

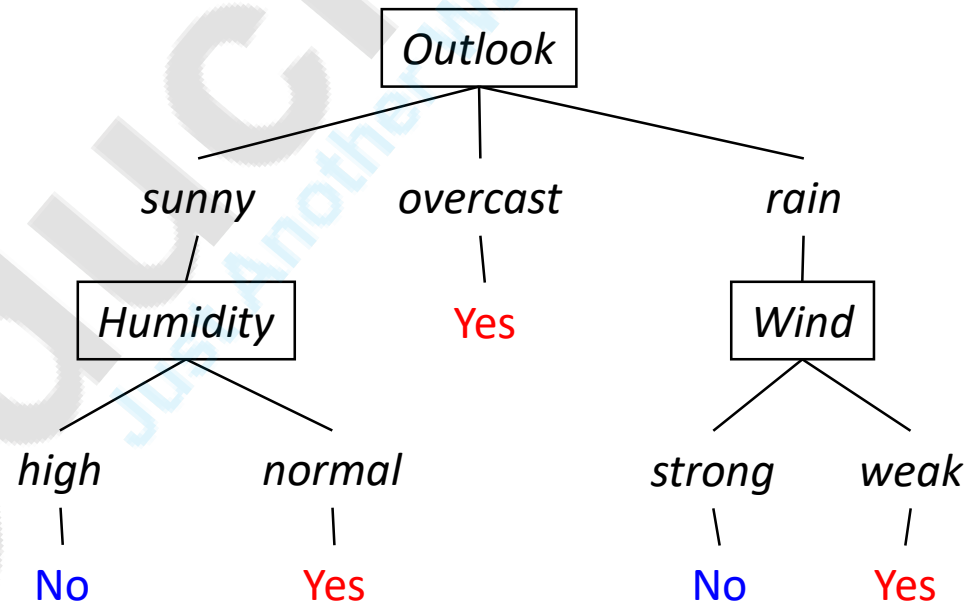
<b>Sr. No.</b>	<b>Module</b>	<b>Detailed Contents</b>	<b>Hrs</b>
<b>1</b>	<b>Business Intelligence-</b>	Introduction and overview of BI-Effective and timely decisions, Data Information and knowledge, BI Architecture, Ethics and BI. BI Applications- Balanced score card, Fraud detection, Telecommunication Industry, Banking and finance, Market segmentation.	<b>06</b>
<b>2</b>	<b>Prediction methods and models for BI</b>	Data preparation, Prediction methods-Mathematical method, Distance methods, Logic method, heuristic method-local optimization technique, stochastic hill climber, evaluation of models	<b>06</b>
<b>3</b>	<b>BI using Data Warehousing</b>	Introduction to DW, DW architecture, ETL Process, Top-down and bottom-up approaches, characteristics and benefits of data mart, Difference between OLAP and OLTP. Dimensional analysis- Define cubes. Drill- down and roll- up – slice and dice or rotation, OLAP models- ROLAP and MOLAP. Define Schemas- Star, snowflake and fact constellations.	<b>08</b>
<b>4</b>	<b>Data Mining and Preprocessing</b>	Data mining- definition and functionalities, KDD Process, Data Cleaning: - Missing values, Noisy data, data integration and transformations. Data Reduction: - Data cube aggregation, dimensionality reduction- data compression, Numerosity reduction- discretization and concept hierarchy.	<b>06</b>
<b>5</b>	<b>Associations and Correlation</b>	Association rule mining:-support and confidence and frequent item sets, market basket analysis, Apriori algorithm, Incremental ARM, Associative classification- Rule Mining.	<b>06</b>
<b>6</b>	<b>Classification and Prediction</b>	Introduction, Classification methods:-Decision Tree- ID3, CART, Bayesian classification- Baye'stheorem( Naïve Bayesian classification),Linear and nonlinear regression.	<b>08</b>
<b>7</b>	<b>Clustering</b>	Introduction, categorization of Major, Clustering Methods:- partitioning methods- K-Means. Hierarchical- Agglomerative and divisive methods, Model- based- Expectation and Maximization.	<b>08</b>
<b>8</b>	<b>Web mining and Text</b>	Text data analysis and Information retrieval, text retrieval methods, dimensionality reduction for text	<b>04</b>

# Decision trees

- Purpose is to expose the structural information contained in the data
- **Classification** trees
  - Response variable (class variable) is nominal or **categorical**
  - E.g. drug A or B
  - ID3 (induction decision tree version 3)
  - C4.5
- **Regression** trees
  - Response variable (class variable) is **continuous**
  - E.g. income
  - CART (classification and regression trees)

# Decision Trees

- Decision tree to represent learned target functions
  - Each internal node tests an attribute
  - Each branch corresponds to attribute value
  - Each leaf node assigns a classification
- Can be represented by logical formulas



# Representation in decision trees

- Example of representing rule in DT' s:

*if* outlook = sunny AND humidity = normal

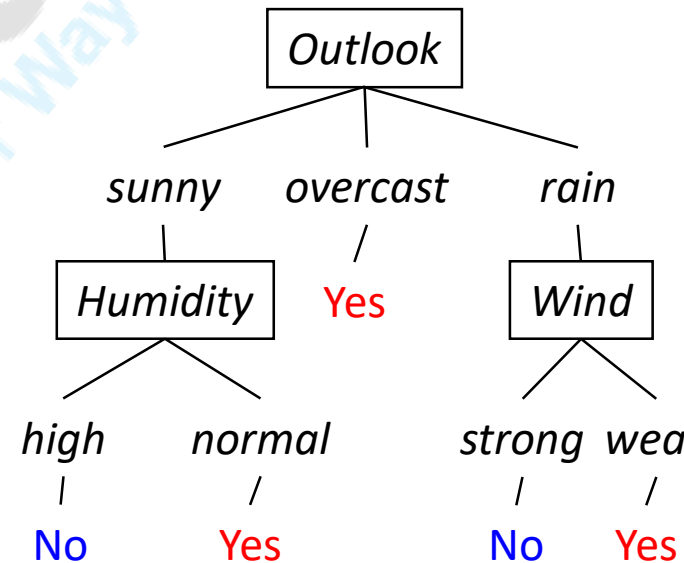
OR

*if* outlook = overcast

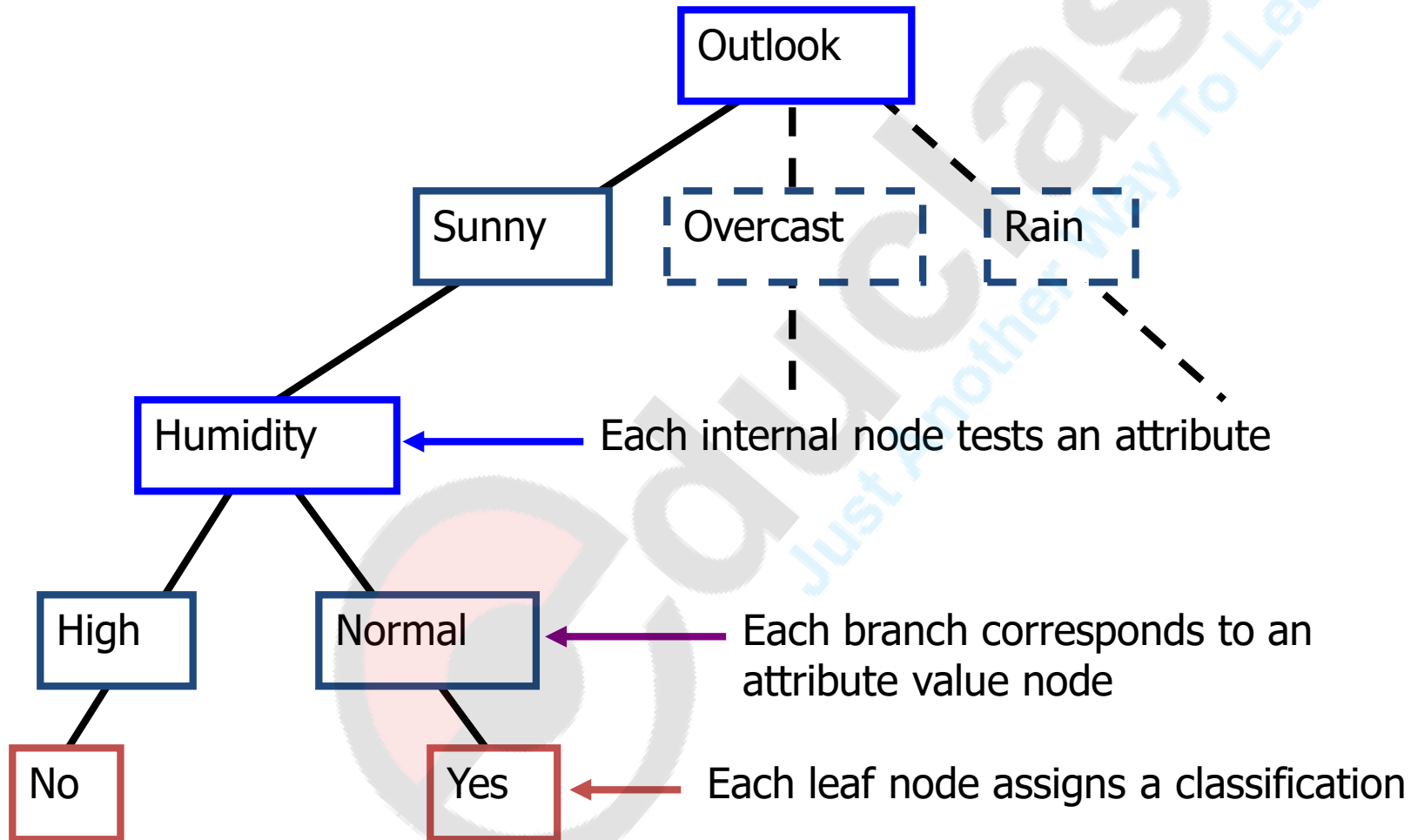
OR

*if* outlook = rain AND wind = weak

*then* playtennis

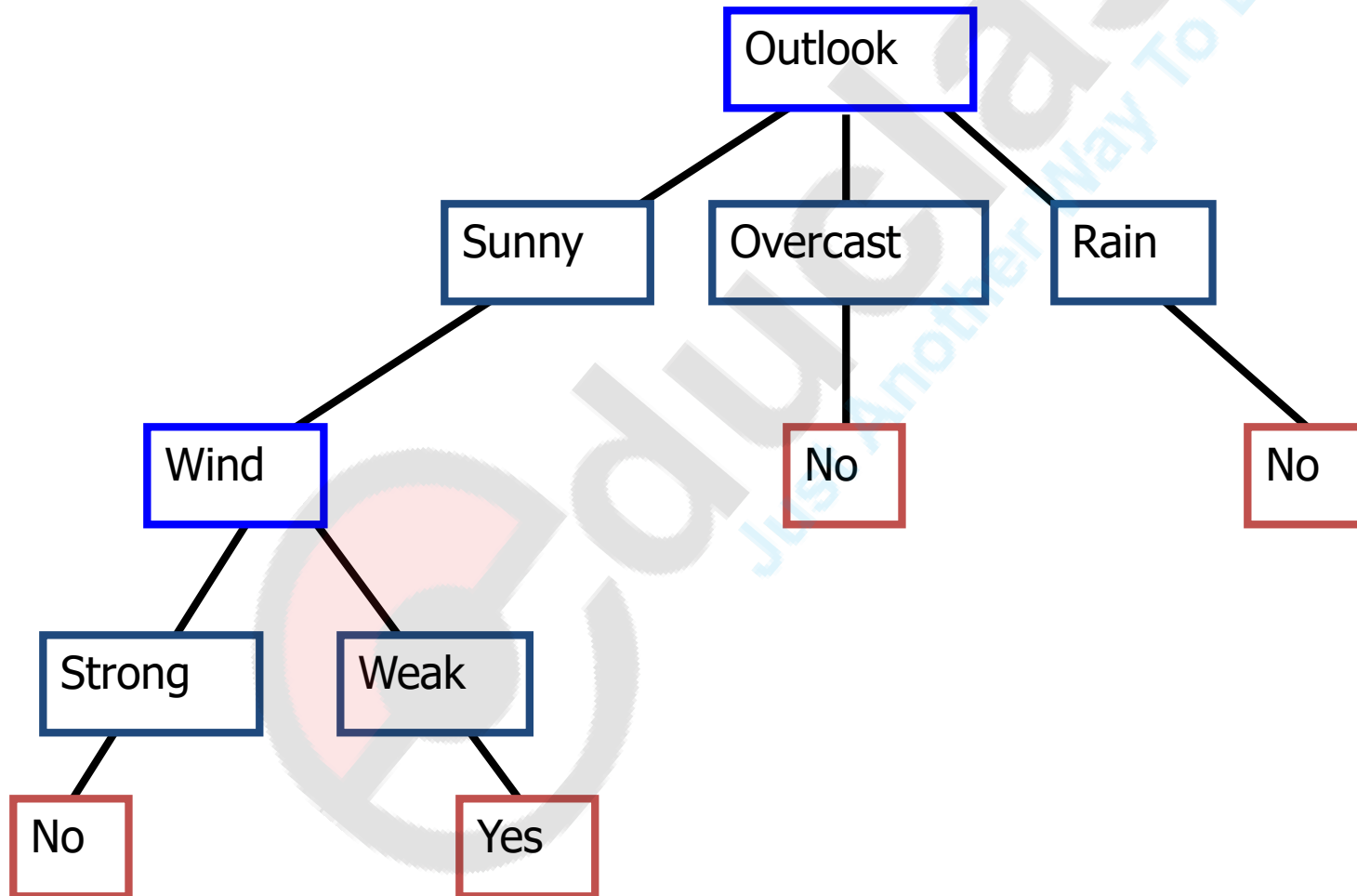


# Decision Tree for PlayTennis



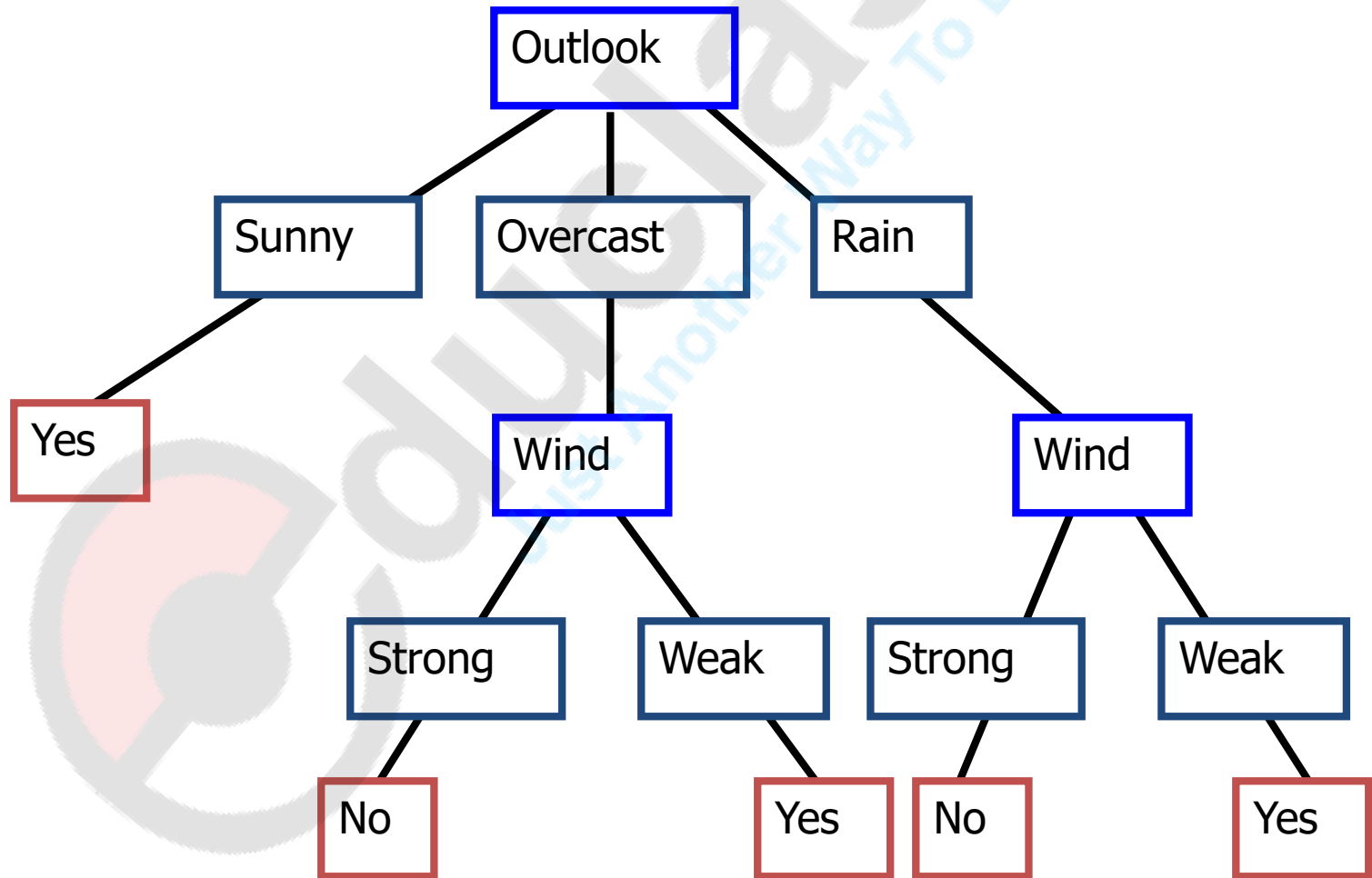
# Decision Tree

Outlook=Sunny  $\wedge$  Wind=Weak



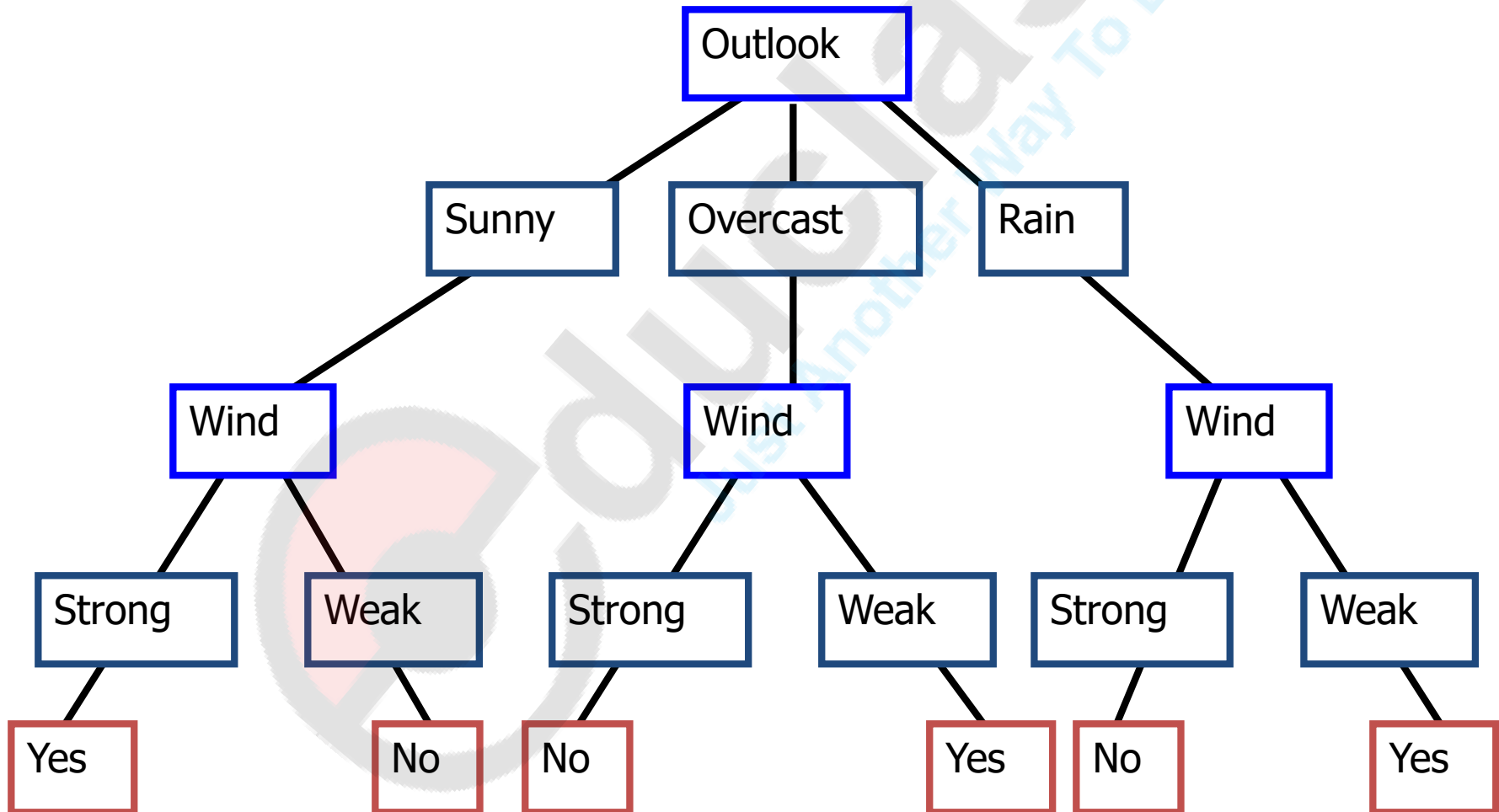
# Decision Tree

Outlook=Sunny  $\vee$  Wind=Weak



# Decision Tree for XOR

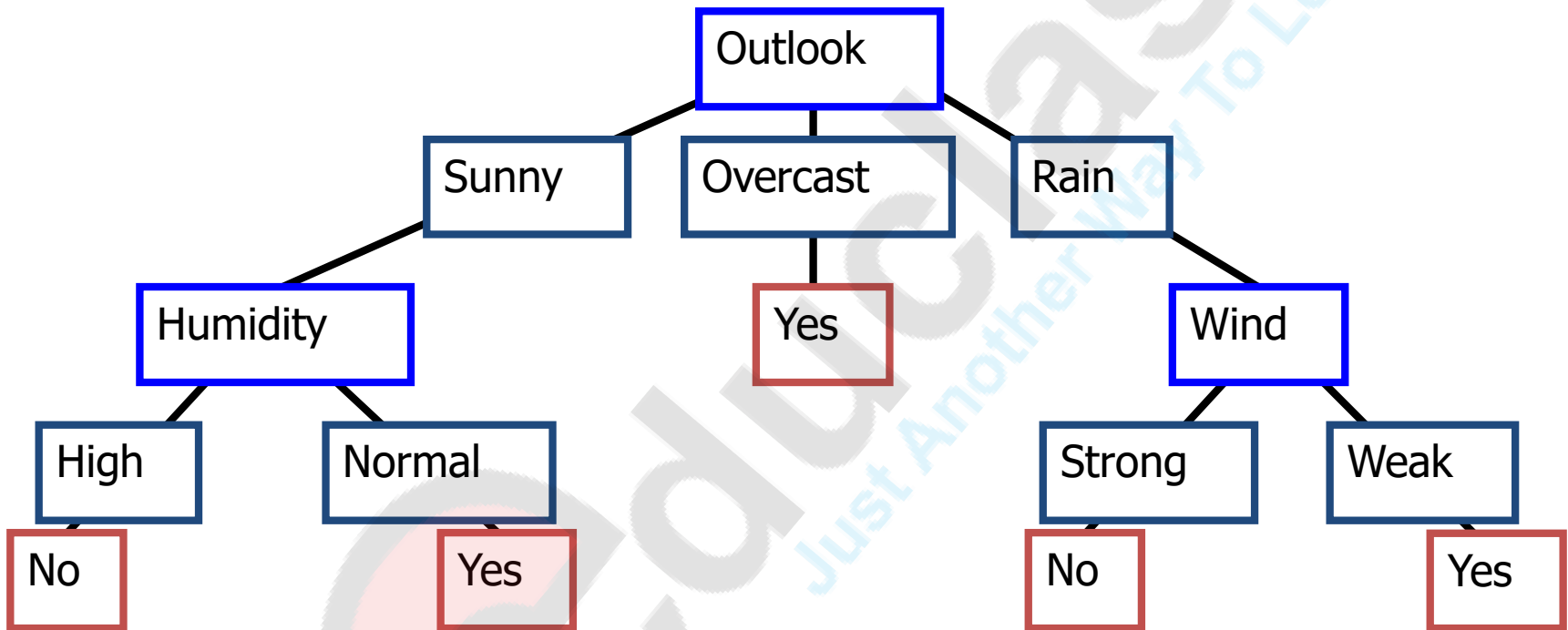
Outlook=Sunny XOR Wind=Weak





# Decision Tree

- decision trees represent disjunctions of conjunctions



(Outlook=Sunny  $\wedge$  Humidity=Normal)

∨ (Outlook=Overcast)

∨ (Outlook=Rain  $\wedge$  Wind=Weak)

# ID3

- Recursive procedure to construct a decision tree from data.
- ID3 algorithm assumes that a good decision tree is the simplest decision tree
  - We should always accept the simplest answer that correctly fits our data
  - The smallest decision tree that correctly classifies all given examples
- The simplest decision tree that covers all examples should be the least likely to include unnecessary constraints

# Important values

- **Information gain**

$$I(p,n) = \frac{-p}{p+n} \log_2 \left[ \frac{p}{p+n} \right] - \frac{n}{p+n} \log_2 \left[ \frac{n}{p+n} \right]$$

- **Entropy**  $E(A) = \sum \frac{p_i + n_i}{p+n} [ I(p,n) ]$   
of the attribute

- **Gain**(A) =  $I(p,n) - E(A)$

p= yes=9

n=no =5

Day	Outlook	Temp	Humidity	Wind	Tennis?
<i>D1</i>	Sunny	Hot	High	Weak	No
<i>D2</i>	Sunny	Hot	High	Strong	No
<i>D3</i>	Overcast	Hot	High	Weak	Yes
<i>D4</i>	Rain	Mild	High	Weak	Yes
<i>D5</i>	Rain	Cool	Normal	Weak	Yes
<i>D6</i>	Rain	Cool	Normal	Strong	No
<i>D7</i>	Overcast	Cool	Normal	Strong	Yes
<i>D8</i>	Sunny	Mild	High	Weak	No
<i>D9</i>	Sunny	Cool	Normal	Weak	Yes
<i>D10</i>	Rain	Mild	Normal	Weak	Yes
<i>D11</i>	Sunny	Mild	Normal	Strong	Yes
<i>D12</i>	Overcast	Mild	High	Strong	Yes
<i>D13</i>	Overcast	Hot	Normal	Weak	Yes
<i>D14</i>	Rain	Mild	High	Strong	No

p= yes=9

n=no =5

- Information gain (IG) of the table

$$I(p,n) = \frac{-p}{p+n} \log_2 \left[ \frac{p}{p+n} \right] - \frac{n}{p+n} \log_2 \left[ \frac{n}{p+n} \right]$$

$$I(9,5) = \frac{-9}{14} \log_2 \left[ \frac{9}{14} \right] - \frac{5}{14} \log_2 \left[ \frac{5}{14} \right]$$

$$I(9,5) = \frac{-9}{14} \log_2 \left[ 0.642 \right] - \frac{5}{14} \log_2 \left[ 0.357 \right]$$

$$I(9,5) = \frac{-9}{14} \frac{\log (0.642)}{\log 2} - \frac{5}{14} \frac{\log (0.357)}{\log 2}$$

$$= [ -9(-0.639) ] / 14 - [ 5 (-1.485) ] / 14 = 0.941$$

- $P = \text{yes} = 9$
- $N = \text{no} = 5$

Day	Outlook	Temp	Humidity	Wind	Tennis?
<i>D1</i>	Sunny	Hot	High	Weak	No
<i>D2</i>	Sunny	Hot	High	Strong	No
<i>D3</i>	Overcast	Hot	High	Weak	Yes
<i>D4</i>	Rain	Mild	High	Weak	Yes
<i>D5</i>	Rain	Cool	Normal	Weak	Yes
<i>D6</i>	Rain	Cool	Normal	Strong	No
<i>D7</i>	Overcast	Cool	Normal	Strong	Yes
<i>D8</i>	Sunny	Mild	High	Weak	No
<i>D9</i>	Sunny	Cool	Normal	Weak	Yes
<i>D10</i>	Rain	Mild	Normal	Weak	Yes
<i>D11</i>	Sunny	Mild	Normal	Strong	Yes
<i>D12</i>	Overcast	Mild	High	Strong	Yes
<i>D13</i>	Overcast	Hot	Normal	Weak	Yes
<i>D14</i>	Rain	Mild	High	Strong	No

$$I(p,n) = \frac{-p}{p+n} \log_2 \left[ \frac{p}{p+n} \right] - \frac{n}{p+n} \log_2 \left[ \frac{n}{p+n} \right]$$

**p** = yes = 9  
**n** = no = 5

• **Entropy**  $E(A) = \sum \frac{p_i + n_i}{p+n} [ I(p_i, n_i) ]$   
of the attribute

Outlook	Tennis?
Sunny	No
Sunny	No
Overcast	Yes
Rain	Yes
Rain	Yes
Rain	No
Overcast	Yes
Sunny	No
Sunny	Yes
Rain	Yes
Sunny	Yes
Overcast	Yes
Overcast	Yes
Rain	No

Outlook	p	n	I(p,n)
Sunny	2	3	0.970
Overcast	4	0	0
Rain	3	2	0.970

$I(2,3) = 0.970$   
 $I(4,0) = 0$ , if p or n is zero  
 $I(3,2) = 0.970$ , same for 3,2 or 2,3

$$E(\text{outlook}) = \left[ \frac{(2+3)}{(9+5)} \right] (0.970) + 0 + \left[ \frac{(3+2)}{14} \right] (0.970) = 0.692$$

$$\begin{aligned} \text{Gain (outlook)} &= IG - E(\text{outlook}) \\ &= 0.940 - 0.692 \\ &= 0.248 \end{aligned}$$

$E(\text{humidity}) = ?$  Gain (humidity) = ?  
 $E(\text{wind}) = ?$  Gain(wind) = ?  
 $E(\text{temp}) = ?$  Gain (temp) = ?

$$I(p,n) = \frac{-p}{p+n} \log_2 \left[ \frac{p}{p+n} \right] - \frac{n}{p+n} \log_2 \left[ \frac{n}{p+n} \right]$$

p = yes = 9  
n = no = 5

• **Entropy**

$$E(A) = \sum \frac{p_i + n_i}{p+n} [ I(p_i, n_i) ]$$

of the attribute

p+n

High → p=3 n=4 I(3,4)= ?  
Normal → p=6 n=1 I(6,1)= ?

E(humidity) = ?

Humidity	Tennis ?
High	No
High	No
High	Yes
High	Yes
Normal	Yes
Normal	No
Normal	Yes
High	No
Normal	Yes
Normal	Yes
Normal	Yes
High	Yes
Normal	Yes
High	No

Gain (humidity) = IG - E(humidity)  
= 0.151

gain(wind) = 0.045  
gain(temp) = 0.029

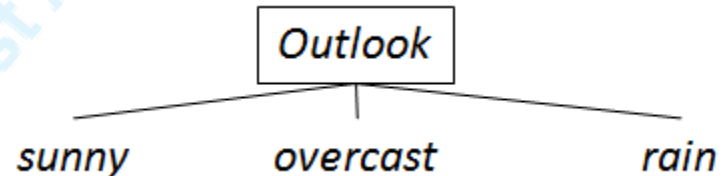


$$\begin{aligned}\text{Gain (outlook)} &= \text{IG} - \mathbf{E}(\text{outlook}) \\ &= 0.940 - 0.692 \\ &= 0.248\end{aligned}$$

Maximum gain is of outlook therefore tree starts from outlook ie. root

$$\begin{aligned}\text{Gain (humidity)} &= \text{IG} - \mathbf{E}(\text{humidity}) \\ &= 0.151\end{aligned}$$

$$\begin{aligned}\mathbf{gain}(\text{wind}) &= 0.045 \\ \mathbf{gain}(\text{temp}) &= 0.029\end{aligned}$$



T

p= yes=9

n=no =5

Day	Outlook	Temp	Humidity	Wind	Tennis?
<i>D1</i>	Sunny	Hot	High	Weak	No
<i>D2</i>	Sunny	Hot	High	Strong	No
<i>D3</i>	Overcast	Hot	High	Weak	Yes
<i>D4</i>	Rain	Mild	High	Weak	Yes
<i>D5</i>	Rain	Cool	Normal	Weak	Yes
<i>D6</i>	Rain	Cool	Normal	Strong	No
<i>D7</i>	Overcast	Cool	Normal	Strong	Yes
<i>D8</i>	Sunny	Mild	High	Weak	No
<i>D9</i>	Sunny	Cool	Normal	Weak	Yes
<i>D10</i>	Rain	Mild	Normal	Weak	Yes
<i>D11</i>	Sunny	Mild	Normal	Strong	Yes
<i>D12</i>	Overcast	Mild	High	Strong	Yes
<i>D13</i>	Overcast	Hot	Normal	Weak	Yes
<i>D14</i>	Rain	Mild	High	Strong	No

T  
sunny

p= yes=2

n=no =3

Day	Outlook	Temp	Humidity	Wind	Tennis?
<i>D1</i>	Sunny	Hot	High	Weak	No
<i>D2</i>	Sunny	Hot	High	Strong	No
<i>D8</i>	Sunny	Mild	High	Weak	No
<i>D9</i>	Sunny	Cool	Normal	Weak	Yes
<i>D11</i>	Sunny	Mild	Normal	Strong	Yes

# Entropy

$$E(A) = \sum \frac{p_i + n_i}{p+n} [ I(p_i, n_i) ]$$

TEMP	p	n	I(p,n)
Hot	0	2	
Mild	1	1	
cold	1	0	

humidity	p	n	I(p,n)
High	0	3	
normal	2	0	

wind	p	n	I(p,n)
Weak	1	2	
strong	1	1	

E(temp)=

E(humidity)=

E(wind)=

Gain(temp) =  $IG_{\text{sunny}} - E(\text{temp}) = 0.571$

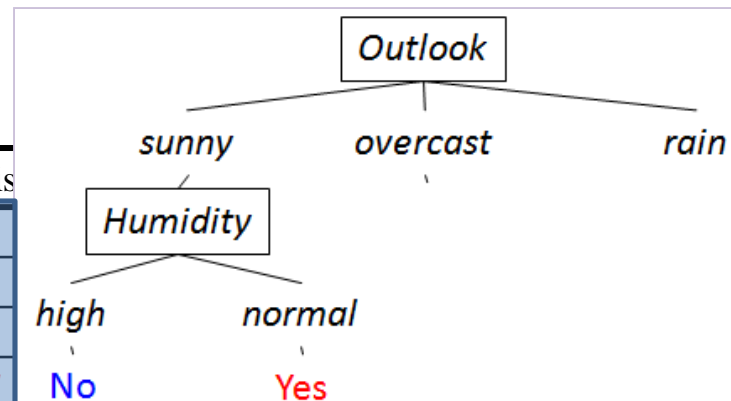
Gain(humidity) =  $0.971$

Gain(wind) =  $0.020$

We find that gain for humidity is max

$IG(T_{\text{sunny}}) =$

Day	Outlook	Temp	Humidity	Wind	Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes

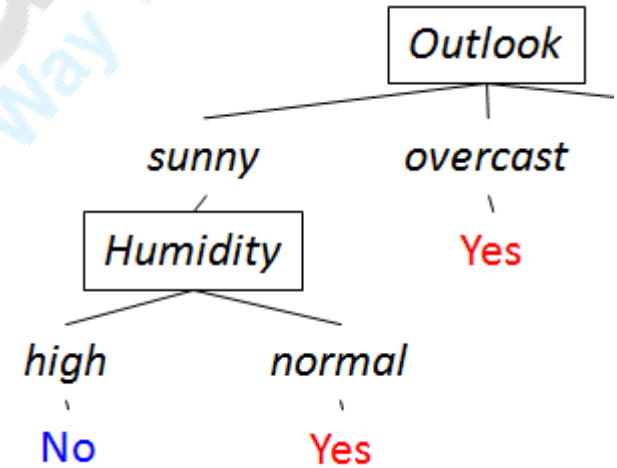


T

p= yes=9  
n=no =5

Day	Outlook	Temp	Humidity	Wind	Tennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

All "yes"  
Thus leaf node found



T overcast

Day	Outlook	Temp	Humidity	Wind	Tennis?
D3	Overcast	Hot	High	Weak	Yes
D7	Overcast	Cool	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes

p= yes=4  
n=no =0

T

p= yes=9

n=no =5

Day	Outlook	Temp	Humidity	Wind	Tennis?
<i>D1</i>	Sunny	Hot	High	Weak	No
<i>D2</i>	Sunny	Hot	High	Strong	No
<i>D3</i>	Overcast	Hot	High	Weak	Yes
<i>D4</i>	Rain	Mild	High	Weak	Yes
<i>D5</i>	Rain	Cool	Normal	Weak	Yes
<i>D6</i>	Rain	Cool	Normal	Strong	No
<i>D7</i>	Overcast	Cool	Normal	Strong	Yes
<i>D8</i>	Sunny	Mild	High	Weak	No
<i>D9</i>	Sunny	Cool	Normal	Weak	Yes
<i>D10</i>	Rain	Mild	Normal	Weak	Yes
<i>D11</i>	Sunny	Mild	Normal	Strong	Yes
<i>D12</i>	Overcast	Mild	High	Strong	Yes
<i>D13</i>	Overcast	Hot	Normal	Weak	Yes
<i>D14</i>	Rain	Mild	High	Strong	No

As humidity is already considered we will exclude it.

T  
Rain

Day	Outlook	Temp	Humidity	Wind	Tennis?
<i>D4</i>	Rain	Mild	High	Weak	Yes
<i>D5</i>	Rain	Cool	Normal	Weak	Yes
<i>D6</i>	Rain	Cool	Normal	Strong	No
<i>D10</i>	Rain	Mild	Normal	Weak	Yes
<i>D14</i>	Rain	Mild	High	Strong	No

p= yes=3

n=no =2

# Entropy

$$E(A) = \sum \frac{p_i + n_i}{p+n} [ I(p_i, n_i) ]$$

TEMP	p	n	I(p,n)
Hot	0	0	0
Mild	2	1	
cold	1	1	

wind	p	n	I(p,n)
Weak	3	0	0
strong	0	2	0

E(temp)=0.950978

E(wind)= 0

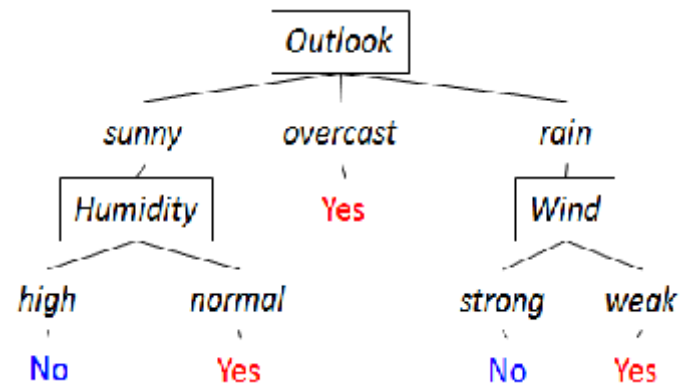
Gain(temp)= IG<sub>sunny</sub> - E(temp)= 0.019973

Gain(wind)= IG<sub>sunny</sub> - E(wind)= 0.970951

Wind is the next branch

T<sub>Rain</sub>

Day	Outlook	Temp	Humidity	Wind	Tennis?
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D10	Rain	Mild	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



# Example 1:

	age	competition	Type	Profit
1	old	Yes	software	down
2	old	No	Software	Down
3	old	No	Hardware	Down
4	Mid	Yes	Software	Down
5	Mid	Yes	Hardware	Down
6	Mid	No	Hardware	Up
7	Mid	No	Software	Up
8	new	Yes	Software	Up
9	new	No	Hardware	Up
10	new	No	Software	Up

TABLE 4.1: Data for Height Classification

Name	Gender	Height	Output1
Kristina	F	1.6 m	Short
Jim	M	2 m	Tall
Maggie	F	1.9 m	Medium
Martha	F	1.88 m	Medium
Stephanie	F	1.7 m	Short
Bob	M	1.85 m	Medium
Kathy	F	1.6 m	Short
Dave	M	1.7 m	Short
Worth	M	2.2 m	Tall
Steven	M	2.1 m	Tall
Debbie	F	1.8 m	Medium
Todd	M	1.95 m	Medium
Kim	F	1.9 m	Medium
Amy	F	1.8 m	Medium
Wynette	F	1.75 m	Medium