

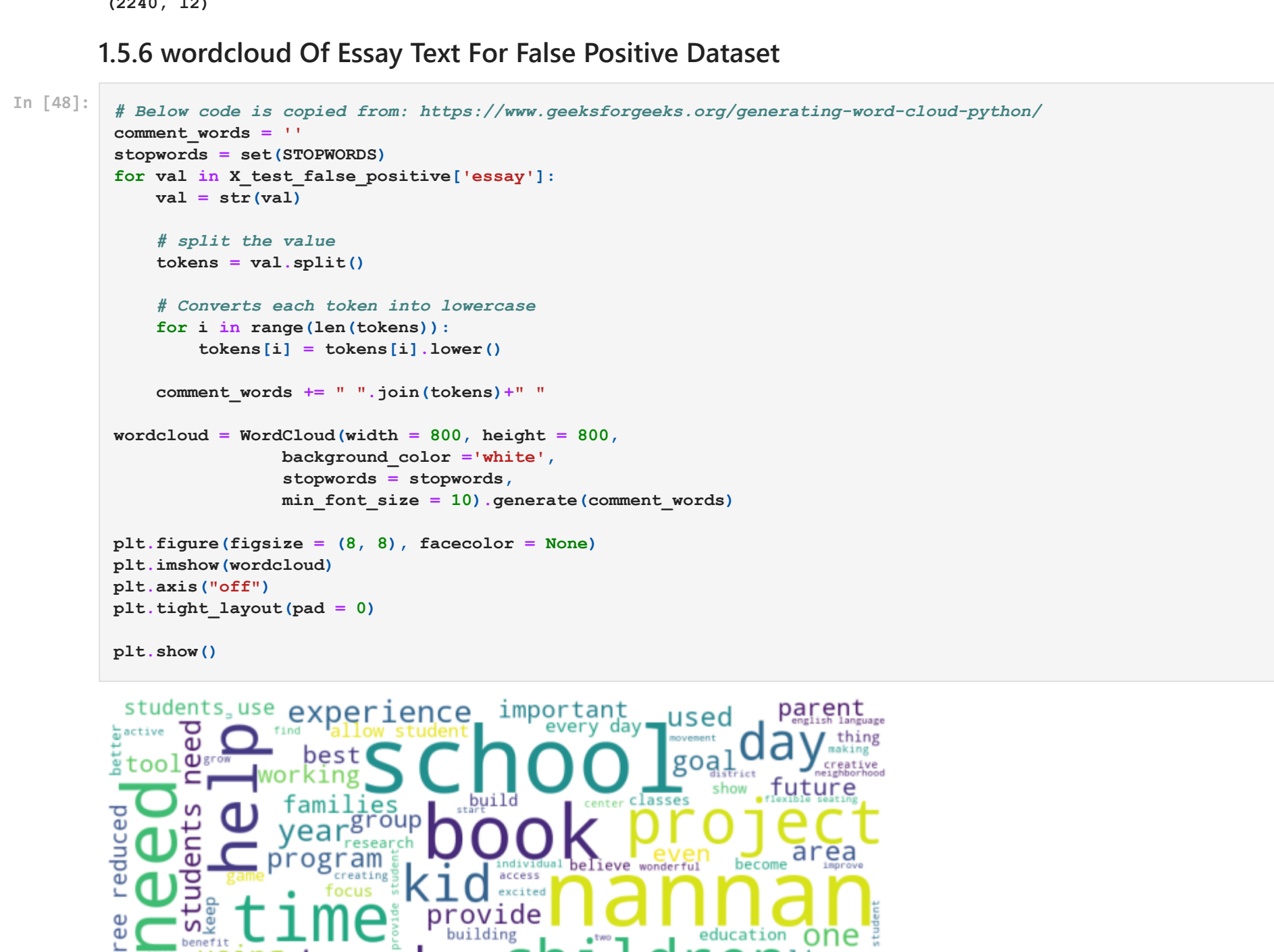
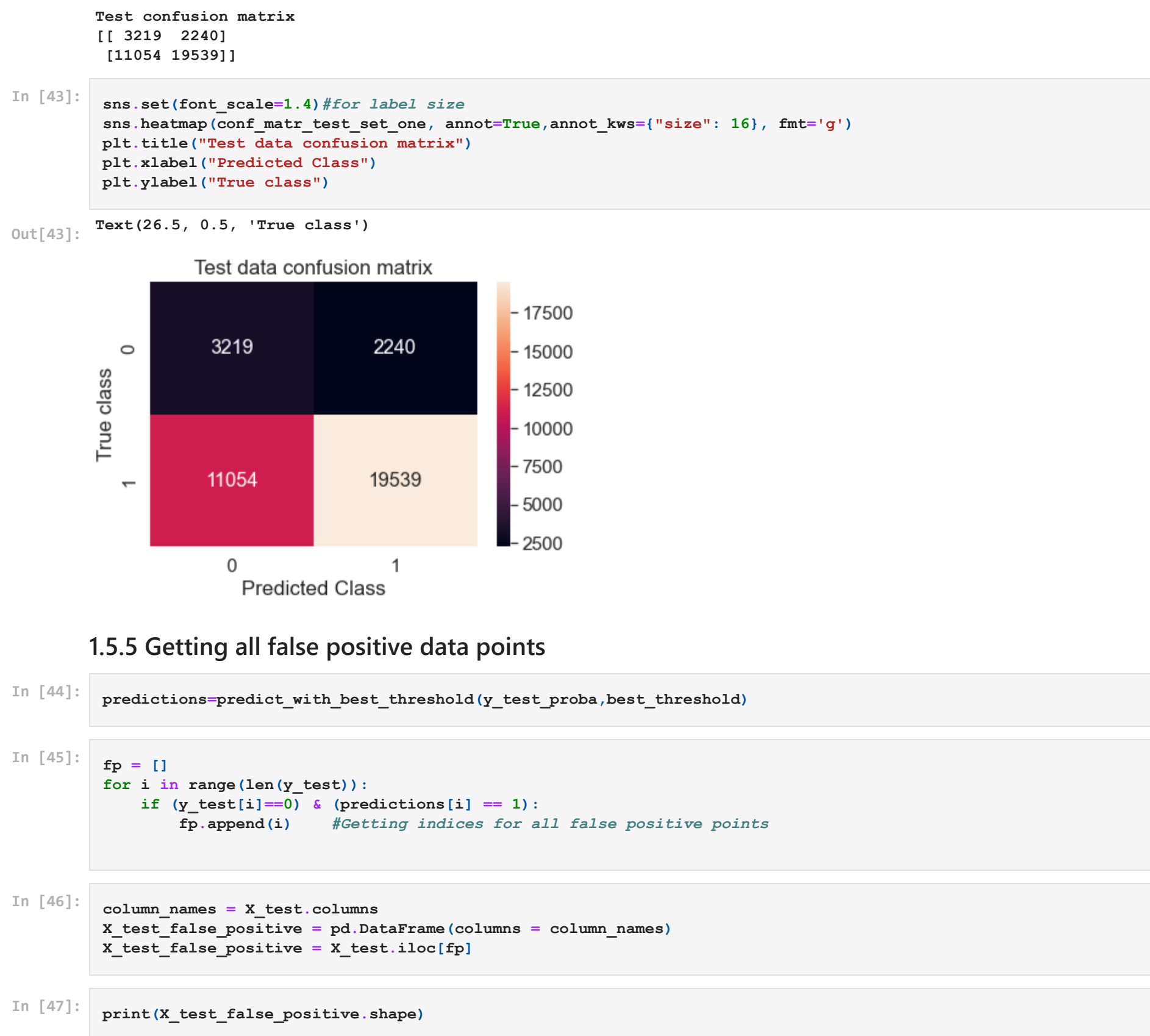
1.5.4 Confusion Matrix for Train and Test data

```
In [38]: def find_best_threshold(probs, threshold, fpr, tpr):
        t = threshold*np.argmax((tpr*(1-fpr)))
        # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
        print("the maximum value of fpr*(1-tpr) = ", max(tpr*(1-fpr)), "for threshold, np.round(t,3))
        return t

In [39]: def predict_with_best_threshold(probs, t):
        predictions = []
        for i in probs:
            if i>=t:
                predictions.append(1)
            else:
                predictions.append(0)
        return predictions

In [40]: print("t==100")
best_threshold = find_best_threshold(y_train_proba, tr_thresholds, train_fpr, train_tpr)
conf_matr_train_set_one = confusion_matrix(y_train, predict_with_best_threshold(y_train_proba, best_threshold))
print("Train confusion matrix")
print(conf_matr_train_set_one)

=====
the maximum value of fpr*(1-tpr) = 0.4027835324942495 for threshold 0.855
Train confusion matrix
[[ 3219 2240]
 [11054 4189]]
(21899 40214)
```



1.5.5 Getting all false positive data points

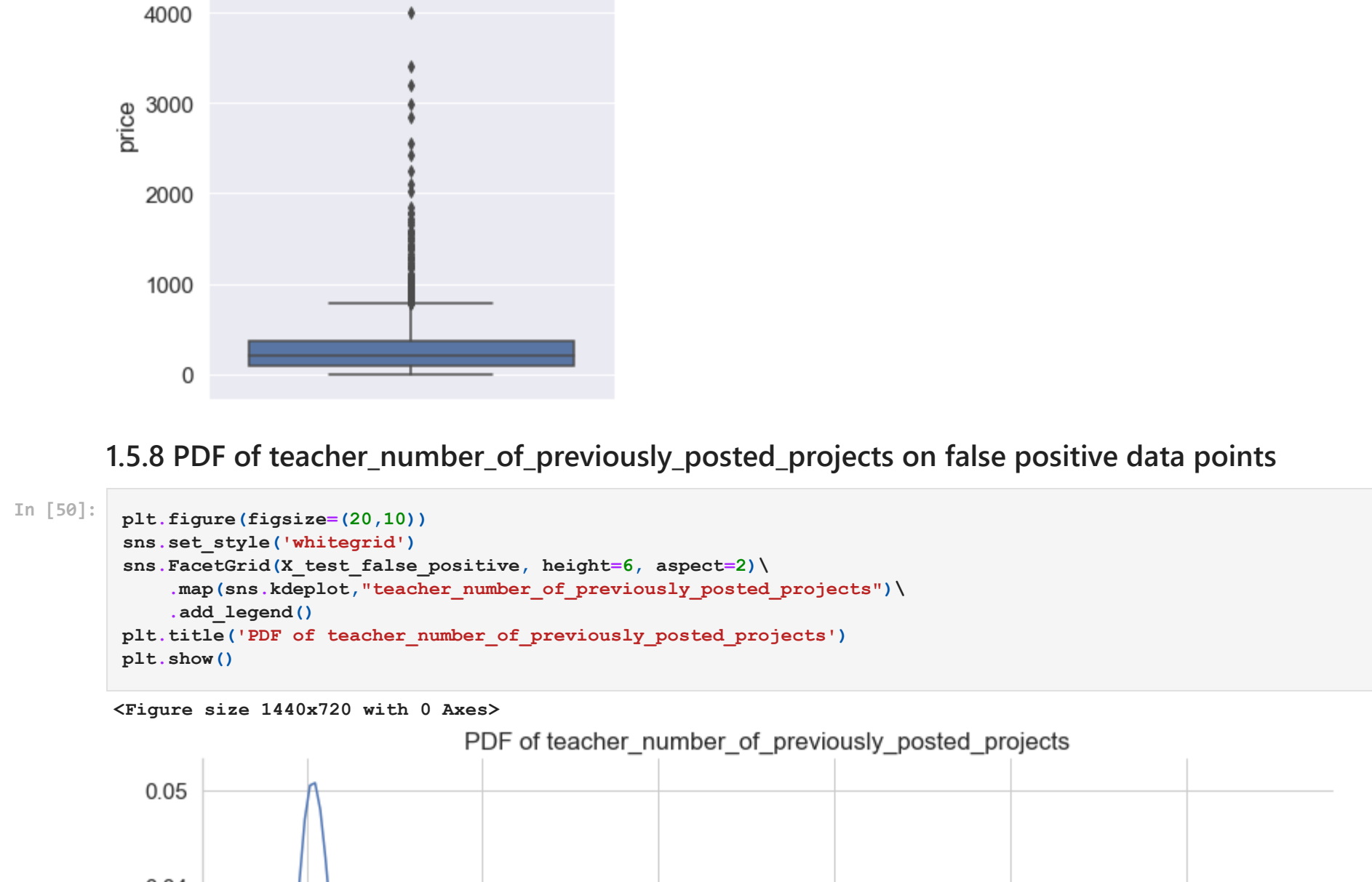
```
In [44]: predictions=predict_with_best_threshold(y_test_proba,best_threshold)

In [45]: fp = []
        for i in range(len(y_test)):
            if (y_test[i]==0) & (predictions[i] == 1):
                fp.append(i) #Getting indices for all false positive points

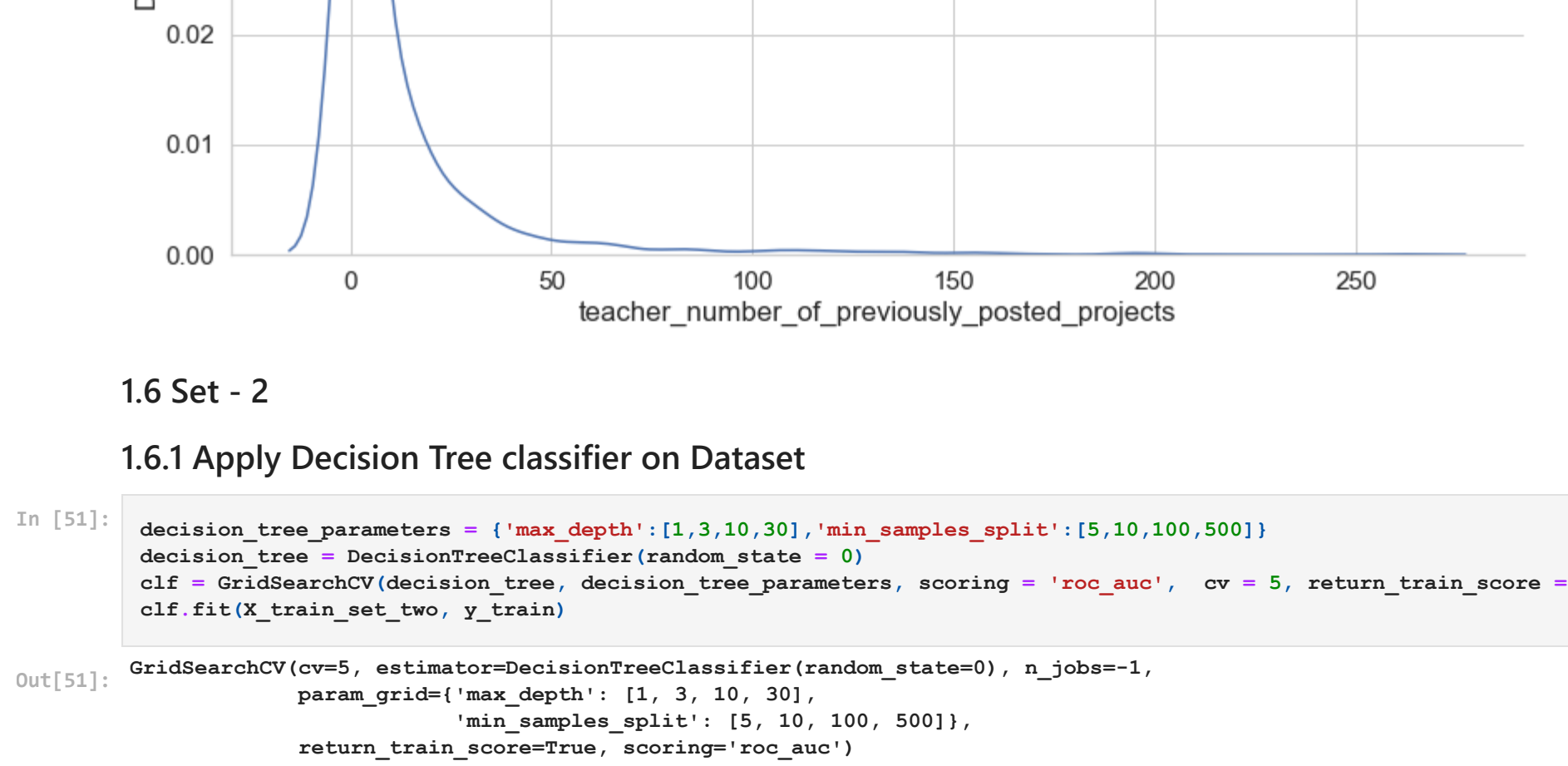
In [46]: column_names = X_test.columns
        X_test_false_positive = pd.DataFrame(columns = column_names)
        X_test_false_positive = X_test.iloc[fp]

In [47]: print(X_test_false_positive.shape)
(2240, 12)
```

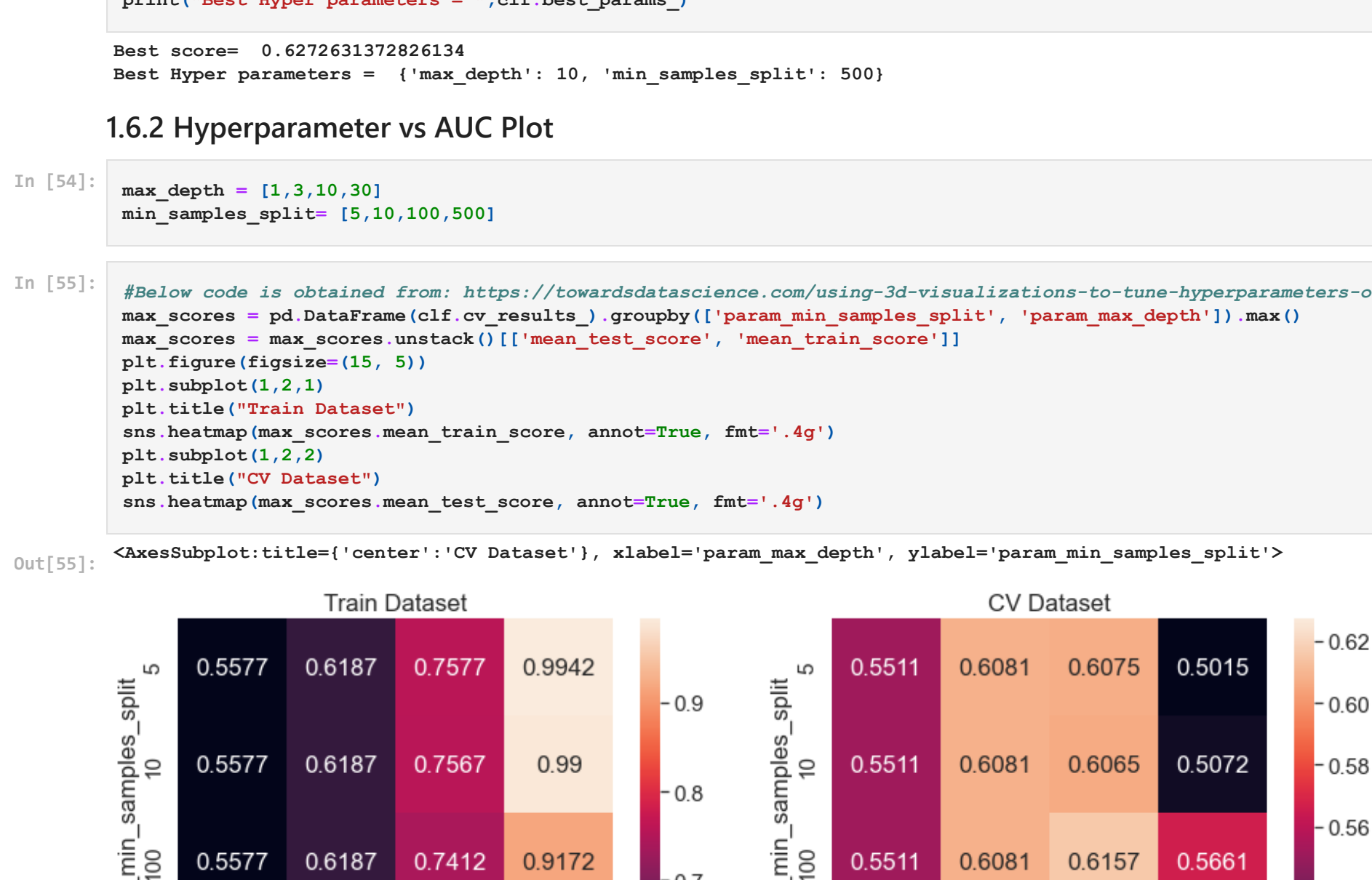
1.5.6 wordcloud Of Essay Text For False Positive Dataset



1.5.7 Box plot of price on false positive data points



1.5.8 PDF of teacher_number_of_previously_posted_projects on false positive data points



1.6 Set - 2

1.6.1 Apply Decision Tree classifier on Dataset

```
In [51]: decision_tree_parameters = {'max_depth':[1,3,10,30], 'min_samples_split':[5,10,100,500]}
        decision_tree = DecisionTreeClassifier(random_state = 0)
        clf = GridSearchCV(decision_tree, decision_tree_parameters, scoring = 'roc_auc', cv = 5, return_train_score =
        GridSearchCV(cv=5, estimator=DecisionTreeClassifier(random_state=0), n_jobs=-1,
                    param_grid={'max_depth': [1, 3, 10, 30],
                                'min_samples_split': [5, 10, 100, 500]},
                    return_train_score=True, scoring='roc_auc')

In [52]: train_auc = clf.cv_results_['mean_train_score']
        cv_auc = clf.cv_results_['mean_test_score']

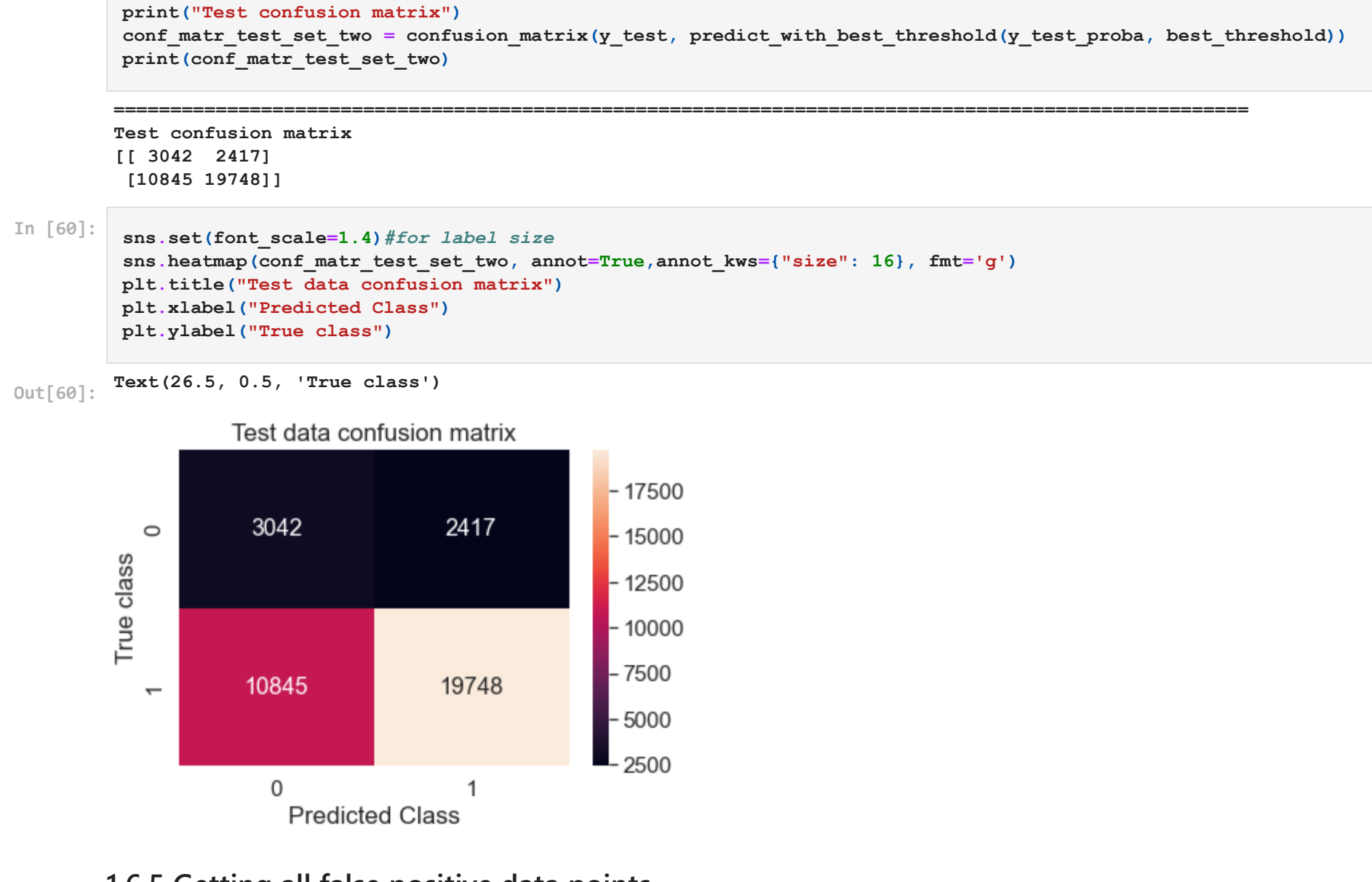
In [53]: print("Best score = ",clf.best_score_)
        print("Best Hyper parameters = ",clf.best_params_)

Best score = 0.6272631372826134
Best Hyper parameters = {'max_depth': 10, 'min_samples_split': 500}
```

1.6.2 Hyperparameter vs AUC Plot



1.6.3 Plot ROC-AUC curve of model with best hyperparameters



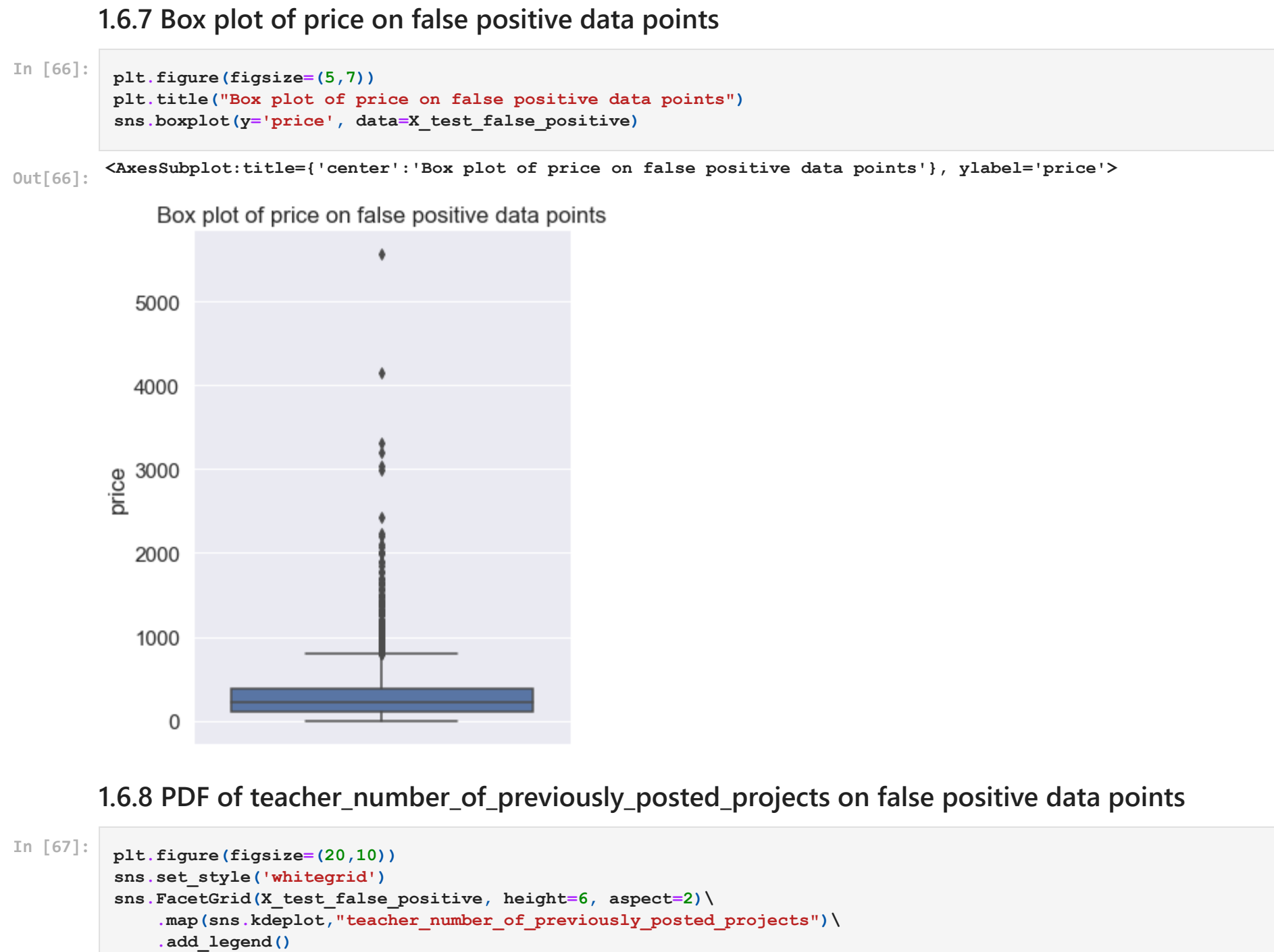
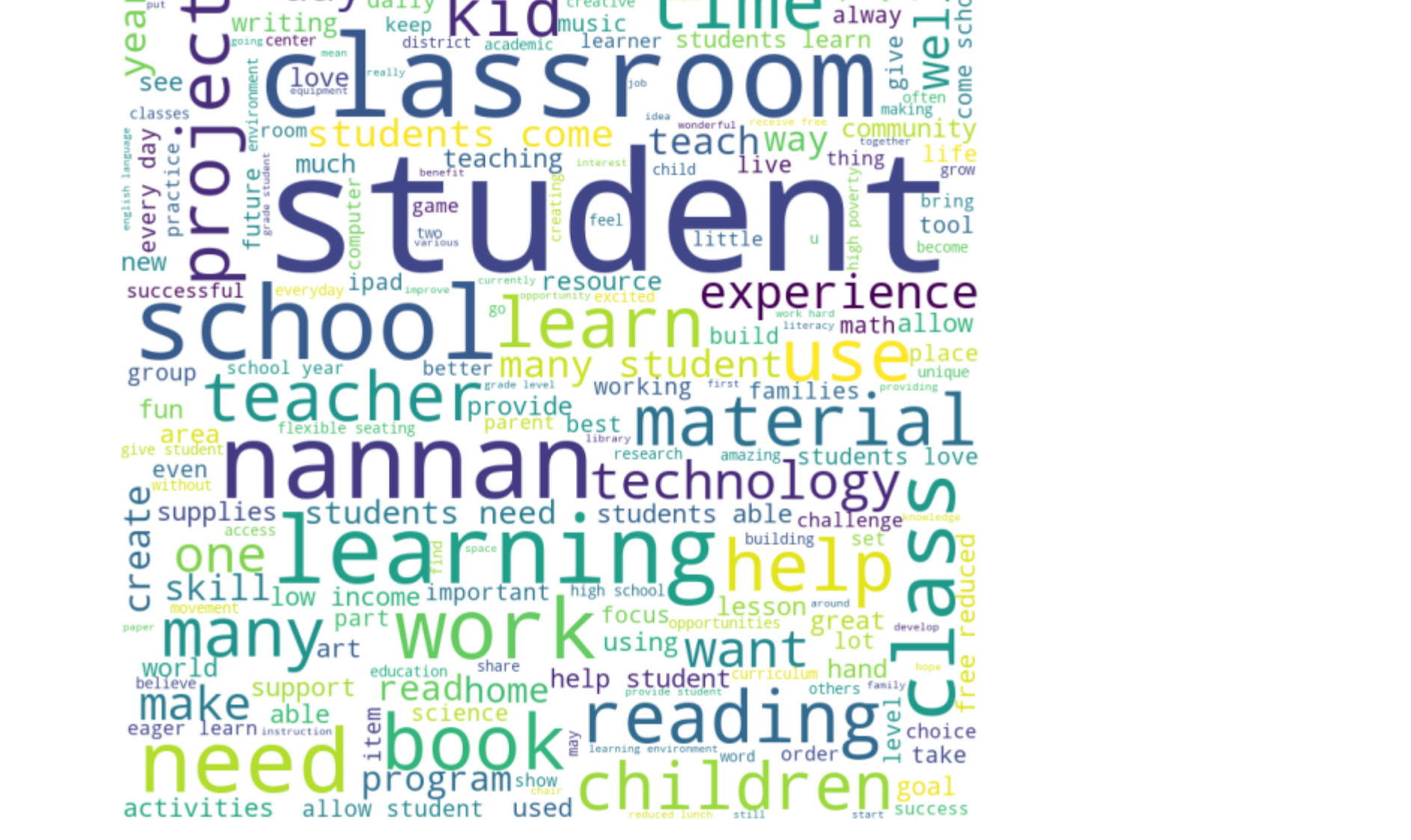
1.6.4 Confusion Matrix for Train and Test data

```
In [57]: print("t==100")
best_threshold = find_best_threshold(y_train_proba, tr_thresholds, train_fpr, train_tpr)
conf_matr_train_set_two = confusion_matrix(y_train, predict_with_best_threshold(y_train_proba, best_threshold))
print("Train confusion matrix")
print(conf_matr_train_set_two)

=====
the maximum value of fpr*(1-tpr) = 0.43386367416184746 for threshold 0.859
Train confusion matrix
[[ 1223 3860]
 [20763 41350]]
(22763 41350)

In [58]: sns.set(font_scale=1.4)#For label size
sns heatmap(conf_matr_train_set_two, annot=True,annot_kws={"size": 16}, fmt='g')
plt title("Train data confusion matrix")
plt xlabel("Predicted Class")
plt ylabel("True Class")

Out[58]: Text(26.5, 0.5, 'True class')
```



1.6.5 Getting all false positive data points

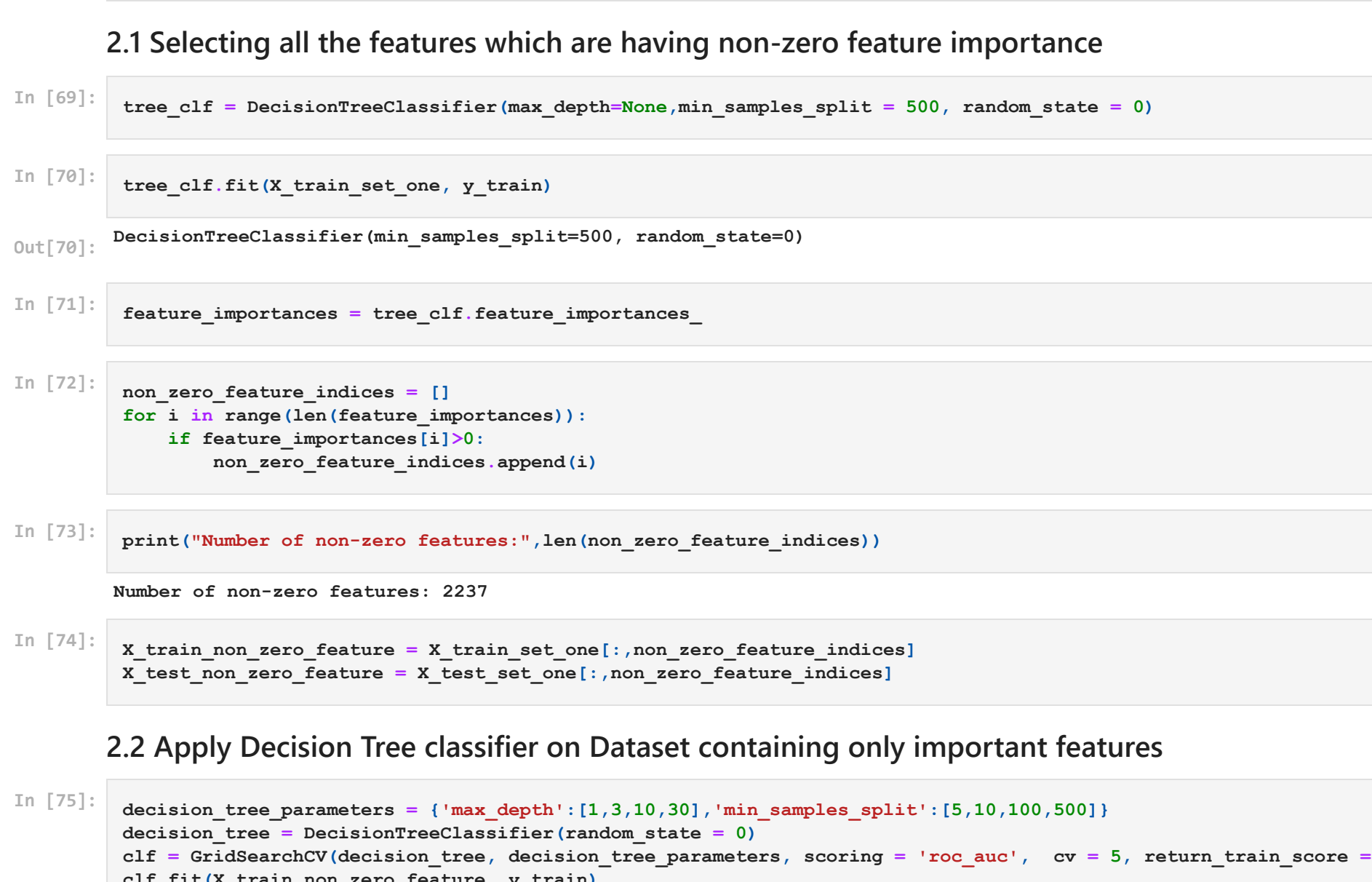
```
In [61]: predictions=predict_with_best_threshold(y_test_proba,best_threshold)

In [62]: fp = []
        for i in range(len(y_test)):
            if (y_test[i]==0) & (predictions[i] == 1):
                fp.append(i) #Getting indices for all false positive points

In [63]: column_names = X_test.columns
        X_test_false_positive = pd.DataFrame(columns = column_names)
        X_test_false_positive = X_test.iloc[fp]

In [64]: print(X_test_false_positive.shape)
(2417, 12)
```

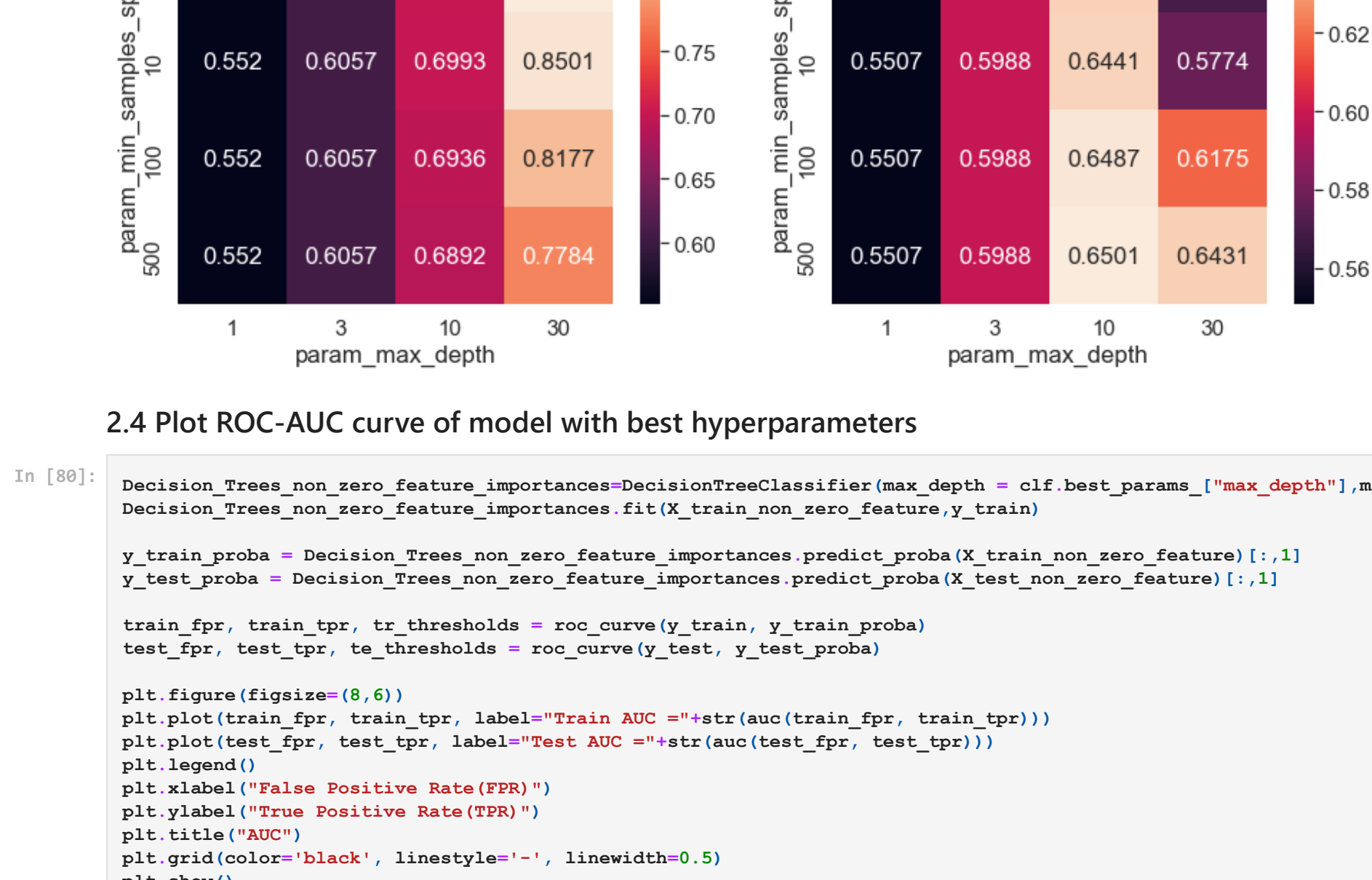
1.6.6 wordcloud Of Essay Text For False Positive Dataset



1.6.7 Box plot of price on false positive data points



1.6.8 PDF of teacher_number_of_previously_posted_projects on false positive data points



Task - 2

```
In [68]: # 1. write your code in following steps for task 2
        # 2. select all non zero features
        # 3. Update your dataset i.e. X_train,X_test and X_cv so that it contains all rows and only non zero features
        # 4. perform hyperparameter tuning and plot either heatmap or 3d plot
        # 5. fit the best model. Plot ROC AUC curve and confusion matrix similar to rowd 1.
```

2.1 Selecting all the features which are having non-zero feature importance

```
In [69]: tree_clf = DecisionTreeClassifier(max_depth=None,min_samples_split = 500, random_state = 0)

In [70]: tree_clf.fit(X_train_set_one, y_train)

Out[70]: DecisionTreeClassifier(min_samples_split=500, random_state=0)

In [71]: feature_importances_ = tree_clf.feature_importances_

In [72]: non_zero_feature_indices = []
        for i in range(len(feature_importances_)):
            if feature_importances_[i]>0:
                non_zero_feature_indices.append(i)

In [73]: print("Number of non-zero features: ",len(non_zero_feature_indices))

Number of non-zero features: 2237

In [74]: X_train_non_zero_feature = X_train_set_one[:,non_zero_feature_indices]
        X_test_non_zero_feature = X_test_set_one[:,non_zero_feature_indices]
```

2.2 Apply Decision Tree classifier on Dataset containing only important features

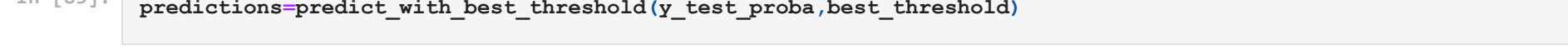
```
In [75]: decision_tree_parameters = {'max_depth':[1,3,10,30], 'min_samples_split':[5,10,100,500]}
        decision_tree = DecisionTreeClassifier(random_state = 0)
        clf = GridSearchCV(decision_tree, decision_tree_parameters, scoring = 'roc_auc', cv = 5, return_train_score =
        GridSearchCV(cv=5, estimator=DecisionTreeClassifier(random_state=0), n_jobs=-1,
                    param_grid={'max_depth': [1, 3, 10, 30],
                                'min_samples_split': [5, 10, 100, 500]},
                    return_train_score=True, scoring='roc_auc')

In [76]: train_auc = clf.cv_results_['mean_train_score']
        cv_auc = clf.cv_results_['mean_test_score']

In [77]: print("Best score = ",clf.best_score_)
        print("Best Hyper parameters = ",clf.best_params_)

Best score = 0.650548549667108
Best Hyper parameters = {'max_depth': 10, 'min_samples_split': 500}
```

2.3 Hyperparameter vs AUC Plot



2.4 Plot ROC-AUC curve of model with best hyperparameters

2.5 Confusion Matrix for Train and Test data

```
In [81]: print("t==100")
best_threshold = find_best_threshold(y_train_proba, tr_thresholds, train_fpr, train_tpr)
conf_matr_train_non_zero_features = confusion_matrix(y_train, predict_with_best_threshold(y_train_proba, best_
print("Train confusion matrix")
print(conf_matr_train_non_zero_features)

=====
the maximum value of fpr*(1-tpr) = 0.4027251144936237 for threshold 0.855
Train confusion matrix
[[ 3216 2243]
 [11048 19549]]
(21899 40214)
```


2.6 Getting all false positive data points

```
In [85]: predictions=predict_with_best_threshold(y_test_proba,best_threshold)
```