

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>	Title of the project. Examples: Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	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
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: Applied Learning Care & Hunger Health & Sports History & Civics Literacy & Language Math & Science Music & The Arts Special Needs Warmth Examples: Music & The Arts Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. Examples: Literacy Literature & Writing, Social Sciences
<code>project_resource_summary</code>	An explanation of the resources needed for the project. Example: My students need hands on literacy materials to manage sensory needs!
<code>project_essay_1</code>	First application essay*
<code>project_essay_2</code>	Second application essay*
<code>project_essay_3</code>	Third application essay*

Feature	Description
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_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> nan Dr. Mr. Mrs. Ms. Teacher.
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the 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
id	A <code>project_id</code> value from the <code>train.csv</code> file. Example: p036502
description	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
quantity	Quantity of the resource required. Example: 3
price	Price of the resource required. Example: 9.95

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

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

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

Notes on the Essay Data

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

- __project_essay_1__: "Introduce us to your classroom"
- __project_essay_2__: "Tell us more about your students"
- __project_essay_3__: "Describe how your students will use the materials you're requesting"
- __project_essay_3__: "Close by sharing why your project will make a difference"

Starting on May 17, 2016, the number of essays was reduced from 4 to 2, and the prompts for the first 2 essays were changed to the following:

- __project_essay_1__: "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- __project_essay_2__: "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

For all projects with `project_submitted_datetime` of 2016-05-17 and later, the values of `project_essay_3` and `project_essay_4` will be NaN.

In [1]:

```
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
```

```

import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

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

```

C:\Users\ADMIN\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
 warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")

1.1 Reading Data

In [2]:

```

project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')

```

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 (109248, 17)
-----
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
 'project_submitted_datetime' 'project_grade_category'
 'project_subject_categories' 'project_subject_subcategories'
 'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
 'project_essay_4' 'project_resource_summary'
 'teacher_number_of_previously_posted_projects' 'project_is_approved']

```

In [4]:

```

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

```

```

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

```

Out[4]:

id	description	quantity	price
10652	Lakeshore Double Space Mobile Drives		

id	description	quantity	price
0 p233245	Mobile Drying Rack	1	149.00
1 p069063	Bouncy Bands for Desks (Blue support pipes)	3	14.95

1.2 preprocessing of project_subject_categories

In [5]:

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

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

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

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

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

1.3 preprocessing of project_subject_subcategories

In [6]:

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

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

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

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

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

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

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

1.3 Text preprocessing

In [7]:

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

In [8]:

```
project_data.head(2)
```

Out[8]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_cat
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades P
1	140945 p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Grade

In [9]:

```
price_data = resource_data.groupby('id').agg({'price': 'sum', 'quantity': 'sum'}).reset_index()
price_data.head(2)
project_data = pd.merge(project_data, price_data, on='id', how='left')
```

Assignment 5: Logistic Regression

1. [Task-1] Logistic Regression(either SGDClassifier with log loss, or LogisticRegression) on these feature sets

- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_eassay ('BOW with bi-grams` with `min_df=10` and `max_features=5000`)
- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay ('TFIDF with bi-grams` with `min_df=10` and `max_features=5000`)
- **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)
- **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

2. Hyper paramter tuning (find best hyper parameters corresponding the algorithm that you choose)

- Find the best hyper parameter which will give the maximum [AUC](#) value
- Find the best hyper paramter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.
- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.
- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).

4. [\[Task-2\] Apply Logistic Regression on the below feature set Set 5 by finding the best hyper parameter as suggested in step 2 and step 3.](#)

5. [Consider these set of features Set 5 :](#)

- [school_state](#) : categorical data
- [clean_categories](#) : categorical data
- [clean_subcategories](#) : categorical data
- [project_grade_category](#) : categorical data
- [teacher_prefix](#) : categorical data
- [quantity](#) : numerical data
- [teacher_number_of_previously_posted_projects](#) : numerical data
- [price](#) : numerical data
- [sentiment score's of each of the essay](#) : numerical data
- [number of words in the title](#) : numerical data
- [number of words in the combine essays](#) : numerical data

[And apply the Logistic regression on these features by finding the best hyper paramter as suggested in step 2 and step 3](#)

6. [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 \[prettytable library link\]\(#\)](#)

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-leakage, 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](#).

2. Logistic Regression

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

In [10]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label

sample_data=project_data.sample(50000)
sample_data.shape
from sklearn.model_selection import train_test_split
data=sample_data
data.head()
y=data.project_is_approved
x=data.drop('project_is_approved',axis=1)
x_train,x_test,y_train,y_test=train_test_split(x,y,train_size=0.8,test_size=0.2,stratify=y)
x_train,x_cv,y_train,y_cv=train_test_split(x_train,y_train,train_size=0.8,test_size=0.2,stratify=y_train)
```

```

print("shape of train data ")
print(x_train.shape)
print(y_train.shape)
print("shape of test data ")
print(x_test.shape)
print(y_test.shape)
print("shape of crossvalidation data ")
print(x_cv.shape)
print(y_cv.shape)

```

```

shape of train data
(32000, 19)
(32000,)
shape of test data
(10000, 19)
(10000,)
shape of crossvalidation data
(8000, 19)
(8000,)

```

2.2 Make Data Model Ready: encoding numerical, categorical features

In [12]:

```

# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
def veccat(x):
    from sklearn.feature_extraction.text import CountVectorizer
    from sklearn.preprocessing import Normalizer
    from scipy.sparse import hstack
    vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
    vectorizer.fit(x['clean_categories'].values)
    print(vectorizer.get_feature_names())

    categories_one_hot = vectorizer.transform(x['clean_categories'].values)
    print("Shape of matrix after one hot encodig ",categories_one_hot.shape)

    vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
    vectorizer.fit(x['clean_subcategories'].values)
    print(vectorizer.get_feature_names())

    subcategories_one_hot = vectorizer.transform(x['clean_subcategories'].values)
    print("Shape of matrix after one hot encodig ",subcategories_one_hot.shape)

    vectorizer = CountVectorizer(lowercase=False, binary=True)
    vectorizer.fit(x['school_state'].values)
    print(vectorizer.get_feature_names())

    state_one_hot = vectorizer.transform(x['school_state'].values)
    print("Shape of matrix after one hot encodig ",state_one_hot.shape)

    x = x.replace(np.nan, '', regex=True)
    vectorizer = CountVectorizer( vocabulary=['Mrs.', 'Ms.', 'Mr.', 'Teacher', 'Dr.'], lowercase=False)
    vectorizer.fit(x['teacher_prefix'].values)
    print(vectorizer.get_feature_names())

```

```

prefix_one_hot = vectorizer.transform(x['teacher_prefix'].values)
print("Shape of matrix after one hot encodig ",prefix_one_hot.shape)

vocab=[]
for i in x['project_grade_category'].values:
    vocab.append(i)
v_set=set(vocab)
vocab=list(v_set)
vectorizer = CountVectorizer(vocabulary=vocab,lowercase=False)
vectorizer.fit(x['project_grade_category'].values)

grade_one_hot = vectorizer.transform(x['project_grade_category'].values)
print("Shape of matrix after one hot encodig ",grade_one_hot.shape)

price_scalar = Normalizer(copy=False,norm='l2')
price_scalar.fit(x['price'].values.reshape(1,-1)) # finding the mean and standard deviation of
this data

# Now standardize the data with above maen and variance.
price_standardized = price_scalar.transform(x['price'].values.reshape(1, -1))
price_standardized=np.transpose(price_standardized)

projects_scalar = Normalizer(copy=False,norm='l2')
projects_scalar.fit(x['teacher_number_of_previously_posted_projects'].values.reshape(1,-1)) # f
inding the mean and standard deviation of this data

# Now standardize the data with above maen and variance.
projects_standardized =
projects_scalar.transform(x['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
projects_standardized =np.transpose(projects_standardized)

qty_scalar= Normalizer(copy=False,norm='l2')
qty_scalar.fit(x['quantity'].values.reshape(1,-1)) # finding the mean and standard deviation of
this data

# Now standardize the data with above maen and variance.
qty_standardized = qty_scalar.transform(x['quantity'].values.reshape(1, -1))
qty_standardized=np.transpose(qty_standardized)

X1 = hstack((categories_one_hot, subcategories_one_hot,state_one_hot,prefix_one_hot,grade_one_h
ot, price_standardized,projects_standardized,qty_standardized))
print(X1.shape)
return(X1)

```

In [13]:

```

x=veccat(x_train)
t=veccat(x_test)
cv=veccat(x_cv)

```

```

['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encodig (32000, 9)
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular',
'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger',
'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']
Shape of matrix after one hot encodig (32000, 30)
['AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DC', 'DE', 'FL', 'GA', 'HI', 'IA', 'ID', 'IL', 'IN', 'K
S', 'KY', 'LA', 'MA', 'MD', 'ME', 'MI', 'MN', 'MO', 'MS', 'MT', 'NC', 'ND', 'NE', 'NH', 'NJ', 'NM',
'NV', 'NY', 'OH', 'OK', 'OR', 'PA', 'RI', 'SC', 'SD', 'TN', 'TX', 'UT', 'VA', 'VT', 'WA', 'WI', 'WV
', 'WY']
Shape of matrix after one hot encodig (32000, 51)
['Mrs.', 'Ms.', 'Mr.', 'Teacher', 'Dr.']
Shape of matrix after one hot encodig (32000, 5)

```



```

Shape of matrix after one hot encoding (32000, 4)
(32000, 102)
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding (10000, 9)
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular',
'Citizens_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger',
'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']
Shape of matrix after one hot encoding (10000, 30)
['AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DC', 'DE', 'FL', 'GA', 'HI', 'IA', 'ID', 'IL', 'IN', 'KS',
'S', 'KY', 'LA', 'MA', 'MD', 'ME', 'MI', 'MN', 'MO', 'MS', 'MT', 'NC', 'ND', 'NE', 'NH', 'NJ', 'NM',
'NV', 'NY', 'OH', 'OK', 'OR', 'PA', 'RI', 'SC', 'SD', 'TN', 'TX', 'UT', 'VA', 'VT', 'WA', 'WI', 'WV',
'WY']
Shape of matrix after one hot encoding (10000, 51)
['Mrs.', 'Ms.', 'Mr.', 'Teacher', 'Dr.']
Shape of matrix after one hot encoding (10000, 5)
Shape of matrix after one hot encoding (10000, 4)
(10000, 102)
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding (8000, 9)
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular',
'Citizens_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger',
'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']
Shape of matrix after one hot encoding (8000, 30)
['AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DC', 'DE', 'FL', 'GA', 'HI', 'IA', 'ID', 'IL', 'IN', 'KS',
'S', 'KY', 'LA', 'MA', 'MD', 'ME', 'MI', 'MN', 'MO', 'MS', 'MT', 'NC', 'ND', 'NE', 'NH', 'NJ', 'NM',
'NV', 'NY', 'OH', 'OK', 'OR', 'PA', 'RI', 'SC', 'SD', 'TN', 'TX', 'UT', 'VA', 'VT', 'WA', 'WI', 'WV',
'WY']
Shape of matrix after one hot encoding (8000, 51)
['Mrs.', 'Ms.', 'Mr.', 'Teacher', 'Dr.']
Shape of matrix after one hot encoding (8000, 5)
Shape of matrix after one hot encoding (8000, 4)
(8000, 102)

```

In [14]:

```

def features(x):
    from sklearn.feature_extraction.text import CountVectorizer
    from sklearn.preprocessing import MinMaxScaler
    from scipy.sparse import hstack
    vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
    vectorizer.fit(x['clean_categories'].values)

    feature_list.extend(vectorizer.get_feature_names())
    vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False,
    binary=True)
    vectorizer.fit(x['clean_subcategories'].values)

    feature_list.extend(vectorizer.get_feature_names())
    vectorizer = CountVectorizer(lowercase=False, binary=True)
    vectorizer.fit(x['school_state'].values)

    feature_list.extend(vectorizer.get_feature_names())
    x = x.replace(np.nan, '', regex=True)
    vectorizer = CountVectorizer(vocabulary=['Mrs.', 'Ms.', 'Mr.', 'Teacher', 'Dr.'], lowercase=False)
    vectorizer.fit(x['teacher_prefix'].values)

    feature_list.extend(vectorizer.get_feature_names())
    vocab=[]
    for i in x['project_grade_category'].values:
        vocab.append(i)
    v_set=set(vocab)
    vocab=list(v_set)
    vectorizer = CountVectorizer(vocabulary=vocab, lowercase=False)
    vectorizer.fit(x['project_grade_category'].values)
    feature_list.extend(vectorizer.get_feature_names())
    feature_list.append("price")
    feature_list.append("teacher_number_of_previously_posted_projects")

```

```
feature_list.append("quantity")
```

In [15]:

```
feature_list=[]
features(x_train)
print(len(feature_list))
```

102

2.3 Make Data Model Ready: encoding eassay, and project_title

Vectorizing essay and project_title

In [16]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly
```

```
# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label
```

```
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 [17]:

```
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',
'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them',
'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll",
'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having',
'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', '
while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during',
'before', 'after',\
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under'
, 'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'e
ach', 'few', 'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll'
, 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "dc
esn't", 'hadn',\
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn'
```

```
"mightn't", 'mustn', \
    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn',
    "wasn't", 'weren', "weren't", \
    'won', "won't", 'wouldn', "wouldn't"]
```

In [18]:

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

In [19]:

```
train_essay=[]
test_essay=[]
cv_essay=[]
train_title=[]
test_title=[]
cv_title=[]
train_essay=preprocessing(x_train['essay'])
test_essay=preprocessing(x_test['essay'])
cv_essay=preprocessing(x_cv['essay'])

train_title=preprocessing(x_train['project_title'])
test_title=preprocessing(x_test['project_title'])
cv_title=preprocessing(x_cv['project title'])
```

[illegible]

In [26]:

```
feature_bow = feature_list.copy()
feature_tfidf = feature_list.copy()
```

Encoding using BOW

In [27]:

```
from scipy.sparse import hstack
from sklearn.preprocessing import StandardScaler
def bow_text(train_essay, test_essay, cv_essay, train_title, test_title, cv_title, x1, t1, cv1):
    from scipy.sparse import hstack

    vectorizer = CountVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)

    vectorizer.fit(train_essay)

    text_bow=vectorizer.transform(train_essay)

    text_bow1 = vectorizer.transform(test_essay)

    text_bow2 = vectorizer.transform(cv_essay)
    feature_bow.extend(vectorizer.get_feature_names())
    vectorizer= CountVectorizer()
    vectorizer.fit(train_title)
    title_bow=vectorizer.transform(train_title)

    title_bow1 = vectorizer.transform(test_title)

    title_bow2 = vectorizer.transform(cv_title)

    x1 = hstack((x1, text_bow, title_bow)).tocsr()
    t1= hstack((t1, text_bow1, title_bow1)).tocsr()
    cv1 = hstack((cv1, text_bow2, title_bow2)).tocsr()
    feature_bow.extend(vectorizer.get_feature_names())

    return x1, t1, cv1
```

In [28]:

```
# Data matrix using BOW
x1, t1, cv1=bow_text(train_essay, test_essay, cv_essay, train_title, test_title, cv_title, x, t, cv)
print(x1.shape)
print(t1.shape)
print(cv1.shape)
# Feature list for BOW encoding
print(len(feature_bow))
```

```
(32000, 14572)
(10000, 14572)
(8000, 14572)
14572
```

Encoding using TFIDF

In [29]:

```
def tfidf_text(train_essay, test_essay, cv_essay, train_title, test_title, cv_title, x1, t1, cv1):
    from sklearn.feature_extraction.text import TfidfVectorizer
    vectorizer = TfidfVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)
    vectorizer.fit(train_essay)
    feature_tfidf.extend(vectorizer.get_feature_names())
    text_tfidf=vectorizer.transform(train_essay)
    text_tfidf1 = vectorizer.transform(test_essay)
    text_tfidf2 = vectorizer.transform(cv_essay)
    vectorizer = TfidfVectorizer()
    vectorizer.fit(train_title)
    feature_tfidf.extend(vectorizer.get_feature_names())
    title_tfidf=vectorizer.transform(train_title)
    title_tfidf1 = vectorizer.transform(test_title)
    title_tfidf2 = vectorizer.transform(cv_title)
    x1 = hstack((x1, text_tfidf, title_tfidf))
    t1= hstack((t1, text_tfidf1, title_tfidf1))
    cv1 = hstack((cv1, text_tfidf2, title_tfidf2))
```

```
return x1,t1,cv1
```

In [30]:

```
x2,t2,cv2=tfidf_text(train_essay, test_essay, cv_essay,train_title,test_title,cv_title,x,t,cv)
print(x2.shape)
print(t2.shape)
print(cv2.shape)
```

```
(32000, 14572)
(10000, 14572)
(8000, 14572)
```

In [31]:

```
print(len(feature_tfidf))
```

```
14572
```

Encoding using AVG W2V

In [48]:

```
def avgword2vec(text):

    with open('glove_vectors', 'rb') as f:
        model = pickle.load(f)
        glove_words = set(model.keys())
    avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in tqdm(text):# 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.append(vector)
    return avg_w2v_vectors
```

In [49]:

```
def w2v(train_essay, test_essay, cv_essay,train_title,test_title,cv_title,x,t,cv):

    text_w2v=avgword2vec(train_essay)
    text_w2v1 = avgword2vec(test_essay)
    text_w2v2 = avgword2vec(cv_essay)

    title_w2v=avgword2vec(train_title)
    title_w2v1 = avgword2vec(test_title)
    title_w2v2 = avgword2vec(cv_title)
    x1 = hstack((x,text_w2v,title_w2v)).tocsr()
    t1= hstack((t,text_w2v1,title_w2v1)).tocsr()
    cv1 = hstack((cv,text_w2v2,title_w2v2)).tocsr()
    return x1,t1,cv1
```

In [50]:

```
x3,t3,cv3=w2v(train_essay, test_essay, cv_essay,train_title,test_title,cv_title,x,t,cv)
print(x3.shape)
print(t3.shape)
print(cv3.shape)
```

```
100%|████████████████████████████████████████| 32000/32000 [00:11<00:00, 2792.41it/s]
100%|████████████████████████████████████████| 10000/10000 [00:03<00:00, 2748.61it/s]
100%|████████████████████████████████████████| 8000/8000 [00:02<00:00, 2733.02it/s]
100%|████████████████████████████████████████| 32000/32000 [00:00<00:00, 52456.01it/s]
100%|████████████████████████████████████████| 10000/10000 [00:00<00:00, 51810.51it/s]
```

```
100%|███████████| 8000/8000 [00:00<00:00, 42325.66it/s]
```

(32000, 702)
(10000, 702)
(8000, 702)

Encoding using TFIDFW2V

In [51]:

```
def tfidf_w2v(text, dictionary, tfidf_words):

    with open('glove_vectors', 'rb') as f:
        model = pickle.load(f)
        glove_words = set(model.keys())

    tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in tqdm(text): # for each review/sentence
        vector = np.zeros(300) # as word vectors are of zero length
        tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
        for word in sentence.split(): # for each word in a review/sentence
            if (word in glove_words) and (word in tfidf_words):
                vec = model[word] # getting the vector for each word
                # here we are multiplying idf value(dictionary[word]) and the tf
                value((sentence.count(word)/len(sentence.split())))
                tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting
                the tfidf value for each word
                vector += (vec * tf_idf) # calculating tfidf weighted w2v
                tf_idf_weight += tf_idf
            if tf_idf_weight != 0:
                vector /= tf_idf_weight
        tfidf_w2v_vectors.append(vector)
    return tfidf_w2v_vectors
```

In [52]:

```
def tfidf2v(train_essay, test_essay, cv_essay, train_title, test_title, cv_title, x, t, cv):

    tfidf_model = TfidfVectorizer()
    tfidf_model.fit(train_essay)
    # we are converting a dictionary with word as a key, and the idf as a value
    dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
    tfidf_words = set(tfidf_model.get_feature_names())

    text_tfw2v=tfidf2v(train_essay,dictionary,tfidf_words)
    text_tfw2v1=tfidf2v(test_essay,dictionary,tfidf_words)
    text_tfw2v2=tfidf2v(cv_essay,dictionary,tfidf_words)

    tfidf_model.fit(train_title)
    # we are converting a dictionary with word as a key, and the idf as a value
    dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
    tfidf_words = set(tfidf_model.get_feature_names())
    title_tfw2v=tfidf2v(train_title,dictionary,tfidf_words)
    title_tfw2v1=tfidf2v(test_title,dictionary,tfidf_words)
    title_tfw2v2=tfidf2v(cv_title,dictionary,tfidf_words)
    x1 = hstack((x,text_tfw2v,title_tfw2v )).tocsr()
    t1= hstack((t,text_tfw2v1,title_tfw2v1 )).tocsr()
    cv1 = hstack((cv,text_tfw2v2,title_tfw2v2 )).tocsr()
    return x1,t1,cv1
```

In [53]:

```
x4,t4,cv4=tfiw2v(train_essay, test_essay, cv_essay,train_title,test_title,cv_title,x,t,cv)
print(x4.shape)
print(t4.shape)
print(cv4.shape)
```

```
100%|██████████| 32000/32000 [01:28<00:00, 360.47it/s]
100%|██████████| 10000/10000 [00:29<00:00, 333.94it/s]
100%|██████████| 8000/8000 [00:22<00:00, 351.98it/s]
100%|██████████| 32000/32000 [00:01<00:00, 21475.28it/s]
```

```
100%|████████████████████| 10000/10000 [00:00<00:00, 21320.74it/s]
100%|████████████████████| 8000/8000 [00:00<00:00, 22159.40it/s]
```

```
(32000, 702)
(10000, 702)
(8000, 702)
```

2.4 Applying Logistic Regression on different kind of featurization as mentioned in the instructions

Apply Logistic Regression 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

In [90]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label

# Function for finding the optimal lambda
def alpha_from_roc(x_train,y_train,x_cv,y_cv):
    from sklearn.linear_model import SGDClassifier
    from sklearn.model_selection import cross_val_score
    from sklearn.metrics import roc_curve
    from sklearn.metrics import roc_auc_score

    a_range=[0.0001,0.001,0.01,0.1,1,10,100,1000,10000]
    scores_auc=[]

    for a in a_range:
        clf=SGDClassifier(loss='log',alpha=a,class_weight='balanced')
        clf.fit(x_train,y_train)
        score_roc = clf.decision_function(x_cv)
        fpr, tpr, thresholds = metrics.roc_curve(y_cv, score_roc)
        roc_auc = auc(fpr, tpr)
        scores_auc.append(roc_auc)
    # Chosen optimal k with the highest AUC score
    optimal_alpha = a_range[scores_auc.index(max(scores_auc))]
```

```
clf =SGDClassifier(loss='log',alpha=optimal_alpha,class_weight='balanced')
clf.fit(x_train,y_train)

score_roc = clf.decision_function(x_cv)
fpr, tpr, thresholds = metrics.roc_curve(y_cv, score_roc)

score_roc = clf.decision_function(x_train)
fpr1, tpr1, thresholds = metrics.roc_curve(y_train, score_roc)

# plot no skill
plt.plot([0, 1], [0, 1], linestyle='--')
# plot the roc curve for the model
plt.plot(fpr, tpr, marker='.',label="CrossValidated or test ROC")
plt.plot(fpr1, tpr1, marker='.', label='Train ROC')
plt.legend()
# show the plot
plt.title("Receiver Operating Characteristics")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.show()
print("Optimal lambda ",optimal_alpha)
print("AUC: ", max(scores_auc))
plt.plot(np.log10(a_range), scores_auc, label='CV or TEST AUC')
plt.scatter(np.log10(a_range), scores_auc, label='CV or TEST AUC points')
plt.xlabel("LOG(lambda): hyperparameter")
```

```
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.legend()

plt.show()
```

In [91]:

```
# K- Fold cross validation
def kfold_cv(x_train,y_train):
    from sklearn.model_selection import cross_val_score
    from sklearn.linear_model import SGDClassifier
    a_range=[0.0001,0.001,0.01,0.1,1,10,100,1000,10000]
    scores_cv=[]
    for a in a_range:
        clf=SGDClassifier(loss='log',alpha=a,class_weight='balanced')
        scores_cv.append(np.mean(cross_val_score(clf, x_train, y_train, cv=3, n_jobs=1)))
    for i in range(len(scores_cv)):
        if scores_cv[i]==max(scores_cv):
            print("Optimal lambda :",a_range[i])
            print("CV SCORE ",max(scores_cv))
```

In [92]:

```
# Gridsearchcv function

def grid_search(x_train,y_train,x_cv,y_cv):
    from sklearn.linear_model import SGDClassifier
    from sklearn.model_selection import GridSearchCV
    from sklearn.metrics import roc_curve
    from sklearn.metrics import roc_auc_score
    from sklearn import preprocessing

    a_range=[0.0001,0.001,0.01,0.1,1,10,100,1000,10000]

    parameters = dict(alpha=a_range)
    # Fitted the model on train data
    clf1 = GridSearchCV(SGDClassifier(loss='log',class_weight='balanced'), parameters, cv=3,verbose
=15, scoring='roc_auc')
    clf1.fit(x_train, y_train)
    train_auc= clf1.cv_results_['mean_train_score']
    cv_auc = clf1.cv_results_['mean_test_score']
    max_score= clf1.best_score_

    opt_par=clf1.best_params_
    print("Optimal lambda ",opt_par)
    print("AUC: ", max_score)
    plt.plot(parameters['alpha'], cv_auc, label='CV AUC')
    plt.plot(parameters['alpha'], train_auc, label='TRAIN AUC')
    plt.xlabel("Alpha: hyperparameter")
    plt.ylabel("AUC")
    plt.title("ERROR PLOTS")
    plt.legend()
    plt.show()
    # Found out the score for crossvalidated data
    print("Accuracy on crossvalidated data: " , clf1.score(x_cv,y_cv))
```

In [75]:

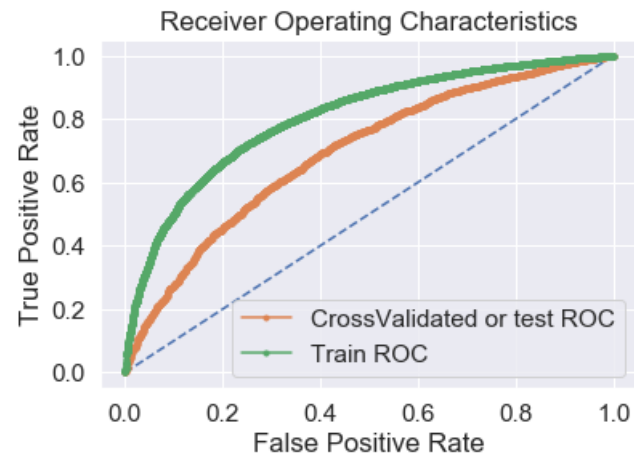
```
# Predictions on test data
def predict(x_train,y_train,x_test,y_test,a):
    from sklearn.linear_model import SGDClassifier
    clf=SGDClassifier(loss='log',alpha=a,class_weight='balanced')
    clf.fit(x_train,y_train)
    score_roc = clf.decision_function(x_test)
    fpr, tpr, thresholds = metrics.roc_curve(y_test, score_roc)
    t = thresholds[np.argmax(tpr*(1-fpr))]
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    print("train or test AUC =",str(roc_auc(fpr, tpr)))
    predictions=[]
    predictions=clf.predict(x_test)
    return predictions
```


Applying Logistic regression on BOW

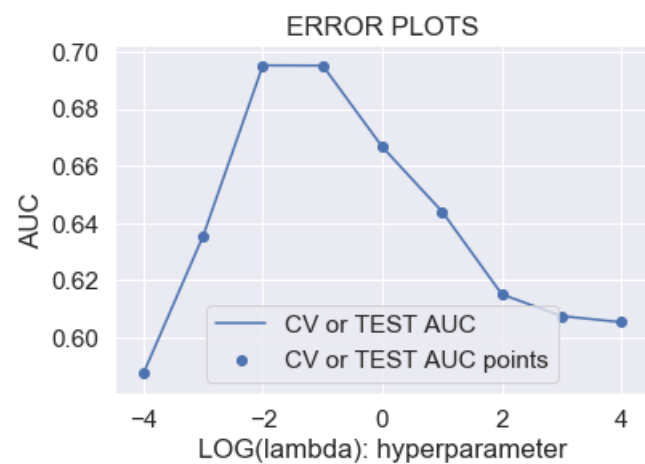
ROC SCORE on CV data

In [93]:

```
alpha_from_roc(x1,y_train,cv1,y_cv)
```



Optimal lambda 0.01
AUC: 0.6951340449761201



K FOLD Cross Validation

In [94]:

```
kfold_cv(x1,y_train)
```

Optimal lambda : 1000
CV SCORE 0.849093751658438

GridsearchCV

In [95]:

```
grid_search(x1,y_train,cv1,y_cv)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```
[CV] alpha=0.0001 .....  
[CV] ..... alpha=0.0001, score=0.6184119623337907, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.1s remaining: 0.0s
```

```
[CV] alpha=0.0001 .....  
[CV] ..... alpha=0.0001, score=0.601254031575042, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.3s remaining: 0.0s
```

```
[CV] alpha=0.0001 .....  
[CV] ..... alpha=0.0001, score=0.6254606945328447, total= 0.0s
```

```
[Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 0.5s remaining: 0.0s
```

```
[CV] alpha=0.001 .....  
[CV] ..... alpha=0.001, score=0.6410854786490254, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 4 out of 4 | elapsed: 0.7s remaining: 0.0s
```

```
[CV] alpha=0.001 .....  
[CV] ..... alpha=0.001, score=0.6266763911377013, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 0.8s remaining: 0.0s
```

```
[CV] alpha=0.001 .....  
[CV] ..... alpha=0.001, score=0.6392984027610122, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 6 out of 6 | elapsed: 1.0s remaining: 0.0s
```

```
[CV] alpha=0.01 .....  
[CV] ..... alpha=0.01, score=0.68076214341606, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 7 out of 7 | elapsed: 1.2s remaining: 0.0s
```

```
[CV] alpha=0.01 .....  
[CV] ..... alpha=0.01, score=0.667654269680567, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 8 out of 8 | elapsed: 1.4s remaining: 0.0s
```

```
[CV] alpha=0.01 .....  
[CV] ..... alpha=0.01, score=0.6871210603998859, total= 0.0s
```

```
[Parallel(n_jobs=1)]: Done 9 out of 9 | elapsed: 1.6s remaining: 0.0s
```

```
[CV] alpha=0.1 .....  
[CV] ..... alpha=0.1, score=0.685159551961113, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 10 out of 10 | elapsed: 1.8s remaining: 0.0s
```

```
[CV] alpha=0.1 .....  
[CV] ..... alpha=0.1, score=0.6690095921139889, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 11 out of 11 | elapsed: 2.0s remaining: 0.0s
```

```
[CV] alpha=0.1 .....  
[CV] ..... alpha=0.1, score=0.6953060185896451, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 12 out of 12 | elapsed: 2.2s remaining: 0.0s
```

```
[CV] alpha=1 .....  
[CV] ..... alpha=1, score=0.659418232491666, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 13 out of 13 | elapsed: 2.4s remaining: 0.0s
```

```
[CV] alpha=1 .....  
[CV] ..... alpha=1, score=0.6411369813129681, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 14 out of 14 | elapsed: 2.6s remaining: 0.0s
```

```
[CV] alpha=1 .....  
[CV] ..... alpha=1, score=0.6659111422645567, total= 0.1s  
[CV] alpha=10 .....  
[CV] ..... alpha=10, score=0.6414658165641072, total= 0.1s  
[CV] alpha=10 .....  
[CV] ..... alpha=10, score=0.627608719654747, total= 0.1s  
[CV] alpha=10 .....  
[CV] ..... alpha=10, score=0.6471653946184404, total= 0.1s  
[CV] alpha=100 .....  
[CV] ..... alpha=100, score=0.6200132082730698, total= 0.1s  
[CV] alpha=100 .....  
[CV] ..... alpha=100, score=0.6087802783887004, total= 0.0s  
[CV] alpha=100 .....  
[CV] ..... alpha=100, score=0.6222662165572943, total= 0.1s  
[CV] alpha=1000 .....  
[CV] ..... alpha=1000, score=0.613637713391447, total= 0.0s  
[CV] alpha=1000 .....  
[CV] ..... alpha=1000, score=0.6016036461965866, total= 0.1s  
[CV] alpha=1000 .....  
[CV] ..... alpha=1000, score=0.616077807102254, total= 0.1s  
[CV] alpha=10000 .....  
[CV] ..... alpha=10000, score=0.6126899546488527, total= 0.1s  
[CV] alpha=10000 .....  
[CV] ..... alpha=10000, score=0.6024184306843408, total= 0.0s  
[CV] alpha=10000 .....  
[CV] ..... alpha=10000, score=0.6151066929306674, total= 0.0s
```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 5.0s finished
```

```
Optimal lambda {'alpha': 0.1}  
AUC: 0.6831580079414458
```



Accuracy on crossvalidated data: 0.6949963477925947

Train Confusion Matrix

In [96]:

```
from sklearn.metrics import confusion_matrix  
from sklearn.metrics import accuracy_score  
from sklearn.linear_model import SGDClassifier  
labeler=[0, 1]
```

```

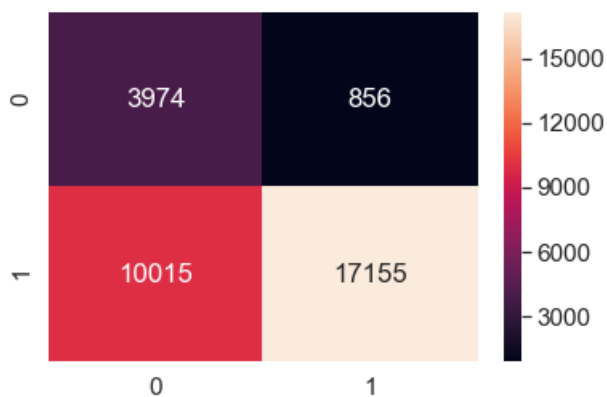
labels=[0,1]
print("Train confusion matrix")
cm=confusion_matrix(y_train, predict(x1, y_train, x1, y_train,0.01),labels)
sns.set(font_scale=1.4)
sns.heatmap(cm,fmt='d',annot=True)

```

Train confusion matrix
the maximum value of $tpr*(1-fpr)$ 0.5380048250757633 for threshold -0.235
train or test AUC = 0.8055649080134206

Out[96]:

<matplotlib.axes._subplots.AxesSubplot at 0x2c62fcf8>



Test Confusion matrix

In [97]:

```

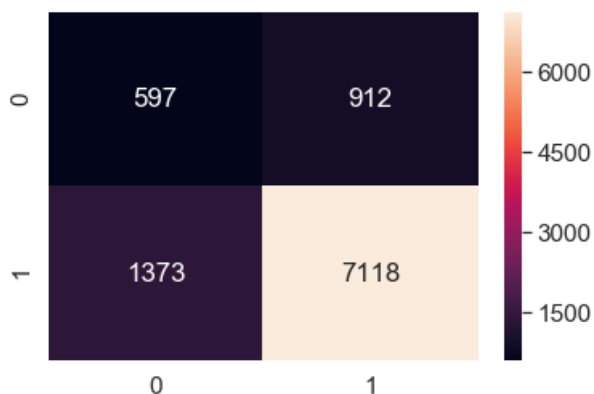
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.linear_model import SGDClassifier
labels=[0,1]
print("Test confusion matrix")
cm=confusion_matrix(y_test, predict(x1, y_train, t1, y_test,0.01),labels)
sns.set(font_scale=1.4)
sns.heatmap(cm,fmt='d',annot=True)

```

Test confusion matrix
the maximum value of $tpr*(1-fpr)$ 0.41206184164592 for threshold 0.492
train or test AUC = 0.6866287845884299

Out[97]:

<matplotlib.axes._subplots.AxesSubplot at 0x2f995eb8>

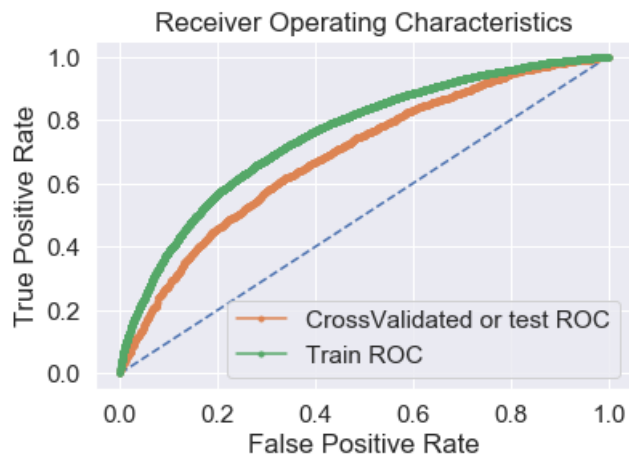


Applying Logistic regression on TFIDF

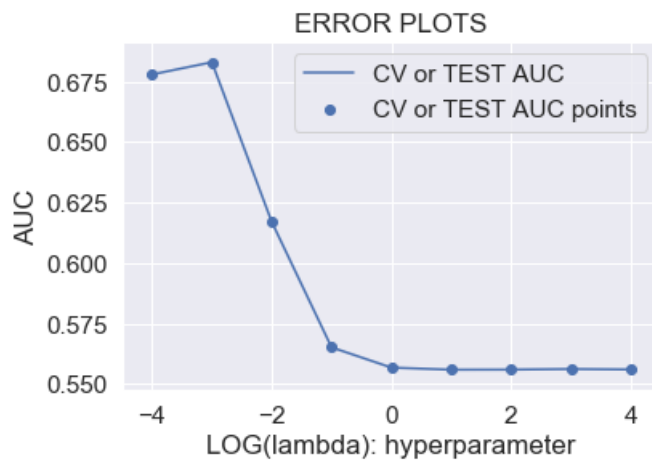
ROC SCORE on CV data

In [98]:

```
alpha_from_roc(x2,y_train,cv2,y_cv)
```



Optimal lambda 0.001
AUC: 0.6830868220380379



K FOLD Cross Validation

In [99]:

```
kfold_cv(x2,y_train)
```

Optimal lambda : 1000
CV SCORE 0.849062499705191

GridsearchCV

In [100]:

```
grid_search(x2,y_train,cv2,y_cv)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[CV] alpha=0.0001
[CV] alpha=0.0001, score=0.6560637014573677, total= 0.1s

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.1s remaining: 0.0s

[CV] alpha=0.0001
[CV] alpha=0.0001, score=0.6679551247893774, total= 0.1s

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.3s remaining: 0.0s

[CV] alpha=0.0001
[CV] alpha=0.0001, score=0.6713970217062091, total= 0.1s

[Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 0.5s remaining: 0.0s

[CV] alpha=0.001
[CV] alpha=0.001, score=0.658641852120833, total= 0.1s

[Parallel(n_jobs=1)]: Done 4 out of 4 | elapsed: 0.7s remaining: 0.0s

[CV] alpha=0.001
[CV] alpha=0.001, score=0.6627244154859115, total= 0.1s

[Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 0.9s remaining: 0.0s

[CV] alpha=0.001
[CV] alpha=0.001, score=0.6704304342339178, total= 0.1s

[Parallel(n_jobs=1)]: Done 6 out of 6 | elapsed: 1.0s remaining: 0.0s

[CV] alpha=0.01
[CV] alpha=0.01, score=0.5931978079478691, total= 0.1s

[Parallel(n_jobs=1)]: Done 7 out of 7 | elapsed: 1.3s remaining: 0.0s

[CV] alpha=0.01
[CV] alpha=0.01, score=0.6246317148055414, total= 0.1s

[Parallel(n_jobs=1)]: Done 8 out of 8 | elapsed: 1.5s remaining: 0.0s

[CV] alpha=0.01
[CV] alpha=0.01, score=0.5942127521234335, total= 0.1s

[Parallel(n_jobs=1)]: Done 9 out of 9 | elapsed: 1.7s remaining: 0.0s

[CV] alpha=0.1
[CV] alpha=0.1, score=0.5597694930039357, total= 0.0s

[Parallel(n_jobs=1)]: Done 10 out of 10 | elapsed: 1.8s remaining: 0.0s

[CV] alpha=0.1
[CV] alpha=0.1, score=0.576518625653813, total= 0.1s

[Parallel(n_jobs=1)]: Done 11 out of 11 | elapsed: 2.0s remaining: 0.0s

[CV] alpha=0.1
[CV] alpha=0.1, score=0.5544584558742839, total= 0.1s

[Parallel(n_jobs=1)]: Done 12 out of 12 | elapsed: 2.3s remaining: 0.0s

[CV] alpha=1
[CV] alpha=1, score=0.5517901461893857, total= 0.1s

[Parallel(n_jobs=1)]: Done 13 out of 13 | elapsed: 2.5s remaining: 0.0s

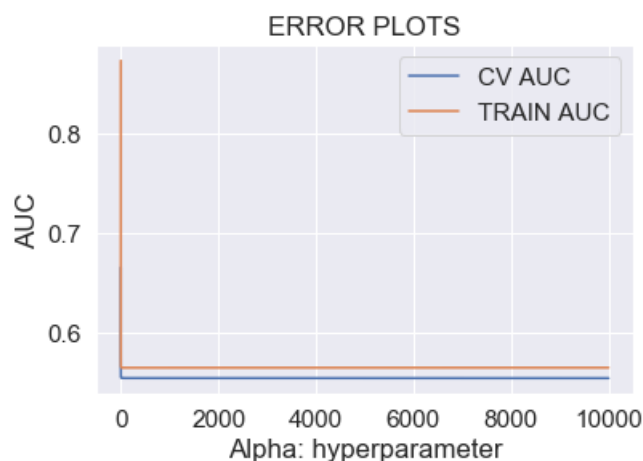
```
[CV] alpha=1 .....
[CV] ..... alpha=1, score=0.5657892697525746, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 14 out of 14 | elapsed: 2.7s remaining: 0.0s
```

```
[CV] alpha=1 .....
[CV] ..... alpha=1, score=0.5471276035379584, total= 0.0s
[CV] alpha=10 .....
[CV] ..... alpha=10, score=0.5506665514543159, total= 0.0s
[CV] alpha=10 .....
[CV] ..... alpha=10, score=0.5654467187453923, total= 0.1s
[CV] alpha=10 .....
[CV] ..... alpha=10, score=0.546251275706165, total= 0.1s
[CV] alpha=100 .....
[CV] ..... alpha=100, score=0.5506416573570972, total= 0.1s
[CV] alpha=100 .....
[CV] ..... alpha=100, score=0.5653660700998576, total= 0.1s
[CV] alpha=100 .....
[CV] ..... alpha=100, score=0.5459751470491407, total= 0.1s
[CV] alpha=1000 .....
[CV] ..... alpha=1000, score=0.5506459778202508, total= 0.0s
[CV] alpha=1000 .....
[CV] ..... alpha=1000, score=0.5653797858559009, total= 0.1s
[CV] alpha=1000 .....
[CV] ..... alpha=1000, score=0.5459042973465312, total= 0.1s
[CV] alpha=10000 .....
[CV] ..... alpha=10000, score=0.5507683909429377, total= 0.1s
[CV] alpha=10000 .....
[CV] ..... alpha=10000, score=0.5654037198501964, total= 0.1s
[CV] alpha=10000 .....
[CV] ..... alpha=10000, score=0.5455964817944384, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 5.1s finished
```

```
Optimal lambda {'alpha': 0.0001}
AUC: 0.6651384204091393
```



Accuracy on crossvalidated data: 0.6641415678281812

Train Confusion Matrix

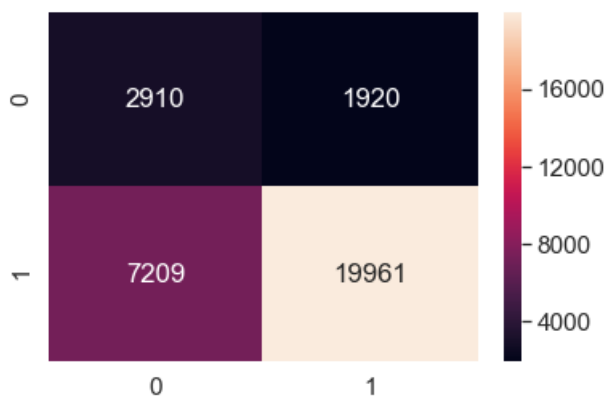
In [79]:

```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.linear_model import SGDClassifier
labels=[0,1]
print("Train confusion matrix")
cm=confusion_matrix(y_train, predict(x2, y_train, x2, y_train,0.001),labels)
sns.set(font_scale=1.4)
sns.heatmap(cm,fmt='d',annot=True)
```

Train confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.4565381681628821 for threshold 0.073
train or test AUC = 0.7374855503001956

Out[79]:

<matplotlib.axes._subplots.AxesSubplot at 0x2c5446a0>



Test Confusion matrix

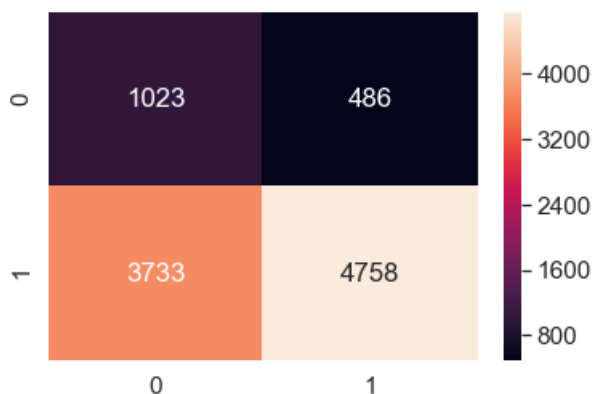
In [83]:

```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.linear_model import SGDClassifier
labels=[0,1]
print("Test confusion matrix")
cm=confusion_matrix(y_test, predict(x2, y_train, t2, y_test,0.001),labels)
sns.set(font_scale=1.4)
sns.heatmap(cm,fmt='d',annot=True)
```

Test confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.3912343471460329 for threshold -0.061
train or test AUC = 0.6632008678116205

Out[83]:

<matplotlib.axes._subplots.AxesSubplot at 0x2c906828>

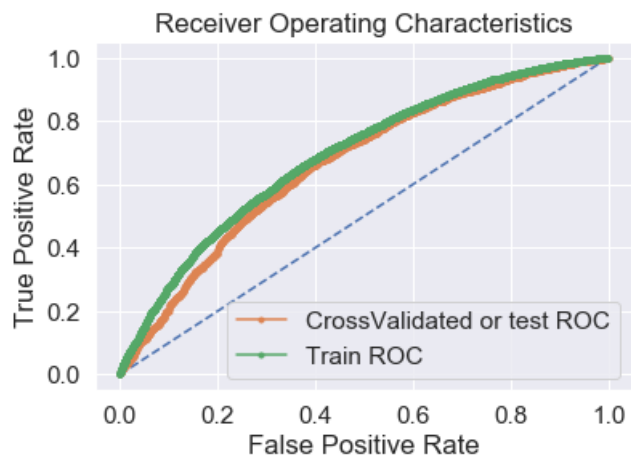


Applying Logistic regression on AVGW2V

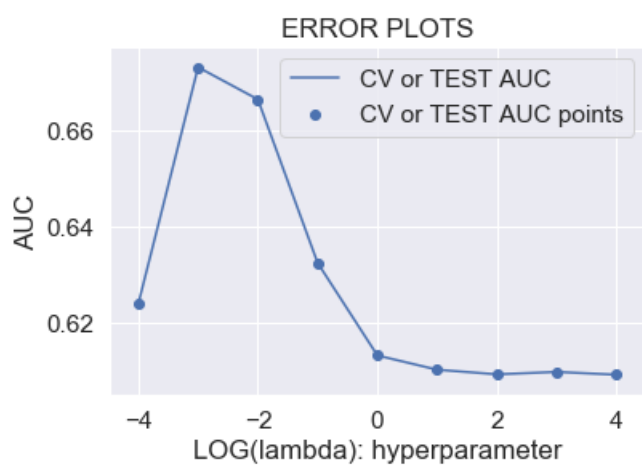
ROC SCORE on CV data

In [101]:

```
alpha_from_roc(x3,y_train,cv3,y_cv)
```

Optimal lambda 0.001
AUC: 0.6730995684797121



K FOLD Cross Validation

In [102]:

```
kfold_cv(x3,y_train)
```

Optimal lambda : 1000
CV SCORE 0.849062499705191
Optimal lambda : 10000
CV SCORE 0.849062499705191

GridsearchCV

In [103]:

```
grid_search(x3,y_train,cv3,y_cv)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[CV] alpha=0.0001
[CV] alpha=0.0001, score=0.6154577256396171, total= 0.3s

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.3s remaining: 0.0s

[CV] alpha=0.0001

[CV] alpha=0.0001, score=0.603192685112987, total= 0.3s

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.8s remaining: 0.0s

[CV] alpha=0.0001

[CV] alpha=0.0001, score=0.6362245681803217, total= 0.3s

[Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 1.2s remaining: 0.0s

[CV] alpha=0.001

[CV] alpha=0.001, score=0.6534009245791148, total= 0.3s

[Parallel(n_jobs=1)]: Done 4 out of 4 | elapsed: 1.7s remaining: 0.0s

[CV] alpha=0.001

[CV] alpha=0.001, score=0.6514580191567965, total= 0.3s

[Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 2.2s remaining: 0.0s

[CV] alpha=0.001

[CV] alpha=0.001, score=0.6521499078199416, total= 0.3s

[Parallel(n_jobs=1)]: Done 6 out of 6 | elapsed: 2.6s remaining: 0.0s

[CV] alpha=0.01

[CV] alpha=0.01, score=0.6398479059812355, total= 0.3s

[Parallel(n_jobs=1)]: Done 7 out of 7 | elapsed: 3.1s remaining: 0.0s

[CV] alpha=0.01

[CV] alpha=0.01, score=0.6429838078642031, total= 0.3s

[Parallel(n_jobs=1)]: Done 8 out of 8 | elapsed: 3.5s remaining: 0.0s

[CV] alpha=0.01

[CV] alpha=0.01, score=0.6436301796413757, total= 0.3s

[Parallel(n_jobs=1)]: Done 9 out of 9 | elapsed: 4.0s remaining: 0.0s

[CV] alpha=0.1

[CV] alpha=0.1, score=0.6081818599525298, total= 0.3s

[Parallel(n_jobs=1)]: Done 10 out of 10 | elapsed: 4.4s remaining: 0.0s

[CV] alpha=0.1

[CV] alpha=0.1, score=0.612247552937675, total= 0.3s

[Parallel(n_jobs=1)]: Done 11 out of 11 | elapsed: 4.8s remaining: 0.0s

[CV] alpha=0.1

[CV] alpha=0.1, score=0.6125310010315388, total= 0.3s

[Parallel(n_jobs=1)]: Done 12 out of 12 | elapsed: 5.3s remaining: 0.0s

[CV] alpha=1

[CV] alpha=1, score=0.5924099063419599, total= 0.3s

[Parallel(n_jobs=1)]: Done 13 out of 13 | elapsed: 5.8s remaining: 0.0s

[CV] alpha=1

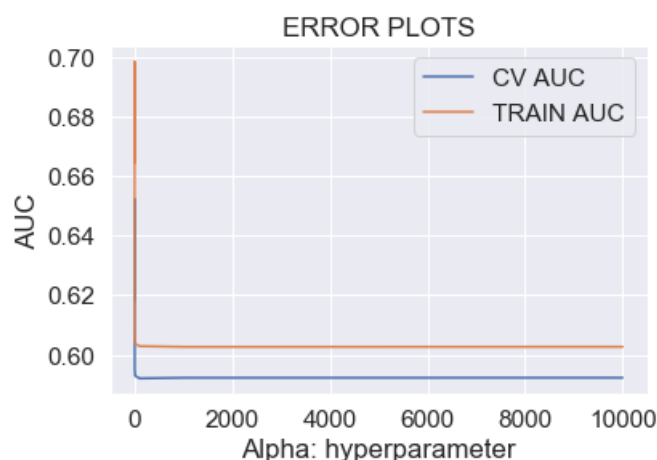
[CV] alpha=1, score=0.5961763901090197, total= 0.3s

```
[Parallel(n_jobs=1)]: Done 14 out of 14 | elapsed: 6.3s remaining: 0.0s
```

```
[CV] alpha=1 .....  
[CV] ..... alpha=1, score=0.5959888643197331, total= 0.3s  
[CV] alpha=10 .....  
[CV] ..... alpha=10, score=0.5929872025138239, total= 0.3s  
[CV] alpha=10 .....  
[CV] ..... alpha=10, score=0.5910359990590991, total= 0.3s  
[CV] alpha=10 .....  
[CV] ..... alpha=10, score=0.5948126083664376, total= 0.3s  
[CV] alpha=100 .....  
[CV] ..... alpha=100, score=0.5922873560617127, total= 0.3s  
[CV] alpha=100 .....  
[CV] ..... alpha=100, score=0.5889573762307319, total= 0.3s  
[CV] alpha=100 .....  
[CV] ..... alpha=100, score=0.5950449789302723, total= 0.3s  
[CV] alpha=1000 .....  
[CV] ..... alpha=1000, score=0.5932157755882859, total= 0.3s  
[CV] alpha=1000 .....  
[CV] ..... alpha=1000, score=0.5886825810584038, total= 0.3s  
[CV] alpha=1000 .....  
[CV] ..... alpha=1000, score=0.5950180244935583, total= 0.3s  
[CV] alpha=10000 .....  
[CV] ..... alpha=10000, score=0.5932093977617258, total= 0.3s  
[CV] alpha=10000 .....  
[CV] ..... alpha=10000, score=0.5887303118894345, total= 0.3s  
[CV] alpha=10000 .....  
[CV] ..... alpha=10000, score=0.5950207679476769, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 11.9s finished
```

```
Optimal lambda {'alpha': 0.001}  
AUC: 0.6523362896762019
```



```
Accuracy on crossvalidated data: 0.6771526710509418
```

Train Confusion Matrix

```
In [81]:
```

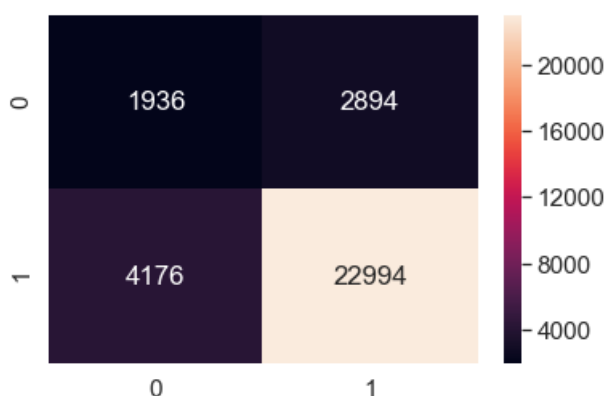
```
from sklearn.metrics import confusion_matrix  
from sklearn.metrics import accuracy_score  
from sklearn.linear_model import SGDClassifier  
labels=[0,1]  
print("Train confusion matrix")  
cm=confusion_matrix(y_train, predict(x3, y_train, x3, y_train,0.001),labels)  
sns.set(font_scale=1.4)  
sns.heatmap(cm,fmt='d',annot=True)
```

```
Train confusion matrix  
the maximum value of tpr*(1-fpr) 0.4180167658428528 for threshold 0.443  
train accuracy: 0.6771526710509418
```

train or test AUC = 0.700124833214078

Out[81]:

<matplotlib.axes._subplots.AxesSubplot at 0x2f2a4358>



Test Confusion matrix

In [84]:

```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.linear_model import SGDClassifier
labels=[0,1]
print("Test confusion matrix")
cm=confusion_matrix(y_test, predict(x3, y_train, t3, y_test,0.001),labels)
sns.set(font_scale=1.4)
sns.heatmap(cm,fmt='d',annot=True)
```

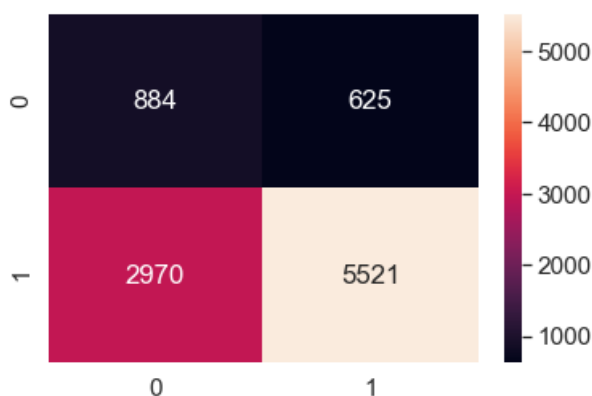
Test confusion matrix

the maximum value of $tpr \cdot (1 - fpr)$ 0.3852314995513512 for threshold 0.056

train or test AUC = 0.6578599302781825

Out[84]:

<matplotlib.axes._subplots.AxesSubplot at 0x2bf902b0>



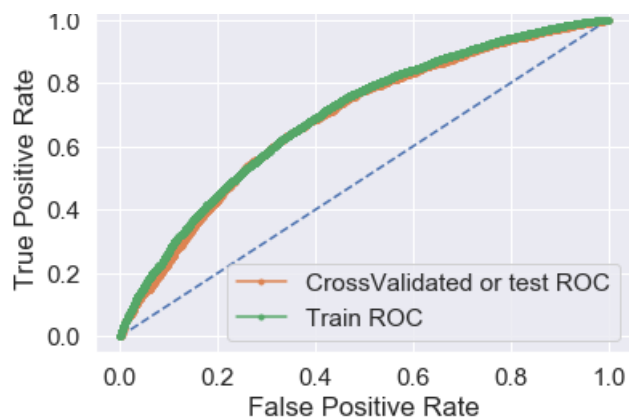
Applying Logistic regression on TFIDF Weighted W2V

ROC SCORE on CV data

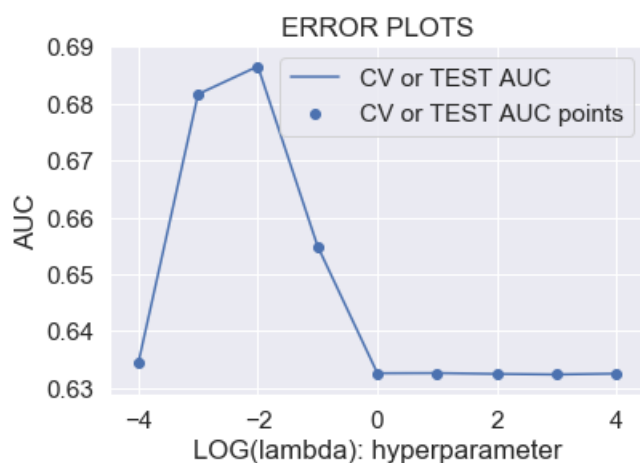
In [104]:

```
alpha_from_roc(x4,y_train,cv4,y_cv)
```

Receiver Operating Characteristics



Optimal lambda 0.01
AUC: 0.6865818180443317



K FOLD Cross Validation

In [105]:

```
kfold_cv(x4,y_train)
```

Optimal lambda : 0.0001
CV SCORE 0.6918394777802422

GridsearchCV

In [106]:

```
grid_search(x4,y_train,cv4,y_cv)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[CV] alpha=0.0001
[CV] alpha=0.0001, score=0.6135880623545702, total= 0.3s

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.3s remaining: 0.0s

[CV] alpha=0.0001
[CV] alpha=0.0001, score=0.6030510013530592, total= 0.3s

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.8s remaining: 0.0s

```
[CV] alpha=0.0001 .....  
[CV] ..... alpha=0.0001, score=0.6444844226675153, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 1.3s remaining: 0.0s
```

```
[CV] alpha=0.001 .....  
[CV] ..... alpha=0.001, score=0.666802864124177, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 4 out of 4 | elapsed: 1.7s remaining: 0.0s
```

```
[CV] alpha=0.001 .....  
[CV] ..... alpha=0.001, score=0.6664459115731491, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 2.1s remaining: 0.0s
```

```
[CV] alpha=0.001 .....  
[CV] ..... alpha=0.001, score=0.6681627636459408, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 6 out of 6 | elapsed: 2.6s remaining: 0.0s
```

```
[CV] alpha=0.01 .....  
[CV] ..... alpha=0.01, score=0.6718773509663093, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 7 out of 7 | elapsed: 3.1s remaining: 0.0s
```

```
[CV] alpha=0.01 .....  
[CV] ..... alpha=0.01, score=0.6644770147931287, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 8 out of 8 | elapsed: 3.5s remaining: 0.0s
```

```
[CV] alpha=0.01 .....  
[CV] ..... alpha=0.01, score=0.6676533042161403, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 9 out of 9 | elapsed: 4.0s remaining: 0.0s
```

```
[CV] alpha=0.1 .....  
[CV] ..... alpha=0.1, score=0.6365236181890127, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 10 out of 10 | elapsed: 4.4s remaining: 0.0s
```

```
[CV] alpha=0.1 .....  
[CV] ..... alpha=0.1, score=0.6349251839797228, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 11 out of 11 | elapsed: 4.9s remaining: 0.0s
```

```
[CV] alpha=0.1 .....  
[CV] ..... alpha=0.1, score=0.6334553255931348, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 12 out of 12 | elapsed: 5.3s remaining: 0.0s
```

```
[CV] alpha=1 .....  
[CV] ..... alpha=1, score=0.6133868522134144, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 13 out of 13 | elapsed: 5.8s remaining: 0.0s
```

```
[CV] alpha=1 .....  
[CV] ..... alpha=1, score=0.6138560682276568, total= 0.3s
```

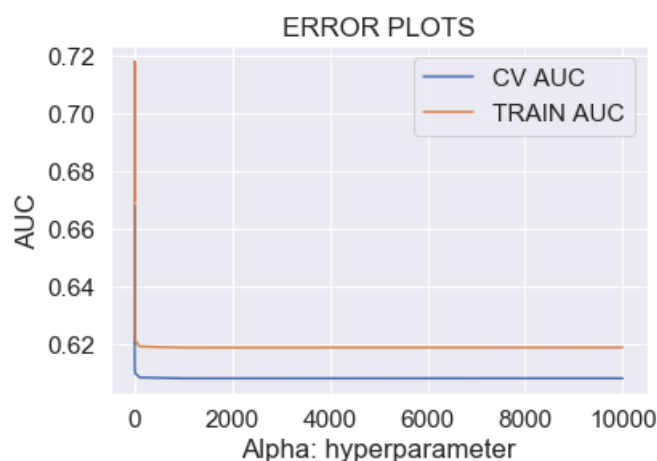
```
[Parallel(n_jobs=1)]: Done 14 out of 14 | elapsed: 6.2s remaining: 0.0s
```

```
[CV] alpha=1 .....  
[CV] ..... alpha=1, score=0.6138560682276568, total= 0.3s
```

```
[CV] ..... alpha=1, score=0.6096004433421855, total= 0.3s
[CV] alpha=10 .....
[CV] ..... alpha=10, score=0.6138590856939863, total= 0.3s
[CV] alpha=10 .....
[CV] ..... alpha=10, score=0.6083603019386535, total= 0.3s
[CV] alpha=10 .....
[CV] ..... alpha=10, score=0.6080663037991353, total= 0.3s
[CV] alpha=100 .....
[CV] ..... alpha=100, score=0.6139197093356978, total= 0.3s
[CV] alpha=100 .....
[CV] ..... alpha=100, score=0.6046695291449529, total= 0.3s
[CV] alpha=100 .....
[CV] ..... alpha=100, score=0.6072835963391349, total= 0.3s
[CV] alpha=1000 .....
[CV] ..... alpha=1000, score=0.6137961989525277, total= 0.3s
[CV] alpha=1000 .....
[CV] ..... alpha=1000, score=0.604053623119827, total= 0.3s
[CV] alpha=1000 .....
[CV] ..... alpha=1000, score=0.6070410749950618, total= 0.3s
[CV] alpha=10000 .....
[CV] ..... alpha=10000, score=0.6138066915059008, total= 0.3s
[CV] alpha=10000 .....
[CV] ..... alpha=10000, score=0.6041279625175818, total= 0.3s
[CV] alpha=10000 .....
[CV] ..... alpha=10000, score=0.6070600048284793, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 11.9s finished
```

```
Optimal lambda {'alpha': 0.01}
AUC: 0.668002567572665
```



```
Accuracy on crossvalidated data: 0.6865355937462305
```

Train Confusion Matrix

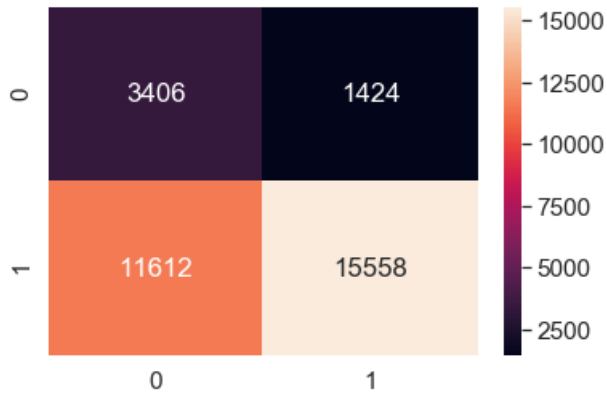
```
In [82]:
```

```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.linear_model import SGDClassifier
labels=[0,1]
print("Train confusion matrix")
cm=confusion_matrix(y_train, predict(x4, y_train, x4, y_train,0.01),labels)
sns.set(font_scale=1.4)
sns.heatmap(cm,fmt='d',annot=True)
```

```
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.4150901272640403 for threshold -0.066
train or test AUC = 0.6944286072432525
```

```
Out[82]:
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x305c4400>
```



Test Confusion matrix

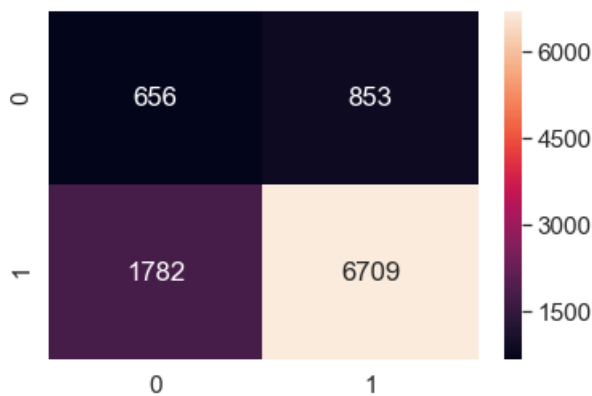
In [85]:

```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.linear_model import SGDClassifier
labels=[0,1]
print("Test confusion matrix")
cm=confusion_matrix(y_test, predict(x4, y_train, t4, y_test,0.01),labels)
sns.set(font_scale=1.4)
sns.heatmap(cm,fmt='d',annot=True)
```

Test confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.38385577868712045 for threshold 0.203
train or test AUC = 0.6640556301027111

Out[85]:

<matplotlib.axes._subplots.AxesSubplot at 0x20e8eda0>

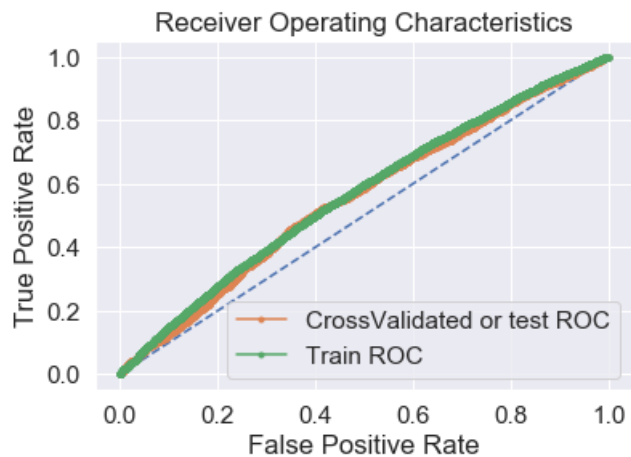


2.5 Logistic Regression with added Features `Set 5`

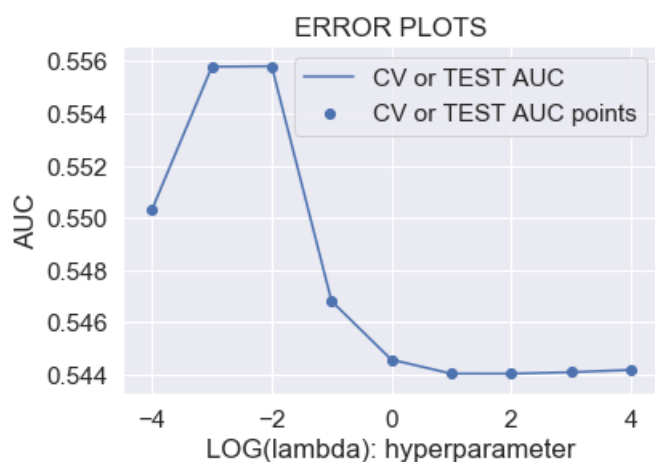
In [107]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label

alpha_from_roc(x,y_train,cv,y_cv)
kfold_cv(x,y_train)
grid_search(x,y_train,cv,y_cv)
```

Optimal lambda 0.01
AUC: 0.5558077903431708



Optimal lambda : 100
CV SCORE 0.849062499705191
Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[CV] alpha=0.0001
[CV] alpha=0.0001, score=0.5292020104555208, total= 0.0s

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.0s remaining: 0.0s

[CV] alpha=0.0001
[CV] alpha=0.0001, score=0.5177165049236135, total= 0.0s

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.0s remaining: 0.0s

[CV] alpha=0.0001
[CV] alpha=0.0001, score=0.5375247939665957, total= 0.0s

[Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 0.1s remaining: 0.0s

[CV] alpha=0.001
[CV] alpha=0.001, score=0.5294385386684881, total= 0.0s

[Parallel(n_jobs=1)]: Done 4 out of 4 | elapsed: 0.2s remaining: 0.0s

[CV] alpha=0.001
[CV] alpha=0.001, score=0.5558077903431708, total= 0.0s

[CV] alpha=0.001, score=0.56521243536521722, total= 0.0s

[Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 0.2s remaining: 0.0s

[CV] alpha=0.001
[CV] alpha=0.001, score=0.5361539585299475, total= 0.0s

[Parallel(n_jobs=1)]: Done 6 out of 6 | elapsed: 0.3s remaining: 0.0s

[CV] alpha=0.01
[CV] alpha=0.01, score=0.5476930441229013, total= 0.0s

[Parallel(n_jobs=1)]: Done 7 out of 7 | elapsed: 0.3s remaining: 0.0s

[CV] alpha=0.01
[CV] alpha=0.01, score=0.5704828014706034, total= 0.0s

[Parallel(n_jobs=1)]: Done 8 out of 8 | elapsed: 0.4s remaining: 0.0s

[CV] alpha=0.01
[CV] alpha=0.01, score=0.5487252883370278, total= 0.0s

[Parallel(n_jobs=1)]: Done 9 out of 9 | elapsed: 0.5s remaining: 0.0s

[CV] alpha=0.1
[CV] alpha=0.1, score=0.5450982973946235, total= 0.0s

[Parallel(n_jobs=1)]: Done 10 out of 10 | elapsed: 0.5s remaining: 0.0s

[CV] alpha=0.1
[CV] alpha=0.1, score=0.5622212529754618, total= 0.0s

[Parallel(n_jobs=1)]: Done 11 out of 11 | elapsed: 0.6s remaining: 0.0s

[CV] alpha=0.1
[CV] alpha=0.1, score=0.5389358210060795, total= 0.0s

[Parallel(n_jobs=1)]: Done 12 out of 12 | elapsed: 0.7s remaining: 0.0s

[CV] alpha=1
[CV] alpha=1, score=0.5425667117229254, total= 0.0s

[Parallel(n_jobs=1)]: Done 13 out of 13 | elapsed: 0.7s remaining: 0.0s

[CV] alpha=1
[CV] alpha=1, score=0.558103440117352, total= 0.0s

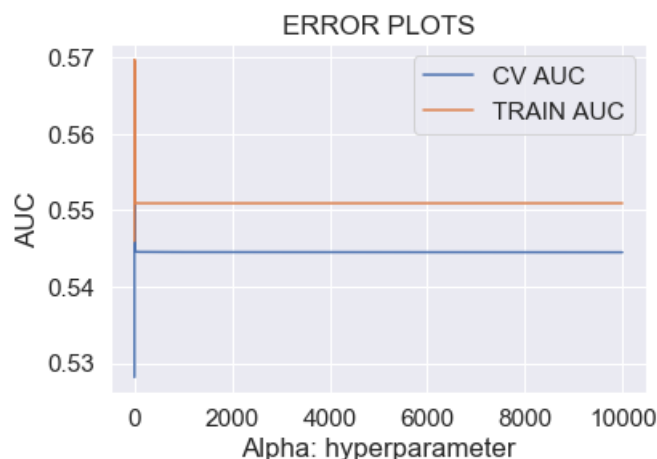
[Parallel(n_jobs=1)]: Done 14 out of 14 | elapsed: 0.8s remaining: 0.0s

[CV] alpha=1
[CV] alpha=1, score=0.5341434867655774, total= 0.0s
[CV] alpha=10
[CV] alpha=10, score=0.5423651586878685, total= 0.0s
[CV] alpha=10
[CV] alpha=10, score=0.5566947633929215, total= 0.0s
[CV] alpha=10
[CV] alpha=10, score=0.5345930703092422, total= 0.0s
[CV] alpha=100
[CV] alpha=100, score=0.5424497849026559, total= 0.0s
[CV] alpha=100
[CV] alpha=100, score=0.5566749441254388, total= 0.0s
[CV] alpha=100
[CV] alpha=100, score=0.5344210557360138, total= 0.0s
[CV] alpha=1000
[CV] alpha=1000, score=0.5424524185780073, total= 0.0s

```
[CV] ..... alpha=1000, score=0.5424534193760073, total= 0.0s
[CV] alpha=1000 .....
[CV] ..... alpha=1000, score=0.5566907172448886, total= 0.0s
[CV] alpha=1000 .....
[CV] ..... alpha=1000, score=0.5343386149397538, total= 0.0s
[CV] alpha=10000 .....
[CV] ..... alpha=10000, score=0.5424534881567875, total= 0.0s
[CV] alpha=10000 .....
[CV] ..... alpha=10000, score=0.5566730239195927, total= 0.0s
[CV] alpha=10000 .....
[CV] ..... alpha=10000, score=0.5342797678489125, total= 0.0s
```

[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 1.5s finished

Optimal lambda {'alpha': 0.01}
AUC: 0.5556339271983954



Accuracy on crossvalidated data: 0.5575033317473967

Train Confusion Matrix

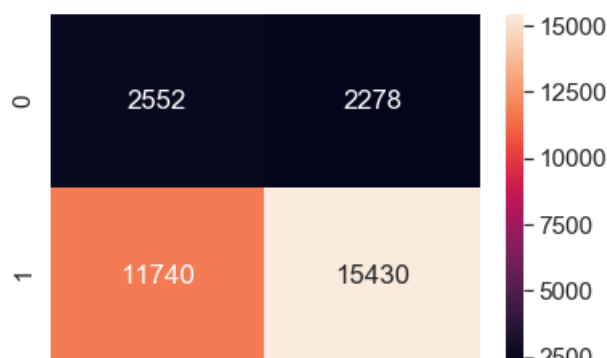
In [88]:

```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.linear_model import SGDClassifier
labels=[0,1]
print("Train confusion matrix")
cm=confusion_matrix(y_train, predict(x, y_train, x, y_train,0.01),labels)
sns.set(font_scale=1.4)
sns.heatmap(cm,fmt='d',annot=True)
```

Train confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.3031962316859342 for threshold 0.029
train or test AUC = 0.5663615979748703

Out[88]:

<matplotlib.axes._subplots.AxesSubplot at 0x2bf88b38>





Test Confusion matrix

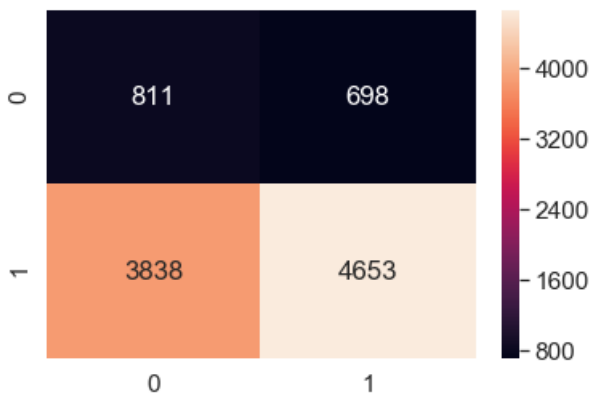
In [89]:

```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.linear_model import SGDClassifier
labels=[0,1]
print("Test confusion matrix")
cm=confusion_matrix(y_test, predict(x, y_train, t, y_test,0.01),labels)
sns.set(font_scale=1.4)
sns.heatmap(cm,fmt='d',annot=True)
```

Test confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.296003978484528 for threshold 0.003
train or test AUC = 0.5637424227843788

Out[89]:

<matplotlib.axes._subplots.AxesSubplot at 0x2bfc8748>



3. Conclusion

In [108]:

```
# Please compare all your models using Prettytable library

from prettytable import PrettyTable
x=PrettyTable()
x.field_names=["Vectorizer","Model","Features","Hyperparameter","Test AUC"]
x.add_row(["BOW","Brute","All","0.01","0.686"])
x.add_row(["TFIDF","Brute","All","0.001","0.663"])
x.add_row(["W2V","Brute","All","0.001","0.657"])
x.add_row(["TFIDF W2V","Brute","All","0.01","0.664"])
x.add_row(["-","Brute","Numerical and Categorical","0.01","0.563"])
print(x)
## Thus we see that for BOW and TFIDF Weighted W2V encoding , the text has almost no effect in the determination of lambda.
```

Vectorizer	Model	Features	Hyperparameter	Test AUC
BOW	Brute	All	0.01	0.686
TFIDF	Brute	All	0.001	0.663
W2V	Brute	All	0.001	0.657
TFIDF W2V	Brute	All	0.01	0.664
-	Brute	Numerical and Categorical	0.01	0.563

In []:

