

Agenda

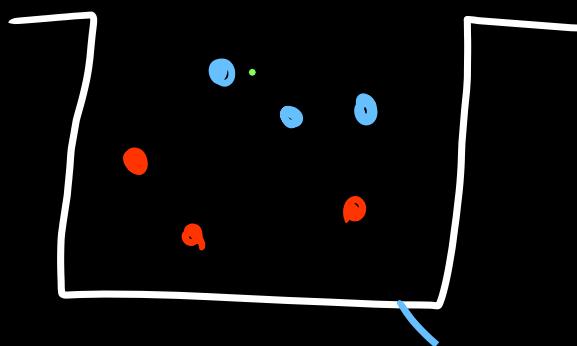
- Gini Index
- Underfitting & Overfitting
- Feature Scaling
- Feature Importance

Entropy

$$H(y) = - \sum_{i=1}^C p_i \log_2 (p_i)$$

→ Gini Impurity

$$GI(y) = 1 - \sum_{i=1}^C p_i^2$$

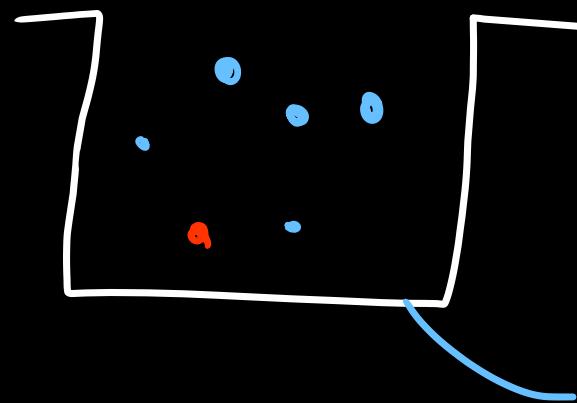


$$\checkmark$$

$$P_b = \frac{3}{6} = \frac{1}{2}$$

$$P_n = \frac{3}{6} = \frac{1}{2}$$

$$H(y) = - \left[\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] = 1$$



$$P_b = \frac{5}{6}$$

$$P_n = \frac{1}{6}$$

$$H(y) = - \left[\frac{5}{6} \log_2 \left(\frac{5}{6} \right) + \frac{1}{6} \log_2 \left(\frac{1}{6} \right) \right]$$

$$= 0.65$$

$$G_I = 1 - \sum P_i^2 = 1 - \left[P_b^2 + P_n^2 \right]$$

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

$$= 0.5$$

$$\underbrace{G(y)}_{\substack{P_b \\ P_n}} = 1 - \left[\left(\frac{5}{6} \right)^2 + \left(\frac{1}{6} \right)^2 \right]$$

$$= 0.277$$

$$\begin{bmatrix} & \ddots \\ \vdots & \ddots \end{bmatrix} \quad \begin{bmatrix} p_b = 1 \\ p_n = 0 \end{bmatrix}$$

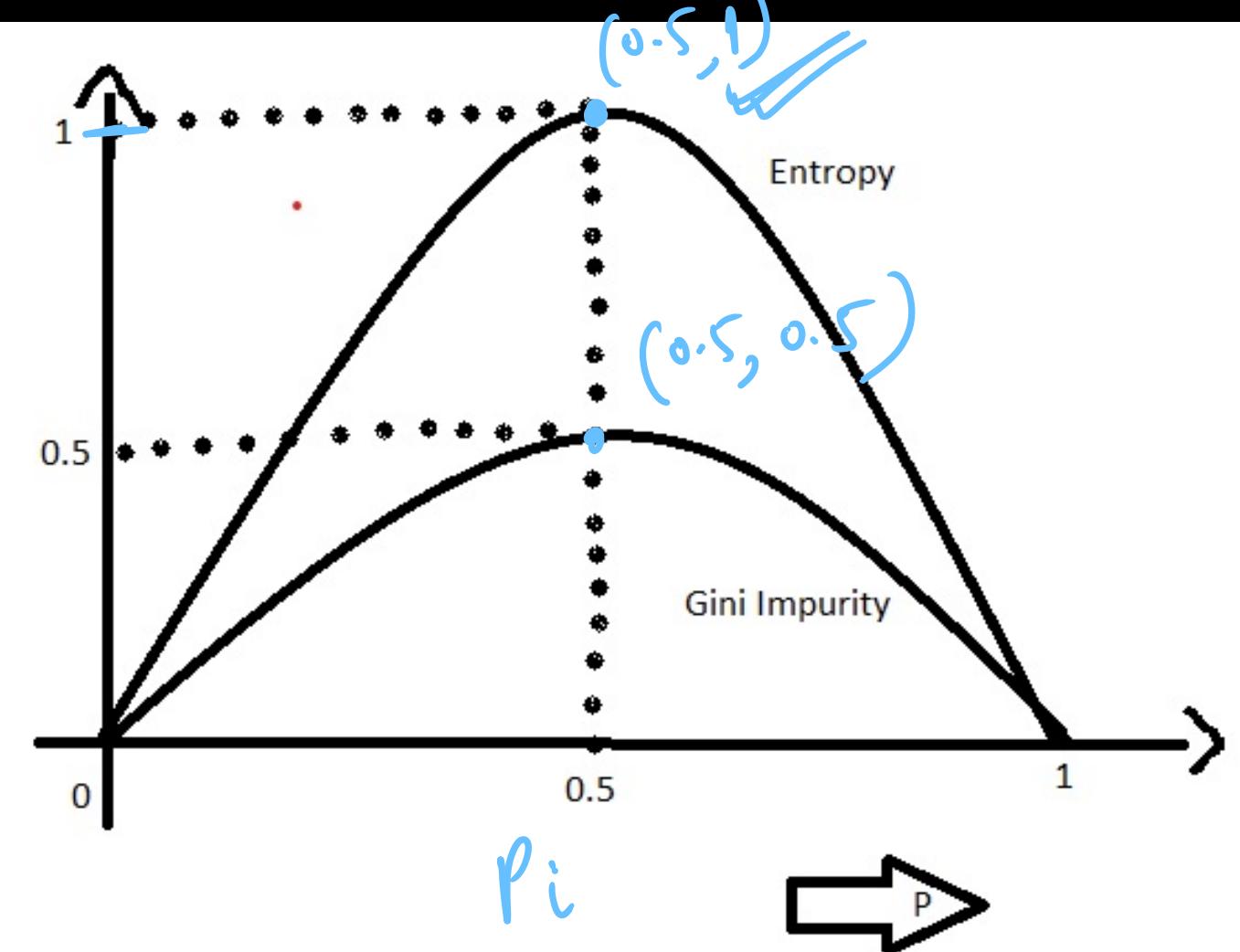
$$G_I(y) = 1 - \left[(1)^2 + (0)^2 \right]$$

$$H(y) = - \left[1 \log_2(1) + 0 \log_2(0) \right] = 0 \quad \checkmark$$

$$\begin{bmatrix} & \ddots \\ \vdots & \ddots \end{bmatrix} \quad \begin{bmatrix} p_b = 0 \\ p_n = 1 \end{bmatrix}$$

$$G_I(y) = 1 - \left[(0)^2 + (1)^2 \right]$$

$$H(y) = - \left[0 \log_2(0) + 1 \log_2(1) \right] = 0 \quad \checkmark$$



$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

$$Gini(E) = 1 - \sum_{j=1}^c p_j^2$$

→ Range of Entropy

$$[0, 1]$$

→ Range of Gini

$$[0, 0.5]$$

1> Entropy

$$-\sum_{i=1}^C p_i \log_2 p_i$$

2> Homogeneous $\rightarrow 0$

3> $p_i = 0.5 \rightarrow H(y) = 1$

4> Slightly slower
due \log_2

Gini Index

$$1 - \sum_{i=1}^C p_i^2$$

Homogeneous $\rightarrow 0$

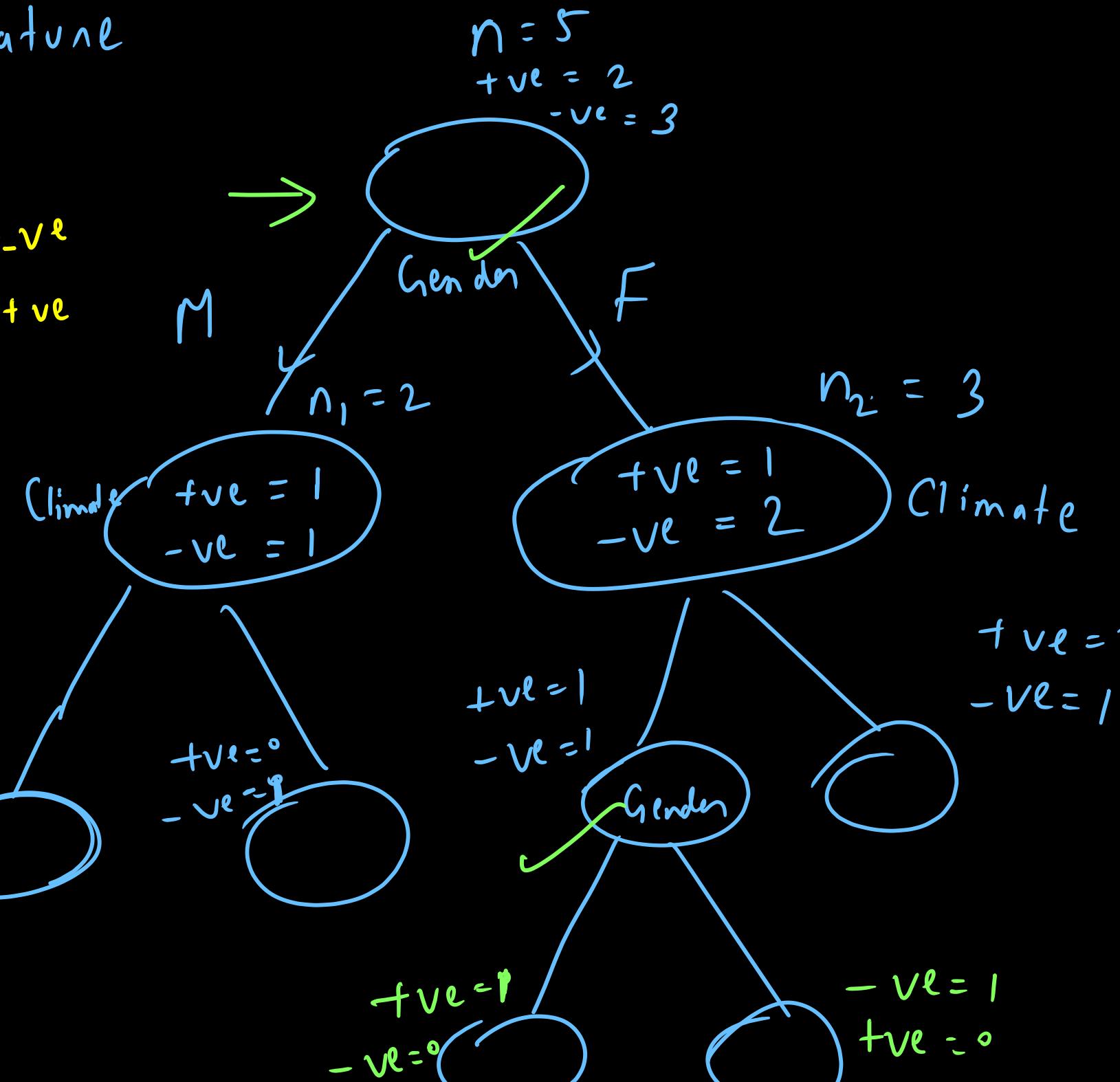
$$p_i = 0.5 \rightarrow GI(y) = 0.5$$

Slightly faster

> Categorical feature

Climate
W → Gender
S → M
F
F
S → M
F

Gender
M
F
F
M
F



Numerical feature

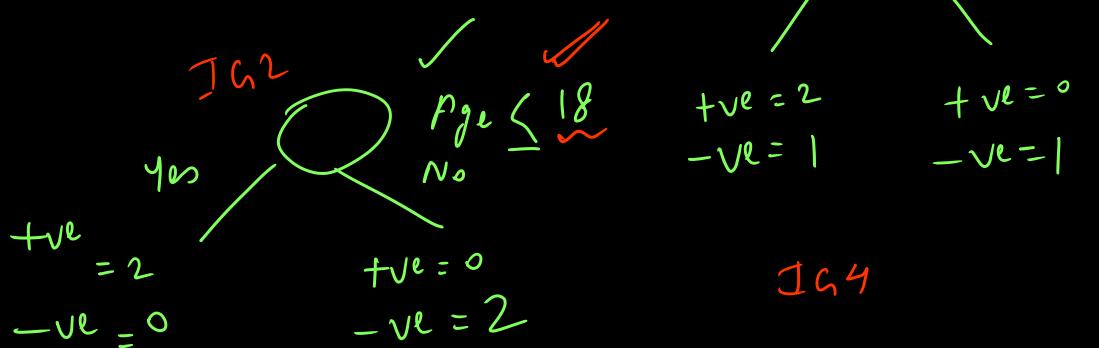
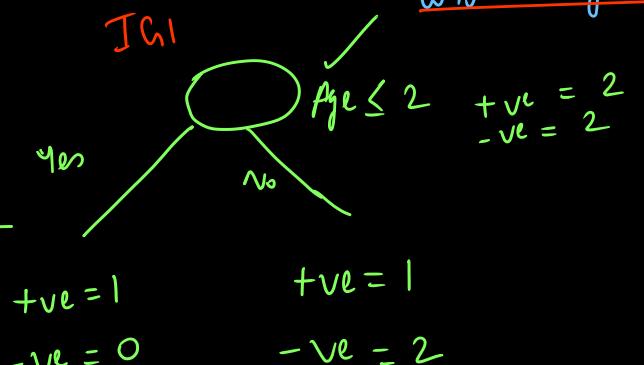
1) Sort the numerical column in ascending order

Age	y
18	1
33	0
41	0
2	1

2) Take each unique value in the numerical feature as threshold and calculate IG_i

3) Choose the threshold with highest IG_i

Age	y
2	1
18	1
33	0
41	0



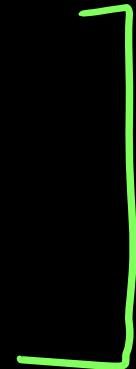
IG4

$$IG_i = H(P) - \left[\frac{n_1}{n} H(c_1) + \frac{n_2}{n} H(c_2) \right]$$

$$IG_i = G_I(P) - \left[\frac{n_1}{n} G_I(c_1) + \frac{n_2}{n} G_I(c_2) \right]$$

Bining

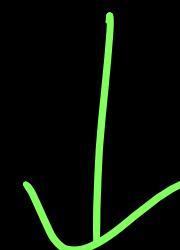
0
5
10



11
14
20



→ Overfitting & Underfitting



Fits training
data perfectly

(noise)

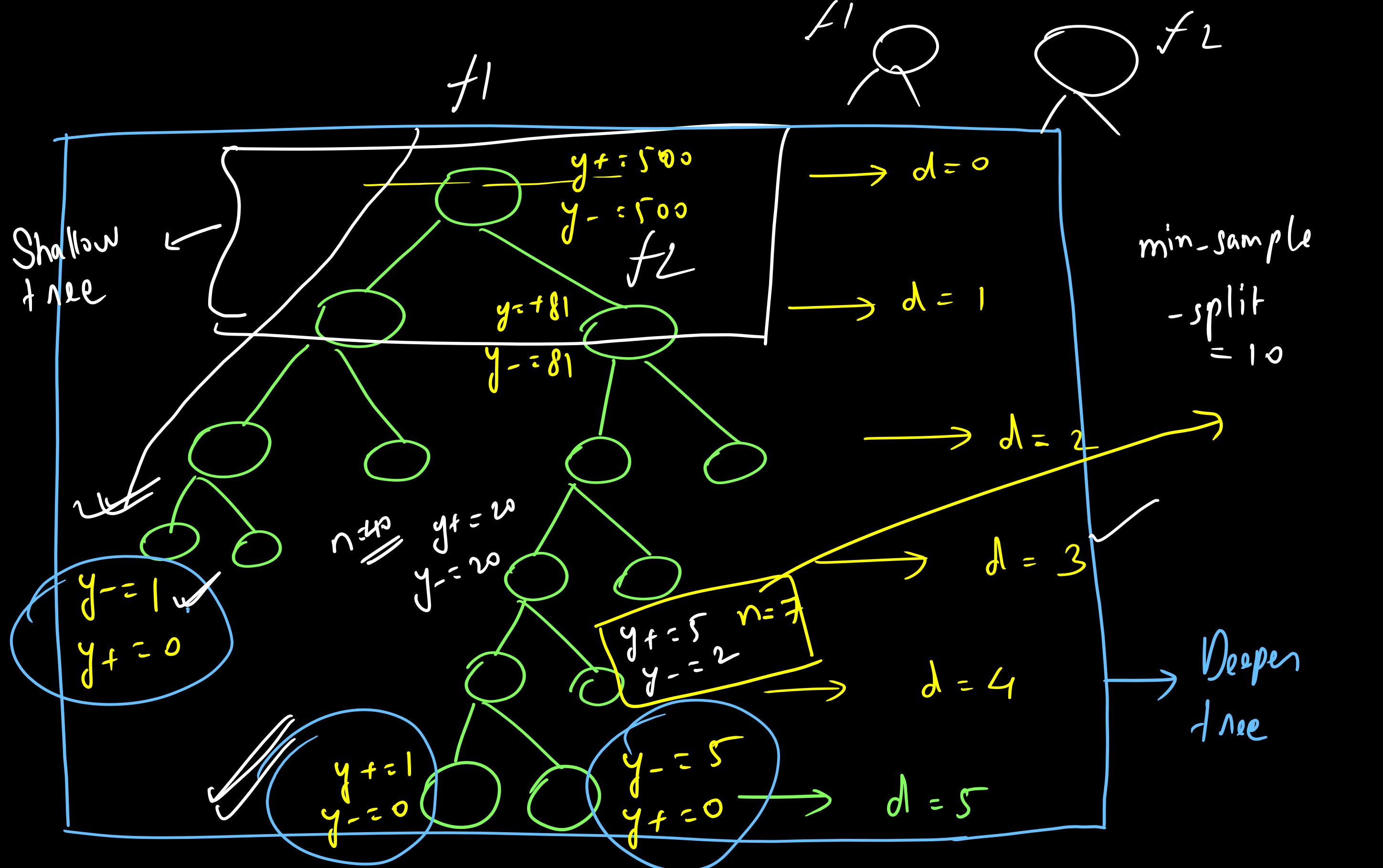
& it
is unable to

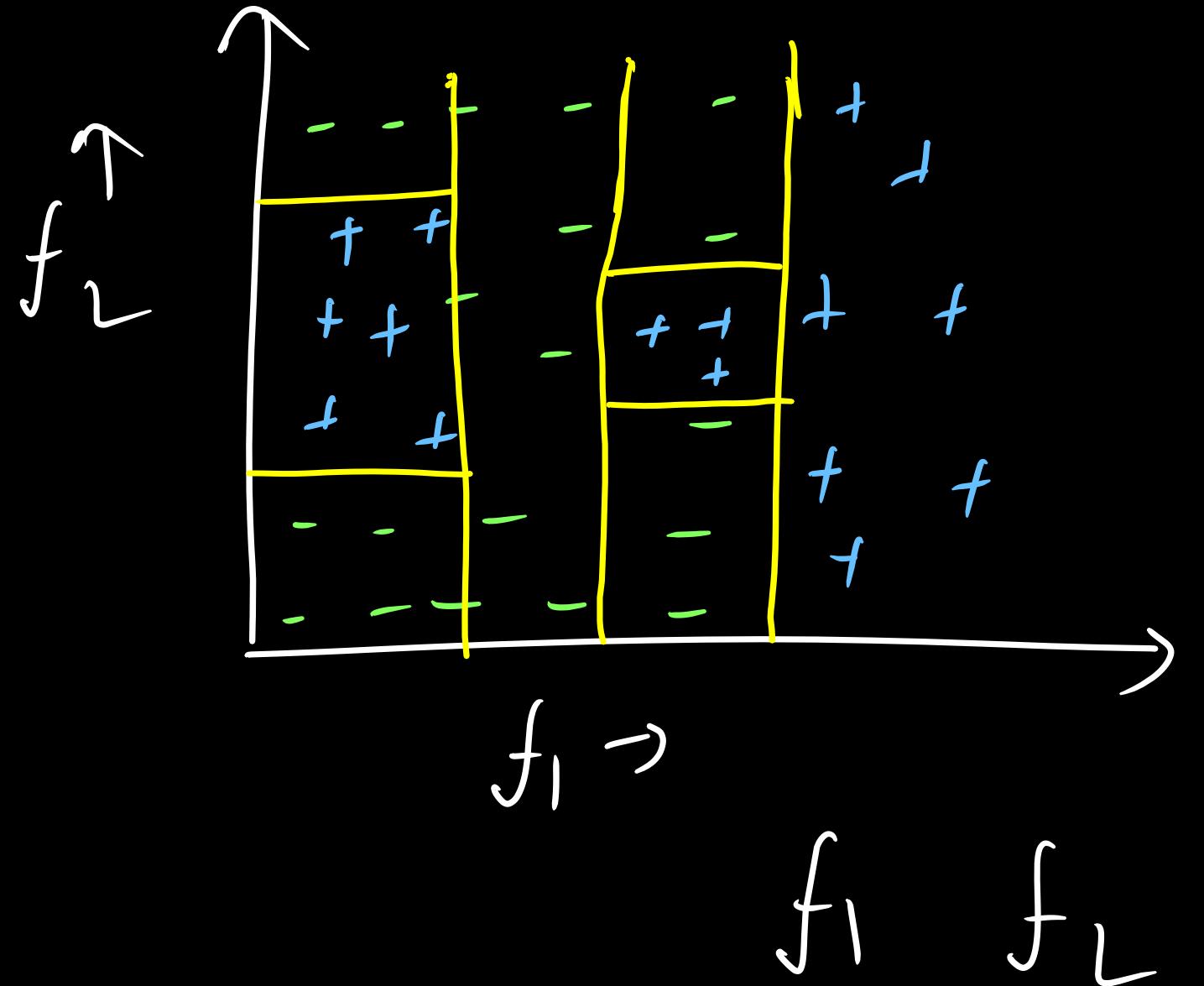
generalize to test data

Does not fit training data
properly

→ Simple model

→ High model complexity





for depth in $[5, 10, 20]$:

model = Decision (depth)

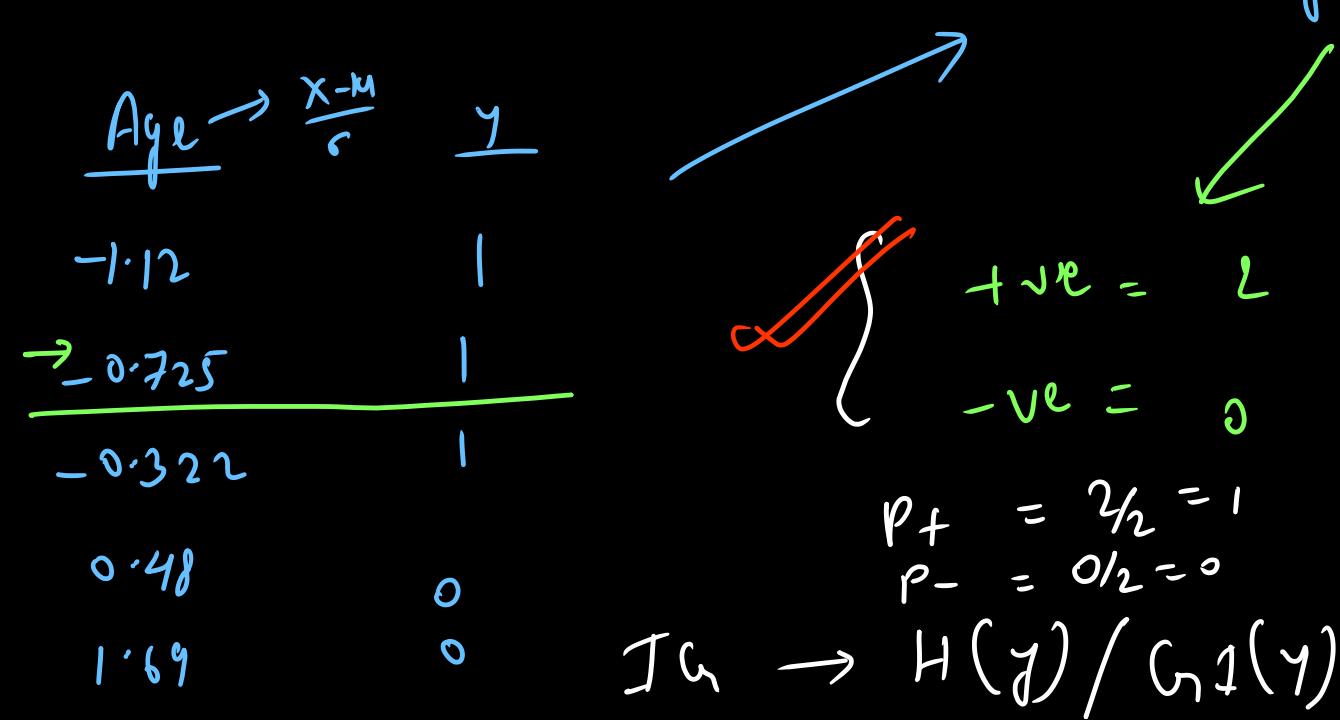
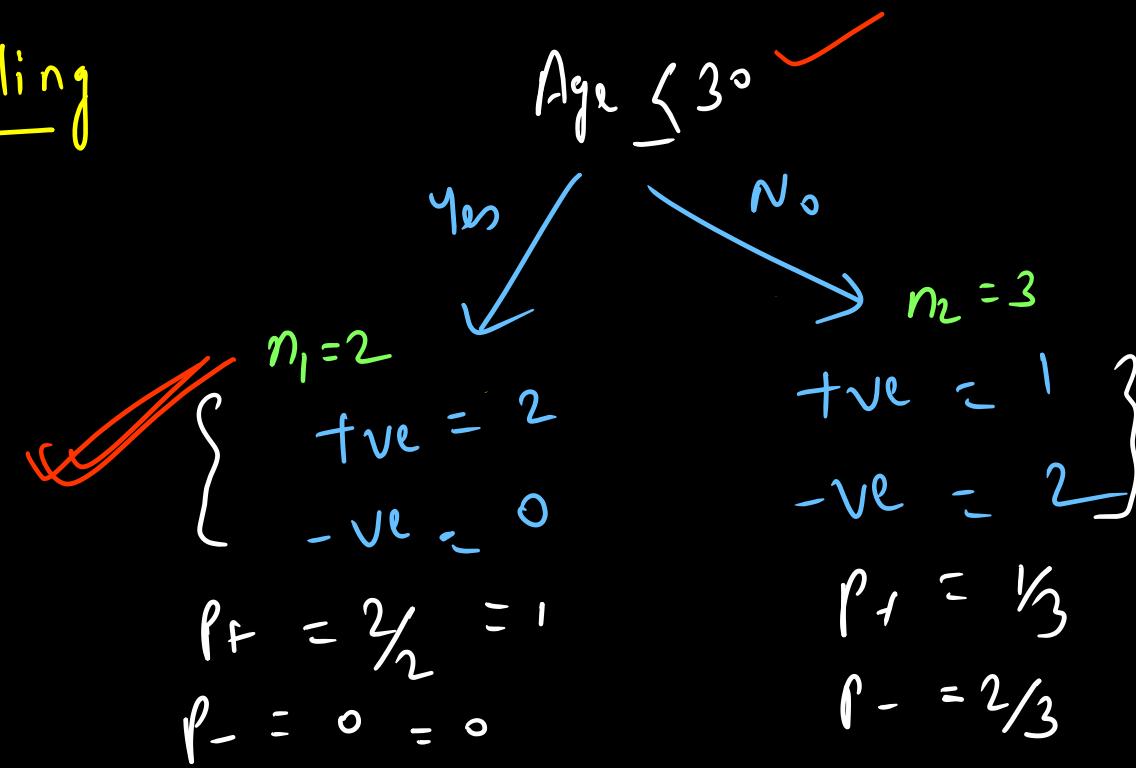
\rightarrow model.fit (X_{train} , y_{train})

\rightarrow model.score (X_{val} , y_{val})

→ Break until 22:24 PM

→ Feature Scaling

<u>Age</u>	<u>y</u>
25	1
30	1
35	1
45	0
60	0



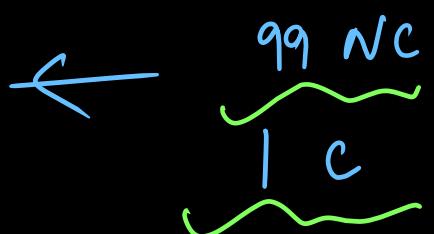
$G_I \rightarrow P_i$

Data Imbalance

$$\frac{y}{N_C}$$

$$N_C$$

C



$$P_C = 0.01$$

$$P_{NC} = 0.99$$

$$\begin{aligned}GI &= 1 - \sum P_i^2 \\&= 1 - [(0.01)^2 + (0.99)^2] \\&= 0.0198 \approx 0\end{aligned}$$



→ For homogeneous → $GI = 0$

- Class weight
- SMOTF
- Under/Overs

→ Class weight

→ NC

$P_{NC} = \frac{4}{5} = 0.8$

$P_C = \frac{1}{5} = 0.2$

$GI = 1 - [0.8^2 + 0.2^2] = 0.32$

$\begin{cases} \text{weight}_C = 4 \\ \text{weight}_{NC} = 1 \end{cases}$

~~Weight NC count
Weight C count~~

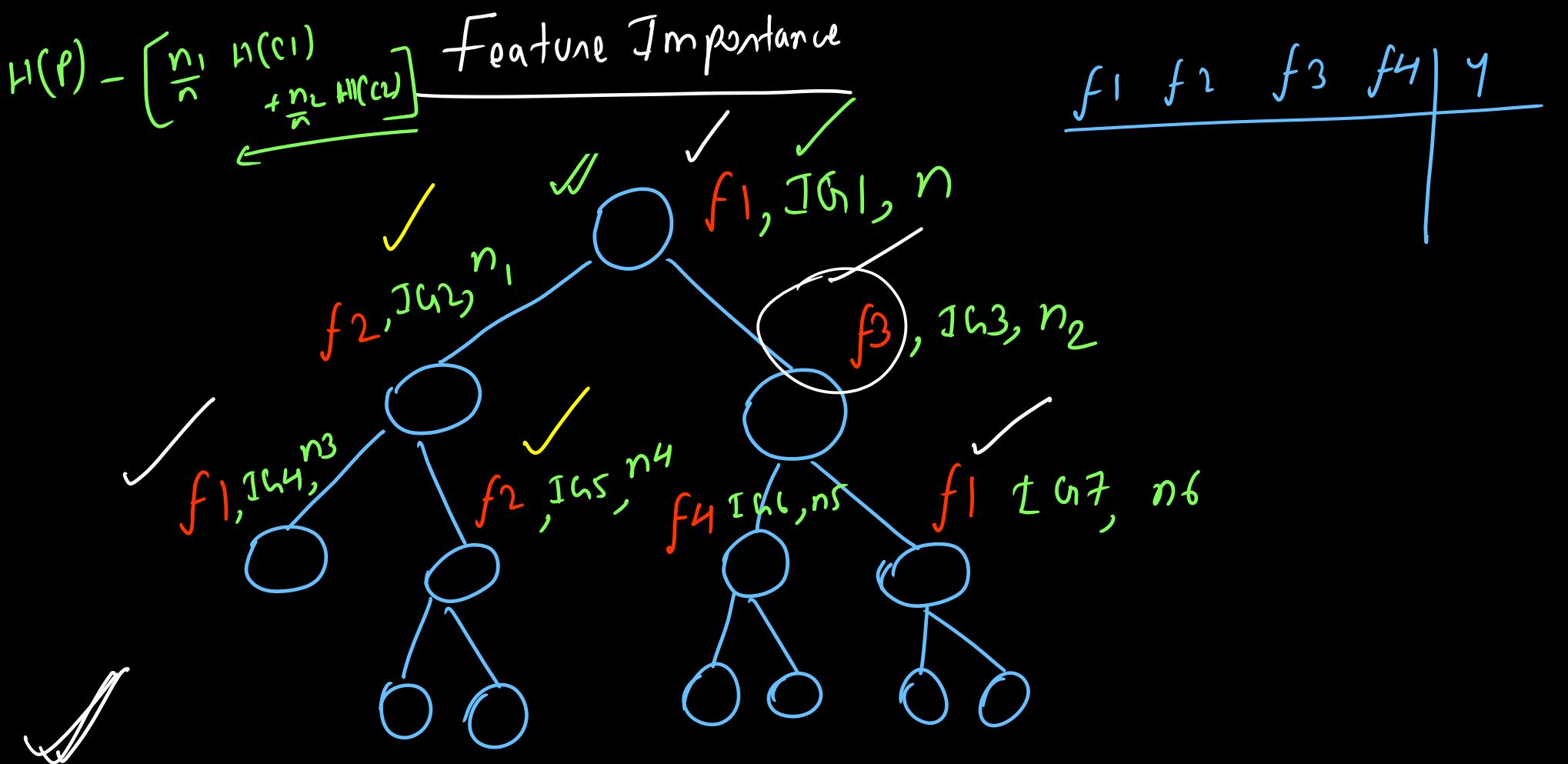
$$= \frac{4 \times 1 = 4}{1 \times 4 = 4}$$

$$P_{NC} = \frac{4}{4+1} = \frac{1}{2}$$

$R = 0.6$
 $G = 0.3$
 $B = 0.1$

$$P_C = \frac{4}{4+4} = \frac{1}{2}$$

$GI = 1 - [0.5^2 + 0.5^2] = 0.5$



$$\text{Feature Importance } f_1 = \frac{n_3}{n} IG_4 + \frac{n_6}{n} IG_7 + \frac{n}{n} IG_1$$

↳ Weight-Avg of IG_i 's where f_1 is used for split

$$\text{Feature Importance } f_2 = \frac{n_4}{n} IG_5 + \frac{n_1}{n} IG_2$$

$$\text{Feature Importance } f_3 = \frac{n_2}{n} IG_3$$

