

# Assignment : DT

In [1]:

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

In [2]:

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

In [3]:

```
1 print("Number of data points in train data", project_data.shape)
2 print('-'*50)
3 print("The attributes of data :", project_data.columns.values)
4 project_data.head(2)
```

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']

Out[3]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL	

In [4]:

```
1 project_data.project_is_approved.value_counts()
```

Out[4]:

```
1    92706
0    16542
Name: project_is_approved, dtype: int64
```

In [5]:

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

Number of data points in train data (1541272, 4)

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

Out[5]:

	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

In [6]:

```
1 project_data['project_subject_categories'] = project_data['project_subject_categories']
2 project_data['project_subject_categories'] = project_data['project_subject_categories']
3 project_data['project_subject_categories'] = project_data['project_subject_categories']
4 project_data['project_subject_categories'] = project_data['project_subject_categories']
5 project_data['project_subject_categories'] = project_data['project_subject_categories']
6 project_data['clean_categories']=project_data['project_subject_categories']
7 project_data.drop(['project_subject_categories'], axis=1, inplace=True)
8 project_data['clean_categories'].value_counts()
```

Out[6]:

literacy_language	23655
math_science	17072
literacy_language_math_science	14636
health_sports	10177
music_arts	5180
specialneeds	4226
literacy_language_specialneeds	3961
appliedlearning	3771
math_science_literacy_language	2289
appliedlearning_literacy_language	2191
history_civics	1851
math_science_specialneeds	1840
literacy_language_music_arts	1757
math_science_music_arts	1642
appliedlearning_specialneeds	1467
history_civics_literacy_language	1421
health_sports_specialneeds	1391
warmth_care_hunger	1309
math_science_appliedlearning	1220
appliedlearning_math_science	1052
literacy_language_history_civics	809
health_sports_literacy_language	803
appliedlearning_music_arts	758
math_science_history_civics	652
literacy_language_appliedlearning	636
appliedlearning_health_sports	608
math_science_health_sports	414
history_civics_math_science	322
history_civics_music_arts	312
specialneeds_music_arts	302
health_sports_math_science	271
history_civics_specialneeds	252
health_sports_appliedlearning	192
appliedlearning_history_civics	178
health_sports_music_arts	155
music_arts_specialneeds	138
literacy_language_health_sports	72
health_sports_history_civics	43
history_civics_appliedlearning	42
specialneeds_health_sports	42
specialneeds_warmth_care_hunger	23
health_sports_warmth_care_hunger	23
music_arts_health_sports	19
music_arts_history_civics	18
history_civics_health_sports	13
math_science_warmth_care_hunger	11
music_arts_appliedlearning	10
appliedlearning_warmth_care_hunger	10

literacy\_language\_warmth\_care\_hunger9

music\_arts\_warmth\_care\_hunger2

history\_civics\_warmth\_care\_hunger1

Name: clean\_categories, dtype: int64

In [7]:

```
1 project_data.head(1)
```

Out[7]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_s
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	:

In [8]:

```
1 project_data['project_subject_subcategories'] = project_data['project_subject_subcateg
2 project_data['project_subject_subcategories'] = project_data['project_subject_subcateg
3 project_data['project_subject_subcategories'] = project_data['project_subject_subcateg
4 project_data['project_subject_subcategories'] = project_data['project_subject_subcateg
5 project_data['project_subject_subcategories'] = project_data['project_subject_subcateg
6 project_data['clean_subcategories']=project_data['project_subject_subcategories']
7 project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
8 project_data['clean_subcategories'].value_counts()
```

Out[8]:

literacy9486

literacy\_mathematics8325

literature\_writing\_mathematics5923

literacy\_literature\_writing5571

mathematics5379

...

civics\_government\_nutritioneducation1

civics\_government\_foreignlanguages1

economics\_foreignlanguages1

parentinvolvement\_warmth\_care\_hunger1

esl\_economics1

Name: clean\_subcategories, Length: 401, dtype: int64

In [9]:

```
1 project_data.head(1)
```

Out[9]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_s
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	:

In [10]:

```
1 project_data['teacher_prefix'].value_counts()
```

Out[10]:

```
Mrs.      57269
Ms.       38955
Mr.       10648
Teacher   2360
Dr.        13
Name: teacher_prefix, dtype: int64
```

In [11]:

```
1 # check if we have any nan values are there
2 print(project_data['teacher_prefix'].isnull().values.any())
3 print("number of nan values",project_data['teacher_prefix'].isnull().values.sum())
```

```
True
number of nan values 3
```

In [12]:

```
1 project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('Mrs.')
```

In [13]:

```
1 # Remove '.'
2 # convert all the chars to small
3 project_data['teacher_prefix'] = project_data['teacher_prefix'].str.replace('.', '')
4 project_data['teacher_prefix'] = project_data['teacher_prefix'].str.lower()
5 project_data['teacher_prefix'].value_counts()
```

Out[13]:

```
mrs      57272
ms       38955
mr       10648
teacher   2360
dr        13
Name: teacher_prefix, dtype: int64
```

In [14]:

```
1 # We need to get rid of The spaces between the text and the hyphens because they're spe
2 #Rmoving multiple characters from a string in Python
3 #https://stackoverflow.com/questions/3411771/multiple-character-replace-with-python
4
5 project_grade_category = []
6
7 for i in range(len(project_data)):
8     a = project_data["project_grade_category"][i].replace(" ", "_").replace("-", "_")
9     project_grade_category.append(a)
```

In [15]:

```
1 project_data.drop(['project_grade_category'], axis = 1, inplace = True)
2 project_data["project_grade_category"] = project_grade_category
3 print("After removing the special characters ,Column values: ")
4 np.unique(project_data["project_grade_category"].values)
```

After removing the special characters ,Column values:

Out[15]:

```
array(['Grades_3_5', 'Grades_6_8', 'Grades_9_12', 'Grades_PreK_2'],
      dtype=object)
```

In [16]:

```
1 # convert all of them into small letters
2 project_data['school_state'] = project_data['school_state'].str.lower()
3 project_data['school_state'].value_counts()
```

Out[16]:

ca	15388
tx	7396
ny	7318
fl	6185
nc	5091
il	4350
ga	3963
sc	3936
mi	3161
pa	3109
in	2620
mo	2576
oh	2467
la	2394
ma	2389
wa	2334
ok	2276
nj	2237
az	2147
va	2045
wi	1827
al	1762
ut	1731
tn	1688
ct	1663
md	1514
nv	1367
ms	1323
ky	1304
or	1242
mn	1208
co	1111
ar	1049
id	693
ia	666
ks	634
nm	557
dc	516
hi	507
me	505
wv	503
nh	348
ak	345
de	343
ne	309
sd	300
ri	285
mt	245
nd	143
wy	98
vt	80

Name: school\_state, dtype: int64

# Preprocessing Categorical Features:

In [17]:

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

## Stopword Removal

In [18]:

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

## Preprocessing Categorical Features: essay



In [19]:

```
1 # Combining all the above students
2 from tqdm import tqdm
3 def preprocess_text(text_data):
4     preprocessed_text = []
5     # tqdm is for printing the status bar
6     for sentance in tqdm(text_data):
7         sent = decontracted(sentance)
8         sent = sent.replace('\\r', ' ')
9         sent = sent.replace('\\n', ' ')
10        sent = sent.replace('\\\"', ' ')
11        sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12        # https://gist.github.com/sebleier/554280
13        sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
14        preprocessed_text.append(sent.lower().strip())
15    return preprocessed_text
```

In [20]:

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

In [21]:

```
1 print("printing some random essay before Preprocessing ")
2 print(9, project_data['essay'].values[9])
3 print('-'*50)
4 print(34, project_data['essay'].values[34])
5 print('-'*50)
6 print(147, project_data['essay'].values[147])
```

printing some random essay before Preprocessing

9 Over 95% of my students are on free or reduced lunch. I have a few who are homeless, but despite that, they come to school with an eagerness to learn. My students are inquisitive eager learners who embrace the challenge of not having great books and other resources every day. Many of them are not afforded the opportunity to engage with these big colorful pages of a book on a regular basis at home and they don't travel to the public library. \r\nIt is my duty as a teacher to do all I can to provide each student an opportunity to succeed in every aspect of life. \r\nReading is Fundamental! My students will read these books over and over again while boosting their comprehension skills. These books will be used for read alouds, partner reading and for Independent reading. \r\nThey will engage in reading to build their "Love for Reading" by reading for pure enjoyment. They will be introduced to some new authors as well as some old favorites. I want my students to be ready for the 21st Century and know the pleasure of holding a good hard back book in hand. There's nothing like a good book to read! \r\nMy students will soar in Reading, and more because of your consideration and generous funding contribution. This will help build stamina and prepare for 3rd grade. Thank you so much for reading our proposal!nannan

-----

34 My students mainly come from extremely low-income families, and the majority of them come from homes where both parents work full time. Most of my students are at school from 7:30 am to 6:00 pm (2:30 to 6:00 pm in the after-school program), and they all receive free and reduced meals for breakfast and lunch. \r\n\r\n\r\nI want my students to feel as comfortable in my classroom as they do at home. Many of my students take on multiple roles both at home as well as in school. They are sometimes the caretakers of younger siblings, cooks, babysitters, academics, friends, and most of all, they are developing who they are going to become as adults. I consider it an essential part of my job to model helping others gain knowledge in a positive manner. As a result, I have a community of students who love helping each other in and outside of the classroom. They consistently look for opportunities to support each other's learning in a kind and helpful way. I am excited to be experimenting with alternative seating in my classroom this school year. Studies have shown that giving students the option of where they sit in a classroom increases focus as well as motivation. \r\n\r\n\r\nBy allowing students choice in the classroom, they are able to explore and create in a welcoming environment. Alternative classroom seating has been experimented with more frequently in recent years. I believe (along with many others), that every child learns differently. This does not only apply to how multiplication is memorized, or a paper is written, but applies to the space in which they are asked to work. I have had students in the past ask "Can I work in the library? Can I work on the carpet?" My answer was always, "As long as you're learning, you can work wherever you want!" \r\n\r\n\r\nWith the yoga balls and the lap-desks, I will be able to increase the options for seating in my classroom and expand its imaginable space.nannan

-----

147 My students are eager to learn and make their mark on the world.\r\n\r\nThey come from a Title 1 school and need extra love.\r\n\r\n\r\nMy fourth grade students are in a high poverty area and still come to school every day

In [22]:

```
1 preprocessed_essays = preprocess_text(project_data['essay'].values)
```

```
100%|███████████| 109248/109248 [04:41<00:00, 388.54it/s]
```

In [23]:

```
1 print("printing some random essay After preprocessed_essays")
2 print(9, preprocessed_essays[9])
3 print('-'*50)
4 print(34, preprocessed_essays[34])
5 print('-'*50)
6 print(147, preprocessed_essays[147])
```

printing some random essay After preprocessed\_essays

9 95 students free reduced lunch homeless despite come school eagerness learn students inquisitive eager learners embrace challenge not great books resources every day many not afforded opportunity engage big colorful pages book regular basis home not travel public library duty teacher provide student opportunity succeed every aspect life reading fundamental students read books boosting comprehension skills books used read alouds partner reading independent reading engage reading build love reading reading pure enjoyment introduced new authors well old favorites want students ready 21st century know pleasure holding good hard back book hand nothing like good book read students soar reading consideration generous funding contribution help build stamina prepare 3rd grade thank much reading proposal nannan

-----  
34 students mainly come extremely low income families majority come homes parents work full time students school 7 30 6 00 pm 2 30 6 00 pm school program receive free reduced meals breakfast lunch want students feel comfortable classroom home many students take multiple roles home well school sometimes caretakers younger siblings cooks babysitters academics friends developing going become adults consider essential part job model helping others gain knowledge positive manner result community students love helping outside classroom consistently look opportunities support learning kind helpful way excited experimenting alternative seating classroom school year studies shown giving students option sit classroom increases focus well motivation allowing students choice classroom able explore create welcoming environment alternative classroom seating experimented frequently recent years believe along many others every child learns differently not apply multiplication memorized paper written applies space asked work students past ask work library work carpenter answer always long learning work wherever want yoga balls lap desks able increase options seating classroom expand imaginable space nannan

-----  
147 students eager learn make mark world come title 1 school need extra love fourth grade students high poverty area still come school every day get education trying make fun educational get schooling created caring environment students bloom deserve best thank requesting 1 chromebook access online interventions differentiate instruction get extra practice chromebook used supplement ela math instruction students play ela math games engaging fun well participate assignments online turn help students improve skills chromebook classroom would not allow students use programs pace would ensure students getting adequate time use programs online programs especially beneficial students special needs able work level well challenged different materials making students confident abilities chromebook would allow students daily access computers increase computing skills change lives better become successful school access technology classroom would help bridge achievement gap nannan

In [24]:

```
1 #creating a new column with the preprocessed essays and replacing it with the original
2 project_data['preprocessed_essays'] = preprocessed_essays
3 project_data.drop(['project_essay_1'], axis=1, inplace=True)
4 project_data.drop(['project_essay_2'], axis=1, inplace=True)
5 project_data.drop(['project_essay_3'], axis=1, inplace=True)
6 project_data.drop(['project_essay_4'], axis=1, inplace=True)
```

In [25]:

```
1 import nltk
2 nltk.download('vader_lexicon')
```

```
[nltk_data] Downloading package vader_lexicon to
[nltk_data] C:\Users\imran\AppData\Roaming\nltk_data...
[nltk_data] Package vader_lexicon is already up-to-date!
```

Out[25]:

True

In [26]:

```
1 import nltk
2 from nltk.sentiment.vader import SentimentIntensityAnalyzer
3
4 sid = SentimentIntensityAnalyzer()
5
6 for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest
7 for learning my students learn in many different ways using all of our senses and multi
8 of techniques to help all my students succeed students in my class come from a variety
9 for wonderful sharing of experiences and cultures including native americans our school
10 learners which can be seen through collaborative student project based learning in and
11 in my class love to work with hands on materials and have many different opportunities
12 mastered having the social skills to work cooperatively with friends is a crucial aspect
13 montana is the perfect place to learn about agriculture and nutrition my students love
14 in the early childhood classroom i have had several kids ask me can we try cooking with
15 and create common core cooking lessons where we learn important math and writing concep
16 food for snack time my students will have a grounded appreciation for the work that wer
17 of where the ingredients came from as well as how it is healthy for their bodies this p
18 nutrition and agricultural cooking recipes by having us peel our own apples to make hon
19 and mix up healthy plants from our classroom garden in the spring we will also create c
20 shared with families students will gain math and literature skills as well as a life le
21 nannan'
22 ss = sid.polarity_scores(for_sentiment)
23
24 for k in ss:
25     print('{0}: {1}, '.format(k, ss[k]), end='')
26
27 # we can use these 4 things as features/attributes (neg, neu, pos, compound)
28 # neg: 0.0, neu: 0.753, pos: 0.247, compound: 0.93
```

neg: 0.01, neu: 0.745, pos: 0.245, compound: 0.9975,

In [27]:

```
1 # Sentiment Analysis on 'essay'
2 sid = SentimentIntensityAnalyzer()
3 negative_sentiments = []
4 positive_sentiments = []
5 neutral_sentiments = []
6 compound_sentiments = []
```

In [28]:

```
1 for i in tqdm(project_data['preprocessed_essays']):
2     sid_sentiments = sid.polarity_scores(i)
3     negative_sentiments.append(sid_sentiments['neg'])
4     positive_sentiments.append(sid_sentiments['pos'])
5     neutral_sentiments.append(sid_sentiments['neu'])
6     compound_sentiments.append(sid_sentiments['compound'])
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
████████████████████████████████████████████████████████████████████████████████ 109248/109248 [13:34<00:00, 134.20it/s]
```

In [29]:

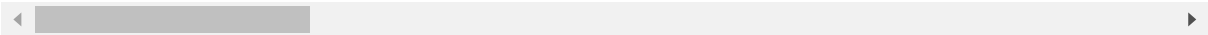
```
1 # Now append these sentiments columns/features to original preprocessed dataframe
2 project_data['negative_sent'] = negative_sentiments
3 project_data['positive_sent'] = positive_sentiments
4 project_data['neutral_sent'] = neutral_sentiments
5 project_data['compound_sent'] = compound_sentiments
```

In [30]:

```
1 project_data.head()
```

Out[30]:

Unnamed: 0		id	teacher_id	teacher_prefix	school_state	project
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	mrs	in	
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	mr	fl	
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	ms	az	
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	mrs	ky	
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	mrs	tx	



# Preprocessing Categorical Features: project\_title

In [31]:

```
1 # Combining all the above students
2 from tqdm import tqdm
3 def preprocess_text(text_data):
4     preprocessed_text = []
5     # tqdm is for printing the status bar
6     for sentence in tqdm(text_data):
7         sent = decontracted(sentence)
8         sent = sent.replace('\\r', ' ')
9         sent = sent.replace('\\n', ' ')
10        sent = sent.replace('\\\"', ' ')
11        sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12        # https://gist.github.com/sebleier/554280
13        sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
14        preprocessed_text.append(sent.lower().strip())
15    return preprocessed_text
```

In [32]:

```
1 print("printing some random reviews Before preprocessed_titles")
2 print(9, project_data['project_title'].values[9])
3 print(34, project_data['project_title'].values[34])
4 print(147, project_data['project_title'].values[147])
```

```
printing some random reviews Before preprocessed_titles
9 Just For the Love of Reading--\r\nPure Pleasure
34 \"Have A Ball!!!\"
147 Who needs a Chromebook?\r\nWE DO!!
```

In [33]:

```
1 preprocessed_titles = preprocess_text(project_data['project_title'].values)
```

```
100%|████████████████████████████████████████████████████████████████████████████████|
109248/109248 [00:14<00:00, 7700.81it/s]
```

In [34]:

```
1 print("printing some random reviews After preprocessed_titles ")
2 print(9, preprocessed_titles[9])
3 print(34, preprocessed_titles[34])
4 print(147, preprocessed_titles[147])
```

```
printing some random reviews After preprocessed_titles
9 love reading pure pleasure
34 ball
147 needs chromebook
```



In [35]:

```
1 project_data['project_title'].head(5)
```

Out[35]:

```
0    Educational Support for English Learners at Home
1                Wanted: Projector for Hungry Learners
2    Soccer Equipment for AWESOME Middle School Stu...
3                Techie Kindergarteners
4                Interactive Math Tools
Name: project_title, dtype: object
```

In [36]:

```
1 #creating a new column with the preprocessed titles,useful for analysis
2 project_data['preprocessed_titles'] = preprocessed_titles
```

In [37]:

```
1 project_data.head(1)
```

Out[37]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_s
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	mrs	in	4

## Merging price with project\_data

In [38]:

```
1 price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_i
2 project_data = pd.merge(project_data, price_data, on='id', how='left')
3 project_data.head(1)
```

Out[38]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_s
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	mrs	in	4

1 rows × 22 columns

In [39]:

```
1 # train test split using sklearn.model selection
2 from sklearn.model_selection import train_test_split
3 # Splitting data into Train and cross validation(or test): Stratified
4 X_train, X_test, y_train, y_test = train_test_split(project_data, project_data['project_is_approved'], stratify=project_data['project_is_approved'], random_state=42)
5 print(X_train.shape, y_train.shape)
6 print(X_test.shape, y_test.shape)
```

```
(73196, 22) (73196,)
(36052, 22) (36052,)
```

In [40]:

```
1 X_train.drop(['project_is_approved'], axis=1, inplace=True)
2 X_test.drop(['project_is_approved'], axis=1, inplace=True)
```

In [41]:

```
1 print('X_train_shape' , X_train.shape)
```

```
X_train_shape (73196, 21)
```

## Vectorization ( convert text to word count vectors with CountVectorizer ) of below Categorical features:

- teacher\_prefix
- project\_grade\_category
- school\_state
- project\_subject\_categories
- project\_subject\_subcategories

In [42]:

```
1 ### Vectorizing Categorical data: teacher_prefix
2
3 vectorizer_teacher_prefix = CountVectorizer(lowercase=False, binary=True)
4
5 train_vectorized_ohe_teacher_prefix = vectorizer_teacher_prefix.fit_transform(X_train['teacher_prefix'])
6
7 test_df_vectorized_ohe_teacher_prefix = vectorizer_teacher_prefix.transform(X_test['teacher_prefix'])
8 print("After vectorizations")
9 print("Shape of matrix of Train data after one hot encoding", train_vectorized_ohe_teacher_prefix.shape)
10 print("Shape of matrix of Test data after one hot encoding", test_df_vectorized_ohe_teacher_prefix.shape)
11 print(vectorizer_teacher_prefix.get_feature_names())
```

After vectorizations

```
Shape of matrix of Train data after one hot encoding (73196, 5) (73196,)
```

```
Shape of matrix of Test data after one hot encoding (36052, 5) (36052,)
```

```
['dr', 'mr', 'mrs', 'ms', 'teacher']
```

In [43]:

```
1 # Vectorizing Categorical data: project_subject_categories
2
3 vectorizer_clean_categories = CountVectorizer(lowercase=False, binary=True)
4
5 train_vectorized_ohe_clean_categories = vectorizer_clean_categories.fit_transform(X_train['project_subject_categories'])
6
7
8 test_df_vectorized_ohe_clean_categories = vectorizer_clean_categories.transform(X_test['project_subject_categories'])
9
10 print("Shape of matrix of Train data after one hot encoding ", train_vectorized_ohe_clean_categories.shape)
11 print("Shape of matrix of Test data after one hot encoding ", test_df_vectorized_ohe_clean_categories.shape)
```

Shape of matrix of Train data after one hot encoding (73196, 51)

Shape of matrix of Test data after one hot encoding (36052, 51)

In [44]:

```
1 ### Vectorizing Categorical data: project_subject_subcategories
2
3 vectorizer_clean_subcategories = CountVectorizer(lowercase=False, binary=True)
4
5 train_vectorized_ohe_clean_subcategories = vectorizer_clean_subcategories.fit_transform(X_train['project_subject_subcategories'])
6
7 test_df_vectorized_ohe_clean_subcategories = vectorizer_clean_subcategories.transform(X_test['project_subject_subcategories'])
8
9 print("Shape of matrix of Train data after one hot encoding ", train_vectorized_ohe_clean_subcategories.shape)
10 print("Shape of matrix of Test data after one hot encoding ", test_df_vectorized_ohe_clean_subcategories.shape)
```

Shape of matrix of Train data after one hot encoding (73196, 393)

Shape of matrix of Test data after one hot encoding (36052, 393)

In [45]:

```
1 ## we use count vectorizer to convert the values into one hot encoded features
2
3 vectorizer_grade = CountVectorizer(lowercase=False, binary=True)
4
5 project_grade_categories_one_hot_train = vectorizer_grade.fit_transform(X_train['project_grade_categories'])
6
7 project_grade_categories_one_hot_test = vectorizer_grade.transform(X_test['project_grade_categories'])
8
9 print(vectorizer_grade.get_feature_names())
10 print("Shape of matrix of Train data after one hot encoding ", project_grade_categories_one_hot_train.shape)
11 print("Shape of matrix of Test data after one hot encoding ", project_grade_categories_one_hot_test.shape)
```

['Grades\_3\_5', 'Grades\_6\_8', 'Grades\_9\_12', 'Grades\_PreK\_2']

Shape of matrix of Train data after one hot encoding (73196, 4)

Shape of matrix of Test data after one hot encoding (36052, 4)

In [46]:

```
1  ### Vectorizing Categorical data: school_state
2
3  vectorizer_school_state = CountVectorizer(lowercase=False, binary=True)
4
5  train_vectorized_ohe_school_state = vectorizer_school_state.fit_transform(X_train['school_state'])
6
7  test_df_vectorized_ohe_school_state = vectorizer_school_state.transform(X_test['school_state'])
8
9  print("Shape of matrix of Train data after one hot encoding", train_vectorized_ohe_school_state.shape)
10 print("Shape of matrix of Test data after one hot encoding", test_df_vectorized_ohe_school_state.shape)
```

Shape of matrix of Train data after one hot encoding (73196, 51)

Shape of matrix of Test data after one hot encoding (36052, 51)

## Vectorizing Numerical features:

Various numerical features are :

1.Price

2.Number of Projects previously proposed by Teacher

In [47]:

```
1  from sklearn.preprocessing import Normalizer
2  # Our first Numerical feature - 'price'
3  normalizer = Normalizer()
4
5  # As explainEd above first I will reshape(1, -1)
6  normalizer.fit(X_train['price'].values.reshape(1, -1))
7
8  train_normalized_price = normalizer.transform(X_train['price'].values.reshape(1, -1))
9
10 test_df_normalized_price = normalizer.transform(X_test['price'].values.reshape(1, -1))
11
12 # After normalization reshape again to (-1, 1) i.e. this time unknown rows (i.e. leaving 0s)
13 train_normalized_price = train_normalized_price.reshape(-1, 1)
14
15
16 test_df_normalized_price = test_df_normalized_price.reshape(-1, 1)
17 print(train_normalized_price.shape)
18 print(test_df_normalized_price.shape)
```

(73196, 1)

(36052, 1)

In [48]:

```
1  ### Normalizing next numerical feature: teacher_number_of_previously_posted_projects
2
3  # Second Numerical feature - 'teacher_number_of_previously_posted_projects'
4  normalizer = Normalizer()
5
6  normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
7
8  train_normalized_teacher_number_of_previously_posted_projects = normalizer.transform(X_train)
9
10 test_df_normalized_teacher_number_of_previously_posted_projects = normalizer.transform(X_test)
11
12 # After normalization reshape again to (-1, 1) i.e. this time unknown rows (i.e. leaving out the last row)
13 train_normalized_teacher_number_of_previously_posted_projects = train_normalized_teacher_number_of_previously_posted_projects.reshape(-1, 1)
14
15 test_df_normalized_teacher_number_of_previously_posted_projects = test_df_normalized_teacher_number_of_previously_posted_projects.reshape(-1, 1)
16 print(train_normalized_teacher_number_of_previously_posted_projects.shape)
17
18 print(test_df_normalized_teacher_number_of_previously_posted_projects.shape)
```

(73196, 1)

(36052, 1)

In [49]:

```
1  ### Normalizing next numerical feature: negative_sent
2
3  # Second Numerical feature - 'negative_sent'
4  normalizer = Normalizer()
5
6  normalizer.fit(X_train['negative_sent'].values.reshape(1, -1))
7
8  train_normalized_negative_sent = normalizer.transform(X_train['negative_sent'].values.reshape(-1, 1))
9
10 test_df_normalized_negative_sent = normalizer.transform(X_test['negative_sent'].values.reshape(-1, 1))
11
12 # After normalization reshape again to (-1, 1) i.e. this time unknown rows (i.e. leaving out the last row)
13 train_normalized_negative_sent = train_normalized_negative_sent.reshape(-1, 1)
14
15 test_df_normalized_negative_sent = test_df_normalized_negative_sent.reshape(-1, 1)
16 print(train_normalized_negative_sent.shape)
17
18 print(test_df_normalized_negative_sent.shape)
```

(73196, 1)

(36052, 1)

In [50]:

```
1  ### Normalizing next numerical feature: negative_sent
2
3  # Second Numerical feature - 'positive_sent'
4  normalizer = Normalizer()
5
6  normalizer.fit(X_train['positive_sent'].values.reshape(1, -1))
7
8  train_normalized_positive_sent= normalizer.transform(X_train['positive_sent'].values.re
9
10 test_df_normalized_positive_sent = normalizer.transform(X_test['positive_sent'].values.
11
12 # After normalization reshape again to (-1, 1) i.e. this time unknown rows (i.e. leavin
13 train_normalized_positive_sent = train_normalized_positive_sent.reshape(-1, 1)
14
15 test_df_normalized_positive_sent = test_df_normalized_positive_sent.reshape(-1, 1)
16 print(train_normalized_positive_sent.shape)
17
18 print(test_df_normalized_positive_sent.shape)
```

(73196, 1)

(36052, 1)

In [51]:

```
1  ### Normalizing next numerical feature: negative_sent
2
3  # Second Numerical feature - 'neutral_sent'
4  normalizer = Normalizer()
5
6  normalizer.fit(X_train['neutral_sent'].values.reshape(1, -1))
7
8  train_normalized_neutral_sent= normalizer.transform(X_train['neutral_sent'].values.rest
9
10 test_df_normalized_neutral_sent = normalizer.transform(X_test['neutral_sent'].values.re
11
12 # After normalization reshape again to (-1, 1) i.e. this time unknown rows (i.e. leavin
13 train_normalized_neutral_sent = train_normalized_neutral_sent.reshape(-1, 1)
14
15 test_df_normalized_neutral_sent = test_df_normalized_neutral_sent.reshape(-1, 1)
16 print(train_normalized_neutral_sent.shape)
17
18 print(test_df_normalized_neutral_sent.shape)
```

(73196, 1)

(36052, 1)

In [52]:

```
1  ### Normalizing next numerical feature: negative_sent
2
3  # Second Numerical feature - 'compound_sent'
4  normalizer = Normalizer()
5
6  normalizer.fit(X_train['compound_sent'].values.reshape(1, -1))
7
8  train_normalized_compound_sent= normalizer.transform(X_train['compound_sent'].values.re
9
10 test_df_normalized_compound_sent = normalizer.transform(X_test['compound_sent'].values.
11
12 # After normalization reshape again to (-1, 1) i.e. this time unknown rows (i.e. leavin
13 train_normalized_compound_sent = train_normalized_compound_sent.reshape(-1, 1)
14
15 test_df_normalized_compound_sent = test_df_normalized_compound_sent.reshape(-1, 1)
16 print(train_normalized_compound_sent.shape)
17
18 print(test_df_normalized_compound_sent.shape)
```

(73196, 1)

(36052, 1)

## preprocessed\_eassay (TFIDF)

In [53]:

```
1  ## Encoding Essay column using tfidf
2  from sklearn.feature_extraction.text import TfidfVectorizer
3
4  vectorizer_essay_tfidf = TfidfVectorizer(min_df=10)
5
6  train_vectorized_tfidf_essay = vectorizer_essay_tfidf.fit_transform(X_train['essay'])
7
8  test_df_vectorized_tfidf_essay = vectorizer_essay_tfidf.transform(X_test['essay'])
9
10 print(train_vectorized_tfidf_essay.shape)
11 print(test_df_vectorized_tfidf_essay.shape)
```

(73196, 14784)

(36052, 14784)

## Merging all categorical, text, numerical vectors and sentiment score based on tfidf

In [54]:

```
1 from scipy.sparse import hstack
2 X_train_hstacked_all_tfidf_features_vectorized = hstack((train_vectorized_ohe_clean_cat
3
4 print('X_train_hstacked_all_tfidf_features_vectorized.shape is ', X_train_hstacked_all_
5
6 test_df_hstacked_all_tfidf_features_vectorized = hstack((test_df_vectorized_ohe_clean_c
7
8 print('test_df_hstacked_all_tfidf_features_vectorized.shape is ', test_df_hstacked_all_
```

```
X_train_hstacked_all_tfidf_features_vectorized.shape is (73196, 15294)
test_df_hstacked_all_tfidf_features_vectorized.shape is (36052, 15294)
```

## Applying Decision Tree on TFIDF(Set1):

**Set 1: categorical, numerical features + preprocessed\_eassay (TFIDF)+Sentiment scores(preprocessed\_essay)**

In [55]:

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.model_selection import GridSearchCV
4 import seaborn as sea
```

In [56]:

```
1 params = {'max_depth': [1, 5, 10, 50], 'min_samples_split': [5, 10, 100, 500]}#as given
2 clf = GridSearchCV(DecisionTreeClassifier(random_state=42), params, cv=3, scoring='roc_
3 clf.fit(X_train_hstacked_all_tfidf_features_vectorized, y_train)
```

Fitting 3 folds for each of 16 candidates, totalling 48 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 48 out of 48 | elapsed: 23.7min finished
```

Out[56]:

```
GridSearchCV(cv=3, estimator=DecisionTreeClassifier(random_state=42), n_jobs
=-1,
            param_grid={'max_depth': [1, 5, 10, 50],
                        'min_samples_split': [5, 10, 100, 500]},
            return_train_score=True, scoring='roc_auc', verbose=1)
```

In [58]:

```
1 clf.best_estimator_
```

Out[58]:

```
DecisionTreeClassifier(max_depth=10, min_samples_split=500, random_state=42)
```



In [61]:

```

1 from sklearn.metrics import roc_auc_score
2 max_depth = [1,5,10,50]
3 min_samples_split = [5,10,100,500]
4
5 def get_auc_matrix(X_train, X_test, y_train, y_test ):
6
7     train_auc_final_arr = []
8     test_auc_final_arr = []
9
10    for depth in tqdm(max_depth):
11        train_auc_batch = []
12        test_auc_batch = []
13
14        for split in min_samples_split:
15            df_clf = DecisionTreeClassifier(class_weight='balanced',min_samples_split=split)
16
17            df_clf.fit(X_train, y_train)
18
19            # I have to predict probabilities (clf.predict_proba) instead of classes for
20            y_train_predicted = df_clf.predict_proba(X_train)[:,-1]
21            y_test_predicted = df_clf.predict_proba(X_test)[:,-1]
22
23            train_auc = roc_auc_score(y_train, y_train_predicted)
24            test_auc = roc_auc_score(y_test, y_test_predicted)
25
26            train_auc_batch.append(train_auc)
27            test_auc_batch.append(test_auc)
28
29        train_auc_final_arr.append(train_auc_batch)
30        test_auc_final_arr.append(test_auc_batch)
31
32    return train_auc_final_arr, test_auc_final_arr
33
34
35 train_auc_final_arr_s1, test_auc_final_arr_s1 = get_auc_matrix(X_train_hstacked_all_tf, y_train_hstacked_all_tf)

```

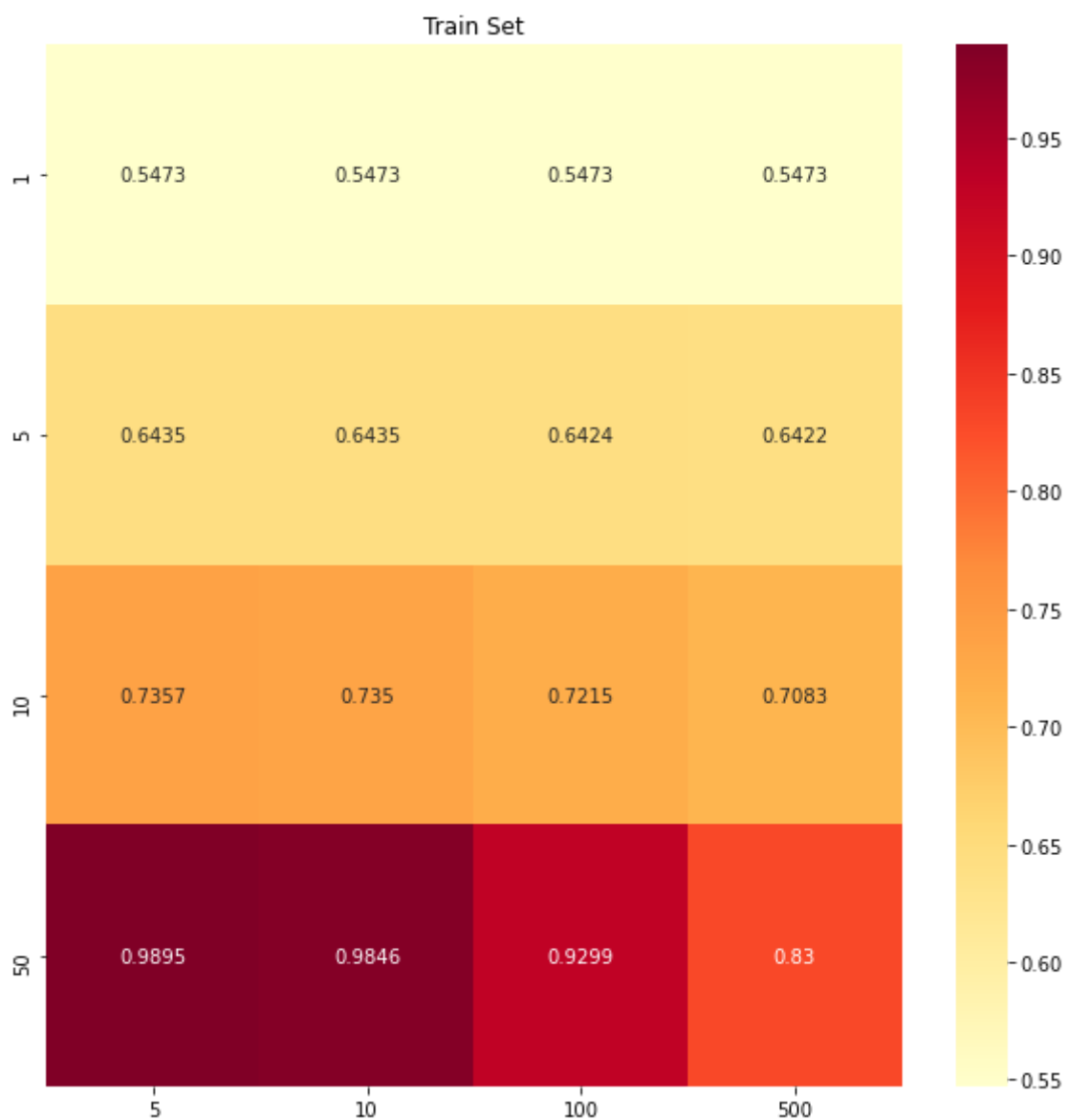
```
100%|██████████| 4/4 [31:42<00:00, 475.54s/it]
```

**plot the performance of model both on train data and cross validation data for each hyper parameter**

- seaborn heat maps with rows as min\_sample\_split, columns as max\_depth, and values inside the cell representing AUC Score.

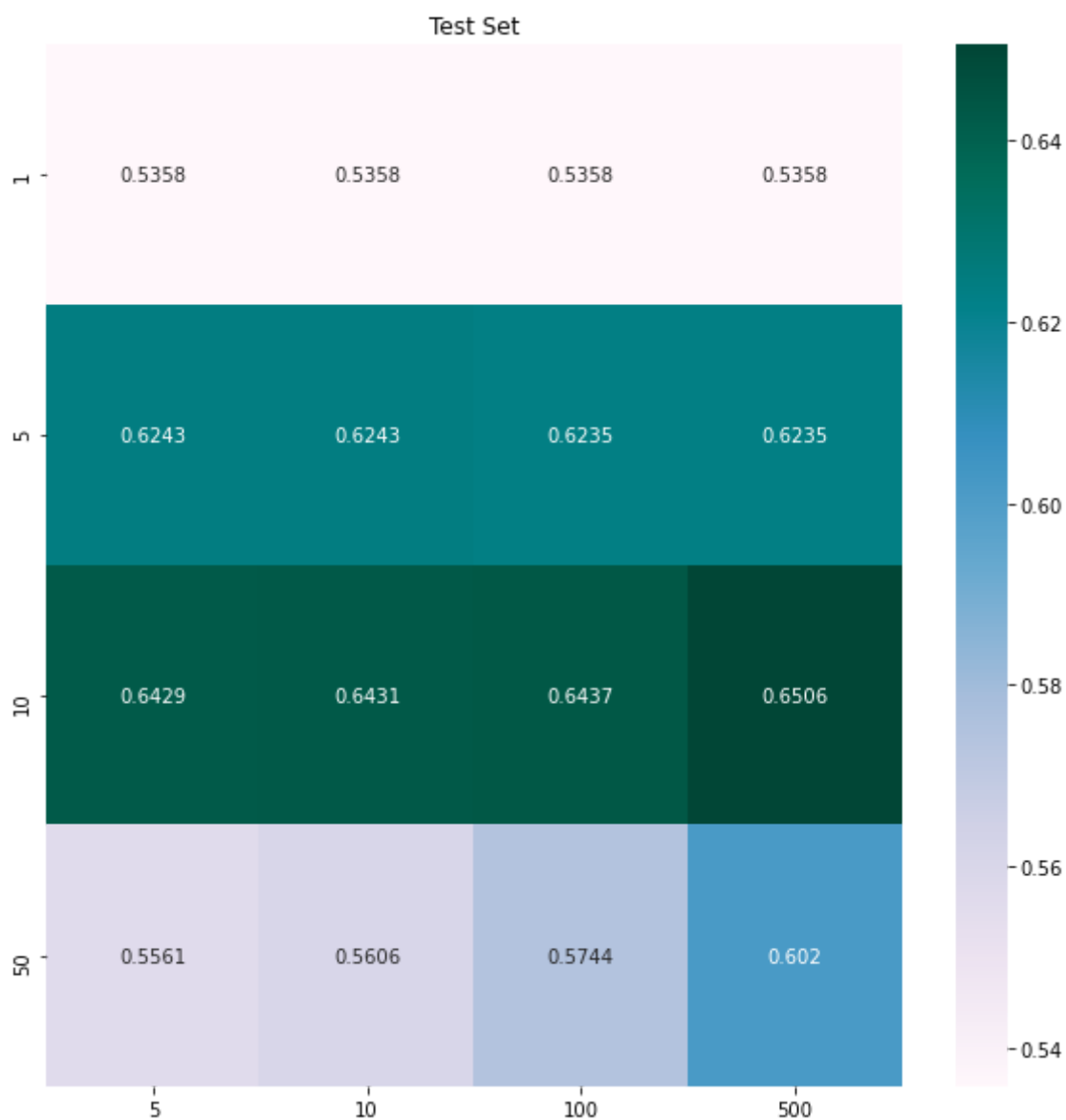
In [65]:

```
1  ## Heatmap for Set S1
2
3  train_auc_final_df_s1 = pd.DataFrame(train_auc_final_arr_s1, columns=min_samples_split,
4  fig, ax = plt.subplots(1, figsize=(10,10))
5  sns.heatmap(train_auc_final_df_s1, annot=True, fmt='.4g', ax=ax, cmap="YlOrRd")
6  ax.set_title('Train Set')
7  plt.show()
8  # train_auc_final_df_s1
```



In [69]:

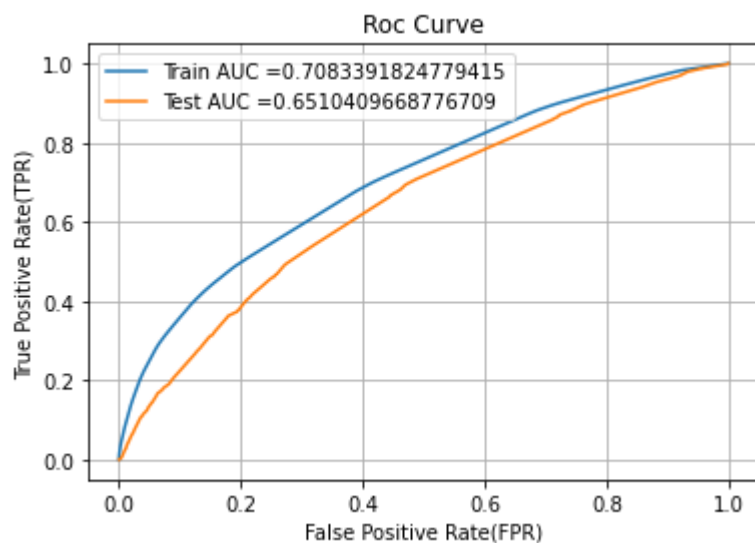
```
1  ## Heatmap for Set S1
2
3  test_auc_final_df_s1 = pd.DataFrame(test_auc_final_arr_s1, columns=min_samples_split, index=
4  fig, ax = plt.subplots(1, figsize=(10,10))
5  sns.heatmap(test_auc_final_df_s1, annot=True, fmt='.4g', ax=ax, cmap="PuBuGn")
6  ax.set_title('Test Set')
7  plt.show()
8  # test_auc_final_df_s1
```



plot the ROC curve on both train and test

In [70]:

```
1 #https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn
2 from sklearn.metrics import roc_curve, auc
3
4 dt_tfidf_Model = DecisionTreeClassifier(class_weight='balanced',min_samples_split=500,r
5 dt_tfidf_Model.fit(X_train_hstacked_all_tfidf_features_vectorized, y_train)
6 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of t
7 # not the predicted outputs
8 # y_train_pred = batch_predict(mnb_bow_testModel, x_train_onehot_bow)
9 y_train_pred=dt_tfidf_Model.predict_proba(X_train_hstacked_all_tfidf_features_vectorize
10 predictions_train_set1=dt_tfidf_Model.predict(X_train_hstacked_all_tfidf_features_vecto
11
12 y_test_pred=dt_tfidf_Model.predict_proba(test_df_hstacked_all_tfidf_features_vectorize
13 predictions_test_set1=dt_tfidf_Model.predict(test_df_hstacked_all_tfidf_features_vecto
14
15
16 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
17 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
18
19 plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
20 plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
21 plt.legend()
22 plt.xlabel("False Positive Rate(FPR)")
23 plt.ylabel("True Positive Rate(TPR)")
24 plt.title("Roc Curve")
25 plt.grid()
26 plt.show()
```



**Confusion Matrix:**

In [71]:

```
1 def predict(proba, threshold, fpr, tpr):
2
3     t = threshold[np.argmax(fpr*(1-tpr))]
4     print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.rou
5     predictions = []
6     global predictions1
7     for i in proba:
8         if i>=t:
9             predictions.append(1)
10        else:
11            predictions.append(0)
12    predictions1 = predictions
13    return predictions
```

In [72]:

```
1 #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
2 import seaborn as sns; sns.set()
3
4 con_m_train = confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr,
5 con_m_test = confusion_matrix(y_test, predict(y_test_pred, te_thresholds, test_fpr, tes
6
7 key = (np.asarray(['TN', 'FP'], ['FN', 'TP'])))
8 fig, ax = plt.subplots(1,2, figsize=(12,5))
9
10 labels_train = (np.asarray(["{0} = {1:.2f}" .format(key, value) for key, value in zip(k
11 labels_test = (np.asarray(["{0} = {1:.2f}" .format(key, value) for key, value in zip(k
12
13 sns.heatmap(con_m_train, linewidths=.5, xticklabels=['PREDICTED : NO', 'PREDICTED : YES
14 sns.heatmap(con_m_test, linewidths=.5, xticklabels=['PREDICTED : NO', 'PREDICTED : YES'
15
16 ax[0].set_title('Train Set')
17 ax[1].set_title('Test Set')
18
19 plt.show()
```

the maximum value of  $tpr*(1-fpr)$  0.4144710515072262 for threshold 0.486  
the maximum value of  $tpr*(1-fpr)$  0.37021130114598005 for threshold 0.486



false positive data points

In [73]:

```
1 #https://www.google.com/search?q=geeks+for+geeks+false+positive&rlz=1C1SQJL_enIN849IN84
2 #https://github.com/pskadasi/DecisionTrees_DonorsChoose/blob/master/Copy_of_8_DonorsCho
3 fpi = []
4 for i in range(len(y_test)):
5     if (y_test.values[i] == 0) & (predictions1[i] == 1) :
6         fpi.append(i)
7 fp_essay1 = []
8 for i in fpi:
9     fp_essay1.append(X_test['essay'].values[i])
```

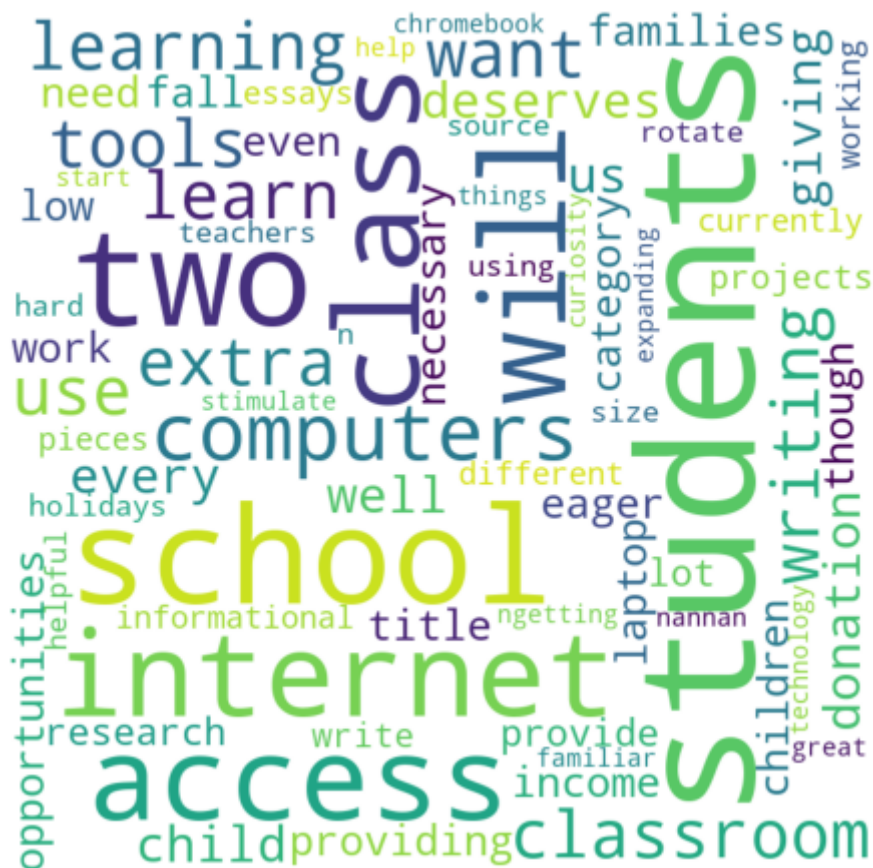
**Plot the WordCloud of essay**

In [74]:

```

1 # Word cloud of essay
2 from wordcloud import WordCloud, STOPWORDS
3 comment_words = ''
4 stopwords = set(STOPWORDS)
5 for val in fp_essay1 :
6     val = str(val)
7     tokens = val.split()
8     for i in range(len(tokens)):
9         tokens[i] = tokens[i].lower()
10    for words in tokens :
11        comment_words = comment_words + words + ' '
12
13 wordcloud = WordCloud(width = 800, height = 800, background_color ='white', stopwords =
14 min_font_size = 10).generate(comment_words)
15
16 plt.figure(figsize = (6, 6), facecolor = None)
17 plt.imshow(wordcloud)
18 plt.axis("off")
19 plt.tight_layout(pad = 0)
20 plt.show()

```



**Plot the box plot with the price of these false positive data points**

In [75]:

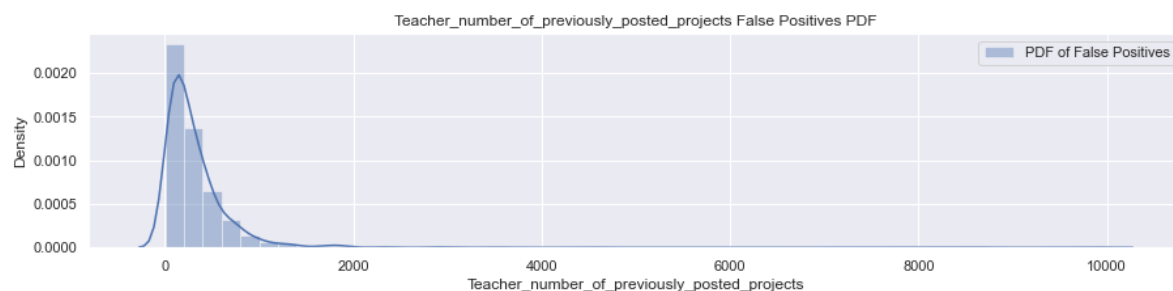
```
1 df = pd.DataFrame(X_test['price'])
2 plt.figure(figsize=(15,3))
3 df1 = df.iloc[fpi, : ]
4 sea.boxplot(df1.values)
5 plt.title("Rejected Projects that Predicted as Positive")
6 plt.ylabel("Box plot for False Positives")
7 plt.xlabel("Price")
8 plt.show()
```



Plot the pdf with the teacher\_number\_of\_previously\_posted\_projects of these false positive data points

In [76]:

```
1 plt.figure(figsize=(15,3))
2 sns.distplot(df1.values, label="PDF of False Positives")
3 plt.title('Teacher_number_of_previously_posted_projects False Positives PDF')
4 plt.xlabel('Teacher_number_of_previously_posted_projects')
5 plt.legend()
6 plt.show()
```



## Preprocessed\_eassay (TFIDF\_W2V)

In [77]:

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



In [78]:

```

1 # In the TF-IDF Word2Vec vectorization, we have to fit the TfidfVectorizer only on X_train
2 # extract 'dictionary' (dictionary with features as the keys and IDF scores as the values)
3 # 'tfidf_words' (a set of all the features extracted from the vectorizer).
4 # We have to use the same 'dictionary' and 'tfidf_words' in vectorizing both X_train['essay'] and X_test['essay']
5
6 # Now, at the very top section of this Notebook, we already have this code of Vectorizer
7 # vectorizer_essay_tfidf = TfidfVectorizer(min_df=10)
8 # vectorizer_essay_tfidf.fit(X_train['essay'].values)
9
10 # Hence we are now converting a dictionary with word as a key, and the idf as a value
11 dictionary = dict(zip(vectorizer_essay_tfidf.get_feature_names(), list(vectorizer_essay_tfidf.idf_)))
12 tfidf_words = set(vectorizer_essay_tfidf.get_feature_names())

```

In [79]:

```

1 # Function to generate Word2Vec weighted by tf-idf
2 def generate_w2v_from_text(essays_text_arr):
3     # compute average word2vec for each review.
4     tfidf_w2v_vectors = []
5     # the avg-w2v for each sentence/review is stored in this list
6
7     for sentence in tqdm(essays_text_arr): # for each sentence
8         vector = np.zeros(300) # as word vectors are of zero length
9         tf_idf_weight = 0
10        # num of words with a valid vector in the sentence
11        for word in sentence.split(): # for each word in a sentence
12            if (word in glove_words) and (word in tfidf_words):
13                vec = model[word] # getting the vector for each word
14                # here we are multiplying idf value(dictionary[word]) and the tf value
15                tf_idf = dictionary[word] * (
16                    sentence.count(word) / len(sentence.split())
17                ) # getting the tfidf value for each word
18                vector += vec * tf_idf # calculating tfidf weighted w2v
19                tf_idf_weight += tf_idf
20        if tf_idf_weight != 0:
21            vector /= tf_idf_weight
22        tfidf_w2v_vectors.append(vector)
23    return tfidf_w2v_vectors
24
25 X_train_vectorized_tfidf_w2v_essay = generate_w2v_from_text(X_train['essay'].values)
26 X_test_vectorized_tfidf_w2v_essay = generate_w2v_from_text(X_test['essay'].values)

```

```
100%|██████████| 73196/73196 [40:40<00:00, 30.00it/s]
100%|██████████| 36052/36052 [13:25<00:00, 44.78it/s]
```

## Merging all categorical, text, numerical vectors and sentiment score based on tfidf\_w2v

In [80]:

```
1 from scipy.sparse import hstack
2 X_train_hstacked_all_tfidf_w2v_vectors_features_vectorized = hstack((train_vectorized_c
3
4 print('X_train_hstacked_all_tfidf_w2v_vectors_features_vectorized.shape is ', X_train_h
5
6 test_df_hstacked_all_tfidf_w2v_vectors_features_vectorized = hstack((test_df_vectorized
7
8 print('test_df_hstacked_all_tfidf_w2v_vectors_features_vectorized.shape is ', test_df_h
```

```
X_train_hstacked_all_tfidf_w2v_vectors_features_vectorized.shape is (73196,
810)
test_df_hstacked_all_tfidf_w2v_vectors_features_vectorized.shape is (36052,
810)
```

## Apply Decision Tree Classifier(DecisionTreeClassifier) on these feature sets

Set 2: categorical, numerical features + preprocessed\_essay (TFIDF W2V) + Sentiment scores(preprocessed\_essay)

In [81]:

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.model_selection import GridSearchCV
4 import seaborn as sea
```

In [82]:

```
1 DT = DecisionTreeClassifier()
2
3 parameters = {'max_depth': [1, 5, 10, 50], 'min_samples_split': [5, 10, 100, 500]}#as c
4
5 classifier = GridSearchCV(DT, parameters, cv=3, scoring='roc_auc', verbose=1, return_train
6 classifier.fit(X_train_hstacked_all_tfidf_w2v_vectors_features_vectorized, y_train)
```

Fitting 3 folds for each of 16 candidates, totalling 48 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 48 out of 48 | elapsed: 48.4min finished
```

Out[82]:

```
GridSearchCV(cv=3, estimator=DecisionTreeClassifier(), n_jobs=-1,
             param_grid={'max_depth': [1, 5, 10, 50],
                         'min_samples_split': [5, 10, 100, 500]},
             return_train_score=True, scoring='roc_auc', verbose=1)
```

In [83]:

```
1 classifier.best_estimator_
```

Out[83]:

```
DecisionTreeClassifier(max_depth=5, min_samples_split=5)
```

In [85]:

```

1 from sklearn.metrics import roc_auc_score
2 max_depth = [1,5,10,50]
3 min_samples_split = [5,10,100,500]
4
5 def get_auc_matrix(X_train, X_test, y_train, y_test ):
6
7     train_auc_final_arr = []
8     test_auc_final_arr = []
9
10    for depth in tqdm(max_depth):
11        train_auc_batch = []
12        test_auc_batch = []
13
14        for split in min_samples_split:
15            df_clf = DecisionTreeClassifier(class_weight='balanced',min_samples_split=split)
16
17            df_clf.fit(X_train, y_train)
18
19            # I have to predict probabilities (clf.predict_proba) instead of classes for
20            y_train_predicted = df_clf.predict_proba(X_train)[: , 1]
21            y_test_predicted = df_clf.predict_proba(X_test)[: , 1]
22
23            train_auc = roc_auc_score(y_train, y_train_predicted)
24            test_auc = roc_auc_score(y_test, y_test_predicted)
25
26            train_auc_batch.append(train_auc)
27            test_auc_batch.append(test_auc)
28
29        train_auc_final_arr.append(train_auc_batch)
30        test_auc_final_arr.append(test_auc_batch)
31
32    return train_auc_final_arr, test_auc_final_arr
33
34
35 train_auc_final_arr_s2, test_auc_final_arr_s2 = get_auc_matrix(X_train_hstacked_all_tf1, y_train_hstacked_all_tf1, X_test_hstacked_all_tf1, y_test_hstacked_all_tf1)

```

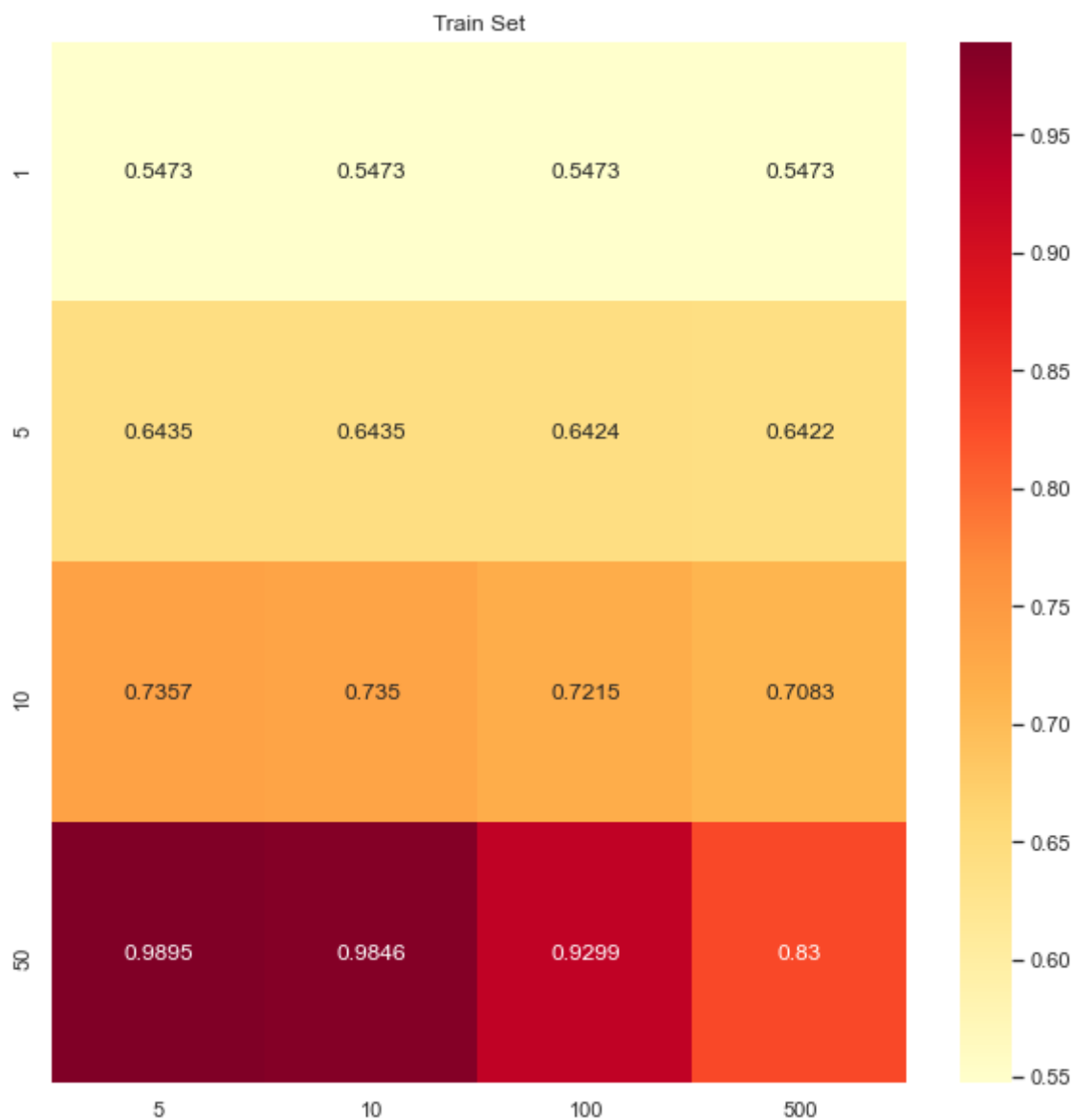
```
100%|██████████| 4/4 [44:00<00:00, 660.10s/it]
```

**plot the performance of model both on train data and cross validation data for each hyper parameter**

- seaborn heat maps with rows as min\_sample\_split, columns as max\_depth, and values inside the cell representing AUC Score.

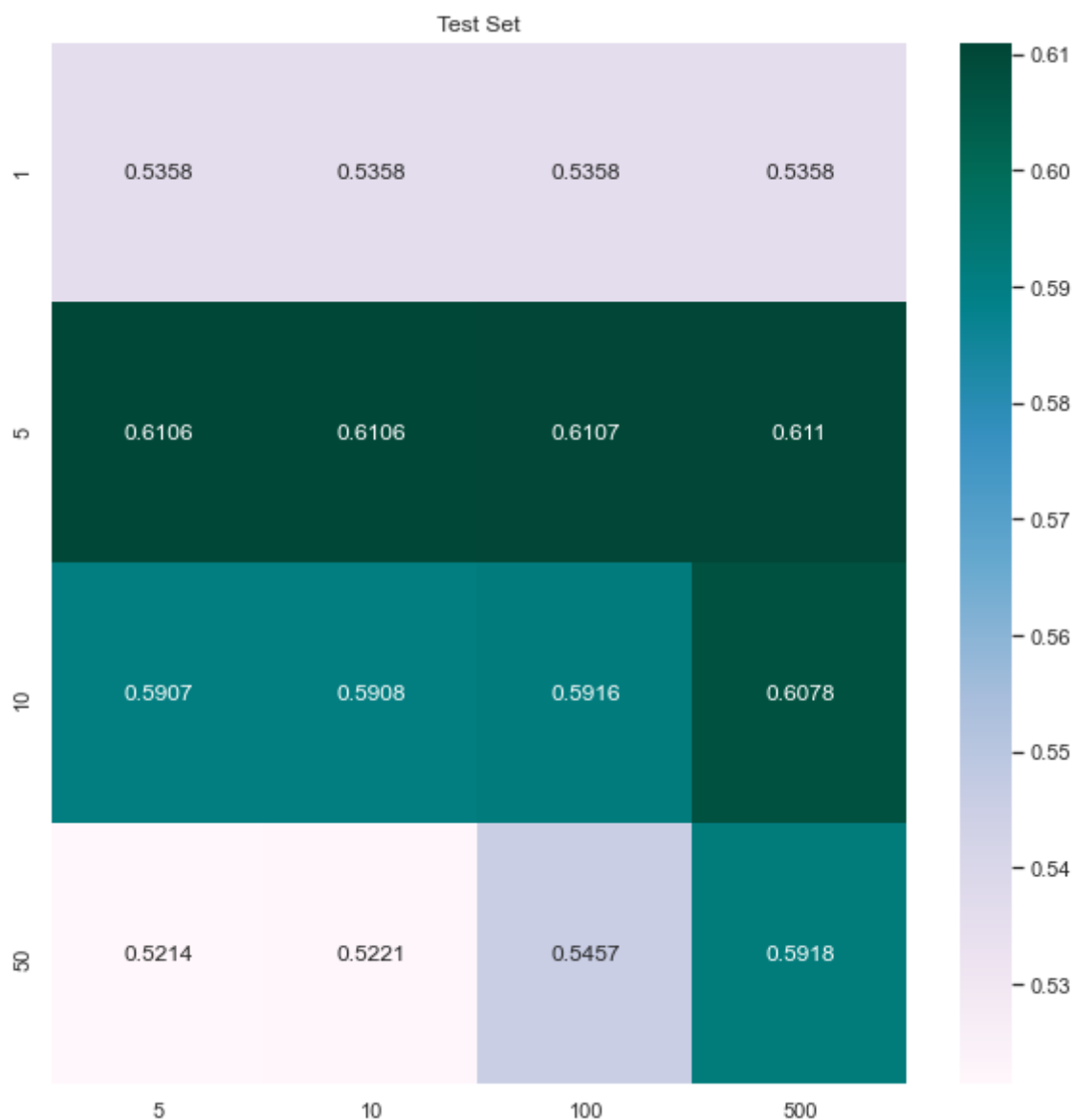
In [86]:

```
1  ## Heatmap for Set S2
2
3  train_auc_final_df_s2 = pd.DataFrame(train_auc_final_arr_s2, columns=min_samples_split,
4  fig, ax = plt.subplots(1, figsize=(10,10))
5  sns.heatmap(train_auc_final_df_s1, annot=True, fmt='.4g', ax=ax, cmap="YlOrRd")
6  ax.set_title('Train Set')
7  plt.show()
8  # train_auc_final_df_s2
```



In [87]:

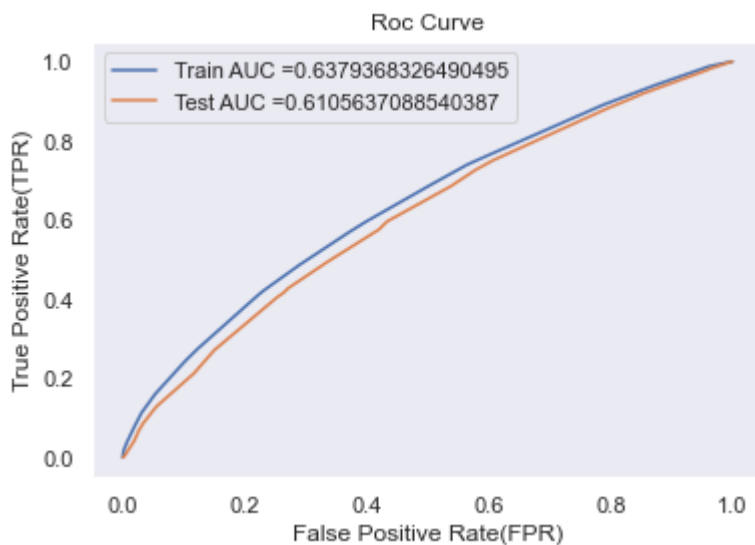
```
1  ## Heatmap for Set S2
2
3  test_auc_final_df_s2 = pd.DataFrame(test_auc_final_arr_s2, columns=min_samples_split, index=
4  fig, ax = plt.subplots(1, figsize=(10,10))
5  sns.heatmap(test_auc_final_df_s2, annot=True, fmt='.4g', ax=ax, cmap="PuBuGn")
6  ax.set_title('Test Set')
7  plt.show()
8  # test_auc_final_df_s2
```



plot roc curve

In [88]:

```
1 #https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn
2 from sklearn.metrics import roc_curve, auc
3
4 dt_tfidf_Model = DecisionTreeClassifier(class_weight='balanced',min_samples_split=5,max
5 dt_tfidf_Model.fit(X_train_hstacked_all_tfidf_w2v_vectors_features_vectorized, y_train)
6 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of t
7 # not the predicted outputs
8 # y_train_pred = batch_predict(mnb_bow_testModel, x_train_onehot_bow)
9 y_train_pred=dt_tfidf_Model.predict_proba(X_train_hstacked_all_tfidf_w2v_vectors_featur
10 predictions_train_set1=dt_tfidf_Model.predict(X_train_hstacked_all_tfidf_w2v_vectors_fe
11
12 y_test_pred=dt_tfidf_Model.predict_proba(test_df_hstacked_all_tfidf_w2v_vectors_featur
13 predictions_test_set1=dt_tfidf_Model.predict(test_df_hstacked_all_tfidf_w2v_vectors_fea
14
15
16 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
17 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
18
19 plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
20 plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
21 plt.legend()
22 plt.xlabel("False Positive Rate(FPR)")
23 plt.ylabel("True Positive Rate(TPR)")
24 plt.title("Roc Curve")
25 plt.grid()
26 plt.show()
```



## Confusion Matrix

In [89]:

```
1 def predict(proba, threshold, fpr, tpr):
2
3     t = threshold[np.argmax(fpr*(1-tpr))]
4     print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.rou
5     predictions = []
6     global predictions1
7     for i in proba:
8         if i>=t:
9             predictions.append(1)
10        else:
11            predictions.append(0)
12    predictions1 = predictions
13    return predictions
```

In [90]:

```
1 #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
2 import seaborn as sns; sns.set()
3
4 con_m_train = confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr,
5 con_m_test = confusion_matrix(y_test, predict(y_test_pred, te_thresholds, test_fpr, tes
6
7 key = (np.asarray(['TN', 'FP'], ['FN', 'TP'])))
8 fig, ax = plt.subplots(1,2, figsize=(12,5))
9
10 labels_train = (np.asarray(["{0} = {1:.2f}" .format(key, value) for key, value in zip(k
11 labels_test = (np.asarray(["{0} = {1:.2f}" .format(key, value) for key, value in zip(k
12
13 sns.heatmap(con_m_train, linewidths=.5, xticklabels=['PREDICTED : NO', 'PREDICTED : YES
14 sns.heatmap(con_m_test, linewidths=.5, xticklabels=['PREDICTED : NO', 'PREDICTED : YES'
15
16 ax[0].set_title('Train Set')
17 ax[1].set_title('Test Set')
18
19 plt.show()
```

the maximum value of  $tpr*(1-fpr)$  0.35805001510615114 for threshold 0.484

the maximum value of  $tpr*(1-fpr)$  0.3375928246728687 for threshold 0.504



False positive point

In [91]:

```
1 #https://www.google.com/search?q=geeks+for+geeks+false+positive&rlz=1C1SQJL_enIN849IN84
2 #https://github.com/pskadasi/DecisionTrees_DonorsChoose/blob/master/Copy_of_8_DonorsCho
3 fpi = []
4 for i in range(len(y_test)):
5     if (y_test.values[i] == 0) & (predictions1[i] == 1) :
6         fpi.append(i)
7 fp_essay1 = []
8 for i in fpi:
9     fp_essay1.append(X_test['essay'].values[i])
```

plot Word cloud of essay

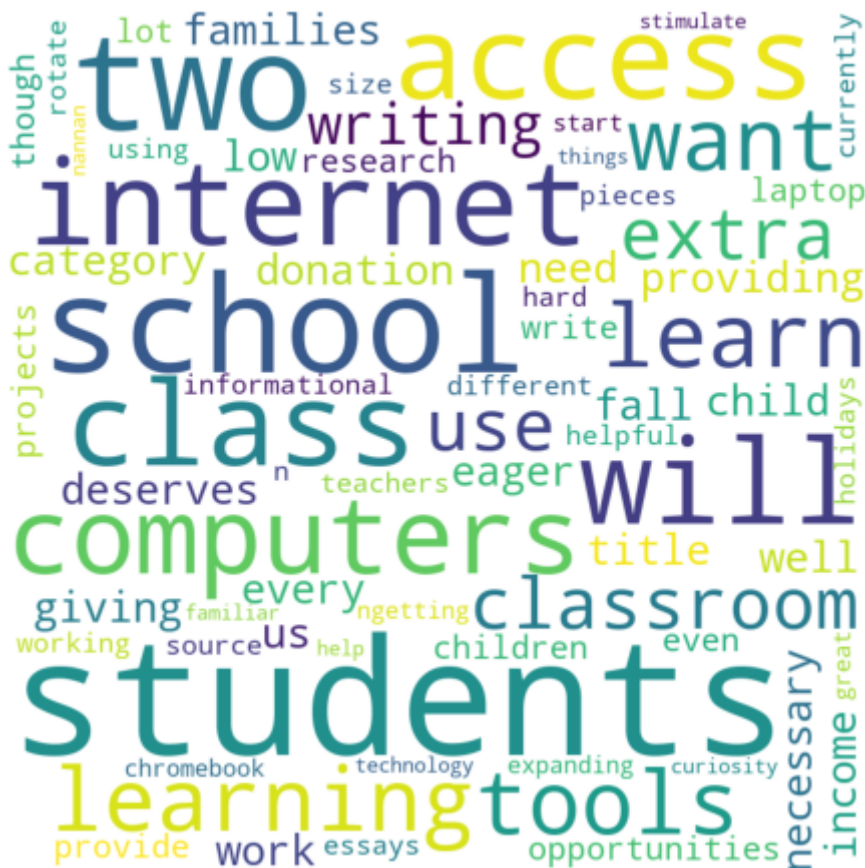


In [92]:

```

1 # Word cloud of essay
2 from wordcloud import WordCloud, STOPWORDS
3 comment_words = ''
4 stopwords = set(STOPWORDS)
5 for val in fp_essay1 :
6     val = str(val)
7     tokens = val.split()
8     for i in range(len(tokens)):
9         tokens[i] = tokens[i].lower()
10    for words in tokens :
11        comment_words = comment_words + words + ' '
12
13 wordcloud = WordCloud(width = 800, height = 800, background_color ='white', stopwords =
14 min_font_size = 10).generate(comment_words)
15
16 plt.figure(figsize = (6, 6), facecolor = None)
17 plt.imshow(wordcloud)
18 plt.axis("off")
19 plt.tight_layout(pad = 0)
20 plt.show()

```



### Box plot for False Positives of price

In [93]:

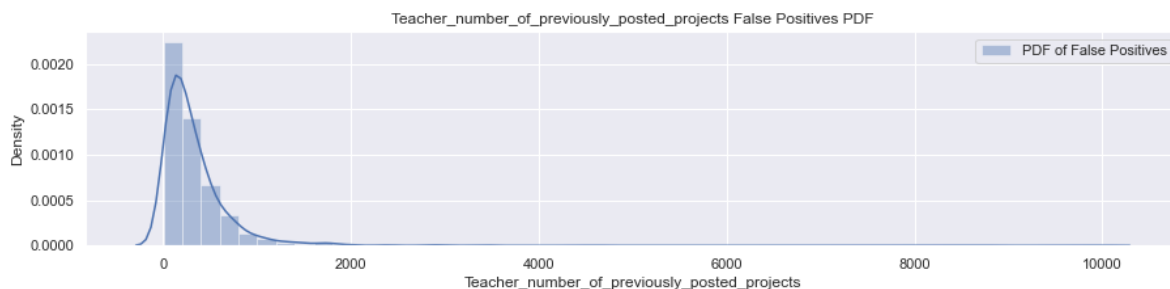
```
1 df = pd.DataFrame(X_test['price'])
2 plt.figure(figsize=(15,3))
3 df1 = df.iloc[fpi, : ]
4 sea.boxplot(df1.values)
5 plt.title("Rejected Projects that Predicted as Positive")
6 plt.ylabel("Box plot for False Positives")
7 plt.xlabel("Price")
8 plt.show()
```



plot Teacher\_number\_of\_previously\_posted\_projects False Positives PDF

In [94]:

```
1 plt.figure(figsize=(15,3))
2 sns.distplot(df1.values, label="PDF of False Positives")
3 plt.title('Teacher_number_of_previously_posted_projects False Positives PDF')
4 plt.xlabel('Teacher_number_of_previously_posted_projects')
5 plt.legend()
6 plt.show()
```



## TASK-2

For this task consider set-1 features.

- Select all the features which are having non-zero feature importance. You can get the feature importance using 'feature\_importances\_' (<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>), discard the all other

remaining features and then apply any of the model of your choice i.e. (Decision tree, Logistic Regression, Linear SVM).

In [95]:

```
1 from scipy.sparse import hstack
2 X_train1 = hstack((train_vectorized_ohe_clean_categories, train_vectorized_ohe_clean_su
3
4 print('X_train1.shape is ', X_train1.shape)
5 print(type(X_train1))
6 X_test1 = hstack((test_df_vectorized_ohe_clean_categories, test_df_vectorized_ohe_clear
7
8 print('X_test1.shape is ', X_test1.shape)
9 print(type(X_test1))
10
```

```
X_train1.shape is (73196, 15294)
<class 'scipy.sparse.csr.csr_matrix'>
X_test1.shape is (36052, 15294)
<class 'scipy.sparse.csr.csr_matrix'>
```

In [96]:

```
1 #https://stackoverflow.com/questions/47111434/randomforestregressor-and-feature-importance
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.ensemble import RandomForestRegressor
4 from sklearn.model_selection import GridSearchCV
5 def selectKImportance(model, X):
6     return X[:,model.best_estimator_.feature_importances_.argsort()]
```

In [97]:

```
1 X_train1_best = selectKImportance(clf, X_train1)
2 X_test1_best = selectKImportance(clf, X_test1)
3
4 print(X_train1_best.shape)
5 print(X_test1_best.shape)
```

```
(73196, 15294)
(36052, 15294)
```

In [98]:

```
1 from sklearn.metrics import roc_auc_score
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import train_test_split
4 from sklearn.model_selection import GridSearchCV
5 from sklearn.model_selection import cross_val_score
6 from sklearn.tree import DecisionTreeClassifier
7 params = {'max_depth': [1, 5, 10, 50], 'min_samples_split': [5, 10, 100, 500]}#as given
8 clf_1 = GridSearchCV(DecisionTreeClassifier(random_state=42), params, cv=3, scoring='roc_auc')
9 clf_1.fit(X_train1_best, y_train)
```

Fitting 3 folds for each of 16 candidates, totalling 48 fits

[Parallel(n\_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.

[Parallel(n\_jobs=-1)]: Done 48 out of 48 | elapsed: 26.0min finished

Out[98]:

```
GridSearchCV(cv=3, estimator=DecisionTreeClassifier(random_state=42), n_jobs=-1,
             param_grid={'max_depth': [1, 5, 10, 50],
                          'min_samples_split': [5, 10, 100, 500]},
             return_train_score=True, scoring='roc_auc', verbose=1)
```

In [99]:

```
1 clf_1.best_estimator_
```

Out[99]:

```
DecisionTreeClassifier(max_depth=10, min_samples_split=500, random_state=42)
```

In [100]:

```

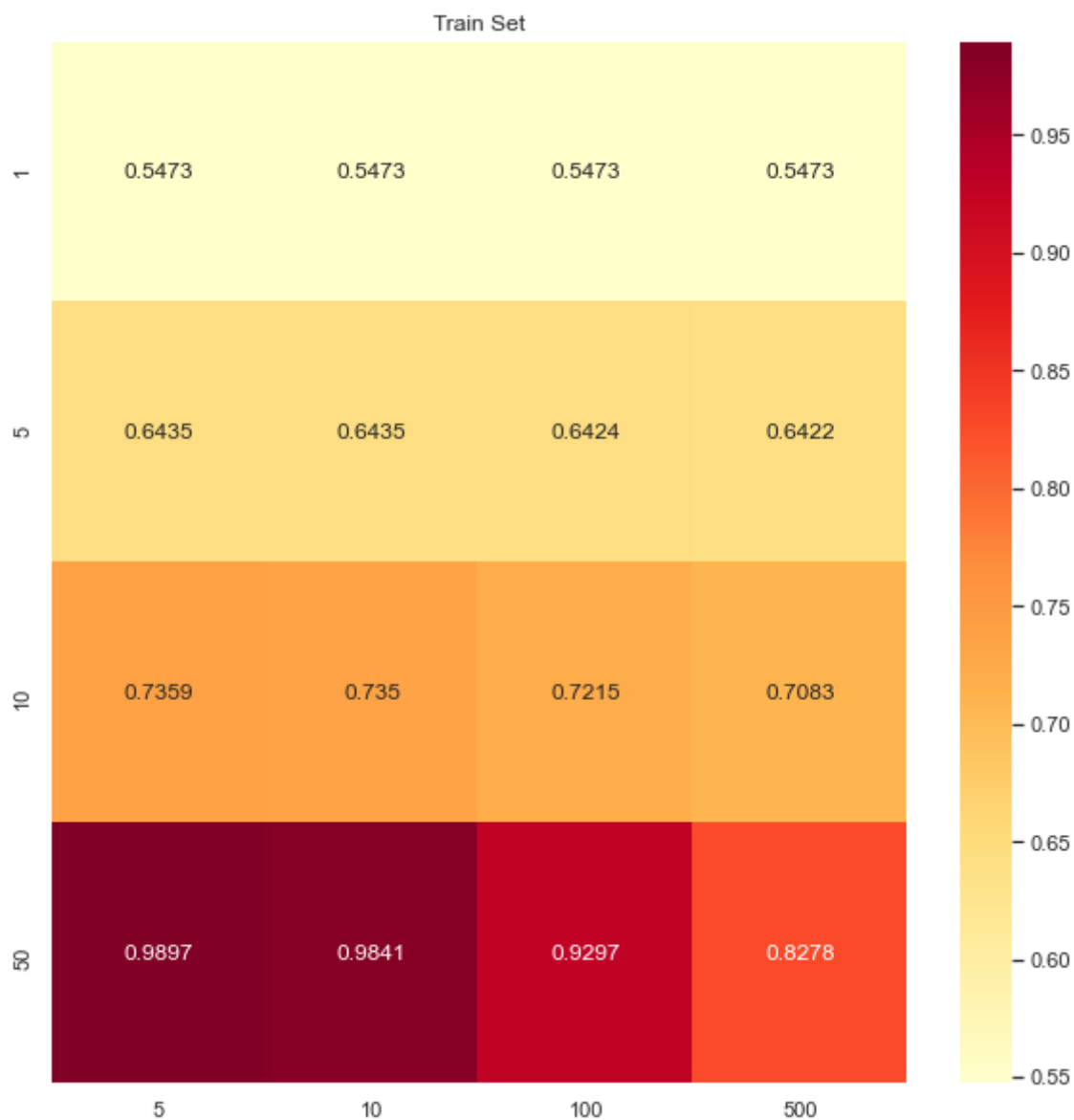
1 from sklearn.metrics import roc_auc_score
2 max_depth = [1,5,10,50]
3 min_samples_split = [5,10,100,500]
4
5 def get_auc_matrix(X_train, X_test, y_train, y_test ):
6
7     train_auc_final_arr = []
8     test_auc_final_arr = []
9
10    for depth in tqdm(max_depth):
11        train_auc_batch = []
12        test_auc_batch = []
13
14        for split in min_samples_split:
15            df_clf = DecisionTreeClassifier(class_weight='balanced',min_samples_split=split)
16
17            df_clf.fit(X_train, y_train)
18
19            # I have to predict probabilities (clf.predict_proba) instead of classes for
20            y_train_predicted = df_clf.predict_proba(X_train)[: , 1]
21            y_test_predicted = df_clf.predict_proba(X_test)[: , 1]
22
23            train_auc = roc_auc_score(y_train, y_train_predicted)
24            test_auc = roc_auc_score(y_test, y_test_predicted)
25
26            train_auc_batch.append(train_auc)
27            test_auc_batch.append(test_auc)
28
29        train_auc_final_arr.append(train_auc_batch)
30        test_auc_final_arr.append(test_auc_batch)
31
32    return train_auc_final_arr, test_auc_final_arr
33
34
35 train_auc_final_arr_task2, test_auc_final_arr_task2 = get_auc_matrix(X_train1_best, X_train2_best, y_train1, y_train2)

```

```
100%|███████████|  
██████████ | 4/4 [31:23<00:00, 470.87s/it]
```

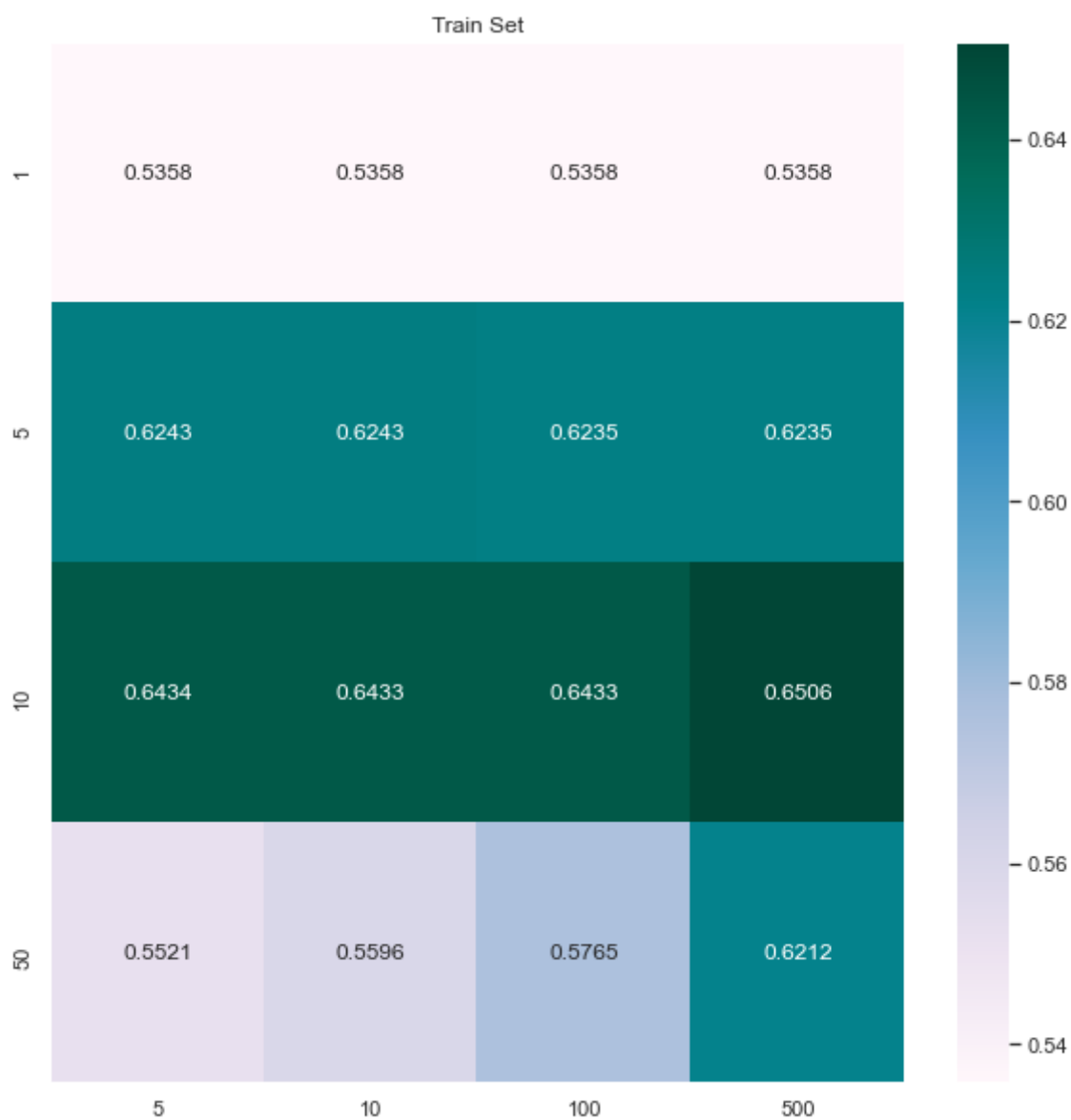
In [101]:

```
1  ## Heatmap for task2
2
3  train_auc_final_df_task2 = pd.DataFrame(train_auc_final_arr_task2, columns=min_samples_s
4  fig, ax = plt.subplots(1, figsize=(10,10))
5  sns.heatmap(train_auc_final_df_task2, annot=True, fmt='.4g', ax=ax, cmap="YlOrRd")
6  ax.set_title('Train Set')
7  plt.show()
8  # train_auc_final_df_task2
```



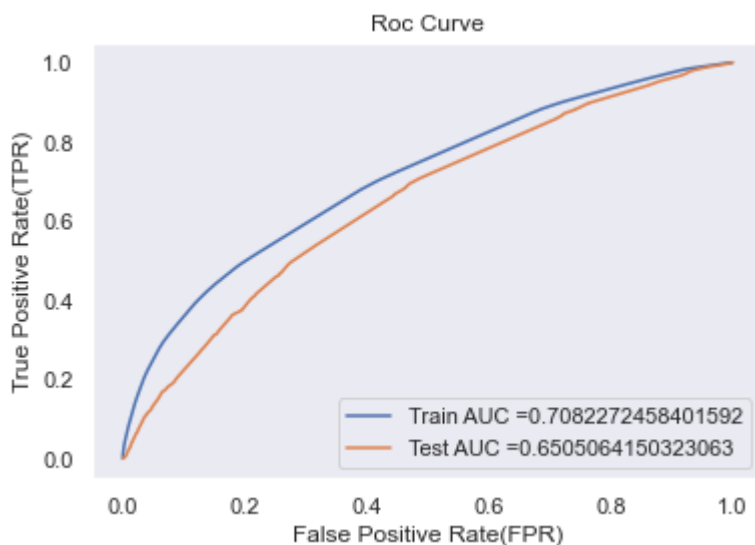
In [105]:

```
1  ## Heatmap for task2
2
3  test_auc_final_df_task2 = pd.DataFrame(test_auc_final_arr_task2, columns=min_samples_sp1
4  fig, ax = plt.subplots(1, figsize=(10,10))
5  sns.heatmap(test_auc_final_df_task2, annot=True, fmt='.4g', ax=ax, cmap="PuBuGn")
6  ax.set_title('Train Set')
7  plt.show()
8  # test_auc_final_df_task2
```



In [104]:

```
1 #https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn
2 from sklearn.metrics import roc_curve, auc
3
4 dt_tfidf_Model = DecisionTreeClassifier(class_weight='balanced',min_samples_split=500,r
5 dt_tfidf_Model.fit(X_train1_best, y_train)
6 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of t
7 # not the predicted outputs
8 # y_train_pred = batch_predict(mnb_bow_testModel, x_train_onehot_bow)
9 y_train_pred=dt_tfidf_Model.predict_proba(X_train1_best)[: ,1]
10 predictions_train_set1=dt_tfidf_Model.predict(X_train1_best)
11
12 y_test_pred=dt_tfidf_Model.predict_proba(X_test1_best)[: ,1]
13 predictions_test_set1=dt_tfidf_Model.predict(X_test1_best)
14
15
16 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
17 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
18
19 plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
20 plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
21 plt.legend()
22 plt.xlabel("False Positive Rate(FPR)")
23 plt.ylabel("True Positive Rate(TPR)")
24 plt.title("Roc Curve")
25 plt.grid()
26 plt.show()
```



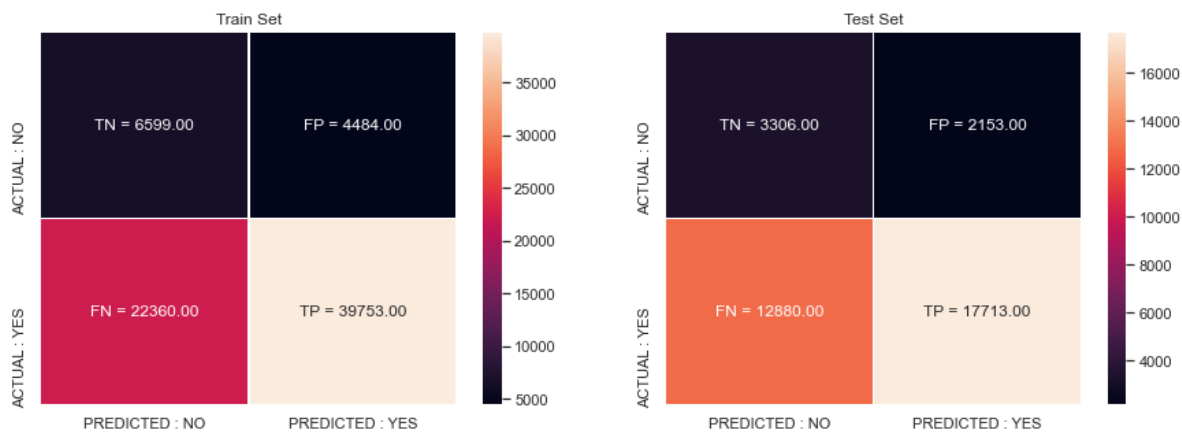


In [194]:

```
1 #CONFUSION MATRIX
2 #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
3 import seaborn as sns; sns.set()
4 con_m_train = confusion_matrix(y_train, predict(y_train_pred1, tr_thresholds1, train_fpr1,
5 con_m_test = confusion_matrix(y_test, predict(y_test_pred1, te_thresholds1, test_fpr1,
6 key = (np.asarray(['TN', 'FP'], ['FN', 'TP'])))
7 fig, ax = plt.subplots(1,2, figsize=(15,5))
8 labels_train = (np.asarray(["{0} = {1:.2f}" .format(key, value) for key, value in zip(key,
9 labels_test = (np.asarray(["{0} = {1:.2f}" .format(key, value) for key, value in zip(key,
10 con_m_test.flatten()))).reshape(2,2)
11 sns.heatmap(con_m_train, linewidths=.5, xticklabels=['PREDICTED : NO', 'PREDICTED : YES',
12 sns.heatmap(con_m_test, linewidths=.5, xticklabels=['PREDICTED : NO', 'PREDICTED : YES',
13 ax[0].set_title('Train Set')
14 ax[1].set_title('Test Set')
15 plt.show()
```

the maximum value of  $tpr*(1-fpr)$  0.3824625740439171 for threshold 0.486

the maximum value of  $tpr*(1-fpr)$  0.3506386704184174 for threshold 0.486



## Conclusion:

In [106]:

```
1 # Please compare all your models using Prettytable library
2 #how to use pretty table http://zetcode.com/python/prettytable/
3 from prettytable import PrettyTable
4 tb = PrettyTable()
5 tb.field_names= ( " Vectorizer ", " Max_depth ", " Min_sample_split ", " Test -AUC ")
6 tb.add_row([ " Tf - Idf", 10 , 500 ,65.10 ])
7 tb.add_row([ "AVG - Tf - Idf", 5 , 5 ,61.07])
8 tb.add_row([ "best Features", 10, 500 ,65.05 ])
9 print(tb.get_string(titles = "Decision trees- Observations"))
```

Vectorizer	Max_depth	Min_sample_split	Test -AUC
Tf - Idf	10	500	65.1
AVG - Tf - Idf	5	5	61.07
best Features	10	500	65.05

### Refrence:

- From naive bayes Assignment.(<http://localhost:8888/notebooks/naive%20bayes.ipynb>  
(<http://localhost:8888/notebooks/naive%20bayes.ipynb>))
- <https://www.kaggle.com/onuragmaji/assignment-08-decision-tree>  
(<https://www.kaggle.com/onuragmaji/assignment-08-decision-tree>)