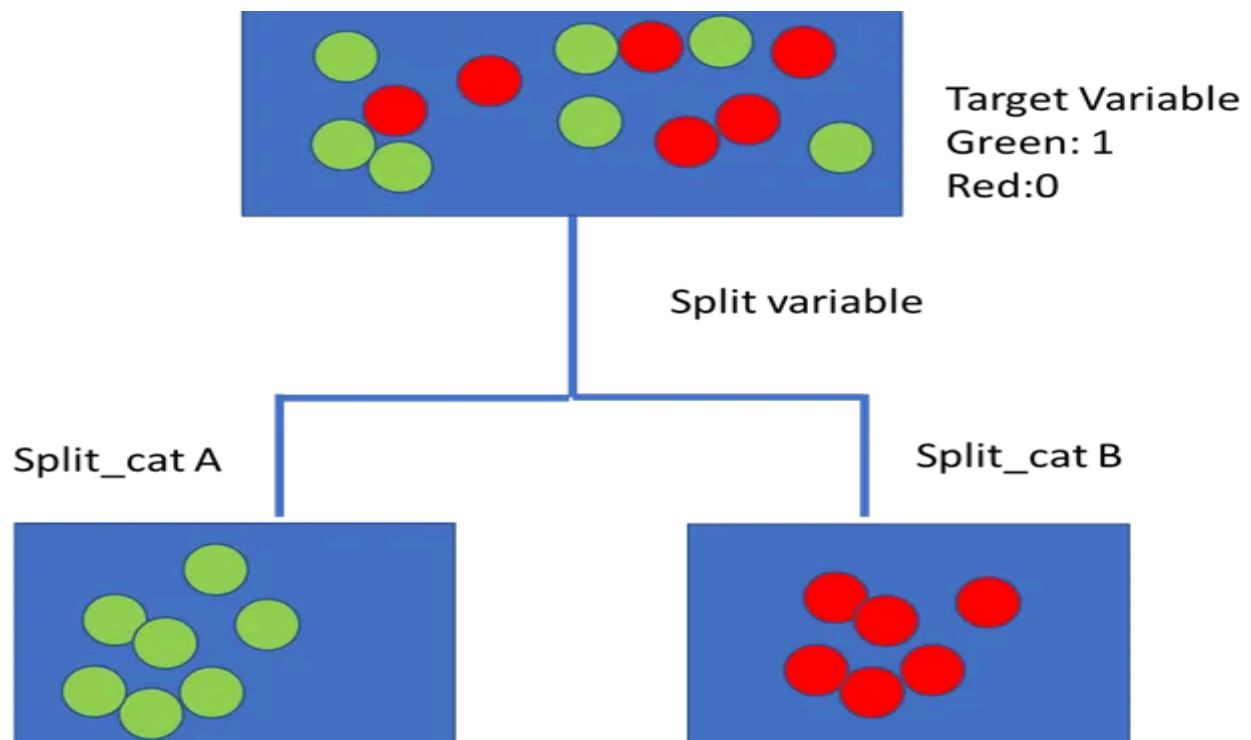


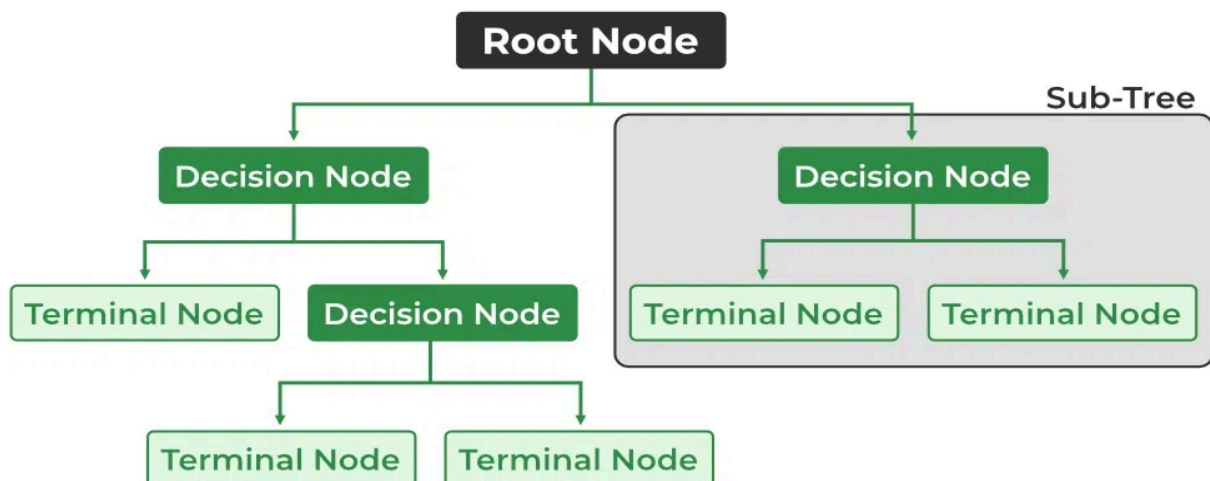
Decision Tree

- A decision tree is a flowchart-like tree structure where each internal node denotes the feature, branches denote the rules and the leaf nodes denote the result of the algorithm.
- Decision tree is a hierarchical data structure that represents data through a divide and conquer strategy.
- A decision tree is a simple model for supervised classification. It is used for classifying a single discrete target feature.

What does a Decision Tree do?



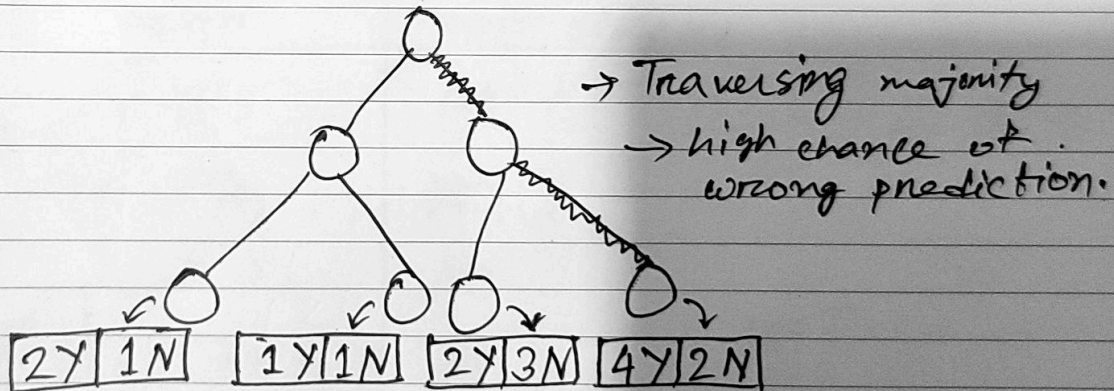
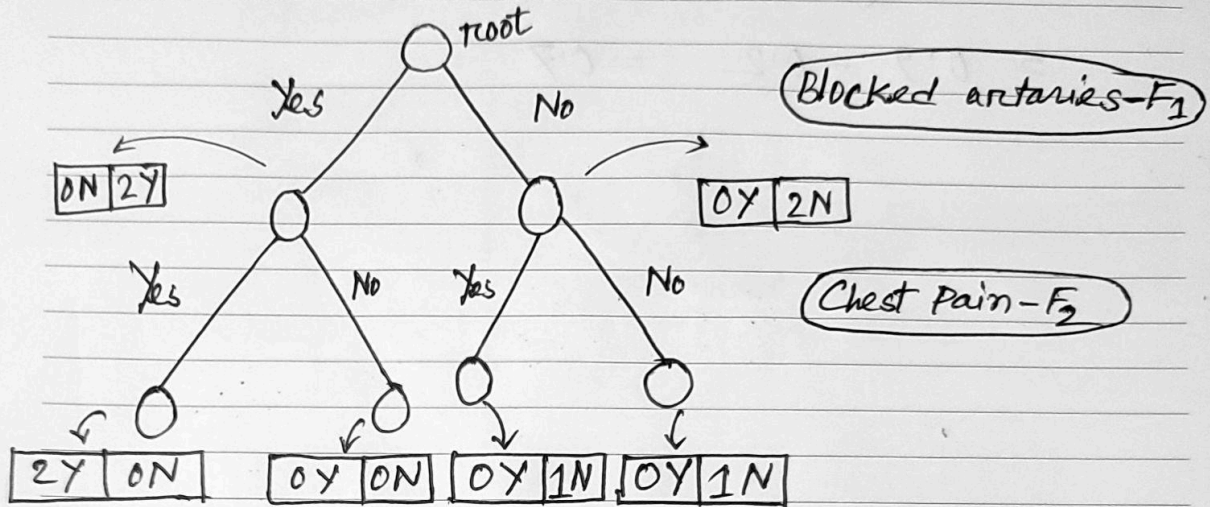
Decision Tree Look like:

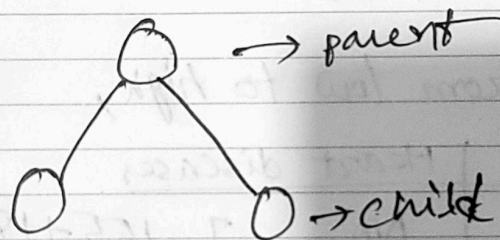
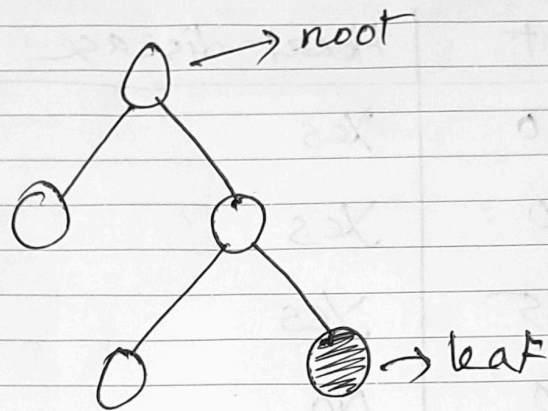


X

→ F_1 means feature 01, etc.

Date: 25 02 2025





→ Decision tree's prediction comes from leaf node ..

Impurity:

for each feature (calculation)

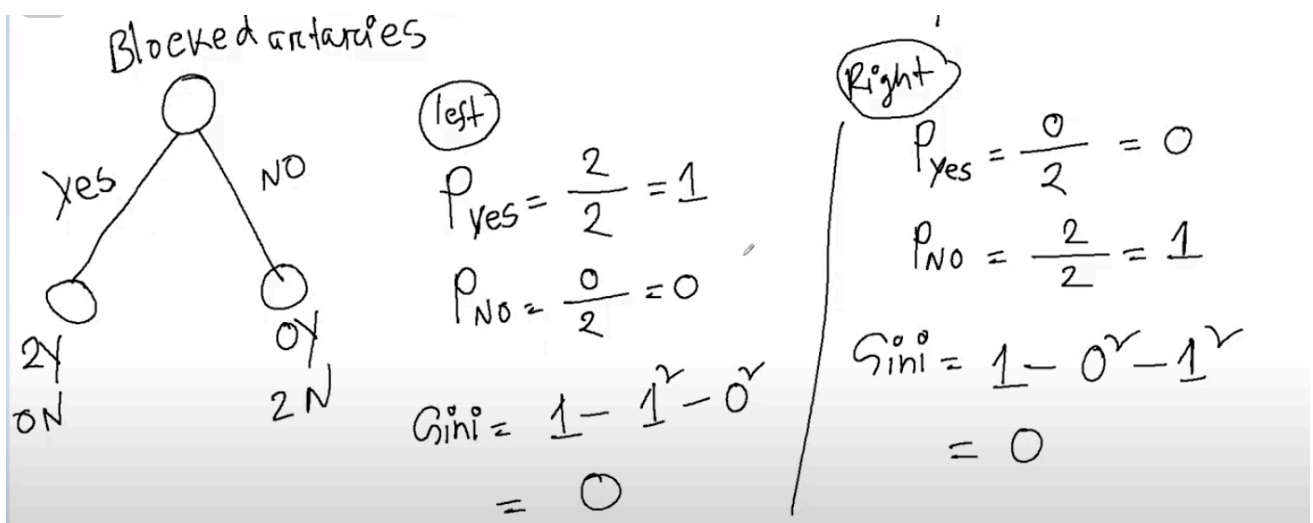
↓
lowest impurity feature (good sign)

↓
make decision tree

✓ Gini = $1 - \sum p(i)^2$ (usually used)

Entropy = $1 - \sum p_i \log_2(p_i)$

chest pain	good blood circulation	Blocked arteries	Heart disease	weight	Heart disease
NO	NO	NO	NO	220	Yes
Yes	yes	yes	yes	180	Yes
Yes	Yes	NO	NO	225	Yes
Yes	NO	Yes	Yes	190	NO
				155	NO



$$Gini = \left(\frac{3}{4}\right) 0.455 + \frac{1}{4} 0$$

$$= \boxed{0.33} \text{ perfect}$$

$$Gini = \frac{2}{4} (0) + \frac{2}{4} (0)$$

$$= 0$$

Decision based on Information Gain:

Decision Trees:

- Think of a decision tree like a flowchart, where each decision leads to more decisions, and finally, to an outcome. In each step, we make a decision based on certain features to get closer to our goal.

Entropy - Measure of Disorder:

- Entropic Playground: Entropy is like a measure of messiness or disorder in our data.
- Low Entropy: Low entropy means our data is more organized or homogeneous.
- High Entropy: High entropy means our data is a bit messy or diverse.

Information Gain - Seeking Order:

- Information Gain is our guide. It helps us decide which features bring more order to our data. We pick the feature that reduces the chaos in our dataset the most.

Classification or Regression:

Once the decision tree is built, it can be used for classification (for categorical target variables) or regression (for numerical target variables). Each path from the root to a leaf node represents a decision rule.

chest pain	good blood circulation	Blocked arteries	Heart disease	weight	Heart disease
NO	NO	NO	NO	220	Yes
Yes	yes	yes	yes	180	Yes
Yes	Yes	NO	NO	225	Yes
Yes	NO	Yes	Yes	190	NO
				155	NO

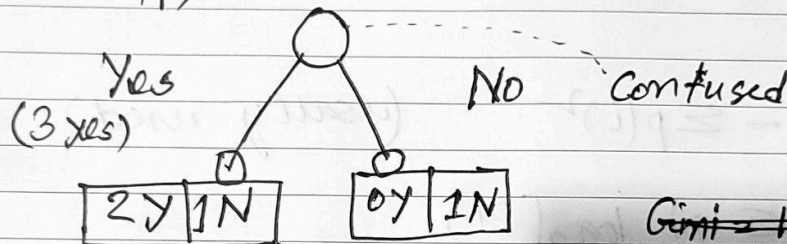
Feature: Chest pain, good blood circulation,
Blocked arteries

Target: Heart disease

Dataset Split:

$F_1 \Rightarrow$ Chest pain

F_1



~~$Gini = 1 - \sum p(i)^2$~~

$\therefore Gini = 1 - \sum p(i)^2$

$\therefore Gini = 1 - \sum p(i)^2$

Left node:

$P_{yes} = \frac{2}{3} = 0.66$

$\therefore Gini_L = 1 - (0.66)^2 - (0.33)^2$
 $= 0.455$

$P_{No} = \frac{1}{3} = 0.33$

Right node:

$P_{yes} = \frac{0}{1} = 0$

$\therefore Gini_R = 1 - (0)^2 - (1)^2$
 $= 0$

$P_{No} = \frac{1}{1} = 1$

$\therefore Gini = Gini_{(Left)} \times \text{sample} + Gini_{(Right)} \times \text{sample}$

$= \frac{3}{4} \times 0.455 + 0 \times \frac{1}{4}$

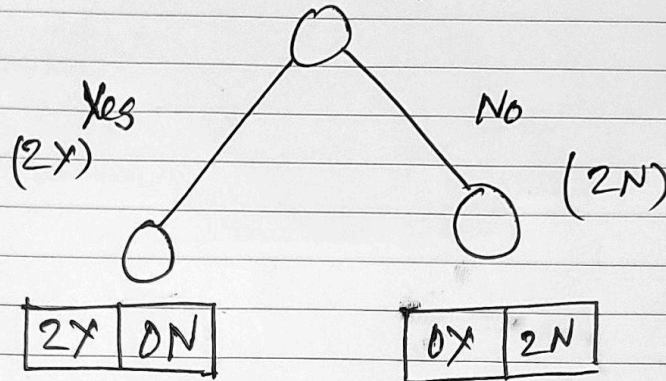
$= 0.34125$

STOLL

KARL MAYER

KMON

(Blocked arteries) $\frac{1}{2}$



Left

$C_{ini} =$

$$P_{yes} = \frac{2}{2} = 1$$

$$P_{no} = \frac{0}{2} = 0$$

$$\therefore C_{ini_L} = 1 - (1)^2 - (0)^2 = 0$$

Right:

$C_{ini} =$

$$P_{yes} = \frac{0}{2} = 0 \quad \therefore C_{ini_R} = 1 - (0)^2 - (1)^2 = 0$$

$$P_{no} = \frac{2}{2} = 1 \quad = 0$$

$$\therefore C_{ini_T} = \frac{1}{4} \times 0 + \frac{1}{4} \times 0 = \boxed{0}$$

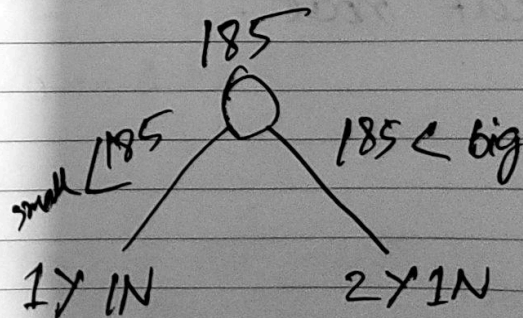
For Numerical value,

Date :

Weight	Heart disease
220	Yes
180	Yes
225	Yes
190	No
155	No

After sorting from low to high,

	Weight	Heart diseases	
167.5	155	No] $\frac{155+180}{2} = 167.5$
	180	Yes	
185	190	No	
205	220	Yes	
222.5	225	Yes	



(2) Average value of Gini is 2.17, after this the scheme draw into 2, 1

Decision Tree Colab:

https://colab.research.google.com/drive/1Kvm_faeSyZOfKop6WuSrQ_BTpChS2xrw?usp=sharing

[3. SVM classifier, Decision Tree Regressor-Classification.ipynb](#)