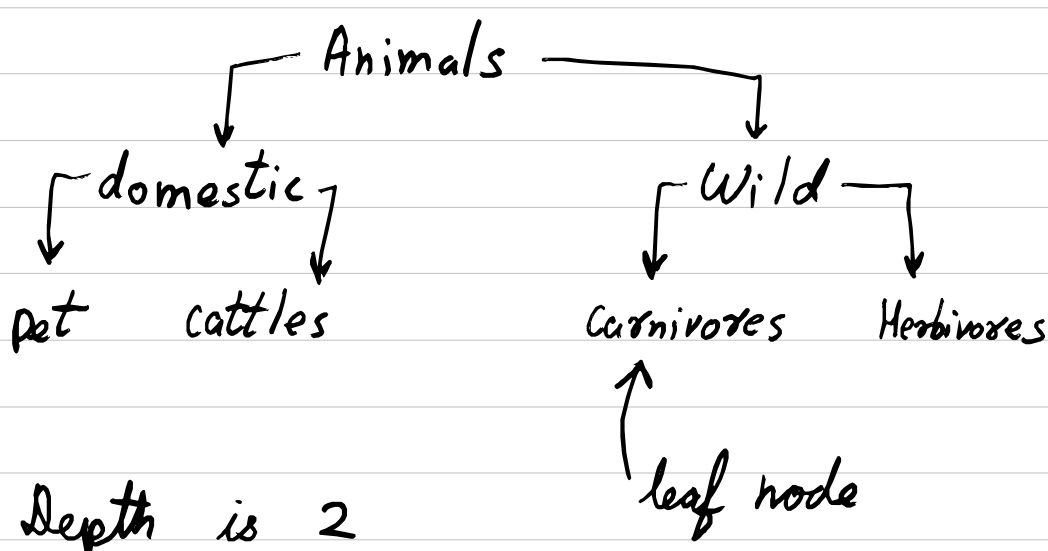


# Decision Tree

\* Used for classification

A dataset can be classified by various parameters. Decision Trees are used to select such parameter

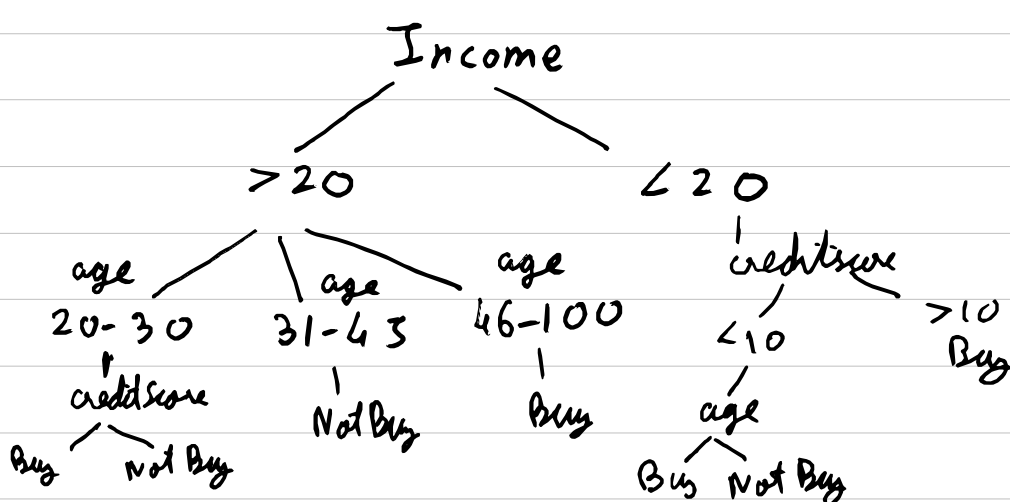


Example if below data is provided to machine. Machine has to decide if person can buy BMW or not

| Income | age | credit Score | buy BMW |
|--------|-----|--------------|---------|
| 20     | 40  | 100          | Y       |
| 10     | 25  | 20           | N       |
| 30     | 45  | 90           | Y       |
| 12     | 45  | 60           | N       |
| 40     | 46  | 80           | Y       |
| 40     | 46  | 80           | Y       |

Machine has to decide which criterion to select for classification

eg one such criterion can be



Machine uses entropy & gain for classification to decide which is the best criterion for such a classification

There are huge number of calculation for this model.

Aim of the model is to decide which tree is the best.

If there are multiple parameters then calculations increase rapidly.

Decision tree & CART both are Top down Non backtracking greedy approaches.

Recursive Partitioning algorithms

# Entropy

Entropy is randomness or "impurity" of an attribute

$$H = - \sum_{i=1}^n P(x_i) \log_2(x_i)$$

when  $P(x) = 0.5$ ,  $H = 1$

① calculate entropy for Gender

| Gender | count |
|--------|-------|
| Male   | 9     |
| Female | 5     |

$$\rightarrow P(\text{male}) = \frac{9}{14}$$

$$P(\text{Female}) = \frac{5}{14}$$

$$H = - \left( \frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14} \right)$$
$$= -0.94$$

Entropy is 0 if outcome is certain

Entropy is maximum if we have no knowledge of the system, or if any outcome is equally possible

# Information gain

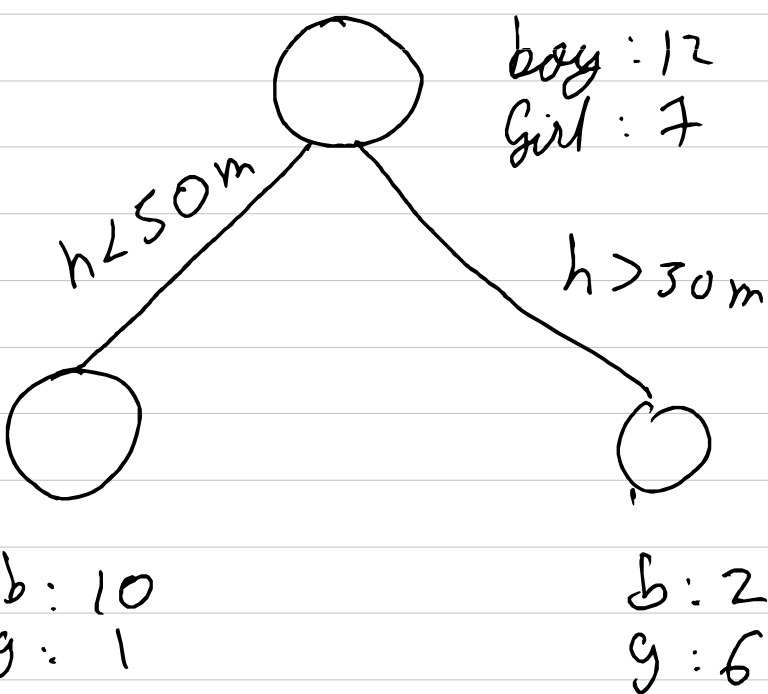
Information is in form of entropy.

$$\text{Info gain} = \text{Info Parent} - \text{Info attribute}$$

$\downarrow$   
 $H(\text{output})$

$\downarrow$   
weighted sum of  $H$   
of all unique values of  
that attribute (feature)

Calculate Info gain for split of height  
Target  $\rightarrow$  gender



$$\text{Info parent} = H(\text{Parent})$$

$$= \frac{12}{19} \log\left(\frac{12}{19}\right) + \frac{7}{19} \log\left(\frac{7}{19}\right)$$
$$= -0.949$$

| $h < 50$  | $h > 50$  |
|---|---|
| $H = \frac{10}{11} \log\left(\frac{10}{11}\right) + \frac{1}{11} \log\left(\frac{1}{11}\right)$ | $H = \frac{2}{8} \log\left(\frac{2}{8}\right) + \frac{6}{8} \log\left(\frac{6}{8}\right)$ |
| $= +0.439$  | $= +0.811$  |

$$\text{Info split} = \frac{11}{19} (0.439) + \frac{8}{19} (0.811)$$
$$= +0.5956$$

$$\text{Info gain} = 0.949 - 0.5956$$
$$= 0.3533$$

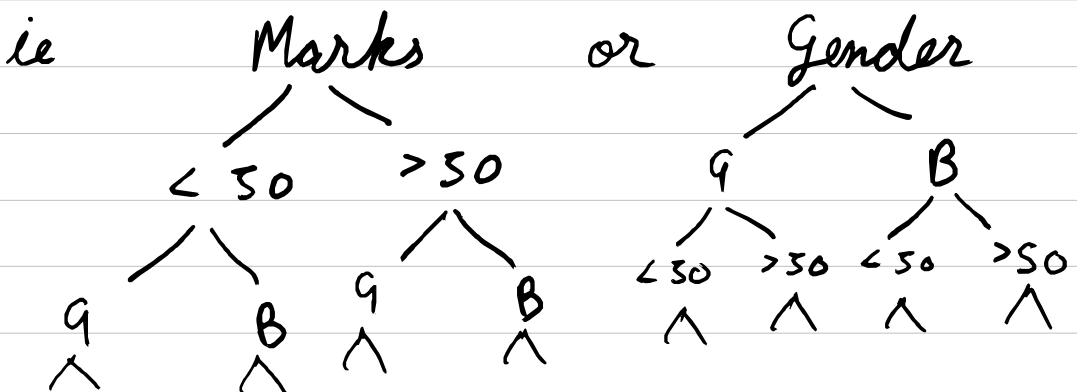
# Decision tree

which attribute would you decide to split the dataset into a decision tree

| Gender | Marks |
|--------|-------|
|--------|-------|

|   |    |
|---|----|
| G | 65 |
| G | 46 |
| B | 56 |
| B | 43 |
| B | 53 |
| B | 49 |
| G | 42 |
| B | 84 |
| B | 44 |
| G | 42 |
| G | 40 |

should data be split on the basis of gender or marks?



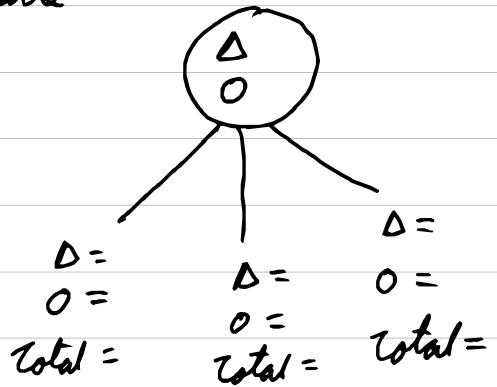
These decisions will be taken by the decision tree based on the info gain

split with the highest info gain is chosen

# Decision tree algorithm

$$H = - \sum P \log_2 P$$

- step 1. find  $H$  of Parent  
step 2. select an attribute &  
split data in that attribute  
count the number of Target classes  
for the attribute



- step 3. calculate entropy for every split  
step 4. Entropy for the attribute is the weighted sum of the attributes ( $\sum p \cdot H$ )  
step 5. Calculate Information gain  
step 6. Repeat steps 2-4 for all attributes  
select lowest value of entropy is the highest value of info gain  
step 7. split the dataset based on the parameter into parts  
step 8. Repeat the process for all parts  
note  $\rightarrow$  dataset will be smaller for a part  
Exclude the attribute for the part

Generally entropy values will be in range 0.8 - 0.99  
or 0 or 1

Don't forget to take  $\log_2$

The lower the value of  $H$ , that means more accurate or skewed the split

eg If split is like 6:4  $H$  will be towards 1  
9:1  $H$  will be lesser

Also  $0/\log_2 0 = 0$  consider

Class labelled Training tuples

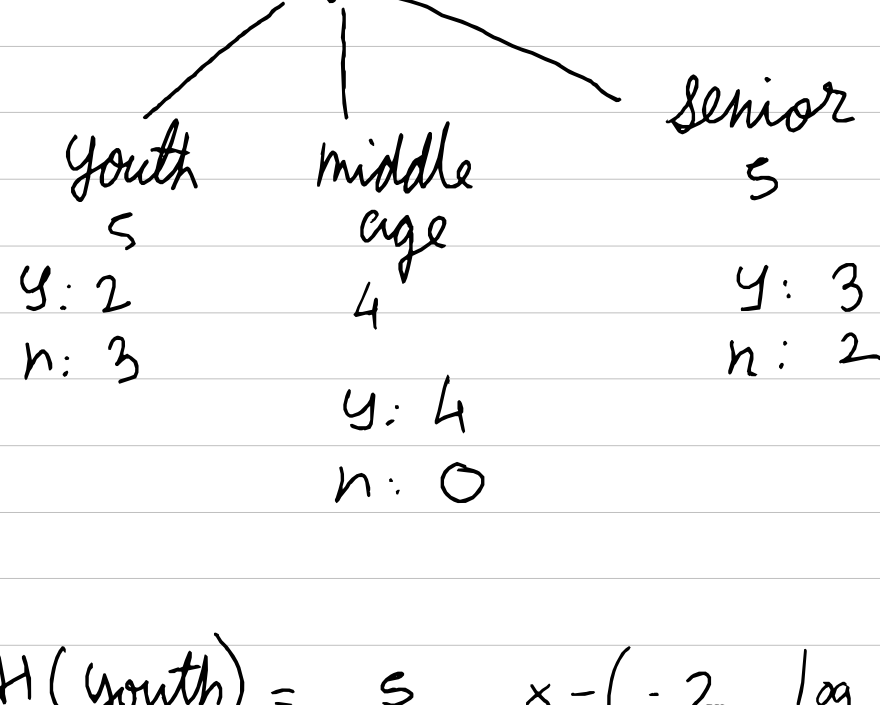
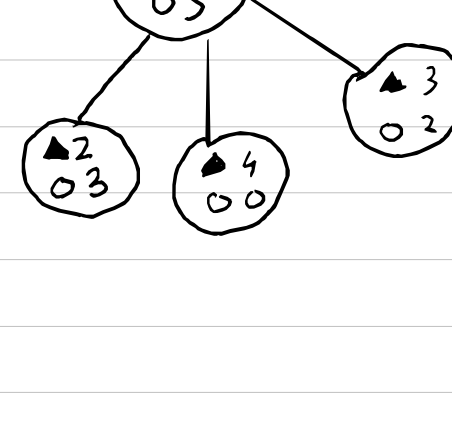
| id | age   | income | student | crediting | buy Car |
|----|-------|--------|---------|-----------|---------|
| 1  | YOUTH | high   | no      | fair      | no      |
| 2  | y     | h      | no      | excellent | no      |
| 3  | m-a   | h      | no      | f         | y       |
| 4  | S     | medium | no      | f         | y       |
| 5  | S     | low    | yes     | f         | y       |
| 6  | S     | L      | yes     | e         | n       |
| 7  | m-a   | L      | yes     | e         | y       |
| 8  | y     | m      | no      | f         | n       |
| 9  | y     | L      | y       | f         | y       |
| 10 | S     | m      | y       | f         | y       |
| 11 | y     | m      | y       | e         | y       |
| 12 | ma    | m      | n       | e         | y       |
| 13 | ma    | h      | y       | f         | y       |
| 14 | S     | m      | n       | e         | n       |

info(D) on the target value

$$\text{info}(D) = -\left(\frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14}\right) \quad \begin{matrix} \blacktriangle 9 \\ \bullet 5 \end{matrix}$$

$$= -0.940$$

info(D)<sub>age</sub> → trying root, age



$$H(\text{youth}) = \frac{5}{14} \times -\left(-\frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5}\right)$$

↑ weighted sum

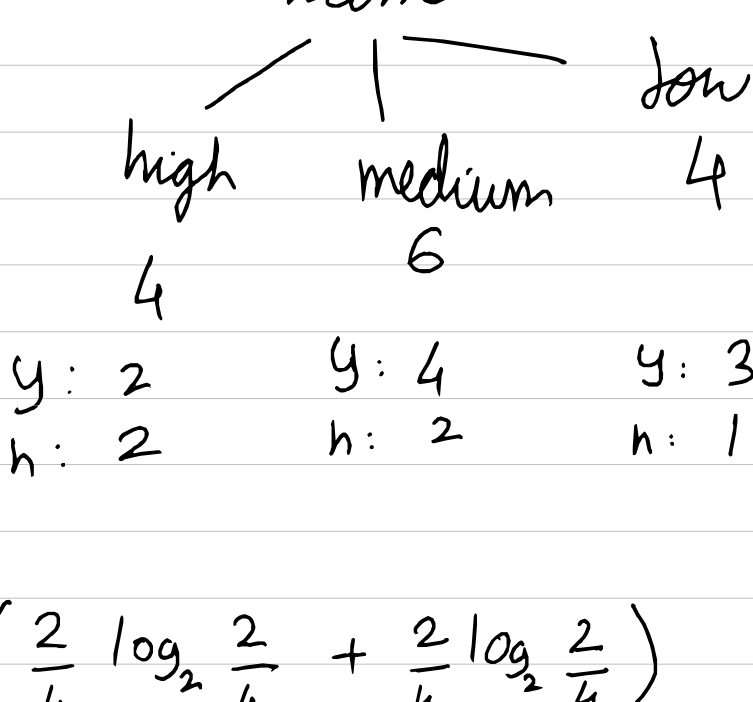
$$H(\text{middle}) = \frac{4}{14} \left(-\frac{4}{4} \log_2 \frac{4}{4} - 0 \log_2 0\right) = 0$$

$$H(\text{senior}) = \frac{5}{14} \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}\right) \quad \text{marked with *}$$

$$H(\text{age}) = H_1 + H_2 + H_3 = 0.693$$

$$\text{Info gain} = 0.940 - 0.693 = \underline{0.247}$$

for Income



$$H = \frac{4}{14} \left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4}\right) \rightarrow \text{for high}$$

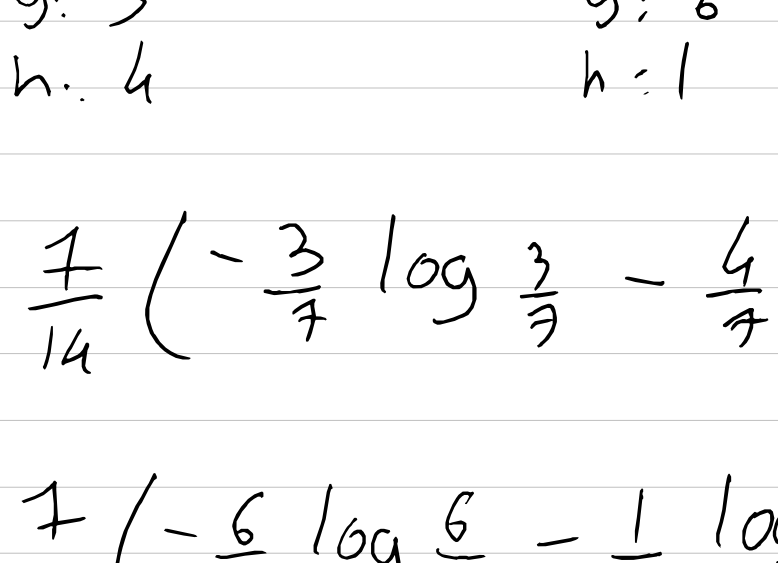
$$+ \frac{6}{14} \left(\frac{4}{6} \log_2 \frac{4}{6} + \frac{2}{6} \log_2 \frac{2}{6}\right) \rightarrow \text{for medium}$$

$$+ \frac{4}{14} \left(\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4}\right) \rightarrow \text{for low}$$

$$= 0.911$$

$$\text{Info gain} = 0.940 - 0.911 = \underline{0.029}$$

for student



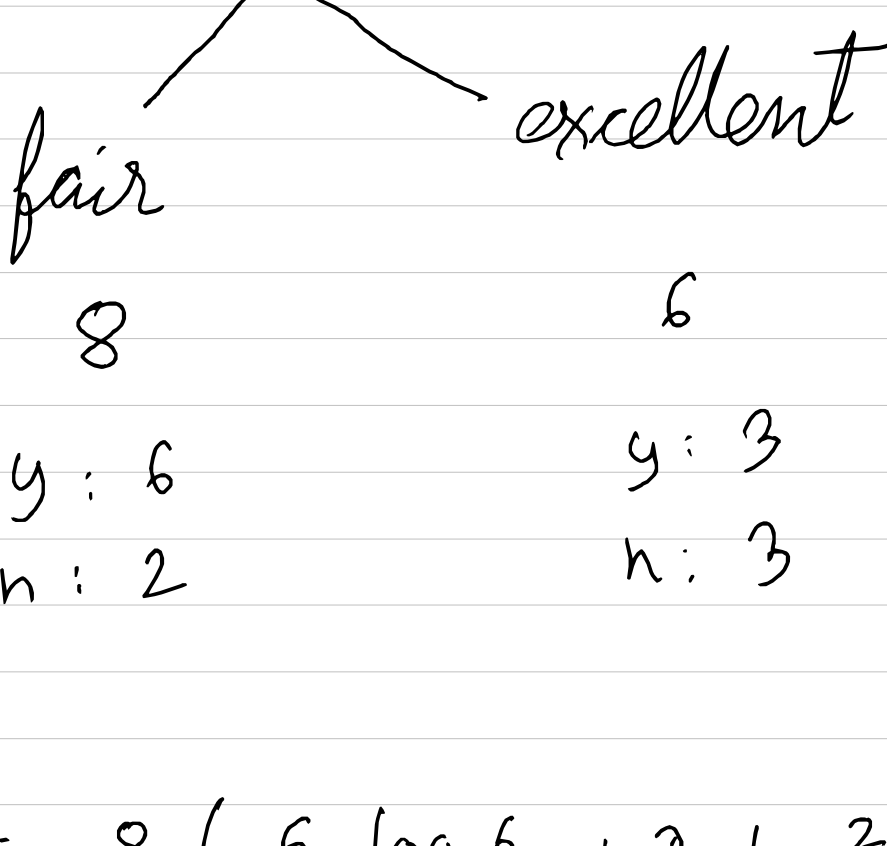
$$H = \frac{7}{14} \left(-\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7}\right)$$

$$+ \frac{7}{14} \left(-\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7}\right)$$

$$= 0.788$$

$$\text{Info gain} = 0.940 - 0.788 = 0.15$$

for credit score



$$H = \frac{8}{14} \left(\frac{6}{8} \log_2 \frac{6}{8} + \frac{2}{8} \log_2 \frac{2}{8}\right)$$

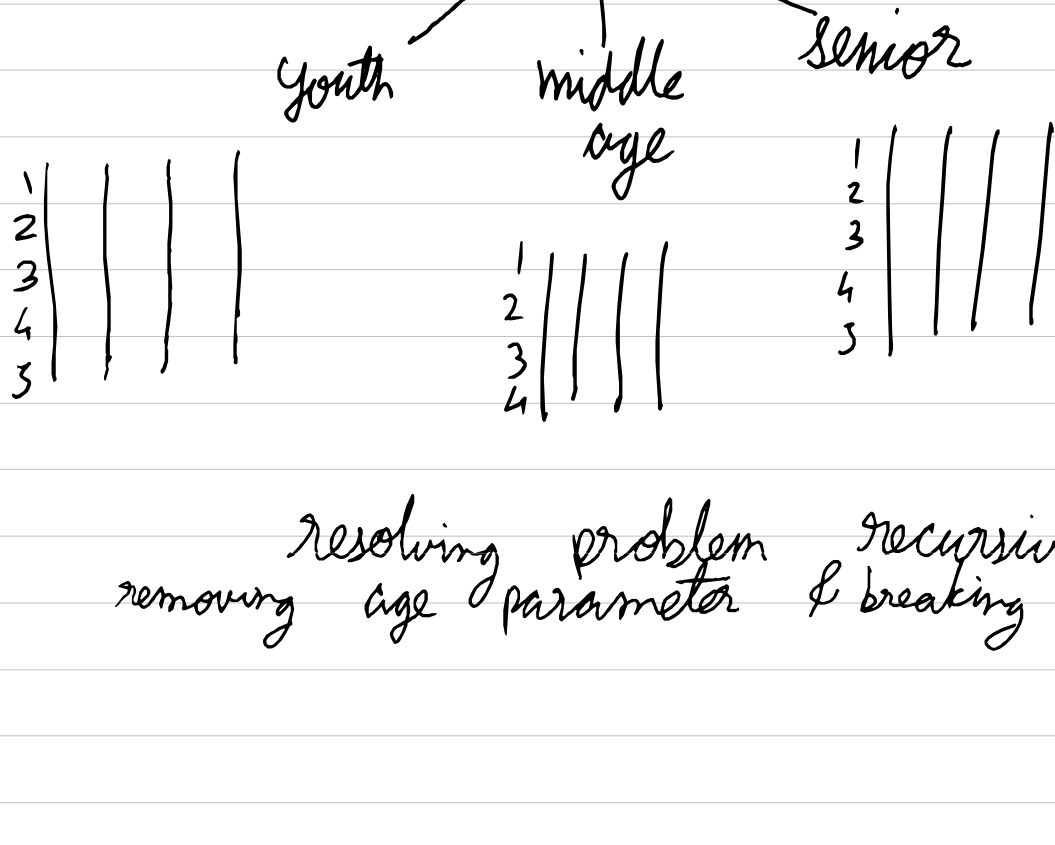
$$+ \frac{6}{14} \left(\frac{3}{6} \log_2 \frac{3}{6} + \frac{3}{6} \log_2 \frac{3}{6}\right)$$

$$= 0.892$$

$$\text{Info gain} = 0.94 - 0.892 = 0.048$$

|                  |         |         |
|------------------|---------|---------|
| Information Gain | age     | → 0.247 |
|                  | Income  | → 0.029 |
|                  | student | → 0.15  |
|                  | credit  | → 0.048 |

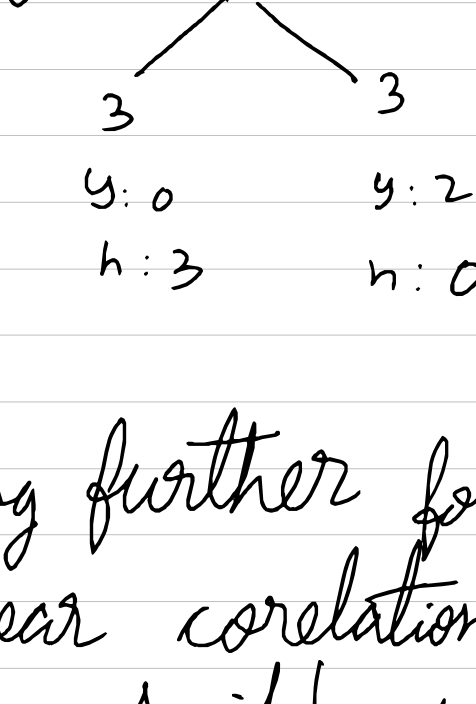
Highest value is by age. Hence age is root node



resolving problem recursively after removing age parameter & breaking into 3 parts

Value was 0 ∴ leaf

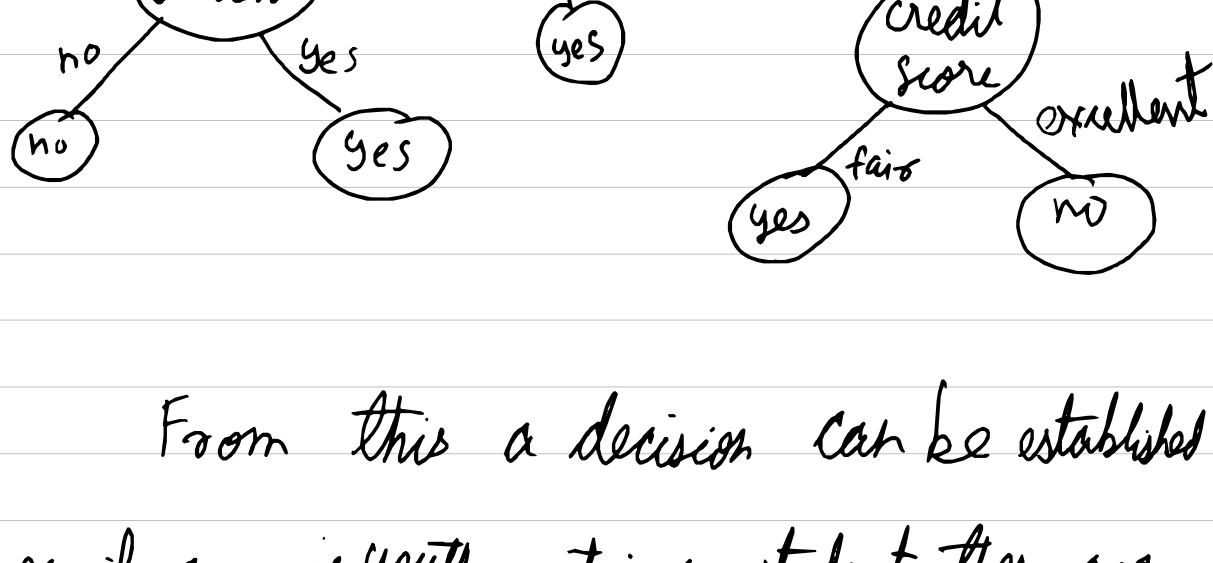
when there is 0 distribution for a node, leaf is found max info gain found



(for middle age, all values are yes ∴ terminal node found)

on solving further for senior, we can see clear correlation with it & the ans so it becomes leaf node

Final Answer →



From this a decision can be established eg if age is youth, it is a student then car can be bought

Q2

| Day | Outlook  | Temp | Humidity | Wind   | play Tennis? |
|-----|----------|------|----------|--------|--------------|
| 1   | Sunny    | Hot  | High     | Weak   | No           |
| 2   | S        | H    | H        | Strong | N            |
| 3   | Overcast | H    | H        | W      | Yes          |
| 4   | Rain     | Mild | H        | W      | Y            |
| 5   | R        | Cool | Normal   | W      | Y            |
| 6   | R        | C    | N        | S      | N            |
| 7   | O        | C    | N        | S      | Y            |
| 8   | S        | M    | H        | W      | N            |
| 9   | S        | C    | N        | W      | Y            |
| 10  | R        | M    | N        | W      | Y            |
| 11  | S        | M    | N        | S      | Y            |
| 12  | O        | M    | H        | S      | Y            |
| 13  | O        | H    | N        | W      | Y            |
| 14  | R        | M    | H        | S      | N            |

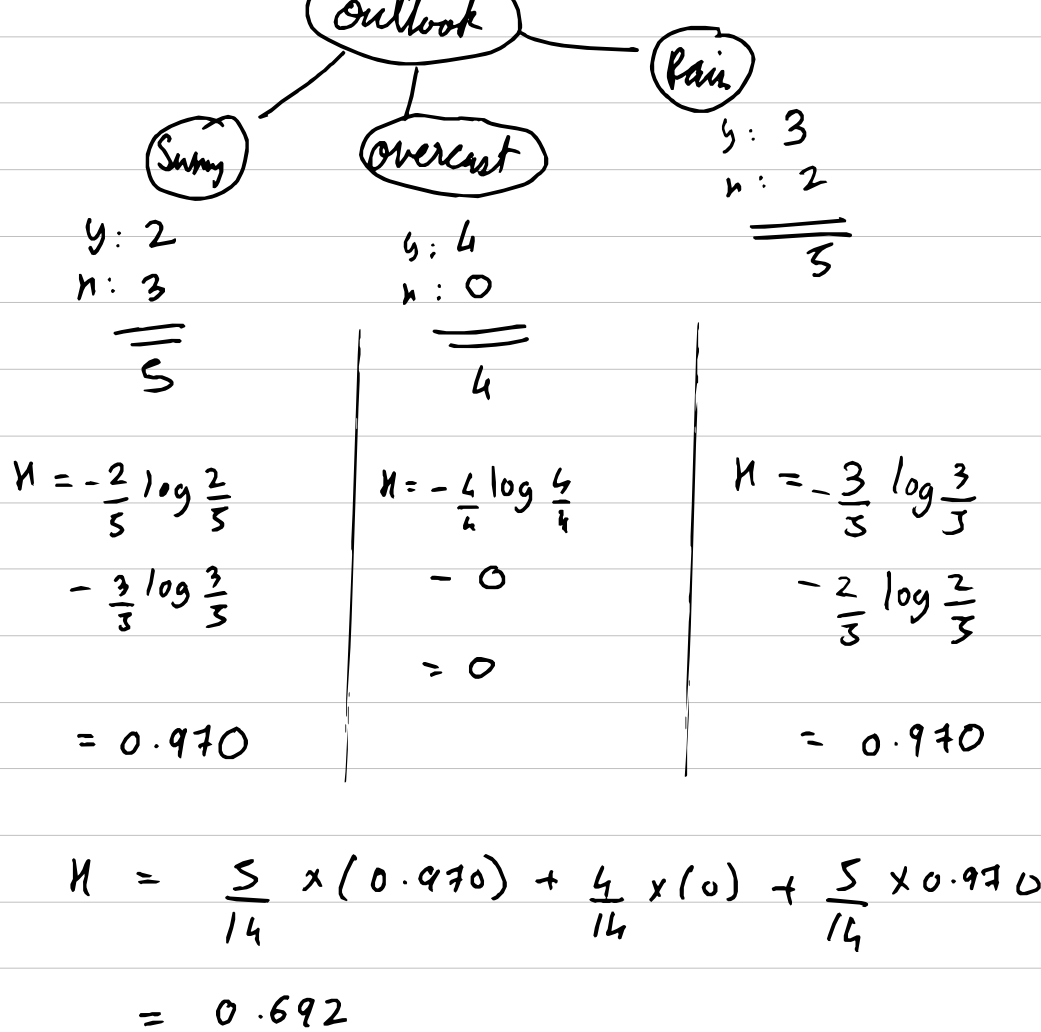
Step 1 → Entropy for parent → Tennis

Yes : 9  
No : 5

$$H = -\frac{9}{14} \log \frac{9}{14} - \frac{5}{14} \log \frac{5}{14}$$

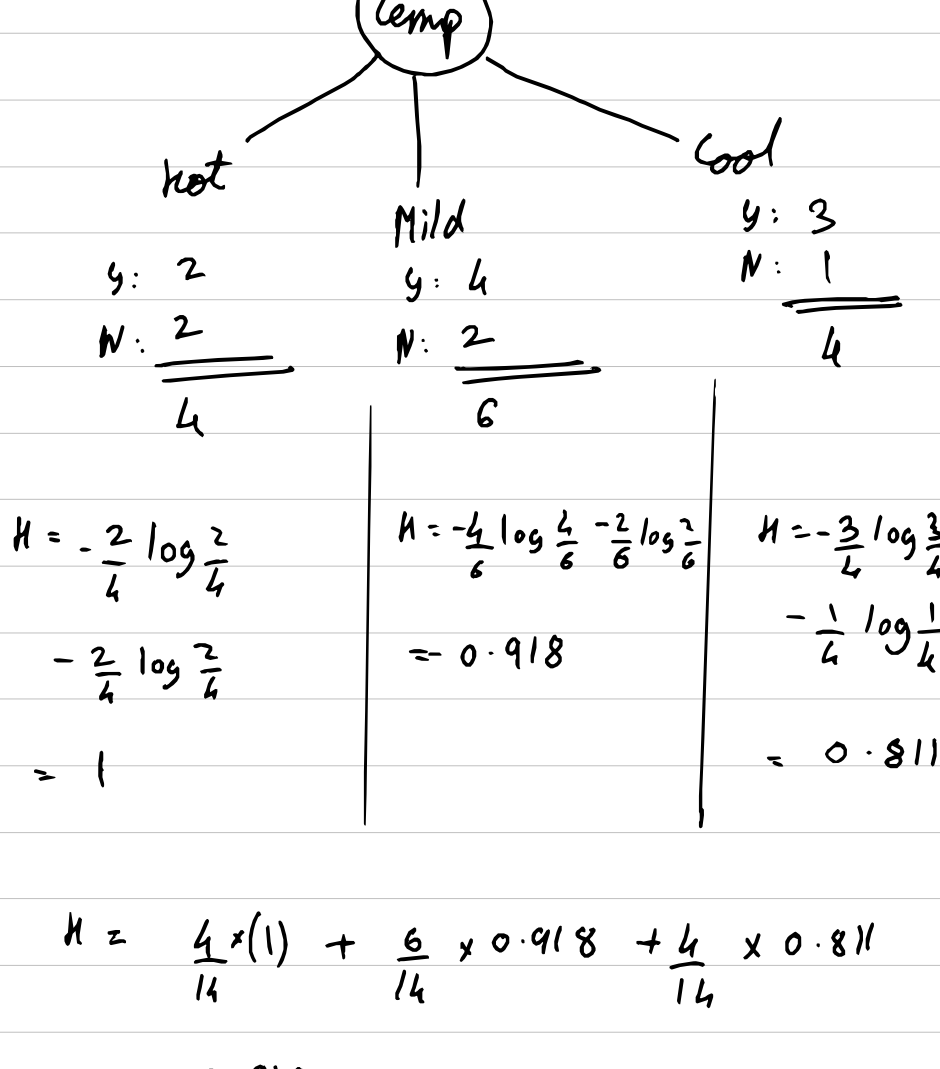
$$= 0.940$$

Step 2 → Entropy for attribute Outlook



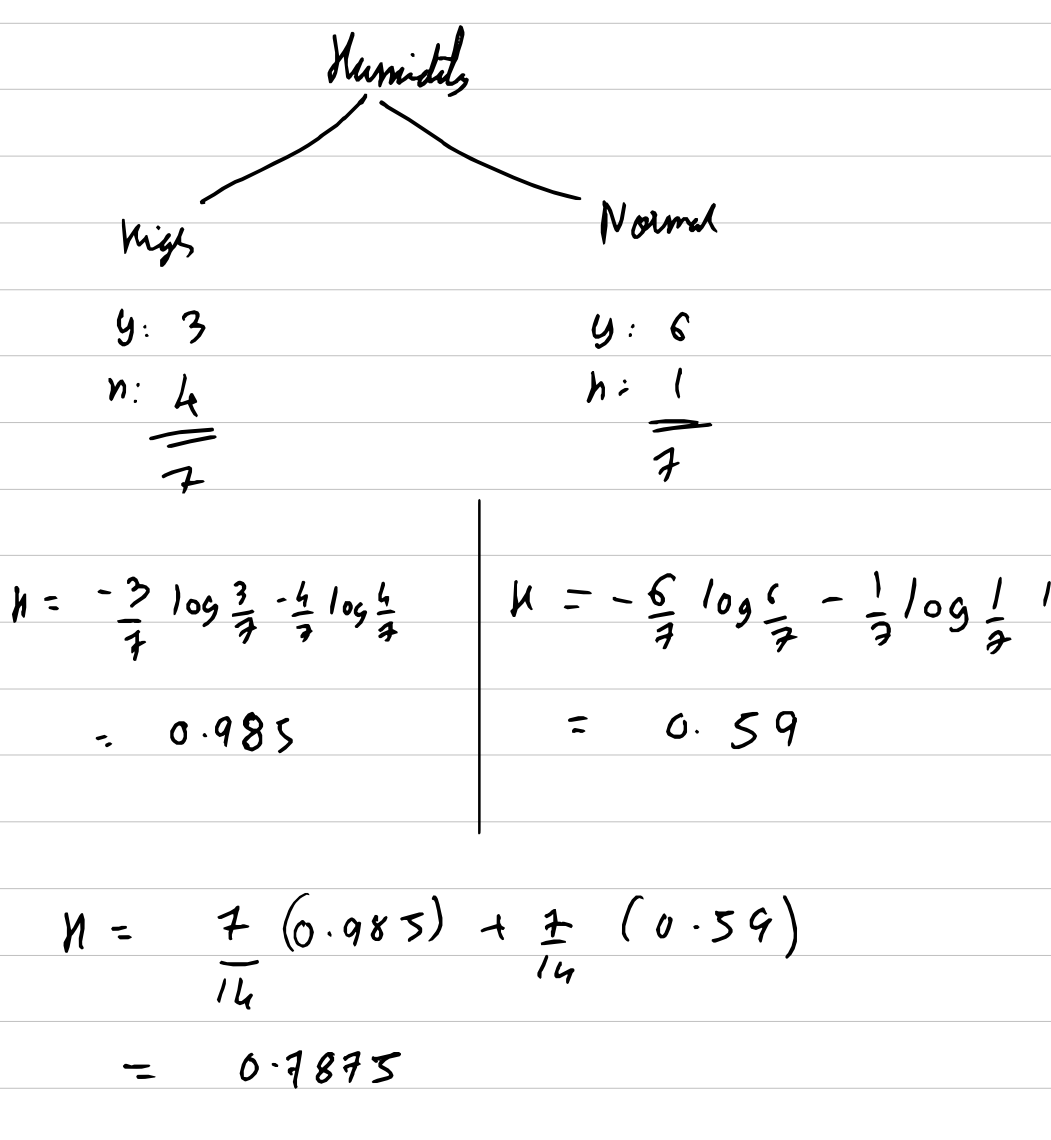
$$\text{Info gain} = 0.940 - 0.692 = 0.248$$

Attribute : Temp.



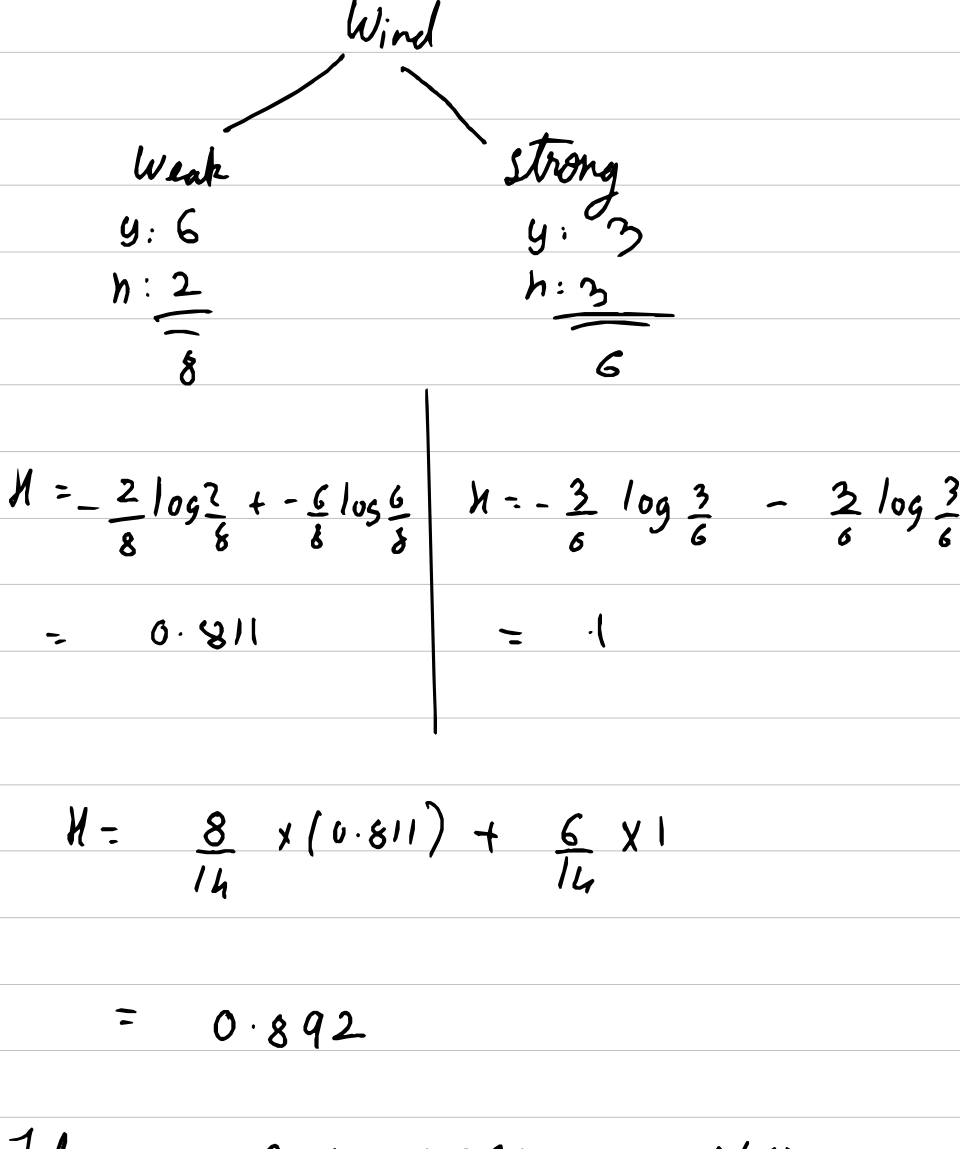
$$\text{Info gain} = 0.940 - 0.910 = 0.03$$

Attribute Humidity



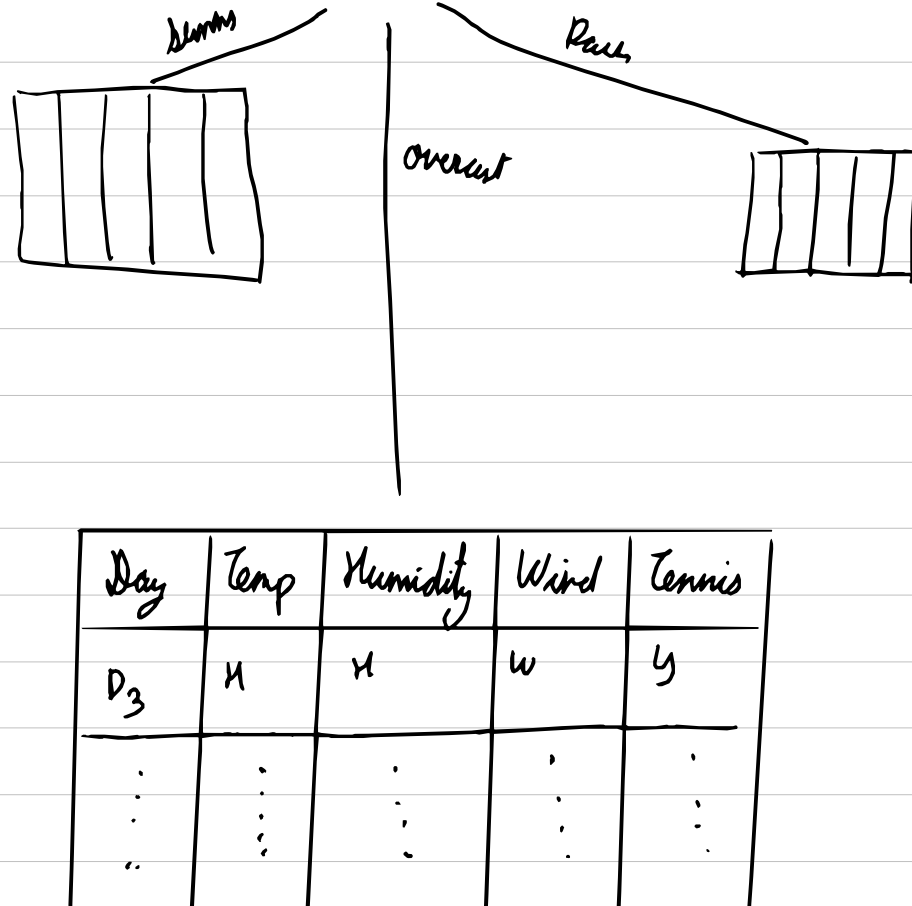
$$\text{Info gain} = 0.940 - 0.7875 = 0.152$$

Attribute : Wind



$$\text{Info gain} = 0.940 - 0.892 = 0.048$$

Best attribute for the split is with the Most Info gain that is outlook



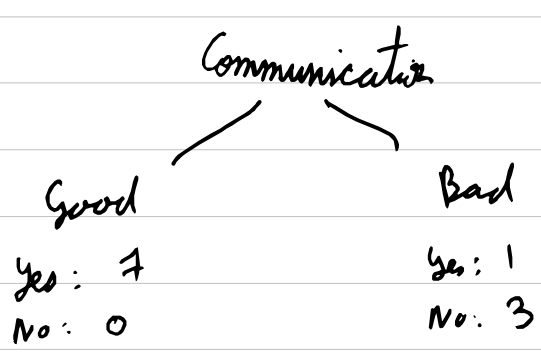
Q3

| CGPA   | Communication | Skills | Employable |
|--------|---------------|--------|------------|
| High   | Good          | Good   | Yes        |
| Medium | G             | G      | Y          |
| H      | G             | Bad    | Y          |
| H      | G             | G      | Y          |
| H      | Bad           | G      | Y          |
| M      | Good          | G      | Y          |
| L      | Bad           | B      | No         |
| L      | Bad           | B      | No         |
| M      | Good          | B      | Y          |
| M      | Bad           | G      | No         |
| M      | Good          | B      | Y          |

→ Parent

Yes: 8  
No: 3

$$H = -\frac{8}{11} \log \frac{8}{11} - \frac{3}{11} \log \frac{3}{11} = 0.845$$

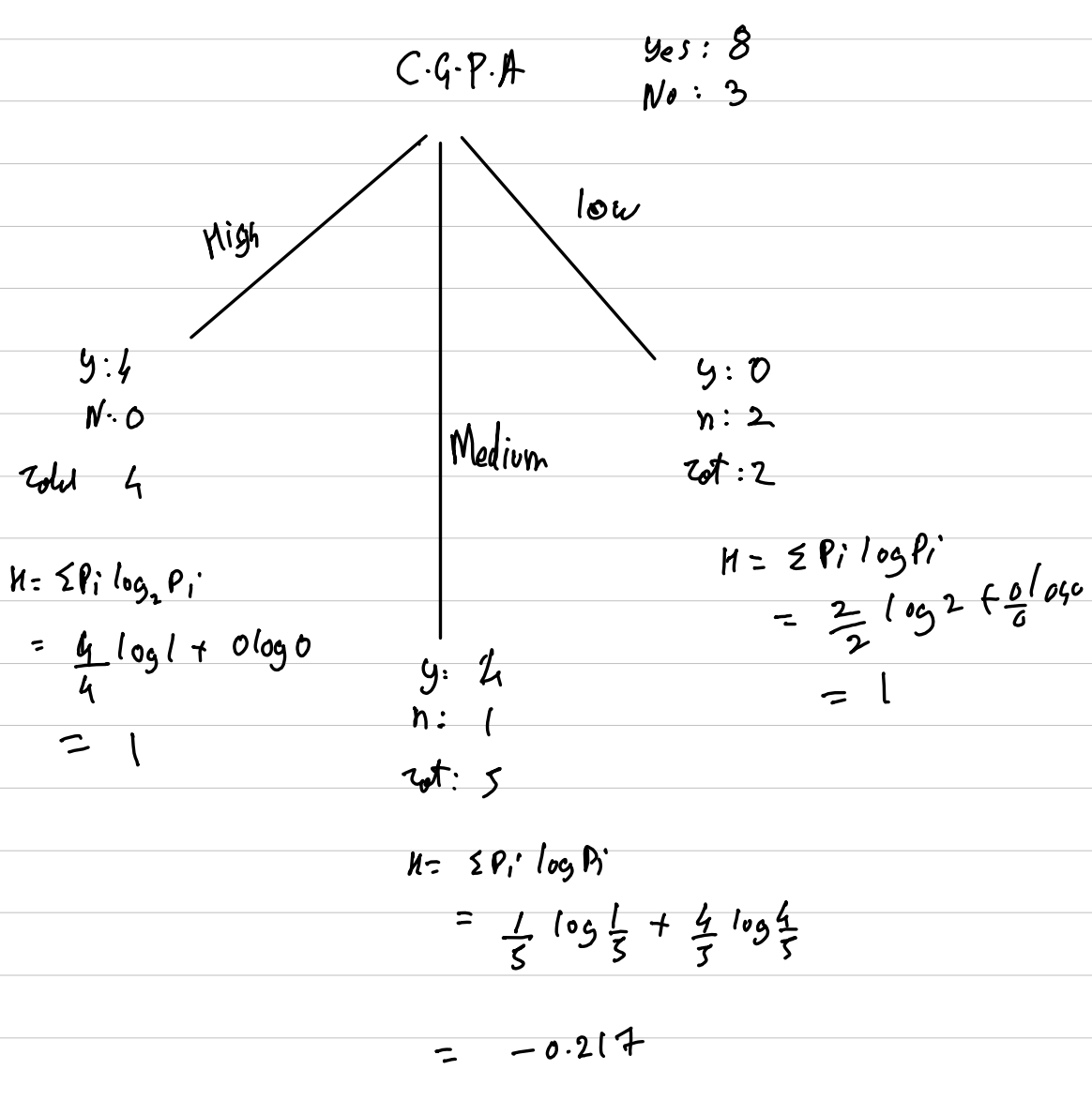


$$H = -\frac{7}{7} \log \left(\frac{7}{7}\right) - 0 = 0$$

$$H = -\frac{1}{4} \log \frac{1}{4} - \frac{3}{4} \log \frac{3}{4} = +0.8112$$

$$H = \frac{7}{11} \times 0 + \frac{4}{11} \times 0.8112 = 0.295$$

Info gain = 0.55



$$H = \sum P_i \log_2 P_i = \frac{4}{4} \log 1 + 0 \log 0 = 1$$

$$H = \sum P_i \log P_i = \frac{2}{2} \log 2 + \frac{0}{2} \log 0 = 1$$

$$H = \sum P_i \log P_i = \frac{1}{5} \log \frac{1}{5} + \frac{4}{5} \log \frac{4}{5} = -0.217$$

$$\Phi = \frac{4}{11} \times 1 + \frac{5}{11} \times -0.217 + \frac{2}{11} \times 1 = 0.4468$$

Info gain = 0.845 - 0.4468 = 0.398



# CART

## Classification & Regression Trees.

CART is a variation of the decision tree algorithm.

CART builds a tree like structure consisting of nodes and branches. Nodes represent decision points while branches represent the possible outcomes of the decisions.

Best split at every point is calculated using Gini impurity.

Gini requires binary split. Hence if the split is more than 2 then parts must be made.

## CART

### (Classification And Regression tree)

\* Based on Gini Index

$$* \text{Gini Index (D)} = 1 - \sum_{i=1}^m p_i^2$$

↑  
partition

$p_i$  is the probability that a tuple belongs to class  $C_i$

\* Gini Needs Binary split

eg Income } won't work with Gini  
low      High      Medium

Data must be divided into only two categories  
In order to do that, consider all possible subsets of the sets

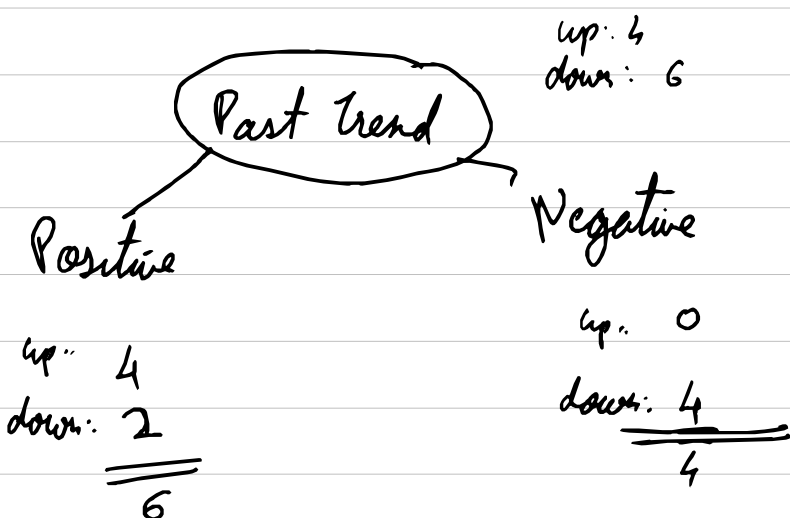
{ high, low } medium  $\rightarrow P_1$

{ high, medium } low  $\rightarrow P_2$

{ low, medium } high  $\rightarrow P_3$

3 combinations  ${}^3C_2$

Example  $\rightarrow$  calculate Gini value of the below split & info gain



$\rightarrow$  Gini Parent

$$G = 1 - p_+^2 - p_-^2$$

$$= 1 - \left(\frac{4}{10}\right)^2 - \left(\frac{6}{10}\right)^2 = 0.48$$

Gini for +ve

$$G = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.44$$

Gini for -ve

$$G = 1 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2 = 0$$

$$\text{Gini for split} = \frac{6}{10} \times (0.44) + 0 \times \frac{4}{10}$$

$$= 0.264$$

$$\text{Info gain} = 0.48 - 0.264 = 0.216$$

Q3) Make decision tree for gini Indexing on basis of Income

| id | age   | income | student | crediting | buy Car |
|----|-------|--------|---------|-----------|---------|
| 1  | Youth | high   | no      | fair      | no      |
| 2  | y     | h      | no      | excellent | no      |
| 3  | m-a   | h      | no      | f         | y       |
| 4  | s     | medium | no      | f         | y       |
| 5  | s     | low    | yes     | f         | y       |
| 6  | s     | L      | yes     | e         | n       |
| 7  | m-a   | L      | yes     | e         | y       |
| 8  | y     | m      | no      | f         | n       |
| 9  | y     | L      | y       | f         | y       |
| 10 | s     | m      | y       | f         | y       |
| 11 | y     | m      | y       | e         | y       |
| 12 | ma    | m      | n       | e         | y       |
| 13 | ma    | h      | y       | f         | y       |
| 14 | s     | m      | n       | e         | n       |

→

Parent      yes = 9  
                  no = 5

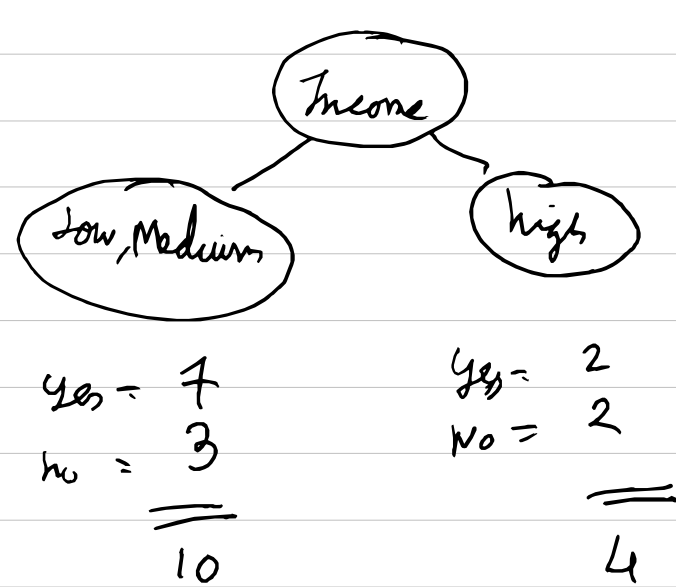
$$\begin{aligned} \text{Gini (Parent)} &= 1 - \sum p_i^2 = 1 - (p_{\text{yes}}^2 + p_{\text{no}}^2) \\ &= 1 - \left[ \left( \frac{9}{14} \right)^2 + \left( \frac{5}{14} \right)^2 \right] \\ &= 0.459 \end{aligned}$$

Possible splits → i) {L, m}, h  
                                 ii) {L, h}, m  
                                 iii) {m, h}, L

Consider the Binary split as {L, m}, h

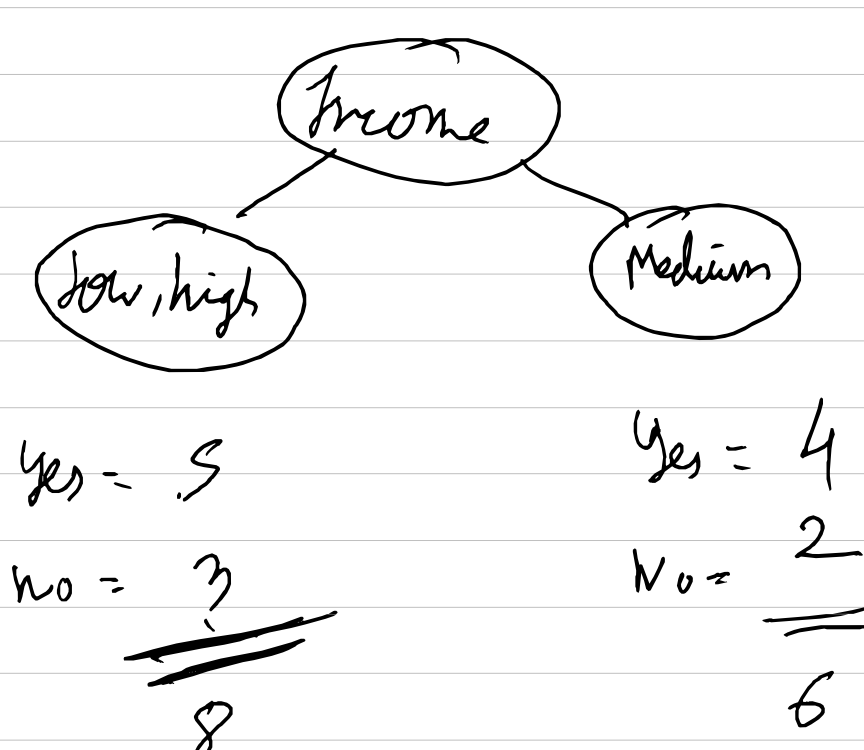
$D_1 = \{L, m\}$       } These are the partitions  
 $D_2 = h$

$$\text{Gini}_A(D) = \left( \frac{D_1}{D} \right) \text{Gini}(D_1) + \left( \frac{D_2}{D} \right) \text{Gini}(D_2)$$



$$\begin{aligned} \text{Gini}(D) &= \frac{10}{14} \left( 1 - \left( \frac{7}{10} \right)^2 - \left( \frac{3}{10} \right)^2 \right) + \frac{2}{4} \left( 1 - \left( \frac{2}{4} \right)^2 - \left( \frac{2}{4} \right)^2 \right) \\ &= 0.443 \end{aligned}$$

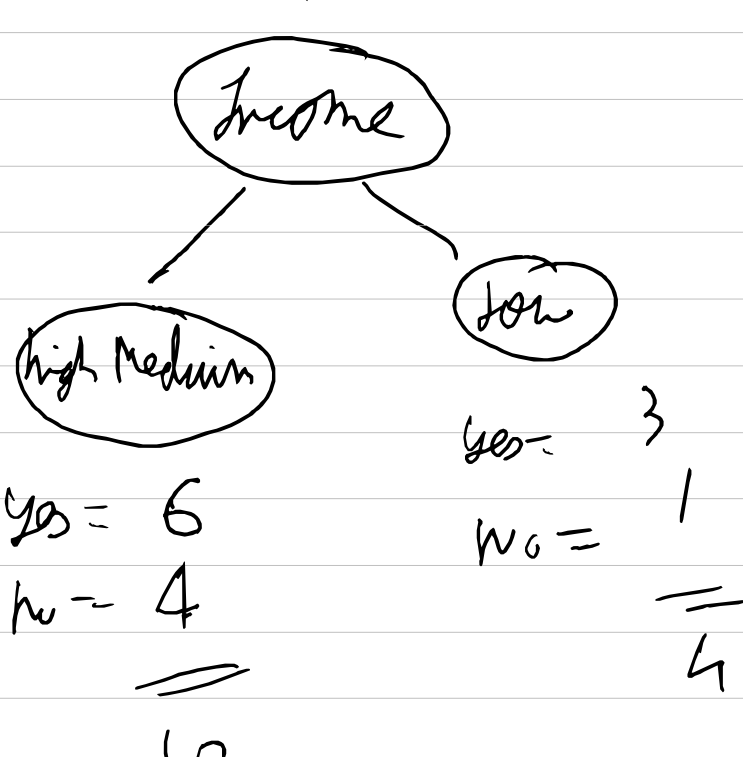
Consider split as {L, h}, m



$$\begin{aligned} \text{Gini} &= 1 - \left( \frac{5}{8} \right)^2 - \left( \frac{3}{8} \right)^2 & \text{Gini} &= 1 - \left( \frac{4}{6} \right)^2 - \left( \frac{2}{6} \right)^2 \\ &= 0.468 & &= 0.44 \end{aligned}$$

$$\begin{aligned} \text{Gini}(D) &= \frac{8}{14} \times (0.468) + \frac{6}{14} \times (0.44) \\ &= 0.457 \end{aligned}$$

Consider split as {h, m}, L



$$\begin{aligned} \text{Gini} &= 1 - \left( \frac{6}{10} \right)^2 - \left( \frac{4}{10} \right)^2 & \text{Gini} &= 1 - \left( \frac{3}{4} \right)^2 - \left( \frac{1}{4} \right)^2 \\ &= 0.48 & &= 0.375 \end{aligned}$$

$$\begin{aligned} \text{Gini} &= \frac{10}{14} (0.48) + \left( \frac{4}{14} \right) (0.375) \\ &= 0.43 \end{aligned}$$

Among all the three, the smallest gini index is of {L, m}, h hence consider that partition.

Q2

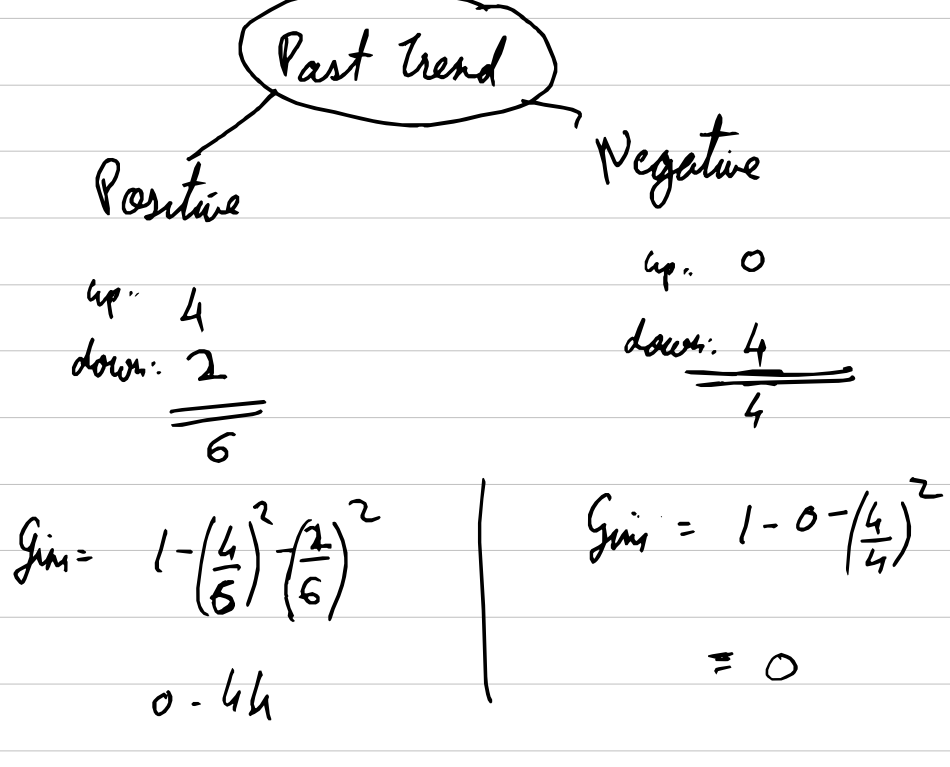
| Past trend | Open interest | Trading Volume | Return |
|------------|---------------|----------------|--------|
| +ve        | Low           | High           | Up     |
| -ve        | High          | Low            | Down   |
| +ve        | L             | H              | U      |
| +ve        | H             | H              | U      |
| -ve        | L             | H              | D      |
| +ve        | L             | L              | D      |
| -ve        | H             | H              | D      |
| -ve        | L             | H              | D      |
| +ve        | L             | L              | D      |
| +ve        | H             | H              | U      |

Parent

Total up: 4  
down: 6

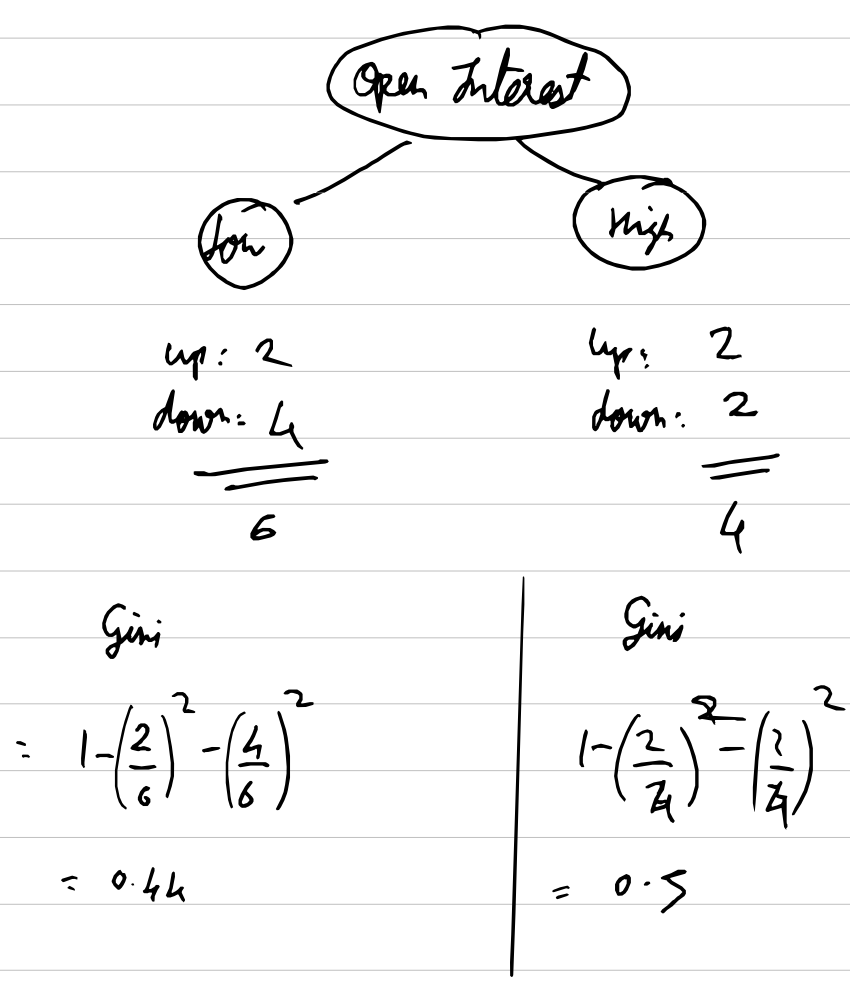
$$\text{Gini Parent} = 1 - \left(\frac{4}{10}\right)^2 - \left(\frac{6}{10}\right)^2 = 0.48$$

For Past trend



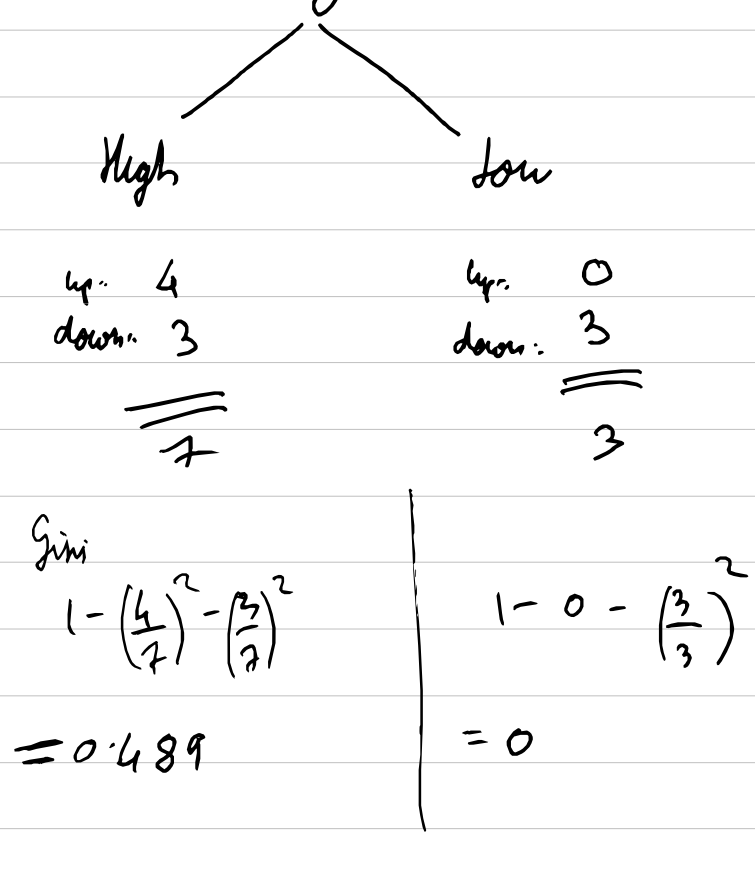
$$\text{Gini} = \frac{6}{10} \times 0.44 = 0.266$$

Gini Open Interest



$$\text{Gini} = \frac{6}{10} \times 0.44 + \frac{4}{10} \times 0.5 = 0.466$$

Gini for Trading volume



$$\text{Gini} = \frac{7}{10} \times 0.489 + \frac{3}{10} \times 0 = 0.3423$$

Smallest Gini value is of Past trend

## Advantages of decision trees

- 1) Understandable rules (Interpretable)
- 2) Easy calculation
- 3) Handle both continuous & categorical variables

## Disadvantages of decision trees

- 1) No global optimization
- 2) Error prone
- 3) Overfitting (especially deep trees)

## Methods to avoid overfitting

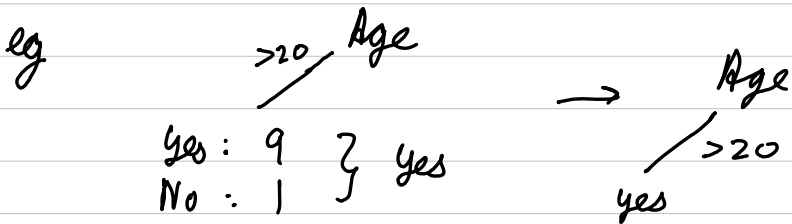
- ① Prepruning - stop growing when data split is not statistically significant
- ② Postpruning - Grow full tree & remove some nodes.

## Prepruning

### Early stopping

If some condition is met, the current node will not be split, even if it is not 100% pure

It will be a leaf node with the label of the majority class in current set



### Common stopping criterion

- ① Entropy / Gini impurity
- ② No of samples
- ③ Depth of tree

## Post Pruning

Prune nodes in bottom up manner if it decreases validation error

↓

unseen dataset

# XOR problem & linearity

XOR is a non linearly separable problem

| $x_1$ | $x_2$ | output |
|-------|-------|--------|
| 0     | 0     | 0      |
| 1     | 0     | 1      |
| 0     | 1     | 1      |
| 1     | 1     | 0      |

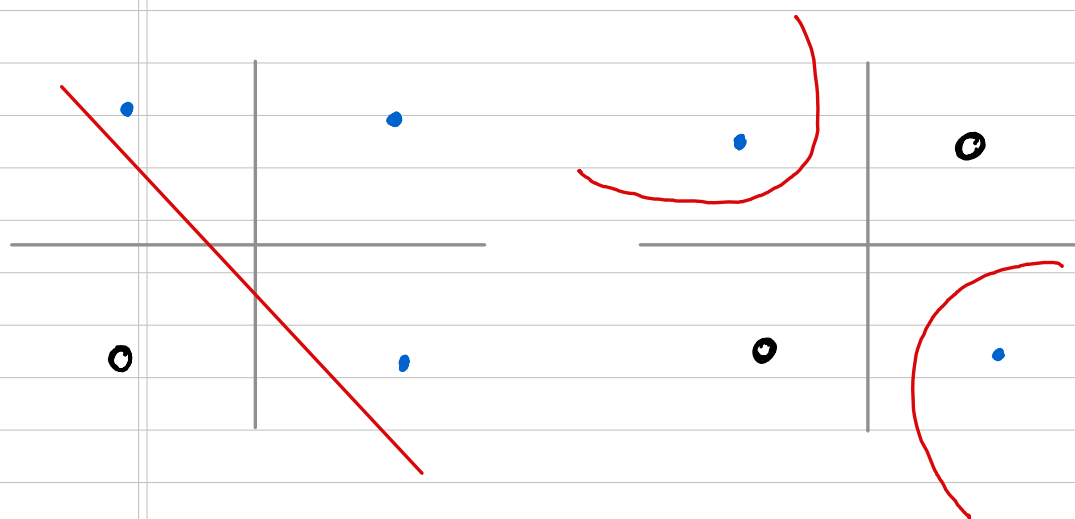
No straight line can separate 1 & 0

when there is such relationship in the data, it is a non linear data.

1) logistic regression is a linear classifier,

meaning it tries to find a single decision boundary in higher dimensions to separate the classes.

Since XOR requires at least two decision boundaries [for (0,0) and (1,1) from (0,1) and (1,0)], logistic regression cannot model it



• → output 1  
 ○ → output 0  
 — → classification line

2) Naive Bayes → Naive Bayes assumes conditional independence between the features and class labels

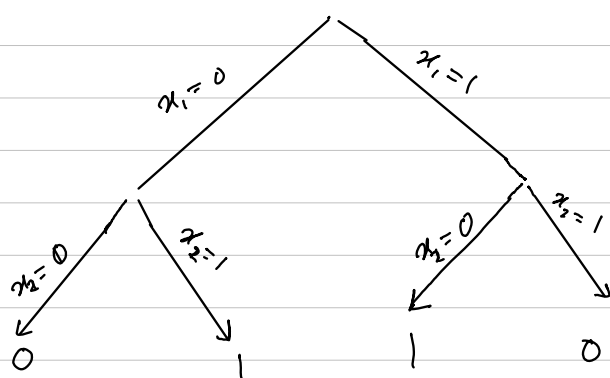
However in XOR, the two features are not independent, their interaction determines the output

| $x_1$ | $x_2$ | $y$ |                       |
|-------|-------|-----|-----------------------|
| 0     | 0     | 0   | $P(y=0 x_1=0) = 0.25$ |
| 1     | 0     | 1   | $P(y=0 x_1=1) = 0.25$ |
| 0     | 1     | 1   | $P(y=1 x_1=0) = 0.25$ |
| 0     | 0     | 0   | $P(y=1 x_1=1) = 0.25$ |
|       |       |     | $P(y=0 x_2=0) = 0.25$ |
|       |       |     | $P(y=0 x_2=1) = 0.25$ |
|       |       |     | $P(y=1 x_2=0) = 0.25$ |
|       |       |     | $P(y=1 x_2=1) = 0.25$ |

No decision can be made in this manner

3) Linear SVM → fails for the same reason as logistic regression

4) Decision trees → split feature space into regions non linearly



can represent the XOR and such non linear data well.

5) Random forest → since it is ensemble of D.T., it can also represent the data well

6) Kernelized SVM → works well as the non linear boundaries are supported

7) Neural networks → can do non linear boundaries