# Decision Tree

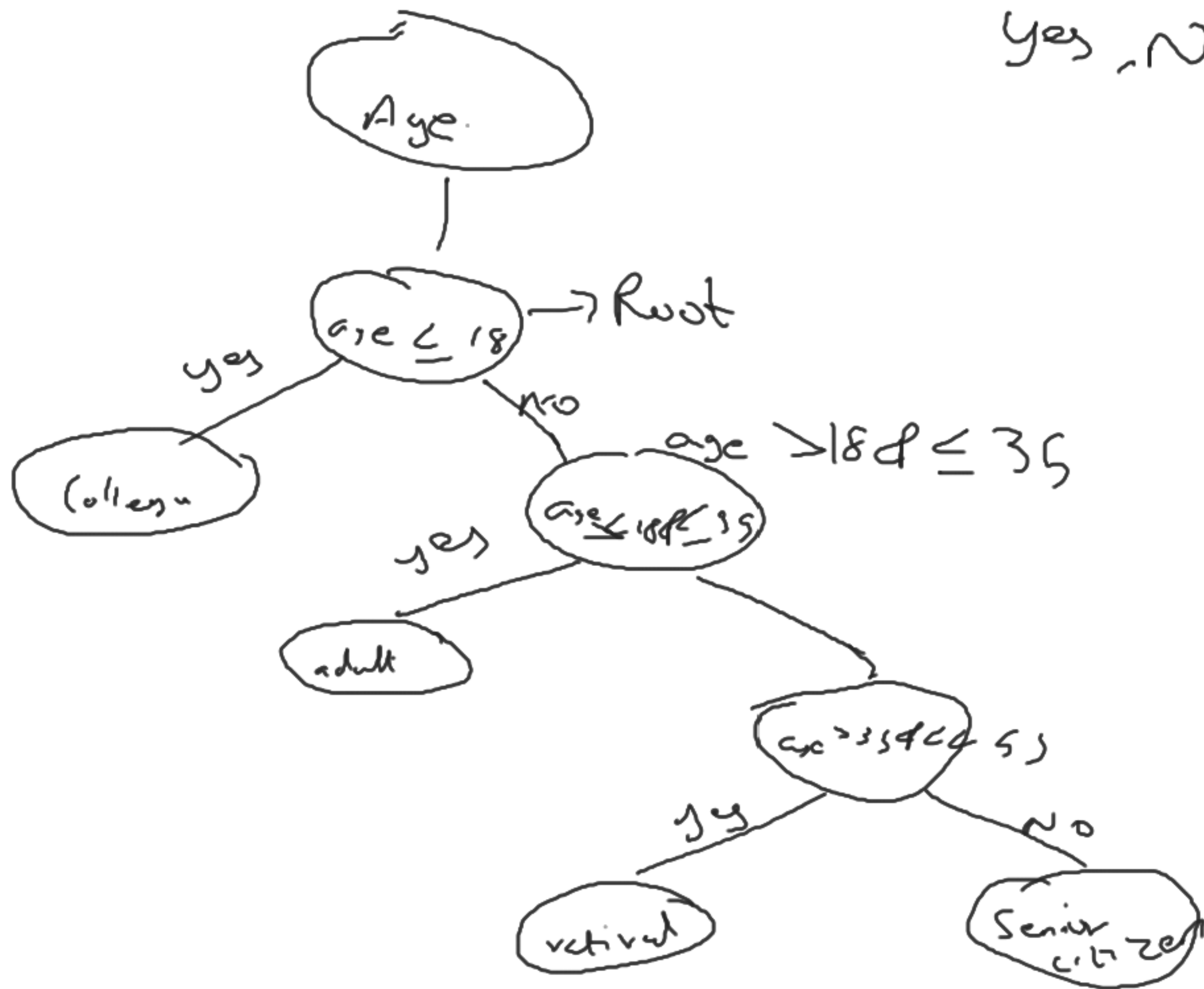- Regression
- Classification

## Age

```
if (age ≤18)
    print (college)
elif (age >18 & ≤35)
    print (adult)
elif (age >35),
    print (retired)
```

age
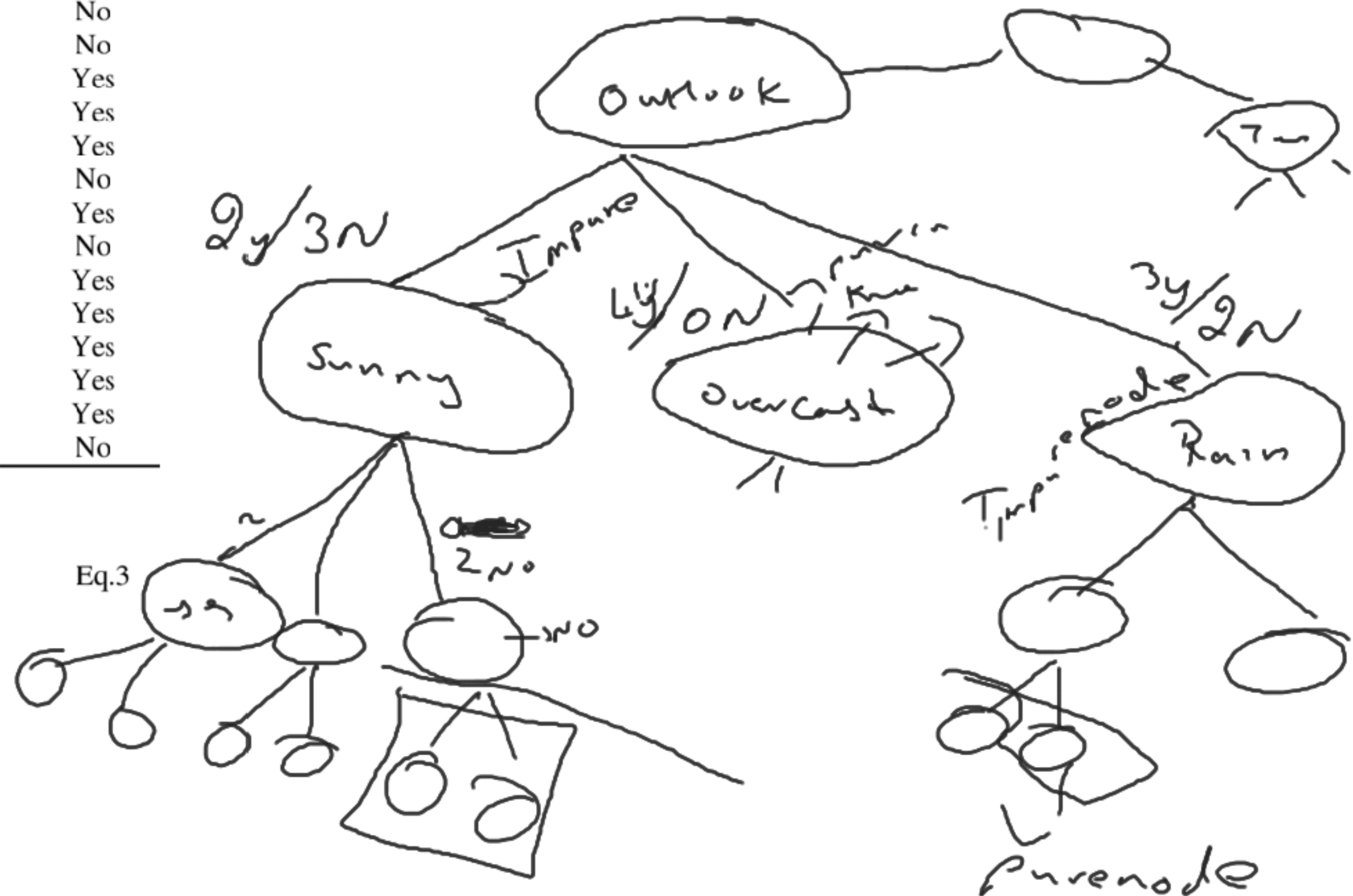
age ≤18 → rool

NO

e...

>18 & ≤

Age

age ≤ 18 → Root

yes — College

no — age >18 & ≤ 35

age >18 & ≤ 35

yes — adult

age >35 & ≤ 55

yes — retired

No — Senior citizen

yes , No

Pure Split

Impure split

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|---|---|---|---|---|---|
| 1 | Sunny | Hot | High | Weak | No |
| 2 | Sunny | Hot | High | Strong | No |
| 3 | Overcast | Hot | High | Weak | Yes |
| 4 | Rain | Mild | High | Weak | Yes |
| 5 | Rain | Cool | Normal | Weak | Yes |
| 6 | Rain | Cool | Normal | Strong | No |
| 7 | Overcast | Cool | Normal | Strong | Yes |
| 8 | Sunny | Mild | High | Weak | No |
| 9 | Sunny | Cool | Normal | Weak | Yes |
| 10 | Rain | Mild | Normal | Weak | Yes |
| 11 | Sunny | Mild | Normal | Strong | Yes |
| 12 | Overcast | Mild | High | Strong | Yes |
| 13 | Overcast | Hot | Normal | Weak | Yes |
| 14 | Rain | Mild | High | Strong | No |

In this example,

$$Entropy(S) = -\left(\frac{9}{14}\right)\log_2\left(\frac{9}{14}\right) - \left(\frac{5}{14}\right)\log_2\left(\frac{5}{14}\right) = 0.9450$$

Eq.3



Purity
- → Entropy → /
- → Gini Impurities (or )Gini Coefficients)

Gain information $\rightarrow$ This Select Particular Column which has good knowlege (or) huge Regarding data and the Output.

Entropy:

$$H(s) = -p^+ \log_2 p^+ - p^- \log_2 p^-$$

$$H(S) = -P_+ \log_2 P_+ + \ominus P_- \log_2 P_-$$

$$= $$

$$H(S) = -\frac{18}{3} + \log_2 \frac{3}{3}$$
$$-\frac{0}{3} - \log_2 \frac{0}{3}$$

Pure node ,,

$f_1$  $6y/3N$ $0$

$3y/3N$ $0$

$3y/0n$

$c_1$

$c_2$

$c_1$

$$H(S) = -\frac{3(}{62} + \log_2 \frac{31}{62} - \frac{31}{62} \log_2 \frac{31}{62}$$

$$= \log_2 = \boxed{1} \quad \text{impure node}$$

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

$c_2$
$$= - -1 \log_2 1$$

$$= 0 \quad \longrightarrow \text{pure node,,}$$

feature Selection $_{it}(1) = -p_+ \log_2 P_+ - P_- \log_2 P_-$

Grain information

Gain $(S, d_1) =$

$\frac{g_y?}{5n}$

No. of sample data

$f_1$

$H(S) - \sum_{v \in \text{val}} \frac{(S_v)}{(S)} c_i$    $(S_v)$

$3y/3n$

$C_2$

overall No. of sample

$6y/2n$

$C_1$

$H(S) \leftarrow - \sum_{v \in \text{val}} \frac{|S_v|}{|S|}$ value

$|S|$ overall sample

$$H(S) = -\frac{9}{14} + \log_2 \frac{9}{14} - \frac{5}{14} - \log_2 \frac{5}{14}$$

$H(S) = \simeq 0.94 \,'''\,'''$

$$G\left(^{s\,:\,\theta\cdots}_{\phantom{s}}\right) = H(s) - \sum_{v\,\in\,Evalue} \frac{|s_v|}{|s|}$$



$$H(s_v\,C_1) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$H(s_v\ C_1) = 0.81$$

$$\approx -\frac{3}{6} \log_2 \frac{3}{6} = -\frac{-3}{6} \log_2 \frac{3}{6}$$

$$H(s_v\ C_2) = 1$$

$$H(s) = 0.94$$

$$H(s \cup c_1) = 0.81$$

$$H(s \cup c_2) = 1$$

$$= H(s) - \sum_{v \in val} \frac{|S_v|}{|S|} \left( \widehat{S_v} \right)$$

feature

$$= \underline{\underline{0.94}} - \left[ \frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

Gain $(s, f_1)$ $= \underbrace{(0.049)}_{Gain} =$

$\text{Gain}(s, f_2) = \underbrace{(0.072)}_{} \longrightarrow$ more knowledge

column

Nodes.

# Gini Impurity

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

0 - r?
|→I~~~

Gin
0 - p~~
0.5 → Imp~~~

$$\frac{z}{w} 2 \qquad G.I = 1 - \left[ (P_+)^2 + (P_-)^2 \right]$$

2y/2N

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

0 - pure,
Gin = 0.5

$$= 1 - 0.5$$
$$= 0.5 \;\; \longrightarrow \; \text{Impure node}_{,,}$$

large dataset

Entropy $\rightarrow$ log $\rightarrow$ runtime

$\rightarrow$ Entropy

min impurity $\rightarrow$ gini impurity

Regression

Continuous

Hyperparameters = Decision
$$6$$
Max depth, Max leaf

Post pruning

pre pruned

3/2 ON

5Y/2 NO