

Table 8.1 Class-Labeled Training Tuples from the *AllElectronics* Customer Database

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class: buys_computer</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

Example 8.1 Induction of a decision tree using information gain. Table 8.1 presents a training set, D , of class-labeled tuples randomly selected from the *AllElectronics* customer database. (The data are adapted from Quinlan [Qui86]. In this example, each attribute is discrete-valued. Continuous-valued attributes have been generalized.) The class label attribute, *buys_computer*, has two distinct values (namely, {yes, no}); therefore, there are two distinct classes (i.e., $m = 2$). Let class C_1 correspond to yes and class C_2 correspond to no. There are nine tuples of class yes and five tuples of class no. A (root) node N is created for the tuples in D . To find the splitting criterion for these tuples, we must compute the information gain of each attribute. We first use Eq. (8.1) to compute the expected information needed to classify a tuple in D :

$$\text{Info}(D) = -\frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right) = 0.940 \text{ bits.}$$

Next, we need to compute the expected information requirement for each attribute. Let's start with the attribute *age*. We need to look at the distribution of yes and no tuples for each category of *age*. For the *age* category “youth,” there are two yes tuples and three no tuples. For the category “middle_aged,” there are four yes tuples and zero no tuples. For the category “senior,” there are three yes tuples and two no tuples. Using Eq. (8.2), the expected information needed to classify a tuple in D if the tuples are partitioned according to *age* is

$$\text{Info}_{\text{age}}(D) = \frac{5}{14} \times \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right)$$

$$\begin{aligned}
 & + \frac{4}{14} \times \left(-\frac{4}{4} \log_2 \frac{4}{4} \right) \\
 & + \frac{5}{14} \times \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) \\
 & = 0.694 \text{ bits.}
 \end{aligned}$$

Hence, the gain in information from such a partitioning would be

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D) = 0.940 - 0.694 = 0.246 \text{ bits.}$$

Similarly, we can compute $\text{Gain}(\text{income}) = 0.029$ bits, $\text{Gain}(\text{student}) = 0.151$ bits, and $\text{Gain}(\text{credit_rating}) = 0.048$ bits. Because age has the highest information gain among the attributes, it is selected as the splitting attribute. Node N is labeled with age , and branches are grown for each of the attribute's values. The tuples are then partitioned accordingly, as shown in Figure 8.5. Notice that the tuples falling into the partition for $\text{age} = \text{middle_aged}$ all belong to the same class. Because they all belong to class "yes," a leaf should therefore be created at the end of this branch and labeled "yes." The final decision tree returned by the algorithm was shown earlier in Figure 8.2. ■

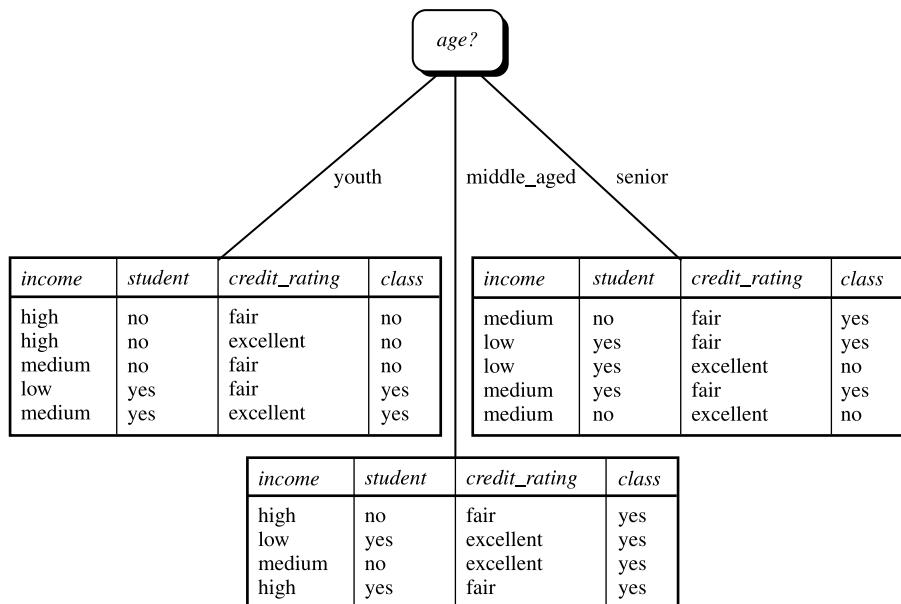


Figure 8.5 The attribute age has the highest information gain and therefore becomes the splitting attribute at the root node of the decision tree. Branches are grown for each outcome of age . The tuples are shown partitioned accordingly.

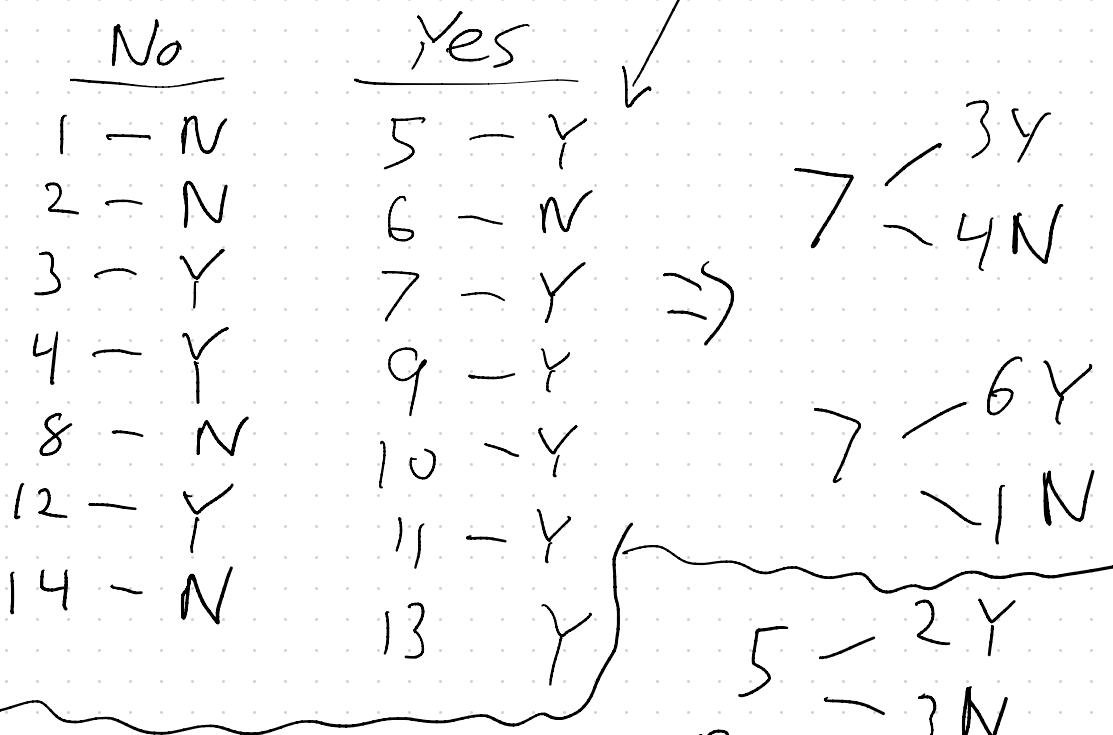
Table 8.1 Class-Labeled Training Tuples from the *AllElectronics* Customer Database

RID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

γ - class Yes buys
Computer.

N - class No

if I use student as split attribute:



VS age \Rightarrow ^{Pure} 4 - 4 Y
5 - 3 Y
2 N

Table 8.1 Class-Labeled Training Tuples from the *AllElectronics* Customer Database

RID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

If I use income as split attr:

low	medium	high
5 - Y	4 - Y	1 - N
6 - N	8 - N	2 - N
7 - Y	10 - Y	3 - Y
9 - Y	12 - Y	13 - Y
	14 - N	

$$\begin{aligned} 4 &\leftarrow 3Y \\ &\quad \sim 1N \\ 6 &\leftarrow 4Y \\ &\quad \sim 2N \\ 4 &\leftarrow 2Y \\ &\quad \sim 2N \end{aligned}$$

VS age \Rightarrow

$$\begin{aligned} 5 &= 2Y \\ &= 3N \end{aligned}$$

$$\begin{aligned} 4 &= 4Y \\ 5 &= 3Y \\ &= 2N \end{aligned}$$

VS Student \Rightarrow

$$\begin{aligned} 7 &\leftarrow 3Y \\ &\sim 4N \\ 7 &\leftarrow 6Y \\ &\sim 1N \end{aligned}$$

Table 8.1 Class-Labeled Training Tuples from the *AllElectronics* Customer Database

RID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

if I use credit as split attr:

<u>fair</u>	<u>excellent</u>
1 - n	2 - n
3 - y	6 - n
4 - y	7 - y
5 - y	11 - y
8 - n	12 - y
9 - y	13 - n
10 - y	
13 - y	



8 -	6 y
1 -	2 n
6 -	3 y
1 -	3 n

VS age \Rightarrow

5 - 2 Y
5 - 3 N
4 - 4 Y
5 - 3 Y
2 N

VS income \Rightarrow

4 - 3 Y
1 N
6 - 4 Y
2 N
4 - 2 Y
2 N

VS student \Rightarrow

7 - 3 Y
7 - 4 N
7 - 6 Y
1 N

VS credit \Rightarrow

8	\leftarrow	6Y
6	\leftarrow	2N
6	\leftarrow	3Y
\leftarrow 3N		

• 048

VS student \Rightarrow

7	\leftarrow	3Y
7	\leftarrow	4N
7	\leftarrow	6Y
\leftarrow 1N		

• 151

best

* VS age \Rightarrow

5	\leftarrow	2Y
4	\leftarrow	3N
5	\leftarrow	4Y
\leftarrow 3N		

• 246

VS income \Rightarrow

6	\leftarrow	3Y
6	\leftarrow	4Y
4	\leftarrow	2Y
\leftarrow 2N		

• 029

* Using age as splitter means

Only need to continue with

10/14, all others require going

on with all

$$Gain(age) = Info(D) - Info_{age}(D) = 0.940 - 0.694 = 0.246 \text{ bits.}$$

Similarly, we can compute $Gain(income) = 0.029$ bits, $Gain(student) = 0.151$ bits, and $Gain(credit_rating) = 0.048$ bits. Because *age* has the highest information gain among the attributes, it is selected as the splitting attribute. Node N is labeled with *age*,

How calculate:

VS credit \Rightarrow

$$.048$$

$$\begin{array}{r} 8 - 6Y \\ 6 - 3Y \\ \hline 3Y \end{array}$$

VS student \Rightarrow

$$.151$$

$$\begin{array}{r} 7 - 3Y \\ 7 - 4N \\ 7 - 6Y \\ \hline 1N \end{array}$$

VS age \Rightarrow

$$\underline{.246}$$

$$\begin{array}{r} 4 - 4Y \\ 5 - 3Y \\ \hline 2Y \end{array}$$

VS income \Rightarrow

$$\underline{.029}$$

$$\begin{array}{r} 6 - 4Y \\ 4 - 2Y \\ \hline 2Y \end{array}$$

credit: $Info_{\text{credit}}(D) =$

$$\begin{aligned} Info(D) - Info_{\text{credit}}(D) &= \\ .940 - .892 &= \end{aligned}$$

$$= .048$$

$$= \frac{8}{14} \left(\frac{1}{4}(-2) + \frac{3}{4}(-.415) \right) + \frac{6}{14} \left(\frac{1}{2}(-1) + \frac{1}{2}(-1) \right)$$

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j).$$

$$Info(D) = - \sum_{i=1}^m p_i \log_2(p_i), \quad \underline{=.892}$$

$$Gain(\text{age}) = Info(D) - Info_{\text{age}}(D) = 0.940 - 0.694 = 0.246 \text{ bits.}$$