

## Decision Tree : >Supervised

	<u>Day</u>	<u>outlook</u>	<u>Temp</u>	<u>Humidity</u>	<u>wind</u>	<u>Play</u>
1	Sunny	Hot	High	Weak	No	①
2	Sunny	Hot	High	Strong	No.	②
3	Overcast	Hot	High	Weak	Yes.	③
4	Rain	Mild	-	-	Yes	④
5.	Rain	-	-	-	Yes	⑤
6.	Rain	-	-	-	No	⑥
7.	Overcast	-	-	-	-	⑦
8.	Sunny	-	-	-	-	⑧
9.	Sunny	-	-	-	-	⑨
10.	Rain	-	-	-	-	⑩
11.	Overcast	-	-	-	-	⑪

Step-1:

P = no. of yes values

n = no. of no. values.

$$P+n = \text{total}$$

$\checkmark P = 9$

$\checkmark n = 5$

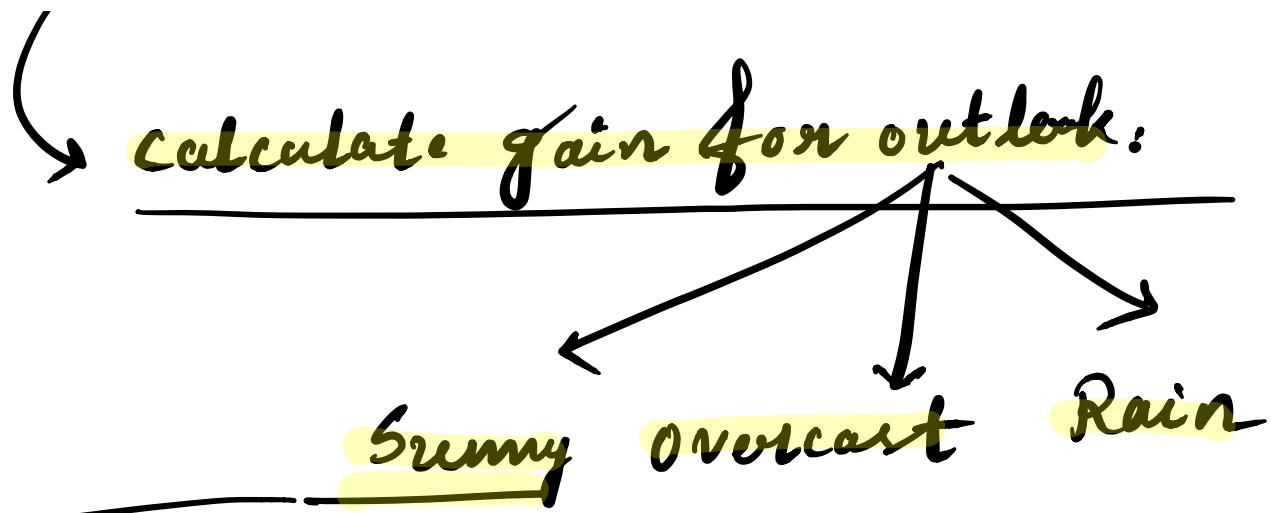
$P+n = 14$

Information Gain: [step-1]

$$I(P, n) = \frac{-P}{P+n} \log_2 \frac{P}{P+n} - \frac{n}{P+n} \log_2 \frac{n}{P+n}$$

$$\begin{aligned} I(P, n) &= \frac{-9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \\ &= 0.9403 \checkmark \end{aligned}$$

Calculate gain for each attribute  
[step-2]



$$\rightarrow P_i = 2 \checkmark$$

$$n_i = 3 \checkmark$$

$$\underline{I(P, n)} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= \underline{0.971}$$

Repeat the same for overcast and Rain.

outlook	$P_i$	$n_i$	$I(P, n)$
Sunny	2	3	0.971
Overcast	4	0	0
Rain	3	2	0.971

Calculate the entropy (Step-3)



$$E(A) = \sum_{i=1}^n \frac{P_i + n_i}{f+n} I(P_i, n_i)$$

$$\begin{aligned} E(\text{outlook}) &= \frac{5}{14} \times 0.971 + \frac{5}{14} \times 0.971 \\ &= 0.99 \end{aligned}$$

Step-4: Gain of outlook

$$\text{Gain}(\text{outlook}) = I(P, n) - E(\text{outlook})$$

$$= 0.9403 - 0.694$$

$$= 0.246$$

Repeat the same steps for  
Temp, Humidity and wind.

51

Step-5:

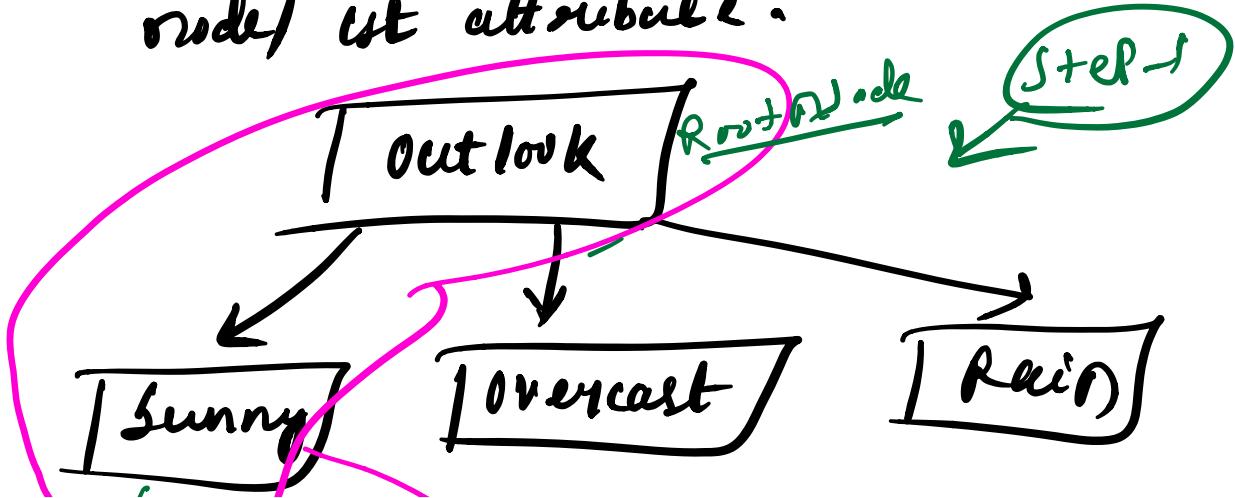
$$\text{Gain}(\text{Temp}) = 0.029$$

$$\text{Gain}(\text{Humidity}) = 0.151$$

$$\text{Gain}(\text{Wind}) = 0.048$$

$$\text{Gain}(\text{Outlook}) = 0.246 \rightarrow \begin{matrix} \text{highest} \\ \text{value} \\ \text{among} \\ \text{all} \end{matrix}$$

Hence choose outlook as the root node/ 1st attribute.



~~D~~

## Iteration-2:

Now it's turn to calculate the decision tree for this part

→ Create a subset from the main dataset containing all values for outlook and sunny.

Doz	outlook	Temp	H	w	P
8	Sunny	mid	High	w	No
9	Sunny	Gold	N	w	Yes
11	Sunny	Hot	N	S	Yes

Repeat all the above steps for this new table

# Random Forest :

why not only decision tree?



over fitting may be a prob.

Solution  $\Rightarrow$  Create a bunch of  
decision tree based on  
random sample of data.

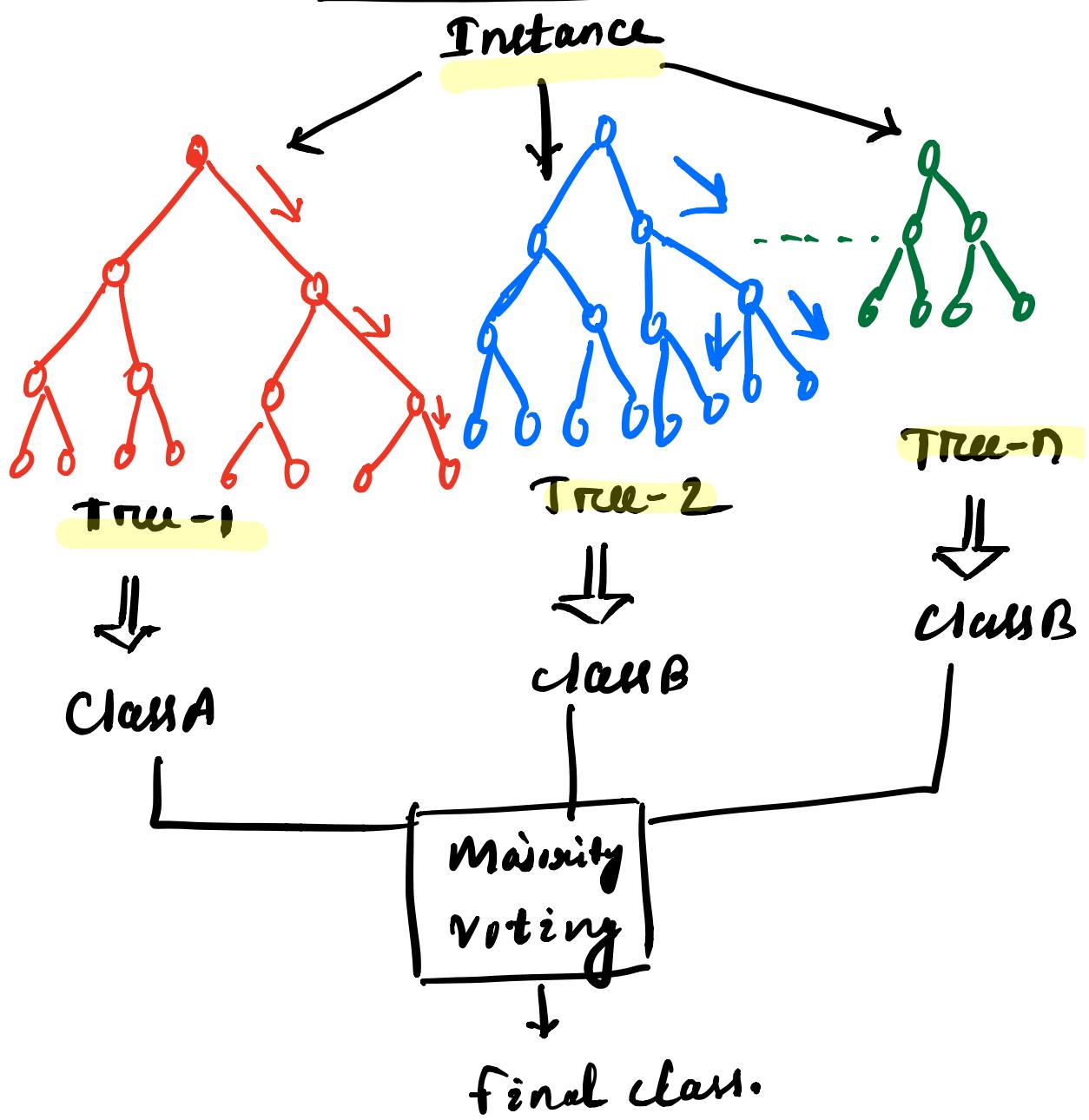
- $\rightarrow$  Define some set of subsample
- $\rightarrow$  for each subsample create a decision tree.
- $\rightarrow$  Then when we have a bunch of decision tree that we

generated, when we have a new data point we will use all these decision tree to predict its class and then we take majority of vote.

- Hence Random Forest gives us a higher accuracy than using a decision tree alone.
- Both Decision Tree and Random Forest can be used for both classification as well as for regression.

→ **most used machine learning technique.**

## Random Forest:



→ The more trees we add the more accuracy score we gets.

