

SVM Algorithm on Donors_Chose dataset

```
In [1]: %matplotlib inline
import warnings
warnings.filterwarnings("ignore")
import sqlite3
import pandas as pd
import numpy as np
import nltk
import gc
gc.enable()
gc.DEBUG_SAVEALL
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os
import math
from plotly import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
gc.set_threshold(2, 1, 1)
```

2.1 Loading Input Data

```
In [2]: # %Load_ext memory_profiler

s=0
# We are taking samples of 0's and 1's and appending them to overcome memory error

project_data = pd.read_csv('train_data.csv')
# project_data=project_data.dropna(how='any')
project_data_1 = project_data[project_data['project_is_approved'] == s+1]
project_data_0 = project_data[project_data['project_is_approved'] == s]
project_data=project_data.fillna("")
project_data_1=project_data_1.head(34000)
project_data_0=project_data_0.tail(16000)
project_data_1=project_data_1.append(project_data_0)
project_data=project_data_1
resource_data = pd.read_csv('resources.csv')

#Sorting them by columns to spread the zeros and one's unevenly in the 'project_'
project_data.sort_values(by=['project_essay_1'])
project_data.sort_values(by=['project_essay_2'], ascending=False)
project_data.sort_values(by=['project_essay_3'])
project_data.sort_values(by=['project_essay_4'], ascending=False)
project_data_1=None
project_data_0=None
```

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

Number of data points in train data (50000, 17)
-----
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
'project_submitted_datetime' 'project_grade_category'
'project_subject_categories' 'project_subject_subcategories'
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

```
In [4]: print("Number of data points in resource data", resource_data.shape)
print(resource_data.columns.values)
resource_data.head(1)
# project_data.head(2)
```

```
Number of data points in resource data (1541272, 4)
['id' 'description' 'quantity' 'price']
```

Out[4]:

	id	description	quantity	price
0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.0

```
In [5]: y = project_data['project_is_approved'].values
X = project_data.drop(['project_is_approved'], axis=1)
X.head(1)
project_data=None
gc.collect()
gc.enable()
gc.DEBUG_SAVEALL
```

Out[5]: 32

2.2 Getting the Data Model Ready:Preprocessing and Vectorizing categorical features

2.2.1 Preprocessing:project_grade_category

```
In [6]: sub_categories = list(X['project_grade_category'].values)
# remove special characters from list of strings python: https://stackoverflow.co

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

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

```
In [7]: sub_categories=None
sub_cat_list=None
temp=None
i=None
j=None
catogories=None
cat_list=None
temp=None
my_counter=None
word=None
cat_dict=None
gc.collect()
gc.enable()
gc.DEBUG_SAVEALL
```

Out[7]: 32

2.2.2 Preprocessing:project subject categories

```

In [8]: categories = list(X['project_subject_categories'].values)
# remove special characters from list of strings python: https://stackoverflow.co

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

cat_list = []
for i in categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science",
        if 'The' in j.split(): # this will split each of the category based on sp
            j=j.replace('The','') # if we have the words "The" we are going to re
        j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty)
        temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trail
        temp = temp.replace('&','_') # we are replacing the & value into
    cat_list.append(temp.strip())

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

```

Out[8]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_sul
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY	

2.2.3 Preprocessing:project_subject_subcategories

```

In [9]: sub_categories = list(X['project_subject_subcategories'].values)
# remove special characters from list of strings python: https://stackoverflow.co

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

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

```

```

In [10]: X['clean_subcategories'] = sub_cat_list
X.drop(['project_subject_subcategories'], axis=1, inplace=True)
X.head(2)

```

Out[10]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_sul
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY	

2.2.4 New Column:digits in summary

```
In [11]: # Creating a new column 'digits_in_summary' which contains flags of 1 for /
# 'project_resource_summary' containing numeric specification in their requiremn
project_resource_summary = []
new=[]

project_resource_summary = list(X['project_resource_summary'].values)

for i in project_resource_summary:
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(' '):
        if j.isdigit():
            new.append(1)
            break
        else:
            continue
    else:
        new.append(0)

X['digits_in_summary']=new
```

```
In [12]: #To make best use of the memory we are setting the variable names to 'None' and p
project_resource_summary=None
new=None
new1=None
i=None
j=None
a=None

gc.collect()
gc.enable()
gc.DEBUG_SAVEALL
```

Out[12]: 32

2.2.5 Preprocessing:Text features (Project Essay's)

```
In [13]: # merge two column text dataframe:
X["essay"] = X["project_essay_1"].map(str) + \
            X["project_essay_2"].map(str) + \
            X["project_essay_3"].map(str) + \
            X["project_essay_4"].map(str)
```

```
In [14]: X = X.drop(['project_essay_1', 'project_essay_2', 'project_essay_3', 'project_ess
X.shape
```

Out[14]: (50000, 14)

2.2.6 Adding column Cost per project in dataset

```
In [15]: # https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indexes
price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'})
price_data.head(2)
type(price_data)
```

Out[15]: pandas.core.frame.DataFrame

```
In [16]: # join two dataframes in python:
X = pd.merge(X, price_data, on='id', how='left')
X.head(2)
```

Out[16]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_sul
0	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	
1	45	p246581	f3cb9bffba169bef1a77b243e620b60	Mrs.	KY	

```
In [17]: #To make best use of the memory we are setting the variable names to 'None' and
resource_data=None
price_data=None
gc.collect()
gc.enable()
gc.DEBUG_SAVEALL
```

Out[17]: 32

2.2.7 Text Preprocessing:Essay Text

In [18]: [# https://stackoverflow.com/a/47091490/4084039](https://stackoverflow.com/a/47091490/4084039)

```
import re

def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"'re", " are", phrase)
    phrase = re.sub(r"'s", " is", phrase)
    phrase = re.sub(r"'d", " would", phrase)
    phrase = re.sub(r"'ll", " will", phrase)
    phrase = re.sub(r"'t", " not", phrase)
    phrase = re.sub(r"'ve", " have", phrase)
    phrase = re.sub(r"'m", " am", phrase)
    return phrase
```

In [19]: `sent = decontracted(X['essay'].values[99])`
`print(sent)`
`print("="*50)`

```
What does \"The message in the music\" mean to you? To my kids, it means so much more! Music is not just for entertainment! My kids are learning how music can be used as a human rights tool! They are learning to express their feelings of the world around them in positive ways.\n\nMy kids are every stereotype you can think of: poor, underprivileged, rough, angry. Turn that same prism around and there is also joy, happiness, intelligence and creativity. They are learning to navigate the world around them and we are working on showing them how to use this energy in a positive way. Our piano lab is aging. We have 3 classes of 35 kids that are bursting at the gills, which is a great problem to have! BUT keyboards are beginning to have the keys break off making them unuseable. We also do not have enough keyboards for every student in class, even if you counted the slightly broken ones!\n\nWe are in desperate need of more keyboards for our labs to meet the needs of our students. A keyboard for each student in class will facilitate their love of music and help them grow and appreciate the world around them. We are down to 24 semi working keyboards and 35 kids in each piano class. These keyboards are needed sooner rather than later for our kids to have a successful year in class.nannan\n\n=====
```



```
In [20]: # \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks
sent = sent.replace('\r', ' ')
sent = sent.replace('\n', ' ')
sent = sent.replace('\t', ' ')
print(sent)
```

What does The message in the music mean to you? To my kids, it means so much more! Music is not just for entertainment! My kids are learning how music can be used as a human rights tool! They are learning to express their feelings of the world around them in positive ways. My kids are every stereotype you can think of: poor, underprivileged, rough, angry. Turn that same prism around the re is also joy, happiness, intelligence and creativity. They are learning to navigate the world around them and we are working on showing them how to use this energy in a positive way. Our piano lab is aging. We have 3 classes of 35 kids that are bursting at the gills, which is a great problem to have! BUT keyboards are beginning to have the keys break off making them unuseable. We also do not have enough keyboards for every student in class, even if you counted the slightly broken ones! We are in desperate need of more keyboards for our labs to meet the needs of our students. A keyboard for each student in class will facilitate their love of music and help them grow and appreciate the world around them. We are down to 24 semi working keyboards and 35 kids in each piano class. These keyboards are needed sooner rather than later for our kids to have a successful year in class. nannan

```
In [21]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

What does The message in the music mean to you To my kids it means so much more Music is not just for entertainment My kids are learning how music can be used as a human rights tool They are learning to express their feelings of the world around them in positive ways My kids are every stereotype you can think of poor underprivileged rough angry Turn that same prism around there is also joy happiness intelligence and creativity They are learning to navigate the world around them and we are working on showing them how to use this energy in a positive way Our piano lab is aging We have 3 classes of 35 kids that are bursting at the gills which is a great problem to have BUT keyboards are beginning to have the keys break off making them unuseable We also do not have enough keyboards for every student in class even if you counted the slightly broken ones We are in desperate need of more keyboards for our labs to meet the needs of our students A keyboard for each student in class will facilitate their love of music and help them grow and appreciate the world around them We are down to 24 semi working keyboards and 35 kids in each piano class These keyboards are needed sooner rather than later for our kids to have a successful year in class nannan


```
In [26]: preprocessed_project_title[4999]
# after preprocesing

# X['project_title'] = None
X['project_title'] = preprocessed_project_title

X.head(2)
```

Out[26]:

Unnamed: 0		id	teacher_id	teacher_prefix	school_state	project_sul
0	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	
1	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY	

2.2.8 Splitting the data into Train and Test

```
In [27]: # train test split(67:33)
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify=y)
# X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33)
```

```
In [28]: # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()
for word in X_train['clean_categories'].values:
    my_counter.update(word.split())
print(my_counter)
# dict sort by value python: https://stackoverflow.com/a/613218/4084039
cat_dict = dict(my_counter)
sorted_cat_dict_train = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
```

```
Counter({'Literacy_Language': 15726, 'Math_Science': 12827, 'Health_Sports': 4373, 'SpecialNeeds': 4305, 'AppliedLearning': 3801, 'Music_Arts': 3148, 'History_Civics': 1762, 'Warmth': 379, 'Care_Hunger': 379})
```

2.2.9 Vectorizing Categorical data: clean_categories(Project subject categories)

```

In [29]: print(X_train.shape, y_train.shape)
# print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

Bow_features_names1=[]

from sklearn.feature_extraction.text import CountVectorizer

vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict_train.keys()), lower
vectorizer.fit(X_train['clean_categories'].values) # fit has to happen only on t

# we use the fitted Countvectorizer to convert the text to vector
X_train_clean_cat_ohe = vectorizer.transform(X_train['clean_categories'].values)
# X_cv_clean_cat_ohe = vectorizer.transform(X_cv['clean_categories'].values)
X_test_clean_cat_ohe = vectorizer.transform(X_test['clean_categories'].values)

print("After vectorizations")
print(X_train_clean_cat_ohe.shape, y_train.shape)
# print(X_cv_clean_cat_ohe.shape, y_cv.shape)
print(X_test_clean_cat_ohe.shape, y_test.shape)

print(vectorizer.get_feature_names())
# print(vectorizer_test.get_feature_names())
print("="*100)

```

```
(33500, 16) (33500,)
```

```
(16500, 16) (16500,)
```

```
=====
```

```
=====
```

```
After vectorizations
```

```
(33500, 9) (33500,)
```

```
(16500, 9) (16500,)
```

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

```
=====
```

```
=====
```

2.2.10 Vectorizing Categorical data: clean_subcategories(Project subject subcategories)

```

In [30]: # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4
from collections import Counter
my_counter = Counter()
for word in X_train['clean_subcategories'].values:
    my_counter.update(word.split())
sub_cat_dict = dict(my_counter)
sorted_sub_cat_dict_train = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]

```

```
In [31]: print(X_train.shape, y_train.shape)
# print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict_train.keys()),
vectorizer.fit(X_train['clean_subcategories'].values) # fit has to happen only on train

# we use the fitted Countvectorizer to convert the text to vector
X_train_clean_sub_cat_ohe = vectorizer.transform(X_train['clean_subcategories'].values)
# X_cv_clean_sub_cat_ohe = vectorizer.transform(X_cv['clean_subcategories'].values)
X_test_clean_sub_cat_ohe = vectorizer.transform(X_test['clean_subcategories'].values)

print("After vectorizations")
print(X_train_clean_sub_cat_ohe.shape, y_train.shape)
# print(X_cv_clean_sub_cat_ohe.shape, y_cv.shape)
print(X_test_clean_sub_cat_ohe.shape, y_test.shape)

print(vectorizer.get_feature_names())
# print(vectorizer_test.get_feature_names())
print("="*100)
```

```
(33500, 16) (33500,)
```

```
(16500, 16) (16500,)
```

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

```
After vectorizations
```

```
(33500, 30) (33500,)
```

```
(16500, 30) (16500,)
```

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

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

2.2.11 Vectorizing Categorical data: school_state

```
In [32]: # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4
from collections import Counter
my_counter = Counter()
for word in X_train['school_state'].values:
    my_counter.update(word.split())
school_state_dict = dict(my_counter)
sorted_school_state_dict_train = dict(sorted(school_state_dict.items(), key=lambda
# sorted_school_state_dict
```

```
In [33]: print(X_train.shape, y_train.shape)
# print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer = CountVectorizer(vocabulary=list(sorted_school_state_dict_train.keys))
vectorizer.fit(X_train['school_state'].values) # fit has to happen only on train

# we use the fitted Countvectorizer to convert the text to vector
X_train_School_state_ohe = vectorizer.transform(X_train['school_state'].values)
# X_cv_School_state_ohe = vectorizer.transform(X_cv['school_state'].values)
X_test_School_state_ohe = vectorizer.transform(X_test['school_state'].values)

print("After vectorizations")
print(X_train_School_state_ohe.shape, y_train.shape)
# print(X_cv_School_state_ohe.shape, y_cv.shape)
print(X_test_School_state_ohe.shape, y_test.shape)

print(vectorizer.get_feature_names())
# print(vectorizer_test.get_feature_names())
print("="*100)
```

```
(33500, 16) (33500,)
(16500, 16) (16500,)
=====
=====
After vectorizations
(33500, 51) (33500,)
(16500, 51) (16500,)
['VT', 'WY', 'ND', 'MT', 'RI', 'NH', 'AK', 'DE', 'SD', 'NE', 'HI', 'WV', 'ME',
'NM', 'DC', 'IA', 'ID', 'KS', 'AR', 'CO', 'OR', 'MN', 'MS', 'KY', 'MD', 'NV',
'UT', 'AL', 'TN', 'CT', 'WI', 'VA', 'OK', 'NJ', 'AZ', 'WA', 'LA', 'MA', 'OH',
'MO', 'IN', 'PA', 'MI', 'SC', 'GA', 'IL', 'NC', 'FL', 'NY', 'TX', 'CA']
=====
=====
```

2.2.12 Vectorizing Categorical data: project_grade_category

```
In [34]: # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4
from collections import Counter
my_counter = Counter()
for word in X_train['project_grade_category'].values:
    my_counter.update(word.split())
project_grade_dict = dict(my_counter)
sorted_project_grade_dict_train = dict(sorted(project_grade_dict.items(), key=lambda
```

```
In [35]: print(X_train.shape, y_train.shape)
# print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer= CountVectorizer(vocabulary=list(sorted_project_grade_dict_train.keys)
vectorizer.fit(X_train['project_grade_category'].values) # fit has to happen only once

# we use the fitted Countvectorizer_pro_gradeto convert the text to vector
X_train_project_grade_category_ohe = vectorizer.transform(X_train['project_grade_category'])
# X_cv_project_grade_category_ohe = vectorizer.transform(X_cv['project_grade_category'])
X_test_project_grade_category_ohe = vectorizer.transform(X_test['project_grade_category'])

print("After vectorizations")
print(X_train_project_grade_category_ohe.shape, y_train.shape)
# print(X_cv_project_grade_category_ohe.shape, y_cv.shape)
print(X_test_project_grade_category_ohe.shape, y_test.shape)

print(vectorizer.get_feature_names())
# print(vectorizer_test.get_feature_names())
print("="*100)
```

```
(33500, 16) (33500,)
```

```
(16500, 16) (16500,)
```

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

```
After vectorizations
```

```
(33500, 4) (33500,)
```

```
(16500, 4) (16500,)
```

```
['Grades9-12', 'Grades6-8', 'Grades3-5', 'GradesPreK-2']
```

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

2.2.13 Vectorizing Categorical data: teacher_prefix

```
In [36]: #To overcome the blanks in the teacher_prefix category the .fillna is used
X_train['teacher_prefix']=X_train['teacher_prefix'].fillna("")
# project_data1=project_data.dropna()
```

```
In [37]: #To overcome the blanks in the teacher_prefix category the .fillna is used
X_test['teacher_prefix']=X_test['teacher_prefix'].fillna("")
# project_data1=project_data.dropna()
```



```
In [38]: # count of all the words in corpus python: https://stackoverflow.com/a/22898595/
from collections import Counter
my_counter = Counter()
my_counter1=[]
# project_data['teacher_prefix']=str(project_data['teacher_prefix'])
for word in X_train['teacher_prefix'].values:
    my_counter.update(word.split())
teacher_prefix_dict = dict(my_counter)
sorted_teacher_prefix_dict_train = dict(sorted(teacher_prefix_dict.items(), key=
# teacher_prefix_dict
```

```
In [39]: print(X_train.shape, y_train.shape)
# print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer = CountVectorizer(vocabulary=list(sorted_teacher_prefix_dict_train.keys))
vectorizer.fit(X_train['teacher_prefix'].values) # fit has to happen only on training data

# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_prefix_ohe = vectorizer.transform(X_train['teacher_prefix'].values)
# X_cv_teacher_prefix_ohe = vectorizer.transform(X_cv['teacher_prefix'].values)
X_test_teacher_prefix_ohe = vectorizer.transform(X_test['teacher_prefix'].values)

print("After vectorizations")
print(X_train_teacher_prefix_ohe.shape, y_train.shape)
# print(X_cv_teacher_prefix_ohe.shape, y_cv.shape)
print(X_test_teacher_prefix_ohe.shape, y_test.shape)

print(vectorizer.get_feature_names())
# print(vectorizer_test.get_feature_names())
print("="*100)
```

```
(33500, 16) (33500,)
(16500, 16) (16500,)
```

```
=====
=====
After vectorizations
(33500, 5) (33500,)
(16500, 5) (16500,)
['Dr.', 'Teacher', 'Mr.', 'Ms.', 'Mrs.']
=====
=====
```

2.3 Make Data Model Ready: Vectorizing Numerical features

2.3.1 Vectorizing Numerical features--Price

```
In [40]: from sklearn.preprocessing import Normalizer
from sklearn.preprocessing import StandardScaler
normalizer = StandardScaler() #Normalizer()
# normalizer_test = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['price'].values.reshape(-1,1))
# normalizer_test.fit(X_test['price'].values.reshape(1,-1))

X_train_price_norm = normalizer.transform(X_train['price'].values.reshape(-1,1))
# X_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(-1,1))
X_test_price_norm = normalizer.transform(X_test['price'].values.reshape(-1,1))

# X_train_price_norm=np.reshape(X_train_price_norm,(1,-1))
# X_test_price_norm=np.reshape(X_test_price_norm,(1,-1))

print("After vectorizations")

# np.reshape(X_train_price_norm,
print(X_train_price_norm.shape, y_train.shape)
# print(X_cv_price_norm.shape, y_cv.shape)
print(X_test_price_norm.shape, y_test.shape)
print("=*100)
```

After vectorizations

(33500, 1) (33500,)

(16500, 1) (16500,)

=====

=====

2.3.2 Vectorizing Numerical features-- teacher_number_of_previously_posted_projects

```
In [41]: from sklearn.preprocessing import Normalizer
normalizer_train = StandardScaler() #Normalizer()
normalizer_test = StandardScaler() #Normalizer()
# normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.re

X_train_teacher_number_of_previously_posted_projects_norm = normalizer.transform
# X_cv_teacher_number_of_previously_posted_projects_norm = normalizer.transform(
X_test_teacher_number_of_previously_posted_projects_norm = normalizer.transform(

# X_train_teacher_number_of_previously_posted_projects_norm=np.reshape(X_train_te
# X_test_teacher_number_of_previously_posted_projects_norm=np.reshape(X_test_tea

print("After vectorizations")
print(X_train_teacher_number_of_previously_posted_projects_norm.shape, y_train.sh
# print(X_cv_teacher_number_of_previously_posted_projects_norm.shape, y_cv.shape,
print(X_test_teacher_number_of_previously_posted_projects_norm.shape, y_test.shap
print("=*100)
```

After vectorizations

(33500, 1) (33500,)

(16500, 1) (16500,)

=====
=====

2.3.3 Vectorizing Numerical features--digits_in_summary

```
In [42]: X_train['digits_in_summary'].fillna(X_train['digits_in_summary'].mean(), inplace=True)
# X_cv['digits_in_summary'].fillna(X_cv['digits_in_summary'].mean(), inplace=True)
X_test['digits_in_summary'].fillna(X_test['digits_in_summary'].mean(), inplace=True)
```

```
In [43]: from sklearn.preprocessing import Normalizer
normalizer_train = StandardScaler() #Normalizer()
normalizer_test = StandardScaler() #Normalizer()
# normalizer.fit(X_train['digits_in_summary'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['digits_in_summary'].values.reshape(-1,1))

X_train_digits_in_summary_norm = normalizer.transform(X_train['digits_in_summary'])
# X_cv_digits_in_summary_norm = normalizer.transform(X_cv['digits_in_summary'].values.reshape(-1,1))
X_test_digits_in_summary_norm = normalizer.transform(X_test['digits_in_summary'].values.reshape(-1,1))

# X_train_digits_in_summary_norm=np.reshape(X_train_digits_in_summary_norm,(1,-1))
# X_test_digits_in_summary_norm=np.reshape(X_test_digits_in_summary_norm,(1,-1))

print("After vectorizations")
print(X_train_digits_in_summary_norm.shape, y_train.shape)
# print(X_cv_digits_in_summary_norm.shape, y_cv.shape)
print(X_test_digits_in_summary_norm.shape, y_test.shape)
print("=="*100)
```

After vectorizations

(33500, 1) (33500,)

(16500, 1) (16500,)

=====

2.4 Make Data Model Ready: Vectorizing Essay and Project_title feature into BOW & TFIDF

Vectorizing Text data

2.4.1 Bag of words:Essays

```

In [44]: print(X_train.shape, y_train.shape)
          # print(X_cv.shape, y_cv.shape)
          print(X_test.shape, y_test.shape)

          print("="*100)

          # We are considering only the words which appeared in at least 10 documents(rows

          vectorizer = CountVectorizer(min_df=10,ngram_range=(2,2), max_features=5000)
          vectorizer.fit(X_train['essay'].values) # fit has to happen only on train data

          # we use the fitted Countvectorizer to convert the text to vector
          X_train_essay_bow = vectorizer.transform(X_train['essay'].values)
          # X_cv_essay_bow = vectorizer.transform(X_cv['essay'].values)
          X_test_essay_bow = vectorizer.transform(X_test['essay'].values)

          print("After vectorizations")
          print(X_train_essay_bow.shape, y_train.shape)
          # print(X_cv_essay_bow.shape, y_cv.shape)
          print(X_test_essay_bow.shape, y_test.shape)
          print("="*100)

```

```
(33500, 16) (33500,)
```

```
(16500, 16) (16500,)
```

```

=====
=====
After vectorizations
(33500, 5000) (33500,)
(16500, 5000) (16500,)
=====
=====

```

2.4.2 Bag of words:Project Title

In [45]:

```

print(X_train.shape, y_train.shape)
# print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer = CountVectorizer(min_df=10,ngram_range=(2,2), max_features=5000)
vectorizer.fit(X_train['project_title'].values) # fit has to happen only on train

# we use the fitted Countvectorizer to convert the text to vector
X_train_project_title_bow = vectorizer.transform(X_train['project_title'].values)
# X_cv_project_title_bow = vectorizer.transform(X_cv['project_title'].values)
X_test_project_title_bow = vectorizer.transform(X_test['project_title'].values)

print("After vectorizations")
print(X_train_project_title_bow.shape, y_train.shape)
# print(X_cv_project_title_bow.shape, y_cv.shape)
print(X_test_project_title_bow.shape, y_test.shape)

print("="*100)

```

(33500, 16) (33500,)

(16500, 16) (16500,)

```

=====
=====
After vectorizations
(33500, 1041) (33500,)
(16500, 1041) (16500,)
=====
=====

```

In [46]:

```

from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense
X_BOW_TRAIN = hstack((X_train_digits_in_summary_norm,X_train_teacher_number_of_projects))
X_BOW_TRAIN=X_BOW_TRAIN.todense()
X_BOW_TRAIN=np.array(X_BOW_TRAIN)

# X_BOW_cv = hstack((X_cv_project_title_bow,X_cv_essay_bow ,X_cv_digits_in_summary_norm))
# X_BOW_cv=X_BOW_cv.todense()
# X_BOW_cv=np.array(X_BOW_cv)

X_BOW_test = hstack((X_test_digits_in_summary_norm,X_test_teacher_number_of_projects))
X_BOW_test=X_BOW_test.todense()
X_BOW_test=np.array(X_BOW_test)

X_train_project_title_bow=None
X_train_essay_bow =None

X_test_project_title_bow=None
X_test_essay_bow =None

```

2.5 Applying SVM on BOW , SET 1

2.5.1 Applying SVM & GridSearchCV on Train data to obtain the best C

```
In [59]: # # Selecting 2000 best features from Tfidf to see the variation in the AUC  
# from sklearn.feature_selection import SelectKBest, f_classif  
# X_BOW_TRAIN = SelectKBest(f_classif, k=5000).fit_transform(X_BOW_TRAIN,y_train,  
# X_BOW_TRAIN.shape  
# X_BOW_test = SelectKBest(f_classif, k=5000).fit_transform(X_BOW_test,y_test)  
# X_BOW_test.shape
```

```
In [48]: #code source: http://occam.olin.edu/sites/default/files/DataScienceMaterials/mach
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.linear_model import LogisticRegression
from sklearn.linear_model import SGDClassifier
from sklearn.metrics import f1_score
from sklearn.calibration import CalibratedClassifierCV
from sklearn.svm import LinearSVC
from sklearn.svm import SVC

# data = load_breast_cancer() #refer: http://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html
tuned_parameters = {'alpha': [10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2]}
# X_train, X_test, y_train, y_test = train_test_split(data.data, data.target, train_size=0.7, random_state=0)

# LR=SGDClassifier(loss='log',class_weight="balanced")
#Using GridSearchCV
# SGDClassifier(loss='log',class_weight="balanced",max_iter=1500)
# LogisticRegression(penalty = 'l1', C = i,class_weight="balanced")
model = GridSearchCV(SGDClassifier(penalty = 'l2',loss='hinge',class_weight="balanced"),
# model=GridSearchCV(SVC(kernel='rbf',class_weight="balanced",max_iter = 2000),
model.fit(X_BOW_TRAIN, y_train)

print(model.best_estimator_)
print(model.score(X_BOW_test, y_test))
```

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:
1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:
1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:
1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:
1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:
1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:
1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

```
SGDClassifier(alpha=100, average=False, class_weight='balanced',
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1500, n_iter_no_change=5, n_jobs=None, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
0.8095238095238095
```

```

In [49]: # from sklearn.model_selection import GridSearchCV

# from sklearn.naive_bayes import MultinomialNB
# # nb = MultinomialNB(class_prior=[0.5,0.5])

# parameters = [0.00001, 0.0001,0.001, 0.01, 0.1,0.5,0.8, 1, 10, 100, 1000]
alpha=[*tuned_parameters.values()]
alpha=alpha[0]
log_alphas=[math.log10(num) for num in alpha]

# clf = GridSearchCV(nb, parameters, cv= 10, scoring='roc_auc',return_train_score=False)

# clf.fit(X_BOW_TRAIN, y_train)

#Selecting the best parameter
best_alpha1=model.best_params_
# print(best_alpha1)
train_auc= model.cv_results_['mean_train_score']
# print(train_auc)
train_auc_std= model.cv_results_['std_train_score']
# print(train_auc_std)
cv_auc = model.cv_results_['mean_test_score']
# print(cv_auc)
cv_auc_std= model.cv_results_['std_test_score']
# print(cv_auc_std)
# print(log_alphas)

plt.figure(figsize=(20,10))

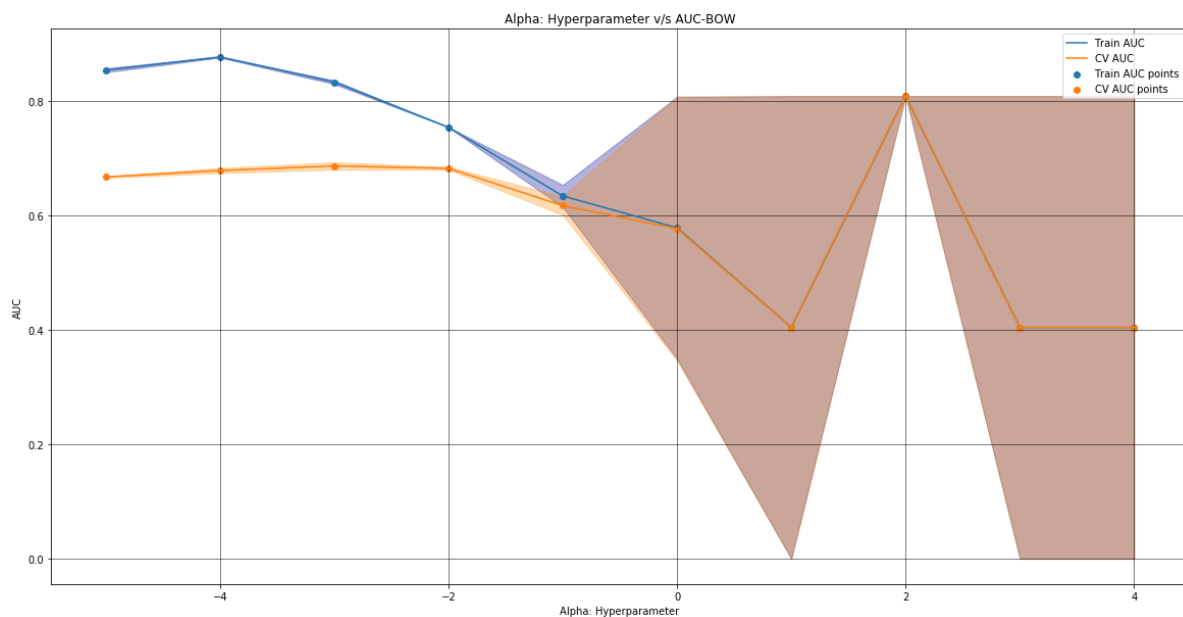
plt.plot(log_alphas, train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_alphas,train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.5)

plt.plot(log_alphas, cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_alphas,cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.5)

plt.scatter(log_alphas, train_auc, label='Train AUC points')
plt.scatter(log_alphas, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("Alpha: Hyperparameter")
plt.ylabel("AUC")
plt.title("Alpha: Hyperparameter v/s AUC-BOW")
plt.grid(color='black', linestyle='--', linewidth=0.5)
plt.show()

```



```
In [50]: print("The best alpha from the above graph is ",best_alpha1)
          best_alpha1['alpha']
```

The best alpha from the above graph is {'alpha': 100}

Out[50]: 100

2.5.2 Receiver Operating Characteristic- (BOW)

```

In [51]: # https://scikit-Learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import SGDClassifier
from sklearn.calibration import CalibratedClassifierCV
from sklearn.svm import SVC
# train model on the best k
# SVM=CalibratedClassifierCV(SGDClassifier(alpha=best_alpha1['C'],penalty = 'l2',

SVM=CalibratedClassifierCV(SGDClassifier(alpha=best_alpha1['alpha'],penalty = 'l2',

SVM.fit(X_BOW_TRAIN, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimate
# not the predicted outputs

probs = SVM.predict_proba(X_BOW_test)
# print(len(probs[:,1]))
probs1 = SVM.predict_proba(X_BOW_TRAIN)
# print(len(probs1[:,1]))
preds = probs[:,1]
# print(preds)
preds1 = probs1[:,1]
# print(preds1)
fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
# print(fpr,tpr,threshold)
fpr1, tpr1, threshold_1 = metrics.roc_curve(y_train, preds1)
# print(fpr1,tpr1,threshold)
roc_auc = metrics.auc(fpr, tpr)
print(roc_auc)
roc_auc1 = metrics.auc(fpr1, tpr1)
print(roc_auc1)

# https://www.programcreek.com/python/example/81207/skLearn.metrics.roc_curve

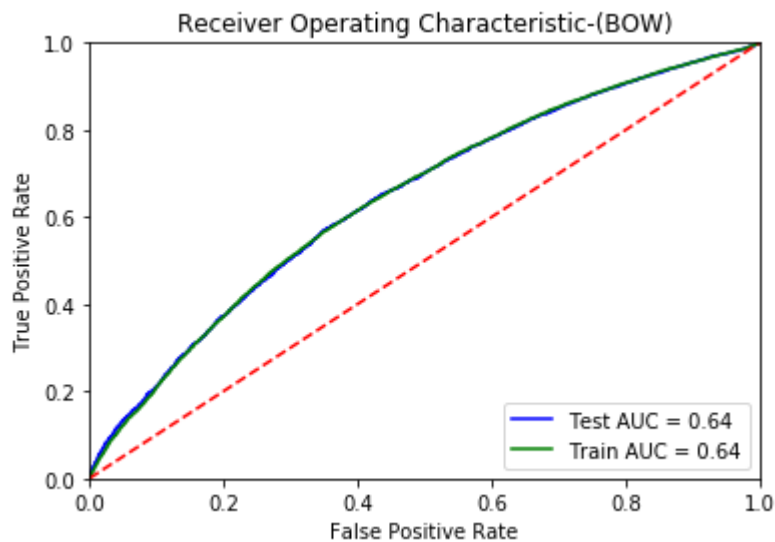
import matplotlib.pyplot as plt
plt.title('Receiver Operating Characteristic-(BOW)')
plt.plot(fpr, tpr, 'b', label = 'Test AUC = %0.2f' % roc_auc)
plt.plot(fpr1, tpr1, 'g', label = 'Train AUC = %0.2f' % roc_auc1)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'r--')
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()

# print(model.best_estimator_)

```

0.6427751951331497

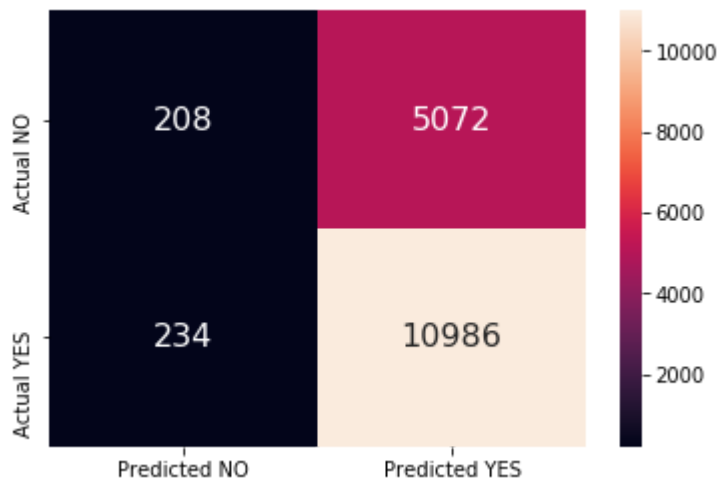
0.6427527870415264



2.5.3 Confusion matrix- BOW

```
In [52]: # https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#function to get heatmap confusion matrix
def get_confusion_matrix(clf,X_te,y_test):
    y_pred = clf.predict(X_te)
    df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), index=['Actual NO','Actual YES'], columns=['Predicted NO','Predicted YES'])
    sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')

# %%time
get_confusion_matrix(SVM,X_BOW_test,y_test)
```



```
In [53]: #To make best use of the memory we are setting the variable names to 'None' and deleting them
X_BOW_TRAIN=None
y_BOW_train=None
gc.collect()
```

Out[53]: 5784

2.6 TFIDF vectorizer

2.6.1 TFIDF vectorizer:Essay

```
In [45]: print(X_train.shape, y_train.shape)
# print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

from sklearn.feature_extraction.text import TfidfVectorizer

vectorizer = TfidfVectorizer(min_df=10,ngram_range=(2,2), max_features=5000)
vectorizer.fit(X_train['essay'].values) # fit has to happen only on train data

# we use the fitted Countvectorizer to convert the text to vector
X_train_essay_Tfidf = vectorizer.transform(X_train['essay'].values)
# X_cv_essay_Tfidf = vectorizer.transform(X_cv['essay'].values)
X_test_essay_Tfidf = vectorizer.transform(X_test['essay'].values)

print("After vectorizations")
print(X_train_essay_Tfidf.shape, y_train.shape)
# print(X_cv_essay_Tfidf.shape, y_cv.shape)
print(X_test_essay_Tfidf.shape, y_test.shape)
print("="*100)

(33500, 16) (33500,)
(16500, 16) (16500,)
=====
=====
After vectorizations
(33500, 5000) (33500,)
(16500, 5000) (16500,)
=====
=====
```

2.6.2 TFIDF vectorizer:Project Title

```
In [55]: # print(X_train.shape, y_train.shape)
# print(X_cv.shape, y_cv.shape)
# print(X_test.shape, y_test.shape)

print("="*100)

# We are considering only the words which appeared in at least 10 documents(rows
vectorizer = TfidfVectorizer(min_df=10,ngram_range=(2,2), max_features=5000)
vectorizer.fit(X_train['project_title'].values) # fit has to happen only on train

# we use the fitted Countvectorizer to convert the text to vector
X_train_project_title_tfidf = vectorizer.transform(X_train['project_title'].values)
# X_cv_project_title_tfidf = vectorizer.transform(X_cv['project_title'].values)
X_test_project_title_tfidf = vectorizer.transform(X_test['project_title'].values)

print("After vectorizations")
print(X_train_project_title_tfidf.shape, y_train.shape)
# print(X_cv_project_title_tfidf.shape, y_cv.shape)
print(X_test_project_title_tfidf.shape, y_test.shape)
print("="*100)
```

```
=====
=====
After vectorizations
(33500, 1034) (33500,)
(16500, 1034) (16500,)
=====
=====
```

```
In [56]: X_Tfidf_train = hstack(( X_train_project_title_tfidf,X_train_essay_Tfidf,X_train_digits_in
X_Tfidf_train=X_Tfidf_train.todense()
X_Tfidf_train=np.array(X_Tfidf_train)

# X_Tfidf_cv = hstack(( X_cv_project_title_tfidf,X_cv_essay_Tfidf,X_cv_digits_in
# X_Tfidf_cv=X_Tfidf_cv.todense()
# X_Tfidf_cv=np.array(X_Tfidf_cv)

X_Tfidf_test = hstack(( X_test_project_title_tfidf,X_test_essay_Tfidf,X_test_digits_in
X_Tfidf_test=X_Tfidf_test.todense()
X_Tfidf_test=np.array(X_Tfidf_test)
```

2.6.3 Applying SVM on TFIDF , SET 2

Applying SVM & GridSearchCV on Train data to obtain the best C

```
In [57]: # # Selecting 2000 best features from Tfidf to see the variation in the AUC  
# from sklearn.feature_selection import SelectKBest, f_classif  
# X_Tfidf_train = SelectKBest(f_classif, k=5000).fit_transform(X_Tfidf_train,y_train)  
# X_Tfidf_train.shape  
# X_Tfidf_test = SelectKBest(f_classif, k=5000).fit_transform(X_Tfidf_test,y_test)  
# X_Tfidf_test.shape
```



```
In [58]: #code source: http://occam.olin.edu/sites/default/files/DataScienceMaterials/mach
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.linear_model import LogisticRegression

# data = load_breast_cancer() #refer: http://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html

tuned_parameters = {'alpha': [10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2]}
# X_train, X_test, y_train, y_test = train_test_split(data.data, data.target, train_size=0.7, random_state=0)

#Using GridSearchCV
model = GridSearchCV(SGDClassifier(penalty = 'l2', loss='hinge', class_weight="balanced"), tuned_parameters, cv=5, scoring='f1')
model.fit(X_Tfidf_train, y_train)

print(model.best_estimator_)
print(model.score(X_Tfidf_test, y_test))
```

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1437: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

```
SGDClassifier(alpha=10, average=False, class_weight='balanced',  
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,  
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',  
              max_iter=1500, n_iter_no_change=5, n_jobs=None, penalty='l2',  
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,  
              validation_fraction=0.1, verbose=0, warm_start=False)  
0.8095238095238095
```

```

In [59]: # from sklearn.model_selection import GridSearchCV

# from sklearn.naive_bayes import MultinomialNB
# # nb = MultinomialNB(class_prior=[0.5,0.5])

# parameters = [0.00001, 0.0001,0.001, 0.01, 0.1,0.5,0.8, 1, 10, 100, 1000]
alpha=[*tuned_parameters.values()]
alpha=alpha[0]
log_alphas=[math.log10(num) for num in alpha]

# clf = GridSearchCV(nb, parameters, cv= 10, scoring='roc_auc',return_train_score=False)

# clf.fit(X_BOW_TRAIN, y_train)

#Selecting the best parameter
best_alpha1=model.best_params_
# print(best_alpha1)
train_auc= model.cv_results_['mean_train_score']
# print(train_auc)
train_auc_std= model.cv_results_['std_train_score']
# print(train_auc_std)
cv_auc = model.cv_results_['mean_test_score']
# print(cv_auc)
cv_auc_std= model.cv_results_['std_test_score']
# print(cv_auc_std)
print(log_alphas)

plt.figure(figsize=(20,10))

plt.plot(log_alphas, train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_alphas,train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.5)

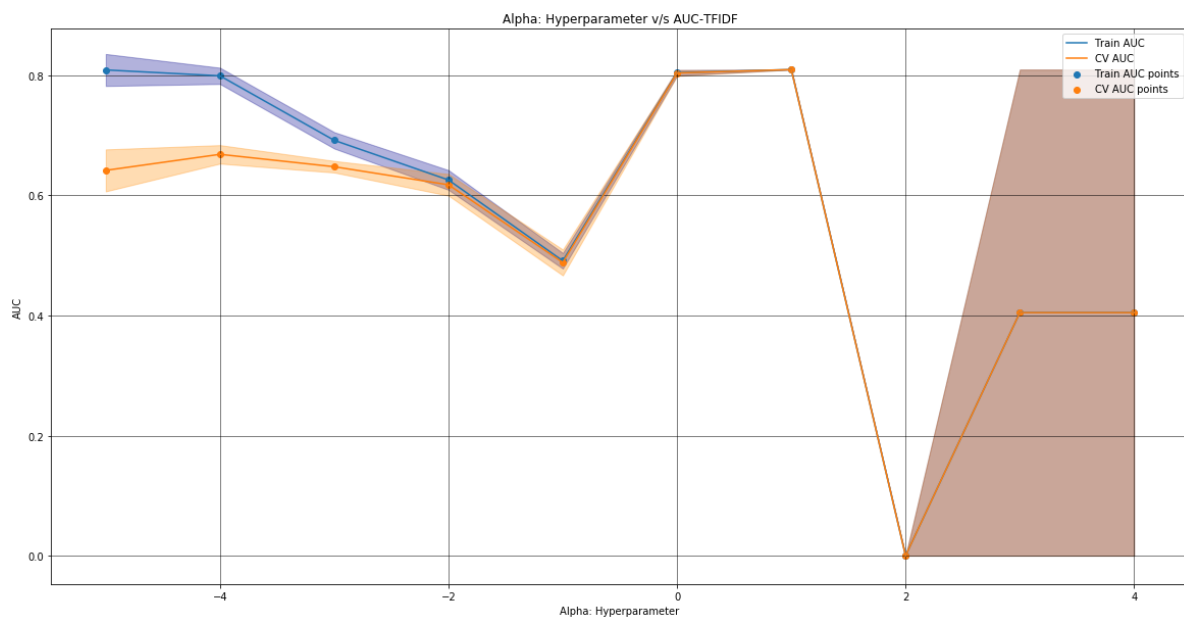
plt.plot(log_alphas, cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_alphas,cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.5)

plt.scatter(log_alphas, train_auc, label='Train AUC points')
plt.scatter(log_alphas, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("Alpha: Hyperparameter")
plt.ylabel("AUC")
plt.title("Alpha: Hyperparameter v/s AUC-TFIDF")
plt.grid(color='black', linestyle='--', linewidth=0.5)
plt.show()

```

```
[-5.0, -4.0, -3.0, -2.0, -1.0, 0.0, 1.0, 2.0, 3.0, 4.0]
```



```
In [60]: # best_alpha1=0.01  
print("The best alpha from the above graph is ",best_alpha1['alpha'])
```

The best alpha from the above graph is 10

2.6.4 Receiver Operating Characteristic- (TFIDF)

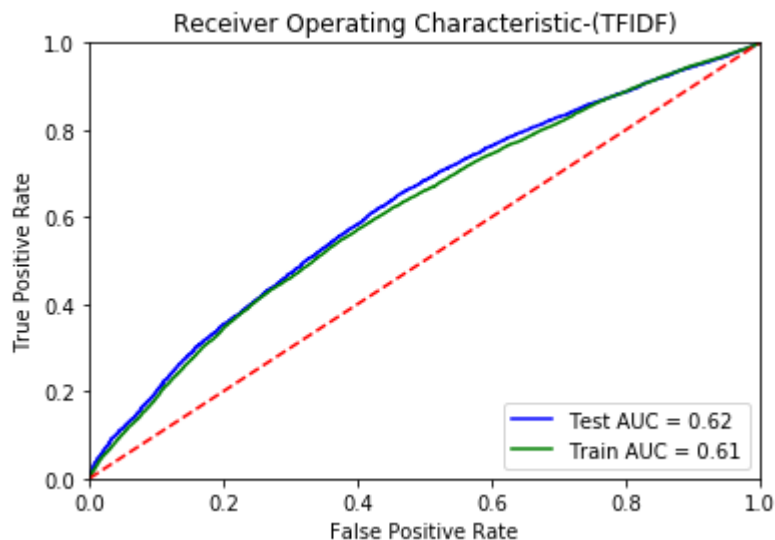
```
In [61]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html
from sklearn.metrics import roc_curve, auc

# train model on the best k
SVM=CalibratedClassifierCV(SGDClassifier(alpha=best_alpha1['alpha'],penalty = 'l2'))
SVM.fit(X_Tfidf_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates
# not the predicted outputs

probs = SVM.predict_proba(X_Tfidf_test)
# print(len(probs[:,1]))
probs1 = SVM.predict_proba(X_Tfidf_train)
# print(len(probs1[:,1]))
preds = probs[:,1]
# print(preds)
preds1 = probs1[:,1]
# print(preds1)
fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
# print(fpr,tpr,threshold)
fpr1, tpr1, threshold = metrics.roc_curve(y_train, preds1)
# print(fpr1,tpr1,threshold)
roc_auc = metrics.auc(fpr, tpr)
# print(roc_auc)
roc_auc1 = metrics.auc(fpr1, tpr1)
# print(roc_auc1)

# https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve

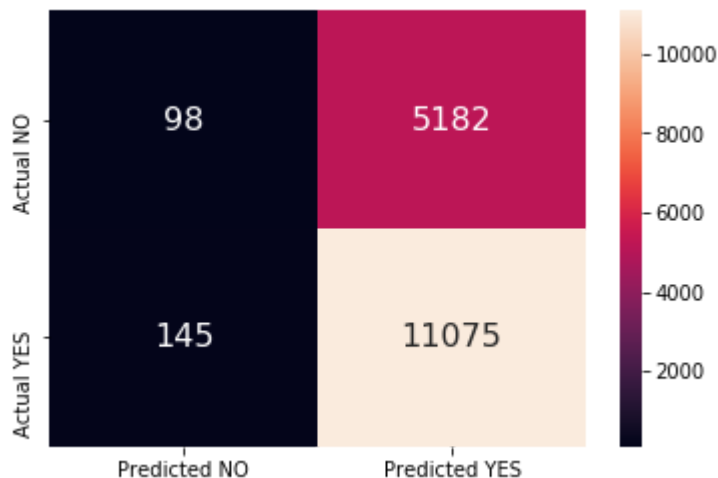
import matplotlib.pyplot as plt
plt.title('Receiver Operating Characteristic-(TFIDF)')
plt.plot(fpr, tpr, 'b', label = 'Test AUC = %0.2f' % roc_auc)
plt.plot(fpr1, tpr1, 'g', label = 'Train AUC = %0.2f' % roc_auc1)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'r--')
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```



2.6.5 Confusion matrix- TFIDF

```
In [62]: # https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#function to get heatmap confusion matrix
def get_confusion_matrix(clf,X_te,y_test):
    y_pred = clf.predict(X_te)
    df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), index=['Actual NO','Actual YES'], columns=['Predicted NO','Predicted YES'])
    sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')

# %%time
get_confusion_matrix(SVM,X_Tfidf_test,y_test)
```



```
In [63]: #To make best use of the memory we are setting the variable names to 'None' and deleting them
X_Tfidf_test=None
X_Tfidf_train=None
gc.collect()
```

Out[63]: 3158

2.7 AVG_W2V

2.7.1 Using Pretrained Models: Avg W2V-Essays

```
In [79]: # train test split
from sklearn.model_selection import train_test_split

avg_w2v_vectors_essay_train_, avg_w2v_vectors_essay_test_, y_train, y_test = tra

# avg_w2v_vectors_essay_train, avg_w2v_vectors_essay_cv, y_train, y_cv = train_te
```

```
In [80]: # stronging variables into pickle files python: http://www.jessicayung.com/how-to
# make sure you have the glove_vectors file
import pickle
with open('glove_vectors', 'rb') as f:
    model_ = pickle.load(f)
    glove_words = set(model_.keys())
```

```
In [81]: # average Word2Vec
# compute average word2vec for each review.
from tqdm import tqdm
# sentence=[]
avg_w2v_vectors_essay_train = []; # the avg-w2v for each sentence/review is store
for sentence in tqdm(avg_w2v_vectors_essay_train_): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model_[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_essay_train.append(vector)

print(len(avg_w2v_vectors_essay_train))
print(len(avg_w2v_vectors_essay_train[0]))
```

100%|██| 33500/33500 [01:12<00:00, 462.71it/s]

33500

300


```
In [88]: import scipy
avg_w2v_vectors_Pro_title_train=scipy.sparse.csr_matrix(avg_w2v_vectors_Pro_title_train)
type(avg_w2v_vectors_Pro_title_train)

import scipy
avg_w2v_vectors_Pro_title_test=scipy.sparse.csr_matrix(avg_w2v_vectors_Pro_title_test)
type(avg_w2v_vectors_Pro_title_test)

# import scipy
# avg_w2v_vectors_Pro_title_cv=scipy.sparse.csr_matrix(avg_w2v_vectors_Pro_title_cv)
# type(avg_w2v_vectors_Pro_title_cv)
```

Out[88]: scipy.sparse.csr.csr_matrix

```
In [89]: X_avg_w2v_train = hstack(( avg_w2v_vectors_Pro_title_train,avg_w2v_vectors_essay_train))
X_avg_w2v_train=X_avg_w2v_train.todense()
X_avg_w2v_train=np.array(X_avg_w2v_train)

# X_avg_w2v_cv = hstack(( avg_w2v_vectors_Pro_title_cv,avg_w2v_vectors_essay_cv))
# X_avg_w2v_cv=X_avg_w2v_cv.todense()
# X_avg_w2v_cv=np.array(X_avg_w2v_cv)

X_avg_w2v_test = hstack(( avg_w2v_vectors_Pro_title_test,avg_w2v_vectors_essay_test))
X_avg_w2v_test=X_avg_w2v_test.todense()
X_avg_w2v_test=np.array(X_avg_w2v_test)

avg_w2v_vectors_Pro_title_train=None
avg_w2v_vectors_essay_train=None

avg_w2v_vectors_Pro_title_test=None
avg_w2v_vectors_essay_test=None
```

2.7.3 Applying SVM on AVG_W2V , SET 3

Applying SVM & GridSearchCV on Train data to obtain the best C

```
In [96]: #code source: http://occam.olin.edu/sites/default/files/DataScienceMaterials/mach
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.linear_model import LogisticRegression

# data = load_breast_cancer() #refer: http://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html

tuned_parameters = {'alpha': [10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2]}
# X_train, X_test, y_train, y_test = train_test_split(data.data, data.target, train_size=0.7, random_state=0)

#Using GridSearchCV
model = GridSearchCV(SGDClassifier(penalty = 'l2', loss='hinge', class_weight="balanced"), tuned_parameters, cv=5, scoring='accuracy')
model.fit(X_avg_w2v_train, y_train)

print(model.best_estimator_)
print(model.score(X_avg_w2v_test, y_test))
```

```
SGDClassifier(alpha=1000, average=False, class_weight='balanced',
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1500, n_iter_no_change=5, n_jobs=None, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
```

0.68

```

In [97]: # from sklearn.model_selection import GridSearchCV

# from sklearn.naive_bayes import MultinomialNB
# # nb = MultinomialNB(class_prior=[0.5,0.5])

# parameters = [0.00001, 0.0001,0.001, 0.01, 0.1,0.5,0.8, 1, 10, 100, 1000]
alpha=[*tuned_parameters.values()]
alpha=alpha[0]
log_alphas=[math.log10(num) for num in alpha]

# clf = GridSearchCV(nb, parameters, cv= 10, scoring='roc_auc',return_train_score=False)

# clf.fit(X_BOW_TRAIN, y_train)

#Selecting the best parameter
best_alpha1=model.best_params_
# print(best_alpha1)
train_auc= model.cv_results_['mean_train_score']
# print(train_auc)
train_auc_std= model.cv_results_['std_train_score']
# print(train_auc_std)
cv_auc = model.cv_results_['mean_test_score']
# print(cv_auc)
cv_auc_std= model.cv_results_['std_test_score']
# print(cv_auc_std)
print(log_alphas)

plt.figure(figsize=(20,10))

plt.plot(log_alphas, train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_alphas,train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.5)

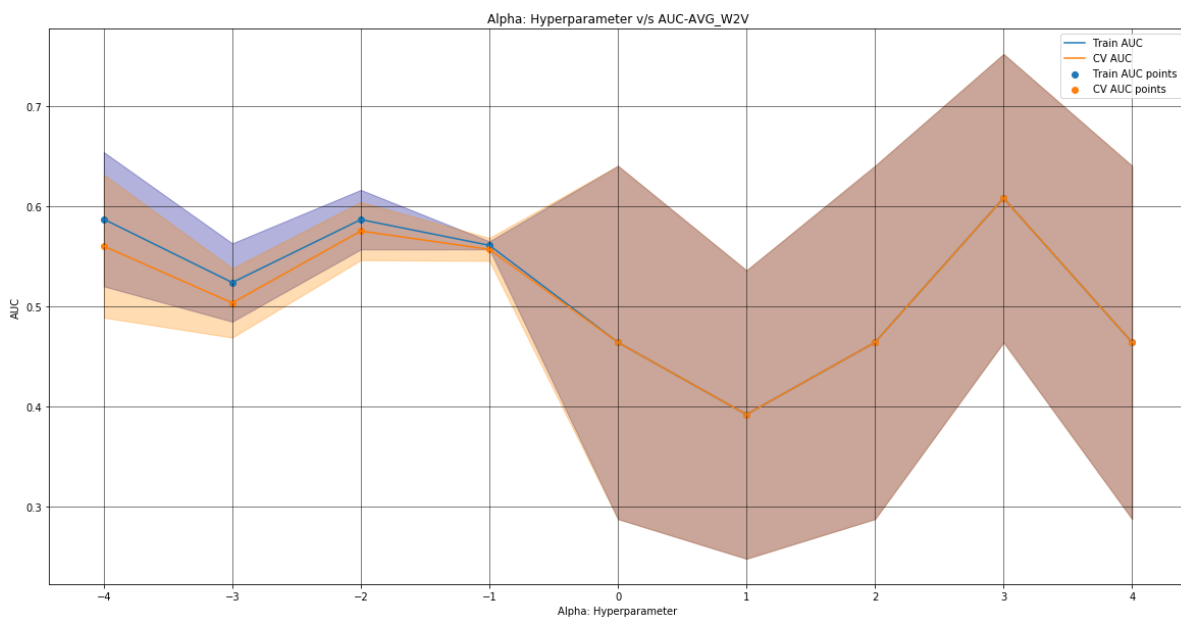
plt.plot(log_alphas, cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_alphas,cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.5)

plt.scatter(log_alphas, train_auc, label='Train AUC points')
plt.scatter(log_alphas, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("Alpha: Hyperparameter")
plt.ylabel("AUC")
plt.title("Alpha: Hyperparameter v/s AUC-AVG_W2V")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()

```

```
[-4.0, -3.0, -2.0, -1.0, 0.0, 1.0, 2.0, 3.0, 4.0]
```



```
In [98]: print("The best alpha from the above graph is ",best_alpha1['alpha'])  
# best_alpha1['alpha']
```

The best alpha from the above graph is 1000

2.7.4 Receiver Operating Characteristic- (AVG_W2V)

```

In [99]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html
from sklearn.metrics import roc_curve, auc

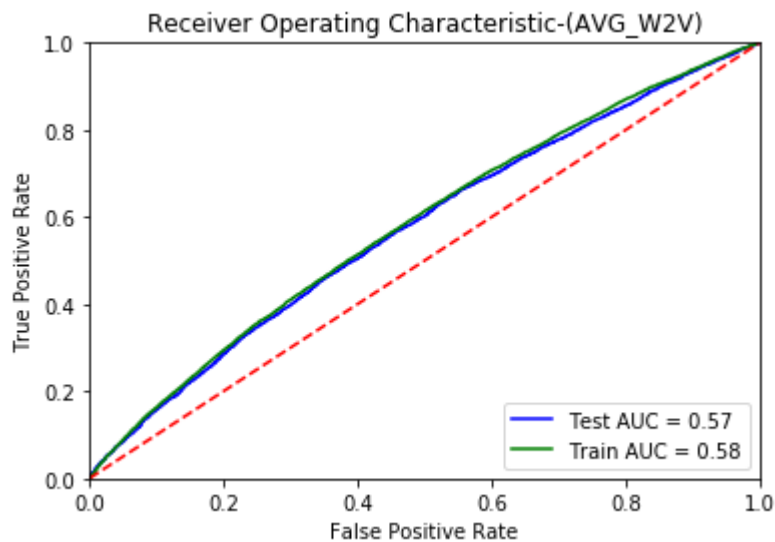
# train model on the best k
SVM=CalibratedClassifierCV(SGDClassifier(alpha=best_alpha1['alpha'],penalty = 'l2'))
SVM.fit(X_avg_w2v_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates
# not the predicted outputs

probs = SVM.predict_proba(X_avg_w2v_test)
# print(len(probs[:,1]))
probs1 = SVM.predict_proba(X_avg_w2v_train)
# print(len(probs1[:,1]))
preds = probs[:,1]
# print(preds)
preds1 = probs1[:,1]
# print(preds1)
fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
# print(fpr,tpr,threshold)
fpr1, tpr1, threshold = metrics.roc_curve(y_train, preds1)
# print(fpr1,tpr1,threshold)
roc_auc = metrics.auc(fpr, tpr)
# print(roc_auc)
roc_auc1 = metrics.auc(fpr1, tpr1)
# print(roc_auc1)

# https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve

import matplotlib.pyplot as plt
plt.title('Receiver Operating Characteristic-(AVG_W2V)')
plt.plot(fpr, tpr, 'b', label = 'Test AUC = %0.2f' % roc_auc)
plt.plot(fpr1, tpr1, 'g', label = 'Train AUC = %0.2f' % roc_auc1)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'r--')
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()

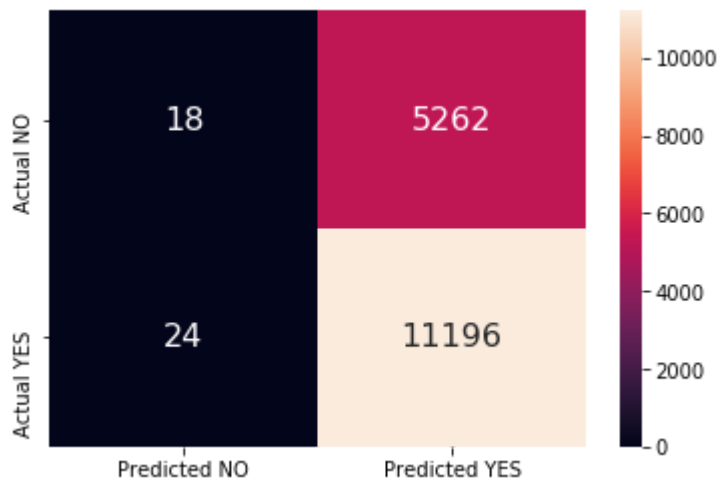
```



2.7.5 Confusion matrix-AVG_W2V

```
In [100]: # https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#function to get heatmap confusion matrix
def get_confusion_matrix(clf,X_te,y_test):
    y_pred = clf.predict(X_te)
    df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), index=['Actual NO','Actual YES'], columns=['Predicted NO','Predicted YES'])
    sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')

# %%time
get_confusion_matrix(SVM,X_avg_w2v_test,y_test)
```



```
In [101]: #To make best use of the memory we are setting the variable names to 'None' and deleting them
X_avg_w2v_test=None
X_avg_w2v_train=None
y_avg_w2v_train=None
gc.collect()
```

Out[101]: 9064

2.8 TFIDF weighted W2V-Essav

2.8.1 Using Pretrained Models: TFIDF weighted W2V-Essay

```
In [102]: # train test split
from sklearn.model_selection import train_test_split

Xtfidf_w2v_vectors_train_, Xtfidf_w2v_vectors_test_, y_train, y_test = train_test_split(Xtfidf_w2v_vectors_, y, test_size=0.2, random_state=42)

# Xtfidf_w2v_vectors_train, Xtfidf_w2v_vectors_cv, y_train, y_cv = train_test_split(Xtfidf_w2v_vectors_, y, test_size=0.2, random_state=42)
```



```
In [109]: # average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_Pro_title_train = []; # the avg-w2v for each sentence/review is
for sentence in tqdm(tfidf_w2v_vectors_Pro_title_train_): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_Pro_title_train.append(vector)

print(len(tfidf_w2v_vectors_Pro_title_train))
print(len(tfidf_w2v_vectors_Pro_title_train[0]))
```

100%|██| 33500/33500 [00:08<00:00, 3828.79it/s]

33500

300

```
In [110]: # average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_Pro_title_test = []; # the avg-w2v for each sentence/review is
for sentence in tqdm(tfidf_w2v_vectors_Pro_title_test_): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_Pro_title_test.append(vector)

print(len(tfidf_w2v_vectors_Pro_title_test))
print(len(tfidf_w2v_vectors_Pro_title_test[0]))
```

100%|██| 16500/16500 [00:04<00:00, 3643.78it/s]

16500

300

```
In [111]: import scipy
tfidf_w2v_vectors_Pro_title_train=scipy.sparse.csr_matrix(tfidf_w2v_vectors_Pro_title_train)
type(tfidf_w2v_vectors_Pro_title_train)

tfidf_w2v_vectors_Pro_title_test=scipy.sparse.csr_matrix(tfidf_w2v_vectors_Pro_title_test)
type(tfidf_w2v_vectors_Pro_title_test)

# tfidf_w2v_vectors_Pro_title_cv=scipy.sparse.csr_matrix(tfidf_w2v_vectors_Pro_title_cv)
# type(tfidf_w2v_vectors_Pro_title_cv)
```

Out[111]: scipy.sparse.csr.csr_matrix

```
In [112]: X_tfidf_w2v_train = hstack((tfidf_w2v_vectors_essay_train,tfidf_w2v_vectors_Pro_title_train))
X_tfidf_w2v_train=X_tfidf_w2v_train.todense()
X_tfidf_w2v_train=np.array(X_tfidf_w2v_train)

# X_tfidf_w2v_cv = hstack((X_tfidf_w2v_vectors_cv,tfidf_w2v_vectors_Pro_title_cv))
# X_tfidf_w2v_cv=X_tfidf_w2v_cv.todense()
# X_tfidf_w2v_cv=np.array(X_tfidf_w2v_cv)

X_tfidf_w2v_test = hstack((tfidf_w2v_vectors_essay_test,tfidf_w2v_vectors_Pro_title_test))
X_tfidf_w2v_test=X_tfidf_w2v_test.todense()
X_tfidf_w2v_test=np.array(X_tfidf_w2v_test)
# X_ALL = hstack((categories_one_hot,sub_categories_one_hot,school_state_one_hot,gender_one_hot))
```

2.8.3 Applying SVM on TFIDF W2V , SET 4

Applying SVM & GridSearchCV on Train data to obtain the best C

```
In [113]: #code source: http://occam.olin.edu/sites/default/files/DataScienceMaterials/mach
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.linear_model import LogisticRegression
from sklearn.linear_model import SGDClassifier

# data = load_breast_cancer() #refer: http://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html

tuned_parameters = {'alpha': [10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2]}
# X_train, X_test, y_train, y_test = train_test_split(data.data, data.target, train_size=0.7, random_state=0)

#Using GridSearchCV
model = GridSearchCV(SGDClassifier(penalty = 'l2', loss='hinge', class_weight="balanced"),
                    tuned_parameters, cv=5, scoring='accuracy')
model.fit(X_tfidf_w2v_train, y_train)

print(model.best_estimator_)
print(model.score(X_tfidf_w2v_test, y_test))
```

```
SGDClassifier(alpha=100, average=False, class_weight='balanced',
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1500, n_iter_no_change=5, n_jobs=None, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
```

0.68

```

In [114]: # from sklearn.model_selection import GridSearchCV

# from sklearn.naive_bayes import MultinomialNB
# # nb = MultinomialNB(class_prior=[0.5,0.5])

# parameters = [0.00001, 0.0001,0.001, 0.01, 0.1,0.5,0.8, 1, 10, 100, 1000]
alpha=[*tuned_parameters.values()]
alpha=alpha[0]
log_alphas=[math.log10(num) for num in alpha]

# clf = GridSearchCV(nb, parameters, cv= 10, scoring='roc_auc',return_train_score=False)

# clf.fit(X_BOW_TRAIN, y_train)

#Selecting the best parameter
best_alpha1=model.best_params_
# print(best_alpha1)
train_auc= model.cv_results_['mean_train_score']
# print(train_auc)
train_auc_std= model.cv_results_['std_train_score']
# print(train_auc_std)
cv_auc = model.cv_results_['mean_test_score']
# print(cv_auc)
cv_auc_std= model.cv_results_['std_test_score']
# print(cv_auc_std)
print(log_alphas)

plt.figure(figsize=(20,10))

plt.plot(log_alphas, train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_alphas,train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.5)

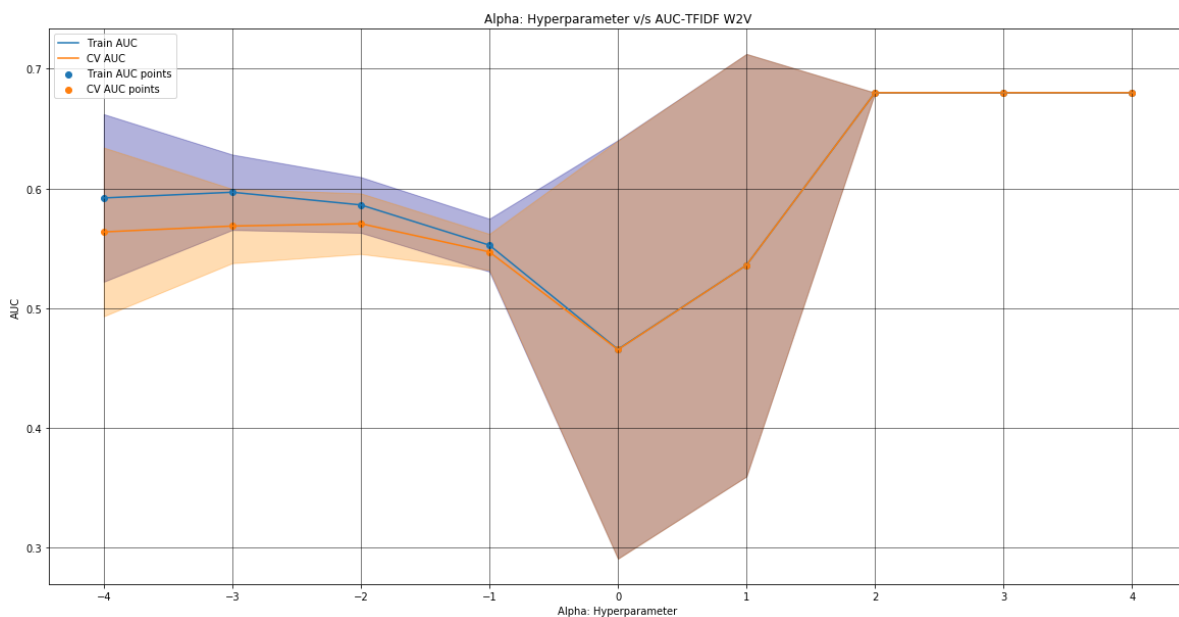
plt.plot(log_alphas, cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_alphas,cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.5)

plt.scatter(log_alphas, train_auc, label='Train AUC points')
plt.scatter(log_alphas, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("Alpha: Hyperparameter")
plt.ylabel("AUC")
plt.title("Alpha: Hyperparameter v/s AUC-TFIDF W2V")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()

```

```
[-4.0, -3.0, -2.0, -1.0, 0.0, 1.0, 2.0, 3.0, 4.0]
```



```
In [115]: print("The best alpha from the above graph is ",best_alpha1['alpha'])  
# best_alpha1['alpha']
```

The best alpha from the above graph is 100

2.8.4 Receiver Operating Characteristic- TFIDF W2V

```

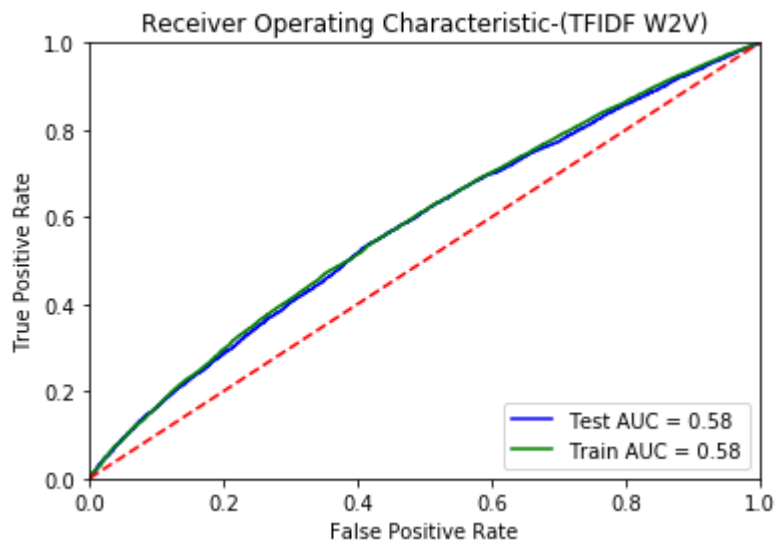
In [116]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html
from sklearn.metrics import roc_curve, auc
from sklearn.calibration import CalibratedClassifierCV
# train model on the best k
SVM=CalibratedClassifierCV(SGDClassifier(alpha=best_alpha1['alpha'],penalty = 'l2'))
SVM.fit(X_tfidf_w2v_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates
# not the predicted outputs

probs = SVM.predict_proba(X_tfidf_w2v_test)
# print(len(probs[:,1]))
probs1 = SVM.predict_proba(X_tfidf_w2v_train)
# print(len(probs1[:,1]))
preds = probs[:,1]
# print(preds)
preds1 = probs1[:,1]
# print(preds1)
fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
# print(fpr,tpr,threshold)
fpr1, tpr1, threshold = metrics.roc_curve(y_train, preds1)
# print(fpr1,tpr1,threshold)
roc_auc = metrics.auc(fpr, tpr)
# print(roc_auc)
roc_auc1 = metrics.auc(fpr1, tpr1)
# print(roc_auc1)

# https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve

import matplotlib.pyplot as plt
plt.title('Receiver Operating Characteristic-(TFIDF W2V)')
plt.plot(fpr, tpr, 'b', label = 'Test AUC = %0.2f' % roc_auc)
plt.plot(fpr1, tpr1, 'g', label = 'Train AUC = %0.2f' % roc_auc1)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'r--')
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()

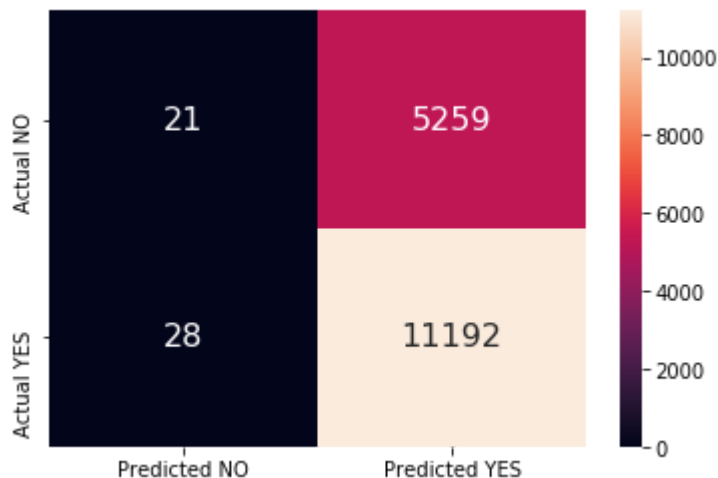
```



2.8.5 Confusion matrix- TFIDF W2V

```
In [117]: # https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#function to get heatmap confusion matrix
def get_confusion_matrix(clf,X_te,y_test):
    y_pred = clf.predict(X_te)
    df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), index=['Actual NO','Actual YES'], columns=['Predicted NO','Predicted YES'])
    sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')

# %%time
get_confusion_matrix(SVM,X_tfidf_w2v_test,y_test)
```



```
In [118]: #To make best use of the memory we are setting the variable names to 'None' and deleting them
X_tfidf_w2v_test=None
X_tfidf_w2v_train=None
y_tfidf_w2v_train=None
gc.collect()
```

Out[118]: 5591


```
In [48]: #At 2500 dimensions we are getting 80% efficiency.
#But we are taking 3000 dimensions since we can approx that it would increase the
svd = TruncatedSVD(n_components= 3000)
svd.fit(X_train_essay_Tfidf)
#Transforms:
#Train SVD
X_train_essay_Tfidf= svd.transform(X_train_essay_Tfidf )
#Test SVD
X_test_essay_Tfidf = svd.transform(X_test_essay_Tfidf )
```

```
In [50]: from scipy.sparse import hstack
X_NO_text_train = hstack((X_train_essay_Tfidf,X_train_digits_in_summary_norm,X_train_digits_in_summary_norm))
X_NO_text_train=X_NO_text_train.todense()
X_NO_text_train=np.array(X_NO_text_train)

# X_tfidf_w2v_cv = hstack((Xtfidf_w2v_vectors_cv,tfidf_w2v_vectors_Pro_title_cv,Xtfidf_w2v_vectors_Pro_title_cv))
# X_tfidf_w2v_cv=X_tfidf_w2v_cv.todense()
# X_tfidf_w2v_cv=np.array(X_tfidf_w2v_cv)

X_NO_text_test = hstack((X_test_essay_Tfidf,X_test_digits_in_summary_norm,X_test_digits_in_summary_norm))
X_NO_text_test=X_NO_text_test.todense()
X_NO_text_test=np.array(X_NO_text_test)
```

2.9.4 Applying SVM on No text features , SET 5

Applying SVM & GridSearchCV on Train data to obtain the best C

```
In [54]: #code source: http://occam.olin.edu/sites/default/files/DataScienceMaterials/mach
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.linear_model import LogisticRegression
from sklearn.linear_model import SGDClassifier

# data = load_breast_cancer() #refer: http://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html

tuned_parameters = {'alpha': [0.00001,0.0001,0.001,0.01,0.1,1,10,100,1000,10000]}
# X_train, X_test, y_train, y_test = train_test_split(data.data, data.target, train_size=0.7, random_state=0)

#Using GridSearchCV
model = GridSearchCV(SGDClassifier(penalty='l2',loss='hinge',class_weight="balanced"), tuned_parameters, cv=5, scoring='f1')
model.fit(X_NO_text_train, y_train)

print(model.best_estimator_)
print(model.score(X_NO_text_test, y_test))
```

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14
37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14
37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14
37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14
37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14
37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14
37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14
37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14

37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14

37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14

37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14

37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

C:\Users\Admin\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:14

37: UndefinedMetricWarning:

F-score is ill-defined and being set to 0.0 due to no predicted samples.

```
SGDClassifier(alpha=0.0001, average=False, class_weight='balanced',
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1500, n_iter_no_change=5, n_jobs=None, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
```

0.7101860920666014

```

In [63]: # from sklearn.model_selection import GridSearchCV

# from sklearn.naive_bayes import MultinomialNB
# # nb = MultinomialNB(class_prior=[0.5,0.5])

# parameters = [0.00001, 0.0001,0.001, 0.01, 0.1,0.5,0.8, 1, 10, 100, 1000]
alpha=[*tuned_parameters.values()]
alpha=alpha[0]
log_alphas=[math.log10(num) for num in alpha]

# clf = GridSearchCV(nb, parameters, cv= 10, scoring='roc_auc',return_train_score=False)

# clf.fit(X_BOW_TRAIN, y_train)

#Selecting the best parameter
best_alpha1=model.best_params_
# print(best_alpha1)
train_auc= model.cv_results_['mean_train_score']
# print(train_auc)
train_auc_std= model.cv_results_['std_train_score']
# print(train_auc_std)
cv_auc = model.cv_results_['mean_test_score']
# print(cv_auc)
cv_auc_std= model.cv_results_['std_test_score']
# print(cv_auc_std)
print(log_alphas)

plt.figure(figsize=(20,10))

plt.plot(log_alphas, train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_alphas,train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.5)

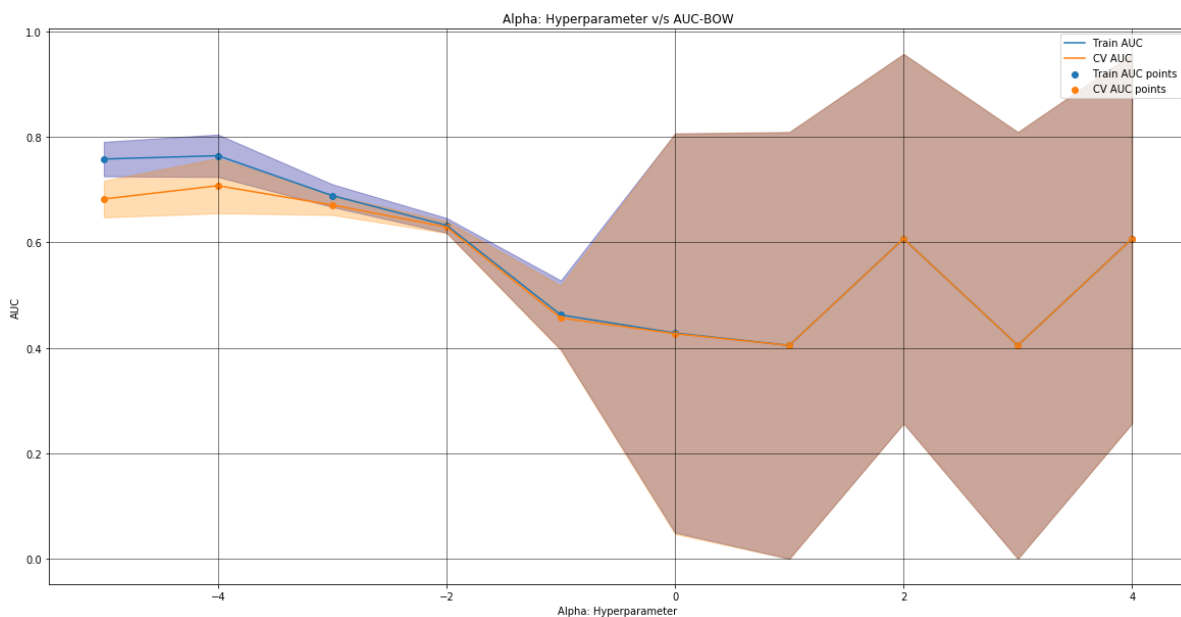
plt.plot(log_alphas, cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_alphas,cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.5)

plt.scatter(log_alphas, train_auc, label='Train AUC points')
plt.scatter(log_alphas, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("Alpha: Hyperparameter")
plt.ylabel("AUC")
plt.title("Alpha: Hyperparameter v/s AUC-BOW")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()

```

```
[-5.0, -4.0, -3.0, -2.0, -1.0, 0.0, 1.0, 2.0, 3.0, 4.0]
```



```
In [56]: print("The best alpha from the above graph is ",best_alpha1)
# best_alpha1['alpha']
```

The best alpha from the above graph is {'alpha': 0.0001}

2.9.5 Receiver Operating Characteristic- (No text features)

```

In [60]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html
from sklearn.metrics import roc_curve, auc

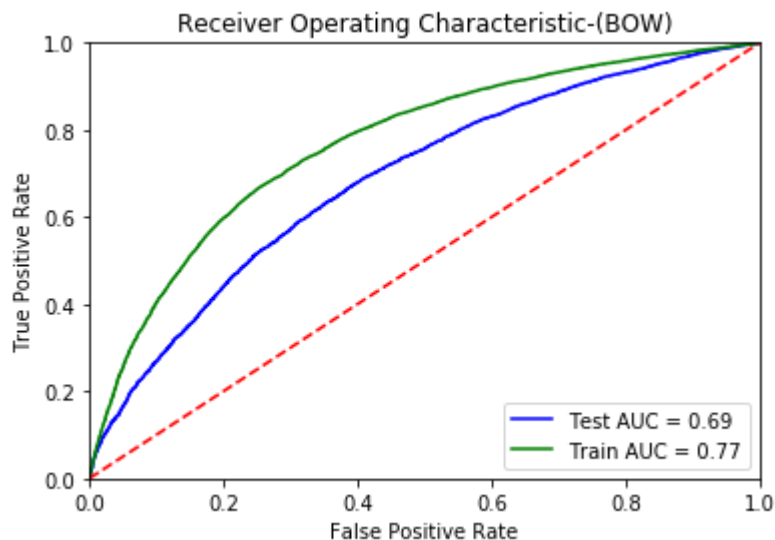
# train model on the best k
SVM=CalibratedClassifierCV(SGDClassifier(alpha=best_alpha1['alpha'],penalty = 'l2'))
SVM.fit(X_NO_text_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates
# not the predicted outputs

probs = SVM.predict_proba(X_NO_text_test)
# print(len(probs[:,1]))
probs1 = SVM.predict_proba(X_NO_text_train)
# print(len(probs1[:,1]))
preds = probs[:,1]
# print(preds)
preds1 = probs1[:,1]
# print(preds1)
fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
# print(fpr,tpr,threshold)
fpr1, tpr1, threshold = metrics.roc_curve(y_train, preds1)
# print(fpr1,tpr1,threshold)
roc_auc = metrics.auc(fpr, tpr)
# print(roc_auc)
roc_auc1 = metrics.auc(fpr1, tpr1)
# print(roc_auc1)

# https://www.programcreek.com/python/example/81207/sklearn.metrics.roc_curve

import matplotlib.pyplot as plt
plt.title('Receiver Operating Characteristic-(BOW)')
plt.plot(fpr, tpr, 'b', label = 'Test AUC = %0.2f' % roc_auc)
plt.plot(fpr1, tpr1, 'g', label = 'Train AUC = %0.2f' % roc_auc1)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'r--')
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()

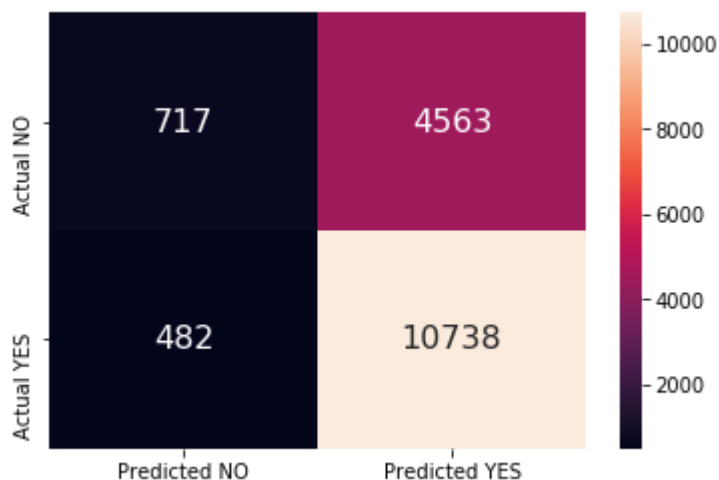
```



2.9.6 Confusion matrix- No text features

```
In [61]: # https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#function to get heatmap confusion matrix
def get_confusion_matrix(clf,X_te,y_test):
    y_pred = clf.predict(X_te)
    df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), index=['Actual NO','Actual YES'], columns=['Predicted NO','Predicted YES'])
    sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')

# %%time
get_confusion_matrix(SVM,X_NO_text_test,y_test)
```



```
In [ ]: #To make best use of the memory we are setting the variable names to 'None' and deleting them
X_NO_text_test=None
X_NO_text_train=None
# y_tfidf_w2v_train=None
gc.collect()
```

2.10 Pretty table summary


```
In [119]: #http://zetcode.com/python/prettytable/
from prettytable import PrettyTable

x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Hyper Parameter-Alpha", "Test-AUC"]
x.add_row(["BOW", "SVM", '0.0001', 0.64])
x.add_row(["Tfidf", "SVM", '0.0001', 0.61])
x.add_row(["avg_w2v", "SVM", '0.0001', 0.57])
x.add_row(["tfidf_w2v", "SVM", '0.0001', 0.58])
x.add_row(["No_Text", "SVM", '0.0001', 0.69])
print(x)
```

Vectorizer	Model	Hyper Parameter-Alpha	Test-AUC
BOW	SVM	0.0001	0.64
Tfidf	SVM	0.0001	0.61
avg_w2v	SVM	0.0001	0.57
tfidf_w2v	SVM	0.0001	0.58
No_Text	SVM	0.0001	0.69

In []: