

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: <ul style="list-style-type: none">• Art Will Make You Happy!• First Grade Fun
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following enumerated values: <ul style="list-style-type: none">• Grades PreK-2• Grades 3-5• Grades 6-8• Grades 9-12
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <ul style="list-style-type: none">• Applied Learning• Care & Hunger• Health & Sports• History & Civics• Literacy & Language• Math & Science• Music & The Arts• Special Needs• Warmth Examples: <ul style="list-style-type: none">• Music & The Arts• Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. Examples: <ul style="list-style-type: none">• Literacy

Feature	Description
<code>project_resource_summary</code>	An explanation of the resources needed for the project. Example: <ul style="list-style-type: none"> 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*
<code>project_essay_4</code>	Fourth application essay*
<code>project_submitted_datetime</code>	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
<code>teacher_id</code>	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
<code>teacher_prefix</code>	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> nan Dr. Mr. Mrs. Ms. Teacher.
<code>teacher_number_of_previously_posted_projects</code>	Number of project applications previously submitted by the same teacher. Example: 2

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

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

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

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

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

Label	Description
<code>project_is_approved</code>	A binary flag indicating whether DonorsChoose approved the project. A value of 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"

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
```

1.1 Reading Data

In [2]:

```
project_data = pd.read_csv("train_new_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
0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.00
1	p069063	Bouncy Bands for Desks (Blue support pipes)	3	14.95

1.2 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"
```

```

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

project_data['clean_subcategories'] = sub_cat_list
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]))

```

preprocessing school state

In [7]:

```

from collections import Counter
my_counter = Counter()
for word in project_data['school_state'].values:
    my_counter.update(word.split())
state_dict = dict(my_counter)
sorted_state_dict = dict(sorted(state_dict.items(), key=lambda kv: kv[1]))

```

preprocessing teacher prefix

In [8]:

```

from collections import Counter
my_counter = Counter()
for word in project_data['teacher_prefix'].values:
    my_counter.update(word.split())
prefix_dict = dict(my_counter)
sorted_prefix_dict = dict(sorted(prefix_dict.items(), key=lambda kv: kv[1]))

```

preprocessing project grade category

In [9]:

```

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

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

```

```

project_data['clean_pgc'] = pgc_list
project_data.drop(['project_grade_category'], axis=1, inplace=True)

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

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

```

1.3 Text preprocessing

In [10]:

```

# 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 [11]:

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

Out[11]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	proj
0	0	p036502	484aaf11257089a66cfedc9461c6bd0a	Ms.	NV	18-11-2016 14:45	Sup Wor Cent
1	3	p185307	525fdbb6ec7f538a48beebaa0a51b24f	Mr.	NC	12-08-2016 15:42	"Kid Insp Equi to In Activ

Decontracting function for sentence

In [12]:

```

# 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 [13]:

```
# https://stackoverflow.com/a/47091490/4084039
```

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you'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', "d
oesn'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 [14]:

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

100%|██████████| 109248/109248 [01:39<00:00, 1096.00it/s]

In [15]:

```
# after preprocessing
preprocessed_essays[2000]
```

Out[15]:

'bilingual first grade students full joy eager learn classroom place daily growth constant challenge discovery students spend year learning foundations reading writing math order succeed lives quickly becoming independent learners taking information learned apply multiple activities all ow use imagination high level thinking skills teacher low income high poverty school district students faced several challenges classroom personal folders used every day reading writing math classes provide students personal space using folders help students focus work not neighbor students able use dividers whole group independent small group time instruction generous donation project improve students self confidence independence donating project not help improve increase student attention focus ultimately help increase academic achievementnannan'

In [16]:

```
project_data["clean_essays"] = preprocessed_essays
project_data.drop(['essay'], axis=1, inplace=True)
```

1.4 Preprocessing of `project_title`

In [17]:

```
preprocessed_pt = []
for titles in tqdm(project_data["project_title"]):
    title = decontracted(titles)
    title = title.replace('\\r', ' ')
    title = title.replace('\\\"', ' ')
    title = title.replace('\\n', ' ')
    title = re.sub('[^A-Za-z0-9]+', ' ', title)
    title = ' '.join(f for f in title.split() if f not in stopwords)
    preprocessed_pt.append(title.lower().strip())
```

100%|██████████| 109248/109248 [00:04<00:00, 25427.20it/s]

In [18]:

```
project_data["clean_pt"] = preprocessed_pt
project_data.drop(['project_title'], axis=1, inplace=True)
```

number of words in title

In [19]:

```
title_word_count = []
for i in project_data["clean_pt"] :
    j = len(i.split())
    title_word_count.append(j)
project_data["title_word_count"] = title_word_count
project_data.head(5)
```

Out[19]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	prc
0	0	p036502	484aaf11257089a66cfedc9461c6bd0a	Ms.	NV	18-11-2016 14:45	Mo kin stu froi
1	3	p185307	525fdbb6ec7f538a48beebaa0a51b24f	Mr.	NC	12-08-2016 15:42	My the stu ...
2	4	p013780	a63b5547a7239eae4c1872670848e61a	Mr.	CA	06-08-2016 09:09	My ath stu ...
3	5	p063374	403c6783e9286e51ab318fba40f8d729	Mrs.	DE	05-11-2016 10:01	My ear the ma
4	6	p103285	4e156c5fb3eea2531601c8736f3751a7	Mrs.	MO	31-08-2016 00:30	Kin the gra stu

number of words in essay

In [160]:

```
essay_word_count = []
for i in project_data["clean_essays"] :
    j = len(i.split())
    essay_word_count.append(j)
project_data["essay_word_count"] = essay_word_count
project_data.head(5)
```

Out[160]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	prc
0	0	p036502	484aaf11257089a66cfedc9461c6bd0a	Ms.	NV	18-11-2016 14:45	Mo kin stu fro
1	3	p185307	525fdbb6ec7f538a48beebaa0a51b24f	Mr.	NC	12-08-2016 15:42	My the stu ...
2	4	p013780	a63b5547a7239eae4c1872670848e61a	Mr.	CA	06-08-2016 09:09	My ath stu ...
3	5	p063374	403c6783e9286e51ab318fba40f8d729	Mrs.	DE	05-11-2016 10:01	My ear the ma
4	6	p103285	4e156c5fb3eea2531601c8736f3751a7	Mrs.	MO	31-08-2016 00:30	Kin the gra stu

5 rows × 26 columns



Calculate Sentiment Scores for the essays

In [21]:

```
import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

analyser = SentimentIntensityAnalyzer()
```

In [22]:

```
neg = []
pos = []
neu = []
compound = []
for i in tqdm(project_data["clean_essays"]) :
    j = analyser.polarity_scores(i)['neg']
    k = analyser.polarity_scores(i)['pos']
    l = analyser.polarity_scores(i)['neu']
    m = analyser.polarity_scores(i)['compound']
    neg.append(j)
```

```
neg.append(j)
pos.append(k)
neu.append(l)
compound.append(m)
```

100%|██████████| 109248/109248 [21:17<00:00, 85.54it/s]

In [23]:

```
project_data["neg"] = neg
project_data["pos"] = pos
project_data["neu"] = neu
project_data["compound"] = compound
```

In [24]:

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

Out[24]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	proj
0	0	p036502	484aaf11257089a66cfedc9461c6bd0a	Ms.	NV	18-11-2016 14:45	Mos kind stud from
1	3	p185307	525fdbb6ec7f538a48beebaa0a51b24f	Mr.	NC	12-08-2016 15:42	My s the c stud ...

2 rows × 24 columns



Splitting data as train ,test and CV

In [25]:

```
from sklearn.model_selection import train_test_split
S_train, S_test, y_train, y_test = train_test_split(project_data,
project_data['project_is_approved'], test_size=0.33, stratify = project_data['project_is_approved']
])
S_train, S_cv, y_train, y_cv = train_test_split(S_train, y_train, test_size=0.30, stratify=y_train)
```

In [26]:

```
S_train.drop(['project_is_approved'], axis=1, inplace=True)
S_test.drop(['project_is_approved'], axis=1, inplace=True)
S_cv.drop(['project_is_approved'], axis=1, inplace=True)
```

1.5 Preparing data for models

In [27]:

```
project_data.columns
```

Out[27]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'project_submitted_datetime', 'project_essay_1', 'project_essay_2',
      'project_essay_3', 'project_essay_4', 'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_pos', 'clean_essay1',
```

```
clean_categories , clean_subcategories , clean_pgc , clean_essays ,
'clean_pt', 'title_word_count', 'essay_word_count', 'neg', 'pos', 'neu',
'compound'],
dtype='object')
```

we are going to consider

```
- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data

- project_title : text data
- text : text data
- project_resource_summary: text data (optinal)

- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical
```

1.5.1 Vectorizing Categorical data

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

VECTORIZING CLEAN CATEGORIES USING ONE HOT ENCODING

In [28]:

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer_clean_cat = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, b
inary=True)
vectorizer_clean_cat.fit(S_train['clean_categories'].values)
categories_one_hot_train = vectorizer_clean_cat.transform(S_train['clean_categories'].values)
categories_one_hot_test = vectorizer_clean_cat.transform(S_test['clean_categories'].values)
categories_one_hot_cv = vectorizer_clean_cat.transform(S_cv['clean_categories'].values)
print(vectorizer_clean_cat.get_feature_names())
print("Shape of matrix of Train data after one hot encoding ",categories_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding ",categories_one_hot_test.shape)
print("Shape of matrix of CV data after one hot encoding ",categories_one_hot_cv.shape)
```

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

```
Shape of matrix of Train data after one hot encoding (51237, 9)
```

```
Shape of matrix of Test data after one hot encoding (36052, 9)
```

```
Shape of matrix of CV data after one hot encoding (21959, 9)
```

VECTORIZING CLEAN SUBCATEGORIES USING ONE HOT ENCODING

In [29]:

```
vectorizer_clean_subcat = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=F
alse, binary=
True)
vectorizer_clean_subcat.fit(S_train['clean_subcategories'].values)
sub_categories_one_hot_train = vectorizer_clean_subcat.transform(S_train['clean_subcategories'].va
lues)
sub_categories_one_hot_test =
vectorizer_clean_subcat.transform(S_test['clean_subcategories'].values)
sub_categories_one_hot_cv = vectorizer_clean_subcat.transform(S_cv['clean_subcategories'].values)
print(vectorizer_clean_subcat.get_feature_names())
print("Shape of matrix of Train data after one hot encoding ",sub_categories_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding ",sub_categories_one_hot_test.shape)
print("Shape of matrix of Cross Validation data after one hot encoding ",sub_categories_one_hot_cv
.shape)
```

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular',
'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger',
'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other',
'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL',
', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences',
'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
Shape of matrix of Train data after one hot encoding (51237, 30)
Shape of matrix of Test data after one hot encoding (36052, 30)
Shape of matrix of Cross Validation data after one hot encoding (21959, 30)
```

VECTORIZING SCHOOL STATE USING ONE HOT ENCODING

In [30]:

```
# you can do the similar thing with state, teacher_prefix and project_grade_category also
vectorizer_school_state= CountVectorizer(vocabulary=list(sorted_state_dict.keys()), lowercase=False, binary=True)
vectorizer_school_state.fit(S_train['school_state'].values)
school_state_one_hot_train = vectorizer_school_state.transform(S_train['school_state'].values)
school_state_one_hot_test = vectorizer_school_state.transform(S_test['school_state'].values)
school_state_one_hot_cv = vectorizer_school_state.transform(S_cv['school_state'].values)
print(vectorizer_school_state.get_feature_names())
print("Shape of matrix of Train data after one hot encoding ",school_state_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding ",school_state_one_hot_test.shape)
print("Shape of matrix of Cross Validation data after one hot encoding ",school_state_one_hot_cv.shape)

['VT', 'WY', 'ND', 'MT', 'RI', 'SD', 'NE', 'DE', 'AK', 'NH', 'WV', 'ME', 'HI', 'DC', 'NM', 'KS', 'IA', 'ID', 'AR', 'CO', 'MN', 'OR', 'KY', 'MS', 'NV', 'MD', 'CT', 'TN', 'UT', 'AL', 'WI', 'VA', 'AZ', 'NJ', 'OK', 'WA', 'MA', 'LA', 'OH', 'MO', 'IN', 'PA', 'MI', 'SC', 'GA', 'IL', 'NC', 'FL', 'NY', 'TX', 'CA']
Shape of matrix of Train data after one hot encoding (51237, 51)
Shape of matrix of Test data after one hot encoding (36052, 51)
Shape of matrix of Cross Validation data after one hot encoding (21959, 51)
```

VECTORIZING TEACHER PREFIX USING ONE HOT ENCODING

In [31]:

```
vectorizer_prefix = CountVectorizer(vocabulary=list(sorted_prefix_dict.keys()), lowercase=False, binary=True)
vectorizer_prefix.fit(S_train['teacher_prefix'].values)
teacher_prefix_one_hot_train = vectorizer_prefix.transform(S_train['teacher_prefix'].values)
teacher_prefix_one_hot_test = vectorizer_prefix.transform(S_test['teacher_prefix'].values)
teacher_prefix_one_hot_cv = vectorizer_prefix.transform(S_cv['teacher_prefix'].values)
print(vectorizer_prefix.get_feature_names())
print("Shape of matrix of Train data after one hot encoding ",teacher_prefix_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding ",teacher_prefix_one_hot_test.shape)
print("Shape of matrix of Cross Validation data after one hot encoding ",teacher_prefix_one_hot_cv.shape)

['Dr.', 'Teacher', 'Mr.', 'Ms.', 'Mrs.']
Shape of matrix of Train data after one hot encoding (51237, 5)
Shape of matrix of Test data after one hot encoding (36052, 5)
Shape of matrix of Cross Validation data after one hot encoding (21959, 5)
```

VECTORIZING PROJECT GRADE CATEGORY USING ONE HOT ENCODING

In [32]:

```
vectorizer_pgc= CountVectorizer(vocabulary=list(sorted_pgc_dict.keys()), lowercase=False, binary=True)
vectorizer_pgc.fit(S_train['clean_pgc'].values)
clean_project_grade_category_one_hot_train = vectorizer_pgc.transform(S_train['clean_pgc'].values)
clean_project_grade_category_one_hot_test = vectorizer_pgc.transform(S_test['clean_pgc'].values)
clean_project_grade_category_one_hot_cv = vectorizer_pgc.transform(S_cv['clean_pgc'].values)
print(vectorizer_pgc.get_feature_names())
print("Shape of matrix of Train data after one hot encoding ")
```

```
print("Shape of matrix of Train data after one hot encoding",clean_project_grade_category_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding",clean_project_grade_category_one_hot_test.shape)
print("Shape of matrix of Cross Validation data after one hot encoding",clean_project_grade_category_one_hot_cv.shape)
```

```
['Grades9-12', 'Grades6-8', 'Grades3-5', 'GradesPreK-2']
Shape of matrix of Train data after one hot encoding (51237, 4)
Shape of matrix of Test data after one hot encoding (36052, 4)
Shape of matrix of Cross Validation data after one hot encoding (21959, 4)
```

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [33]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer_bow = CountVectorizer(min_df=10)
text_bow = vectorizer_bow.fit_transform(S_train["clean_essays"])
print("Shape of matrix after one hot encoding ",text_bow.shape)
```

Shape of matrix after one hot encoding (51237, 12264)

In [34]:

```
text_bow_test = vectorizer_bow.transform(S_test["clean_essays"])
print("Shape of matrix after one hot encoding ",text_bow_test.shape)
```

Shape of matrix after one hot encoding (36052, 12264)

In [35]:

```
text_bow_cv = vectorizer_bow.transform(S_cv["clean_essays"])
print("Shape of matrix after one hot encoding ",text_bow_cv.shape)
```

Shape of matrix after one hot encoding (21959, 12264)

In [36]:

```
vectorizer_title_bow = CountVectorizer( min_df=10)
title_bow_train= vectorizer_title_bow.fit_transform(S_train["clean_pt"])
print("Shape of matrix after one hot encoding ",title_bow_train.shape)
```

Shape of matrix after one hot encoding (51237, 2134)

In [37]:

```
title_bow_test = vectorizer_title_bow.transform(S_test["clean_pt"])
print("Shape of matrix after one hot encoding ",title_bow_test.shape)
```

Shape of matrix after one hot encoding (36052, 2134)

In [38]:

```
title_bow_cv = vectorizer_title_bow.transform(S_cv["clean_pt"])
print("Shape of matrix after one hot encoding ",title_bow_cv.shape)
```

Shape of matrix after one hot encoding (21959, 2134)

1.5.2.2 TFIDF vectorizer

In [161]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_tfidf_essay = TfidfVectorizer(min_df=10)
vectorizer_tfidf_essay.fit(S_train["clean_essays"])
text_tfidf_train = vectorizer_tfidf_essay.transform(S_train["clean_essays"])
print("Shape of matrix after one hot encoding ", text_tfidf_train.shape)
```

Shape of matrix after one hot encoding (51237, 12264)

In [40]:

```
text_tfidf_test = vectorizer_tfidf_essay.transform(S_test["clean_essays"])
print("Shape of matrix after one hot encoding ", text_tfidf_test.shape)
```

Shape of matrix after one hot encoding (36052, 12264)

In [41]:

```
text_tfidf_cv = vectorizer_tfidf_essay.transform(S_cv["clean_essays"])
print("Shape of matrix after one hot encoding ", text_tfidf_cv.shape)
```

Shape of matrix after one hot encoding (21959, 12264)

In [42]:

```
vectorizer_tfidf_title = TfidfVectorizer(min_df=10)
vectorizer_tfidf_title.fit(S_train["clean_pt"])
title_tfidf_train = vectorizer_tfidf_title.transform(S_train["clean_pt"])
print("Shape of matrix after one hot encoding ", title_tfidf_train.shape)
```

Shape of matrix after one hot encoding (51237, 2134)

In [43]:

```
title_tfidf_test = vectorizer_tfidf_title.transform(S_test["clean_pt"])
print("Shape of matrix after one hot encoding ", title_tfidf_test.shape)
```

Shape of matrix after one hot encoding (36052, 2134)

In [44]:

```
title_tfidf_cv = vectorizer_tfidf_title.transform(S_cv["clean_pt"])
print("Shape of matrix after one hot encoding ", title_tfidf_cv.shape)
```

Shape of matrix after one hot encoding (21959, 2134)

1.5.2.3 Using Pretrained Models: Avg W2V

In [45]:

```
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
def loadGloveModel(gloveFile):
    print("Loading Glove Model")
    f = open(gloveFile, 'r', encoding="utf8")
    model = {}
    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding
    print("Done.", len(model), " words loaded!")
    return model
model = loadGloveModel('glove.42B.300d.txt')
```

```

words = []
for i in preprocessed_essays:
    words.extend(i.split(' '))

for i in preprocessed_pt:
    words.extend(i.split(' '))
print("all the words in the coupus", len(words))
words = set(words)
print("the unique words in the coupus", len(words))

inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove vectors and our coupus", \
      len(inter_words), "("np.round(len(inter_words)/len(words)*100,3), "%)"")

words_courpus = {}
words_glove = set(model.keys())
for i in words:
    if i in words_glove:
        words_courpus[i] = model[i]
print("word 2 vec length", len(words_courpus))

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/

import pickle
with open('glove_vectors', 'wb') as f:
    pickle.dump(words_courpus, f)

```

Loading Glove Model

279727it [01:06, 4210.13it/s]

Done. 279727 words loaded!
all the words in the coupus 15565024
the unique words in the coupus 58960
The number of words that are present in both glove vectors and our coupus 44760 (75.916 %)
word 2 vec length 44760

In [46]:

```

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/
# make sure you have the glove_vectors file
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())

```

In [47]:

```

# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_train["clean_essays"]): # 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_train.append(vector)

print(len(avg_w2v_vectors_train))
print(len(avg_w2v_vectors_train[0]))

```

100%|██████████| 51237/51237 [00:54<00:00, 945.53it/s]

51237
300

In [48]:

```
avg_w2v_vectors_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_test["clean_essays"]): # 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_test.append(vector)

print(len(avg_w2v_vectors_test))
print(len(avg_w2v_vectors_test[0]))
```

100%|██████████| 36052/36052 [00:15<00:00, 2390.45it/s]

36052
300

In [49]:

```
avg_w2v_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_cv["clean_essays"]): # 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_cv.append(vector)

print(len(avg_w2v_vectors_cv))
print(len(avg_w2v_vectors_cv[0]))
```

100%|██████████| 21959/21959 [00:09<00:00, 2321.57it/s]

21959
300

In [50]:

```
avg_w2v_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_train["clean_pt"]): # 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_title_train.append(vector)

print(len(avg_w2v_title_train))
print(len(avg_w2v_title_train[0]))
```

100%|██████████| 51237/51237 [00:01<00:00, 40073.48it/s]

51237
300

In [51]:

```
avg_w2v_title_test = []; # the avg-w2v for each sentence/review is stored in this list
```



```

avg_w2v_title_test = []; # the avg_w2v for each sentence/review is stored in this list
for sentence in tqdm(S_test["clean_pt"]): # 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_title_test.append(vector)

print(len(avg_w2v_vectors_test))
print(len(avg_w2v_vectors_test[0]))

```

100%|██████████| 36052/36052 [00:01<00:00, 33533.25it/s]

36052
300

In [52]:

```

avg_w2v_title_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_cv["clean_pt"]): # 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_title_cv.append(vector)

print(len(avg_w2v_title_cv))
print(len(avg_w2v_title_cv[0]))

```

100%|██████████| 21959/21959 [00:00<00:00, 31817.48it/s]

21959
300

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [53]:

```

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

```

In [54]:

```

tfidf_w2v_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_train["clean_essays"]): # 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 tf
            idf 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_train.append(vector)

```

```
print(len(tfidf_w2v_vectors_train))
print(len(tfidf_w2v_vectors_train[0]))
```

```
100%|██████████| 51237/51237 [02:38<00:00, 322.33it/s]
```

```
51237
300
```

In [55]:

```
tfidf_w2v_vectors_test= []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_test["clean_essays"]): # 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 tf
            idf 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_test.append(vector)

print(len(tfidf_w2v_vectors_test))
print(len(tfidf_w2v_vectors_test[0]))
```

```
100%|██████████| 36052/36052 [01:50<00:00, 327.25it/s]
```

```
36052
300
```

In [56]:

```
tfidf_w2v_vectors_cv= []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_cv["clean_essays"]): # 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 tf
            idf 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_cv.append(vector)

print(len(tfidf_w2v_vectors_cv))
print(len(tfidf_w2v_vectors_cv[0]))
```

```
100%|██████████| 21959/21959 [01:07<00:00, 325.35it/s]
```

```
21959
300
```

In [57]:

```
# Similarly you can vectorize for title also
# average Word2Vec
# compute average word2vec for each review.
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(S_train["clean_pt"])
```

```

# 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())
tfidf_w2v_ppt_train= []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_train["clean_pt"]): # 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 tf
            idf 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_ppt_train.append(vector)

print(len(tfidf_w2v_ppt_train))
print(len(tfidf_w2v_ppt_train[0]))

```

```
100%|██████████| 51237/51237 [00:02<00:00, 18230.67it/s]
```

```
51237
300
```

In [58]:

```

# Similarly you can vectorize for title also
# average Word2Vec
# compute average word2vec for each review.
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_w2v_ppt_test= []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_test["clean_pt"]): # 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 tf
            idf 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_ppt_test.append(vector)

print(len(tfidf_w2v_ppt_test))
print(len(tfidf_w2v_ppt_test[0]))

```

```
100%|██████████| 36052/36052 [00:01<00:00, 19774.90it/s]
```

```
36052
300
```

In [59]:

```

tfidf_w2v_ppt_cv= []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(S_cv["clean_pt"]): # 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 tf
            idf 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_ppt_cv.append(vector)

print(len(tfidf_w2v_ppt_cv))
print(len(tfidf_w2v_ppt_cv[0]))

```

```
100%|██████████| 21959/21959 [00:01<00:00, 18722.58it/s]
```

```
21959
300
```

1.5.3 Vectorizing Numerical features

In [60]:

```

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

```

In [61]:

```

S_train = pd.merge(S_train, price_data, on='id', how='left')
S_test = pd.merge(S_test, price_data, on='id', how='left')
S_cv = pd.merge(S_cv, price_data, on='id', how='left')

```

Normalizing Price

In [62]:

```

# check this one: https://www.youtube.com/watch?v=0HOqOcln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html
from sklearn.preprocessing import Normalizer

price_scalar = Normalizer()
price_scalar.fit(S_train['price'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
price_standardized_train = price_scalar.transform(S_train['price'].values.reshape(-1, 1))
price_standardized_test = price_scalar.transform(S_test['price'].values.reshape(-1, 1))
price_standardized_cv = price_scalar.transform(S_cv['price'].values.reshape(-1, 1))

```

In [63]:

```

print(price_standardized_train.shape)
print(price_standardized_test.shape)
print(price_standardized_cv.shape)

```

```

(51237, 1)
(36052, 1)
(21959, 1)

```

Normalizing number of previously posted projects

In [64]:

```

price_scalar.fit(S_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
prev_project_standardized_train =
price_scalar.transform(S_train['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
prev_project_standardized_test =
price_scalar.transform(S_test['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
prev_project_standardized_cv =
price_scalar.transform(S_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))

```

In [65]:

```
print(prev_project_standardized_train.shape)
print(prev_project_standardized_test.shape)
print(prev_project_standardized_cv.shape)
```

```
(51237, 1)
(36052, 1)
(21959, 1)
```

Normalizing Quantity

In [66]:

```
price_scalar.fit(S_train['quantity'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
quantity_standardized_train = price_scalar.transform(S_train['quantity'].values.reshape(-1, 1))
quantity_standardized_test = price_scalar.transform(S_test['quantity'].values.reshape(-1, 1))
quantity_standardized_cv = price_scalar.transform(S_cv['quantity'].values.reshape(-1, 1))
```

In [67]:

```
print(quantity_standardized_train.shape)
print(quantity_standardized_test.shape)
print(quantity_standardized_cv.shape)
```

```
(51237, 1)
(36052, 1)
(21959, 1)
```

normalizing title word count

In [68]:

```
normalizer = Normalizer()
normalizer.fit(S_train['title_word_count'].values.reshape(-1,1))
title_word_count_train = normalizer.transform(S_train['title_word_count'].values.reshape(-1,1))
title_word_count_cv = normalizer.transform(S_cv['title_word_count'].values.reshape(-1,1))
title_word_count_test = normalizer.transform(S_test['title_word_count'].values.reshape(-1,1))
print("After vectorizations")
print(title_word_count_train.shape, y_train.shape)
print(title_word_count_cv.shape, y_cv.shape)
print(title_word_count_test.shape, y_test.shape)
```

```
After vectorizations
(51237, 1) (51237,)
(21959, 1) (21959,)
(36052, 1) (36052,)
```

NORMALIZING ESSAY WORD COUNT

In [69]:

```
normalizer = Normalizer()
normalizer.fit(S_train['essay_word_count'].values.reshape(-1,1))
essay_word_count_train = normalizer.transform(S_train['essay_word_count'].values.reshape(-1,1))
essay_word_count_cv = normalizer.transform(S_cv['essay_word_count'].values.reshape(-1,1))
essay_word_count_test = normalizer.transform(S_test['essay_word_count'].values.reshape(-1,1))
print("After vectorizations")
print(essay_word_count_train.shape, y_train.shape)
print(essay_word_count_cv.shape, y_cv.shape)
print(essay_word_count_test.shape, y_test.shape)
```

```
After vectorizations
(51237, 1) (51237,)
(21959, 1) (21959,)
(36052, 1) (36052,)
```

In [70]:

```
normalizer = Normalizer()
normalizer.fit(S_train['essay_word_count'].values.reshape(-1,1))
essay_word_count_train = normalizer.transform(S_train['essay_word_count'].values.reshape(-1,1))
essay_word_count_cv = normalizer.transform(S_cv['essay_word_count'].values.reshape(-1,1))
essay_word_count_test = normalizer.transform(S_test['essay_word_count'].values.reshape(-1,1))
print("After vectorizations")
print(essay_word_count_train.shape, y_train.shape)
print(essay_word_count_cv.shape, y_cv.shape)
print(essay_word_count_test.shape, y_test.shape)
```

After vectorizations

```
(51237, 1) (51237,)
(21959, 1) (21959,)
(36052, 1) (36052,)
```

NORMALIZING ESSAY SENTIMENT-POS

In [71]:

```
normalizer = Normalizer()
normalizer.fit(S_train['pos'].values.reshape(-1,1))
essay_sent_pos_train = normalizer.transform(S_train['pos'].values.reshape(-1,1))
essay_sent_pos_cv = normalizer.transform(S_cv['pos'].values.reshape(-1,1))
essay_sent_pos_test = normalizer.transform(S_test['pos'].values.reshape(-1,1))
print("After vectorizations")
print(essay_sent_pos_train.shape, y_train.shape)
print(essay_sent_pos_cv.shape, y_cv.shape)
print(essay_sent_pos_test.shape, y_test.shape)
```

After vectorizations

```
(51237, 1) (51237,)
(21959, 1) (21959,)
(36052, 1) (36052,)
```

NORMALIZING ESSAY SENTIMENT-NEG

In [72]:

```
normalizer = Normalizer()
normalizer.fit(S_train['neg'].values.reshape(-1,1))
essay_sent_neg_train = normalizer.transform(S_train['neg'].values.reshape(-1,1))
essay_sent_neg_cv = normalizer.transform(S_cv['neg'].values.reshape(-1,1))
essay_sent_neg_test = normalizer.transform(S_test['neg'].values.reshape(-1,1))
print("After vectorizations")
print(essay_sent_neg_train.shape, y_train.shape)
print(essay_sent_neg_cv.shape, y_cv.shape)
print(essay_sent_neg_test.shape, y_test.shape)
```

After vectorizations

```
(51237, 1) (51237,)
(21959, 1) (21959,)
(36052, 1) (36052,)
```

NORMALIZING ESSAY SENTIMENT-NEU

In [73]:

```
normalizer = Normalizer()
normalizer.fit(S_train['neu'].values.reshape(-1,1))
essay_sent_neu_train = normalizer.transform(S_train['neu'].values.reshape(-1,1))
essay_sent_neu_cv = normalizer.transform(S_cv['neu'].values.reshape(-1,1))
essay_sent_neu_test = normalizer.transform(S_test['neu'].values.reshape(-1,1))
print("After vectorizations")
print(essay_sent_neu_train.shape, y_train.shape)
print(essay_sent_neu_cv.shape, y_cv.shape)
print(essay_sent_neu_test.shape, y_test.shape)
```

After vectorizations

```
(51237, 1) (51237,)  
(21959, 1) (21959,)  
(36052, 1) (36052,)
```

NORMALIZING ESSAY SENTIMEN-COMPOUND

In [74]:

```
normalizer = Normalizer()  
normalizer.fit(S_train['compound'].values.reshape(-1,1))  
essay_sent_comp_train = normalizer.transform(S_train['compound'].values.reshape(-1,1))  
essay_sent_comp_cv = normalizer.transform(S_cv['compound'].values.reshape(-1,1))  
essay_sent_comp_test = normalizer.transform(S_test['compound'].values.reshape(-1,1))  
print("After vectorizations")  
print(essay_sent_comp_train.shape, y_train.shape)  
print(essay_sent_comp_cv.shape, y_cv.shape)  
print(essay_sent_comp_test.shape, y_test.shape)
```

```
After vectorizations  
(51237, 1) (51237,)  
(21959, 1) (21959,)  
(36052, 1) (36052,)
```

1.5.4 Merging all the above features

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

Assignment 7: Decision tree

Apply Decision Tree Classifier(DecisionTreeClassifier) on these feature sets Set 1: categorical, numerical features + project_title(BOW) + preprocessed_eassay (BOW) Set 2: categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay (TFIDF) Set 3: categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V) Set 4: categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_eassay (TFIDF W2V)

Hyper paramter tuning (best depth in range [1, 5, 10, 50, 100, 500, 100], and the best min_samples_split in range [5, 10, 100, 500]) 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

Graphviz Visualize your decision tree with Graphviz. It helps you to understand how a decision is being made, given a new vector. Since feature names are not obtained from word2vec related models, visualize only BOW & TFIDF decision trees using Graphviz Make sure to print the words in each node of the decision tree instead of printing its index. Just for visualization purpose, limit max_depth to 2 or 3 and either embed the generated images of graphviz in your notebook, or directly upload them as .png files.

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 Once after you plot the confusion matrix with the test data, get all the false positive data points Plot the WordCloud WordCloud Plot the box plot with the price of these false positive data points Plot the pdf with the teacher_number_of_previously_posted_projects of these false positive data points

[Task-2] Select 5k best features from features of Set 2 usingfeature importances, discard all the other remaining features and then apply any of the model of you choice i.e. (Dession tree, Logistic Regression, Linear SVM), you need to do hyperparameter tuning corresponding to the model you selected and procedure in step 2 and step 3

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. Decision Tree Classifier

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

Feature set 1 using BOW

In [122]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
S_BOW_train=
hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train,teacher_prefix_one_hot_train,clean_project_grade_category_one_hot_train,text_bow,title_bow_train,price_standardized_train,prev_project_standardized_train,quantity_standardized_train,title_word_count_train,essay_word_count_train,essay_sent_pos_train,essay_sent_neg_train,essay_sent_neu_train,essay_sent_comp_train)).tocsr()
print(S_BOW_train.shape)
```

(51237, 14506)

In [123]:

```
S_BOW_test= hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test,teacher_prefix_one_hot_test,clean_project_grade_category_one_hot_test,text_bow_test,title_bow_test,price_standardized_test,prev_project_standardized_test,quantity_standardized_test,title_word_count_test,essay_word_count_test,essay_sent_pos_test,essay_sent_neg_test,essay_sent_neu_test,essay_sent_comp_test)).tocsr()
print(S_BOW_test.shape)
```

(36052, 14506)

In [124]:

```
S_BOW_cv=
hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv,teacher_prefix_one_hot_cv,clean_project_grade_category_one_hot_cv,text_bow_cv,title_bow_cv,price_standardized_cv,prev_project_standardized_cv,quantity_standardized_cv,title_word_count_cv,essay_word_count_cv,essay_sent_pos_cv,essay_sent_neg_cv,essay_sent_neu_cv,essay_sent_comp_cv)).tocsr()
print(S_BOW_cv.shape)
```

(21959, 14506)

In [78]:

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    y_data_pred = []
    tr_loop = data.shape[0] - data.shape[0]%1000
    # consider you X_tr shape is 49041, then your tr_loop will be 49041 - 49041%1000 = 49000
    # in this for loop we will iterate until the last 1000 multiplier
    for i in range(0, tr_loop, 1000):
        y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
    # we will be predicting for the last data points
    if data.shape[0]%1000 != 0:
        y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
    return y_data_pred
```



```
return y_data_pred
```

finding best hyperparameter-max_depth using CV

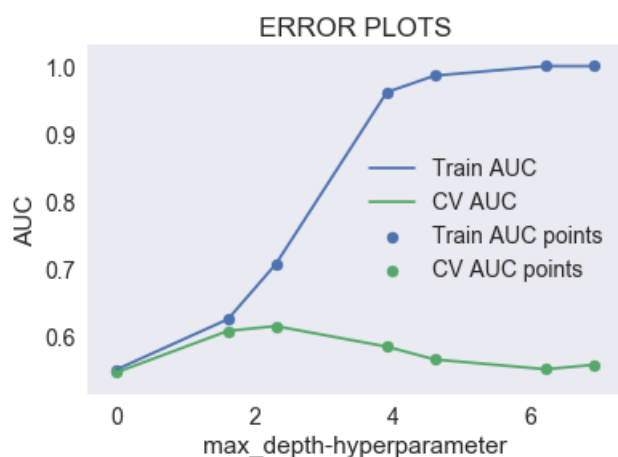
In [226]:

```
from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
train_auc = []
cv_auc = []
a = []
b = []
import math
max_depth=[1, 5, 10, 50, 100, 500, 1000]

for i in tqdm(max_depth):
    dtc= DecisionTreeClassifier(max_depth=i,class_weight="balanced")
    l=dtc.fit(S_BOW_train, y_train)
    y_train_pred = batch_predict(dtc,S_BOW_train)
    y_cv_pred = batch_predict(dtc, S_BOW_cv)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
    # class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
    a.append(y_train_pred)
    b.append(y_cv_pred)

plt.plot([math.log(i) for i in max_depth],train_auc, label='Train AUC')
plt.plot([math.log(i) for i in max_depth],cv_auc, label='CV AUC')
plt.scatter([math.log(i) for i in max_depth],train_auc, label='Train AUC points')
plt.scatter([math.log(i) for i in max_depth],cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("max_depth-hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

100%|██████████| 7/7 [08:50<00:00, 108.45s/it]



using CV for finding best hyperparameter-min_samples_split:

In [227]:

```
from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
train_auc = []
cv_auc = []
a = []
b = []
import math
min_samples_split=[5, 50, 100, 500]
```

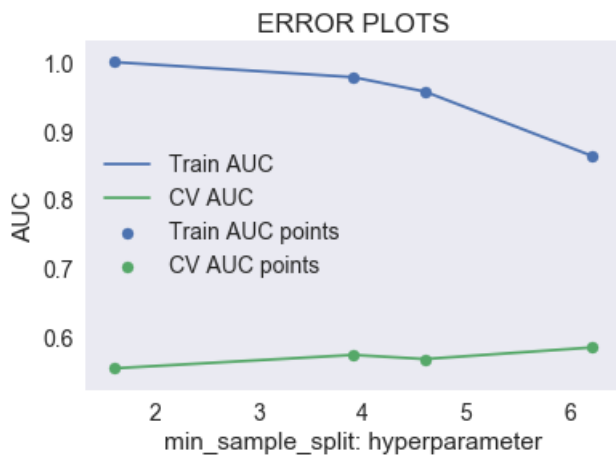
```

for i in tqdm(min_samples_split):
    dtc= DecisionTreeClassifier(min_samples_split=i,class_weight="balanced")
    l=dtc.fit(S_BOW_train, y_train)
    y_train_pred = batch_predict(dtc,S_BOW_train)
    y_cv_pred = batch_predict(dtc, S_BOW_cv)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
    a.append(y_train_pred)
    b.append(y_cv_pred)

plt.plot([math.log(i) for i in min_samples_split],train_auc, label='Train AUC')
plt.plot([math.log(i) for i in min_samples_split],cv_auc, label='CV AUC')
plt.scatter([math.log(i) for i in min_samples_split],train_auc, label='Train AUC points')
plt.scatter([math.log(i) for i in min_samples_split],cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("min_sample_split: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```

100%|██████████| 4/4 [08:29<00:00, 125.40s/it]



we will consider max_depth=5 and min samples split=5

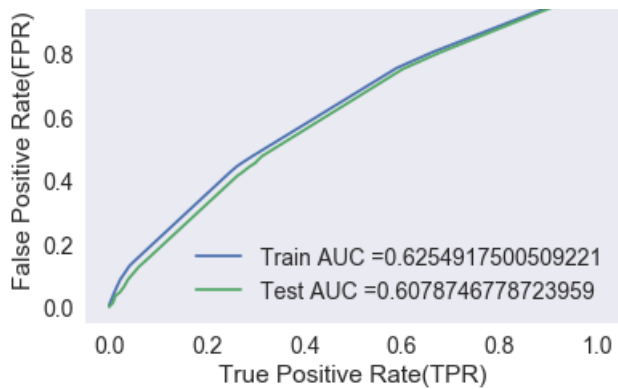
In [125]:

```

#https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
model = DecisionTreeClassifier(max_depth = 5,min_samples_split=5,random_state=0, class_weight='balanced')
model.fit(S_BOW_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
y_train_pred = batch_predict(model, S_BOW_train)
y_test_pred = batch_predict(model, S_BOW_test)
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()

```





In [153]:

```
def prediction(proba, threshold, fpr, tpr):
    t = threshold[np.argmax(fpr*(1-tpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    predictions = []
    global predictions1
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    predictions1=predictions
    return predictions
```

confusion matrix for train data

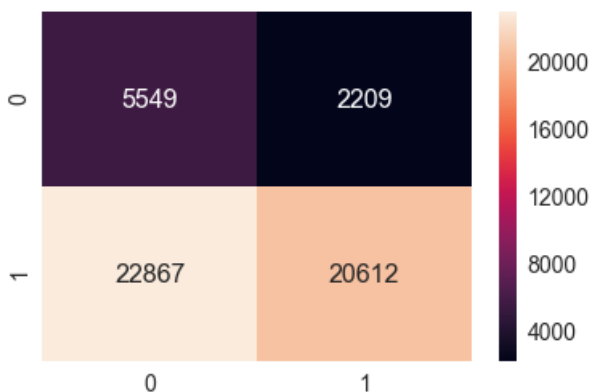
In [154]:

```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train, prediction(y_train_pred, tr_thresholds,
train_fpr, train_tpr)), range(2),range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of tpr*(1-fpr) 0.3390826248709637 for threshold 0.49

Out[154]:

<matplotlib.axes._subplots.AxesSubplot at 0x17bd9c31f98>



Confusion matrix for test data

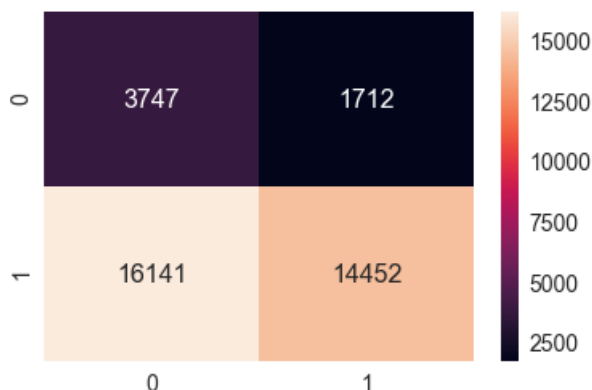
In [155]:

```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test,prediction(y_test_pred, tr_thresholds,
train_fpr, train_tpr)), range(2),range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of $\text{tpr} \times (1 - \text{fpr})$ 0.3390826248709637 for threshold 0.49

Out[155]:

<matplotlib.axes._subplots.AxesSubplot at 0x17bd9875320>



In [127]:

```
features_name_BOW= []
for a in vectorizer_clean_cat.get_feature_names() :
    features_name_BOW.append(a)
for a in vectorizer_clean_subcat.get_feature_names() :
    features_name_BOW.append(a)
for a in vectorizer_school_state.get_feature_names() :
    features_name_BOW.append(a)
for a in vectorizer_pgc.get_feature_names() :
    features_name_BOW.append(a)
for a in vectorizer_prefix.get_feature_names() :
    features_name_BOW.append(a)
features_name_BOW.append("price")
features_name_BOW.append("prev_proposed_projects")
features_name_BOW.append("quantity")
features_name_BOW.append("essay_word_count")
features_name_BOW.append("title_word_count")
features_name_BOW.append("pos")
features_name_BOW.append("neg")
features_name_BOW.append("neu")
features_name_BOW.append("compound")
for a in vectorizer_bow.get_feature_names() :
    features_name_BOW.append(a)
for a in vectorizer_title_bow.get_feature_names() :
    features_name_BOW.append(a)
print(len(features_name_BOW))
```

14506

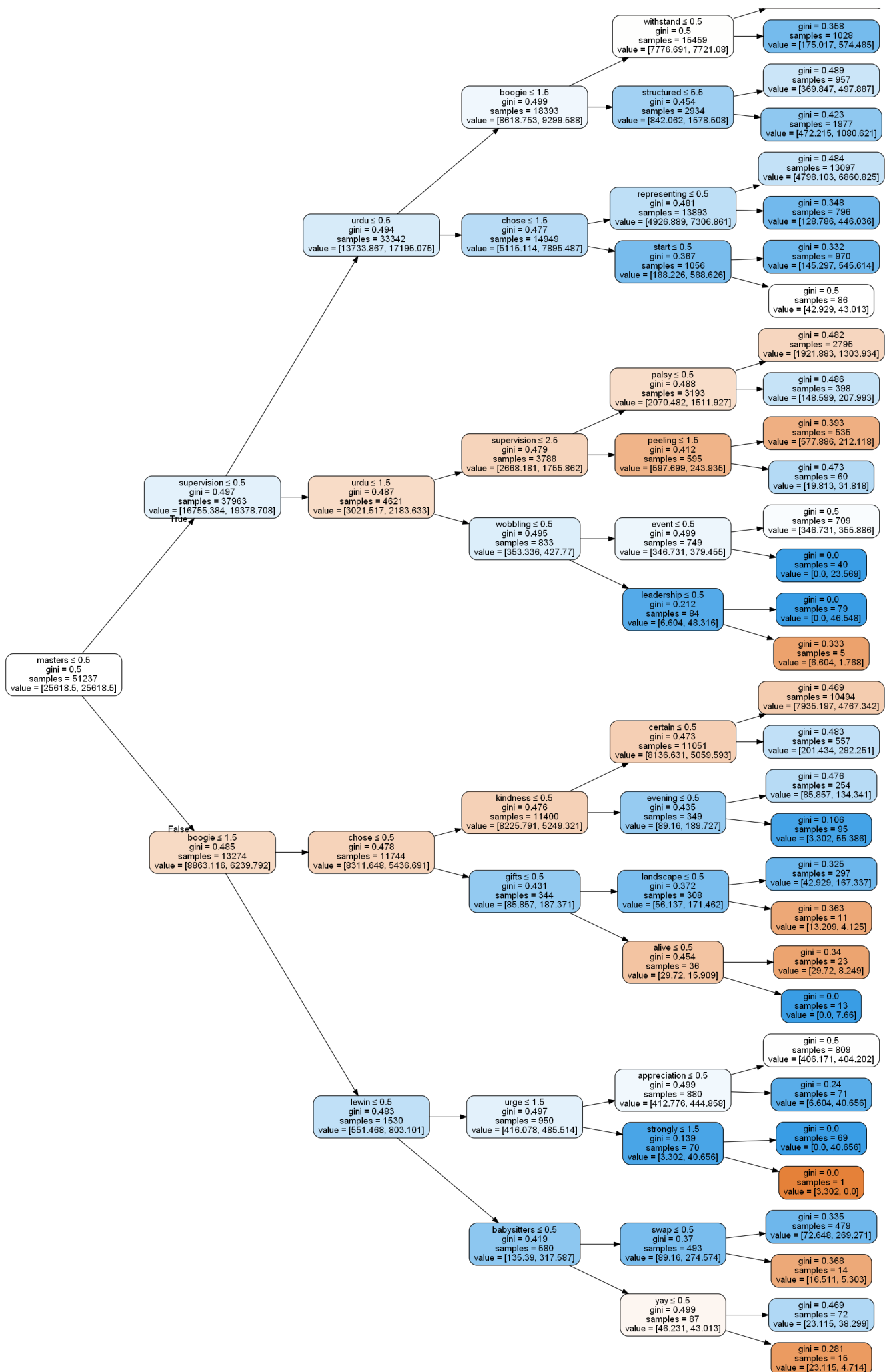
graphviz-tree visualization

In [128]:

```
import warnings
warnings.filterwarnings("ignore")
from sklearn.externals.six import StringIO
from IPython.display import Image
from sklearn.tree import export_graphviz
import pydotplus
dot_data = StringIO()
export_graphviz(model, out_file=dot_data, filled=True, rounded=True, special_characters=True, feature_names=features_name_BOW, rotate=True)
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
Image(graph.create_png())
```

Out[128]:

gini = 0.5
samples = 14431
value = [7601.674, 7146.595]



Finding false positive points

In [162]:

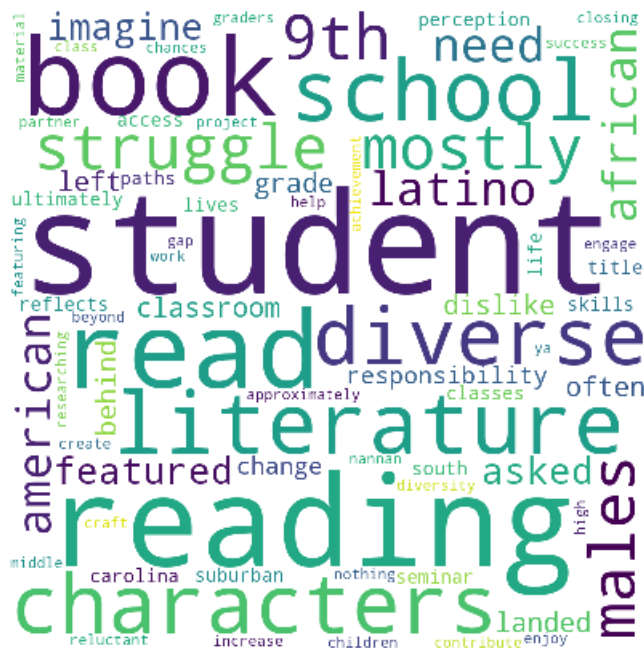
```
#https://www.google.com/search?
q=geeks+for+geeks+false+positive&rlz=1C1SQJL_enIN849IN849&oq=geeks+for+geeks+false+positive&aqs=ch
.69i57j3315.6431j0j7&sourceid=chrome&ie=UTF-8
#https://github.com/pskadasi/DecisionTrees_DonorsChoose/blob/master/Copy_of_8_DonorsChoose_DT_(1).:

fpi = []
for i in range(len(y_test)) :
    if (y_test.values[i] == 0) & (predictions1[i] == 1) :
        fpi.append(i)
fp_essay1 = []
for i in fpi :
    fp_essay1.append(S_test["clean_essays"].values[i])
```

word cloud

In [164]:

```
from wordcloud import WordCloud, STOPWORDS
comment_words = ' '
stopwords = set(STOPWORDS)
for val in fp_essay1 :
    val = str(val)
    tokens = val.split()
for i in range(len(tokens)):
    tokens[i] = tokens[i].lower()
for words in tokens :
    comment_words = comment_words + words + ' '
wordcloud = WordCloud(width = 800, height = 800, background_color ='white', stopwords =
stopwords,min_font_size = 10).generate(comment_words)
plt.figure(figsize = (6, 6), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad = 0)
plt.show()
```



In [165]:

```
cols = S_test.columns
S_test_falsePos1 = pd.DataFrame(columns=cols)
# get the data of the false positives
for i in fpi : # (in fpi all the false positives data points indexes)
    S_test_falsePos1 = S_test_falsePos1.append(S_test.filter(items=[i], axis=0))
```

```
S_test_falsePos1.head(1)
len(S_test_falsePos1)
```

Out[165]:

1712

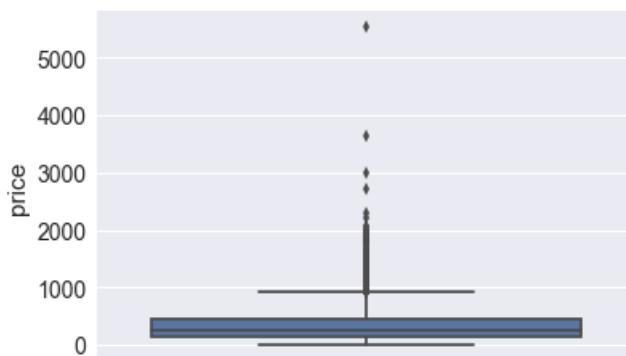
Boxplot:

In [166]:

```
sns.boxplot(y='price', data=S_test_falsePos1)
```

Out[166]:

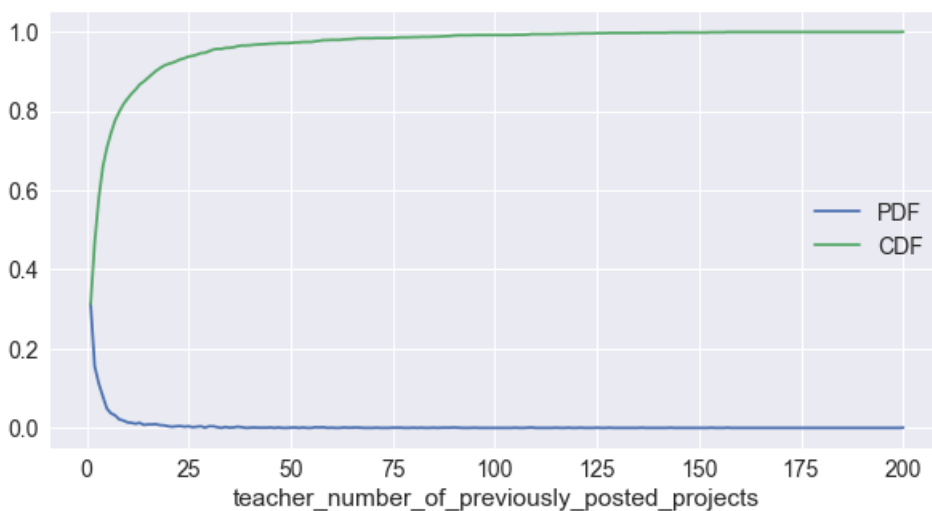
<matplotlib.axes._subplots.AxesSubplot at 0x17bf7e2d080>



PDF (FP ,teacher_number_of_previously_posted_projects)

In [167]:

```
plt.figure(figsize=(10,5))
counts, bin_edges = np.histogram(S_test_falsePos1['teacher_number_of_previously_posted_projects'],
,bins='auto', density=True)
pdf = counts/sum(counts)
cdf = np.cumsum(pdf)
pdf_FP, = plt.plot(bin_edges[1:], pdf)
cdf_FP, = plt.plot(bin_edges[1:], cdf)
plt.legend([pdf_FP, cdf_FP], ["PDF", "CDF"])
plt.xlabel('teacher_number_of_previously_posted_projects')
plt.show()
```



Feature set 2 USING TFIDF_Train

In [89]:

```
# Please write all the code with proper documentation
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
S_TFIDF_train=
hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train,teacher_prefix_one_hot_train,clean_project_grade_category_one_hot_train,text_tfidf_train,title_tfidf_train,price_standardized_train,prev_project_standardized_train,quantity_standardized_train,title_word_count_train,essay_word_count_train,essay_sent_pos_train,essay_sent_neg_train,essay_sent_neu_train,essay_sent_comp_train)).tocsr()
S_TFIDF_train.shape
```

Out[89]:

(51237, 14506)

In [90]:

```
S_TFIDF_test=
hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test,teacher_prefix_one_hot_test,clean_project_grade_category_one_hot_test,text_tfidf_test,title_tfidf_test,price_standardized_test,prev_project_standardized_test,quantity_standardized_test,title_word_count_test,essay_word_count_test,essay_sent_pos_test,essay_sent_neg_test,essay_sent_neu_test,essay_sent_comp_test)).tocsr()
S_TFIDF_test.shape
```

Out[90]:

(36052, 14506)

In [91]:

```
S_TFIDF_cv=
hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv,teacher_prefix_one_hot_cv,clean_project_grade_category_one_hot_cv,text_tfidf_cv,title_tfidf_cv,price_standardized_cv,prev_project_standardized_cv,quantity_standardized_cv,title_word_count_cv,essay_word_count_cv,essay_sent_pos_cv,essay_sent_neg_cv,essay_sent_neu_cv,essay_sent_comp_cv)).tocsr()
S_TFIDF_cv.shape
```

Out[91]:

(21959, 14506)

Finding best parameter using CV

In [228]:

```
from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
train_auc = []
cv_auc = []
a = []
b = []
import math
max_depth=[1, 5, 10, 50, 100, 500, 1000]

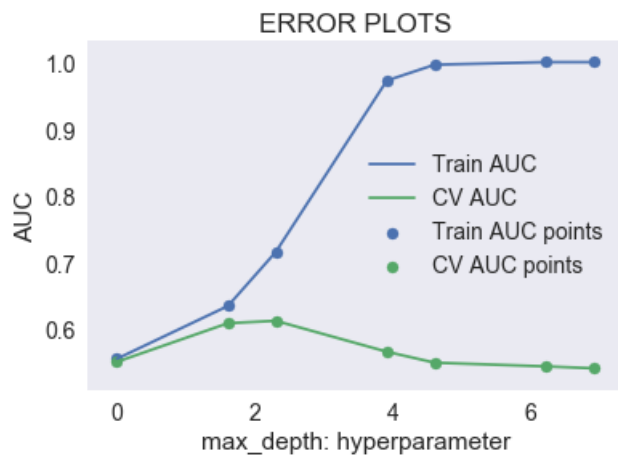
for i in tqdm(max_depth):
    dtc= DecisionTreeClassifier(max_depth=i,class_weight="balanced")
    l=dtc.fit(S_TFIDF_train, y_train)
    y_train_pred = batch_predict(dtc,S_TFIDF_train)
    y_cv_pred = batch_predict(dtc, S_TFIDF_cv)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
    a.append(y_train_pred)
    b.append(y_cv_pred)

plt.plot([math.log(i) for i in max_depth],train_auc, label='Train AUC')
```



```
plt.plot([math.log(i) for i in max_depth],cv_auc, label='CV AUC')
plt.scatter([math.log(i) for i in max_depth],train_auc, label='Train AUC points')
plt.scatter([math.log(i) for i in max_depth],cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("max_depth: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

100%|██████████| 7/7 [08:32<00:00, 98.86s/it]



Finding best hyperparameter using GridSearchCV

In [229]:

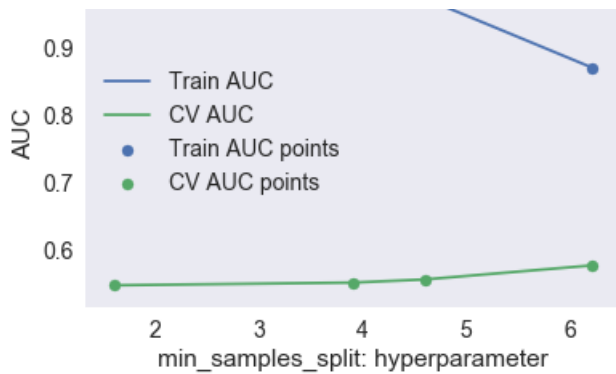
```
from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
train_auc = []
cv_auc = []
a = []
b = []
import math
min_samples_split=[5, 50, 100, 500]

for i in tqdm(min_samples_split):
    dtc= DecisionTreeClassifier(min_samples_split=i,class_weight="balanced")
    l=dtc.fit(S_TFIDF_train, y_train)
    y_train_pred = batch_predict(dtc,S_TFIDF_train)
    y_cv_pred = batch_predict(dtc, S_TFIDF_cv)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
    a.append(y_train_pred)
    b.append(y_cv_pred)

plt.plot([math.log(i) for i in min_samples_split],train_auc, label='Train AUC')
plt.plot([math.log(i) for i in min_samples_split],cv_auc, label='CV AUC')
plt.scatter([math.log(i) for i in min_samples_split],train_auc, label='Train AUC points')
plt.scatter([math.log(i) for i in min_samples_split],cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("min_samples_split: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

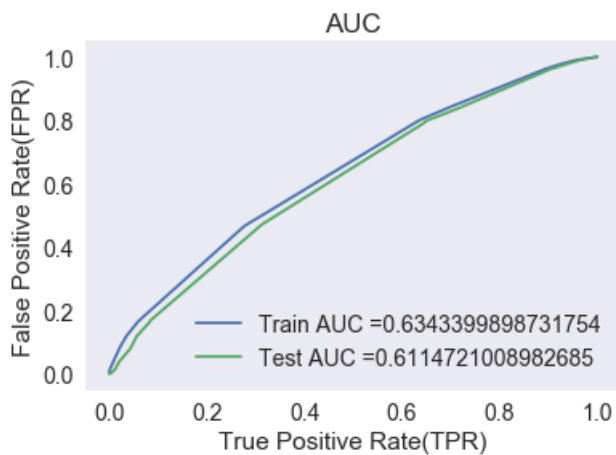
100%|██████████| 4/4 [07:29<00:00, 111.27s/it]





In [168]:

```
#https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
model = DecisionTreeClassifier(max_depth = 5,min_samples_split=5,random_state=0, class_weight='balanced')
model.fit(S_TFIDF_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs
y_train_pred = batch_predict(model, S_TFIDF_train)
y_test_pred = batch_predict(model, S_TFIDF_test)
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



confusion matrix for train data

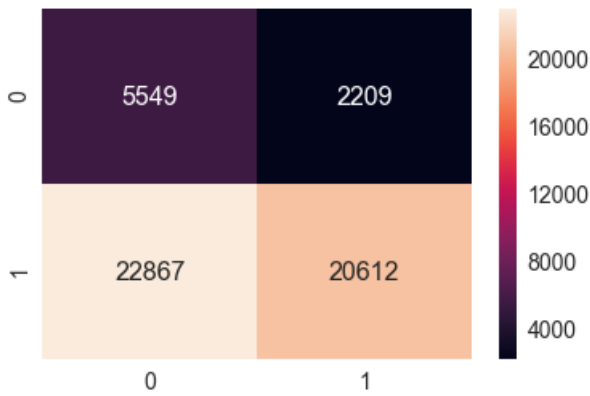
In [103]:

```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train, prediction(y_train_pred, tr_thresholds,
train_fpr, train_tpr)), range(2),range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.3390826248709637 for threshold 0.49

Out[103]:

<matplotlib.axes._subplots.AxesSubplot at 0x17bd9a25a58>



Confusion matrix for test data

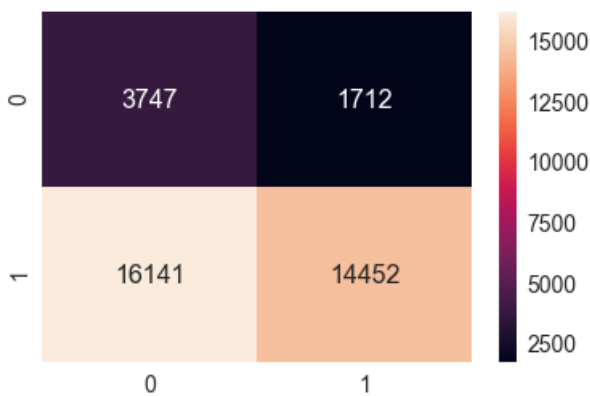
In [104]:

```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test, prediction(y_test_pred, tr_thresholds,
train_fpr, train_tpr)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.3390826248709637 for threshold 0.49

Out[104]:

<matplotlib.axes._subplots.AxesSubplot at 0x17bd9bc6ac8>



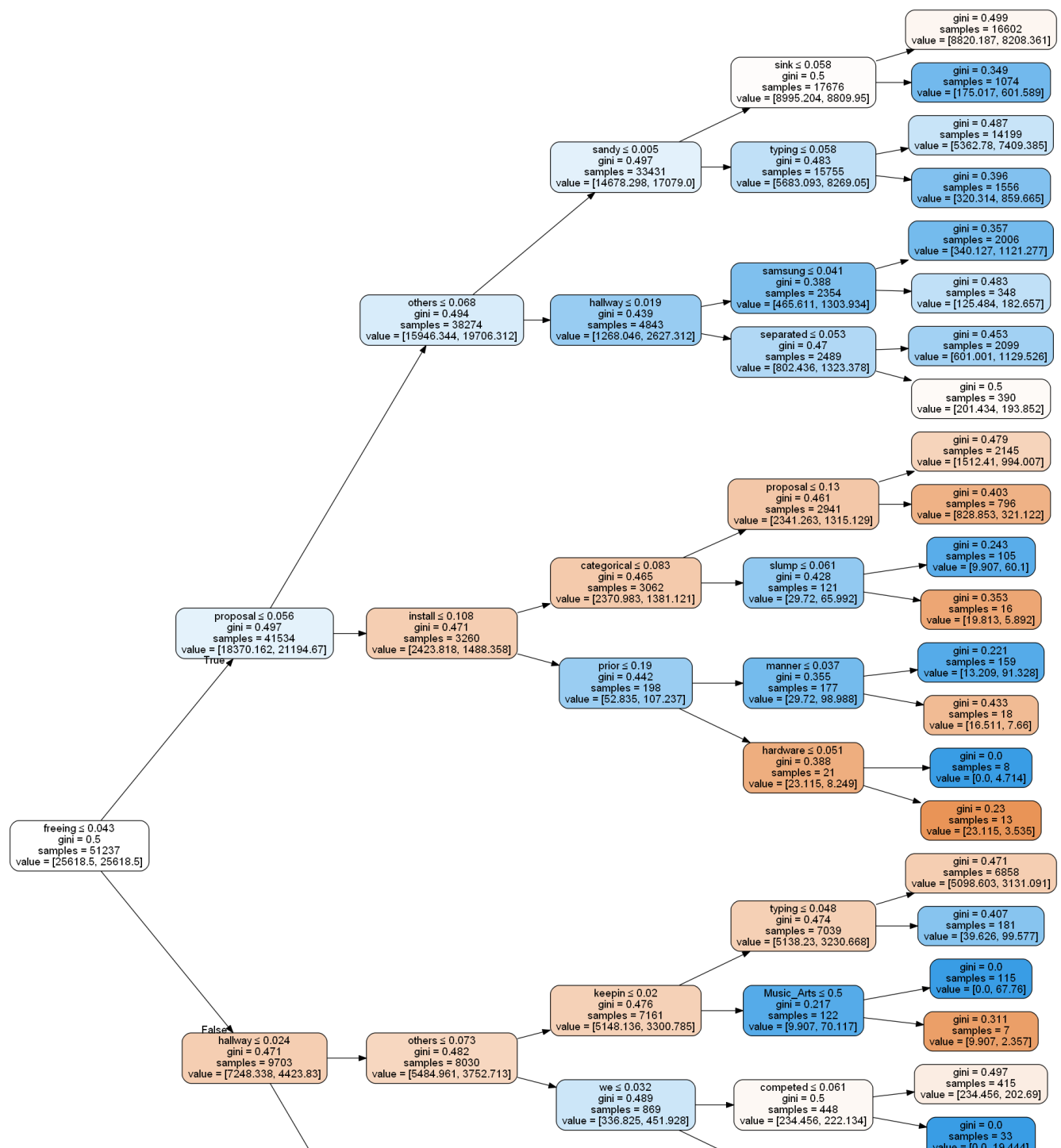
In [129]:

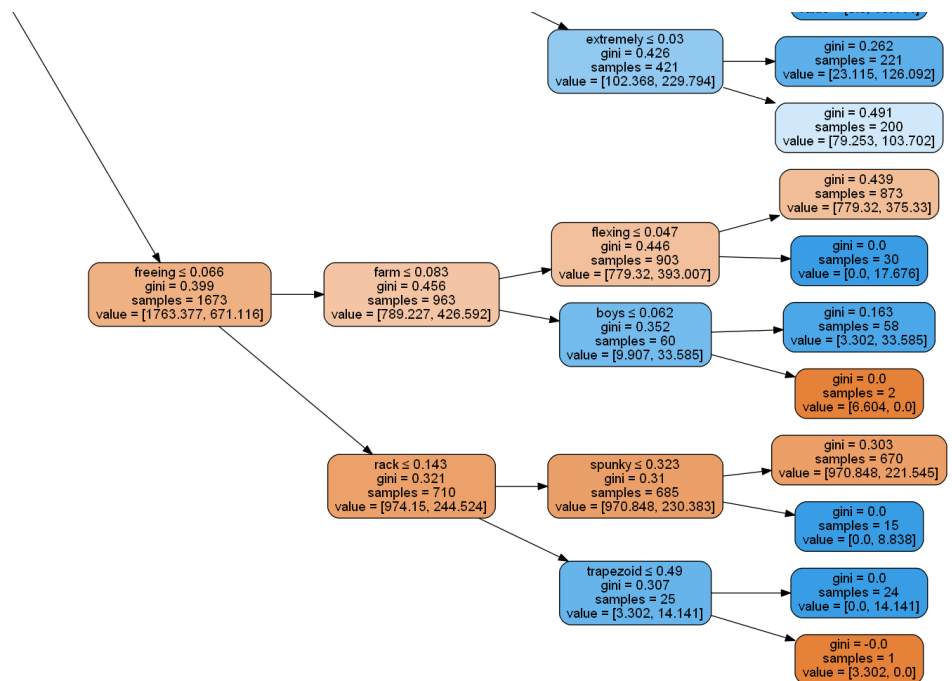
```
features_name_TFIDF= []
for a in vectorizer_clean_cat.get_feature_names() :
    features_name_TFIDF.append(a)
for a in vectorizer_clean_subcat.get_feature_names() :
    features_name_TFIDF.append(a)
for a in vectorizer_school_state.get_feature_names() :
    features_name_TFIDF.append(a)
for a in vectorizer_pgc.get_feature_names() :
    features_name_TFIDF.append(a)
for a in vectorizer_prefix.get_feature_names() :
    features_name_TFIDF.append(a)
features_name_TFIDF.append("price")
features_name_TFIDF.append("prev_proposed_projects")
features_name_TFIDF.append("quantity")
features_name_TFIDF.append("essay_word_count")
features_name_TFIDF.append("title_word_count")
features_name_TFIDF.append("pos")
features_name_TFIDF.append("neg")
features_name_TFIDF.append("neu")
features_name_TFIDF.append("compound")
for a in vectorizer_tfidf_title.get_feature_names() :
    features_name_TFIDF.append(a)
for a in vectorizer_tfidf_essay.get_feature_names() :
    features_name_TFIDF.append(a)
print(len(features_name_TFIDF))
```

14506

In [132]:

Out [132] :





finding false positive points

In [169]:

```

fpi = []
for i in range(len(y_test)) :
    if (y_test.values[i] == 0) & (predictions1[i] == 1) :
        fpi.append(i)
fp_essay1 = []
for i in fpi :
    fp_essay1.append(S_test["clean_essays"].values[i])

```

making a word cloud

In [170]:

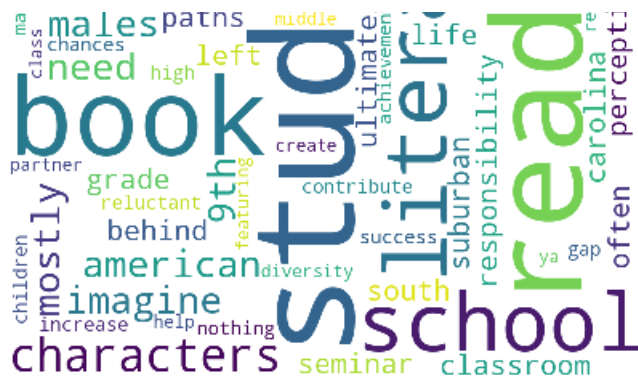
```

from wordcloud import WordCloud, STOPWORDS
comment_words = ' '
stopwords = set(STOPWORDS)
for val in fp_essay1 :
    val = str(val)
    tokens = val.split()
for i in range(len(tokens)):
    tokens[i] = tokens[i].lower()
for words in tqdm(tokens) :
    comment_words = comment_words + words + ' '
wordcloud = WordCloud(width = 800, height = 800, background_color = 'white', stopwords =
stopwords,min_font_size = 10).generate(comment_words)
plt.figure(figsize = (6, 6), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad = 0)
plt.show()

```

100% |██████████| 125/125 [00:00<00:00, 949.52it/s]





In [171]:

```
cols = S_test.columns
S_test_falsePos1 = pd.DataFrame(columns=cols)
# get the data of the false positives
for i in fpi : # (in fpi all the false positives data points indexes)
    S_test_falsePos1 = S_test_falsePos1.append(S_test.filter(items=[i], axis=0))
S_test_falsePos1.head(1)
len(S_test_falsePos1)
```

Out [171]:

1712

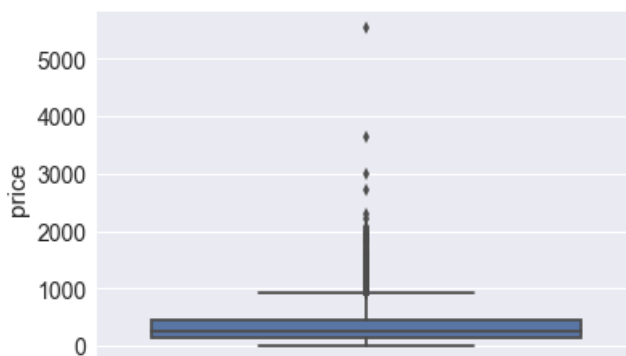
Box plot

In [172]:

```
sns.boxplot(y='price', data=S test falsePos1)
```

Out [172] :

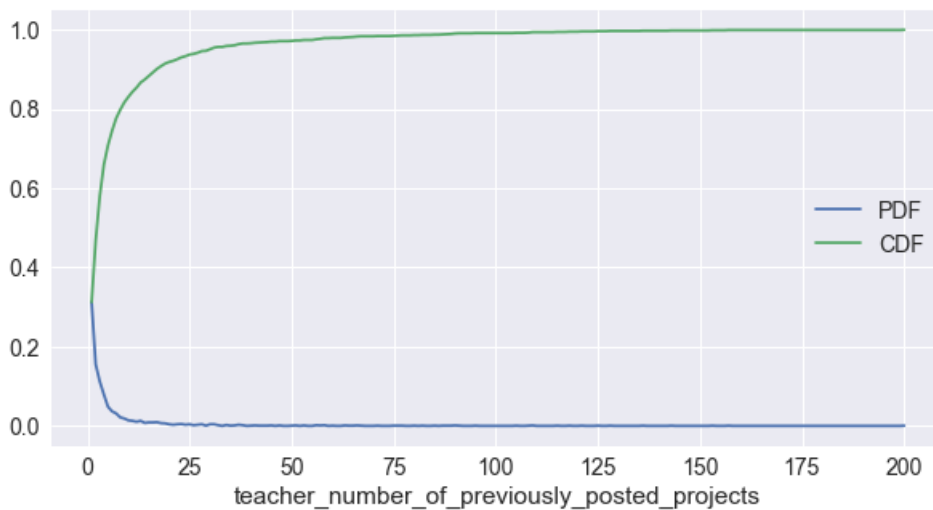
```
<matplotlib.axes. subplots.AxesSubplot at 0x17bf7e2de10>
```



PDF (FP ,teacher_number_of_previously_posted_projects)

In [173]:

```
plt.figure(figsize=(10,5))
counts, bin_edges = np.histogram(S_test_falsePos1['teacher_number_of_previously_posted_projects'],
,bins='auto', density=True)
pdf = counts/sum(counts)
cdf = np.cumsum(pdf)
pdf_FP, = plt.plot(bin_edges[1:], pdf)
cdf_FP, = plt.plot(bin_edges[1:], cdf)
plt.legend([pdf_FP, cdf_FP], ["PDF", "CDF"])
plt.xlabel('teacher_number_of_previously_posted_projects')
plt.show()
```



Feature set 3 USING AVG_W2V

In [94]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
S_avgw2v_train=
hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train,teacher_prefix_one_hot_train,clean_project_grade_category_one_hot_train,avg_w2v_vectors_train,avg_w2v_title_train,price_standardized_train,prev_project_standardized_train,quantity_standardized_train,title_word_count_train,essay_word_count_train,essay_sent_pos_train,essay_sent_neg_train,essay_sent_neu_train,essay_sent_comp_train)).tocsr()
print(S_avgw2v_train.shape)

S_avgw2v_test=
hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test,teacher_prefix_one_hot_test,clean_project_grade_category_one_hot_test,avg_w2v_vectors_test,avg_w2v_title_test,price_standardized_test,prev_project_standardized_test,quantity_standardized_test,title_word_count_test,essay_word_count_test,essay_sent_pos_test,essay_sent_neg_test,essay_sent_neu_test,essay_sent_comp_test)).tocsr()
print(S_avgw2v_test.shape)

S_avgw2v_cv=
hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv,teacher_prefix_one_hot_cv,clean_project_grade_category_one_hot_cv,avg_w2v_vectors_cv,avg_w2v_title_cv,price_standardized_cv,prev_project_standardized_cv,quantity_standardized_cv,title_word_count_cv,essay_word_count_cv,essay_sent_pos_cv,essay_sent_neg_cv,essay_sent_neu_cv,essay_sent_comp_cv)).tocsr()
print(S_avgw2v_cv.shape)
```

```
(51237, 708)
(36052, 708)
(21959, 708)
```

FINDING BEST HYPERPARAMETER USING CV

In [230]:

```
from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
train_auc = []
cv_auc = []
a = []
b = []
import math
max_depth=[1, 5, 10, 50, 100, 500, 1000]

for i in tqdm(max_depth):
    dtc= DecisionTreeClassifier(max_depth=i,class_weight="balanced")
    l=dtc.fit(S_avgw2v_train, y_train)
    y_train_pred = batch_predict(dtc,S_avgw2v_train)
    y_cv_pred = batch_predict(dtc, S_avgw2v_cv)
```

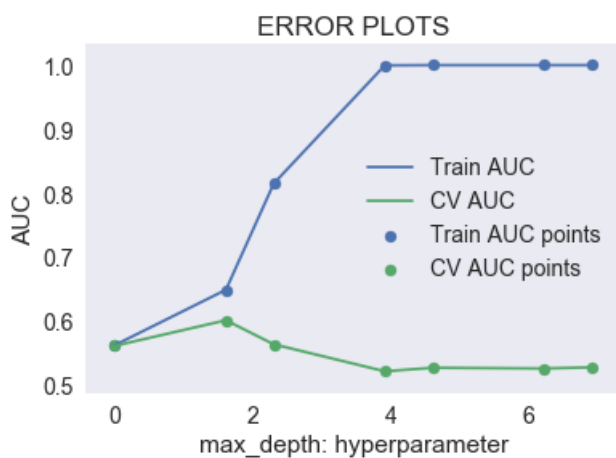
```

# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
train_auc.append(roc_auc_score(y_train,y_train_pred))
cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
a.append(y_train_pred)
b.append(y_cv_pred)

plt.plot([math.log(i) for i in max_depth],train_auc, label='Train AUC')
plt.plot([math.log(i) for i in max_depth],cv_auc, label='CV AUC')
plt.scatter([math.log(i) for i in max_depth],train_auc, label='Train AUC points')
plt.scatter([math.log(i) for i in max_depth],cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("max_depth: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```

100%|██████████| 7/7 [20:24<00:00, 212.30s/it]



FINDING BEST HYPERPARAMETER USING GRIDSEARCHCV

In [231]:

```

from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score

train_auc = []
cv_auc = []
a = []
b = []
import math
min_samples_split=[5, 50, 100, 500]

for i in tqdm(min_samples_split):
    dtc= DecisionTreeClassifier(min_samples_split=i,class_weight="balanced")
    l=dtc.fit(S_avgw2v_train, y_train)
    y_train_pred = batch_predict(dtc,S_avgw2v_train)
    y_cv_pred = batch_predict(dtc, S_avgw2v_cv)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
    class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
    a.append(y_train_pred)
    b.append(y_cv_pred)

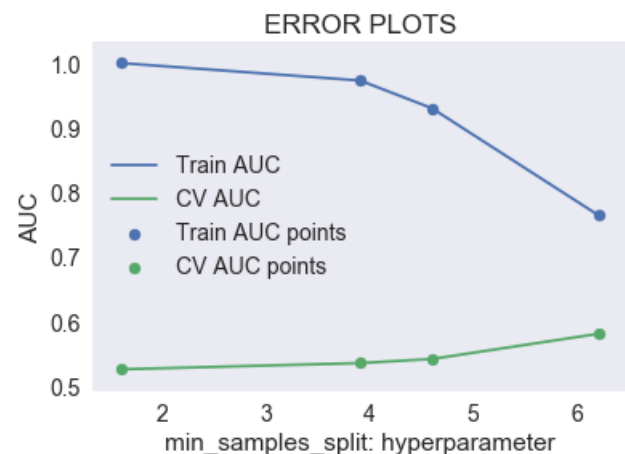
plt.plot([math.log(i) for i in min_samples_split],train_auc, label='Train AUC')
plt.plot([math.log(i) for i in min_samples_split],cv_auc, label='CV AUC')
plt.scatter([math.log(i) for i in min_samples_split],train_auc, label='Train AUC points')
plt.scatter([math.log(i) for i in min_samples_split],cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("min_samples_split: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



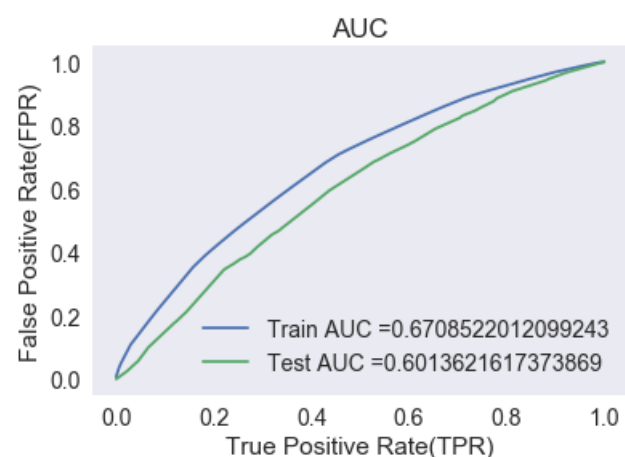
```
plt.grid()
plt.show()
```

100% | 4/4 [13:26<00:00, 196.89s/it]



In [174]:

```
#https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
model = DecisionTreeClassifier(max_depth=6,min_samples_split=7,random_state=0, class_weight='balanced')
model.fit(S_avgw2v_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs
y_train_pred = batch_predict(model, S_avgw2v_train)
y_test_pred = batch_predict(model, S_avgw2v_test)
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



CONFUSION MATRIX FOR TRAIN DATA

In [106]:

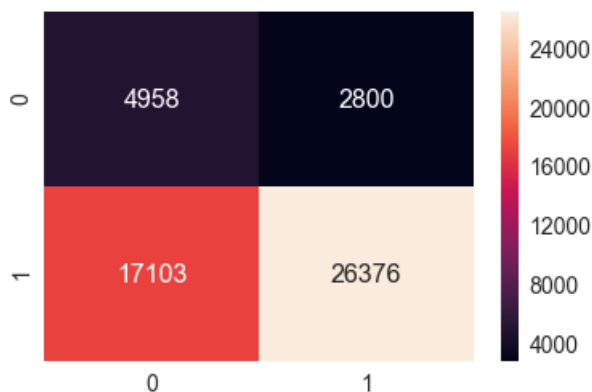
```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train, prediction(y_train_pred, tr_thresholds
, train_fpr, train_tpr)), range(2),range(2))
sns.set(font_scale=1.4)#for label size
```

```
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.38990012459811385 for threshold 0.533

Out[106]:

<matplotlib.axes._subplots.AxesSubplot at 0x17bd998b3c8>



CONFUSION MATRIX FOR TEST DATA

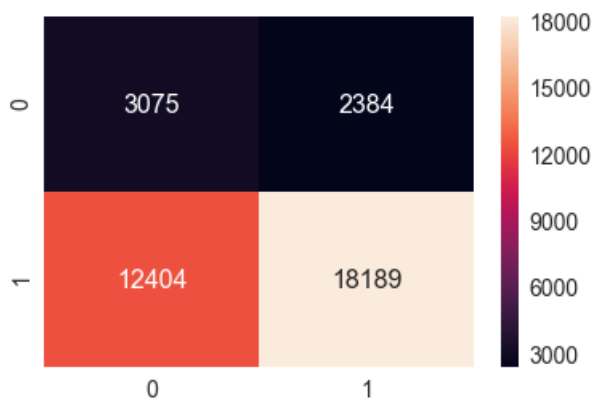
In [107]:

```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test, prediction(y_test_pred, tr_thresholds,
train_fpr, train_tpr)), range(2),range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.38990012459811385 for threshold 0.533

Out[107]:

<matplotlib.axes._subplots.AxesSubplot at 0x17bd980ac18>



finding false positive points

In [175]:

```
fpi = []
for i in range(len(y_test)) :
    if (y_test.values[i] == 0) & (predictions1[i] == 1) :
        fpi.append(i)
fp_essay1 = []
for i in fpi :
    fp_essay1.append(S_test["clean_essays"].values[i])
```

Word cloud

In [176]:

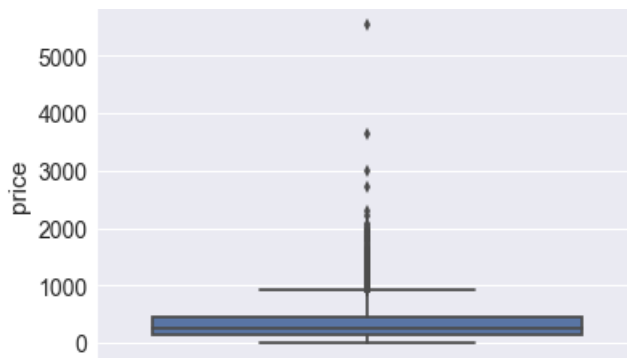
```
from wordcloud import WordCloud, STOPWORDS
comment_words = ' '
stopwords = set(STOPWORDS)
for val in fp_essay1 :
    val = str(val)
    tokens = val.split()
for i in range(len(tokens)):
    tokens[i] = tokens[i].lower()
for words in tqdm(tokens) :
    comment_words = comment_words + words + ' '
wordcloud = WordCloud(width = 800, height = 800, background_color='white', stopwords =
stopwords,min_font_size = 10).generate(comment_words)
plt.figure(figsize = (6, 6), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad = 0)
plt.show()
```

```
cols = S_test.columns
S_test_falsePos1 = pd.DataFrame(columns=cols)
# get the data of the false positives
for i in fpi : # (in fpi all the false positives data points indexes)
    S_test_falsePos1 = S_test_falsePos1.append(S_test.filter(items=[i], axis=0))
S_test_falsePos1.head(1)
len(S_test_falsePos1)
```

1712

```
sns.boxplot(y='price', data=S test falsePos1)
```

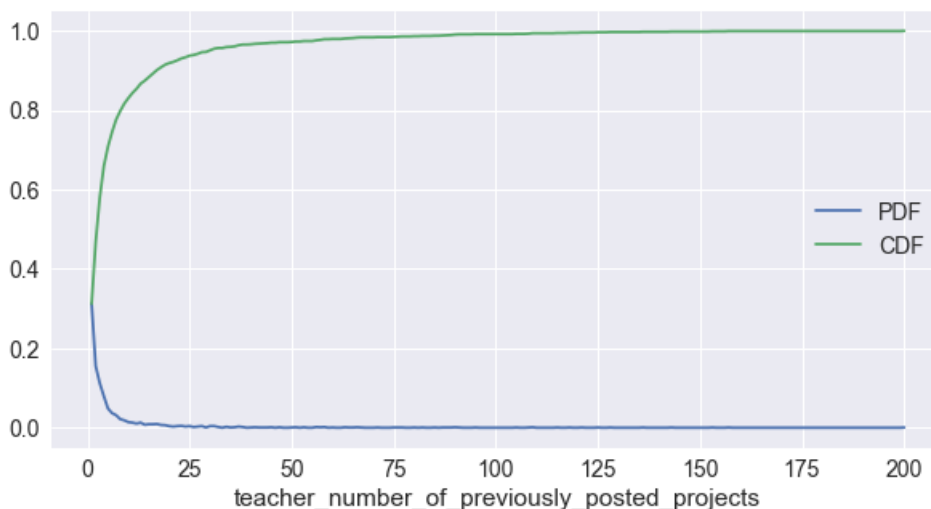
```
<matplotlib.axes. subplots.AxesSubplot at 0x17bf7db5518>
```



PDF (FP ,teacher_number_of_previously_posted_projects)

In [179]:

```
plt.figure(figsize=(10,5))
counts, bin_edges = np.histogram(S_test_falsePos1['teacher_number_of_previously_posted_projects'],
,bins='auto', density=True)
pdf = counts/sum(counts)
cdf = np.cumsum(pdf)
pdf_FP, = plt.plot(bin_edges[1:], pdf)
cdf_FP, = plt.plot(bin_edges[1:], cdf)
plt.legend([pdf_FP, cdf_FP], ["PDF", "CDF"])
plt.xlabel('teacher_number_of_previously_posted_projects')
plt.show()
```



FEATURE SET 4:TFIDF_W2V

In [97]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
S_tfidf_w2v_train=
hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train,teacher_pre
fix_one_hot_train,clean_project_grade_category_one_hot_train,tfidf_w2v_vectors_train,tfidf_w2v_ppt
_train,price_standardized_train,prev_project_standardized_train,quantity_standardized_train,title_w
ord_count_train,essay_word_count_train,essay_sent_pos_train,essay_sent_neg_train,essay_sent_neu_tra
in,essay_sent_comp_train)).tocsr()
print(S_tfidf_w2v_train.shape)

S_tfidf_w2v_test=
hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test,teacher_prefi
x_one_hot_test,clean_project_grade_category_one_hot_test,tfidf_w2v_vectors_test,tfidf_w2v_ppt_test
,price_standardized_test,prev_project_standardized_test,quantity_standardized_test,title_word_count
_test,essay_word_count_test,essay_sent_pos_test,essay_sent_neg_test,essay_sent_neu_test,essay_sent_
comp_test)).tocsr()
print(S_tfidf_w2v_test.shape)
```

```
S_tfidf_w2v_cv= hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv,teacher_prefix_one_hot_cv,clean_project_grade_category_one_hot_cv,tfidf_w2v_vectors_cv,tfidf_w2v_ppt_cv,price_standardized_cv,prev_project_standardized_cv,quantity_standardized_cv,title_word_count_cv,essay_word_count_cv,essay_sent_pos_cv,essay_sent_neg_cv,essay_sent_neu_cv,essay_sent_comp_cv)).to_csr()
print(S_tfidf_w2v_cv.shape)

(51237, 708)
(36052, 708)
(21959, 708)
```

Using CV to find best hyperparameter

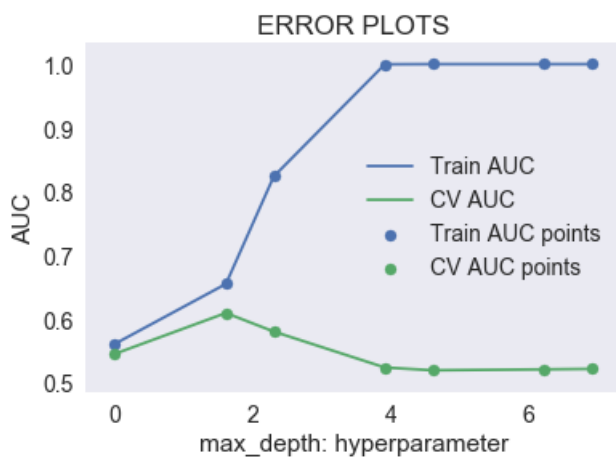
In [232]:

```
from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
train_auc = []
cv_auc = []
a = []
b = []
import math
max_depth=[1, 5, 10, 50, 100, 500, 1000]

for i in tqdm(max_depth):
    dtc= DecisionTreeClassifier(max_depth=i,class_weight="balanced")
    l=dtc.fit(S_tfidf_w2v_train, y_train)
    y_train_pred = batch_predict(dtc,S_tfidf_w2v_train)
    y_cv_pred = batch_predict(dtc, S_tfidf_w2v_cv)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
    a.append(y_train_pred)
    b.append(y_cv_pred)

plt.plot([math.log(i) for i in max_depth],train_auc, label='Train AUC')
plt.plot([math.log(i) for i in max_depth],cv_auc, label='CV AUC')
plt.scatter([math.log(i) for i in max_depth],train_auc, label='Train AUC points')
plt.scatter([math.log(i) for i in max_depth],cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("max_depth: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

100%|██████████| 7/7 [18:58<00:00, 197.95s/it]



In [233]:

```
from sklearn.tree import DecisionTreeClassifier
```

```

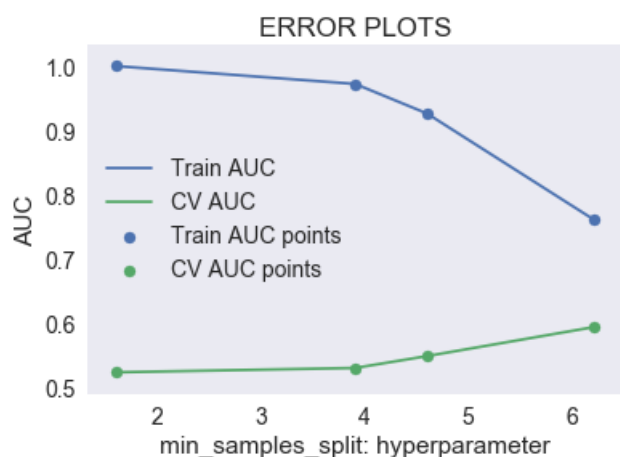
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
train_auc = []
cv_auc = []
a = []
b = []
import math
min_samples_split=[5, 50, 100, 500]

for i in tqdm(min_samples_split):
    dtc= DecisionTreeClassifier(min_samples_split=i,class_weight="balanced")
    l=dtc.fit(S_tfidf_w2v_train, y_train)
    y_train_pred = batch_predict(dtc,S_tfidf_w2v_train)
    y_cv_pred = batch_predict(dtc, S_tfidf_w2v_cv)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
    a.append(y_train_pred)
    b.append(y_cv_pred)

plt.plot([math.log(i) for i in min_samples_split],train_auc, label='Train AUC')
plt.plot([math.log(i) for i in min_samples_split],cv_auc, label='CV AUC')
plt.scatter([math.log(i) for i in min_samples_split],train_auc, label='Train AUC points')
plt.scatter([math.log(i) for i in min_samples_split],cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("min_samples_split: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```

100%|██████████| 4/4 [13:14<00:00, 194.80s/it]



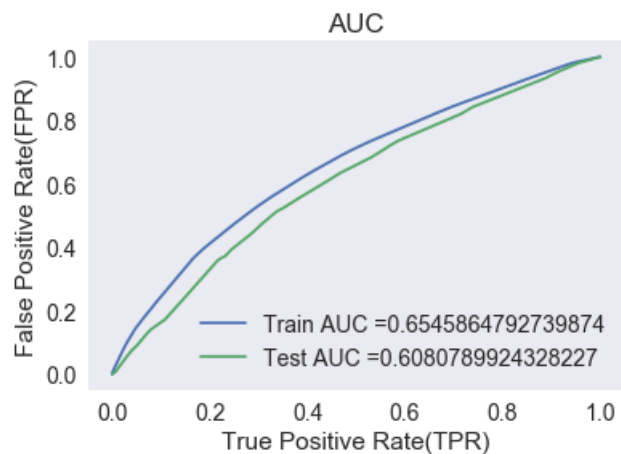
In [180]:

```

#https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
model = DecisionTreeClassifier(max_depth =5,min_samples_split=6,random_state=0, class_weight='balanced')
model.fit(S_tfidf_w2v_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs
y_train_pred = batch_predict(model, S_tfidf_w2v_train)
y_test_pred = batch_predict(model, S_tfidf_w2v_test)
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
plt.plot(train_fpr, train_tpr, label="Train AUC "+str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC "+str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()

```

```
plt.show()
```



confusion matrix for train data

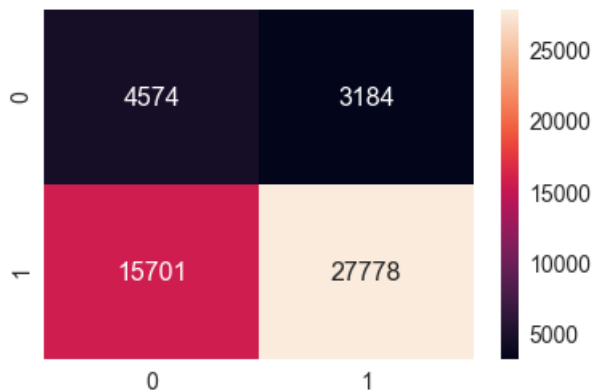
In [109]:

```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train, prediction(y_train_pred, tr_thresholds,
train_fpr, train_tpr)), range(2), range(2))
sns.set(font_scale=1.4) #for label
sns.heatmap(conf_matr_df_train, annot=True, annot_kws={"size": 16}, fmt='g')
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.3785268713076889 for threshold 0.478

Out[109]:

<matplotlib.axes._subplots.AxesSubplot at 0x17bd999e748>



Confusion matrix on test data

In [110]:

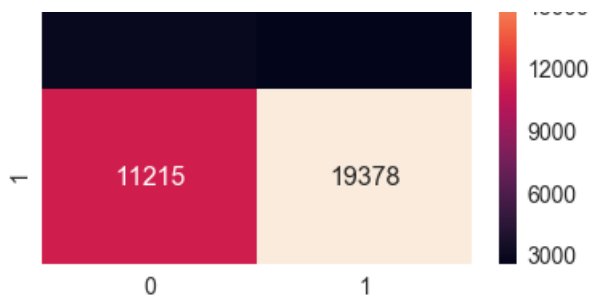
```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test, prediction(y_test_pred, tr_thresholds,
train_fpr, train_tpr)), range(2), range(2))
sns.set(font_scale=1.4) #for label size
sns.heatmap(conf_matr_df_test, annot=True, annot_kws={"size": 16}, fmt='g')
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.3785268713076889 for threshold 0.478

Out[110]:

<matplotlib.axes._subplots.AxesSubplot at 0x17bd961ac18>





finding false positive points

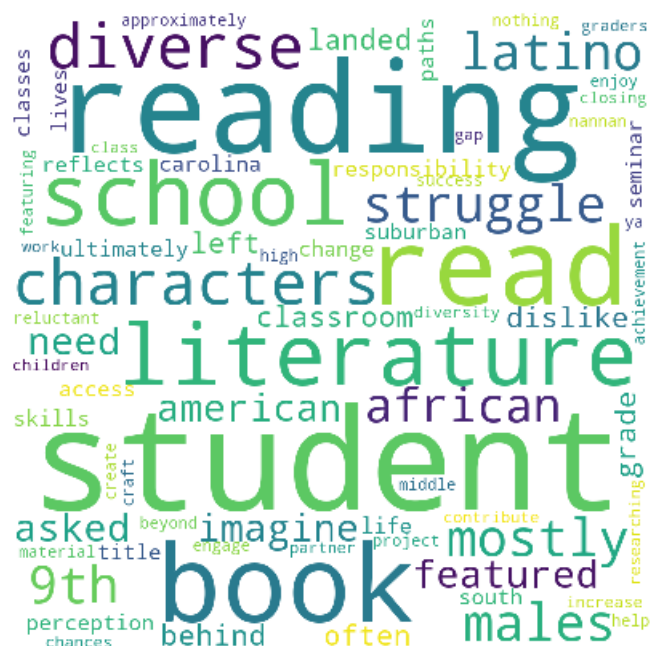
In [181]:

```
fpi = []
for i in range(len(y_test)) :
    if (y_test.values[i] == 0) & (predictions1[i] == 1) :
        fpi.append(i)
fp_essay1 = []
for i in fpi :
    fp_essay1.append(S_test["clean_essays"].values[i])
```

Word cloud

In [182]:

```
from wordcloud import WordCloud, STOPWORDS
comment_words = ' '
stopwords = set(STOPWORDS)
for val in fp_essay1 :
    val = str(val)
    tokens = val.split()
for i in range(len(tokens)):
    tokens[i] = tokens[i].lower()
for words in tokens :
    comment_words = comment_words + words + ' '
wordcloud = WordCloud(width = 800, height = 800, background_color = 'white', stopwords =
stopwords,min_font_size = 10).generate(comment_words)
plt.figure(figsize = (6, 6), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad = 0)
plt.show()
```



In [183]:

```
cols = S_test.columns
S_test_falsePos1 = pd.DataFrame(columns=cols)
# get the data of the false positives
for i in fpi : # (in fpi all the false positives data points indexes)
    S_test_falsePos1 = S_test_falsePos1.append(S_test.filter(items=[i], axis=0))
S_test_falsePos1.head(1)
len(S_test_falsePos1)
```

Out[183]:

1712

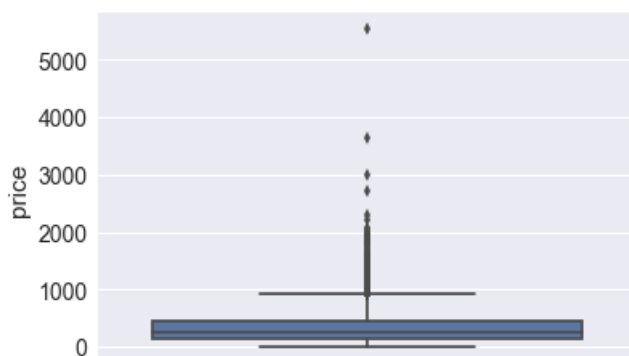
boxplot:

In [184]:

```
sns.boxplot(y='price', data=S_test_falsePos1)
```

Out[184]:

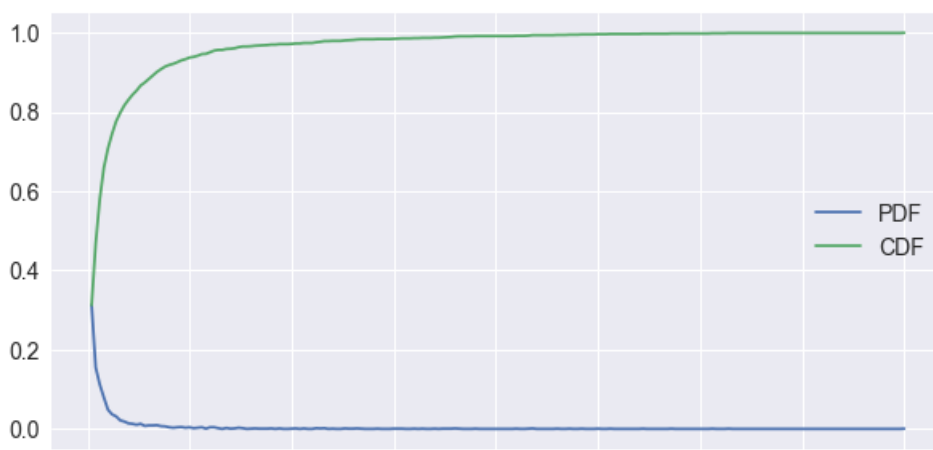
<matplotlib.axes._subplots.AxesSubplot at 0x17bd973f3c8>



PDF (FP ,teacher_number_of_previously_posted_projects)

In [185]:

```
plt.figure(figsize=(10,5))
counts, bin_edges = np.histogram(S_test_falsePos1['teacher_number_of_previously_posted_projects'],
,bins='auto', density=True)
pdf = counts/sum(counts)
cdf = np.cumsum(pdf)
pdf_FP, = plt.plot(bin_edges[1:], pdf)
cdf_FP, = plt.plot(bin_edges[1:], cdf)
plt.legend([pdf_FP, cdf_FP], ["PDF", "CDF"])
plt.xlabel('teacher_number_of_previously_posted_projects')
plt.show()
```



0 25 50 75 100 125 150 175 200
teacher_number_of_previously_posted_projects

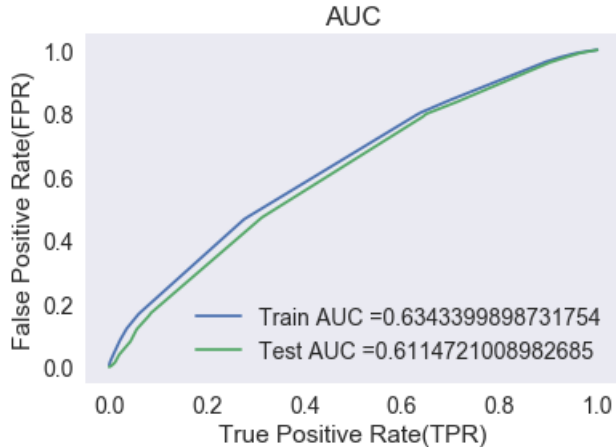
2.5 [Task-2]Getting top 5k features using `feature_importances_`

We have already found max_depth=5 and min_sample_split=5 for TFIDF data earlier so we will use those values as hyperparameters

We will use a decision tree to find 5000 best features

In [186]:

```
#https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
model = DecisionTreeClassifier(max_depth = 5,min_samples_split=5,random_state=0, class_weight='balanced')
model.fit(S_TFIDF_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs
y_train_pred = batch_predict(model, S_TFIDF_train)
y_test_pred = batch_predict(model, S_TFIDF_test)
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



Selecting best 5K features

In [211]:

```
#https://stackoverflow.com/questions/47111434/randomforestregressor-and-feature-importances-error
def selectKImportance(model, X, k=5):
    return X[:,model.feature_importances_.argsort()[::-1][:k]]
```

In [214]:

```
# for tf-idf set 2
S_set5_train = selectKImportance(model,S_TFIDF_train,5000)
S_set5_test = selectKImportance(model,S_TFIDF_test, 5000)
S_set5_cv = selectKImportance(model,S_TFIDF_cv, 5000)
print(S_set5_train.shape)
print(S_set5_test.shape)
print(S_set5_cv.shape)
```

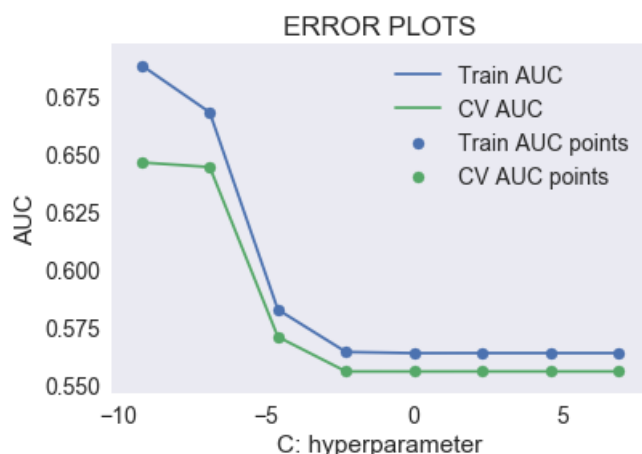
```
(51237, 5000)
(36052, 5000)
(21959, 5000)
```

now training a linear SVM with these 5k features

In [215]:

```
from sklearn.linear_model import SGDClassifier
from sklearn.calibration import CalibratedClassifierCV
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
train_auc = []
cv_auc = []
a = []
b = []
import math
alpha=[10**x for x in range(-4,4)]
for i in tqdm(alpha):
    svm= SGDClassifier(alpha=i,loss='hinge', class_weight='balanced' )
    s=svm.fit(S_set5_train, y_train)
    clfcalibrated=CalibratedClassifierCV(svm,cv='prefit',method='isotonic')
    clfcalibrated.fit(S_set5_cv,y_cv)
    y_train_pred = batch_predict(clfcalibrated,S_set5_train)
    y_cv_pred = batch_predict(clfcalibrated, S_set5_cv)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
    a.append(y_train_pred)
    b.append(y_cv_pred)
plt.plot([math.log(i) for i in alpha],train_auc, label='Train AUC')
plt.plot([math.log(i) for i in alpha],cv_auc, label='CV AUC')
plt.scatter([math.log(i) for i in alpha],train_auc, label='Train AUC points')
plt.scatter([math.log(i) for i in alpha],cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

100%|██████████| 8/8 [00:06<00:00, 1.42it/s]



In [216]:

```
from sklearn.model_selection import GridSearchCV
svm = SGDClassifier(loss='hinge',class_weight='balanced')
alpha_vals=[10**x for x in range(-4,4)]
penalty=['l1','l2']
parameters = {'alpha':alpha_vals,'penalty':penalty}
clf = GridSearchCV(svm, parameters, cv= 10, scoring='roc_auc')
best_model=clf.fit(S_set5_train, y_train)
print('Best alpha:', best_model.best_estimator_.get_params()['alpha'])
```

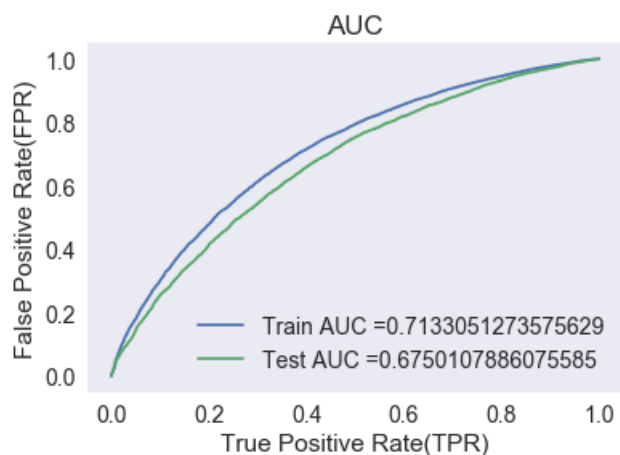
```
print('Best penalty:', best_model.best_estimator_.get_params()['penalty'])
```

Best alpha: 0.0001
Best penalty: l1

we will choose alpha=0.0001 and penalty=l1

In [218]:

```
from sklearn.metrics import roc_curve, auc
svm= SGDClassifier(alpha=0.0001,loss='hinge', penalty='l1', class_weight='balanced', )
s=svm.fit(S_set5_train[0:26237], y_train[0:26237])
clfcalibrated=CalibratedClassifierCV(svm,method='isotonic')
clfcalibrated.fit(S_set5_train[26237:51237],y_train[26237:51237])
y_train_pred = batch_predict(clfcalibrated,S_set5_train)
y_test_pred = batch_predict(clfcalibrated, S_set5_test)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



Confusion matix for train data:

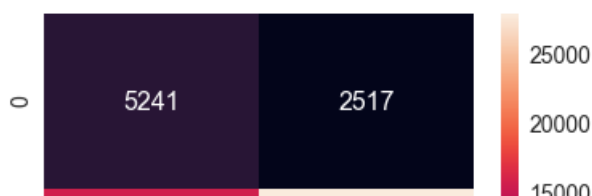
In [219]:

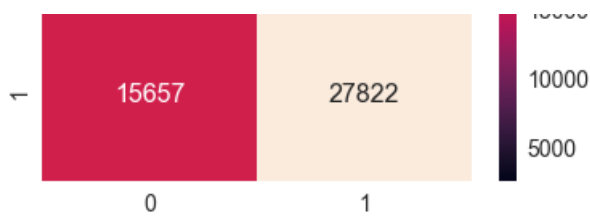
```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train, prediction(y_train_pred, tr_thresholds
,
train_fpr, train_tpr)), range(2),range(2))
sns.set(font_scale=1.4)#for label
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.4349424841680244 for threshold 0.848

Out[219]:

<matplotlib.axes._subplots.AxesSubplot at 0x17bd977e048>





Confusion matrix for test data:

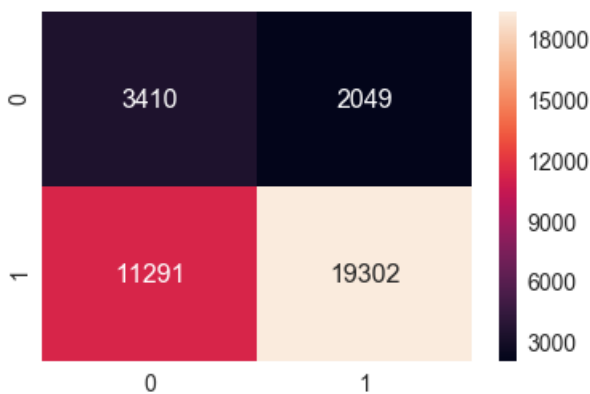
In [220]:

```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test, prediction(y_test_pred, tr_thresholds,
train_fpr, train_tpr)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

the maximum value of tpr*(1-fpr) 0.4349424841680244 for threshold 0.848

Out[220]:

<matplotlib.axes._subplots.AxesSubplot at 0x17bd9bce0f0>



finding false positive points

In [221]:

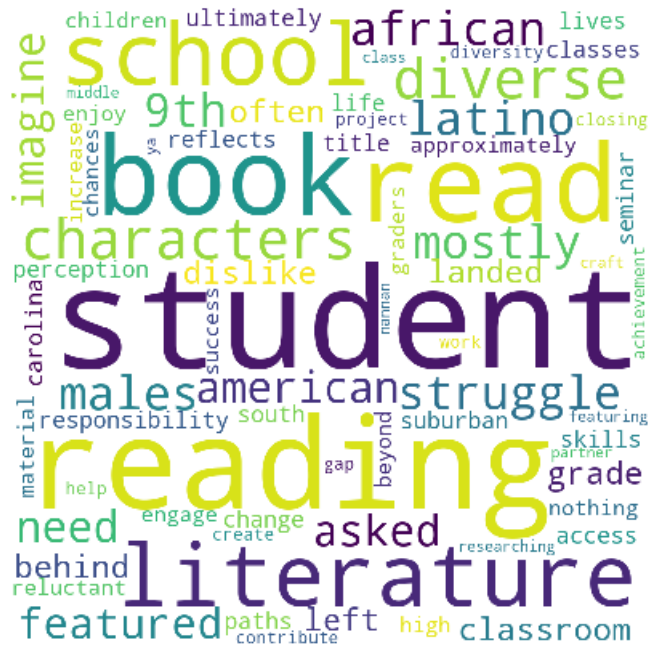
```
fpi = []
for i in range(len(y_test)) :
    if (y_test.values[i] == 0) & (predictions1[i] == 1) :
        fpi.append(i)
fp_essay1 = []
for i in fpi :
    fp_essay1.append(S_test["clean_essays"].values[i])
```

Word cloud

In [222]:

```
from wordcloud import WordCloud, STOPWORDS
comment_words = ' '
stopwords = set(STOPWORDS)
for val in fp_essay1 :
    val = str(val)
    tokens = val.split()
for i in range(len(tokens)):
    tokens[i] = tokens[i].lower()
for words in tokens :
    comment_words = comment_words + words + ' '
wordcloud = WordCloud(width = 800, height = 800, background_color ='white', stopwords =
stopwords,min_font_size = 10).generate(comment_words)
plt.figure(figsize = (6, 6), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
```

```
plt.tight_layout(pad = 0)
plt.show()
```



In [223]:

```
cols = S_test.columns
S_test_falsePos1 = pd.DataFrame(columns=cols)
# get the data of the false positives
for i in fpi : # (in fpi all the false positives data points indexes)
    S_test_falsePos1 = S_test_falsePos1.append(S_test.filter(items=[i], axis=0))
S_test_falsePos1.head(1)
len(S_test_falsePos1)
```

Out [223]:

2049

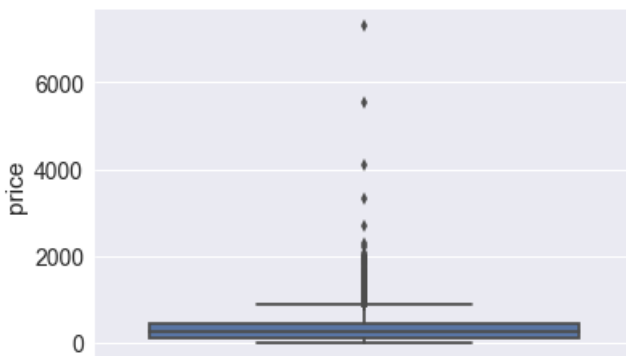
Box plot

In [224]:

```
sns.boxplot(y='price', data=S_test_falsePos1)
```

Out [224] :

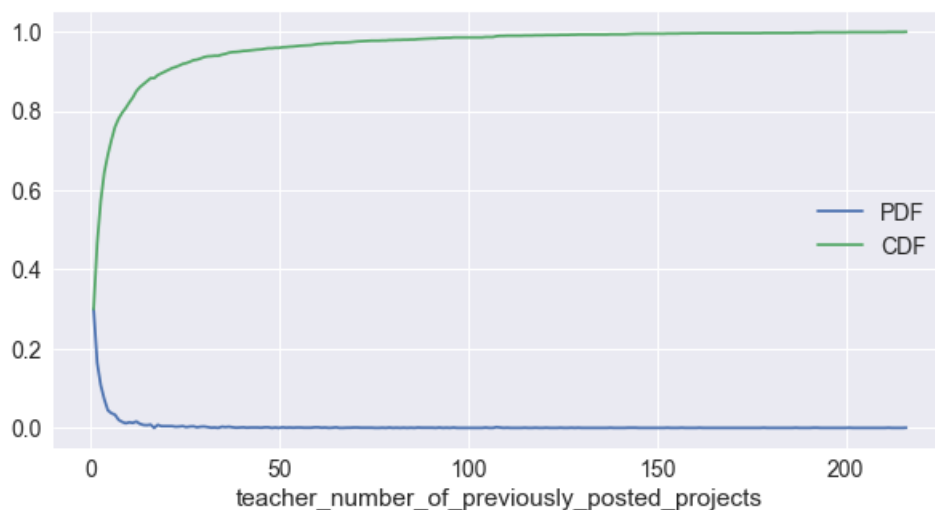
```
<matplotlib.axes._subplots.AxesSubplot at 0x17bf7cc5cc0>
```



PDF (FP ,teacher_number_of_previously_posted_projects)

In [225]:

```
plt.figure(figsize=(10,5))
counts, bin_edges = np.histogram(S_test_falsePos1['teacher_number_of_previously_posted_projects'],
,bins='auto', density=True)
pdf = counts/sum(counts)
cdf = np.cumsum(pdf)
pdf_FP, = plt.plot(bin_edges[1:], pdf)
cdf_FP, = plt.plot(bin_edges[1:], cdf)
plt.legend([pdf_FP, cdf_FP], ["PDF", "CDF"])
plt.xlabel('teacher_number_of_previously_posted_projects')
plt.show()
```



3. Conclusion

In [244]:

```
# Please compare all your models using Prettytable library
from prettytable import PrettyTable
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable
x = PrettyTable()
x.field_names=["Vectorizer","Model","max_depth & AUC","min_split & AUC"]
x.add_row(["BOW","Decision Tree",'5 & 0.6','5 & 0.6'])
x.add_row(["TFIDF","Decision Tree",'5 & 0.59','5 & 0.59'])
x.add_row(["AVG W2V","Decision Tree",'6 & 0.52','7 & 0.58'])
x.add_row(["TFIDF W2V","decision tree",'5 & 0.52','6 & 0.62'])
print(x)
```

Vectorizer	Model	max_depth & AUC	min_split & AUC
BOW	Decision Tree	5 & 0.6	5 & 0.6
TFIDF	Decision Tree	5 & 0.59	5 & 0.59
AVG W2V	Decision Tree	6 & 0.52	7 & 0.58
TFIDF W2V	decision tree	5 & 0.52	6 & 0.62

pretty table for feature set 5:

In [246]:

```
s = PrettyTable()
s.field_names=["Vectorizer","Model","alpha"," AUC","penalty"]
s.add_row(["TFIDF-5K FEATURES","SGD-HINGE LOSS-L1",'ALPHA:0.0001', 'AUC:0.56', '11'])
print(s)
```

Vectorizer	Model	alpha	AUC	penalty
TFIDF-5K FEATURES	SGD-HINGE LOSS-L1	ALPHA:0.0001	AUC:0.56	11

