

Decision Tree

Classification & Regression

- For a given dataset, there can be multiple split conditions
- But then how does decision tree decides which one to choose
 - In human eye it will be best how to mathematically do that
 - We use concept of Entropy and Information Gain

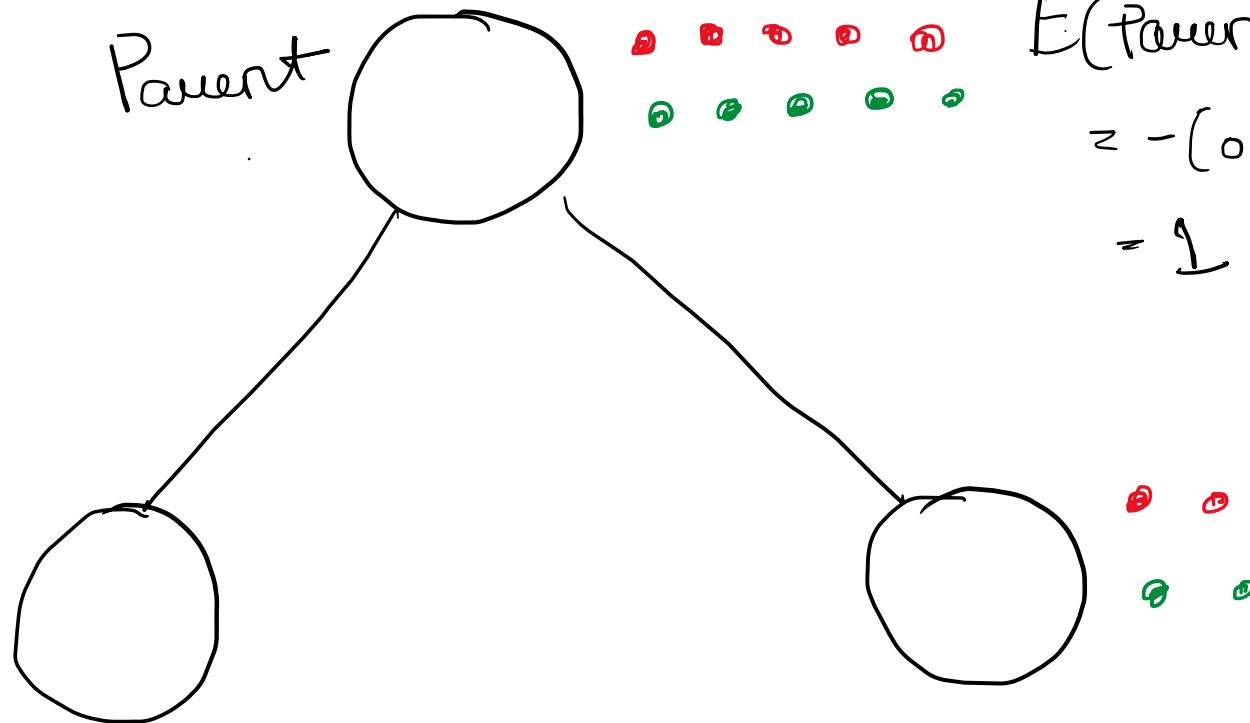
Entropy: The measure of randomness



$$\text{Entropy} = \sum_{i=1}^C -p_i \log(p_i)$$

Example

st possibility



$E(\text{Parent})$

$$\approx -(0.5) \log(0.5) - (0.5) \log(0.5)$$
$$= 1$$

$$E(C_1) = -\frac{3}{6} \log\left(\frac{3}{6}\right) - \frac{3}{6} \log\left(\frac{3}{6}\right)$$

$$= 1$$

:

$$E(C_2) = -\frac{2}{4} \log\left(\frac{2}{4}\right) - \frac{2}{4} \log\left(\frac{2}{4}\right)$$
$$= 1$$

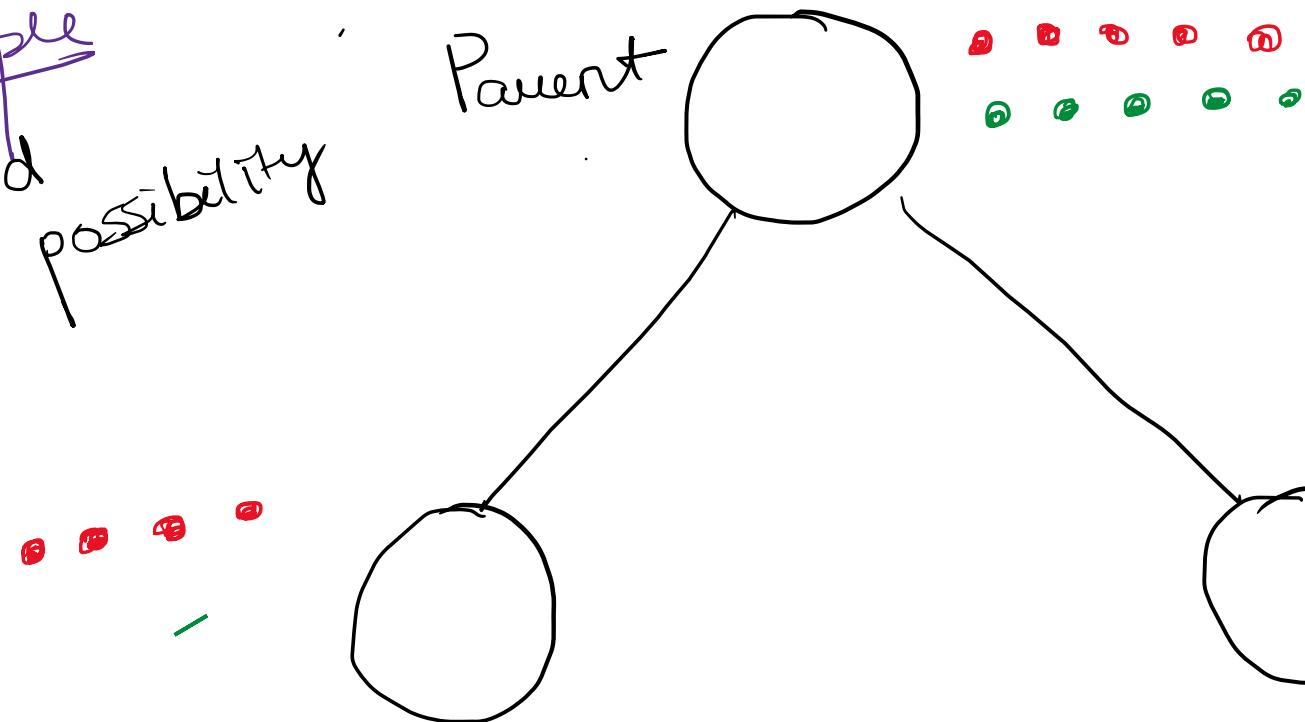
$$\text{Information Gain} = E(\text{parent}) - \sum w_i E(\text{child}_i)$$

$$= 1 - \left[\frac{6}{10} \times 1 + \frac{4}{10} \times 1 \right]$$

$$= 1 - \left[\frac{6}{10} + \frac{4}{10} \right] = 1 - \frac{10}{10} = 1 - 1 = 0$$

Example

- 2^{n d} possibility



$$E(C_1) = -(\log(1))$$

$$= 0$$

$E(\text{Parent})$

$$\approx -(0.5)\log(0.5) - (0.5)\log(0.5)$$

$$= 1$$



$$E(C_2) = -\frac{1}{6}\log\left(\frac{1}{6}\right) - \frac{5}{6}\log\left(\frac{5}{6}\right)$$

$$\approx -(0.17)(-0.48) -$$

$$(0.83)(-0.08)$$

$$= 0.132 + 0.06$$

$$= 0.192$$

$$\text{Information Gain} = E(\text{parent}) - \sum w_i E(\text{child}_i)$$

$$= 1 - \left[0 \times \frac{4}{10} + 0.192 \times \frac{6}{10} \right]$$

$$\geq 1 - 0.1152 = \underline{\underline{0.8848}}$$

$\therefore IG(1^{\text{st}}) < IG(2^{\text{nd}})$ \therefore 2nd split it is
 \equiv

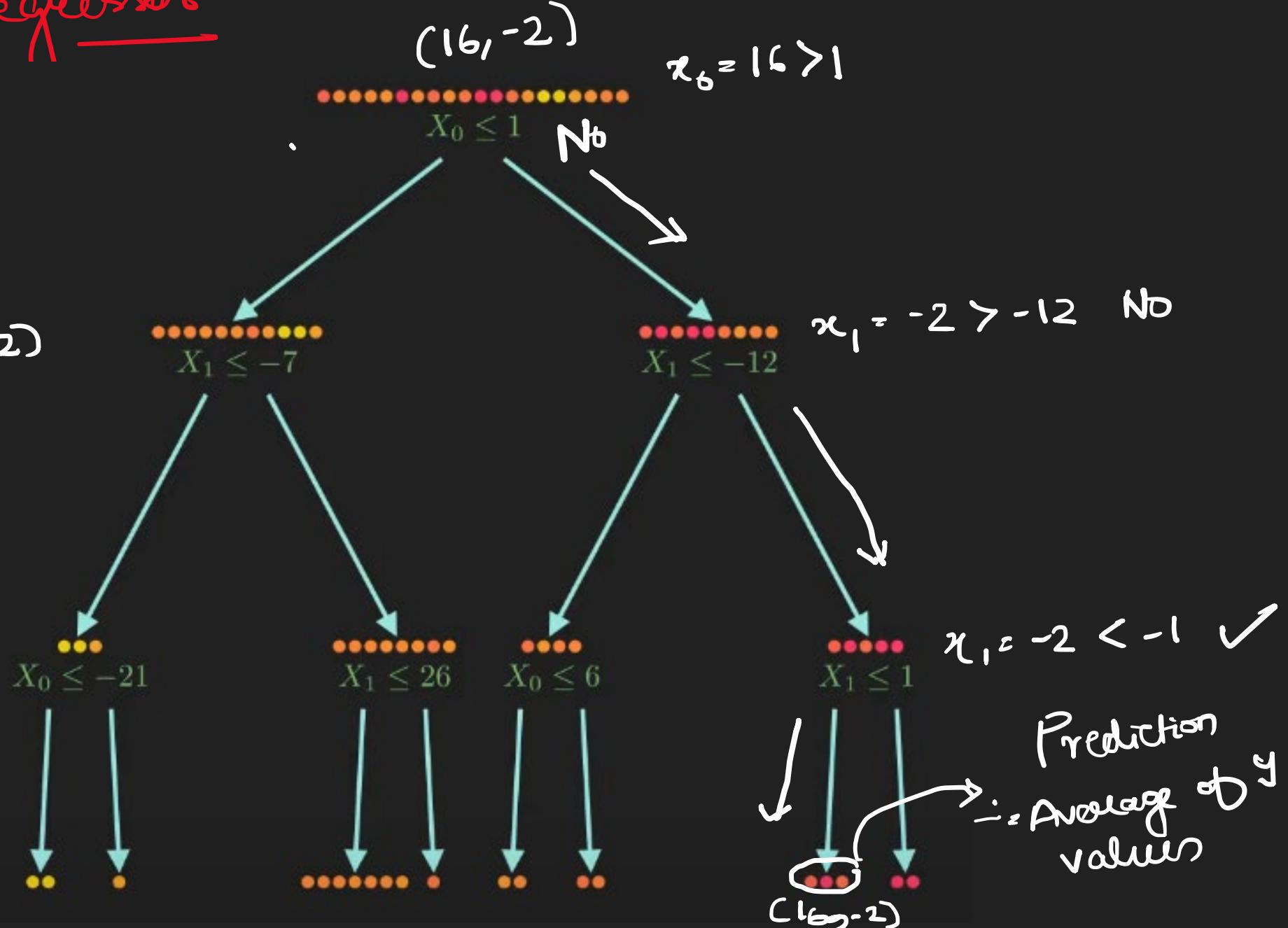
This way the algorithm calculates all the way for all possible splitting possibilities.

- Hence we get our best split and decision tree

Decision tree Regressor

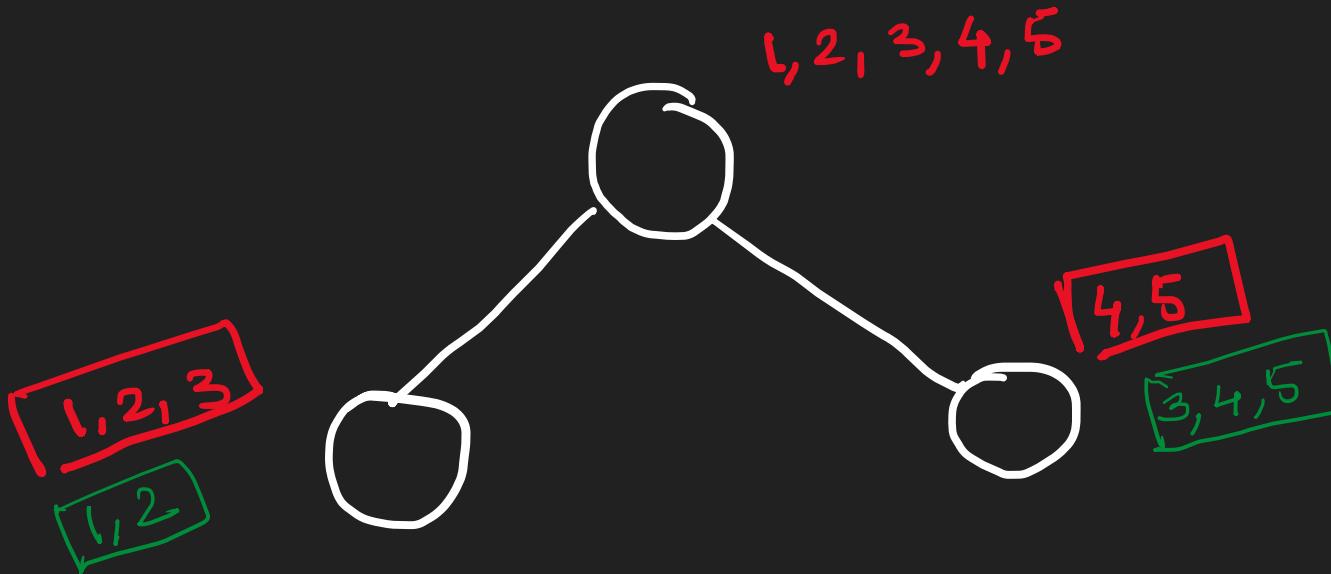
Suppose we have
a new point

$$(x_0, x_1) = (16, -2)$$



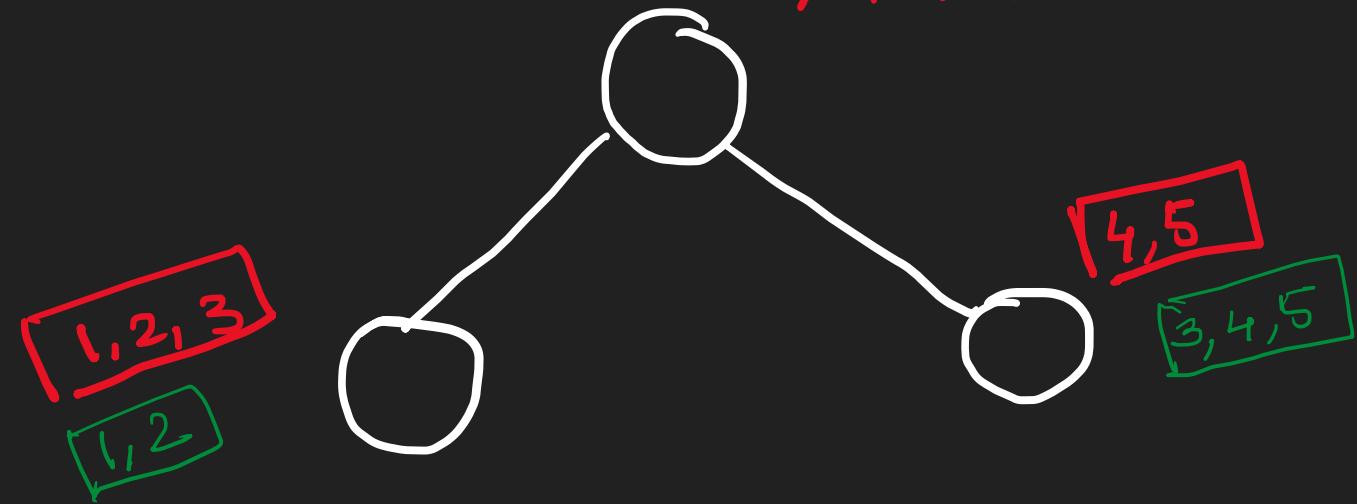
But how to split

We use a parameter called variance reduction



1, 2, 3, 4, 5

$$\bar{y} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$



$$\text{Var}_{\text{Parent}} = \frac{1}{n} \sum (y_i - \bar{y})^2 = \frac{1}{5} [(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2]$$
$$= \frac{1}{5} [4 + 1 + 0 + 1 + 4] = \frac{10}{5} = 2$$

$\boxed{1, 2, 3}$

V_1

$\boxed{4, 5}$

V_2

$$V_1 \text{ mean} = \frac{1+2+3}{3} = 2$$

$$V_2 \text{ mean} = \frac{4+5}{2} = 4.5$$

$$V_1 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{3} \left[(1-2)^2 + (2-2)^2 + (3-2)^2 \right] = \frac{1}{3} [1+0+1] = \frac{2}{3} = 0.67$$

$$V_2 = \frac{1}{2} \left[(4-4.5)^2 + (5-4.5)^2 \right] = \frac{1}{2} [0.25 + 0.25] = \frac{1}{2} \times 0.5 = 0.25$$

$$\boxed{12}$$
$$\sqrt{1}$$

$$\boxed{3 \ 4 \ 5}$$
$$\sqrt{2}$$

$$\sqrt{1}_{\text{mean}} = \frac{1+2+3}{3} = 3/2 = 1.5$$

$$\sqrt{2}_{\text{mean}} = \frac{3+4+5}{3} = 4 \quad \frac{12}{3} = 4$$

$$\sqrt{1} = \frac{1}{2} [(1-1.5)^2 + (2-1.5)^2] = 0.25$$

$$\sqrt{2} = \frac{1}{3} [(3-4)^2 + (4-4)^2 + (5-4)^2] = \frac{1}{3} * 2 = \frac{2}{3} = 0.67$$

For split 1

$$\text{Var}_{\text{Red}} = \text{Var}(\text{Parent}) - \sum w_i V(\text{child}_i)$$

$$= 2 - \left[\frac{3}{5} \times 0.67 + \frac{2}{5} \times 0.25 \right]$$

$$= 2 - [0.402 + 0.1]$$

$$= 2 - 0.502 = 1.498$$

For split 1

$$\text{Var}_{\text{Red}} = \text{Var}(\text{Parent}) - \sum w_i V(\text{child}_i)$$

$$= 2 - \left[\frac{2}{5} \times 0.25 + \frac{3}{5} \times 0.67 \right]$$

$$= 2 - [0.402 + 0.1]$$

$$= 2 - 0.502 = 1.498$$

So Variance Reduction for both splits are
same, So we can Select any one of them

for our split -

