

Sudha Amarnath



**SAN JOSÉ STATE
UNIVERSITY**

Graduate and Extended Studies

FA19: CMPE-297 Sec 01 - Special Topics

Prof. Chandrasekar Vuppalapati

Assignment #2

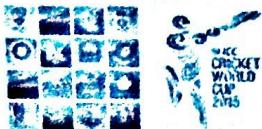
Sudha Amarnath
ID: 013709956

ICC Worldcup 2015 PlayCricket - Decisiontree

Please Develop Decision tree for ICC worldcup 2015 PlayCricket - Decision tree

Deliverable:

- Decision tree calculations



Training example for the target match "PlayCricket"

Match Day	Outlook	Temperature	Humidity	Wind	Play Cricket
Day 1	Sunny	Hot	High	Weak	No
Day 2	Sunny	Hot	High	Strong	No
Day 3	Overcast	Hot	High	Weak	Yes
Day 4	Rain	Mild	High	Weak	Yes
Day 5	Rain	Cool	Normal	Weak	Yes
Day 6	Rain	Cool	Normal	Strong	No
Day 7	Overcast	Cool	Normal	Strong	Yes
Day 8	Sunny	Mild	High	Weak	No
Day 9	Sunny	Cool	Normal	Weak	Yes
Day 10	Rain	Mild	Normal	Weak	Yes
Day 11	Sunny	Mild	Normal	Strong	Yes
Day 12	Overcast	Mild	High	Strong	Yes
Day 13	Overcast	Hot	Normal	Weak	Yes
Day 14	Rain	Mild	High	Strong	No

1. classification using ID3 algorithm to build a decision tree for the given table. This uses Entropy and information gain as metric.

Independent variables	Dependent variable
- outlook - Temperature - Humidity - Wind	- play Cricket

3. To calculate the root for the decision tree, calculate Entropy of the target ($\text{Yes} = 9, \text{No} = 5$)
- The expected information needed to classify a tuple in the data partition D , is given by

$$\text{Entropy}(D) = - \sum_{i=1}^m P_i \log(P_i)$$

$$\begin{aligned}\text{Entropy}(\text{Play Cricket}) &= \text{Entropy}(5, 9) \\ &= \text{Entropy}(0.36, 0.64) \\ &= -(0.36 \log_2 0.36) - (0.64 \log_2 0.64) \\ &= 0.94\end{aligned}$$

4. The dataset is split on the different attributes. The Entropy is calculated for each branch and total entropy is calculated.

- From the table derive the data for the attribute outlook

outlook	play		Total
	Yes	No	
Sunny	2	3	5
Overscast	4	0	4
Rainy	3	2	5
			14

- The average weighted entropy is calculated using:

$$\text{Entropy}(D) = \sum_{j=1}^{|D|} \frac{|D_j|}{|D|} \times \text{Entropy}(D_j)$$

$$\begin{aligned}
 \text{Entropy}(D, \text{outlook}) &= \frac{5}{14} \times \text{Entropy}(2, 3) + \left(\frac{4}{14}\right) \times \\
 &\quad \text{Entropy}(4, 0) + \left(\frac{5}{14}\right) \times \text{Entropy}(3, 2) \\
 &= \left(\frac{5}{14}\right) \left(-\left(\frac{2}{5}\right) \log\left(\frac{2}{5}\right) - \left(\frac{3}{5}\right) \log\left(\frac{3}{5}\right)\right) + \left(\frac{4}{14}\right)(0) + \\
 &\quad \left(\frac{5}{14}\right) \left(-\left(\frac{3}{5}\right) \log\left(\frac{3}{5}\right) - \left(\frac{2}{5}\right) \log\left(\frac{2}{5}\right)\right) \\
 &= 0.693
 \end{aligned}$$

- Information gain for the above data is calculated using

$$\text{Gain}(A) = \text{Entropy}(D) - \text{Entropy}_n(D)$$

$$\text{Gain}(D, \text{outlook}) = 0.940 - 0.693 = 0.247$$

- calculate information gain for each of the attributes

outlook	Playcricket	
	Yes	No
Sunny	2	3
overcast	4	0
Rainy	3	2

$$\text{Gain} = 0.247$$

Temperature	Playcricket	
	Yes	No
Hot	2	2
Mild	4	2
Cool	3	1

$$\text{Gain} = 0.029$$

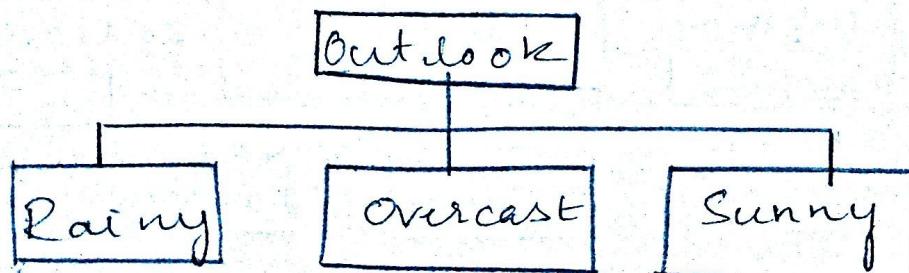
Humidity	Playcricket	
	Yes	No
High	3	4
Normal	6	1

$$\text{Gain} = 0.152$$

Wind	Playcricket	
	Yes	No
Weak	6	2
Strong	3	3

$$\text{Gain} = 0.048$$

5. choose the attribute with largest Information Gain as the decision node (root). Repeat the same step on each branch by dividing the dataset

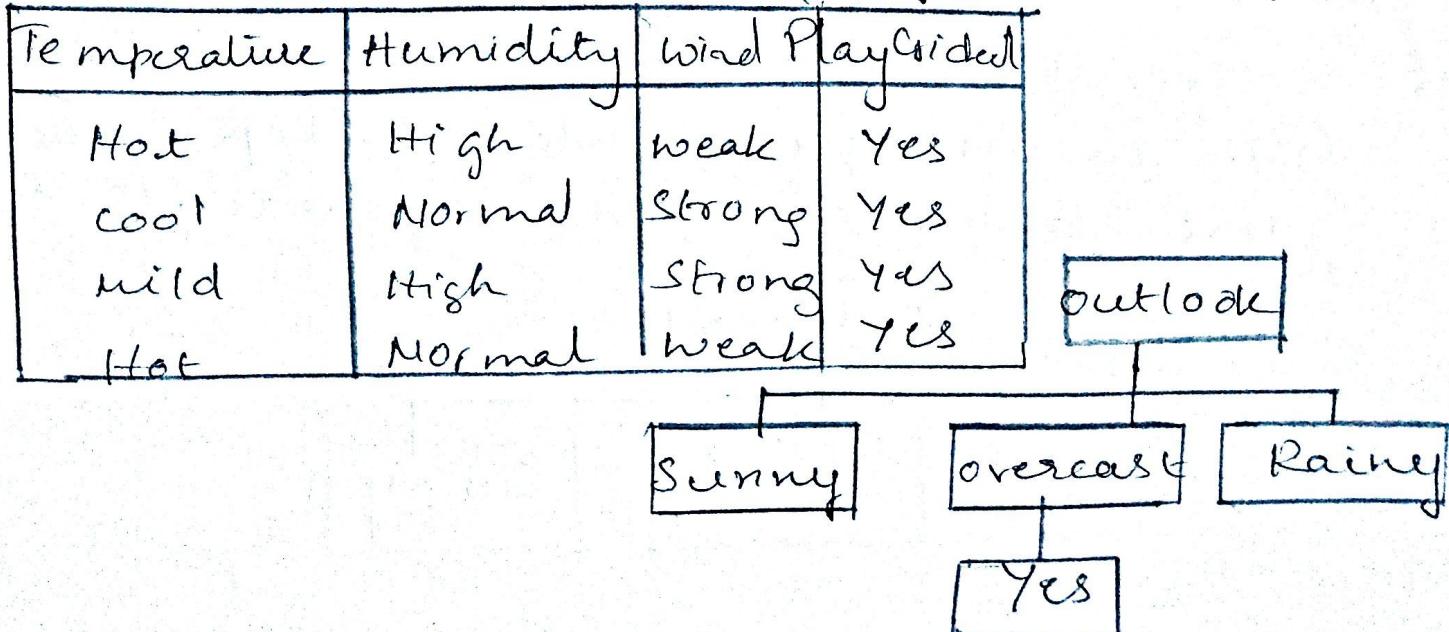


outlook	Temperature	Humidity	wind	PlayCricket
Sceney	Hot	High	weak	NO
Sunny	Hot	High	Strong	NO
Sunny	Mild	High	weak	NO
Sunny	cool	Normal	weak	YES
Sunny	Mild	Normal	Strong	YES

overcast	Hot	High	weak	YES
overcast	cool	Normal	Strong	YES
overcast	Mild	High	Strong	YES
overcast	Hot	Normal	weak	YES

Rain	Mild	High	weak	YES
Rain	cool	Normal	weak	YES
Rain	cool	Normal	Strong	NO
Rain	Mild	Normal	weak	YES
Rain	Mild	High	Strong	NO

- The branch node with Entropy 0 is the leaf node



6. Entropy more than 0 needs further splitting.
Consider the table for attribute weather outlook with tuple sunny.

Temperature	Humidity	Wind	Play/cricket
Hot	High	weak	NO
Hot	High	Strong	NO
Mild	High	weak	NO
Cool	Normal	weak	YES
Mild	Normal	Strong	YES

Temp	play		Total
	Yes	No	
Hot	0	2	2
Cool	1	1	2
Mild	1	0	1
			5

$$E(\text{Sunny}) = \left(-\frac{2}{5} \log \left(\frac{2}{5} \right) - \frac{3}{5} \log \left(\frac{3}{5} \right) \right) = 0.971$$

$$E(\text{Sunny, temperature}) = \left(\frac{2}{5} \right) * E(0,2) + \left(\frac{2}{5} \right) * E(1,1) + \left(\frac{1}{5} \right) * E(1,0) = 2 \cdot \frac{1}{5} = 0.4$$

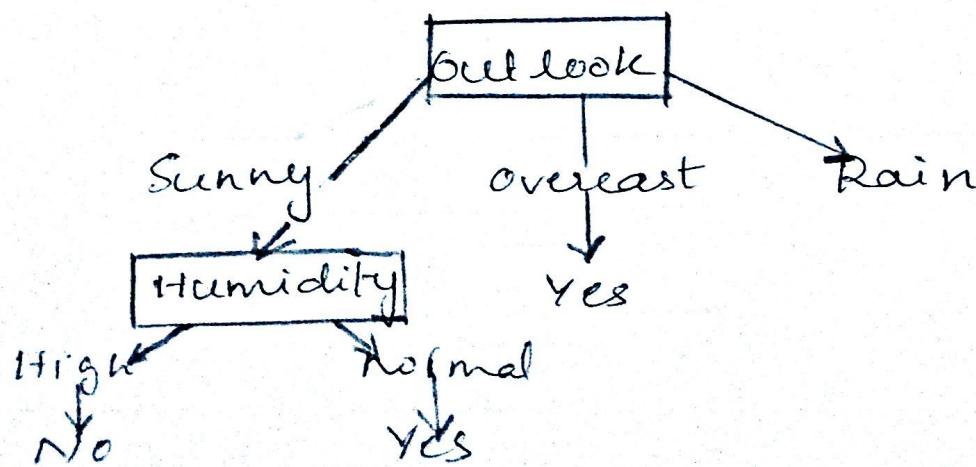
$$\text{Gain}(\text{Sunny, temperature}) = 0.971 - 0.4 = 0.571$$

Similarly for other attributes we get

$$\text{Gain}(\text{Sunny, Humidity}) = 0.971$$

$$\text{Gain}(\text{Sunny, Wind}) = 0.020$$

- Gain(Sunny, Humidity) has the largest value, so it is the node that comes under Sunny



7. The ID3 algorithm is run recursively on non-leaf branches until the data is classified

- Finally the tree looks as below

