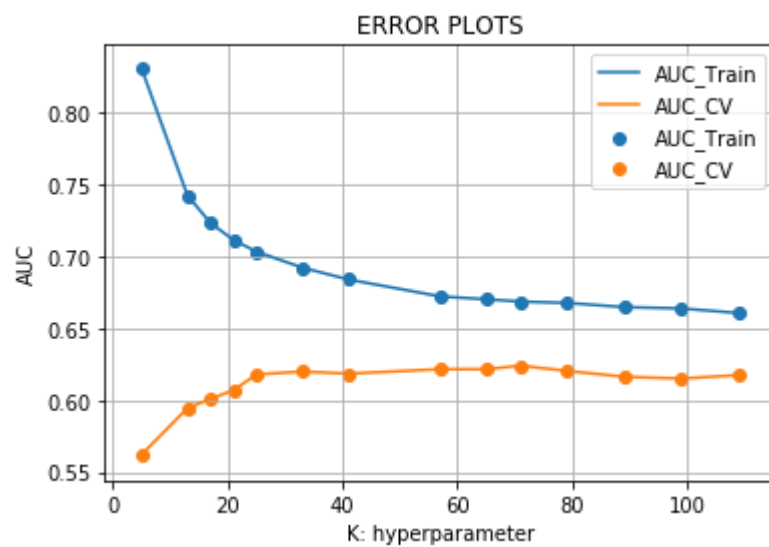
Find the best hyper parameter which results in the maximum AUC value



In [0]:

```
1 %%time
2 BOW_TCV_CSR = BOW_TCV.tocsr()
3 BOW_TR_CSR = BOW_Train.tocsr()
4 BOW_CV_CSR = BOW_CV.tocsr()
5 BOW_Test_CSR = BOW_Test.tocsr()
6
7 k_range = [5,13,17,21,25,33,41,57,65,71,79,89,99,109]
8
9 ACCV = []
10 AUC_TR = []
11 AUC_TS = []
12
13 for i in tqdm(k_range):
14 knn = KNeighborsClassifier(n_neighbors=i, n_jobs=-1,algorithm='brute')
15 knn.fit(BOW_Train, Y_Train)
16 pred = knn.predict(BOW_CV)
17 acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
18 ACCV.append(acc)
19 Train_pred = batch_predict(knn, BOW_TR_CSR)
20 a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
21 AUC_TR.append(auc(a_fpr_train, a_tpr_train))
22
23 Test_pred = batch_predict(knn, BOW_CV_CSR)
24 a_fpr_Test,a_tpr_Test,c = roc_curve(Y_CV, Test_pred)
25 AUC_TS.append(auc(a_fpr_Test, a_tpr_Test))
26
27 # Performance of model on Train data and Test data for each hyper parameter.
28 plt.plot(k_range, AUC_TR, label='AUC_Train')
29 plt.scatter(k_range, AUC_TR, label='AUC_Train')
30 plt.gca()
31 plt.plot(k_range, AUC_TS, label='AUC_CV')
32 plt.scatter(k_range, AUC_TS, label='AUC_CV')
33 plt.gca()
34 plt.legend()
35 plt.xlabel("K: hyperparameter")
36 plt.ylabel("AUC")
37 plt.title("ERROR PLOTS")
38 plt.grid()
39 plt.show()
```

100%|██████████| 14/14 [06:05<00:00, 25.98s/it]



CPU times: user 10min 3s, sys: 1.8 s, total: 10min 5s  
 Wall time: 6min 6s

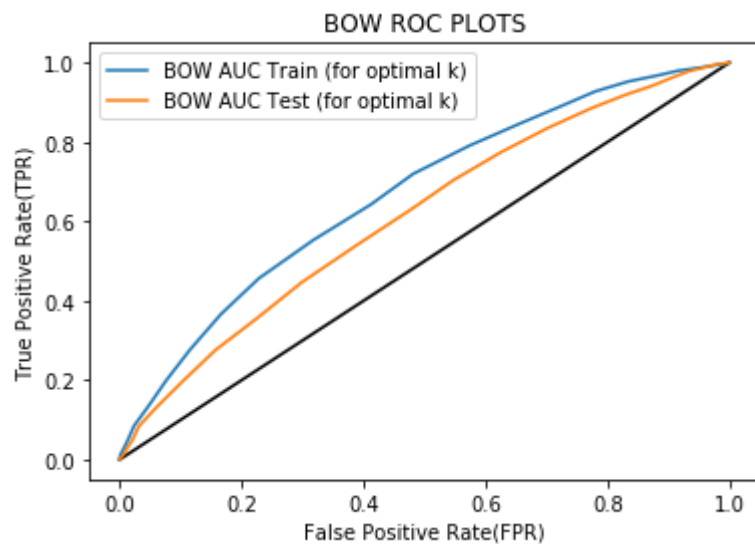
**OBSERVATION:** From the above plot the difference between AUC train and AUC cv started to reduce from the K value 71. From the plot we can see that K value 71 and 79 have same difference. Hence i have choose the optimal K value as 71.

```
In [0]: 1 k_opt=71
2 knn_opt = KNeighborsClassifier(n_neighbors = k_opt, n_jobs=-1,algorithm='brute')
3 knn_opt.fit(BOW_Train, Y_Train)
4 pred = knn.predict(BOW_Test)
5 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
6 print('\nTest accuracy for k = {0} is {1}%'.format(k_opt,acc))
7
8 Y_Train_pred = batch_predict(knn_opt, BOW_TR_CSR)
9 Y_Test_pred = batch_predict(knn_opt, BOW_Test_CSR)
10
11 fpr_Train, tpr_Train, thresholds = roc_curve(Y_Train, Y_Train_pred)
12 fpr_Test, tpr_Test, thresholds = roc_curve(Y_Test, Y_Test_pred)
```

Test accuracy for k = 71 is 84.87272727272727%

## BOW ROC PLOT

```
In [0]: 1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-')
6 plt.plot(fpr_Train, tpr_Train, label="BOW AUC Train (for optimal k)")
7 plt.plot(fpr_Test, tpr_Test, label="BOW AUC Test (for optimal k)")
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("BOW ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal k) =", auc(fpr_Train, tpr_Train))
15 print("AUC Test (for optimal k) =", auc(fpr_Test, tpr_Test))
16 BOW_AUC=round(auc(fpr_Test, tpr_Test)*100)
17 BOW_K=k_opt
18
19
20
```



AUC Train (for optimal k) = 0.6686594927594207

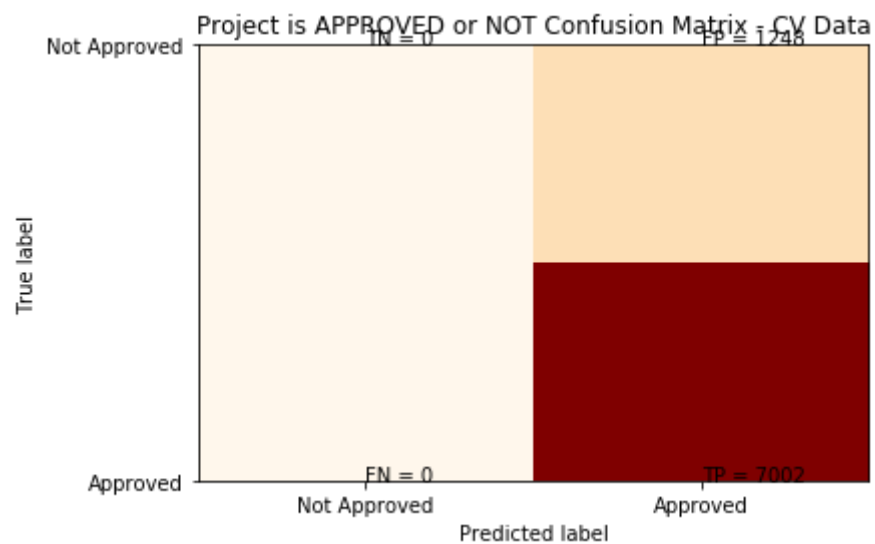
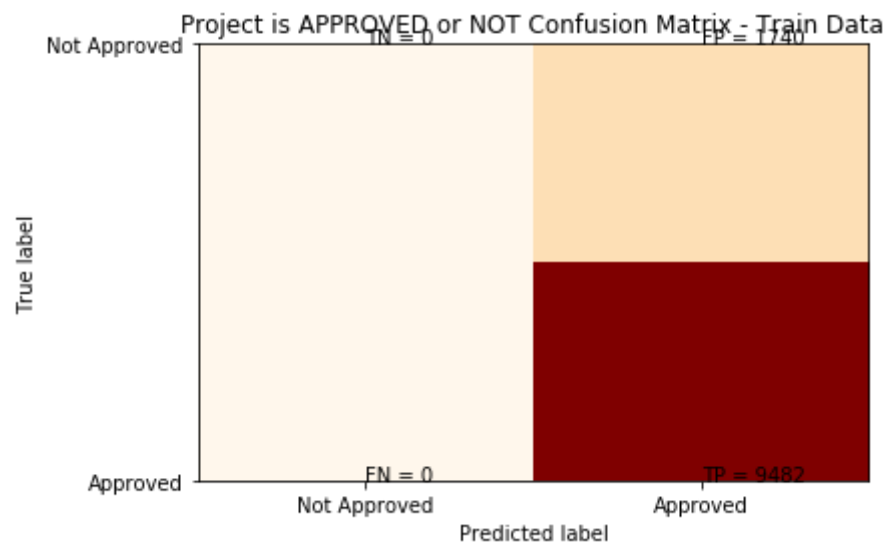
AUC Test (for optimal k) = 0.6138982039929983

CPU times: user 188 ms, sys: 3 ms, total: 191 ms  
Wall time: 201 ms

## BOW CONFUSION MATRIX

In [0]:

```
1 %%time
2 #https://tatwan.github.io/How-To-Plot-A-Confusion-Matrix-In-Python/
3 #https://matplotlib.org/3.1.1/gallery/color/colormap_reference.html
4
5 PR1= knn_opt.predict(BOW_Train)
6 PR2= knn_opt.predict(BOW_Test)
7 #-----Confusion matrix for BOW Train Data-----
8 plt.clf()
9 CM1 = confusion_matrix(Y_Train,PR1)
10 plt.imshow(CM1, interpolation='nearest', cmap='OrRd', aspect='auto')
11 classNames = ['Not Approved', 'Approved']
12 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
13 plt.ylabel('True label')
14 plt.xlabel('Predicted label')
15 tick_marks = np.arange(len(classNames))
16 plt.xticks(tick_marks, classNames, rotation=0)
17 plt.yticks(tick_marks, classNames)
18 s = [['TN', 'FP'], ['FN', 'TP']]
19 for i in range(2):
20 for j in range(2):
21 plt.text(j,i, str(s[i][j])+ " = "+str(CM1[i][j]))
22 plt.show()
23
24 #-----Confusion matrix for BOW Test Data-----
25 plt.clf()
26 CM2 = confusion_matrix(Y_Test,PR2)
27 plt.imshow(CM2, interpolation='nearest', cmap='OrRd', aspect='auto')
28 classNames = ['Not Approved', 'Approved']
29 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
30 plt.ylabel('True label')
31 plt.xlabel('Predicted label')
32 tick_marks = np.arange(len(classNames))
33 plt.xticks(tick_marks, classNames, rotation=0)
34 plt.yticks(tick_marks, classNames)
35 s = [['TN', 'FP'], ['FN', 'TP']]
36 for i in range(2):
37 for j in range(2):
38 plt.text(j,i, str(s[i][j])+ " = "+str(CM2[i][j]))
39 plt.show()
```

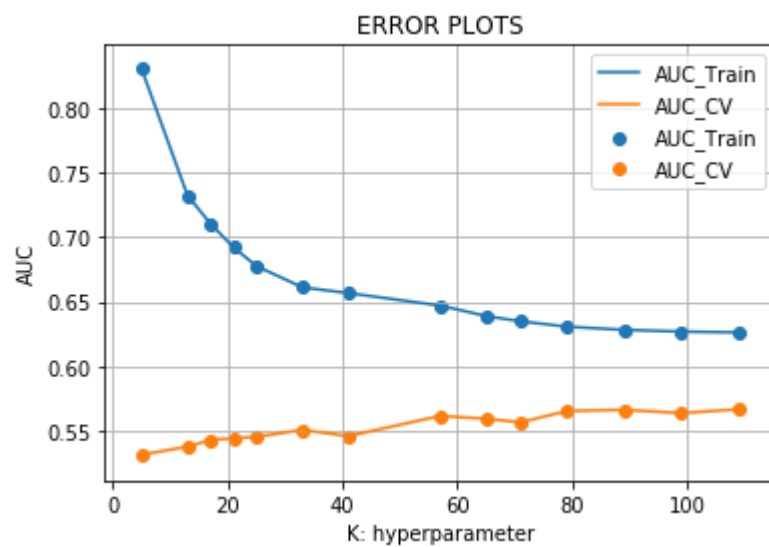


CPU times: user 44.5 s, sys: 86.9 ms, total: 44.6 s  
Wall time: 27.5 s

## 2.4.2 Applying KNN brute force on TFIDF, SET 2

In [0]:

```
1 %%time
2
3 TFIDF_TCV_CSR=TFIDF_TCV.tocsr()
4 TFIDF_TR_CSR = TFIDF_Train.tocsr()
5 TFIDF_CV_CSR = TFIDF_CV.tocsr()
6 TFIDF_Test_CSR = TFIDF_Test.tocsr()
7
8
9 k_range = [5,13,17,21,25,33,41,57,65,71,79,89,99,109]
10
11 ACCV = []
12 AUC_TR = []
13 AUC_TS = []
14
15 for i in tqdm(k_range):
16 knn = KNeighborsClassifier(n_neighbors=i, n_jobs=-1,algorithm='brute')
17 knn.fit(TFIDF_TR_CSR, Y_Train)
18 pred = knn.predict(TFIDF_CV)
19 acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
20 ACCV.append(acc)
21 Train_pred = batch_predict(knn, TFIDF_TR_CSR)
22 a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
23 AUC_TR.append(auc(a_fpr_train, a_tpr_train))
24
25 Test_pred = batch_predict(knn, TFIDF_CV_CSR)
26 a_fpr_Test,a_tpr_Test,c = roc_curve(Y_CV, Test_pred)
27 AUC_TS.append(auc(a_fpr_Test, a_tpr_Test))
28
29 # Performance of model on Train data and Test data for each hyper parameter.
30 plt.plot(k_range, AUC_TR, label='AUC_Train')
31 plt.scatter(k_range, AUC_TR, label='AUC_Train')
32 plt.gca()
33 plt.plot(k_range, AUC_TS, label='AUC_CV')
34 plt.scatter(k_range, AUC_TS, label='AUC_CV')
35 plt.gca()
36 plt.legend()
37 plt.xlabel("K: hyperparameter")
38 plt.ylabel("AUC")
39 plt.title("ERROR PLOTS")
40 plt.grid()
41 plt.show()
```



**OBSERVATION:** From the above plot the difference between AUC train and AUC cv started to reduce from the K value 89. After that difference seems to be same. Hence i have choose Optimal K value as 89.

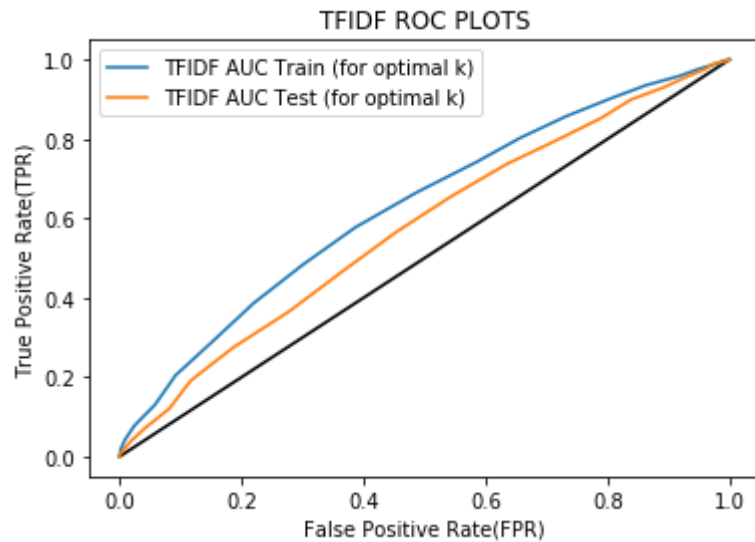
```
In [0]: 1 k_opt=89
2
3 knn_opt = KNeighborsClassifier(n_neighbors = k_opt, n_jobs=-1,algorithm='brute')
4 knn_opt.fit(TFIDF_Train, Y_Train)
5 pred = knn.predict(TFIDF_Test)
6 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
7 print('\nTest accuracy for k = {0} is {1}%'.format(k_opt,acc))
8
9 Y_Train_pred = batch_predict(knn_opt, TFIDF_TR_CSR)
10 Y_Test_pred = batch_predict(knn_opt, TFIDF_Test_CSR)
11
12 fpr_Train, tpr_Train, thresholds = roc_curve(Y_Train, Y_Train_pred)
13 fpr_Test, tpr_Test, thresholds = roc_curve(Y_Test, Y_Test_pred)
```

Test accuracy for k = 89 is 84.87272727272727%

## TFIDF ROC PLOT



```
In [0]: 1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-')
6 plt.plot(fpr_Train, tpr_Train, label="TFIDF AUC Train (for optimal k)")
7 plt.plot(fpr_Test, tpr_Test, label="TFIDF AUC Test (for optimal k)")
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("TFIDF ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal k) =", auc(fpr_Train, tpr_Train))
15 print("AUC Test (for optimal k) =", auc(fpr_Test, tpr_Test))
16 TFIDF_AUC=round(auc(fpr_Test, tpr_Test)*100)
17 TFIDF_K=k_opt
```



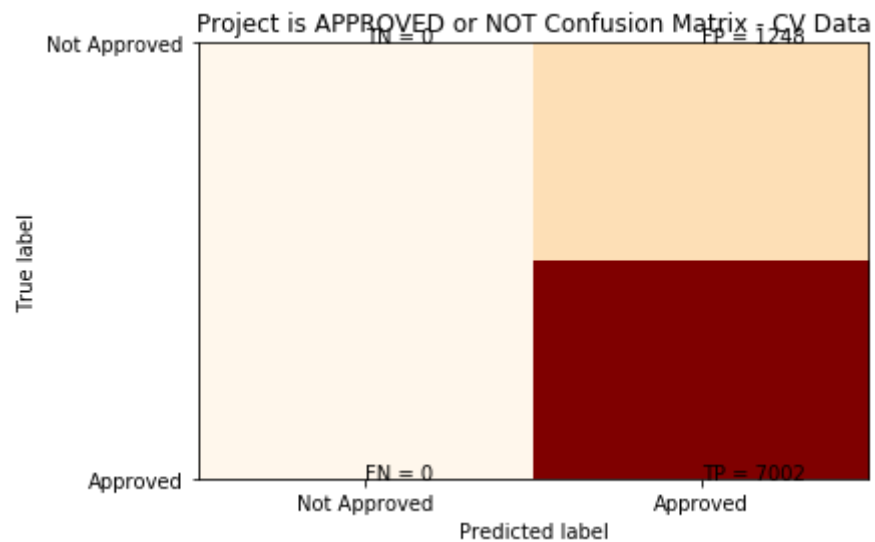
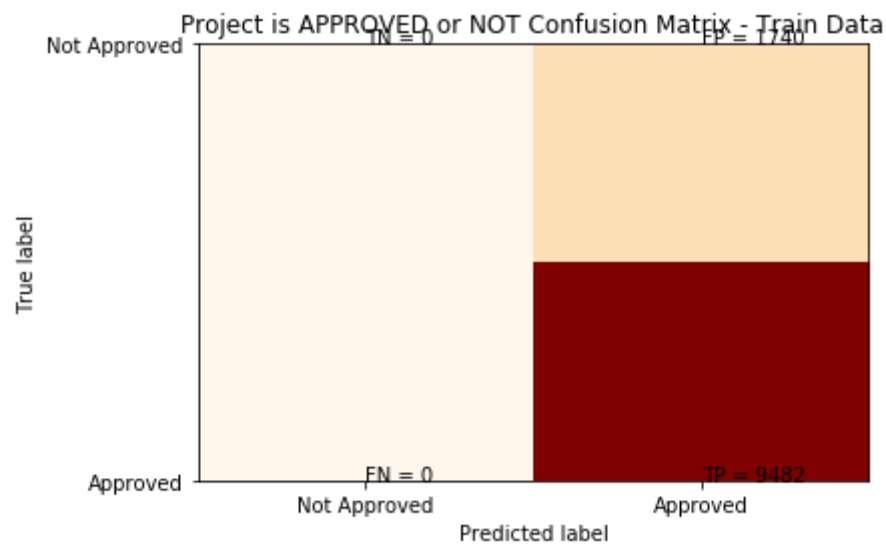

---

AUC Train (for optimal k) = 0.6281867397876679  
AUC Test (for optimal k) = 0.5761106945634581  
CPU times: user 202 ms, sys: 4.97 ms, total: 207 ms  
Wall time: 210 ms

## TFIDF CONFUSION MATRIX



```
In [0]: 1 %%time
2 #https://tatwan.github.io/How-To-Plot-A-Confusion-Matrix-In-Python/
3 #https://matplotlib.org/3.1.1/gallery/color/colormap_reference.html
4
5 PR3= knn_opt.predict(TFIDF_Train)
6 PR4= knn_opt.predict(TFIDF_Test)
7 #-----Confusion matrix for TFIDF Train Data-----
8 plt.clf()
9 CM3 = confusion_matrix(Y_Train,PR3)
10 plt.imshow(CM3, interpolation='nearest', cmap='OrRd', aspect='auto')
11 classNames = ['Not Approved', 'Approved']
12 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
13 plt.ylabel('True label')
14 plt.xlabel('Predicted label')
15 tick_marks = np.arange(len(classNames))
16 plt.xticks(tick_marks, classNames, rotation=0)
17 plt.yticks(tick_marks, classNames)
18 s = [['TN', 'FP'], ['FN', 'TP']]
19 for i in range(2):
20 for j in range(2):
21 plt.text(j,i, str(s[i][j])+ " = "+str(CM3[i][j]))
22 plt.show()
23
24 #-----Confusion matrix for TFIDF Test Data-----
25 plt.clf()
26 CM4 = confusion_matrix(Y_Test,PR4)
27 plt.imshow(CM4, interpolation='nearest', cmap='OrRd', aspect='auto')
28 classNames = ['Not Approved', 'Approved']
29 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
30 plt.ylabel('True label')
31 plt.xlabel('Predicted label')
32 tick_marks = np.arange(len(classNames))
33 plt.xticks(tick_marks, classNames, rotation=0)
34 plt.yticks(tick_marks, classNames)
35 s = [['TN', 'FP'], ['FN', 'TP']]
36 for i in range(2):
37 for j in range(2):
38 plt.text(j,i, str(s[i][j])+ " = "+str(CM4[i][j]))
39 plt.show()
```



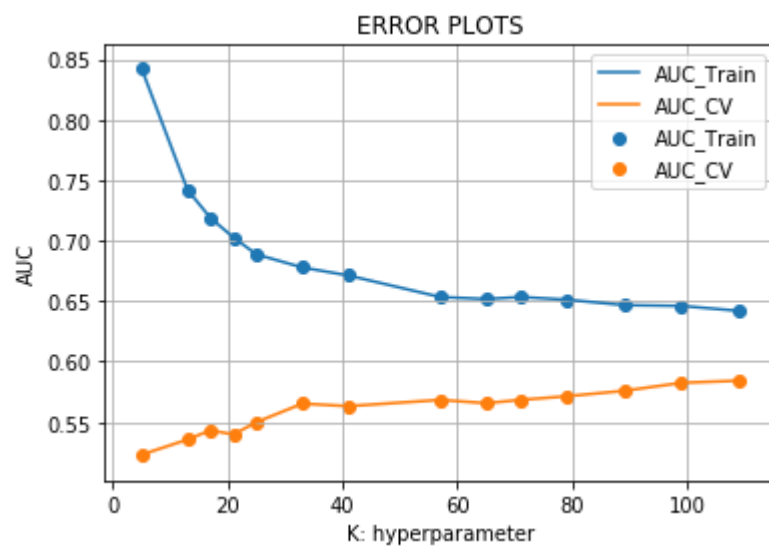
CPU times: user 45.1 s, sys: 89.8 ms, total: 45.2 s  
Wall time: 27.9 s

### 2.4.3 Applying KNN brute force on AVG W2V, SET 3

In [0]:

```
1 %%time
2 AVG_W2V_TCV_CSR = AVG_W2V_TCV.tocsr()
3 AVG_W2V_TR_CSR = AVG_W2V_Train.tocsr()
4 AVG_W2V_CV_CSR = AVG_W2V_CV.tocsr()
5 AVG_W2V_Test_CSR = AVG_W2V_Test.tocsr()
6
7 k_range = [5,13,17,21,25,33,41,57,65,71,79,89,99,109]
8
9 ACCV = []
10 AUC_TR = []
11 AUC_TS = []
12
13 for i in tqdm(k_range):
14 knn = KNeighborsClassifier(n_neighbors=i, n_jobs=-1,algorithm='brute')
15 knn.fit(AVG_W2V_TR_CSR, Y_Train)
16 pred = knn.predict(AVG_W2V_CV)
17 acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
18 ACCV.append(acc)
19 Train_pred = batch_predict(knn, AVG_W2V_TR_CSR)
20 a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
21 AUC_TR.append(auc(a_fpr_train, a_tpr_train))
22
23 Test_pred = batch_predict(knn, AVG_W2V_CV_CSR)
24 a_fpr_Test,a_tpr_Test,c = roc_curve(Y_CV, Test_pred)
25 AUC_TS.append(auc(a_fpr_Test, a_tpr_Test))
26
27 # Performance of model on Train data and Test data for each hyper parameter.
28 plt.plot(k_range, AUC_TR, label='AUC_Train')
29 plt.scatter(k_range, AUC_TR, label='AUC_Train')
30 plt.gca()
31 plt.plot(k_range, AUC_TS, label='AUC_CV')
32 plt.scatter(k_range, AUC_TS, label='AUC_CV')
33 plt.gca()
34 plt.legend()
35 plt.xlabel("K: hyperparameter")
36 plt.ylabel("AUC")
37 plt.title("ERROR PLOTS")
38 plt.grid()
39 plt.show()
```

100%|██████████| 14/14 [1:22:53<00:00, 361.47s/it]



CPU times: user 2h 41min 54s, sys: 8.2 s, total: 2h 42min 2s  
 Wall time: 1h 22min 53s

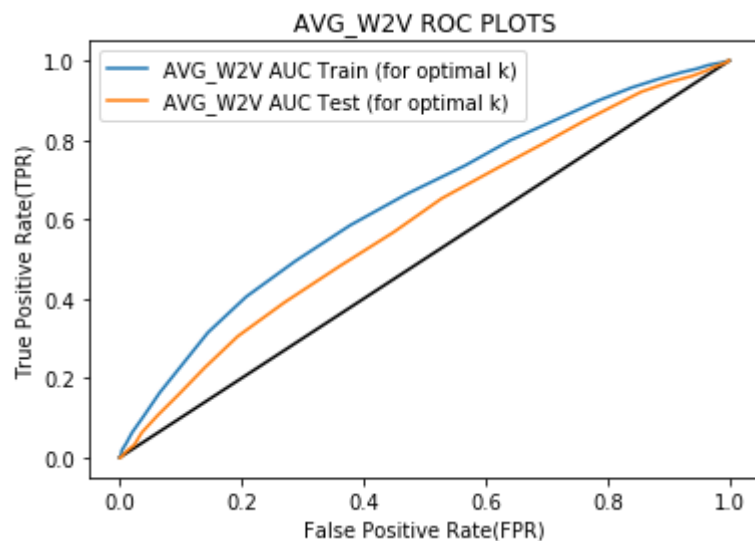
**OBSERVATION:** From the above plot the difference between AUC train and AUC cv have very less value for the K=109. Hence i have choose that as an Optimal K value.

```
In [0]: 1 k_opt=109
2
3 knn_opt = KNeighborsClassifier(n_neighbors = k_opt, n_jobs=-1,algorithm='brute')
4 knn_opt.fit(AVG_W2V_Train, Y_Train)
5 pred = knn.predict(AVG_W2V_Test)
6 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
7 print('\nTest accuracy for k = {0} is {1}%'.format(k_opt,acc))
8
9 Y_Train_pred = batch_predict(knn_opt, AVG_W2V_TR_CSR)
10 Y_Test_pred = batch_predict(knn_opt, AVG_W2V_Test_CSR)
11
12 fpr_Train, tpr_Train, thresholds = roc_curve(Y_Train, Y_Train_pred)
13 fpr_Test, tpr_Test, thresholds = roc_curve(Y_Test, Y_Test_pred)
```

Test accuracy for k = 109 is 84.87272727272727%

## AVG W2V ROC PLOT

```
In [0]: 1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1],'k-')
6 plt.plot(fpr_Train, tpr_Train, label="AVG_W2V AUC Train (for optimal k)")
7 plt.plot(fpr_Test, tpr_Test, label="AVG_W2V AUC Test (for optimal k)")
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("AVG_W2V ROC PLOTS")
12 plt.show()
13 print("-"*120)
14
15 print("AUC Train (for optimal k) =", auc(fpr_Train, tpr_Train))
16 print("AUC Test (for optimal k) =", auc(fpr_Test, tpr_Test))
17 AVG_W2V_AUC=round(auc(fpr_Test, tpr_Test)*100)
18 AVG_W2V_K=k_opt
```




---

```
AUC Train (for optimal k) = 0.6419674483049553
AUC Test (for optimal k) = 0.5878782801983315
CPU times: user 200 ms, sys: 5 ms, total: 205 ms
Wall time: 208 ms
```

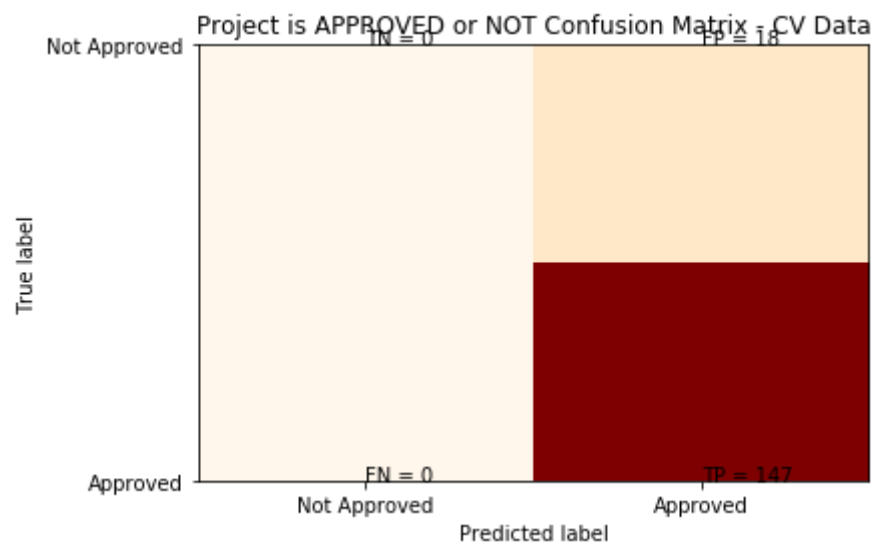
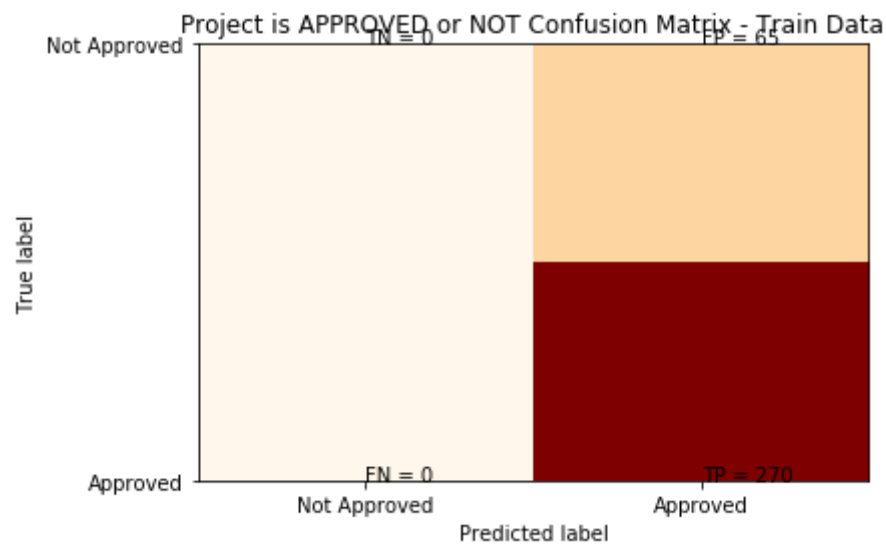
## AVG W2V CONFUSION MATRIX





In [0]:

```
1 %%time
2 #https://tatwan.github.io/How-To-Plot-A-Confusion-Matrix-In-Python/
3 #https://matplotlib.org/3.1.1/gallery/color/colormap_reference.html
4
5 PR5= knn_opt.predict(AVG_W2V_Train)
6 PR6= knn_opt.predict(AVG_W2V_Test)
7 #-----Confusion matrix for AVG_W2V Train Data-----
8 plt.clf()
9 CM5 = confusion_matrix(Y_Train,PR5)
10 plt.imshow(CM5, interpolation='nearest', cmap='OrRd', aspect='auto')
11 classNames = ['Not Approved', 'Approved']
12 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
13 plt.ylabel('True label')
14 plt.xlabel('Predicted label')
15 tick_marks = np.arange(len(classNames))
16 plt.xticks(tick_marks, classNames, rotation=0)
17 plt.yticks(tick_marks, classNames)
18 s = [['TN', 'FP'], ['FN', 'TP']]
19 for i in range(2):
20 for j in range(2):
21 plt.text(j,i, str(s[i][j])+ " = "+str(CM5[i][j]))
22 plt.show()
23
24 #-----Confusion matrix for AVG_W2V Test Data-----
25 plt.clf()
26 CM6 = confusion_matrix(Y_Test,PR6)
27 plt.imshow(CM6, interpolation='nearest', cmap='OrRd', aspect='auto')
28 classNames = ['Not Approved', 'Approved']
29 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
30 plt.ylabel('True label')
31 plt.xlabel('Predicted label')
32 tick_marks = np.arange(len(classNames))
33 plt.xticks(tick_marks, classNames, rotation=0)
34 plt.yticks(tick_marks, classNames)
35 s = [['TN', 'FP'], ['FN', 'TP']]
36 for i in range(2):
37 for j in range(2):
38 plt.text(j,i, str(s[i][j])+ " = "+str(CM6[i][j]))
39 plt.show()
```



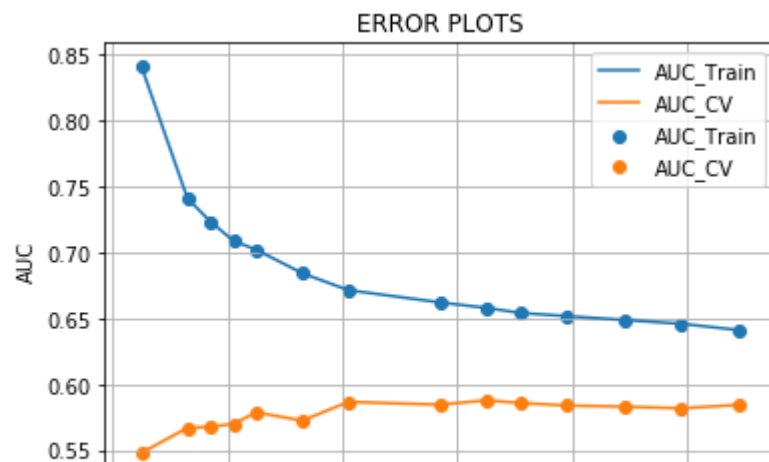
Wall time: 612 ms

## 2.4.4 Applying KNN brute force on TFIDF W2V, SET 4

In [0]:

```
1 %%time
2 TFIDF_W2V_TCV_CSR = TFIDF_W2V_TCV.tocsr()
3 TFIDF_W2V_TR_CSR = TFIDF_W2V_Train.tocsr()
4 TFIDF_W2V_CV_CSR = TFIDF_W2V_CV.tocsr()
5 TFIDF_W2V_Test_CSR = TFIDF_W2V_Test.tocsr()
6
7 k_range = [5,13,17,21,25,33,41,57,65,71,79,89,99,109]
8
9 ACCV = []
10 AUC_TR = []
11 AUC_TS = []
12
13 for i in tqdm(k_range):
14 knn = KNeighborsClassifier(n_neighbors=i, n_jobs=-1,algorithm='brute')
15 knn.fit(TFIDF_W2V_TR_CSR, Y_Train)
16 pred = knn.predict(TFIDF_W2V_CV)
17 acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
18 ACCV.append(acc)
19 Train_pred = batch_predict(knn, TFIDF_W2V_TR_CSR)
20 a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
21 AUC_TR.append(auc(a_fpr_train, a_tpr_train))
22
23 Test_pred = batch_predict(knn, TFIDF_W2V_CV_CSR)
24 a_fpr_Test,a_tpr_Test,c = roc_curve(Y_CV, Test_pred)
25 AUC_TS.append(auc(a_fpr_Test, a_tpr_Test))
26
27 # Performance of model on Train data and Test data for each hyper parameter.
28 plt.plot(k_range, AUC_TR, label='AUC_Train')
29 plt.scatter(k_range, AUC_TR, label='AUC_Train')
30 plt.gca()
31 plt.plot(k_range, AUC_TS, label='AUC_CV')
32 plt.scatter(k_range, AUC_TS, label='AUC_CV')
33 plt.gca()
34 plt.legend()
35 plt.xlabel("K: hyperparameter")
36 plt.ylabel("AUC")
37 plt.title("ERROR PLOTS")
38 plt.grid()
39 plt.show()
```

100%|██████████| 14/14 [42:32<00:00, 183.43s/it]



**OBSERVATION:** From the above plot the difference between AUC train and AUC cv have very less value for the K=109. Hence i have choose that as an Optimal K value.

```
In [0]: 1 k_opt=109
2
3 knn_opt = KNeighborsClassifier(n_neighbors = k_opt, n_jobs=-1,algorithm='brute')
4 knn_opt.fit(TFIDF_W2V_Train, Y_Train)
5 pred = knn.predict(TFIDF_W2V_Test)
6 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
7 print('\nTest accuracy for k = {0} is {1}%'.format(k_opt,acc))
8
9 Y_Train_pred = batch_predict(knn_opt, TFIDF_W2V_TR_CSR)
10 Y_Test_pred = batch_predict(knn_opt, TFIDF_W2V_Test_CSR)
11
12 fpr_Train, tpr_Train, thresholds = roc_curve(Y_Train, Y_Train_pred)
13 fpr_Test, tpr_Test, thresholds = roc_curve(Y_Test, Y_Test_pred)
14
```

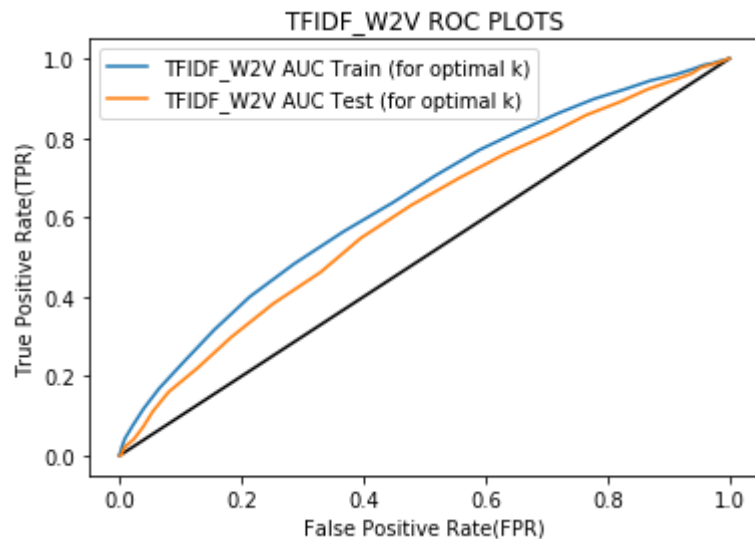
Test accuracy for k = 109 is 84.87272727272727%

## TFIDF W2V ROC PLOT

```

In [0]: 1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-')
6 plt.plot(fpr_Train, tpr_Train, label="TFIDF_W2V AUC Train (for optimal k)")
7 plt.plot(fpr_Test, tpr_Test, label="TFIDF_W2V AUC Test (for optimal k)")
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("TFIDF_W2V ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal k) =", auc(fpr_Train, tpr_Train))
15 print("AUC Test (for optimal k) =", auc(fpr_Test, tpr_Test))
16 TFIDF_W2V_AUC=round(auc(fpr_Test, tpr_Test)*100)
17 TFIDF_W2V_K=k_opt
18

```



-----

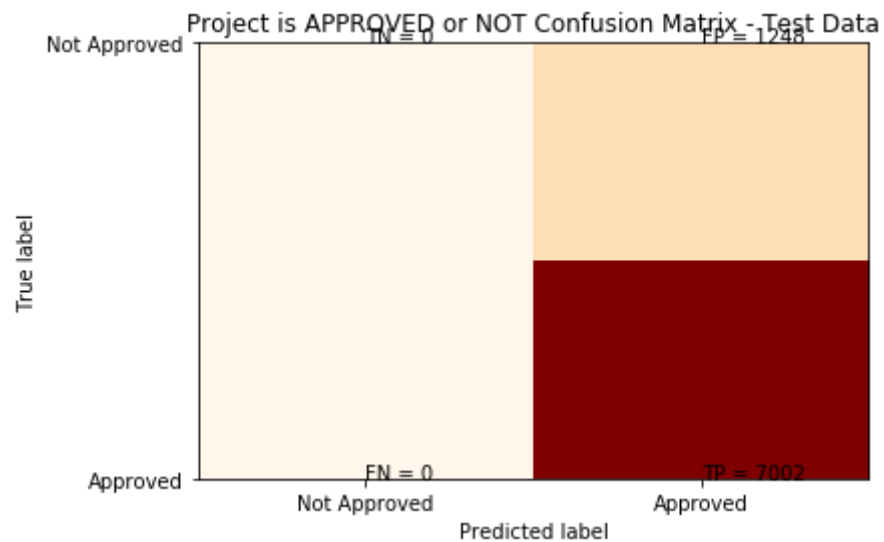
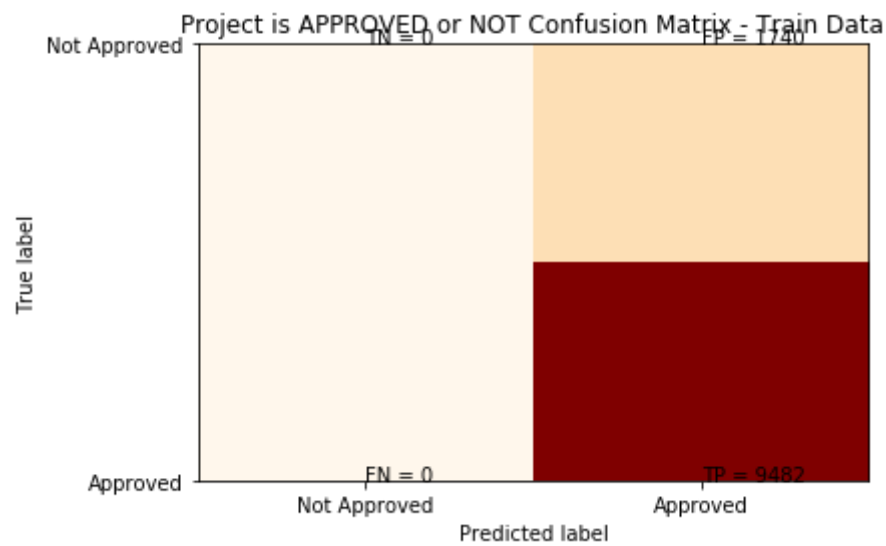
AUC Train (for optimal k) = 0.6408993931635742

AUC Test (for optimal k) = 0.6011983641120852  
CPU times: user 199 ms, sys: 4.99 ms, total: 204 ms  
Wall time: 207 ms

## TFIDF W2V CONFUSION MATRIX

In [0]:

```
1 %%time
2 #https://tatwan.github.io/How-To-Plot-A-Confusion-Matrix-In-Python/
3 #https://matplotlib.org/3.1.1/gallery/color/colormap_reference.html
4
5 PR7= knn_opt.predict(TFIDF_W2V_Train)
6 PR8= knn_opt.predict(TFIDF_W2V_Test)
7 #-----Confusion matrix for TFIDF_W2V Train Data-----
8 plt.clf()
9 CM7 = confusion_matrix(Y_Train,PR7)
10 plt.imshow(CM7, interpolation='nearest', cmap='OrRd', aspect='auto')
11 classNames = ['Not Approved', 'Approved']
12 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
13 plt.ylabel('True label')
14 plt.xlabel('Predicted label')
15 tick_marks = np.arange(len(classNames))
16 plt.xticks(tick_marks, classNames, rotation=0)
17 plt.yticks(tick_marks, classNames)
18 s = [['TN', 'FP'], ['FN', 'TP']]
19 for i in range(2):
20 for j in range(2):
21 plt.text(j,i, str(s[i][j])+ " = "+str(CM7[i][j]))
22 plt.show()
23
24 #-----Confusion matrix for TFIDF_W2V Test Data-----
25 plt.clf()
26 CM8 = confusion_matrix(Y_Test,PR8)
27 plt.imshow(CM8, interpolation='nearest', cmap='OrRd', aspect='auto')
28 classNames = ['Not Approved', 'Approved']
29 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
30 plt.ylabel('True label')
31 plt.xlabel('Predicted label')
32 tick_marks = np.arange(len(classNames))
33 plt.xticks(tick_marks, classNames, rotation=0)
34 plt.yticks(tick_marks, classNames)
35 s = [['TN', 'FP'], ['FN', 'TP']]
36 for i in range(2):
37 for j in range(2):
38 plt.text(j,i, str(s[i][j])+ " = "+str(CM8[i][j]))
39 plt.show()
```



CPU times: user 5min 10s, sys: 572 ms, total: 5min 10s  
Wall time: 2min 41s

## 2.5 Feature selection with SelectKBest



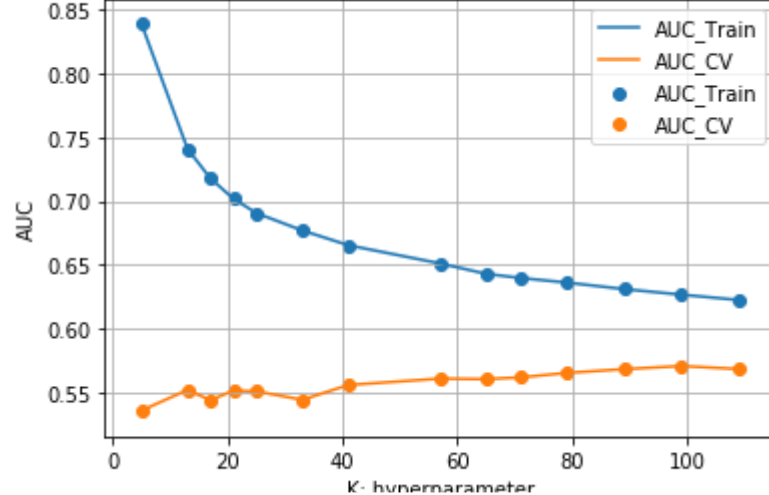
```
In [0]: 1 %%time
2 #https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.f_classif.html#sklearn.feature_selection.f_classif
3 warnings.filterwarnings("ignore")
4
5 TFIDF_Train_KB1 = SelectKBest(f_classif, k=2000).fit(TFIDF_TrainKB, Y_Train)
6 TFIDF_Train_KBST = TFIDF_Train_KB1.transform(TFIDF_TrainKB)
7 TFIDF_CV_KBST = TFIDF_Train_KB1.transform(TFIDF_CVKB)
8 TFIDF_Test_KBST = TFIDF_Train_KB1.transform(TFIDF_Test_KB)
9
10 print("-"*120)
11 print("Selecting the K Best")
12 print("-"*120)
13 print('Shape of Train dataset matrix after one hot encoding is: {0}'.format(TFIDF_Train_KBST.shape))
14 print('Shape of CV dataset matrix after one hot encoding is: {0}'.format(TFIDF_CV_KBST.shape))
15 print('Shape of Test dataset matrix after one hot encoding is: {0}'.format(TFIDF_Test_KBST.shape))
```

-----  
Selecting the K Best  
-----

Shape of Train dataset matrix after one hot encoding is: (11222, 2000)  
Shape of CV dataset matrix after one hot encoding is: (5528, 2000)  
Shape of Test dataset matrix after one hot encoding is: (8250, 2000)  
CPU times: user 212 ms, sys: 1 ms, total: 213 ms  
Wall time: 214 ms

In [0]:

```
1 %%time
2 TFIDF_TCV_CSR_KBST = TFIDF_KB_TCV.tocsr()
3 TFIDF_TR_CSR_KBST = TFIDF_Train_KBST.tocsr()
4 TFIDF_CV_CSR_KBST = TFIDF_CV_KBST.tocsr()
5 TFIDF_Test_CSR_KBST = TFIDF_Test_KB.tocsr()
6
7 k_range = [5,13,17,21,25,33,41,57,65,71,79,89,99,109]
8
9
10 ACCV = []
11 AUC_TR = []
12 AUC_TS = []
13
14 for i in tqdm(k_range):
15 knn = KNeighborsClassifier(n_neighbors=i, n_jobs=-1,algorithm='brute')
16 knn.fit(TFIDF_Train_KBST, Y_Train)
17 pred = knn.predict(TFIDF_CV_KBST)
18 acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
19 ACCV.append(acc)
20 Train_pred = batch_predict(knn, TFIDF_TR_CSR_KBST)
21 a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
22 AUC_TR.append(auc(a_fpr_train, a_tpr_train))
23
24 Test_pred = batch_predict(knn, TFIDF_CV_CSR_KBST)
25 a_fpr_Test,a_tpr_Test,c = roc_curve(Y_CV, Test_pred)
26 AUC_TS.append(auc(a_fpr_Test, a_tpr_Test))
27
28 # Performance of model on Train data and Test data for each hyper parameter.
29 plt.plot(k_range, AUC_TR, label='AUC_Train')
30 plt.scatter(k_range, AUC_TR, label='AUC_Train')
31 plt.gca()
32 plt.plot(k_range, AUC_TS, label='AUC_CV')
33 plt.scatter(k_range, AUC_TS, label='AUC_CV')
34 plt.gca()
35 plt.legend()
36 plt.xlabel("K: hyperparameter")
37 plt.ylabel("AUC")
38 plt.title("ERROR PLOTS")
39 plt.grid()
40 plt.show()
```



**OBSERVATION:** From the above plot the difference between AUC train and AUC cv started to reduce from the K value 99 and difference remain same for the K values 99 and 109. Hence i have choose the optimal K value as 99.

In [0]:

```

1
2 k_opt=99
3
4 knn_opt = KNeighborsClassifier(n_neighbors = k_opt, n_jobs=-1,algorithm='brute')
5 knn_opt.fit(TFIDF_Train_KBST, Y_Train)
6 pred = knn.predict(TFIDF_Test_KBST)
7 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
8 print('\nTest accuracy for k = {0} is {1}%'.format(k_opt,acc))
9
10 Y_Train_pred = batch_predict(knn_opt, TFIDF_TR_CSR_KBST)
11 Y_Test_pred = batch_predict(knn_opt, TFIDF_Test_KBST.tocsr())
12
13 fpr_Train, tpr_Train, thresholds = roc_curve(Y_Train, Y_Train_pred)
14 fpr_Test, tpr_Test, thresholds = roc_curve(Y_Test, Y_Test_pred)
15

```

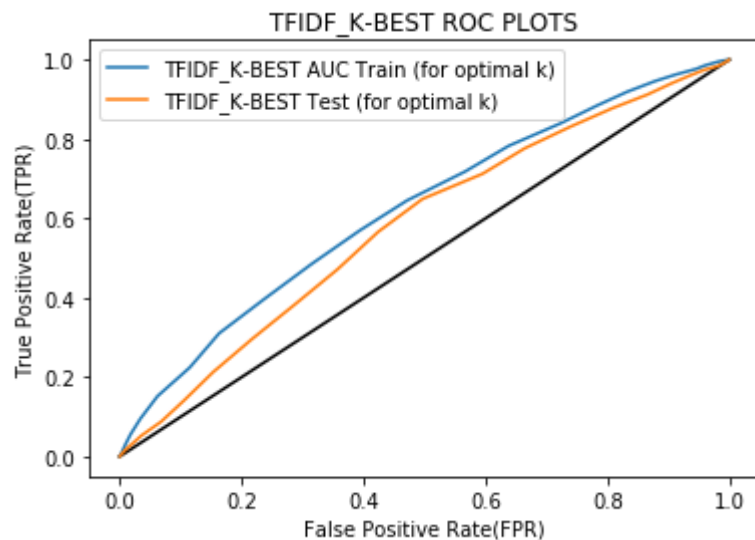
Test accuracy for k = 99 is 84.87272727272727%

## SELECT K BEST TFIDF ROC PLOT

```

In [0]: 1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-')
6 plt.plot(fpr_Train, tpr_Train, label="TFIDF_K-BEST AUC Train (for optimal k)")
7 plt.plot(fpr_Test, tpr_Test, label="TFIDF_K-BEST Test (for optimal k)")
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("TFIDF_K-BEST ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal k) =", auc(fpr_Train, tpr_Train))
15 print("AUC Test (for optimal k) =", auc(fpr_Test, tpr_Test))
16 KBST_AUC=round(auc(fpr_Test, tpr_Test)*100)
17 KBST_K=k_opt
18

```




---

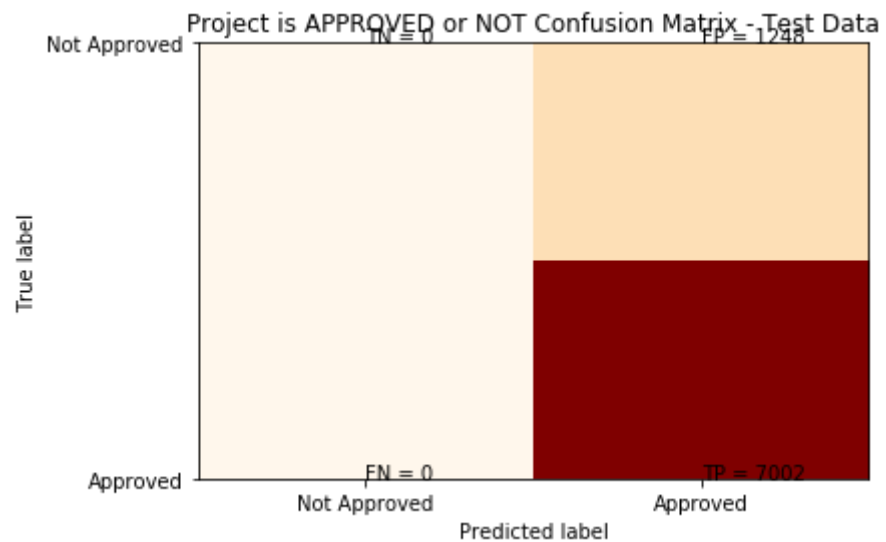
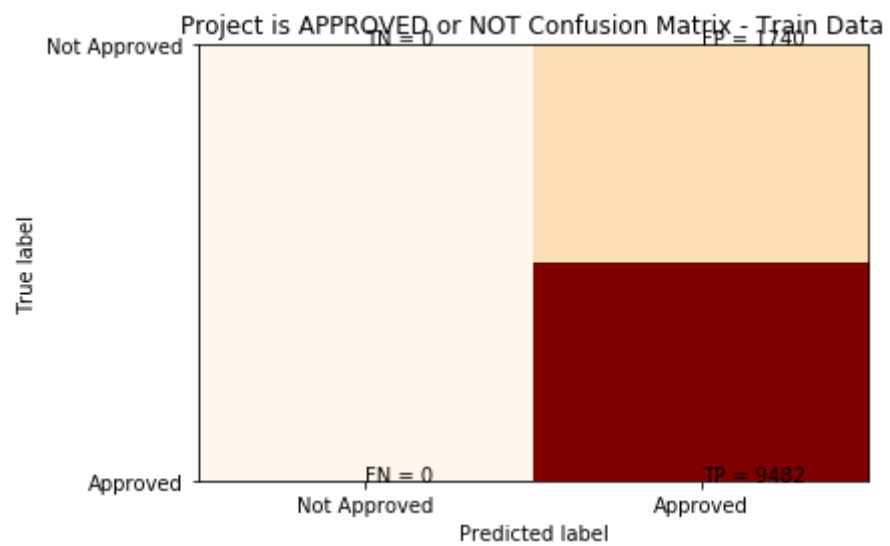
AUC Train (for optimal k) = 0.6223923671469476  
 AUC Test (for optimal k) = 0.5824019945766411

CPU times: user 183 ms, sys: 3.99 ms, total: 187 ms  
Wall time: 189 ms

## **SELECT K BEST TFIDF CONFUSION MATRIX**

In [0]:

```
1 %%time
2 #https://tatwan.github.io/How-To-Plot-A-Confusion-Matrix-In-Python/
3 #https://matplotlib.org/3.1.1/gallery/color/colormap_reference.html
4
5 PR9= knn_opt.predict(TFIDF_Train_KBST)
6 PR10= knn_opt.predict(TFIDF_Test_KBST)
7 #-----Confusion matrix for K-BEST TFIDF_W2V Train Data-----
8 plt.clf()
9 CM9 = confusion_matrix(Y_Train,PR9)
10 plt.imshow(CM9, interpolation='nearest', cmap='OrRd', aspect='auto')
11 classNames = ['Not Approved', 'Approved']
12 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
13 plt.ylabel('True label')
14 plt.xlabel('Predicted label')
15 tick_marks = np.arange(len(classNames))
16 plt.xticks(tick_marks, classNames, rotation=0)
17 plt.yticks(tick_marks, classNames)
18 s = [['TN', 'FP'], ['FN', 'TP']]
19 for i in range(2):
20 for j in range(2):
21 plt.text(j,i, str(s[i][j])+ " = "+str(CM9[i][j]))
22 plt.show()
23
24 #-----Confusion matrix for K-BEST TFIDF_W2V Test Data-----
25 plt.clf()
26 CM10 = confusion_matrix(Y_Test,PR10)
27 plt.imshow(CM8, interpolation='nearest', cmap='OrRd', aspect='auto')
28 classNames = ['Not Approved', 'Approved']
29 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
30 plt.ylabel('True label')
31 plt.xlabel('Predicted label')
32 tick_marks = np.arange(len(classNames))
33 plt.xticks(tick_marks, classNames, rotation=0)
34 plt.yticks(tick_marks, classNames)
35 s = [['TN', 'FP'], ['FN', 'TP']]
36 for i in range(2):
37 for j in range(2):
38 plt.text(j,i, str(s[i][j])+ " = "+str(CM10[i][j]))
39 plt.show()
```



CPU times: user 34.2 s, sys: 1.38 s, total: 35.6 s  
Wall time: 22.5 s

### 3. Conclusions

```
In [0]: 1 # Please compare all your models using Prettytable Library
2
3 pt = PrettyTable()
4 pt.field_names= ("Vectorizer", "Model", "HyperParameter", "AUC")
5 pt.add_row(["BOW", "Brute", BOW_K, BOW_AUC])
6 pt.add_row(["Tf-Idf", "Brute", TFIDF_K, TFIDF_AUC])
7 pt.add_row(["AVG W2V", "Brute", AVG_W2V_K, AVG_W2V_AUC])
8 pt.add_row(["TFIDF W2V", "Brute", TFIDF_W2V_K, TFIDF_W2V_AUC])
9 pt.add_row(["Top 2000 features of Tf-Idf", "Brute", KBST_K, KBST_AUC])
10 print(pt)
```

| Vectorizer                  | Model | HyperParameter | AUC  |
|-----------------------------|-------|----------------|------|
| BOW                         | Brute | 71             | 61.0 |
| Tf-Idf                      | Brute | 89             | 58.0 |
| AVG W2V                     | Brute | 109            | 59.0 |
| TFIDF W2V                   | Brute | 109            | 60.0 |
| Top 2000 features of Tf-Idf | Brute | 99             | 58.0 |