

# Decision tree

*SONG JIANXING*

Data: weather\_norminal.CSV

num	outlook	temperature	humidity	windy	play
1	sunny	hot	high	F	no
2	sunny	hot	high	T	no
3	Overcast	hot	high	F	yes
4	rainy	mild	high	F	yes
5	rainy	cool	normal	F	yes
6	rainy	cool	normal	T	no
7	Overcast	cool	normal	T	yes
8	sunny	mild	high	F	no
9	sunny	cool	normal	F	yes
10	rainy	mild	normal	F	yes
11	sunny	mild	normal	T	yes
12	Overcast	mild	high	T	yes
13	Overcast	hot	normal	F	yes
14	rainy	mild	high	T	no

## Method I: ID3

response variable(play): [yes9,no5]

**Total:**

$$\begin{aligned}\text{Entropy}(s) &= \text{info}([9,5]) = -9/14\log(9/14) - 5/14\log(5/14) \text{ (base number of log is 2)} \\ &= 0.940\end{aligned}$$

Use each of predictor variable (outlook, temperature, humidity, wind) as key to classification discuss.

**Outlook:**

**Sunny:**

Result of sunny: [yes2,no3]

$$\text{Entropy}(\text{sunny}) = \text{Info}([2,3]) = -2/5\log(2/5) - 3/5\log(3/5) = 0.971$$

**Overcast:**

Result of overcast: [yes4,no0]

$$\text{Entropy}(\text{overcast}) = \text{info}([4,0]) = -4/4\log(4/4) - 0\log(0) = 0$$

**Rainy:**

Result of rainy: [yes3,no2]

$$\text{Entropy}(\text{rainy}) = \text{info}([3,2]) = -3/5\log(3/5) - 2/5\log(2/5) = 0.971$$

**Summary outlook:**

$$\begin{aligned}\text{Entropy}(\text{outlook}) &= \text{info}([2,3],[4,0],[3,2]) = 5/14\text{info}[2,3] + 4/14\text{info}[4,0] + \\ &5/14\text{info}[3,2] \\ &= 5/14 \times 0.971 + 4/14 \times 0 + 5/14 \times 0.971 \\ &= 0.694\end{aligned}$$

**information gain:**

$$\begin{aligned}\text{gain}(\text{outlook}) &= \text{info}([9,5]) - \text{info}([2,3],[4,0],[3,2]) \\ &= 0.940 - 0.694 \\ &= 0.246\end{aligned}$$

**Temperature:**

**Hot:**

Result of hot: [yes2,no2]

$$\text{Entropy}(\text{hot}) = \text{Info}([2,2]) = -2/4\log(2/4) - 2/4\log(2/4) = 1$$

**Mild:**

Result of mild: [yes5,no1]

$$\text{Entropy}(\text{mild}) = \text{info}([5,1]) = -5/6\log(5/6) - 1/6\log(1/6) = 0.65$$

**Cool:**

Result of cool: [yes3,no1]

$$\text{Entropy}(\text{cool}) = \text{info}([3,1]) = -3/4\log(3/4) - 1/4\log(1/4) = 0.811$$

**Summary temperature:**

$$\begin{aligned}\text{Entropy}(\text{temperature}) &= \text{info}([2,2],[5,1],[3,1]) = 4/14\text{info}([2,2]) + 6/14\text{info}([5,1]) \\ &\quad + 4/14\text{info}([3,1]) \\ &= 4/14 \times 1 + 6/14 \times 0.65 + 4/14 \times 0.811 \\ &= 0.796\end{aligned}$$

**information gain:**

$$\begin{aligned}\text{gain}(\text{temperature}) &= \text{info}([9,5]) - \text{info}([2,2],[5,1],[3,1]) \\ &= 0.940 - 0.796 \\ &= 0.029\end{aligned}$$

**Humidity:**

**High:**

Result of high: [yes3,no4]

$$\text{Entropy}(\text{high}) = \text{info}([3,4]) = -3/7\log(3/7) - 4/7\log(4/7) = 0.985$$

**Normal:**

Result of normal: [yes6,no1]

$$\text{Entropy}(\text{normal}) = \text{info}([6,1]) = -6/7\log(6/7) - 1/7\log(1/7) = 0.591$$

**Summary humidity:**

$$\text{Entropy}(\text{humidity}) = \text{info}([3,4],[6,1]) = 7/14\text{info}([3,4]) + 7/14\text{info}([6,1]) = 0.788$$

**information gain:**

$$\text{gain}(\text{humidity}) = \text{info}([9,5]) - \text{info}([3,4],[6,1]) = 0.940 - 0.788 = 0.152$$

**Windy:**

**True:**

Result of true: [yes3,no3]

Entropy(ture) = Info([3,3]) =  $-3/6\log 3/6 - 3/6\log 3/6 = 1$

**False:**

Result of false: [yes6,no2]

Entropy(false) = info([6,2]) =  $-6/8\log(6/8) - 2/8\log(2/8) = 0.811$

**Summary windy:**

Entropy(windy) = info([3,3],[6,2]) = 0.892

**information gain:**

gain(windy) = info([9,5]) – info([3,3],[6,2]) = 0.05

gain(outlook) = 0.25

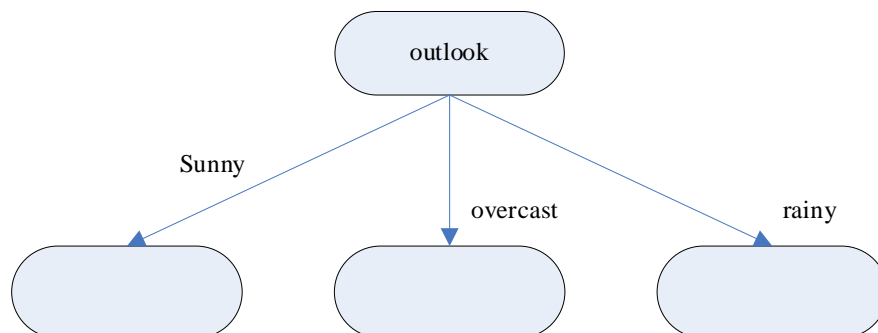
gain(tempearture) = 0.03

gain(humidity) = 0.15

gain(wind) = 0.05

compare them to find gain(outlook) is max

**A: Set ‘outlook’ as Root of the tree, the preliminary is as follows:**



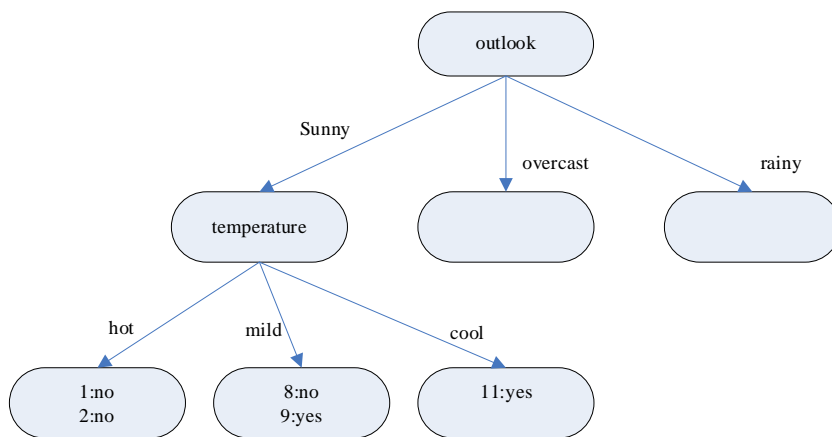
Take ‘sunny’ as total information

Sunny: [yes2,no3]

Info([2,3]) =  $-2/5\log(2/5) - 3/5\log(3/5) = 0.97$

**B: Sunny: All possible results are as follows**

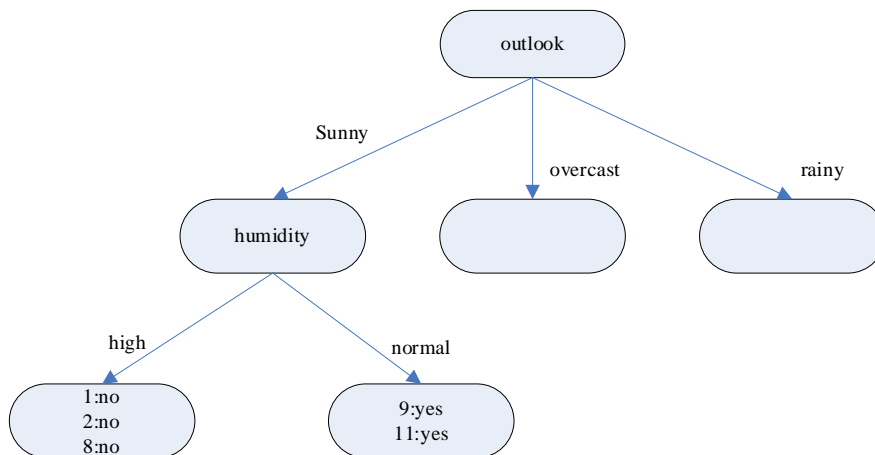
### Result I



$$\text{Entropy}(\text{temperature}) = \text{Info}([0,2],[1,1],[1,0]) = 0.40$$

$$\text{Gain}(\text{temperature}) = \text{Info}([2,3]) - \text{Info}([0,2],[1,1],[1,0]) = 0.57$$

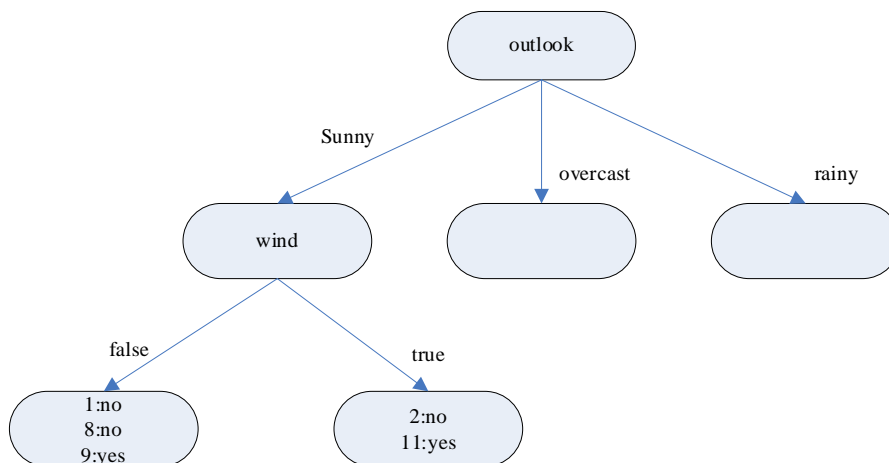
### Result II



$$\text{Entropy}(\text{humidity}) = \text{info}([0,3],[2,0]) = 0$$

$$\text{Gain}(\text{humidity}) = \text{Info}([2,3]) - \text{info}([0,3],[2,0]) = 0.97$$

### Result III



$$\text{Entropy}(\text{wind}) = \text{info}([2,1],[1,1]) = 0.95$$

$$\text{Gain}(\text{wind}) = \text{Info}([2,3]) - \text{info}([2,1],[1,1]) = 0.02$$

$$\text{Gain}(\text{temperature}) = 0.57$$

$$\text{Gain}(\text{humidity}) = 0.97$$

$$\text{Gain}(\text{wind}) = 0.02$$

compare them to find gain(humidity) is max

**C: Set 'humidity' as Internal node of the 'sunny'**

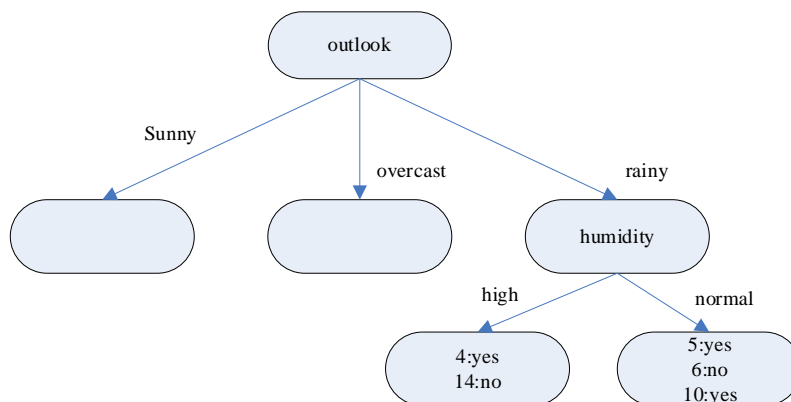
**D: For 'overcast', recursive termination has been formed, do not need to classified**

**F: Rainy: All possible results are as follows**

Rainy:[yes3,no2]

$$\text{Info}([3,2]) = -3/5\log(3/5) - 2/5\log(2/5) = 0.971$$

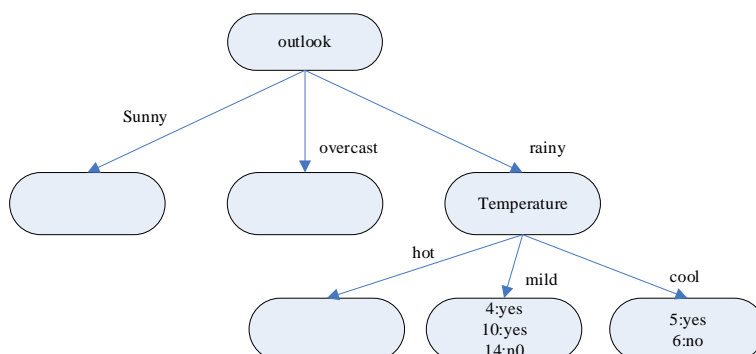
**Result I**



$$\text{Entropy}(\text{humidity}) = \text{info}([1,1],[2,1]) = 0.951$$

$$\text{Gain}(\text{humidity}) = \text{Info}([3,2]) - \text{info}([1,1],[2,1]) = 0.02$$

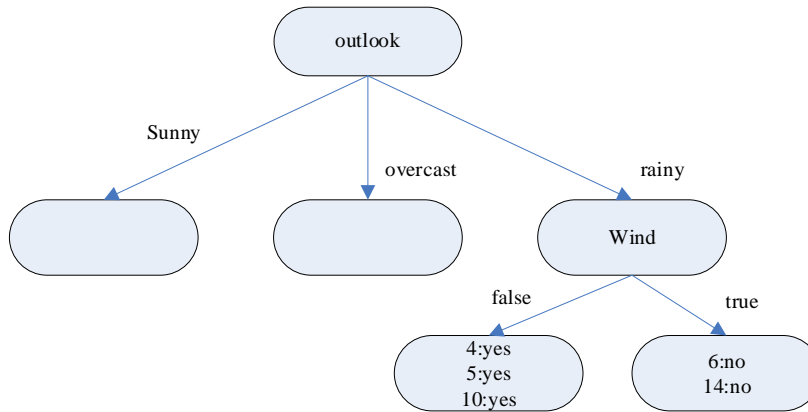
**Result II**



$$\text{Entropy}(\text{temperature}) = \text{info}([2,1],[1,1]) = 0.951$$

$$\text{Gain}(\text{temperature}) = \text{Info}([3,2]) - \text{info}([1,1],[2,1]) = 0.02$$

### Result III



$$\text{Entropy}(\text{wind}) = \text{info}([3,0],[0,2]) = 0$$

$$\text{Gain}(\text{wind}) = \text{Info}([3,2]) - \text{info}([3,0],[0,2]) = 0.971$$

$$\text{Gain}(\text{humidity}) = 0.02$$

$$\text{Gain}(\text{temperature}) = 0.02$$

$$\text{Gain}(\text{wind}) = 0.971$$

compare them to find gain(wind) is max

**G: Set 'wind' as Internal node of the 'rainy'**

**H: To sum up, the decision tree is generated as follows**

