

Consider the following training dataset and the original decision tree induction algorithm (ID3). Risk is the class label attribute. The Height values have been already discretized into disjoint ranges. Calculate the information gain if Gender is chosen as the test attribute. Calculate the information gain if Height is chosen as the test attribute. Draw the final decision tree (without any pruning) for the training dataset. Generate all the “IF-THEN rules from the decision tree.

Gender	Height	Risk
F	(1.5, 1.6)	Low
M	(1.9, 2.0)	High
F	(1.8, 1.9)	Medium
F	(1.8, 1.9)	Medium
F	(1.6, 1.7)	Low
M	(1.8, 1.9)	Medium
F	(1.5, 1.6)	Low
M	(1.6, 1.7)	Low
M	(2.0, 8)	High
M	(2.0, 8)	High
F	(1.7, 1.8)	Medium
M	(1.9, 2.0)	Medium
F	(1.8, 1.9)	Medium
F	(1.7, 1.8)	Medium
F	(1.7, 1.8)	Medium

1. The original entropy is $I_{\text{Risk}} = I(\text{Low, Medium, High}) = I(4, 8, 3) = 1.4566$. Consider Gender .

Gender	entropy
F	$I(3, 6, 0)$
M	$I(1, 2, 3)$

The expected entropy is $\frac{9}{15} \cdot I(3, 6, 0) + \frac{6}{15} \cdot I(1, 2, 3) = 1.1346$. The information gain is $1.4566 - 1.1346 = 0.3220$

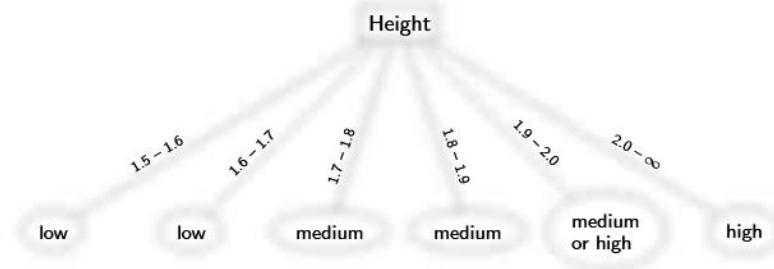
2. Consider Height .

Height	entropy
(1.5, 1.6]	$I(2, 0, 0)$
(1.6, 1.7]	$I(2, 0, 0)$
(1.7, 1.8]	$I(0, 3, 0)$
(1.8, 1.9]	$I(0, 4, 0)$
(1.9, 2.0]	$I(0, 1, 1)$
(2.0, ∞]	$I(0, 0, 2)$

The expected entropy is $\frac{2}{15} \cdot I(2, 0, 0) + \frac{2}{15} \cdot I(2, 0, 0) + \frac{3}{15} \cdot I(0, 3, 0) + \frac{4}{15} \cdot I(0, 4, 0) + \frac{2}{15} \cdot I(0, 1, 1) + \frac{2}{15} \cdot I(0, 0, 2) = 0.1333$. The information gain is $1.4566 - 0.1333 = 1.3233$

3. ID3 decision tree:

- ▶ According to the computation above, we should first choose *Height* to split
- ▶ After split, the only problematic partition is the (1.9, 2.0] one. However, the only remaining attribute *Gender* cannot divide them. As there is a draw, we can take any label.
- ▶ The final tree is show in the figure below.



4. The rules are

- ▶ **IF** $height \in (1.5, 1.6]$, **THEN** $Rish = Low$.
- ▶ **IF** $height \in (1.6, 1.7]$, **THEN** $Rish = Low$.
- ▶ **IF** $height \in (1.7, 1.8]$, **THEN** $Rish = Medium$.
- ▶ **IF** $height \in (1.8, 1.9]$, **THEN** $Rish = Medium$.
- ▶ **IF** $height \in (1.9, 2.0]$, **THEN** $Rish = Medium$ (or High).
- ▶ **IF** $height \in (2.0, \infty]$, **THEN** $Rish = High$.