

Decision Tree

WORLD STAR™

Date: _____

Page: _____

Q1- Is the animal big?
NO | YES

← Parent Nodes

child Node → Q... ?

Q2. Does it have a long neck?
NO | YES

leaf Node → Q3. Does it have a trunk
NO | YES

Guess Giraffe

Q4. Does it live in water?
NO | YES

Guess Elephant

Q5. Is it a feline?
NO | YES

Guess Hippo

Guess Rhino

Q... ?

→ Information Gain = (base entropy) - (new entropy)

• Gini Index Impurity = $1 - (\text{the probability of "yes"})^2 - (\text{the prob "no"})^2$

Independent Variable →

Cholesterol

Dependent Variable →

Heart disease

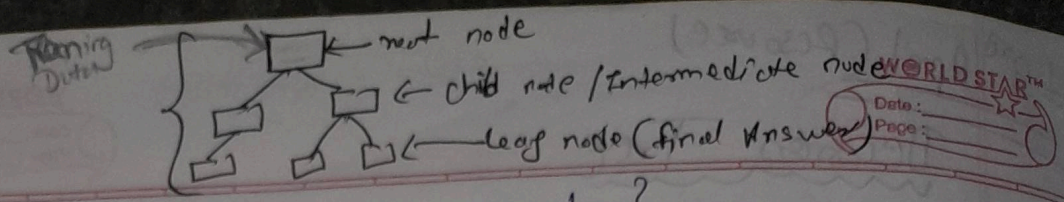
Heart Disease

← dependent variable

Yes	No	Yes	No
105	39	34	125

$$= 1 - \left(\frac{105}{105+39} \right)^2 - \left(\frac{39}{105+39} \right)^2$$

$$= 0.395$$



Q How learning is going over here?

- which root is good
- ↳ means we are classifying the data based on the root based on different-different root i.e. child root, root.
- ↳ best split possibility
- ↳ There are diff-different methods like Gini Indexing, Entropy, Information gain, chi-square method.
- ↳ gini is half of entropy

* Gini Impurity # Gini Impurity is faster due to no log calculation

Formula

→ Gini Impurity = $1 - \sum_{i=1}^n (P_i)^2$

(or) Gini Index

* Entropy (negative) log of prob (ave) log y - prob (ave)

→ $H(A) = P_+ \log(P_+) - P_- \log(P_-)$

* Information Gain

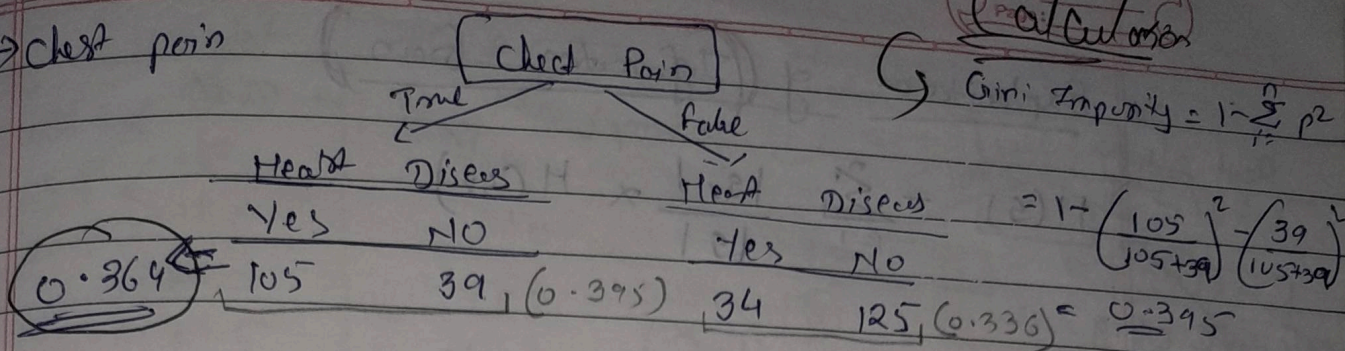
→ $H(S) = \sum \frac{|S_v|}{|S|} \times H(S_v)$

eg: Dataset = chest Pain, Good Blood Circulation, Blocked Arteries, Heart Disease

chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
NO	NO	NO	NO
YES	YES	YES	YES
YES	YES	NO	NO
YES	NO	???	YES
etc...	etc...	etc...	etc.

To determine which separation is best, we need a way to measure and compare "Impurity"

→ chest pain



→

Good Blood Circulation

Less Impure

0.360

True		False	
Heart	Disease	Heart	Disease
Yes	No	Yes	No
37	127	100	33

→

Blocked Arteries

High Impure

0.981

True		False	
Heart	Disease	Heart	Disease
Yes	No	Yes	No
92	31	45	129

Gini Impurity for chest pain
= weighted avg of
Gini impurities for the
leaf node

$$= \left(\frac{144}{144+159} \right) 0.395$$

$$+ \left(\frac{159}{144+159} \right) 0.336$$

$$= 0.364 //$$

Step 1 ⇒ Calculate Gini Impurity

Step 2 ⇒ weighted avg of Gini Impurity

$$(1) \text{ Gini Impurity} = 1 - \sum_{i=1}^n p_i^2$$

(2) Gini impurity for chest pain = weighted avg of Gini Impurity for the leaf nodes

Pre-pruning ⇒ Pre-Pruning = We are not reaching to the leaf node in-between, we are going to stop a construction of decision tree.

(It is hyperparameter)

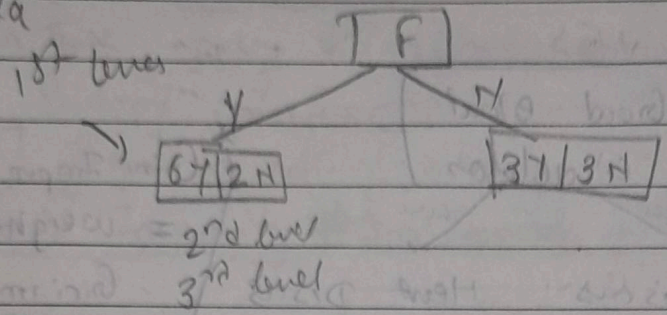
* Explanation of (Information Gain)

Root Entry

$$H(S) = - \sum_{i=1}^n \frac{|S_v|}{|S|} \times H(S_v)$$

Sample 1st level

calculate using
Entropy formula



DTR

Step 1
Sort

$$H(S) = \left(\frac{8}{14} \times H(6) + \frac{6}{14} \times H(5) \right)$$

DTR

Sort (height)	Height	Weight	Sort (weight)	threshold value	height
4.2	5.4	60 kg	80 kg	≤ 4.2	80 kg
4.3	5.6	70 kg	68.4 kg	> 4.2	68.4 + 72.6 + 60 + 70 + 75 + 72
4.4	5.8	75 kg	72.6 kg		
5.4	6.1	72 kg	60 kg		
5.6	4.3	68.4 kg	70 kg		
5.8	4.4	72.6 kg	75 kg		
6.1	4.2	80 kg	72 kg		

mean of height

Then find MSE (or) MAE

(find mean of these values)

67

* How we can find MSE?

→ write each and every data points we will try to calculate mse.

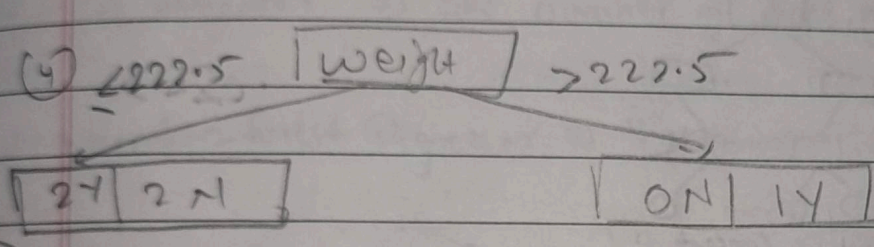
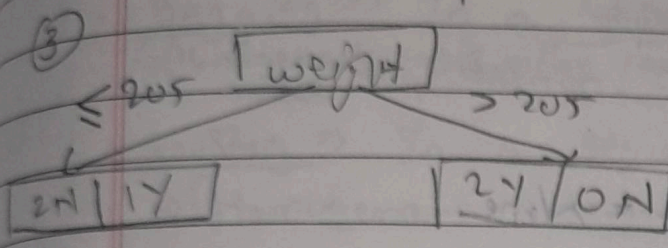
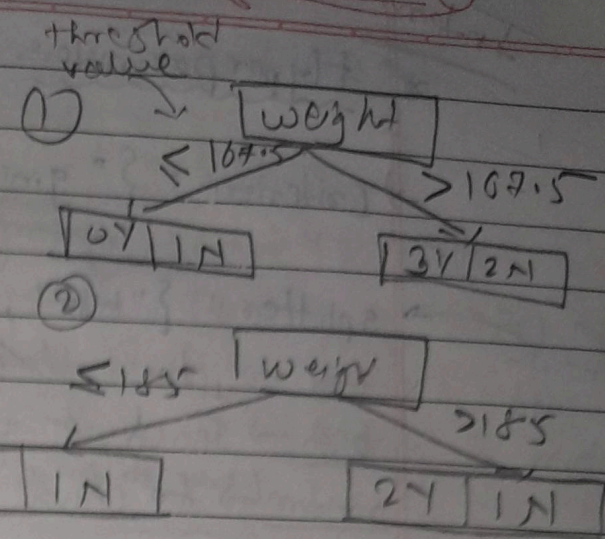
Calculation for MSE: $\sum (y - \hat{y})^2$

$(80-80)^2 + (80-68.4)^2 + (80-72.6)^2 + (80-60)^2 + (80-70)^2 + (80-72)^2$

- Steps ⇒
- Sort height from lowest to highest (mean)
 - Calculate the Average weight for all adjacent

The less mse the best node selection

weight	Heart Disease
155	No
180	YES
190	No
220	YES
225	YES

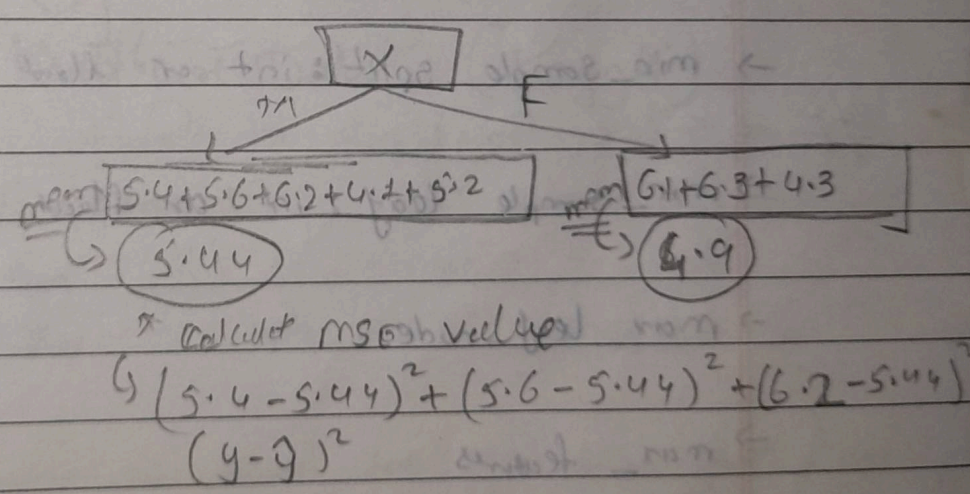


Now we can calculate

- ① Gini Impurity (or)
- ② Entropy

no need to sort

cloth size	Height
X	5.4
M	5.6
L	6.1
M	6.2
L	6.3
L	4.3
M	4.2
M	5.3



Imp { CART ⇒ Gini Index (By default) Impurity we are going to build DT)

IO3 ⇒ Information Gain (Based on Information Gain we are going to build DT)

Interview Ques

* Hyperparameters of DT

→ Criterion: $\{ \text{"gini"}, \text{"entropy"} \}$, default = "gini"
for Information Gain
minimum Gini Impurity we are gaining to consider

→ Splitter: $\{ \text{"best"}, \text{"random"} \}$, default = "best"

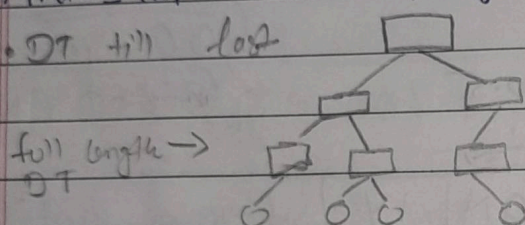
used to choose the split at each node.

best → to choose the best split

random → to choose the best random split

→ max_depth: int, default = None

• DT till last



full length →
DT

• full length $\xrightarrow{\text{two concept comes into picture}}$ $\begin{cases} \text{(1) Pre-pruning} \\ \text{(2) Post-pruning} \end{cases}$

→ min_sample_split: int or float, default = 2

→ min_sample_leaf: int or float, default = 1

→ max_leaf_node

→ max_features

• Splitting:- It is a process of dividing a node into two or more sub-nodes.

• Pruning:- when we remove sub-nodes of a decision node, this process is called pruning. You can say opposite process of Splitting.