MACHINE LEARNING Decision Trees

PROF COWAN
CHARLES

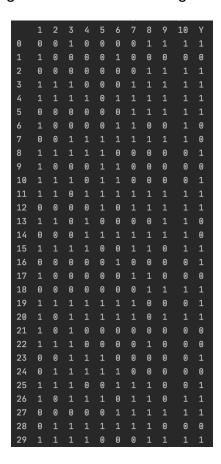
RUTGERS UNIVERSITY

Fall 2020

Part - 1

Generate Dataset

For a given value of k, m, (number of features, number of data points), generated a training data set based on the given scheme in part 1.



Generating ID3 Tree

Entropy, conditional entropies calculation functions are written, and decision tree was built using the max information gain of the Xi values.

```
{1: {0: {6: {0: 1, 1: {3: {0: 1, 1: 0}}}},

1: {4: {0: 0, 1: {2: {0: {7: {0: 0, 1: 1}}}, 1: 1}}}}}

Training error 0.0
```

Training error for the decision tree is 0 which means that the tree fits the dataset.

Finding Training Error and Typical Error for the generated tree

Data for k = 4 and m = 200 was generated and a tree was built.

```
Run: Decision_Trees ×

/Users/rithvikananth/PycharmProjects/Fire_Maze/venv/bin/python /Users/rithvikananth/Desktop/Decision_Trees.py

{8: {0: {7: {0: {1: {0: {2: {0: 1, 1: 0}}}, 1: {3: {0: 0, 1: 1}}}}, 1: 1}},

1: {2: {0: {10: {4: {0: 1, 1: 0}}, 1: {0}}, 1: {1: {0: 0, 1: 1}}}}}

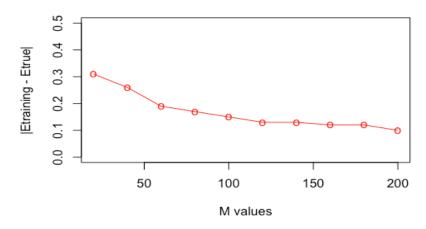
Process finished with exit code 0
```

Answering question 3: Yes, the target variable Y depends on X1 as data is generated as such (dependency of X2, X3. X4, ... Xk on X1 by weights). The weights are also calculated such that X2 has the next higher preference than X3, X4 ... and X3 has higher preference than X4, X5 ... so on. **This dependency of Y on the features is not reflected on the decision tree though.**

Error for 1000 sample data was calculated and the typical error was around 4%.

For K = 10, increasing the m values we observe that the absolute difference between the training error and true error decreases.

Difference as a function of M

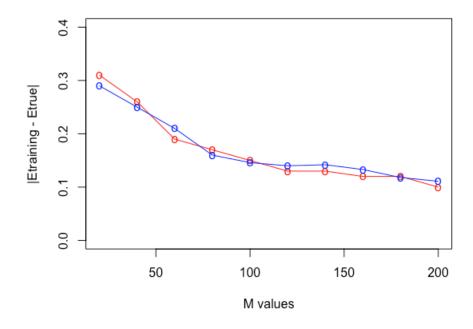


ID3 – Splitting by Gini Tree

An alternative way of splitting the dataset is by calculating the Gini value for all the columns and split on the column with maximum Gini value.

A comparison between information gain and Gini-index split on | Etraining – Etree | as a function of M:

Difference as a function of M



Blue: Gini-Index, Red: Information Gain

Part - 2

Generating Dataset

For a given value of m (number of data points), generated a training data set based on the given scheme in part 2.

Data was generated according to professors' email later.

(Same idea just shrunk down a little bit)

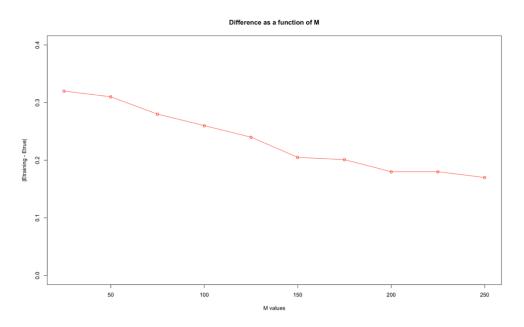
x0 as a fair coinflip, x1 through x10 each dependent on the previous as in the problem, and x11 through x15 as independent random variables - Y is either the majority of x1 through x5, or the majority of x6 through x10, depending on the value of x0.

(**Note:** generating data as per the given question is also written and commented in the code)

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	Υ
0	1	1	1	1	1	1	1	1	0	0	0	0	0	1	0	0	0
1	1	0	0	1	1	1	0	0	0	0	0	0	0	0	1	0	0
2	0	0	1	0	0	0	0	0	1	1	1	0	0	1	1	0	0
3	1	1	1	1	1	1	1	0	1	1	1	1	0	0	1	0	1
4	1	0	0	0	0	0	0	1	1	0	1	1	1	0	1	0	1
5	0	0	1	1	0	1	0	0	0	1	0	1	0	1	1	1	1
6	0	1	0	1	1	0	0	0	0	0	0	1	0	1	1	0	1
7	1	1	1	1	1	0	0	0	0	1	0	0	1	1	0	0	0
8	0	0	1	1	1	0	0	0	0	1	1	1	0	0	0	0	1
9	0	0	1	1	1	1	1	1	0	0	0	0	1	0	0	1	1
10	0	0	1	1	1	1	1	1	1	1	0	0	0	0	1	1	1
11	0	0	1	1	1	1	0	0	0	1	1	0	1	0	1	0	1
12	1	0	0	1	1	0	0	0	0	0	0	0	1	1	0	0	0
13	1	1	0	0	0	1	1	1	1	1	1	1	0	0	0	1	1
14	0	0	0	0	0	0	0	0	1	1	1	0	1	1	0	0	0
15	1	1	1	1	1	1	0	1	1	1	1	1	0	1	0	1	1
16	1	1	1	1	1	1	1	1	0	0	0	1	0	1	0	1	0
17	0	1	1	1	0	1	1	0	0	1	1	0	1	0	0	0	1
18	0	1	1	0	0	0	1	1	1	1	1	1	1	0	1	1	1
19	0	0	0	1	1	0	0	1	1	1	0	0	1	0	1	1	1
20	1	1	1	1	0	0	0	0	0	1	1	0	1	0	0	1	0
21	0	1	1	1	0	0	0	0	0	0	1	0	1	1	0	1	1
22	0	0	0	0	0	0	0	0	1	1	1	1	0	1	1	1	0
23	1	1	0	0	0	0	0	0	0	0	0	1	1	1	0	1	0
24	0	1	1	1	1	1	1	1	1	1	1	1	0	0	1	0	1
25	0	1	0	0	1	1	0	0	0	0	0	0	1	1	0	0	1
26	1	1	1	1	1	1	1	1	1	1	1	0	0	1	1	0	1
27	1	1	1	1	1	0	1	0	0	0	0	0	1	0	1	0	0
28	1	1	1	1	1	0	0	0	1	1	1	0	1	1	1	1	1
29	1	1	1	1	0	1	0	0	1	0	0	0	0	1	0	0	0

Plotting |Etraining - Etrue| as a function of M

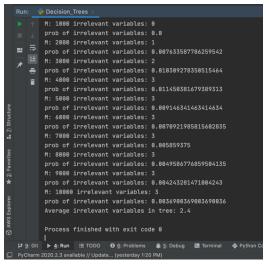
Pruning data was generated and typical error on M data points were plotted. As the m value increases the |Etraining - Etrue| values decrease.

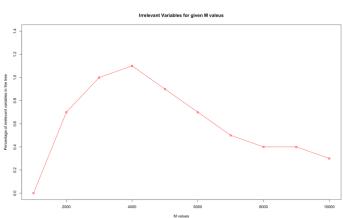


Yes, the above plot agrees with my intuition as there is a marginal decrease in error values with the increase in m values

Irrelevant Variables - M values

For M values in range 1000, 10000 we see that the probability of finding an irrelevant variable in the tree keeps decreasing.



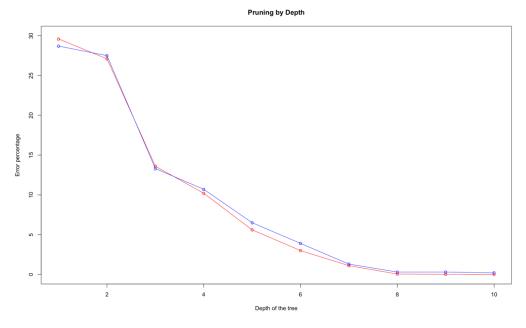


Pruning by Depth - Training and Testing errors

Generating dataset with m = 10000, first 8000 values as training data and remaining 2000 values as testing data.

```
pruning_depth_error()
/Users/rithvikananth/PycharmProjects/Fire_Maze/venv/bin/python /
Depth 1 Error of training data 0.296625
Testing error: 0.287
Depth 2 Error of training data 0.270875
Testing error: 0.275
Depth 3 Error of training data 0.136375
Testing error: 0.133
Depth 4 Error of training data 0.10275
Testing error: 0.107
Depth 5 Error of training data 0.056375
Testing error: 0.065
Depth 6 Error of training data 0.03075
Depth 7 Error of training data 0.0115
Testing error: 0.013
Depth 8 Error of training data 0.0005
Testing error: 0.003
Depth 9 Error of training data 0.000375
Testing error: 0.003
Depth 10 Error of training data 0.0
Testing error: 0.002
Process finished with exit code 0
```

The tree will be similar to ID3 but after a certain given depth we take the maximum Y value from the data at that point.



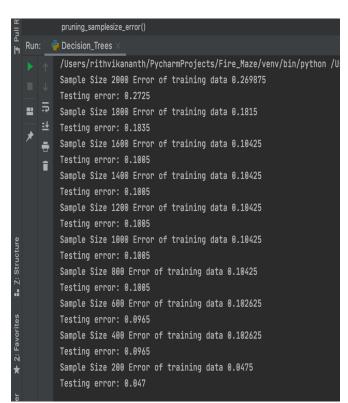
Blue -Testing data, Red – training data.

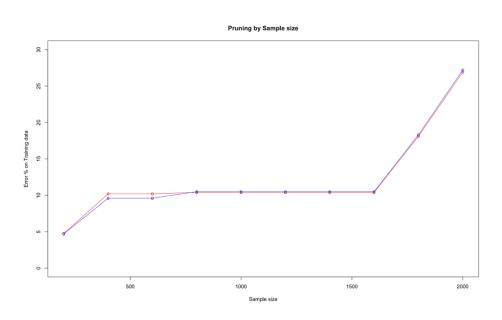
As the depth increases, the error of both training and testing data reduces. At **depth 8** the error on training and testing data is almost 0.

Pruning by Sample Size - Training and Testing errors

The tree will be similar to ID3 but after a sample size shrinks below a given size, we take the maximum Y value from the data at that point.

```
Sample Size 200 Error of training data 0.068625
Testing error: 0.07
Sample Size 190 Error of training data 0.0575
Testing error: 0.0575
Sample Size 180 Error of training data 0.0575
Testing error: 0.0575
Sample Size 170 Error of training data 0.051
Testing error: 0.0545
Sample Size 160 Error of training data 0.051
Testing error: 0.0545
Sample Size 150 Error of training data 0.049875
Testing error: 0.052
Sample Size 140 Error of training data 0.049875
Testing error: 0.052
Sample Size 130 Error of training data 0.049875
Testing error: 0.052
Sample Size 120 Error of training data 0.049875
Testing error: 0.052
Sample Size 110 Error of training data 0.0495
Testing error: 0.0525
Sample Size 100 Error of training data 0.0445
Testing error: 0.0485
Sample Size 90 Error of training data 0.039125
Testing error: 0.046
Sample Size 80 Error of training data 0.0345
Testing error: 0.04
Sample Size 70 Error of training data 0.029625
Testing error: 0.034
Sample Size 60 Error of training data 0.029
```





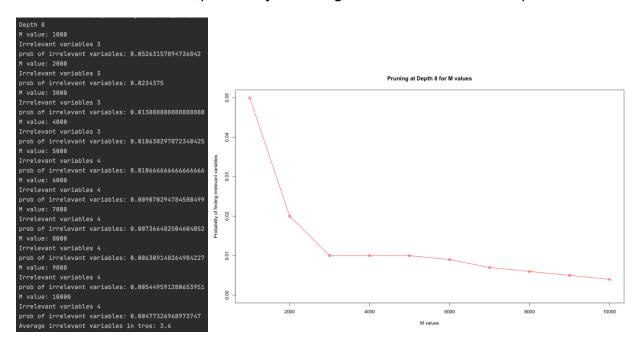
Blue – Testing data, Red – Training data.

The lower s value we round to get Y, the less error in training data. As we can see in the graph, if the rounding off of the data for Y values increases, the error also increases.

We take **sample size as 100 to be ideal** to avoid irrelevant variables.

Irrelevant Features – Pruning by depth

Now for various m values, probability of finding irrelevant variables at depth 8 is found.



At Depth 8 the probability of finding irrelevant variables is less than 0.05 and for M values greater than 7000 the probability of finding irrelevant variables in the tree is almost 0.

Irrelevant Features - Pruning by Sample Size

As mentioned in 2.3b we take sample size as 100.

```
100
M value: 1000
prob of irrelevant variables: 0.0
M value: 2000
prob of irrelevant variables: 0.0
                                                                 Pruning at sample size 0 for M values
M value: 3000
prob of irrelevant variables: 0.0
M value: 4000
prob of irrelevant variables: 0.0
M value: 5000
prob of irrelevant variables: 0.0
M value: 6000
prob of irrelevant variables: 0.0
M value: 7000
prob of irrelevant variables: 0.0
M value: 8000
prob of irrelevant variables: 0.0
M value: 9000
prob of irrelevant variables: 0.0
M value: 10000
```

For different m values the probability of finding irrelevant variables in the tree is always 0 if sample size is 100.

Source Code:

```
import numpy as np
def generate td(k, m):
       header.append(i)
           data.append(d1)
def entropy(data):
```

```
def conditional entropy(data, c):
   entropy conditional = 0
       entropy conditional += (len_xcon/len(data))*e_con_for_x
   return entropy conditional
```

```
def tree value(instance data, tree):
   for node in tree.keys():
   error = 0
def typical error(k, m, tree):
```

```
features = data.columns[:-1]
def gini index(data, target):
   subtree = decision tree(data[data[qini] == 0], features)
```

```
return tree
   header.append(i)
dataset.append(header)
    data.append(x0)
        data.append(d1)
```

```
def typical error pd(m, tree):
def error estimate pd():
error estimate pd()
def tree variables(data, tree, a = []):
        a.append(node)
def irrilavent variables():
```

```
tree variable = tree variables(data, tree)
           if variables > 10:
def decision tree depth(data, features, depth, d):
       info gains.append(entropy(data) - conditional entropy(data, col))
```

```
subtree = decision tree depth(data[data[max gain] == 0], features,
depth, d)
       subtree = decision tree depth(data[data[max gain] == 1], features,
def pruning depth error():
   data = pruning data(10000)
       depth += 1
def decision tree sample size(data, features, sample size):
```

```
info gains.append(entropy(data) - conditional entropy(data, col))
def pruning_samplesize_error():
def irrilavent variables depth(depth):
```

```
if variables > 10:
def irrilavent variables samplesize(sample size):
               b.append(variables)
```

R-code for plotting:

```
plot(mvalue, difference IG, type="o", col="red", pch="o", lty=1, ylim=c(0,0.4),
   ylab="|Etraining - Etrue|", xlab = "M values") +
 title("Difference as a function of M") +
 points(mvalue, difference_GI, type="o", col="blue", pch="o", lty=1) +
 legend(1, 95, legend=c("Information gain", "Gini Index"),
     col=c("red", "blue"), Ity=1:2, cex=0.8)
plot(m, training_error, type="o", col="red", pch="o", lty=1, ylim=c(0,30),
   ylab="Error % on Training data", xlab = "Sample size") +
 title("Pruning by Sample size") +
 points(m, testing error, type="o", col="blue", pch="o", lty=1)
plot(m, prob, type="o", col="red", pch="o", lty=1, ylim=c(0,1),
   ylab="Probability of finding irrelevant variables", xlab = "M values") +
 title("Pruning at sample size 0 for M values")
plot(m, error, type="o", col="red", pch="o", lty=1, ylim=c(0,30),
   ylab="Error % on Training data", xlab = "Sample Size (s)") +
 title("Pruning by sample size")
plot(depth, error, type="o", col="red", pch="o", lty=1, ylim=c(0,30),
   ylab="Error percentage", xlab = "Depth of the tree") +
 title("Pruning by Depth") + points(depth, error2, type="o", col="blue", pch="o", lty=1)
```

```
plot(mvalue, difference_IG, type="o", col="red", pch="o", lty=1, ylim=c(0,0.4),
    ylab="|Etraining - Etrue|", xlab = "M values") +
    title("Difference as a function of M")
```