

# Assignment 8: DT

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

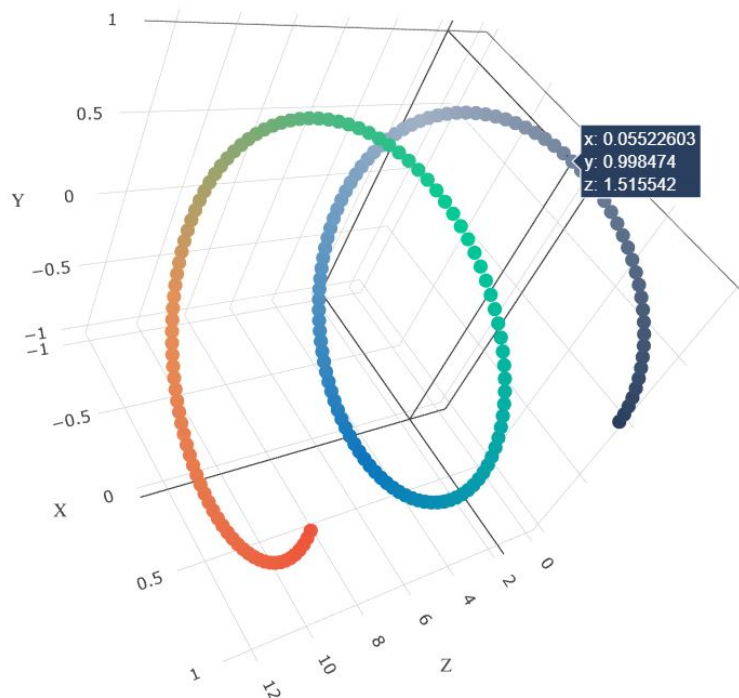
- **Set 1:** categorical, numerical features + preprocessed\_eassay (TFIDF)
- **Set 2:** categorical, numerical features + preprocessed\_eassay (TFIDF W2V)

## 2. The hyper paramter tuning (best depth in range [1, 5, 10, 50], and the best min\_samples\_split in range [5, 10, 100, 500])

- Find the best hyper parameter which will give the maximum [AUC](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/>) value
- find the best hyper paramter using k-fold cross validation(use gridsearch cv or randomsearch cv)/simple cross validation data(you can write your own for loops refer sample solution)

## 3. Representation of results

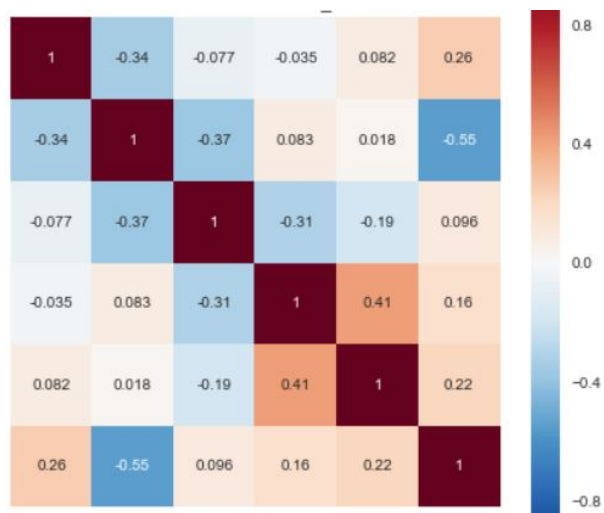
- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure



with X-axis as **min\_sample\_split**, Y-axis as **max\_depth**, and Z-axis as **AUC Score** , we have given the notebook which explains how to plot this 3d plot, you can find it in the same drive *3d\_scatter\_plot.ipynb*

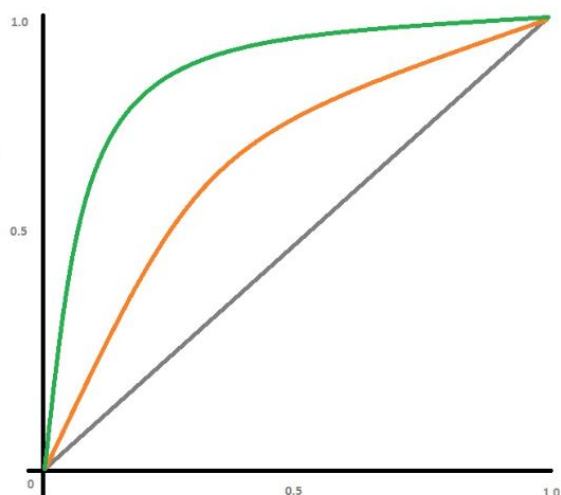
**or**

- 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



[seaborn heat maps](https://seaborn.pydata.org/generated/seaborn.heatmap.html) (<https://seaborn.pydata.org/generated/seaborn.heatmap.html>) with rows as **n\_estimators**, columns as **max\_depth**, and values inside the cell representing **AUC Score**

- You choose either of the plotting techniques out of 3d plot or heat map
- 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](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/>) with predicted and original labels of test data points

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

- Once after you plot the confusion matrix with the test data, get all the false positive data points
  - Plot the WordCloud(<https://www.geeksforgeeks.org/generating-word-cloud-python/>) with the words of essay text of these false positive data points

- 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
4. **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), you need to do hyperparameter tuning corresponding to the model you selected and procedure in step 2 and step 3
- Note: when you want to find the feature importance make sure you don't use max\_depth parameter keep it None.

5. You need to summarize the results at the end of the notebook, summarize it in the table format

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

# 1. Decision Tree

## 1.1 Loading Data

In [1]: 1 ! pip install plotly

Requirement already satisfied: plotly in c:\users\aaa\anaconda3\lib\site-packages (4.6.0)

Requirement already satisfied: six in c:\users\aaa\anaconda3\lib\site-packages (from plotly) (1.11.0)

Requirement already satisfied: retrying>=1.3.3 in c:\users\aaa\anaconda3\lib\site-packages (from plotly) (1.3.3)

distributed 1.21.8 requires msgpack, which is not installed.

You are using pip version 10.0.1, however version 20.1 is available.

You should consider upgrading via the 'python -m pip install --upgrade pip' command.

```

In [2]: 1 %matplotlib inline
        2 import warnings
        3 warnings.filterwarnings("ignore")
        4
        5 import pandas as pd
        6 import numpy as np
        7 import nltk
        8 import matplotlib.pyplot as plt
        9 import seaborn as sns
       10 from sklearn.feature_extraction.text import TfidfVectorizer
       11 from sklearn.feature_extraction.text import CountVectorizer
       12 from sklearn.metrics import confusion_matrix
       13 from sklearn import metrics
       14 from sklearn.metrics import roc_curve, auc
       15
       16 import re
       17 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
       18
       19 import pickle
       20 from tqdm import tqdm
       21 import os
       22
       23 # from plotly import plotly
       24 # import plotly.offline as offline
       25 # import plotly.graph_objs as go
       26 # offline.init_notebook_mode()
       27 from collections import Counter

```

```

In [3]: 1 import pandas as pd
        2 data = pd.read_csv('preprocessed_data.csv',nrows=50000)
        3 data.head(1)

```

```

Out[3]:
  school_state  teacher_prefix  project_grade_category  teacher_number_of_previously_posted_projec

```

```

0          ca          mrs          grades_prek_2

```

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

```
In [4]: 1 # please write all the code with proper documentation, and proper titles for
2 # go through documentations and blogs before you start coding
3 # first figure out what to do, and then think about how to do.
4 # reading and understanding error messages will be very much helpfull in debu
5 # when you plot any graph make sure you use
6 # a. Title, that describes your plot, this will be very helpful to the re
7 # b. Legends if needed
8 # c. X-axis Label
9 # d. Y-axis Label
```

```
In [5]: 1 y = data['project_is_approved'].values
2 X = data.drop(['project_is_approved'], axis=1)
3 X.head(1)
```

```
Out[5]:
```

	school_state	teacher_prefix	project_grade_category	teacher_number_of_previously_posted_projec
0	ca	mrs	grades_prek_2	

```
In [6]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, str
3 # X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size
4
5 print(X_train.shape)
6 print(X_test.shape)
7 print("*****")
8 print(y_train.shape)
9 print(y_test.shape)
```

```
(33500, 8)
(16500, 8)
*****
(33500,)
(16500,)
```

```
In [ ]: 1
```

```
In [ ]: 1
```

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

## TFIDF Vectorization

```
In [7]: 1 # please write all the code with proper documentation, and proper titles for
2 # go through documentations and blogs before you start coding
3 # first figure out what to do, and then think about how to do.
4 # reading and understanding error messages will be very much helpfull in debu
5 # make sure you featurize train and test data separatly
6
7 # when you plot any graph make sure you use
8     # a. Title, that describes your plot, this will be very helpful to the re
9     # b. Legends if needed
10    # c. X-axis label
11    # d. Y-axis label
```

## TFIDF Vectorization Eassy

```
In [8]: 1 vect= TfidfVectorizer(min_df=10,ngram_range=(1,4), max_features=50000)
2 vect.fit(X_train['essay'].values)
3 X_train_essay_tfidf = vect.transform(X_train['essay'].values)
4 # X_cv_essay_tfidf = vect.transform(X_cv['essay'].values)
5 X_test_essay_tfidf = vect.transform(X_test['essay'].values)
6 essay_tfidf_features=vect.get_feature_names()
7
8 print('X_train_essay_tfidf.shape',X_train_essay_tfidf.shape, y_train.shape)
9 print('X_test_essay_tfidf.shape',X_test_essay_tfidf.shape, y_test.shape)
10 # print('X_cv_essay_tfidf.shape',X_cv_essay_tfidf.shape, y_cv.shape)s
11 print(essay_tfidf_features[:10])
```

X\_train\_essay\_tfidf.shape (33500, 50000) (33500,)

X\_test\_essay\_tfidf.shape (16500, 50000) (16500,)

['00', '000', '000 steps', '000 students', '10', '10 000', '10 11', '10 12', '10 15', '10 chromebooks']

## 1.4 Make Data Model Ready: encoding numerical, categorical features

```
In [9]: 1 # please write all the code with proper documentation, and proper titles for
2 # go through documentations and blogs before you start coding
3 # first figure out what to do, and then think about how to do.
4 # reading and understanding error messages will be very much helpfull in debu
5 # make sure you featurize train and test data separatly
6
7 # when you plot any graph make sure you use
8     # a. Title, that describes your plot, this will be very helpful to the re
9     # b. Legends if needed
10    # c. X-axis label
11    # d. Y-axis label
```

## school\_state-OHE

```
In [10]: 1 vect= CountVectorizer(min_df=10,ngram_range=(1,4), max_features=50000)
2 vect.fit(X_train['school_state'].values)
3 X_train_ohe_school_state= vect.transform(X_train['school_state'].values)
4 # X_cv_ohe_school_state= vect.transform(X_cv['school_state'].values)
5 X_test_ohe_school_state= vect.transform(X_test['school_state'].values)
6 school_state_features=vect.get_feature_names()
7
8 print('X_train_ohe_school_state',X_train_ohe_school_state.shape, y_train.shape)
9 print('X_test_ohe_school_state',X_test_ohe_school_state.shape, y_test.shape)
10 # print('X_cv_ohe_school_state',X_cv_ohe_school_state.shape, y_cv.shape)
11 print(school_state_features[:10])
```

```
X_train_ohe_school_state (33500, 51) (33500,)
X_test_ohe_school_state (16500, 51) (16500,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl']
```

## teacher\_prefix-OHE

```
In [11]: 1 vect= CountVectorizer(min_df=10,ngram_range=(1,4), max_features=50000)
2 vect.fit(X_train['teacher_prefix'].values)
3 X_train_ohe_teacher_prefix= vect.transform(X_train['teacher_prefix'].values)
4 # X_cv_ohe_teacher_prefix= vect.transform(X_cv['teacher_prefix'].values)
5 X_test_ohe_teacher_prefix= vect.transform(X_test['teacher_prefix'].values)
6 teacher_prefix_features=vect.get_feature_names()
7
8 print('X_train_ohe_teacher_prefix',X_train_ohe_teacher_prefix.shape, y_train.shape)
9 print('X_test_ohe_teacher_prefix',X_test_ohe_teacher_prefix.shape, y_test.shape)
10 # print('X_cv_ohe_teacher_prefix',X_cv_ohe_teacher_prefix.shape, y_cv.shape)
11 print(teacher_prefix_features[:10])
```

```
X_train_ohe_teacher_prefix (33500, 4) (33500,)
X_test_ohe_teacher_prefix (16500, 4) (16500,)
['mr', 'mrs', 'ms', 'teacher']
```

## project\_grade\_category-OHE

```
In [12]: 1 vect= CountVectorizer(min_df=10,ngram_range=(1,4), max_features=50000)
2 vect.fit(X_train['project_grade_category'].values)
3 X_train_ohe_project_grade_category= vect.transform(X_train['project_grade_category'].values)
4 # X_cv_ohe_project_grade_category= vect.transform(X_cv['project_grade_category'].values)
5 X_test_ohe_project_grade_category= vect.transform(X_test['project_grade_category'].values)
6 project_grade_category_features=vect.get_feature_names()
7
8 print('X_train_ohe_project_grade_category',X_train_ohe_project_grade_category.shape, y_train.shape)
9 print('X_test_ohe_project_grade_category',X_test_ohe_project_grade_category.shape, y_test.shape)
10 # print('X_cv_ohe_project_grade_category',X_cv_ohe_project_grade_category.shape, y_cv.shape)
11 print(project_grade_category_features[:10])
12
```

```
X_train_ohe_project_grade_category (33500, 4) (33500,)
X_test_ohe_project_grade_category (16500, 4) (16500,)
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
```



## clean\_categories-OHE

```
In [13]: 1 vect= CountVectorizer(min_df=10,ngram_range=(1,4), max_features=50000)
2 vect.fit(X_train['clean_categories'].values)
3 X_train_ohe_clean_categories= vect.transform(X_train['clean_categories'].valu
4 # X_cv_ohe_clean_categories= vect.transform(X_cv['clean_categories'].values)
5 X_test_ohe_clean_categories= vect.transform(X_test['clean_categories'].values)
6 clean_categories_features=vect.get_feature_names()
7
8 print('X_train_ohe_clean_categories',X_train_ohe_clean_categories.shape, y_tr
9 print('X_test_ohe_clean_categories',X_test_ohe_clean_categories.shape, y_test
10 # print('X_cv_ohe_clean_categories',X_cv_ohe_clean_categories.shape, y_cv.sha
11 print(clean_categories_features[:10])
12
```

```
X_train_ohe_clean_categories (33500, 40) (33500,)
X_test_ohe_clean_categories (16500, 40) (16500,)
['appliedlearning', 'appliedlearning health_sports', 'appliedlearning history_c
ivics', 'appliedlearning literacy_language', 'appliedlearning math_science', 'a
ppliedlearning music_arts', 'appliedlearning specialneeds', 'health_sports', 'h
ealth_sports appliedlearning', 'health_sports history_civics']
```

## clean\_subcategories-OHE

```
In [14]: 1 vect= CountVectorizer(min_df=10,ngram_range=(1,4), max_features=50000)
2 vect.fit(X_train['clean_subcategories'].values)
3 X_train_ohe_clean_subcategories= vect.transform(X_train['clean_subcategories'
4 # X_cv_ohe_clean_subcategories= vect.transform(X_cv['clean_subcategories'].va
5 X_test_ohe_clean_subcategories= vect.transform(X_test['clean_subcategories'].
6 clean_subcategories_features=vect.get_feature_names()
7
8 print('X_train_ohe_clean_subcategories',X_train_ohe_clean_subcategories.shape
9 print('X_test_ohe_clean_subcategories',X_test_ohe_clean_subcategories.shape,
10 # print('X_cv_ohe_clean_subcategories',X_cv_ohe_clean_subcategories.shape, y_
11 print(clean_subcategories_features[:10])
```

```
X_train_ohe_clean_subcategories (33500, 181) (33500,)
X_test_ohe_clean_subcategories (16500, 181) (16500,)
['appliedsciences', 'appliedsciences charactereducation', 'appliedsciences coll
ege_careerprep', 'appliedsciences earlydevelopment', 'appliedsciences environme
ntalscience', 'appliedsciences esl', 'appliedsciences extracurricular', 'applie
dsciences health_lifescience', 'appliedsciences health_wellness', 'appliedscien
ces history_geography']
```

## Normalize numerical features

### price



```
In [15]: 1 from sklearn.preprocessing import Normalizer
2 norml=Normalizer()
3 norml.fit(X_train['price'].values.reshape(1,-1))
4 X_train_norml_price=norml.transform(X_train['price'].values.reshape(1,-1))
5 X_test_norml_price=norml.transform(X_test['price'].values.reshape(1,-1))
6 # X_cv_norml_price=norml.transform(X_cv['price'].values.reshape(1,-1))
7
8 #X_cv_norm_price=norm.transform(xcv['price'].values.reshape(-1,1))
9 print('X_train_norml_price shape',X_train_norml_price.shape,y_train.shape)
10 print('X_test_norml_price shape',X_test_norml_price.shape,y_test.shape)
11 # print('X_cv_norml_price shape',X_cv_norml_price.shape,y_cv.shape)
12
```

X\_train\_norml\_price shape (1, 33500) (33500,)

X\_test\_norml\_price shape (1, 16500) (16500,)

## teacher\_number\_of\_previously\_posted\_projects

```
In [16]: 1 norml.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
2 X_train_norml_teacher_number_of_previously_posted_projects=norml.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
3 X_test_norml_teacher_number_of_previously_posted_projects=norml.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
4 # X_cv_norml_teacher_number_of_previously_posted_projects=norml.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
5
6
7 # X_cv_norm_price=norm.transform(xcv['price'].values.reshape(-1,1))
8 print('X_train_norml_teacher_number_of_previously_posted_projects shape',X_train_norml_teacher_number_of_previously_posted_projects.shape,y_train.shape)
9 print('X_test_norml_teacher_number_of_previously_posted_projects shape',X_test_norml_teacher_number_of_previously_posted_projects.shape,y_test.shape)
10
```

X\_train\_norml\_teacher\_number\_of\_previously\_posted\_projects shape (1, 33500) (33500,)

X\_test\_norml\_teacher\_number\_of\_previously\_posted\_projects shape (1, 16500) (16500,)

## SET-1

```

In [17]: 1 # https://stackoverflow.com/a/19710648/4084039
2 # combine all features into one single set
3
4 X_train_norml_tnoprepst_projects=X_train_norml_teacher_number_of_previously_p
5 X_test_norml_tnoprepst_projects=X_test_norml_teacher_number_of_previously_pos
6 # X_cv_norml_tnoprepst_projects=X_cv_norml_teacher_number_of_previously_poste
7
8 X_train_norml_price= X_train_norml_price.reshape(33500,1)
9 X_test_norml_price=X_test_norml_price.reshape(16500,1)
10 # X_cv_norml_price=X_cv_norml_price.reshape(11055,1)
11
12 from scipy.sparse import hstack
13 X_tr_set1 = hstack((X_train_essay_tfidf,X_train_norml_price,X_train_norml_tnc
14 X_tst_set1= hstack((X_test_essay_tfidf,X_test_norml_price,X_test_norml_tnopre
15 # X_cr_set1= hstack((X_cv_essay_tfidf,X_cv_norml_price,X_cv_norml_tnoprepst_p
16
17
18 print("Final Data matrix")
19 print(X_tr_set1.shape, y_train.shape)
20 # print(X_cr_set1.shape, y_cv.shape)
21 print(X_tst_set1.shape, y_test.shape)
22

```

```

Final Data matrix
(33500, 50282) (33500,)
(16500, 50282) (16500,)

```

## TFIDF W2V - EASSY

```
In [ ]:
```

```
1
```

```
In [18]:
```

```

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

```

```
In [19]:
```

```

1 tfidf_model = TfidfVectorizer()
2 tfidf_model.fit(data['essay'].values)
3 # we are converting a dictionary with word as a key, and the idf as a value
4 dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_))
5 tfidf_words = set(tfidf_model.get_feature_names())

```



```
In [22]: 1 X_train_norml_tnoprepst_projects=X_train_norml_teacher_number_of_previously_p
2 X_test_norml_tnoprepst_projects=X_test_norml_teacher_number_of_previously_pos
3 # X_cv_norml_tnoprepst_projects=X_cv_norml_teacher_number_of_previously_poste
4
5 X_train_norml_price= X_train_norml_price.reshape(33500,1)
6 X_test_norml_price=X_test_norml_price.reshape(16500,1)
7 # X_cv_norml_price=X_cv_norml_price.reshape(11055,1)
8
9 from scipy.sparse import hstack
10 X_tr_set2 = hstack((tfidf_w2v_vectors_xtrain,X_train_norml_price,X_train_norml_t
11 X_tst_set2= hstack((tfidf_w2v_vectors_xtest,X_test_norml_price,X_test_norml_t
12 # X_cr_set1= hstack((X_cv_essay_tfidf,X_cv_norml_price,X_cv_norml_tnoprepst_p
13
14
15 print("Final Data matrix")
16 print(X_tr_set1.shape, y_train.shape)
17 # print(X_cr_set1.shape, y_cv.shape)
18 print(X_tst_set1.shape, y_test.shape)
```

```
Final Data matrix
(33500, 50282) (33500,)
(16500, 50282) (16500,)
```

## 1.5 Applying Decision Tree on different kind of featurization as mentioned in the instructions

Apply Decision Tree 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

## finding best hyperparameter for set1

```
In [23]: 1
2 from sklearn.model_selection import RandomizedSearchCV
3 import matplotlib.pyplot as plt
4 from sklearn.tree import DecisionTreeClassifier
5
6 from sklearn.metrics import roc_auc_score
7 from sklearn import cross_validation
8 import math
9 import warnings
10 warnings.filterwarnings("ignore")
11
12 model=DecisionTreeClassifier(random_state=0)
13 parameters = {'max_depth':[1,5,10,50], 'min_samples_split':[10,50,100,500]}
14 clf = RandomizedSearchCV(estimator = model,param_distributions = parameters,
15 clf.fit(X_tr_set1, y_train)
16 # clf.estimator
```

C:\Users\aaa\Anaconda3\lib\site-packages\sklearn\cross\_validation.py:41: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model\_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.

"This module will be removed in 0.20.", DeprecationWarning)

```
Out[23]: RandomizedSearchCV(cv=None, error_score='raise',
      estimator=DecisionTreeClassifier(class_weight=None, criterion='gini',
      max_depth=None,
      max_features=None, max_leaf_nodes=None,
      min_impurity_decrease=0.0, min_impurity_split=None,
      min_samples_leaf=1, min_samples_split=2,
      min_weight_fraction_leaf=0.0, presort=False, random_state=0,
      splitter='best'),
      fit_params=None, iid=True, n_iter=10, n_jobs=1,
      param_distributions={'max_depth': [1, 5, 10, 50], 'min_samples_split': [10, 50, 100, 500]},
      pre_dispatch='2*n_jobs', random_state=None, refit=True,
      return_train_score='warn', scoring='roc_auc', verbose=0)
```

```
In [24]: 1 clf.estimator
```

```
Out[24]: DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
      max_features=None, max_leaf_nodes=None,
      min_impurity_decrease=0.0, min_impurity_split=None,
      min_samples_leaf=1, min_samples_split=2,
      min_weight_fraction_leaf=0.0, presort=False, random_state=0,
      splitter='best')
```

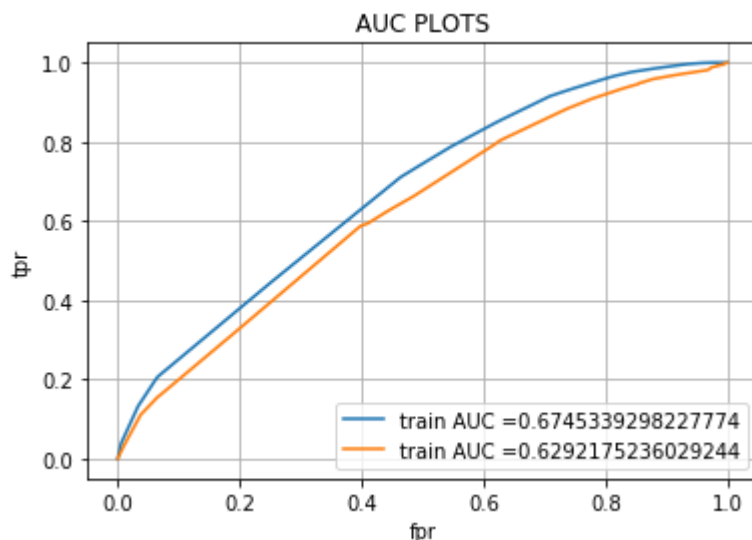
```
In [25]: 1 clf.best_params_
```

```
Out[25]: {'min_samples_split': 500, 'max_depth': 10}
```

```

In [26]: 1 def batch_predict(clf, data):
2         # roc_auc_score(y_true, y_score) the 2nd parameter should be probability
3         # not the predicted outputs
4
5         y_data_pred = []
6         tr_loop = data.shape[0] - data.shape[0]%1000
7         # consider you X_tr shape is 49041, then your tr_loop will be 49041 - 490
8         # in this for loop we will iterate until the last 1000 multiplier
9         for i in range(0, tr_loop, 1000):
10            y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
11            # we will be predicting for the last data points
12            if data.shape[0]%1000 !=0:
13                y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
14
15            return y_data_pred
16
17 clf_set1=DecisionTreeClassifier(max_depth=clf.best_params_['max_depth'],min_s
18 clf_set1.fit(X_tr_set1, y_train)
19 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability esti
20 # not the predicted outputs
21
22 y_train_pred_set1 = batch_predict(clf_set1, X_tr_set1)
23 y_test_pred_set1 = batch_predict(clf_set1, X_tst_set1)
24
25 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_set1)
26 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_set1)
27
28 auc_set1=auc(test_fpr, test_tpr)
29
30 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_t
31 plt.plot(test_fpr, test_tpr, label="train AUC =" +str(auc(test_fpr, test_tpr))
32 plt.legend()
33 plt.xlabel("fpr")
34 plt.ylabel("tpr")
35 plt.title("AUC PLOTS")
36 plt.grid()
37 plt.show()
38
39
40 train_auc=roc_auc_score(y_train,y_train_pred_set1)
41 test_auc=roc_auc_score(y_test,y_test_pred_set1)
42 print("train_Auc:",train_auc)
43 print("test_Auc:",test_auc)

```



train\_Auc: 0.6745339298227774

test\_Auc: 0.6292175236029244

## confusion matrix set1

```
In [27]: 1 def find_best_threshold(threshold, fpr, tpr):
2         t = threshold[np.argmax(tpr*(1-fpr))]
3         # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
4         print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", t)
5         return t
6
7 def predict_with_best_t(proba, threshold):
8     predictions = []
9     for i in proba:
10         if i >= threshold:
11             predictions.append(1)
12         else:
13             predictions.append(0)
14     return predictions
15
16 from sklearn.metrics import confusion_matrix
17 best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
18 # print("Train confusion matrix")
19 # print(confusion_matrix(y_train, predict_with_best_t(y_train_pred_set1, best_t)))
20 print("Test confusion matrix")
21 print(confusion_matrix(y_test, predict_with_best_t(y_test_pred_set1, best_t)))
```

the maximum value of tpr\*(1-fpr) 0.38112492838830747 for threshold 0.868

Test confusion matrix

```
[[1591 1051]
 [5725 8133]]
```

## finding best hyperparameter for set2



```
In [28]: 1 model=DecisionTreeClassifier(random_state=0)
2 parameters = {'max_depth':[1,5,10,50], 'min_samples_split':[10,50,100,500]}
3 clf1 = RandomizedSearchCV(estimator = model,param_distributions = parameters,
4 clf1.fit(X_tr_set2, y_train)
5 # clf.estimator
```

```
Out[28]: RandomizedSearchCV(cv=None, error_score='raise',
      estimator=DecisionTreeClassifier(class_weight=None, criterion='gini',
      max_depth=None,
      max_features=None, max_leaf_nodes=None,
      min_impurity_decrease=0.0, min_impurity_split=None,
      min_samples_leaf=1, min_samples_split=2,
      min_weight_fraction_leaf=0.0, presort=False, random_state=0,
      splitter='best'),
      fit_params=None, iid=True, n_iter=10, n_jobs=1,
      param_distributions={'max_depth': [1, 5, 10, 50], 'min_samples_spli
t': [10, 50, 100, 500]},
      pre_dispatch='2*n_jobs', random_state=None, refit=True,
      return_train_score='warn', scoring='roc_auc', verbose=0)
```

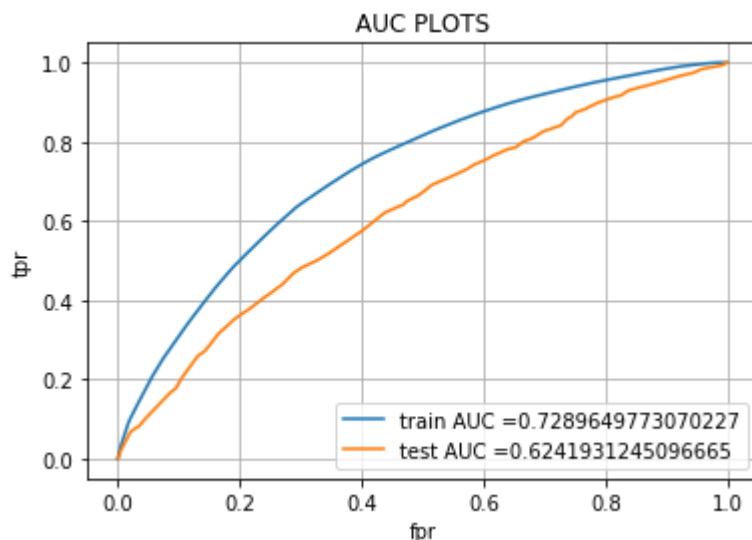
```
In [29]: 1 clf1.best_params_
```

```
Out[29]: {'min_samples_split': 10, 'max_depth': 5}
```

```

In [30]: 1 def batch_predict(clf1, data):
2         # roc_auc_score(y_true, y_score) the 2nd parameter should be probability
3         # not the predicted outputs
4
5         y_data_pred = []
6         tr_loop = data.shape[0] - data.shape[0]%1000
7         # consider you X_tr shape is 49041, then your tr_loop will be 49041 - 490
8         # in this for loop we will iterate until the last 1000 multiplier
9         for i in range(0, tr_loop, 1000):
10            y_data_pred.extend(clf1.predict_proba(data[i:i+1000]))[:,1])
11            # we will be predicting for the last data points
12            if data.shape[0]%1000 !=0:
13                y_data_pred.extend(clf1.predict_proba(data[tr_loop:]))[:,1])
14
15
16        return y_data_pred
17
18 clf_set2=DecisionTreeClassifier(max_depth=clf.best_params_['max_depth'],min_s
19 clf_set2.fit(X_tr_set2, y_train)
20 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability esti
21 # not the predicted outputs
22
23 y_train_pred_set2 = batch_predict(clf_set2, X_tr_set2)
24 y_test_pred_set2 = batch_predict(clf_set2, X_tst_set2)
25
26 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_set2)
27 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_set2)
28
29 # auc_set1=auc(test_fpr, test_tpr)
30
31 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_t
32 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
33 plt.legend()
34 plt.xlabel("fpr")
35 plt.ylabel("tpr")
36 plt.title("AUC PLOTS")
37 plt.grid()
38 plt.show()
39
40 train_auc=roc_auc_score(y_train,y_train_pred_set2)
41 test_auc=roc_auc_score(y_test,y_test_pred_set2)
42 print("train_Auc:",train_auc)
43 print("test_Auc:",test_auc)

```



train\_Auc: 0.7289649773070227

test\_Auc: 0.6241931245096665

## confusion matrix

```
In [31]: 1 def find_best_threshold(threshold, fpr, tpr):
2         t = threshold[np.argmax(tpr*(1-fpr))]
3         # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
4         print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", t)
5         return t
6
7 def predict_with_best_t(proba, threshold):
8     predictions = []
9     for i in proba:
10         if i >= threshold:
11             predictions.append(1)
12         else:
13             predictions.append(0)
14     return predictions
15
16 from sklearn.metrics import confusion_matrix
17 best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
18 print("Train confusion matrix")
19 print(confusion_matrix(y_train, predict_with_best_t(y_train_pred_set2, best_t)))
20 print("Test confusion matrix")
21 print(confusion_matrix(y_test, predict_with_best_t(y_test_pred_set2, best_t)))
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.45199795752439104 for threshold 0.846

Train confusion matrix

```
[[ 3576 1789]
 [ 9056 19079]]
```

Test confusion matrix

```
[[1524 1118]
 [5505 8353]]
```

In [32]: 1 X\_test.shape

Out[32]: (16500, 8)

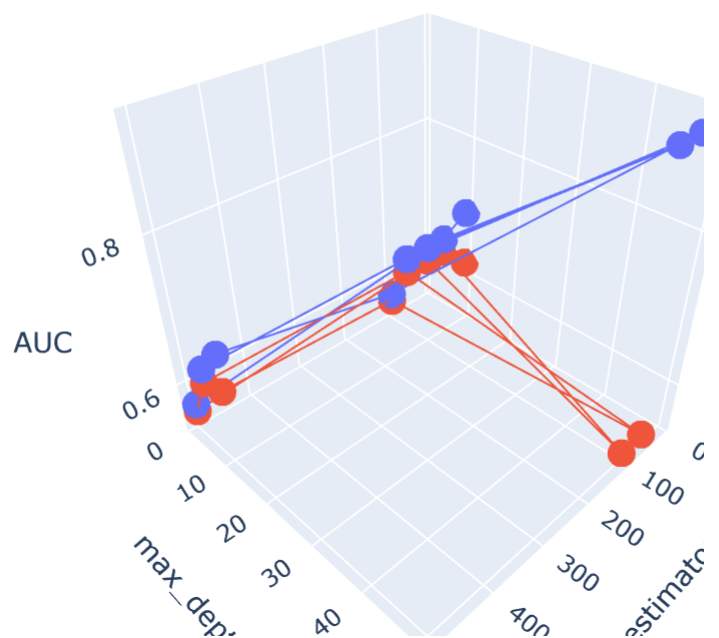
## 3D map for set1

In [33]:

```

1 import plotly.offline as offline
2 import plotly.graph_objs as go
3 offline.init_notebook_mode()
4 import numpy as np
5
6 min_samples_split=clf.cv_results_['param_min_samples_split']
7 max_depth=clf.cv_results_['param_max_depth']
8 train_auc=clf.cv_results_['split0_train_score']
9 test_auc=clf.cv_results_['split0_test_score']
10
11 x = min_samples_split
12 y = max_depth
13 z1 = train_auc
14 z2 = test_auc
15
16 # https://plot.ly/python/3d-axes/
17 trace1 = go.Scatter3d(x=x,y=y,z=z1, name = 'train')
18 trace2 = go.Scatter3d(x=x,y=y,z=z2, name = 'test')
19 data = [trace1,trace2]
20
21 layout = go.Layout(scene = dict(
22     xaxis = dict(title='n_estimators'),
23     yaxis = dict(title='max_depth'),
24     zaxis = dict(title='AUC'),))
25
26 fig = go.Figure(data=data, layout=layout)
27 offline.iplot(fig, filename='3d-scatter-colorscale')
28

```





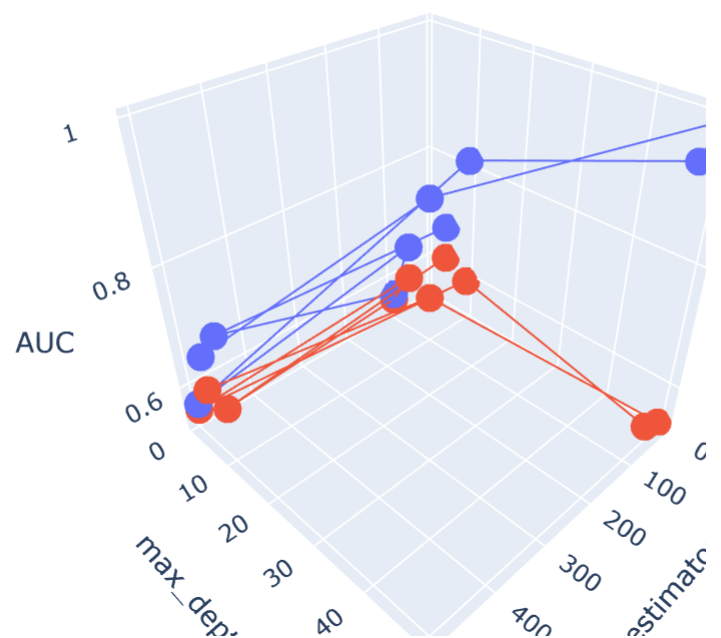
## 3D map for set2

In [34]:

```

1  import plotly.offline as offline
2  import plotly.graph_objs as go
3  offline.init_notebook_mode()
4  import numpy as np
5
6  min_samples_split=clf1.cv_results_['param_min_samples_split']
7  max_depth=clf1.cv_results_['param_max_depth']
8  train_auc=clf1.cv_results_['split0_train_score']
9  test_auc=clf1.cv_results_['split0_test_score']
10
11 x = min_samples_split
12 y = max_depth
13 z1 = train_auc
14 z2 = test_auc
15
16 # https://plot.ly/python/3d-axes/
17 trace1 = go.Scatter3d(x=x,y=y,z=z1, name = 'train')
18 trace2 = go.Scatter3d(x=x,y=y,z=z2, name = 'test')
19 data = [trace1,trace2]
20
21 layout = go.Layout(scene = dict(
22     xaxis = dict(title='n_estimators'),
23     yaxis = dict(title='max_depth'),
24     zaxis = dict(title='AUC'),))
25
26 fig = go.Figure(data=data, layout=layout)
27 offline.iplot(fig, filename='3d-scatter-colorscale')
28

```





## FALSE POSITIVE DATA POINTS SET 1

```
In [35]: 1 cloud1=X_test.copy(deep=True)
          2 cloud2=X_test.copy(deep=True)
```

```
In [ ]: 1
```

```
In [36]: 1 ytp=[]
          2 for i in y_test_pred_set1:
          3     if i<0.5:
          4         ytp.append(0)
          5     else:
          6         ytp.append(1)
          7 cloud1['y_pred'] = ytp
          8 cloud1['y_test'] = y_test
          9 cloud1
         10
         11 # fp_list=[]
         12 # for i in range(len(y_test)):
         13 #     if y_test[i]==0 and y_test_pred_set1[i]==1:
         14 #         fp_list.append(i)
         15 # fp
```

Out[36]:

	school_state	teacher_prefix	project_grade_category	teacher_number_of_previously_posted
36320	wa	mrs	grades_prek_2	
37485	ny	mrs	grades_prek_2	
41428	ca	mrs	grades_prek_2	
30726	va	mr	grades_3_5	

```
In [37]: 1 df1= cloud1[(cloud1['y_pred'] == 1) & (cloud1['y_test']==0)]  
        2 df1
```

Out[37]:

	school_state	teacher_prefix	project_grade_category	teacher_number_of_previously_posted_
30726	va	mr	grades_3_5	
32090	ok	mrs	grades_9_12	
43254	ga	mr	grades_prek_2	
6686	ar	mrs	grades_prek_2	

```
In [38]: 1 df1.shape[0]
```

Out[38]: 2565

## FALSE POSITIVE DATA POINTS SET 2

```

In [39]: 1 # cloud2=X_tst_set2.copy()
2 ytp=[]
3 for i in y_test_pred_set2:
4     if i<0.5:
5         ytp.append(0)
6     else:
7         ytp.append(1)
8 cloud2['y_pred'] = ytp
9 cloud2['y_test'] = y_test
10 cloud2
11
12 df2 = cloud2[(cloud2['y_pred'] == 1) & (cloud2['y_test']==0)]
13 df2
14

```

```

Out[39]:

```

	school_state	teacher_prefix	project_grade_category	teacher_number_of_previously_posted_
30726	va	mr	grades_3_5	
32090	ok	mrs	grades_9_12	
43254	ga	mr	grades_prek_2	
4018	mo	mrs	grades_3_5	

```

In [40]: 1 df2.shape[0]

```

```

Out[40]: 2525

```

## plot wordcloud set 1

```
In [41]: 1 # https://www.geeksforgeeks.org/generating-word-cloud-python/
2 from wordcloud import WordCloud, STOPWORDS
3 comment_words = ''
4 stopwords = set(STOPWORDS)
5
6 # iterate through the csv file
7 for val in df1.essay:
8
9     # typecaste each val to string
10    val = str(val)
11
12    # split the value
13    tokens = val.split()
14
15    # Converts each token into lowercase
16    for i in range(len(tokens)):
17        tokens[i] = tokens[i].lower()
18
19    comment_words += " ".join(tokens)+" "
20
21
22 wordcloud = WordCloud(width = 800, height = 800,
23                       background_color = 'white',
24                       stopwords = stopwords,
25                       min_font_size = 10).generate(comment_words)
26
27 # plot the WordCloud image
28 plt.figure(figsize = (8, 8), facecolor = None)
29 plt.imshow(wordcloud)
30 plt.axis("off")
31 plt.tight_layout(pad = 0)
32
33 plt.show()
```



In [42]:

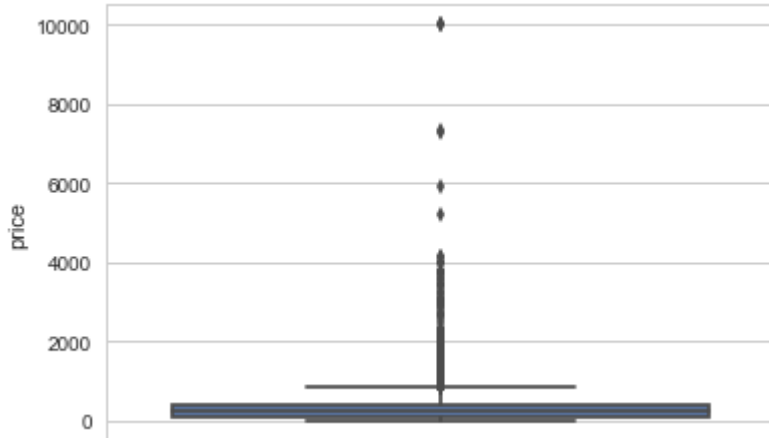
```
1  # https://www.geeksforgeeks.org/generating-word-cloud-python/
2  from wordcloud import WordCloud, STOPWORDS
3  comment_words = ''
4  stopwords = set(STOPWORDS)
5
6  # iterate through the csv file
7  for val in df2.essay:
8
9      # typecaste each val to string
10     val = str(val)
11
12     # split the value
13     tokens = val.split()
14
15     # Converts each token into lowercase
16     for i in range(len(tokens)):
17         tokens[i] = tokens[i].lower()
18
19     comment_words += " ".join(tokens)+" "
20
21
22 wordcloud = WordCloud(width = 800, height = 800,
23                       background_color = 'white',
24                       stopwords = stopwords,
25                       min_font_size = 10).generate(comment_words)
26
27 # plot the WordCloud image
28 plt.figure(figsize = (8, 8), facecolor = None)
29 plt.imshow(wordcloud)
30 plt.axis("off")
31 plt.tight_layout(pad = 0)
32
33 plt.show()
```



## Plot the box plot with the price of these false positive data points set1

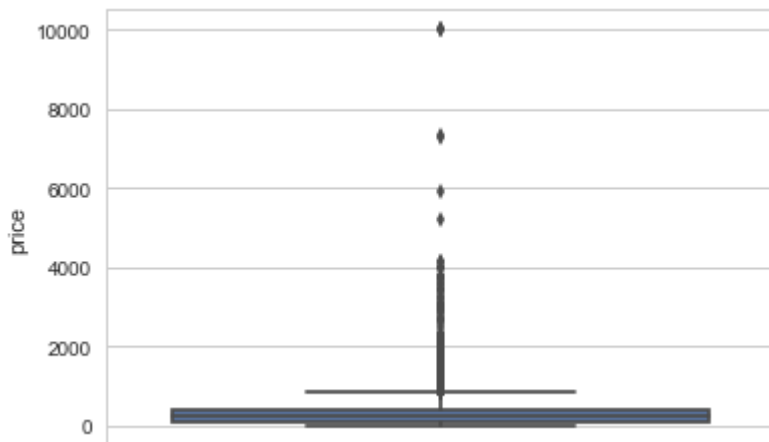


```
In [43]: 1 import seaborn as sns
2 sns.set(style='whitegrid')
3 ax = sns.boxplot(y=cloud1["price"])
4
```



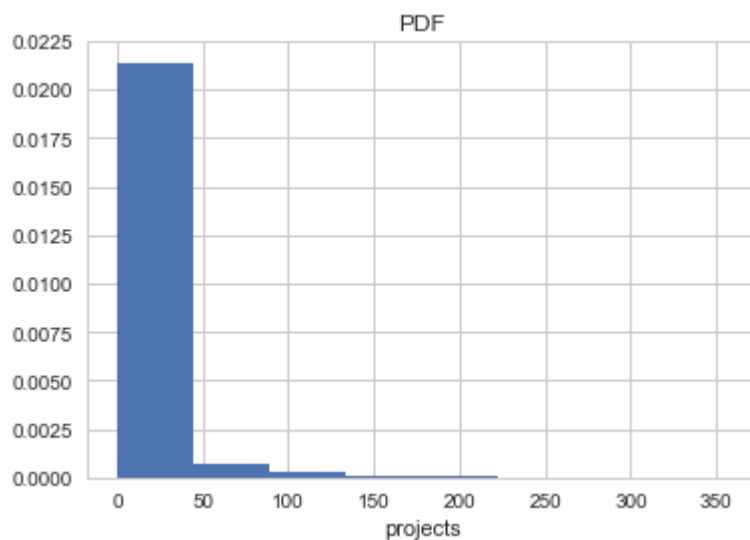
**Plot the box plot with the price of these false positive data points set 2**

```
In [44]: 1 import seaborn as sns
2 sns.set(style='whitegrid')
3 ax = sns.boxplot(y=cloud2["price"])
4
```



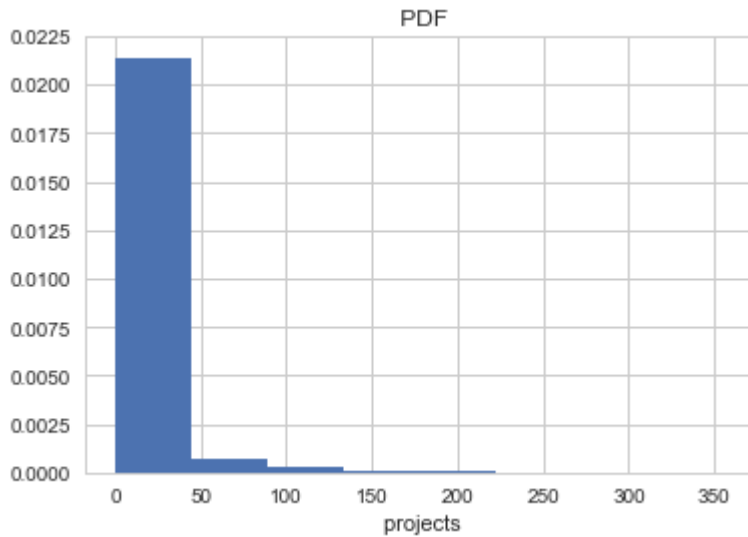
**Plot the pdf set1**

```
In [45]: 1 plt.hist(cloud1['teacher_number_of_previously_posted_projects'], bins=8, dens
2          # pdf = counts/(sum(counts))
3          # print(pdf)
4          plt.title('PDF')
5          plt.xlabel('projects')
6          plt.show()
7          # plt.plot(bin_edges[1:], cdf)
```



## Plot the pdf set2

```
In [46]: 1 plt.hist(cloud2['teacher_number_of_previously_posted_projects'], bins=8, dens
2 # pdf = counts/(sum(counts))
3 # print(pdf)
4 plt.title(' PDF')
5 plt.xlabel(' projects')
6 plt.show()
7 # plt.plot(bin_edges[1:], cdf)
```



```
In [ ]:
```

```
1
```

## 1.6 Getting top features using feature\_importances\_

### TASK 2

```
In [47]: 1 imp=[]
2 for i in clf_set1.feature_importances_:
3     if i!=0.0:
4         imp.append(i)
5 imp
6 len(imp)
7
8
9 X_train_task2=X_tr_set1[:,imp]
10 X_test_task2=X_tst_set1[:,imp]
11
12 print(X_train_task2.shape)
13 print(X_test_task2.shape)
14 print(y_train.shape)
```

```
(33500, 61)
```

```
(16500, 61)
```

```
(33500,)
```

```
In [48]: 1
          2
          3 model=DecisionTreeClassifier(random_state=0)
          4 parameters = {'max_depth':[1,5,10,50], 'min_samples_split':[10,50,100,500]}
          5 clf_featureimp = RandomizedSearchCV(estimator = model,param_distributions = p
          6 clf_featureimp.fit(X_train_task2,y_train)
          7
          8
```

```
Out[48]: RandomizedSearchCV(cv=None, error_score='raise',
                             estimator=DecisionTreeClassifier(class_weight=None, criterion='gini',
                             max_depth=None,
                             max_features=None, max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=1, min_samples_split=2,
                             min_weight_fraction_leaf=0.0, presort=False, random_state=0,
                             splitter='best'),
                             fit_params=None, iid=True, n_iter=10, n_jobs=1,
                             param_distributions={'max_depth': [1, 5, 10, 50], 'min_samples_split': [10, 50, 100, 500]},
                             pre_dispatch='2*n_jobs', random_state=None, refit=True,
                             return_train_score='warn', scoring='roc_auc', verbose=0)
```

```
In [49]: 1 clf_featureimp.best_params_
          2
```

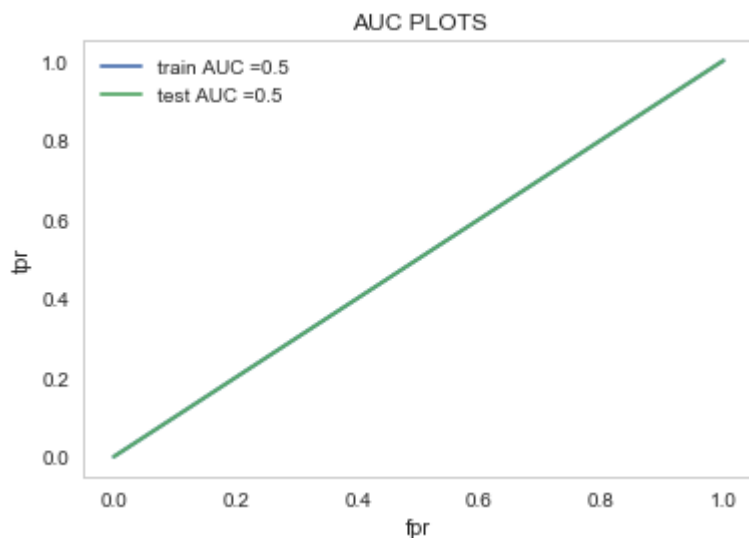
```
Out[49]: {'min_samples_split': 50, 'max_depth': 1}
```

## Decision tree for hyperparameter tuning

```

In [50]: 1 def batch_predict(clf, data):
2         # roc_auc_score(y_true, y_score) the 2nd parameter should be probability
3         # not the predicted outputs
4
5         y_data_pred = []
6         tr_loop = data.shape[0] - data.shape[0]%1000
7         # consider you X_tr shape is 49041, then your tr_loop will be 49041 - 490
8         # in this for loop we will iterate until the last 1000 multiplier
9
10
11        for i in range(0, tr_loop, 1000):
12            y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
13            # we will be predicting for the last data points
14            if data.shape[0]%1000 !=0:
15                y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
16
17        return y_data_pred
18 from sklearn.metrics import auc
19
20 best_model=DecisionTreeClassifier(max_depth=clf.featureimp.best_params_['max_
21 best_model.fit(X_train_task2, y_train)
22 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability esti
23 # not the predicted outputs
24
25 y_train_pred = best_model.predict( X_train_task2)
26 y_test_pred = best_model.predict( X_test_task2)
27
28
29 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
30 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
31
32
33
34 # auc=auc(test_fpr, test_tpr)
35
36 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_t
37 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
38 plt.legend()
39 plt.xlabel("fpr")
40 plt.ylabel("tpr")
41 plt.title("AUC PLOTS")
42 plt.grid()
43
44
45 plt.show()
46
47
48 train_auc=roc_auc_score(y_train,y_train_pred)
49 test_auc=roc_auc_score(y_test,y_test_pred)
50 print("train_Auc:",train_auc)
51 print("test_Auc:",test_auc)

```



train\_Auc: 0.5

test\_Auc: 0.5

## confusion matrix for task 2

```
In [51]: 1 def find_best_threshold(threshold, fpr, tpr):
2         t = threshold[np.argmax(tpr*(1-fpr))]
3         # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very h
4         print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for thresho
5         return t
6
7 def predict_with_best_t(proba, threshold):
8     predictions = []
9     for i in proba:
10        if i>=threshold:
11            predictions.append(1)
12        else:
13            predictions.append(0)
14    return predictions
15
16 from sklearn.metrics import confusion_matrix
17 best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
18 print("Train confusion matrix")
19 print(confusion_matrix(y_train, predict_with_best_t(y_train_pred_set2, best_t
20 print("Test confusion matrix")
21 print(confusion_matrix(y_test, predict_with_best_t(y_test_pred_set2, best_t))
```

the maximum value of tpr\*(1-fpr) 0.0 for threshold 2

Train confusion matrix

```
[[ 5365    0]
 [28135    0]]
```

Test confusion matrix

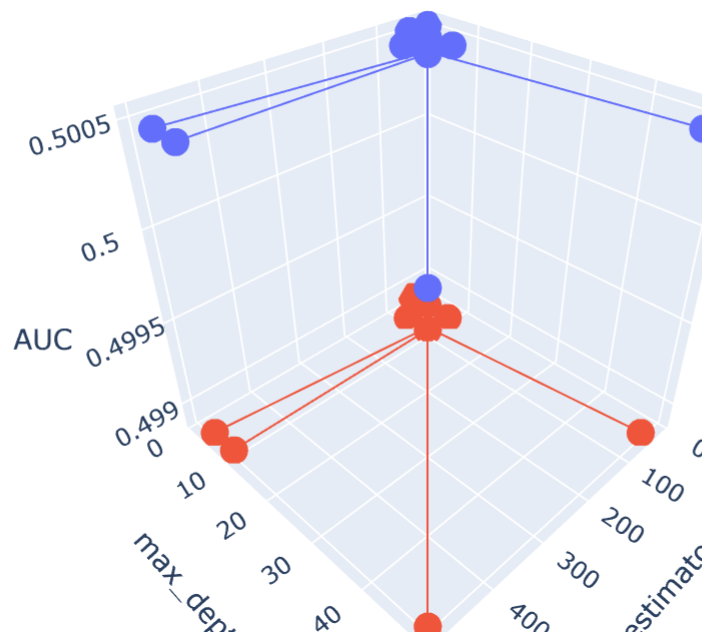
```
[[ 2642    0]
 [13858    0]]
```

## 3D map for task 2

```

In [52]: 1 import plotly.offline as offline
2 import plotly.graph_objs as go
3 offline.init_notebook_mode()
4 import numpy as np
5
6 min_samples_split=clf_featureimp.cv_results_['param_min_samples_split']
7 max_depth=clf_featureimp.cv_results_['param_max_depth']
8 train_auc=clf_featureimp.cv_results_['split0_train_score']
9 test_auc=clf_featureimp.cv_results_['split0_test_score']
10
11 x = min_samples_split
12 y = max_depth
13 z1 = train_auc
14 z2 = test_auc
15
16 # https://plot.ly/python/3d-axes/
17 trace1 = go.Scatter3d(x=x,y=y,z=z1, name = 'train')
18 trace2 = go.Scatter3d(x=x,y=y,z=z2, name = 'test')
19 data = [trace1,trace2]
20
21 layout = go.Layout(scene = dict(
22     xaxis = dict(title='n_estimators'),
23     yaxis = dict(title='max_depth'),
24     zaxis = dict(title='AUC'),))
25
26 fig = go.Figure(data=data, layout=layout)
27 offline.iplot(fig, filename='3d-scatter-colorscale')
28

```





# fp for task 2

```
In [53]: 1 X_test_task2.shape
```

Out[53]: (16500, 61)

```
In [54]: 1 cloud3=X_test.copy()  
2 cloud3
```

Out[54]:

	school_state	teacher_prefix	project_grade_category	teacher_number_of_previously_posted_
36320	wa	mrs	grades_prek_2	
37485	ny	mrs	grades_prek_2	
41428	ca	mrs	grades_prek_2	
30726	va	mr	grades_3_5	

```
In [55]: 1 # cloud3=X_test_task2.copy()
2
3 ytp=[]
4 for i in y_test_pred:
5     if i<0.5:
6         ytp.append(0)
7     else:
8         ytp.append(1)
9 cloud3['y_pred'] = ytp
10 cloud3['y_test'] = y_test
11 cloud3
12
13 df3 = cloud3[(cloud3['y_pred'] == 1) & (cloud3['y_test']==0)]
14 df3
```

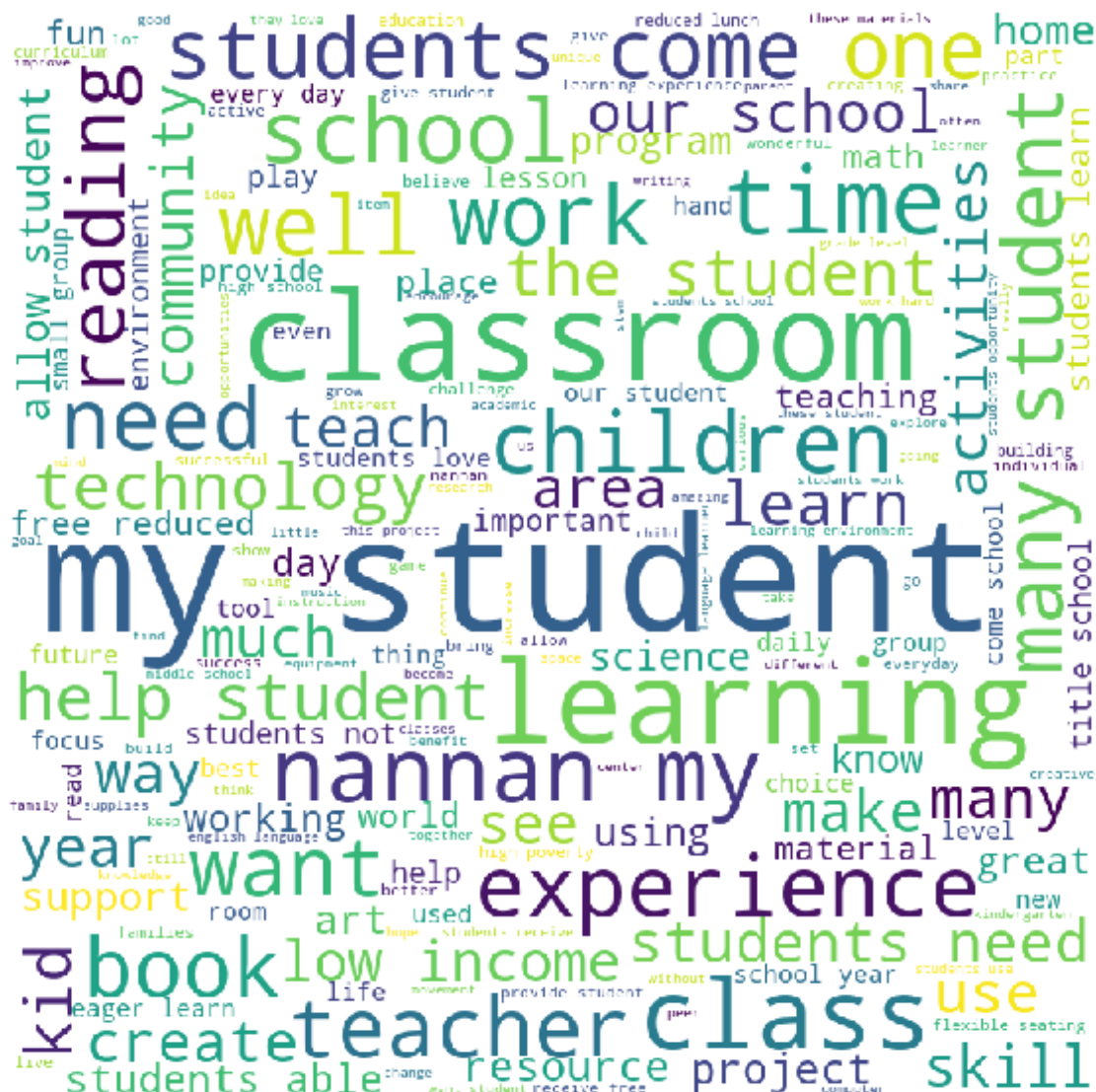
```
Out[55]:
```

	school_state	teacher_prefix	project_grade_category	teacher_number_of_previously_posted_
30726	va	mr	grades_3_5	
32090	ok	mrs	grades_9_12	
43254	ga	mr	grades_prek_2	
6686	ar	mrs	grades_prek_2	

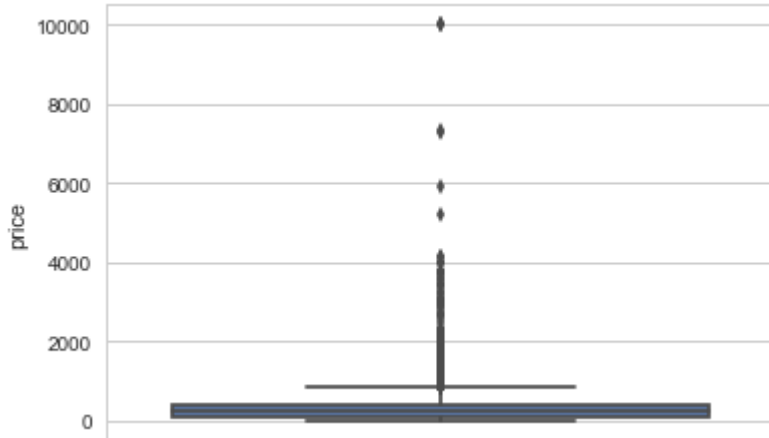
## wordcloud for task 2

In [56]:

```
1  # https://www.geeksforgeeks.org/generating-word-cloud-python/
2  from wordcloud import WordCloud, STOPWORDS
3  comment_words = ''
4  stopwords = set(STOPWORDS)
5
6  # iterate through the csv file
7  for val in df3.essay:
8
9      # typecaste each val to string
10     val = str(val)
11
12     # split the value
13     tokens = val.split()
14
15     # Converts each token into lowercase
16     for i in range(len(tokens)):
17         tokens[i] = tokens[i].lower()
18
19     comment_words += " ".join(tokens)+" "
20
21
22 wordcloud = WordCloud(width = 800, height = 800,
23                       background_color = 'white',
24                       stopwords = stopwords,
25                       min_font_size = 10).generate(comment_words)
26
27 # plot the WordCloud image
28 plt.figure(figsize = (8, 8), facecolor = None)
29 plt.imshow(wordcloud)
30 plt.axis("off")
31 plt.tight_layout(pad = 0)
32
33 plt.show()
```

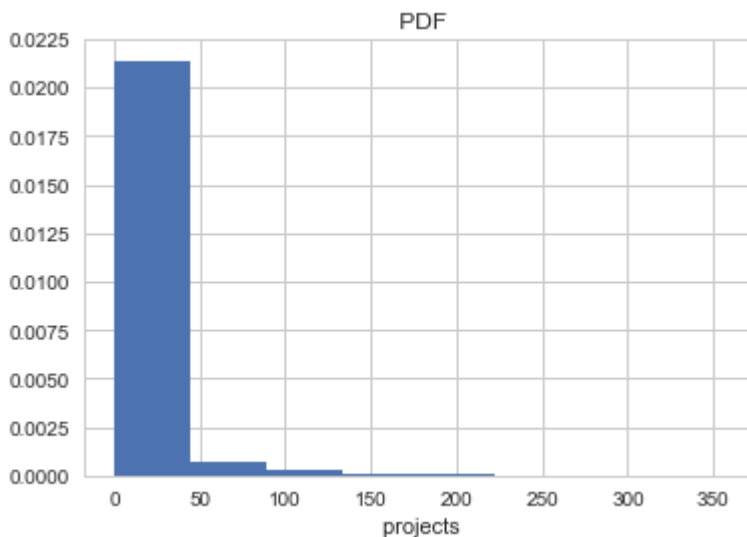


```
In [57]: 1 import seaborn as sns
2 sns.set(style='whitegrid')
3 ax = sns.boxplot(y=cloud3["price"])
4
```



## PDF for task 2

```
In [58]: 1 plt.hist(cloud3['teacher_number_of_previously_posted_projects'], bins=8, dens
2 # pdf = counts/(sum(counts))
3 # print(pdf)
4 plt.title(' PDF')
5 plt.xlabel(' projects')
6 plt.show()
7 # plt.plot(bin_edges[1:], cdf)
```



## 2. Summary

```
In [59]: 1 from prettytable import PrettyTable
2 pryt=PrettyTable()
3 pryt.field_names=['vectorizer','model','best hyper params','train-AUC','test-
4 pryt.add_row(['TF-IDF','Decision Tree',clf.best_params_,train_auc,test_auc])
5 pryt.add_row(['TF-IDF W2V','Decision Tree',clf1.best_params_,train_auc,test_a
6 pryt.add_row(['TF-IDF','Decision Tree',clf_featureimp.best_params_,train_auc,
7 print(pryt)
8
```

```
+-----+-----+-----+-----+
| vectorizer |      model      |      best hyper params      |      test-AUC      |
+-----+-----+-----+-----+
| TF-IDF     | Decision Tree   | {'min_samples_split': 500, 'max_depth': 10} | [0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 |
| 0.50050704 0.50050704 0.50050704 0.50050704 0.50050704 0.50050704 | 0.49889657
| 0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 |
| TF-IDF W2V | Decision Tree   | {'min_samples_split': 10, 'max_depth': 5}   | [0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 |
| 0.50050704 0.50050704 0.50050704 0.50050704 0.50050704 0.50050704 | 0.49889657
| 0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 |
| TF-IDF     | Decision Tree   | {'min_samples_split': 50, 'max_depth': 1}    | [0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 |
| 0.50050704 0.50050704 0.50050704 0.50050704 0.50050704 0.50050704 | 0.49889657
| 0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 0.49889657 |
+-----+-----+-----+-----+
```

```
In [ ]: 1
```