

Tutorial Business Analytics

Tutorial 6 - Solution

Exercise 6.1

Calculate:

- a) $\text{entropy}(0.1, 0.9)$
- b) $\text{entropy}(0.8, 0.2)$
- c) $\text{entropy}(0.3, 0.7)$
- d) $\text{entropy}(0.5, 0.5)$
- e) $\text{entropy}(0.8, 0.1, 0.1)$

Solution 6.1

- a) $\text{entropy}(0.1, 0.9) = -0.1 \cdot \log_2 0.1 - 0.9 \cdot \log_2 0.9 = 0.469.$
- b) $\text{entropy}(0.8, 0.2) = 0.722.$
- c) $\text{entropy}(0.3, 0.7) = 0.881.$
- d) $\text{entropy}(0.5, 0.5) = 1.$
- e) $\text{entropy}(0.8, 0.1, 0.1) = 0.922.$

Exercise 6.2

Calculate:

- a) $\text{info}([2, 3])$
- b) $\text{info}([5, 4])$
- c) $\text{info}([2,3], [5, 4])$
- d) $\text{info}([2, 3], [9, 0])$

Solution 6.2

a) $\text{info}([2, 3]) = \text{entropy}\left(\frac{2}{2+3}, \frac{3}{2+3}\right) = 0.971.$

b) $\text{info}([5, 4]) = 0.991.$

c) $\text{info}([2, 3], [5, 4]) = \frac{2+3}{(2+3)+(5+4)} \cdot \text{info}([2, 3]) + \frac{5+4}{(2+3)+(5+4)} \cdot \text{info}([5, 4]) = 0.984.$

d) $\text{info}([2, 3], [9, 0]) = 0.347.$

Exercise 6.3

Construct a tree:

Temperature	Visibility	Snow depth	Sport
< -5	Clear	≥ 50	Skiing
< -5	Fog	≥ 50	Swimming
< -5	Fog	< 50	Swimming
< -5	Rain	≥ 50	Skiing
< -5	Rain	< 50	Swimming
≥ -5	Clear	≥ 50	Skiing
≥ -5	Clear	< 50	Skiing
≥ -5	Fog	< 50	Swimming
≥ -5	Rain	≥ 50	Skiing

Solution 6.3

Constructing the tree:

	< -5	≥ -5	Clear	Fog	Rain	< 50	≥ 50	Σ
Skiing	2	3	3	0	2	1	4	5
Swimming	3	1	0	3	1	3	1	4

Attribute temperature

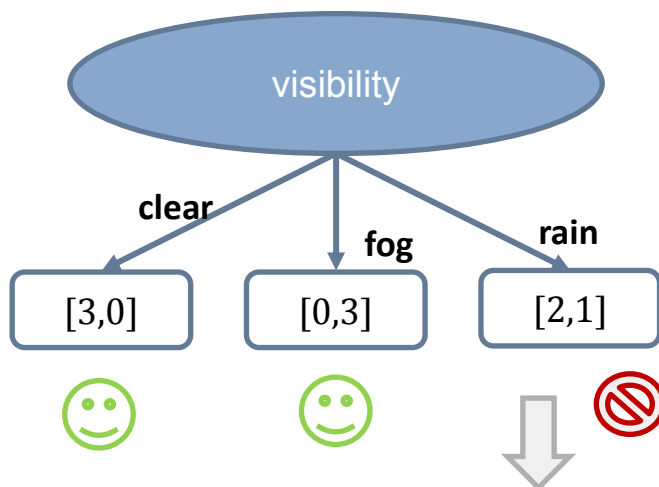
$$\begin{aligned}
 \text{gain}(\text{temperature}) &= \text{info}(\text{root}) - \text{info}(\text{temperature}) \\
 &= \text{info}([5, 4]) - \text{info}([2, 3], [3, 1]) \\
 &= 0.991 - 0.900 = 0.091
 \end{aligned}$$

Attribute visibility

$$\begin{aligned}
 \text{gain}(\text{visibility}) &= \text{info}(\text{root}) - \text{info}(\text{visibility}) \\
 &= \text{info}([5, 4]) - \text{info}([3, 0], [0, 3], [2, 1]) \\
 &= 0.991 - 0.306 = 0.685
 \end{aligned}$$

Attribute snow depth

$$\begin{aligned}
 \text{gain}(\text{snow depth}) &= \text{info}(\text{root}) - \text{info}(\text{snow depth}) \\
 &= \text{info}([5, 4]) - \text{info}([1, 3], [4, 1]) \\
 &= 0.991 - 0.762 = 0.229
 \end{aligned}$$



Temperature	Visibility	Snow depth	Sport
< -5	Rain	≥ 50	Skiing
< -5	Rain	< 50	Swimming
≥ -5	Rain	≥ 50	Skiing

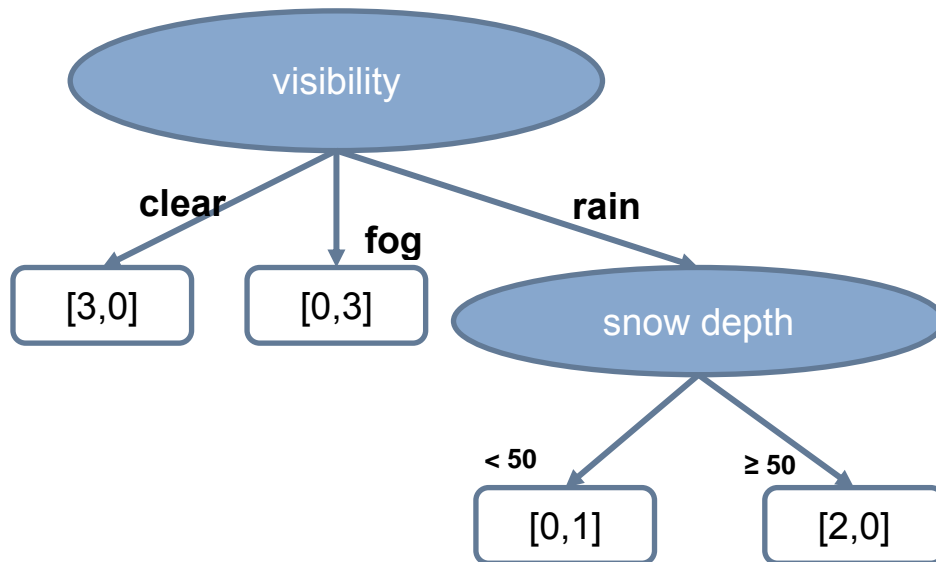
Attribute temperature

$$\begin{aligned}
 \text{gain}(\text{temperature}) &= \text{info}([2, 1]) - \text{info}([1, 1], [1, 0]) \\
 &= 0.918 - 0.667 \\
 &= 0.251
 \end{aligned}$$

Attribute snow depth

$$\begin{aligned}\text{gain}(\text{snow depth}) &= \text{info}([2, 1]) - \text{info}([0, 1], [2, 0]) \\ &= 0.918 - 0 \\ &= 0.918\end{aligned}$$

Result:



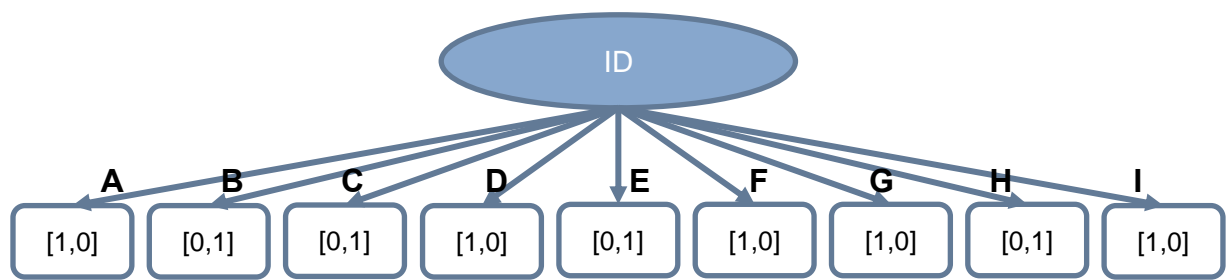
Exercise 6.4

Construct a tree for the following table, that has an additional ID attribute now.

ID	Temperature	Visibility	Snow depth	Sport
A	< -5	Clear	≥ 50	Skiing
B	< -5	Fog	≥ 50	Swimming
C	< -5	Fog	< 50	Swimming
D	< -5	Rain	≥ 50	Skiing
E	< -5	Rain	< 50	Swimming
F	≥ -5	Clear	≥ 50	Skiing
G	≥ -5	Clear	< 50	Skiing
H	≥ -5	Fog	< 50	Swimming
I	≥ -5	Rain	≥ 50	Skiing

Solution 6.4

Construct a tree:



Exercise 6.5

Construct the tree from exercise 6.4 a second time using gain ratio:

ID	Temperature	Visibility	Snow depth	Sport
A	< -5	Clear	≥ 50	Skiing
B	< -5	Fog	≥ 50	Swimming
C	< -5	Fog	< 50	Swimming
D	< -5	Rain	≥ 50	Skiing
E	< -5	Rain	< 50	Swimming
F	≥ -5	Clear	≥ 50	Skiing
G	≥ -5	Clear	< 50	Skiing
H	≥ -5	Fog	< 50	Swimming
I	≥ -5	Rain	≥ 50	Skiing

Solution 6.5

Construct the tree from exercise 6.4 a second time using gain ratio:

	< -5	≥ -5	Clear	Fog	Rain	< 50	≥ 50	Σ
Skiing	2	3	3	0	2	1	4	5
Swimming	3	1	0	3	1	3	1	4

Attribute temperature

- $\text{gain}(\text{temperature}) = \text{info}([5,4]) - \text{info}([2,3], [3,1]) = 0.991 - 0.900 = 0.091$
- $\text{gainRatio}(\text{temperature}) = 0.091 / \text{info}([5,4]) = 0.091 / 0.991 = 0.092$

Attribute visibility

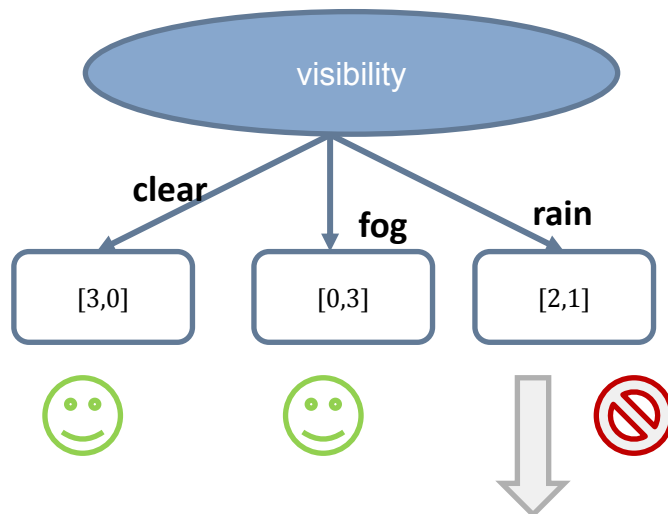
- $\text{gain}(\text{visibility}) = \text{info}([5,4]) - \text{info}([3,0], [0,3], [2,1]) = 0.991 - 0.306 = 0.685$
- $\text{gainRatio}(\text{visibility}) = 0.685 / \text{info}([3,3,3]) = 0.685 / 1.585 = \mathbf{0.432}$

Attribute snow depth

- $\text{gain}(\text{snow depth}) = \text{info}([5,4]) - \text{info}([1,3], [4,1]) = 0.991 - 0.762 = 0.229$
- $\text{gainRatio}(\text{snow depth}) = 0.229 / \text{info}([4,5]) = 0.229 / 0.991 = 0.231$

Attribute ID

- $\text{gain}(\text{ID}) = \text{info}([5,4]) - \text{info}([1,0], [1,0], [1,0], [1,0], [1,0], [1,0], [1,0], [1,0], [1,0], [1,0]) = 0.991$
- $\text{gainRatio}(\text{ID}) = 0.991 / \text{info}([1,1,1,1,1,1,1,1,1,1]) = 0.991 / 3.17 = 0.313$



ID	Temperature	Visibility	Snow depth	Sport
D	< -5	Rain	≥ 50	Skiing
E	< -5	Rain	< 50	Swimming
I	≥ -5	Rain	≥ 50	Skiing

Attribute temperature

- $\text{gain}(\text{temperature}) = \text{info}([2, 1]) - \text{info}([1, 1], [1, 0]) = 0.918 - 0.667 = 0.251$
- $\text{gainRatio}(\text{temperature}) = 0.251 / \text{info}([2,1]) = 0.251 / 0.918 = 0.273$

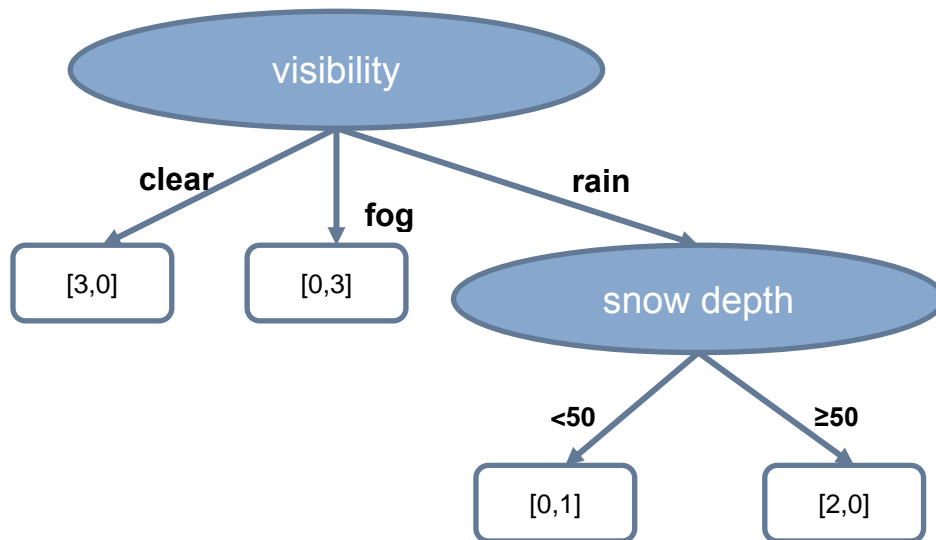
Attribute snow depth

- $\text{gain}(\text{snow depth}) = \text{info}([2, 1]) - \text{info}([0, 1], [2, 0]) = 0.918 - 0 = 0.918$
- $\text{gainRatio}(\text{snow depth}) = 0.981 / \text{info}([1, 2]) = 0.918 / 0.918 = 1$

Attribute ID

- $\text{gain}(\text{ID}) = \text{info}([2, 1]) - \text{info}([1, 0], [0, 1], [1, 0]) = 0.918 - 0 = 0.918$
- $\text{gainRatio}(\text{ID}) = 0.918 / \text{info}([1, 1, 1]) = 0.918 / 1.585 = 0.579$

Result:



Exercise 6.6

Find the optimal binary splits.

a)

60	60	120	120	180	180	180
F	F	T	F	F	T	T

b)

5	5	7	7	7	8	9	9
T	T	T	T	F	T	F	F

Solution 6.6

Find the optimal binary splits.

a)

60	60	120	120	180	180	180
F	F	T	F	F	T	T

All possible splits: [0,2], [1,1], [2,1]

- $\text{info}([0,2],[3,2]) = \mathbf{0.694}$ split at 90
- $\text{info}([1,3],[2,1]) = 0.857$

b)

5	5	7	7	7	8	9	9
T	T	T	T	F	T	F	F

All possible splits: [2,0],[2,1],[1,0],[0,2]

- $\text{info}([2,0],[3,3]) = 0.75$
- $\text{info}([4,1],[1,2]) = 0.796$
- $\text{info}([5,1],[0,2]) = \mathbf{0.488}$ split at 8.5