

# 1 Decision Tree

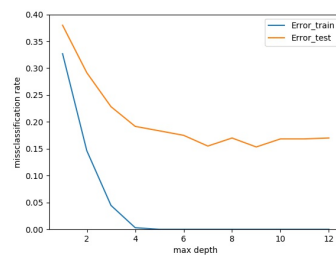
In this assignment, we use xgboost and scikit-learn packages in python to train the datasets separately. Since we don't have labels in test data sets, we have to use valid datasets as test datasets. The scikit-learn DecisionTreeClassifier package has two main splitting criterion: Gini and Entropy. Here, we use entropy as splitting criterion following the lecture notes.

## 1.1 Madelon Dataset

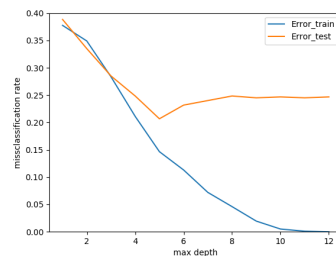
Table 1: Data information

Dataset	Features	Train_obs	Test_obs
Madelon	500	2000	600

In this dataset, this dataset has 500 features. We have 2000 observations as training data. Figure 1 shows after depth=5, the model from xgboost over fits the training data. And the error of test data converges to 0.15. The optimal choice of depth is 3 or 4.



(a) Results from xgboost



(b) Results from Scikit-learn

Figure 1: Results for dataset Madelon

Table 2: Test error with depth

Testerror	Tree depth
0.38	1
0.292	2
0.228	3
0.192	4
0.183	5
0.175	6
0.155	7
0.17	8
0.153	9
0.1683	10
0.168	11
0.17	12

## 1.2 Wilt Dataset

Table 3: Data information

Dataset	Features	Train_obs	Test_obs
Wilt	5	4339	500

In this dataset, there are only 5 features but 4339 observations. After depth=2, the model from xgboost over fits the training data. And the error of test data converges to 0.18. The optimal choice of depth is 7.

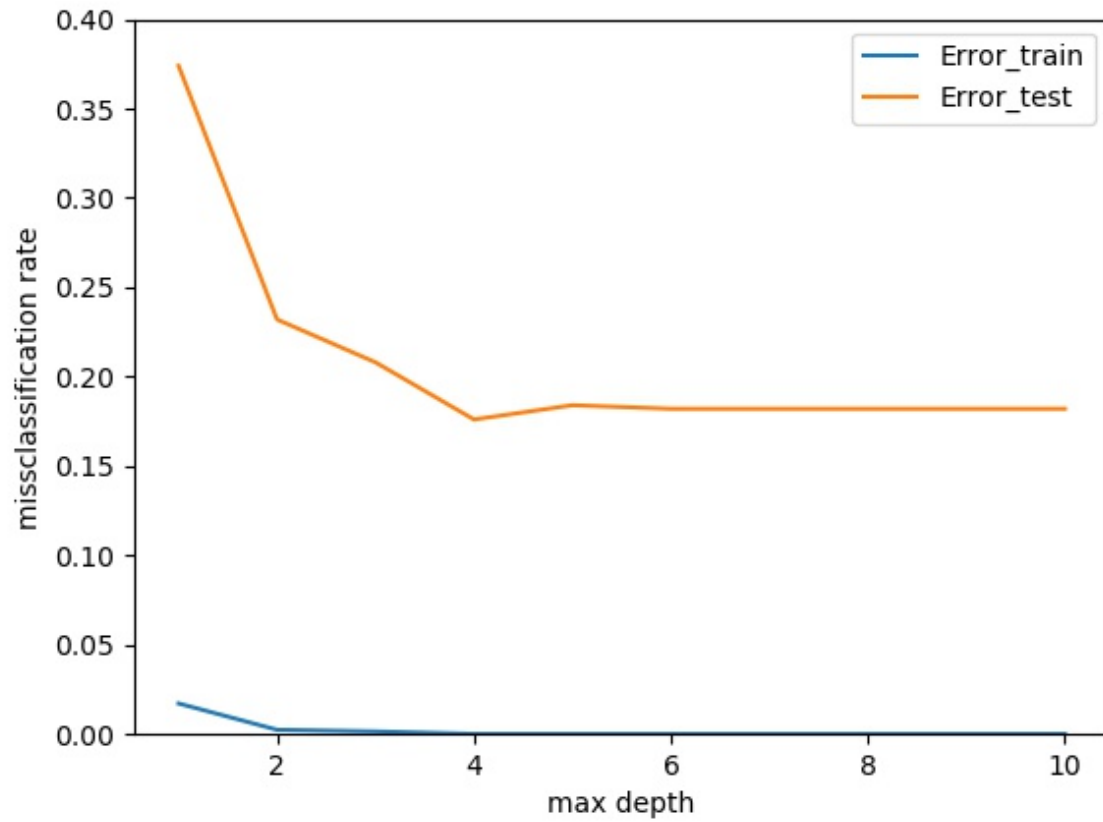


Figure 2

Table 4: Test error with depth

Testerror	Tree depth
0.374	1
0.232	2
0.208	3
0.176	4
0.184	5
0.182	6
0.182	7
0.182	8
0.182	9
0.182	10

### 1.3 Gisette Dataset

Table 5: Data information

Dataset	Features	Train_obs	Test_obs
Gisette	5000	6000	1000

First of all, in this dataset, we have too many features and relative small observations. With decision tree algorithm in xgboost, after depth=4, we can see the model over fits the training data. Therefore, the optimal choice of depth is 4.

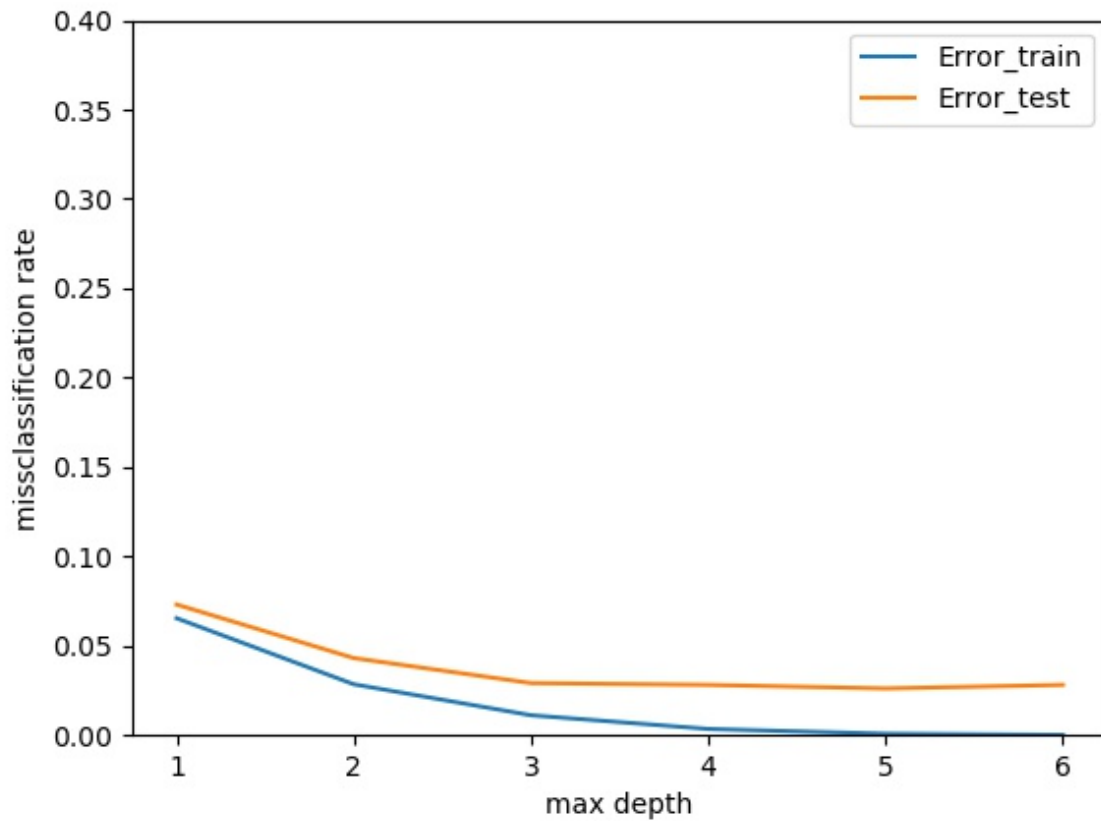


Figure 3

Table 6: Test error with depth

Testerror	Tree depth
0.073	1
0.043	2
0.029	3
0.028	4
0.026	5
0.028	6

## 2 Code

```

c def clearall(): all = [var for var in globals() if (var[:2], var[-2:]) != (" ", " ")]for var in all: del globals()[var]
clearall()
import xgboost as xgb import numpy as np from xgboost import XGBClassifier from sklearn import tree from
sklearn.model_selection import train_test_split from sklearn.metrics import accuracy_score import pandas as pd import matplotlib
def decision_tree(X, Y, XTest, YTest, maxDepth) : rateTrain = np.zeros(maxDepth.shape[0]) rateTest =
np.zeros(maxDepth.shape[0]) for i in range(maxDepth.shape[0]) : fit_model = tree.DecisionTreeClassifier(maxDepth = maxDepth[i])
entropy', maxDepth = maxDepth[i] clf = XGBClassifier(maxDepth = maxDepth[i]) clf.fit(X, Y) print(clf) y =
clf.predict(X) yTest = clf.predict(XTest) py = [round(value) for value in y] pyTest = [round(value) for value in yTest] rateTrain[i] =
1 - np.mean(py == Y) rateTest[i] = 1 - np.mean(pyTest == YTest) accuracy = accuracy_score(pyTest, YTest) print("Accuracy: ", accuracy)
print('train, test misclassification error : ', rateTrain[i], rateTest[i]) return rateTrain, rateTest
def problem1a() : load_data = np.loadtxt('E : material2018_Fall2018_train.txt') dmdl = np.genfromtxt('E :
material2018_Fall2018_train_labels.txt') dmdv = np.loadtxt('E : material2018_Fall2018_valid.txt') dmdlv =
np.loadtxt('E : material2018_Fall2018_valid_labels.txt')
maxDepth = np.arange(12) + 1 train, test = decision_tree(dmd, dmdl, dmdv, dmdlv, maxDepth)
figureIndex = 0 plt.figure(figureIndex) figureIndex += 1 plt.plot(maxDepth, train, label='Error_train') plt.plot(maxDepth, test,
Error_test')
plt.xlabel('max depth') plt.ylabel('misclassification rate') plt.ylim([0, 0.4]) plt.legend() plt.show()
def problem1b() : X = np.genfromtxt('E : material2018_Fall2018_train.csv', delimiter=',') Y = np.loadtxt('E :
material2018_Fall2018_train_labels.txt') XTest = np.genfromtxt('E : material2018_Fall2018_test.csv', delimiter=',')
YTest = np.loadtxt('E : material2018_Fall2018_test_labels.txt') maxDepth = np.arange(10) + 1 train, test =
decision_tree(X, Y, XTest, YTest, maxDepth)
figureIndex = 0 plt.figure(figureIndex) figureIndex += 1 plt.plot(maxDepth, train, label='Error_train') plt.plot(maxDepth, test,
Error_test')
plt.xlabel('max depth') plt.ylabel('misclassification rate') plt.ylim([0, 0.4]) plt.legend() plt.show()
def problem1c() : X = np.loadtxt('E : material2018_Fall2018_train.txt') Y = np.loadtxt('E : material2018_Fall2018_train_labels.txt')
YTest = np.loadtxt('E : material2018_Fall2018_valid_labels.txt') maxDepth = np.arange(6) + 1 train, test = decision_tree(X, Y, XTest, YTest, maxDepth)
figureIndex = 0 plt.figure(figureIndex) figureIndex += 1 plt.plot(maxDepth, train, label='Error_train') plt.plot(maxDepth, test,
Error_test')
plt.xlabel('max depth') plt.ylabel('misclassification rate') plt.ylim([0, 0.4]) plt.legend() plt.show()
problem1a() problem1b() problem1c()

```