

Homework 1 - Decision Trees

Kyle Shaw

Thomas Johansen

This file includes analysis using decision trees for five machine learning datasets. The tree depths producing the most accurate classification for each dataset are summarized in a table near the bottom.

Link to our working directory for this course: https://github.com/johansent/ML_Fall2017
(https://github.com/johansent/ML_Fall2017)

```
In [1]: ### Import a machine Learning dataset
import pandas

def import_data(path, train1, test1, train2, test2, head = 0):
    filepath = 'C:/Users/joh10/Desktop/FSU/FA17/5635/git/Data/' + path + '/'
    #filepath = os.pardir + '\\Data\\' + path + '\\'

    train_data = pandas.read_table(filepath + train1, sep = ' ', header =
head)
    train2_data = pandas.read_table(filepath + train2, sep = ' ', header = head)

    test_data = pandas.read_table(filepath + test1, sep = ' ', header = head)
    test2_data = pandas.read_table(filepath + test2, sep = ' ', header = head)

    return train_data, train2_data, test_data, test2_data
```

```

In [2]: ### train decision trees of depth 1 through 12, computing the classification error at each step for both training and testing sets.
#from import_data import import_data
from sklearn import tree
import matplotlib.pyplot as plt
import sklearn
import numpy as np

def grow_trees(title, legend_loc = 1):
    %matplotlib inline

    errorsTrain = []
    errorsTest = []
    D = range(1,12)

    for d in D:
        clf = tree.DecisionTreeClassifier(max_depth = d)
        clf = clf.fit(X, Y)

        error = 1 - clf.score(X, Y)
        errorsTrain.append(error)

        error = 1-clf.score(Xtest, Ytest)
        errorsTest.append(error)

    plt.plot(D, errorsTrain, label = 'Training Error')
    plt.plot(D, errorsTest, label = 'Test Error')
    plt.xlabel('Tree Depth')
    plt.ylabel('Classification Error')
    plt.title(title)
    plt.legend(loc = legend_loc)
    plt.show()

    #print('min is', min(errorsTest), 'at d =', np.argmin(errorsTest) + 1)

    return errorsTest

```

Madelon data

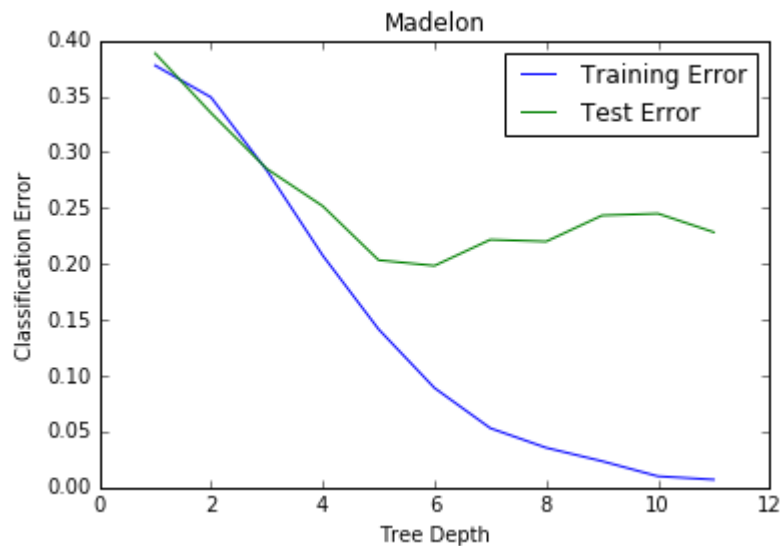
```

In [3]: min_table = {}

X, Y, Xtest, Ytest = import_data('madelon', 'madelon_train.data', 'madelon_val
id.data', 'madelon_train.labels', 'madelon_valid.labels', head = None)
X.drop(X.columns[len(X.columns)-1], axis=1, inplace=True)
Xtest.drop(Xtest.columns[len(Xtest.columns)-1], axis=1, inplace=True) #see ref
[5]

```

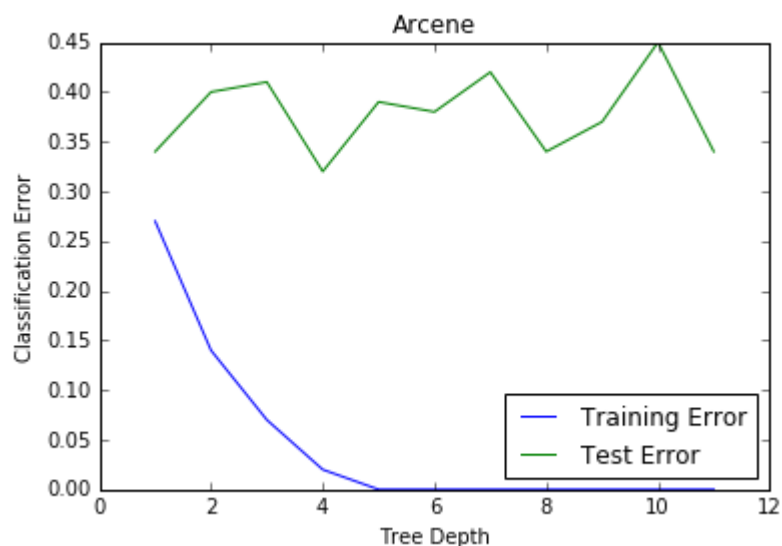
```
In [4]: errorsTest = grow_trees('Madelon')
min_table['Madelon'] = [min(errorsTest), np.argmin(errorsTest) + 1]
```



Arcene data

```
In [5]: X, Y, Xtest, Ytest = import_data('arcene', 'arcene_train.data', 'arcene_valid.
data', 'arcene_train.labels', 'arcene_valid.labels', head = None)
X.drop(X.columns[len(X.columns)-1], axis=1, inplace=True)
Xtest.drop(Xtest.columns[len(Xtest.columns)-1], axis=1, inplace=True)
```

```
In [6]: errorsTest = grow_trees('Arcene', legend_loc = 4)
min_table['Arcene'] = [min(errorsTest), np.argmin(errorsTest) + 1]
```



Hill/Valley data

```
In [7]: X, Y, Xtest, Ytest = import_data('hill_valley', 'X.dat', 'Xtest.dat', 'Y.dat',
      'Ytest.dat', head = None)
```

```
In [8]: errorsTest = grow_trees('Hill/Valley', legend_loc = 3)
      min_table['Hill/Valley '] = [min(errorsTest), np.argmin(errorsTest) + 1]
```



Poker data

```
In [10]: X, Y, Xtest, Ytest = import_data('poker', 'X.dat', 'Xtest.dat', 'Y.dat', 'Ytest.dat', head = None)
```

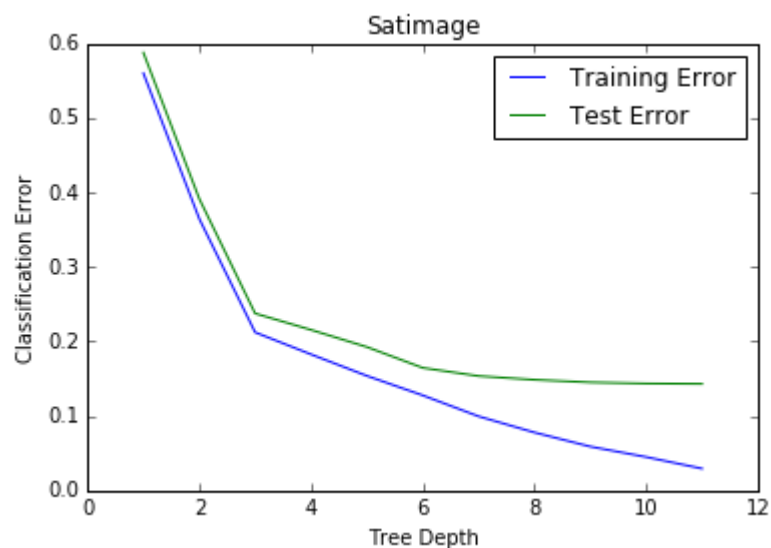
```
In [11]: errorsTest = grow_trees('Poker')
      min_table['Poker      '] = [min(errorsTest), np.argmin(errorsTest) + 1]
```



Satimage

```
In [12]: X, Y, Xtest, Ytest = import_data('satimage', 'X.dat', 'Xtest.dat', 'Y.dat', 'Ytest.dat', head = None)
```

```
In [13]: errorsTest = grow_trees('Satimage')
min_table['Satimage'] = [min(errorsTest), np.argmin(errorsTest) + 1]
```



Summary of Results

```
In [14]: print('{0:12s} {1:15s} {2:12s}'.format('Data name', 'Minimum Value', 'Tree Depth'))
print('-----')
for k in min_table.keys():
    print('{0:12s} {1:13f} {2:12d}'.format(k, min_table[k][0], min_table[k][1]))
```

Data name	Minimum Value	Tree Depth
-----	-----	-----
Madelon	0.198333	6
Satimage	0.143000	11
Hill/Valley	0.463696	10
Poker	0.472820	7
Arcene	0.320000	4

Bibliography

1. Stéfan van der Walt, S. Chris Colbert and Gaël Varoquaux. The NumPy Array: A Structure for Efficient Numerical Computation, Computing in Science & Engineering, 13, 22-30 (2011), DOI:10.1109/MCSE.2011.37
2. John D. Hunter. Matplotlib: A 2D Graphics Environment, Computing in Science & Engineering, 9, 90-95 (2007), DOI:10.1109/MCSE.2007.55
3. Wes McKinney. Data Structures for Statistical Computing in Python, Proceedings of the 9th Python in Science Conference, 51-56 (2010)
4. Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, Édouard Duchesnay. Scikit-learn: Machine Learning in Python, Journal of Machine Learning Research, 12, 2825-2830 (2011)
5. conner.xyz at <https://stackoverflow.com/questions/20517650/how-to-delete-the-last-column-of-data-of-a-pandas-dataframe> (<https://stackoverflow.com/questions/20517650/how-to-delete-the-last-column-of-data-of-a-pandas-dataframe>)