

Homework 2

Due : April 4th (11:59 pm)

Submitted by: Jyotsna Rajaraman

Q1 ALMA3 \rightarrow Ex 18.13, Ex 18.14

18.13) P.T :

$\text{Num}(\text{rules in decision list}) \leq \text{leaves}(\text{decision tree})$

We know that in decision lists, each rule has a binary outcome, YES and NO. If YES \rightarrow there will be a prediction, if no, we will proceed to the next rule.

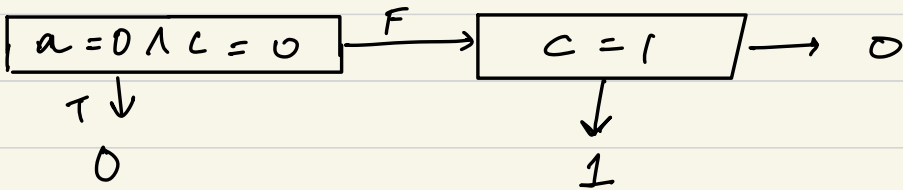
Let us make the assumption that, in worst case, each of the n leaves in the decision tree corresponds to an unique outcome. This means that we would need the same number of rules to create n unique outcomes.

If not, we may actually be able to combine the nodes of a decision tree to create a rule that may reduce the overall number of rules together.

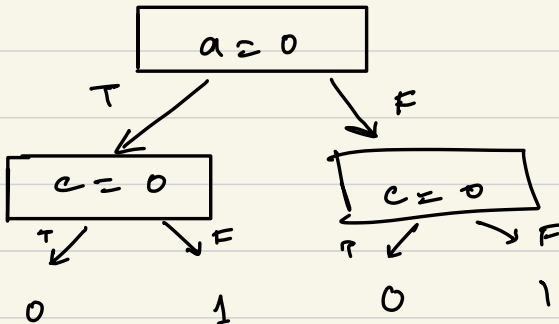
$\therefore \text{Num}(\text{rules})_{\text{in list}} < \text{Num}(\text{leaves})_{\text{in tree}}$

a	b	c	out come.
0	0	0	0
0	0	1	1
1	0	0	0
1	0	1	1

Decision List



Decision Tree

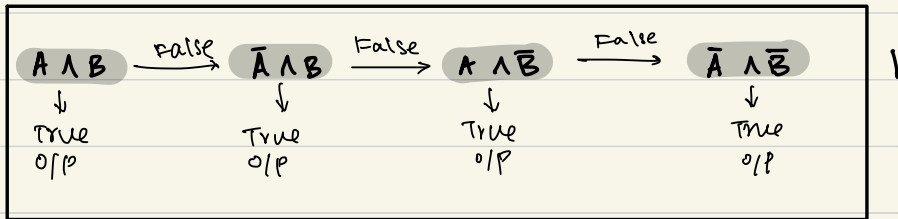


18.4) Expressiveness of decision lists

a)

Each rule in a decision list is a conjunction of Boolean variable and thus we can essentially create an expression for any Boolean function as:

as:



For 2 variables

In this way we can actually write an expression for each outcome and represent it with decision lists.

b)

\therefore The decision tree has depth = k. At most, each path has k nodes. \therefore Taking the longest path, we can convert it into a decision list. Hence, we can use at most k literals to build a decision list for the function.

Q2) Generalization Error:

* $T \rightarrow$ training set ; $V \rightarrow$ validation set
(both i.i.d)

* Max Depth	Err Rate
4	0.45
5	0.35
6	0.40

* Gen Error : $E(\text{misclassification err over all alg.})$

STATEMENT #1:

An unbiased estimate of the generalized error rate for a decision tree w/ max depth 5 = 0.35

Since T and V are i.i.d samples & are random it is true.

We can explain this in detail by considering that V (validation) set of data is independent of the training data. Which means that

the generalized error is
an unbiased estimate

Statement 2:

An unbiased estimate of a
tuned decision tree is 0.35

False, because the tuned decision
tree has a validation set that is
dependent. This means that the tree
is built based on the training &
validation. Therefore cannot be generalized
over the universe U .

\therefore It is NOT an unbiased estimate.

Statement : 3

The unbiased estimate of the random decision tree is 0.35

0.35 is the unbiased estimate specifically for the random case chosen since T and V are not dependent. However, in the case of a random decision tree - it may be the case that the max depth could change and \therefore The generalized error would also vary.

Q3) CROSS VALIDATION:

Task \rightarrow Binary Classification

$U \rightarrow$ Universe

- $\nearrow p = \text{class 1}$
- $\searrow 1-p = \text{class 0}$

$$0.5 < p \leq 1$$

\therefore class 1 is the majority class.

$A \rightarrow$ classification algorithm.

\hookrightarrow always predicts the majority in training set
or if equal \rightarrow picks random.

$\hookrightarrow (R, \text{sample})$

$D \rightarrow$ dataset of m examples. (i.i.d) $\in U$

$C \rightarrow$ accuracy of A on D w k fold cv

$\hookrightarrow (R, V)$

a) $E[C] = ?$

$X_j \rightarrow$ class of example j

$1 \leq j \leq m \Rightarrow D$ has m samples.

training & test for i th fold = $T_i, S_i \rightarrow 1 \leq i \leq k$.
 $\Rightarrow k$ folds.

$C_i \rightarrow$ accuracy of each fold.

For fold i :

$$E[X|Y=y] = \sum_x x f_{X|Y}(x|y)$$

$$E[C_i | T_{i(maj)} = 1] \rightarrow C = \text{accuracy} = \frac{\text{correct}}{\text{total}}.$$

$$\text{If } T_{i(maj)} = 1 \rightarrow C_i | T_{i(maj)} = 1 = \begin{cases} \frac{p}{m} & \text{w prob } 1 \\ \frac{1-p}{m} & \text{w prob } 0 \end{cases}$$

$$C_i | T_{i(maj)} = 0 = \begin{cases} \frac{1-p}{m} & \text{w prob } 1 \\ \frac{p}{m} & \text{w prob } 0 \end{cases}$$

Total no. of correctly classified instances in each fold = no. of instances of majority class in that fold.

$$\therefore \text{Tot no. of correctly classified instances} \left\} = \left[\left(p \times \frac{m}{k} \right) + (1-p) \frac{m}{k} \right] k$$

$\frac{m}{k} \rightarrow$ no. of instances in each fold.

$$E(C) = \frac{\text{Correct}}{\text{Total.}}$$

$$= \frac{k \left[\left(p \times \frac{m}{k} \right) + (1-p) \times \frac{m}{k} \right]}{m}$$

⑥, ⑦ \Rightarrow Python.

(d) if no. of samples, $m \rightarrow \infty$

$$C = \frac{\text{correct}}{\text{total}}$$

\therefore class 1 is $= p$ = majority

Prediction : class 1

$$\begin{aligned} \therefore \text{correct} &= p \\ \text{total} &= m \end{aligned}$$

Weak Law: Sample av converges towards expected value.

$$\therefore \lim_{m \rightarrow \infty} Pr(C) = p$$

⑤ $m \neq 0$, $p \neq 0$ and k .

such that $C \in D$

Suppose there are:

$$m = 10$$

$$p = 6 \Rightarrow \text{Class 1}$$

$$1-p = 4 \Rightarrow \text{Class 0}$$

Say training set = all class 1

test set = all class 0

\therefore accuracy will be 0

(Q4)

i/p \rightarrow integer ; o/p \Rightarrow Boolean

$$y(n) = (a \leq n_1 \leq b) \text{ and } (c \leq n_2 \leq d)$$

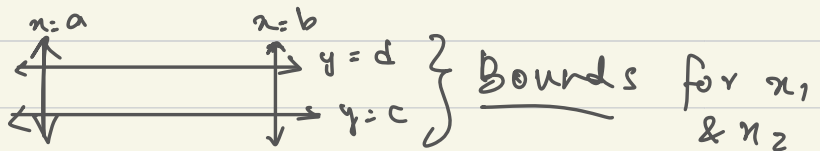
② S.T the no. of hyp is

$$= \binom{n(n+1)}{2}^2$$

No. of hypothesis = no. of distinct rect
restricted to 1 to n int position.

We know $\rightarrow a, b, c, d \in \{1 \dots n\}$.

if $n_1 \in [a, b]$ & $n_2 \in [c, d]$



Number of rectangles:

choices [A] + choices [B] + choices [C] + choices [D]
can't be same can't be same

∴

We pick two lines for a, b & 2 for c, d.
from $\{0, 1, 2, \dots, n\}$.

$$\begin{aligned}\therefore {}^{n+1}C_2 \times {}^{n+1}C_2 &= \frac{n(n+1)}{2} \times \frac{n(n+1)}{2} \\ &= \left[\frac{n(n+1)}{2} \right]^2\end{aligned}$$

Proved.

⑦ How many training samples for 0 training error to have $E_{\text{gen}} \leq 10\%$ with prob 95% \rightarrow prob $\rightarrow 1 - \frac{0.05}{8}$.

$$N \geq \frac{1}{\epsilon} \times \left[\ln |H| + \ln \frac{1}{\delta} \right] \geq \frac{1}{\left(\frac{10}{100}\right)} \times \ln |H| + \ln \frac{1}{0.95}$$

$$N \geq \frac{100}{10} [\ln |H| + \ln 20]$$

$$N \geq 10 \left[\left(\frac{n(n+1)}{2} \right)^2 + 2.99 \right] //$$

Q5. Python code has been submitted

Q6. Discussion:

- * I used 10,000 rows of data

- * For Imputer

 - > for classification | discrete \rightarrow 'most freq'

 - \therefore we don't want to end up with predictions that do not exist.

 - > for continuous values \rightarrow 'median'

 - \therefore to avoid creating a bias within the mean/expectation of data

- * It may be useful to check for which rows have missing values—especially for smaller data sets to eliminate rows with little to no info.

For instance we may have data samples that do not have any cols except for 1 - 2 $\rightarrow \therefore$ introducing error

However for larger datasets it may be a trivial result requiring additional computation & analysis.

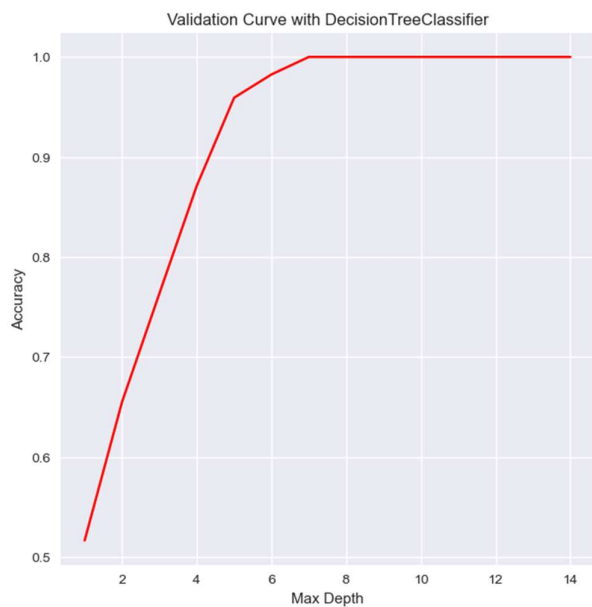
(This has been reflected in the data test)

Learning Curves for DecisionTree Vs LogisticReg.

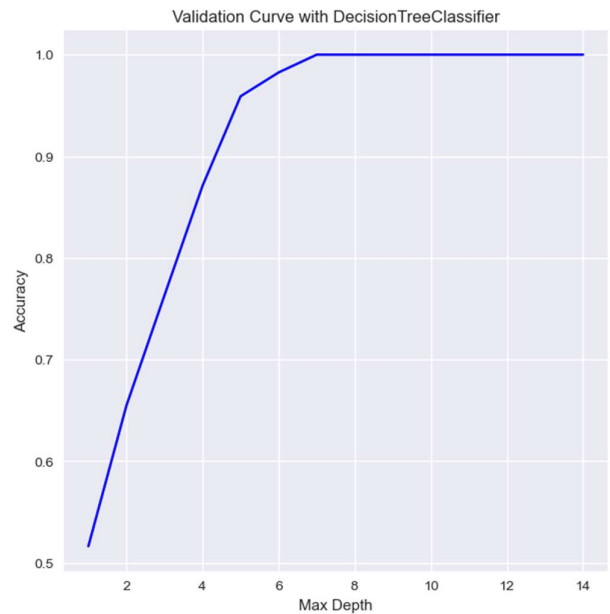
For decision tree, the error goes slightly \uparrow as training set increases but for logistic regression \uparrow in training set size decreases the error.

For Decision Tree Classifier Model:

Validation Curves:



Validation Curve - Training Data



Validation Curve - Testing Data

Learning Curves:



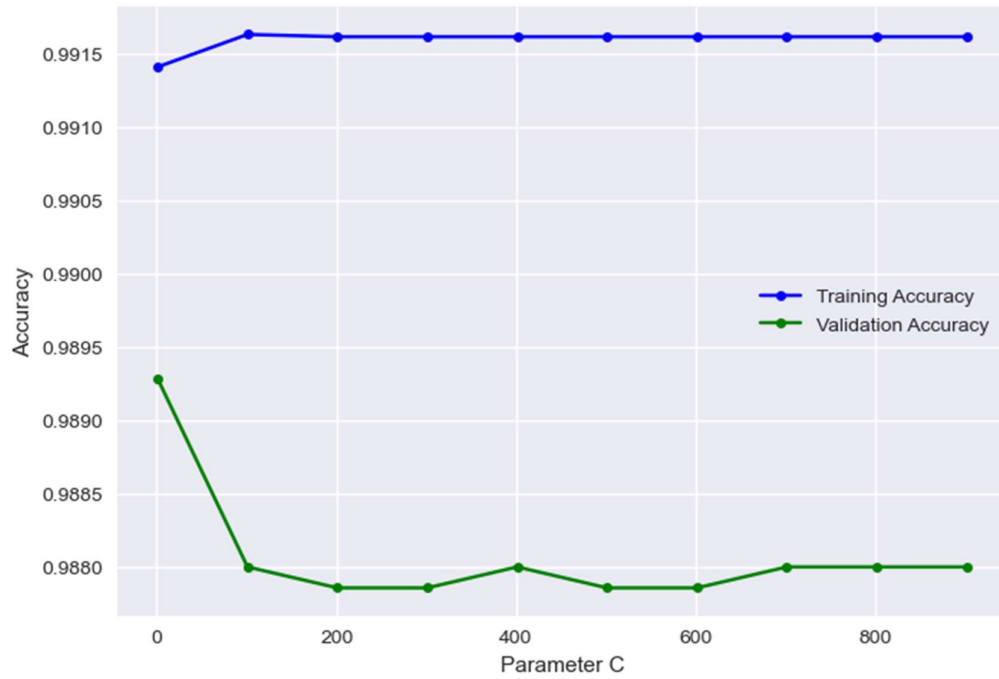
Learning Curve : Depth = 5



Learning Curve : Depth = 7

For Logistic Regression Model:

Validation Curve:



Learning Curve:

