

Day-4 Machine Learning Algorithms

Agenda

- ① Decision Tree Classification
- ② Decision Tree Regression
- ③ Practical Implementation
- ④ Ensemble Technique

Decision Tree of solving many use cases by

↓
→ Regression (It has worst time complexity)
Classification

If (Age ≤ 18):

print("College")

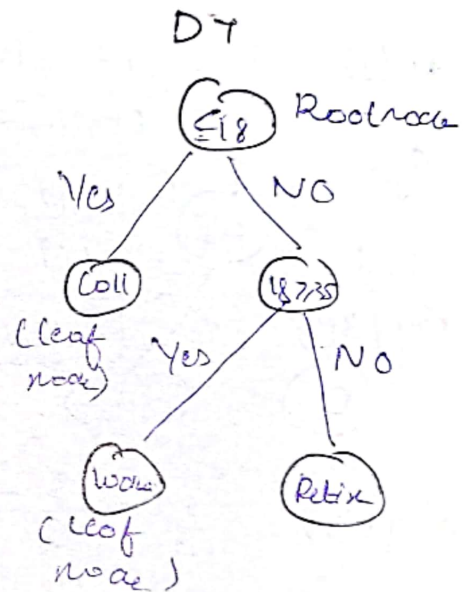
elif (Age > 18 and Age ≤ 35):

print("work")

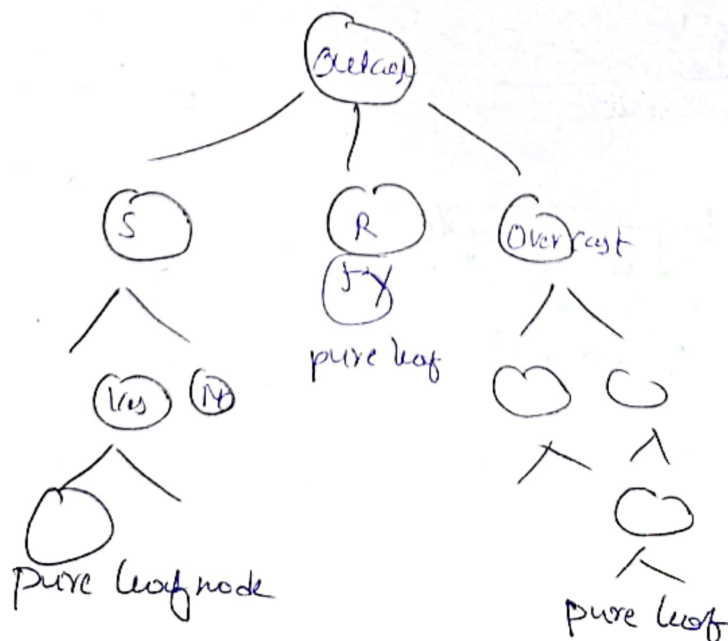
else:

print("retire")

Nested If-else \Rightarrow Decision Tree



eg:-



① purity \rightarrow purity Split?
 (it is used for less number of features)
 \rightarrow Entropy
 \rightarrow Gini Coefficient } ways to find purity
 (it is used for huge features)

② How the features are selected

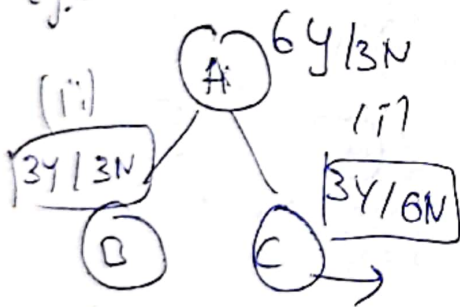
\rightarrow Information Gain?

① Entropy

$$H(S) = -P+ \log_2 P+ - P- \log_2 P$$

(+ = Y, - = No)

Eg:-



$$i) H(S) = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$$

$$= -1 \log_2 1$$

$$= \boxed{0}$$

\downarrow
 Pure split

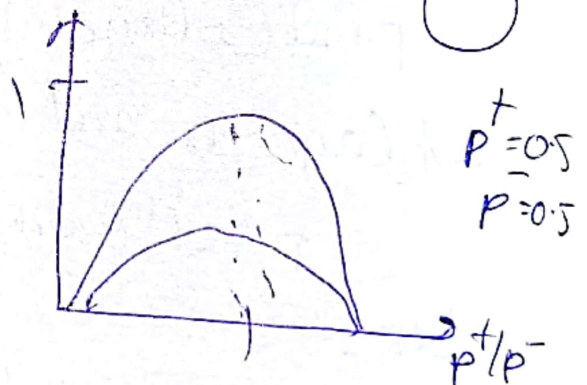
$$ii) H(S) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$= \boxed{1} \text{ Impure split}$$

① GINI Impurity

$$G.I = 1 - \sum_{i=1}^n (p_i)^2$$

$2N/2N$



$$= 1 - [(p_+)^2 + (p_-)^2]$$

$$= 1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right]$$

$$= 1 - \frac{1}{2}$$

$$= 0.5$$

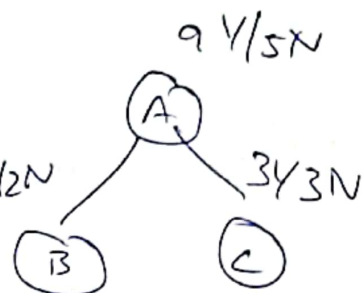
Information Gain

$$\text{Gain}(S, F_1) = H(S) - \sum_{v \in V(F_1)} \frac{|S_v|}{|S|} H(S_v)$$

$$H(S) = -P + \log_2 P + -P - \log_2(P)$$

$$= -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right)$$

$$= \boxed{0.94}$$



$$H(S_{v_B}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$\boxed{H(S_{v_B}) = 0.81}$$

$$H(S_{v_C}) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$\boxed{H(S_{v_C}) = 1}$$

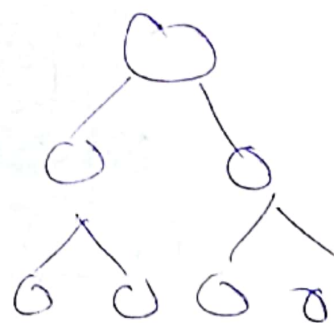
$$\text{Gain}(S, F_1) = 0.94 - \left[\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$\boxed{\text{Gain}(S, F_1) = 0.049}$$

In second feature set

let suppose

$$\boxed{\text{Gain}(S, F_2) = 0.051}$$



$$\boxed{\text{Gain}(S, F_2) > \text{Gain}(S, F_1)}$$

So select 2nd feature because it have high Information Value.