
UNIT-5

Classification and Prediction

Jiawei Han, Micheline Kamber, and Jian Pei

University of Illinois at Urbana-Champaign &

Simon Fraser University

©2011 Han, Kamber & Pei. All rights reserved.

Supervised vs. Unsupervised Learning

- Supervised learning (classification)
 - Supervision: The training data (observations, measurements, etc.) are accompanied by **labels** indicating the class of the observations
 - New data is classified based on the training set
- Unsupervised learning (clustering)
 - The class labels of training data is unknown
 - Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

Prediction Problems: Classification vs. Numeric Prediction

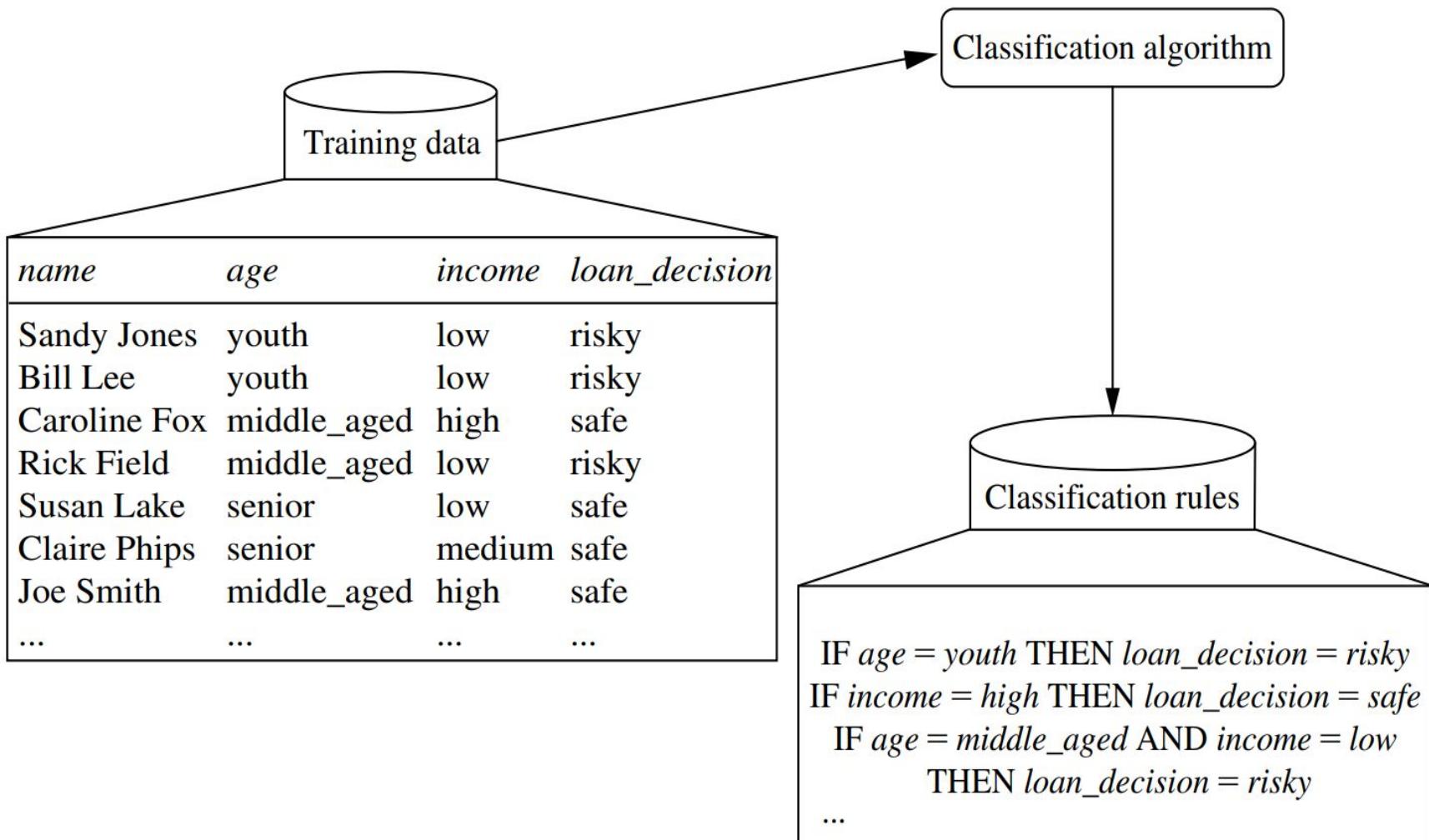
- **Classification**
 - predicts categorical class labels
 - constructs a model based on the training set and the values in a classifying attribute and uses it in classifying new data
- **Numeric Prediction**
 - models continuous-valued functions, i.e., predicts unknown or missing values
 - Regression analysis is a statistical methodology that is most often used for numeric prediction
- **Typical applications**
 - Credit/loan approval:
 - Medical diagnosis: if a tumor is cancerous or benign
 - Fraud detection: if a transaction is fraudulent

Classification—A Two-Step Process

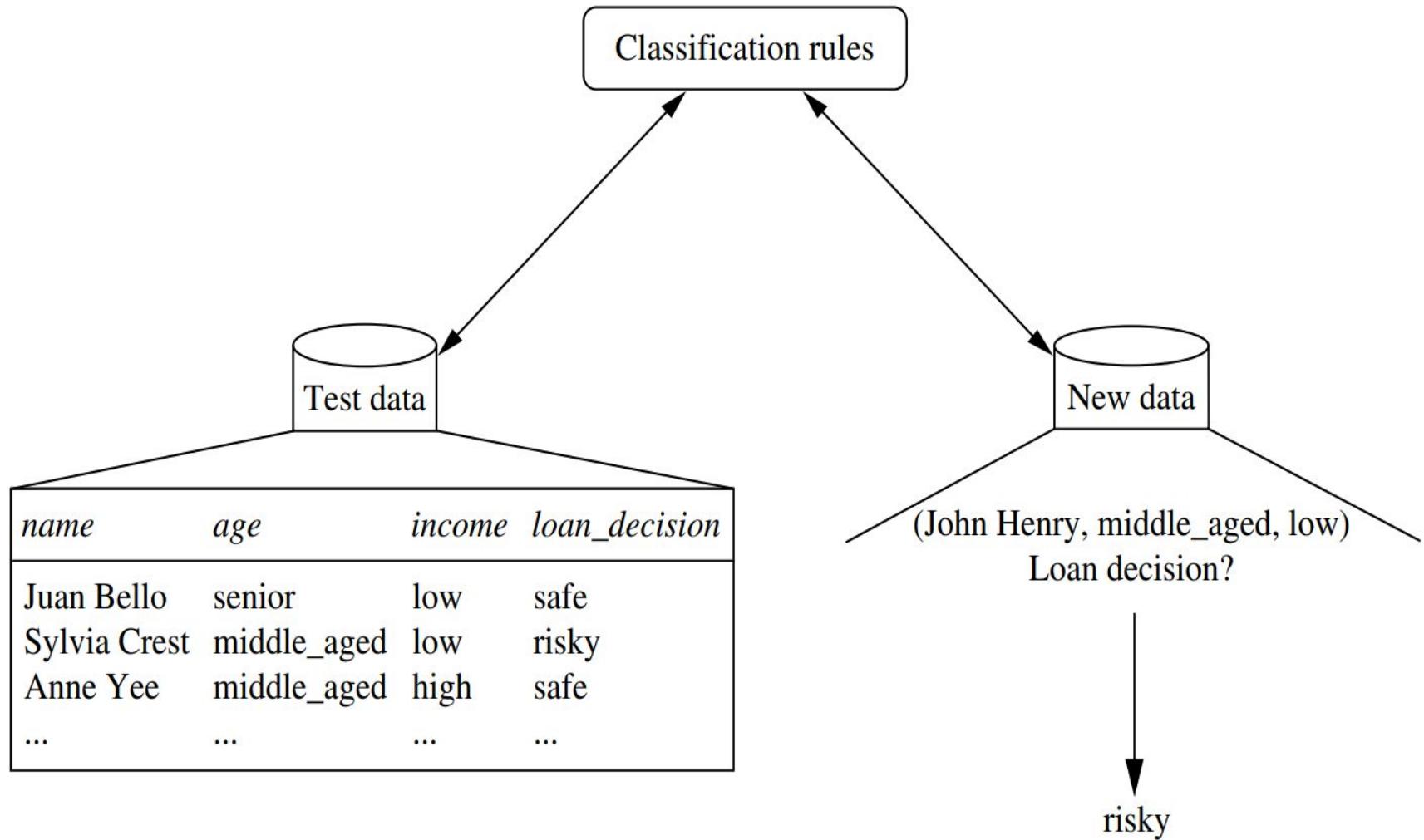
- **Model construction:** describing a set of predetermined classes
 - Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label attribute**
 - The set of tuples used for model construction is **training set**
 - The model is represented as classification rules, decision trees, or mathematical formulae
- **Model usage:** for classifying future or unknown objects
 - **Estimate accuracy** of the model
 - The known label of test sample is compared with the classified result from the model
 - **Accuracy** rate is the percentage of test set samples that are correctly classified by the model
 - **Test set** is independent of training set (otherwise overfitting)

-
- If the accuracy is acceptable, use the model to **classify new data**
 - Note: If *the test set* is used to select models, it is called **validation (test) set**

Process (1): Model Construction



Process (2): Using the Model in Prediction



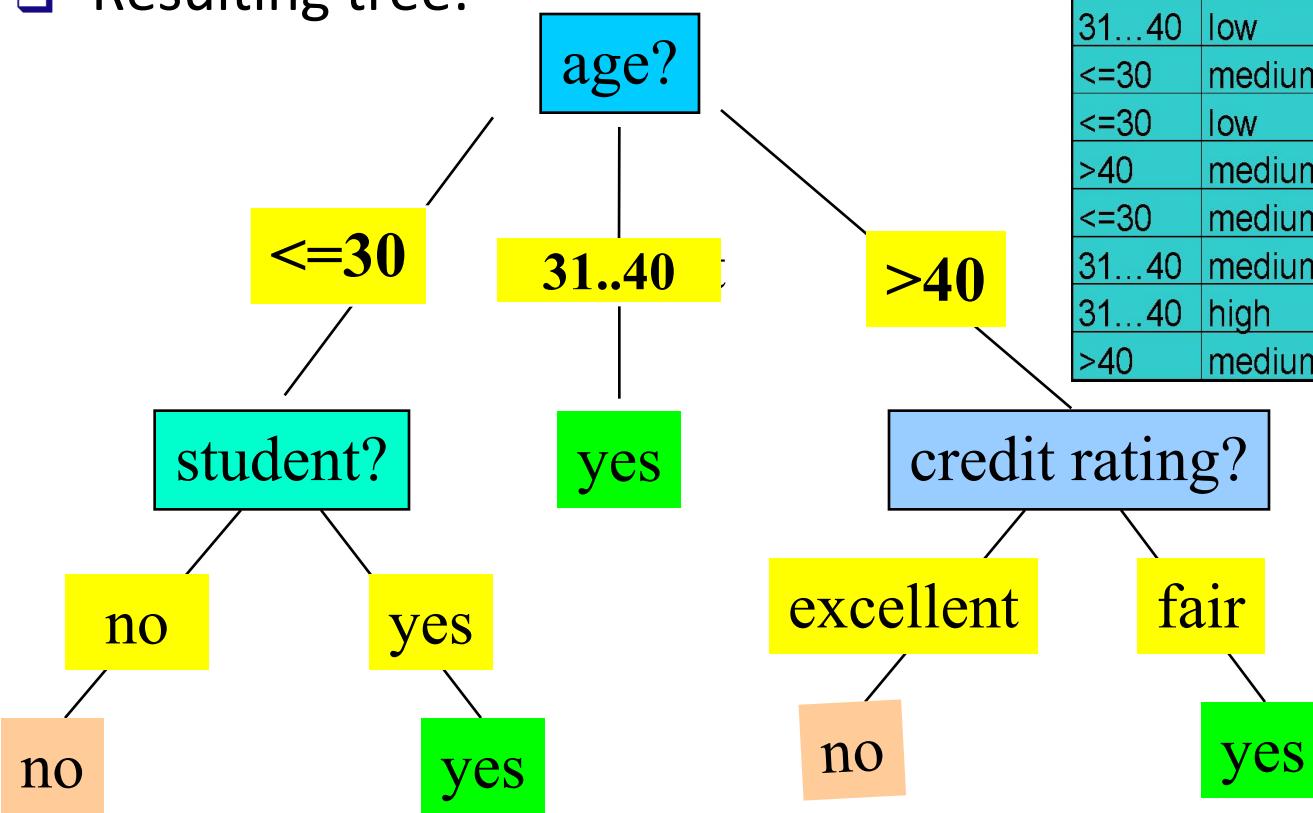
Decision Tree Induction

- A decision tree algorithm known as ID3 (Iterative Dichotomiser).
- A decision tree is a flowchart-like tree structure,
- Where each internal node (nonleaf node) denotes a test on an attribute,
- Each branch represents an outcome of the test
- Each leaf node (or terminal node) holds a class label
- The topmost node in a tree is the root node
- Tree is constructed in a **top-down recursive divide-and-conquer manner**

-
- Attributes are categorical (if continuous-valued, they are discretized in advance)
 - Test attributes are selected on the basis of a heuristic or statistical measure (e.g., **Information gain or Gini Index**)

Decision Tree Induction: An Example

- Training data set: Buys_computer
- The data set follows an example of Quinlan's ID3 (Playing Tennis)
- Resulting tree:



age	income	student	credit rating	buys computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31..40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31..40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31..40	medium	no	excellent	yes
31..40	high	yes	fair	yes
>40	medium	no	excellent	no

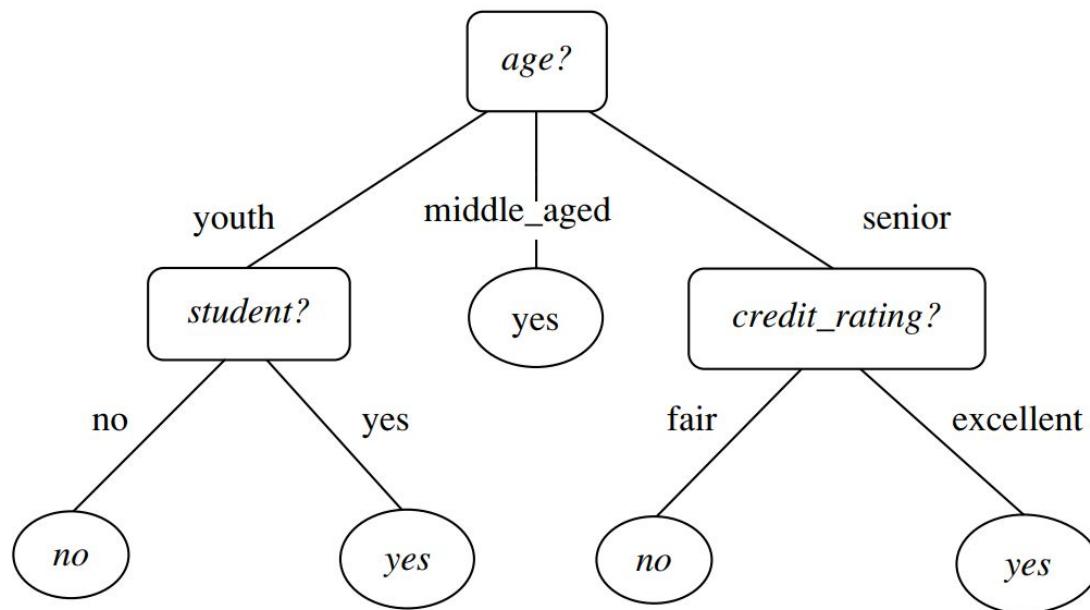


Figure 8.2 A decision tree for the concept *buys_computer*, indicating whether an *AllElectronics* customer is likely to purchase a computer. Each internal (nonleaf) node represents a test on an attribute. Each leaf node represents a class (either *buys_computer* = *yes* or *buys_computer* = *no*).

ENTROPY MEASURES HOMOGENEITY OF EXAMPLES

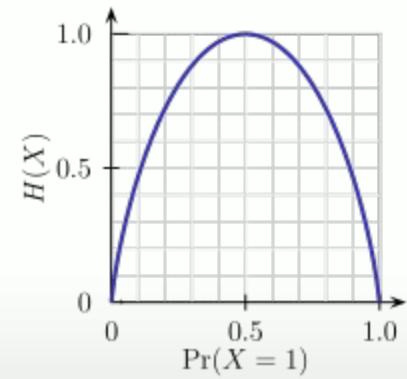
- Entropy measures the *impurity* of a collection of examples. It depends from the distribution of the random variable p .

- S is a collection of training examples

- p_+ the proportion of positive examples in S

- p_- the proportion of negative examples in S

$$\text{Entropy}(S) \equiv -p_+ \log_2 p_+ - p_- \log_2 p_-$$



Examples

$$\text{Entropy}([14+, 0-]) = -14/14 \log_2 (14/14) - 0 \log_2 (0) = 0$$

$$\text{Entropy}([9+, 5-]) = -9/14 \log_2 (9/14) - 5/14 \log_2 (5/14) = 0.94$$

$$\text{Entropy}([7+, 7-]) = -7/14 \log_2 (7/14) - 7/14 \log_2 (7/14) = 1/2 + 1/2 = 1$$

INFORMATION GAIN MEASURES THE EXPECTED REDUCTION IN ENTROPY

- Given entropy as a measure of the impurity in a collection of training examples, the *information gain*, is simply the expected reduction in entropy caused by partitioning the examples according to an attribute.
- More precisely, the information gain, $\text{Gain}(S, A)$ of *an* attribute A , relative to a collection of examples S , is defined as,

$$\text{Gain}(S, A) \equiv \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Values(Wind) = Weak, Strong

S = [9+, 5-]

S_{Weak} ← [6+, 2-]

S_{Strong} ← [3+, 3-]

$$\begin{aligned}
 \text{Gain}(S, \text{Wind}) &= \text{Entropy}(S) - \sum_{v \in \{\text{Weak}, \text{Strong}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v) \\
 &= \text{Entropy}(S) - (8/14)\text{Entropy}(S_{\text{Weak}}) \\
 &\quad - (6/14)\text{Entropy}(S_{\text{Strong}}) \\
 &= 0.940 - (8/14)0.811 - (6/14)1.00 \\
 &= 0.048
 \end{aligned}$$

Instance	a1	a2	a3	Classification
1	True	Hot	High	No
2	True	Hot	High	No
3	False	Hot	High	Yes
4	False	Cool	Normal	Yes
5	False	Cool	Normal	Yes
6	True	Cool	High	No
7	True	Hot	High	No
8	True	Hot	Normal	Yes
9	False	Cool	Normal	Yes
10	False	Cool	High	Yes

Attribute: a1**Values(a1) = True, False**

$$S = [6+, 4-]$$

$$\text{Entropy}(S) = -\frac{6}{10} \log_2 \frac{6}{10} - \frac{4}{10} \log_2 \frac{4}{10} = 0.9709$$

$$S_{True} = [1+, 4-]$$

$$\text{Entropy}(S_{True}) = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} = 0.7219$$

$$S_{False} \leftarrow [5+, 0-]$$

$$\text{Entropy}(S_{False}) = 0.0$$

$$\text{Gain}(S, a1) = \text{Entropy}(S) - \sum_{v \in \{\text{True}, \text{False}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Gain}(S, a1) = \text{Entropy}(S) - \frac{5}{10} \text{Entropy}(S_{True}) - \frac{5}{10} \text{Entropy}(S_{False})$$

$$\text{Gain}(S, a1) = 0.9709 - \frac{5}{10} * 0.7219 - \frac{5}{10} * 0.0 = 0.6099$$

Attribute: a2**Values (a2) = Hot, Cool**

$$S = [6+, 4-]$$

$$\text{Entropy}(S) = -\frac{6}{10} \log_2 \frac{6}{10} - \frac{4}{10} \log_2 \frac{4}{10} = 0.9709$$

$$S_{Hot} = [2+, 3-]$$

$$\text{Entropy}(S_{Hot}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.9709$$

$$S_{Cool} \leftarrow [4+, 1-]$$

$$\text{Entropy}(S_{Cool}) = -\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} = 0.7219$$

$$\text{Gain}(S, a2) = \text{Entropy}(S) - \sum_{v \in \{Hot, Cool\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Gain}(S, a2) = \text{Entropy}(S) - \frac{5}{10} \text{Entropy}(S_{Hot}) - \frac{5}{10} \text{Entropy}(S_{Cool})$$

$$\text{Gain}(S, a2) = 0.9709 - \frac{5}{10} * 0.9709 - \frac{5}{10} * 0.7219 = 0.1245$$

Values(a3) = High, Normal

$$S = [6+, 4-] \quad Entropy(S) = -\frac{6}{10} \log_2 \frac{6}{10} - \frac{4}{10} \log_2 \frac{4}{10} = 0.9709$$

$$S_{High} = [2+, 4-] \quad Entropy(S_{High}) = -\frac{2}{6} \log_2 \frac{2}{6} - \frac{4}{6} \log_2 \frac{4}{6} = 0.9183$$

$$S_{Normal} \leftarrow [4+, 0-] \quad Entropy(S_{Normal}) = 0.0$$

$$Gain(S, a3) = Entropy(S) - \sum_{v \in \{High, Normal\}} \frac{|S_v|}{|S|} Entropy(S_v)$$

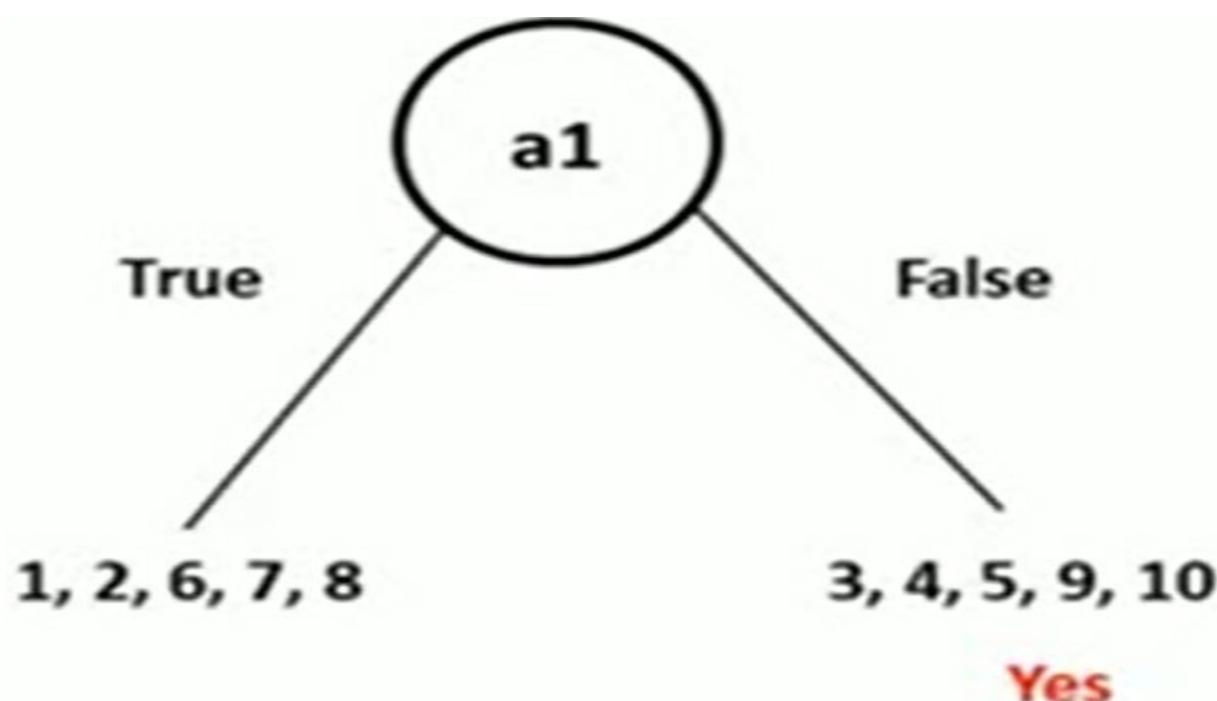
$$Gain(S, a3) = Entropy(S) - \frac{6}{10} Entropy(S_{High}) - \frac{4}{10} Entropy(S_{Normal})$$

$$Gain(S, a3) = 0.9709 - \frac{6}{10} * 0.9183 - \frac{4}{10} * 0.0 = 0.4199$$

$Gain(S, a1) = 0.6099$ – Maximum Gain

$Gain(S, a2) = 0.1245$

$Gain(S, a3) = 0.4199$



Instance	a2	a3	Classification
1	Hot	High	No
2	Hot	High	No
6	Cool	High	No
7	Hot	High	No
8	Hot	Normal	Yes

Attribute: a2

Values (a2) = Hot, Cool

$$S_{a1} = [1+, 4-]$$

$$\text{Entropy}(S_{a1}) = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} = 0.7219$$

$$S_{Hot} = [1+, 3-]$$

$$\text{Entropy}(S_{Hot}) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.8112$$

$$S_{Cool} \leftarrow [0+, 1-]$$

$$\text{Entropy}(S_{Cool}) = 0.0$$

$$\text{Gain}(S, a2) = \text{Entropy}(S) - \sum_{v \in \{Hot, Cool\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Gain}(S, a2) = \text{Entropy}(S) - \frac{4}{5} \text{Entropy}(S_{Hot}) - \frac{1}{5} \text{Entropy}(S_{Cool})$$

$$\text{Gain}(S, a2) = 0.9709 - \frac{4}{5} * 0.8112 - \frac{1}{5} * 0.0 = 0.3219$$

Attribute: a3

Values (a3) = High, Normal

$$S_{a1} = [1+, 4-] \quad \text{Entropy}(S_{a1}) = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} = 0.7219$$

$$S_{High} = [0+, 4-] \quad \text{Entropy}(S_{High}) = 0.0$$

$$S_{Normal} \leftarrow [1+, 0-] \quad \text{Entropy}(S_{Normal}) = 0.0$$

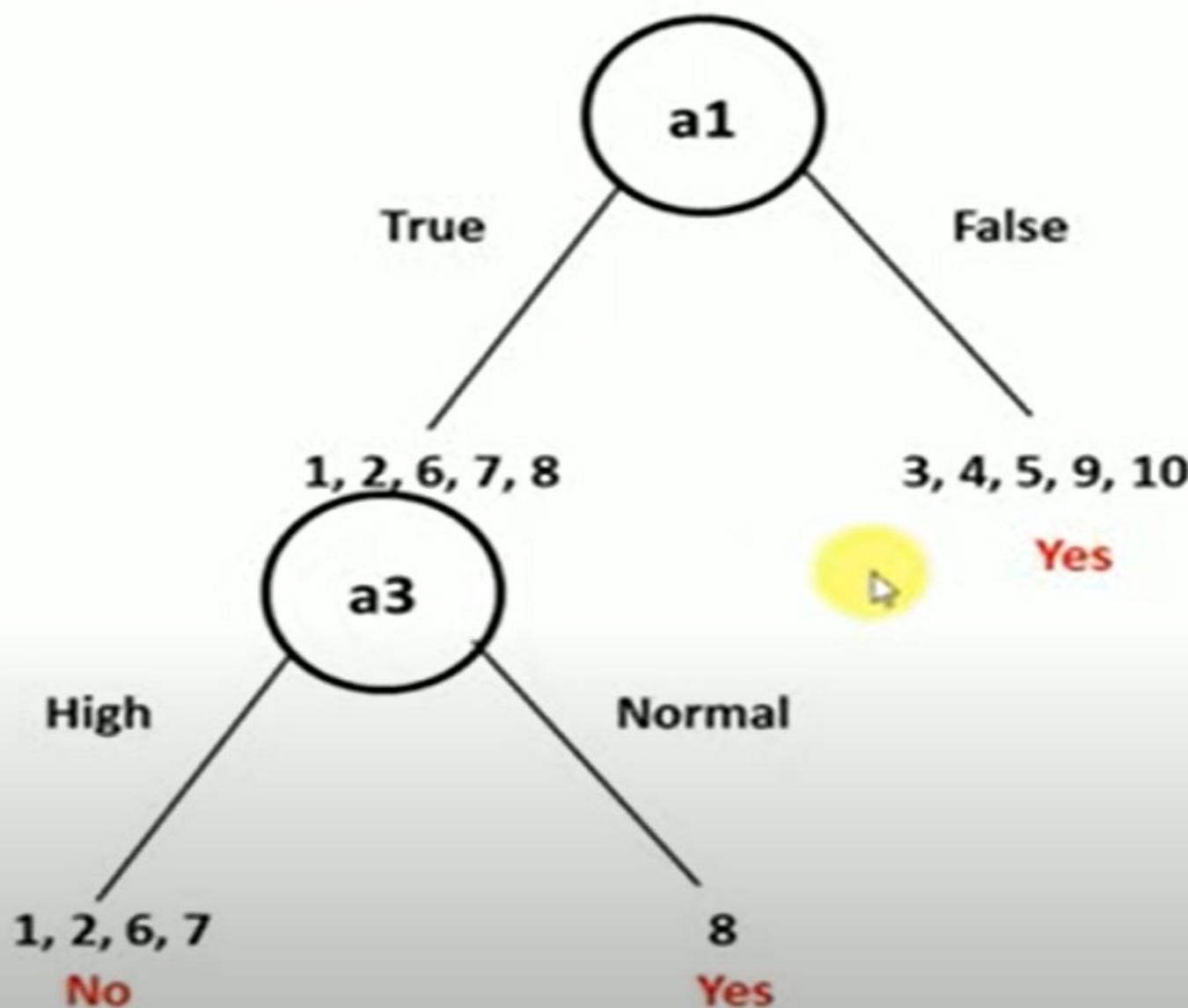
$$\text{Gain}(S, a3) = \text{Entropy}(S) - \sum_{v \in \{High, Normal\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Gain}(S, a3) = \text{Entropy}(S) - \frac{4}{5} \text{Entropy}(S_{High}) - \frac{1}{5} \text{Entropy}(S_{Normal})$$

$$\text{Gain}(S, a3) = 0.9709 - \frac{4}{5} * 0.0 - \frac{1}{5} * 0.0 = 0.7219$$

$$Gain(S_{a1}, a2) = 0.3219$$

$$Gain(S_{a1}, a3) = 0.7219 - \text{Maximum Gain}$$



Gini Index

- If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^n p_j^2$$

where p_j is the relative frequency of class j in D

- If a data set D is split on A into two subsets D_1 and D_2 , the *gini* index $gini(D)$ is defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

- Reduction in Impurity:

$$\Delta gini(A) = gini(D) - gini_A(D)$$

- The attribute provides the smallest $gini_{split}(D)$ (or the largest reduction in impurity) is chosen to split the node (*need to enumerate all the possible splitting points for each attribute*)

Weekend	Weather	Parents	Money	Decision
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay In
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

Overfitting and Tree Pruning

- Overfitting: An induced tree may overfit the training data
 - Too many branches, some may reflect anomalies due to noise or outliers
 - Poor accuracy for unseen samples
- Two approaches to avoid overfitting
 - Prepruning: *Halt tree construction early*-do not split a node if this would result in the goodness measure falling below a threshold
 - Difficult to choose an appropriate threshold
 - Postpruning: *Remove branches* from a “fully grown” tree—get a sequence of progressively pruned trees
 - Use a set of data different from the training data to decide which is the “best pruned tree”

Rainforest: Training Set and Its AVC Sets

Training Examples

age	income	student	credit rating	computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

AVC-set on Age

Age	Buy_Computer	
	yes	no
<=30	2	3
31..40	4	0
>40	3	2

AVC-set on *income*

income	Buy_Computer	
	yes	no
high	2	2
medium	4	2
low	3	1

AVC-set on *Student*

student	Buy_Computer	
	yes	no
yes	6	1
no	3	4

AVC-set on *credit_rating*

Credit rating	Buy_Computer	
	yes	no
fair	6	2
excellent	3	3

Bayesian Classification: Why?

- A statistical classifier: performs *probabilistic prediction*, i.e., predicts class membership probabilities
- Foundation: Based on Bayes' Theorem.
- Performance: A simple Bayesian classifier, *naïve Bayesian classifier*, has comparable performance with decision tree and selected neural network classifiers
- Incremental: Each training example can incrementally increase/decrease the probability that a hypothesis is correct — prior knowledge can be combined with observed data
- Standard: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured

Bayes' Theorem: Basics

- Total probability Theorem: $P(B) = \sum_{i=1}^M P(B|A_i)P(A_i)$
- Bayes' Theorem: $P(H|\mathbf{X}) = \frac{P(\mathbf{X}|H)P(H)}{P(\mathbf{X})} = P(\mathbf{X}|H) \times P(H) / P(\mathbf{X})$
 - Let \mathbf{X} be a data sample (“*evidence*”): class label is unknown
 - Let H be a *hypothesis* that \mathbf{X} belongs to class C
 - Classification is to determine $P(H|\mathbf{X})$, (i.e., *posteriori probability*): the probability that the hypothesis holds given the observed data sample \mathbf{X}
 - $P(H)$ (*prior probability*): the initial probability
 - E.g., \mathbf{X} will buy computer, regardless of age, income, ...
 - $P(\mathbf{X})$: probability that sample data is observed
 - $P(\mathbf{X}|H)$ (*likelihood*): the probability of observing the sample \mathbf{X} , given that the hypothesis holds
 - E.g., Given that \mathbf{X} will buy computer, the prob. that X is 31..40, medium income

Prediction Based on Bayes' Theorem

- Given training data \mathbf{X} , *posteriori probability of a hypothesis H*, $P(H|\mathbf{X})$, follows the Bayes' theorem

$$P(H|\mathbf{X}) = \frac{P(\mathbf{X}|H)P(H)}{P(\mathbf{X})} = P(\mathbf{X}|H) \times P(H) / P(\mathbf{X})$$

- Informally, this can be viewed as
posteriori = likelihood x prior/evidence
- Predicts \mathbf{X} belongs to C_i iff the probability $P(C_i|\mathbf{X})$ is the highest among all the $P(C_k|\mathbf{X})$ for all the k classes
- Practical difficulty: It requires initial knowledge of many probabilities, involving significant computational cost

Bayes' Theorem

$$v_{NB} = \operatorname{argmax}_{v_j \in \{\text{yes}, \text{no}\}} P(v_j) \prod_i P(a_i | v_j)$$

Exercise 4. Naïve Bayes for data with nominal attributes

Given the training data in the table below (*Buy Computer* data), predict the class of the following new example using Naïve Bayes classification: age<=30, income=medium, student=yes, credit-rating=fair

RID	age	income	student	credit_rating	Class: buys_computer
1	<=30	high	no	fair	no
2	<=30	high	no	excellent	no
3	31 ... 40	high	no	fair	yes
4	>40	medium	no	fair	yes
5	>40	low	yes	fair	yes
6	>40	low	yes	excellent	no
7	31 ... 40	low	yes	excellent	yes
8	<=30	medium	no	fair	no
9	<=30	low	yes	fair	yes
10	>40	medium	yes	fair	yes
11	<=30	medium	yes	excellent	yes
12	31 ... 40	medium	no	excellent	yes
13	31 ... 40	high	yes	fair	yes
14	>40	medium	no	excellent	no

Solution:

$E = \text{age} \leq 30, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair}$

E_1 is $\text{age} \leq 30$, E_2 is $\text{income} = \text{medium}$, E_3 is $\text{student} = \text{yes}$, E_4 is $\text{credit-rating} = \text{fair}$

We need to compute $P(\text{yes}|E)$ and $P(\text{no}|E)$ and compare them.

$$P(\text{yes} | E) = \frac{P(E_1 | \text{yes}) P(E_2 | \text{yes}) P(E_3 | \text{yes}) P(E_4 | \text{yes}) P(\text{yes})}{P(E)}$$

$$P(\text{yes}) = 9/14 = 0.643$$

$$P(\text{no}) = 5/14 = 0.357$$

$$P(E_1|\text{yes}) = 2/9 = 0.222$$

$$P(E_1|\text{no}) = 3/5 = 0.6$$

$$P(E_2|\text{yes}) = 4/9 = 0.444$$

$$P(E_2|\text{no}) = 2/5 = 0.4$$

$$P(E_3|\text{yes}) = 6/9 = 0.667$$

$$P(E_3|\text{no}) = 1/5 = 0.2$$

$$P(E_4|\text{yes}) = 6/9 = 0.667$$

$$P(E_4|\text{no}) = 2/5 = 0.4$$

$$P(\text{yes} | E) = \frac{0.222 \ 0.444 \ 0.667 \ 0.668 \ 0.443}{P(E)} = \frac{0.028}{P(E)}$$

$$P(\text{no} | E) = \frac{0.6 \ 0.4 \ 0.2 \ 0.4 \ 0.357}{P(E)} = \frac{0.007}{P(E)}$$

Hence, the Naïve Bayes classifier predicts `buys_computer=yes` for the new example.

Classification Is to Derive the Maximum Posteriori

- Let D be a training set of tuples and their associated class labels, and each tuple is represented by an n -D attribute vector $\mathbf{X} = (x_1, x_2, \dots, x_n)$
- Suppose there are m classes C_1, C_2, \dots, C_m .
- Classification is to derive the maximum posteriori, i.e., the maximal $P(C_i | \mathbf{X})$
- This can be derived from Bayes' theorem

$$P(C_i | \mathbf{X}) = \frac{P(\mathbf{X} | C_i)P(C_i)}{P(\mathbf{X})}$$

- Since $P(\mathbf{X})$ is constant for all classes, only

$$P(C_i | \mathbf{X}) = P(\mathbf{X} | C_i)P(C_i)$$

needs to be maximized

Naïve Bayes Classifier: Training Dataset

Class:

C1:buys_computer = 'yes'

C2:buys_computer = 'no'

Data to be classified:

X = (age <=30,

Income = medium,

Student = yes

Credit_rating = Fair)

age	income	student	credit rating	com
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Naïve Bayes Classifier: An Example

- $P(C_i)$: $P(\text{buys_computer} = \text{"yes"}) = 9/14 = 0.643$
 $P(\text{buys_computer} = \text{"no"}) = 5/14 = 0.357$
- Compute $P(X|C_i)$ for each class
 - $P(\text{age} = \text{"<=30"} | \text{buys_computer} = \text{"yes"}) = 2/9 = 0.222$
 - $P(\text{age} = \text{"<= 30"} | \text{buys_computer} = \text{"no"}) = 3/5 = 0.6$
 - $P(\text{income} = \text{"medium"} | \text{buys_computer} = \text{"yes"}) = 4/9 = 0.444$
 - $P(\text{income} = \text{"medium"} | \text{buys_computer} = \text{"no"}) = 2/5 = 0.4$
 - $P(\text{student} = \text{"yes"} | \text{buys_computer} = \text{"yes"}) = 6/9 = 0.667$
 - $P(\text{student} = \text{"yes"} | \text{buys_computer} = \text{"no"}) = 1/5 = 0.2$
 - $P(\text{credit_rating} = \text{"fair"} | \text{buys_computer} = \text{"yes"}) = 6/9 = 0.667$
 - $P(\text{credit_rating} = \text{"fair"} | \text{buys_computer} = \text{"no"}) = 2/5 = 0.4$
- $X = (\text{age} \leq 30, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit_rating} = \text{fair})$
- $P(X|C_i)$: $P(X|\text{buys_computer} = \text{"yes"}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$
 $P(X|\text{buys_computer} = \text{"no"}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$
- $P(X|C_i) * P(C_i)$: $P(X|\text{buys_computer} = \text{"yes"}) * P(\text{buys_computer} = \text{"yes"}) = 0.028$
 $P(X|\text{buys_computer} = \text{"no"}) * P(\text{buys_computer} = \text{"no"}) = 0.007$
- Therefore, X belongs to class ("buys_computer = yes")

age	income	student	credit rating	com
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

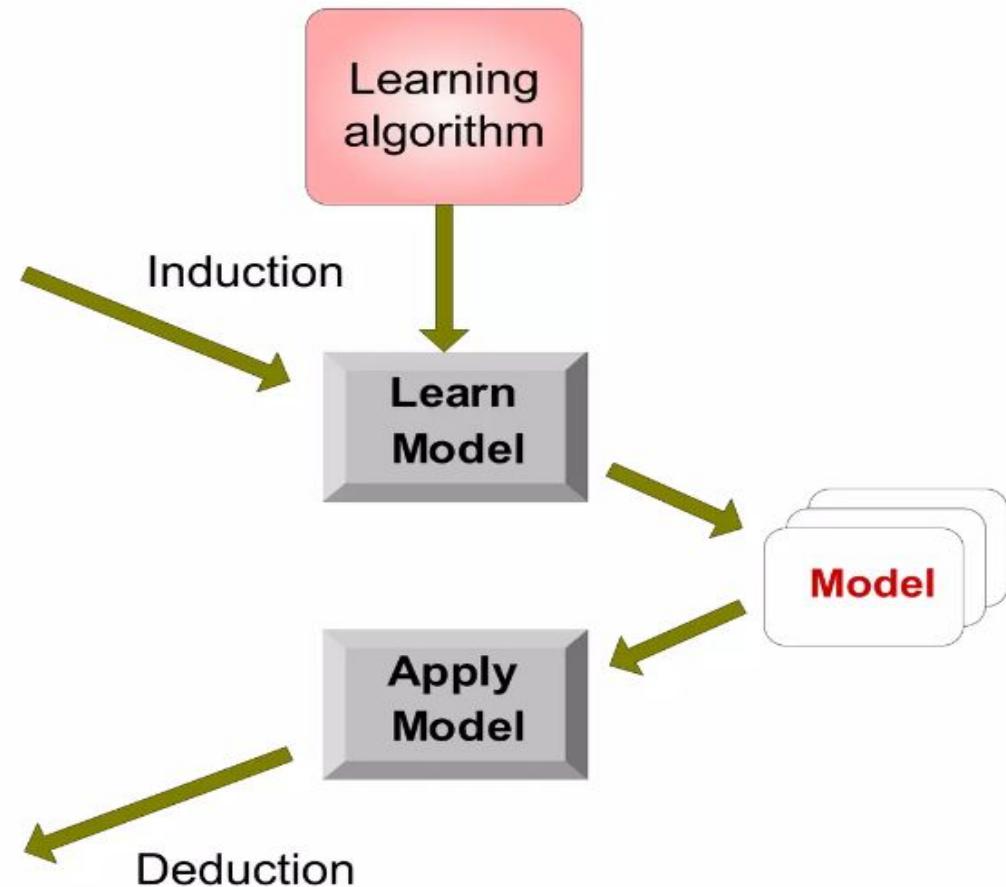
Rule-Based Classification

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



Rule-Based Classification

- A rule-based classifier uses a set of IF-THEN rules for classification.
- An IF-THEN rule is an expression of the form

IF *condition* THEN *conclusion*.

An example is rule *R1*,

R1: IF *age = youth* AND *student = yes* THEN *buys_computer = yes*.

LHS:- rule antecedent or precondition.

RHS:- rule consequent.

- If the condition (i.e., all the attribute tests) in a rule antecedent holds true for a given tuple, we say that the rule antecedent is satisfied (or simply, that the rule is satisfied)
- If the condition in antecedent is available then that the rule **covers** the tuple

Tid	Refund	Marital Status	Taxable Income	Class
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

(Status=Single) → No

Coverage = 40%, Accuracy = 50%

-
- Quality of rule is decided by **coverage and accuracy**
 - A rule R can be assessed by its coverage and accuracy
 - Given a tuple, X, from a class labeled data set, D,
 - let **ncovers** be the number of tuples covered by R;
 - **ncorrect** be the number of tuples correctly classified by R;
 - and $|D|$ be the number of tuples in D.
 - We can define the coverage and accuracy of R as

$$\text{coverage}(R) = \frac{n_{\text{covers}}}{|D|}$$

$$\text{accuracy}(R) = \frac{n_{\text{correct}}}{n_{\text{covers}}}.$$

How rule-based classifier Work

R1: (Give Birth = no) \wedge (Can Fly = yes) \rightarrow Birds

R2: (Give Birth = no) \wedge (Live in Water = yes) \rightarrow Fishes

R3: (Give Birth = yes) \wedge (Blood Type = warm) \rightarrow Mammals

R4: (Give Birth = no) \wedge (Can Fly = no) \rightarrow Reptiles

R5: (Live in Water = sometimes) \rightarrow Amphibians

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
lemur	warm	yes	no	no	?
turtle	cold	no	no	sometimes	?
dogfish shark	cold	yes	no	yes	?

R1: (Give Birth = no) \wedge (Can Fly = yes) \rightarrow Birds

R2: (Give Birth = no) \wedge (Live in Water = yes) \rightarrow Fishes

R3: (Give Birth = yes) \wedge (Blood Type = warm) \rightarrow Mammals

R4: (Give Birth = no) \wedge (Can Fly = no) \rightarrow Reptiles

R5: (Live in Water = sometimes) \rightarrow Amphibians

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
lemur	warm	yes	no	no	?
turtle	cold	no	no	sometimes	?
dogfish shark	cold	yes	no	yes	?

A lemur triggers rule R3, so it is classified as a mammal

A turtle triggers both R4 and R5

A dogfish shark triggers none of the rules

Solution to first conflict(tuple satisfy both rule)

- we need a conflict resolution strategy to figure out which rule gets to fire and assign its class prediction
- There are many possible strategies.
 - Size ordering and
 - Rule ordering

Size Ordering

- The size ordering scheme assigns the highest priority to the triggering rule that has the “toughest” requirements,
- Where toughness is measured by the rule antecedent size.
- That is, the triggering rule with the most attribute tests is fired

Rule Ordering

- scheme prioritizes the rules beforehand
- The ordering may be class-based or rule-based

Rule-Based Ordering

- Individual rules are ranked based on quality of rule
- Rules are unordered

Class-Based Ordering

- Rules that belong to the same class appear together
- Rules are ordered

Solution to second conflict(tuple satisfy no rule)

- Default rule can be set up to specify a default class, based on a training set
- This may be the class in majority or the majority class of the tuples that were not covered by any rule.
- The default rule is evaluated at the end, if and only if no other rule covers X.
- The condition in the default rule is empty. In this way, the rule fires when no other rule is satisfied.

Building Classification Rule

Rule Characteristic:-

Mutually exclusive :-

Mutually exclusive means that we cannot have rule conflicts here because no two rules will be triggered for the same tuple. (We have one rule per leaf, and any tuple can map to only one leaf.)

Exhaustive:-

means there is one rule for each possible attribute–value combination, so that this set of rules does not require a default rule.

Direct Method

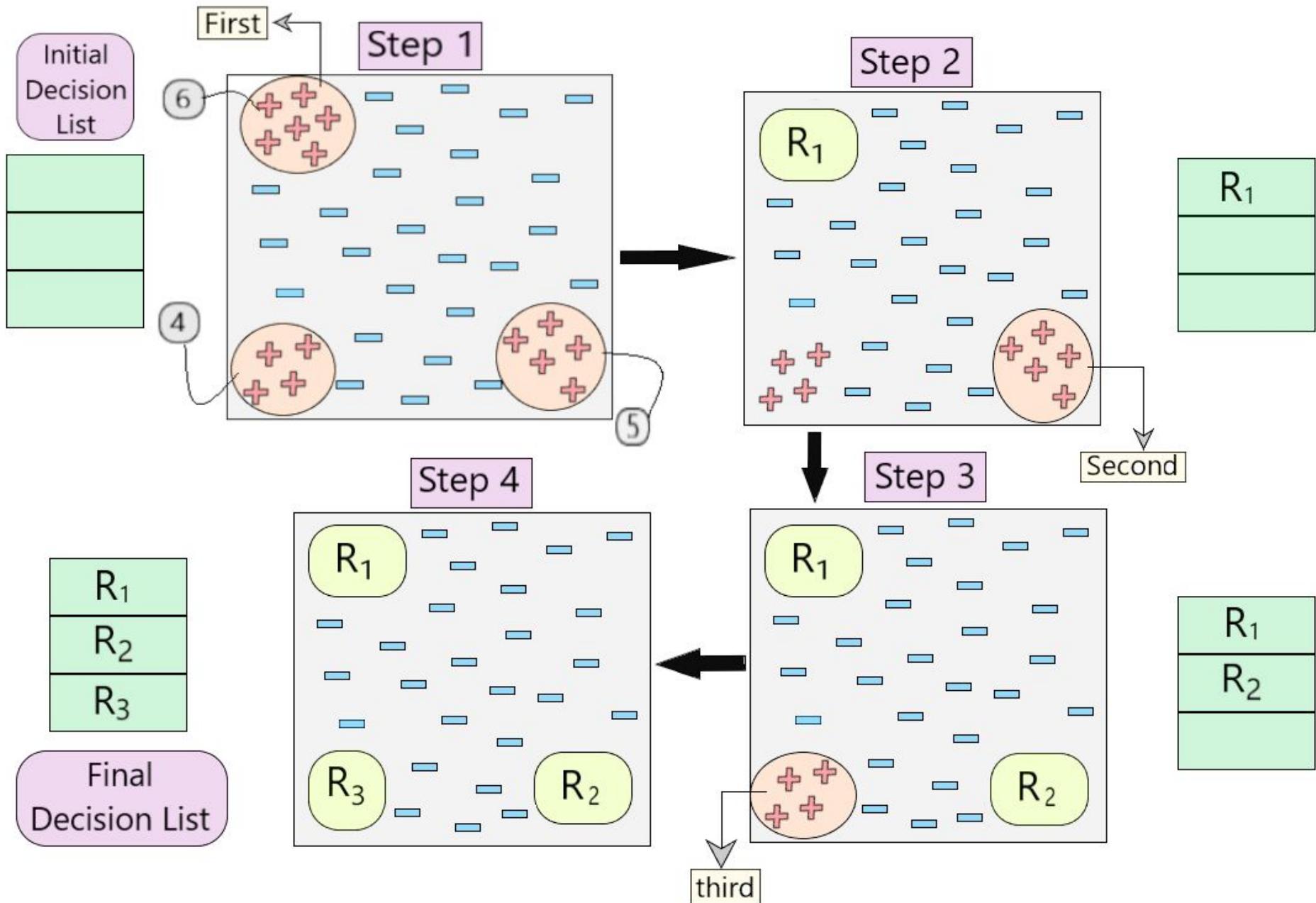
- Rules are extracted from data
- Example:- Sequential Covering Algorithm

Indirect Method

- Extract rules from other classification Method
- Example:-Decision Tree, Neural Network

Sequential Covering Algorithm

- IF-THEN rules can be extracted directly from the training data
- The name comes from the notion that the rules are learned sequentially
- Rules are learned one at a time.
- Each time a rule is learned, the tuples covered by the rule are removed, and the process repeats on the remaining tuples



Algorithm: Sequential covering. Learn a set of IF-THEN rules for classification.

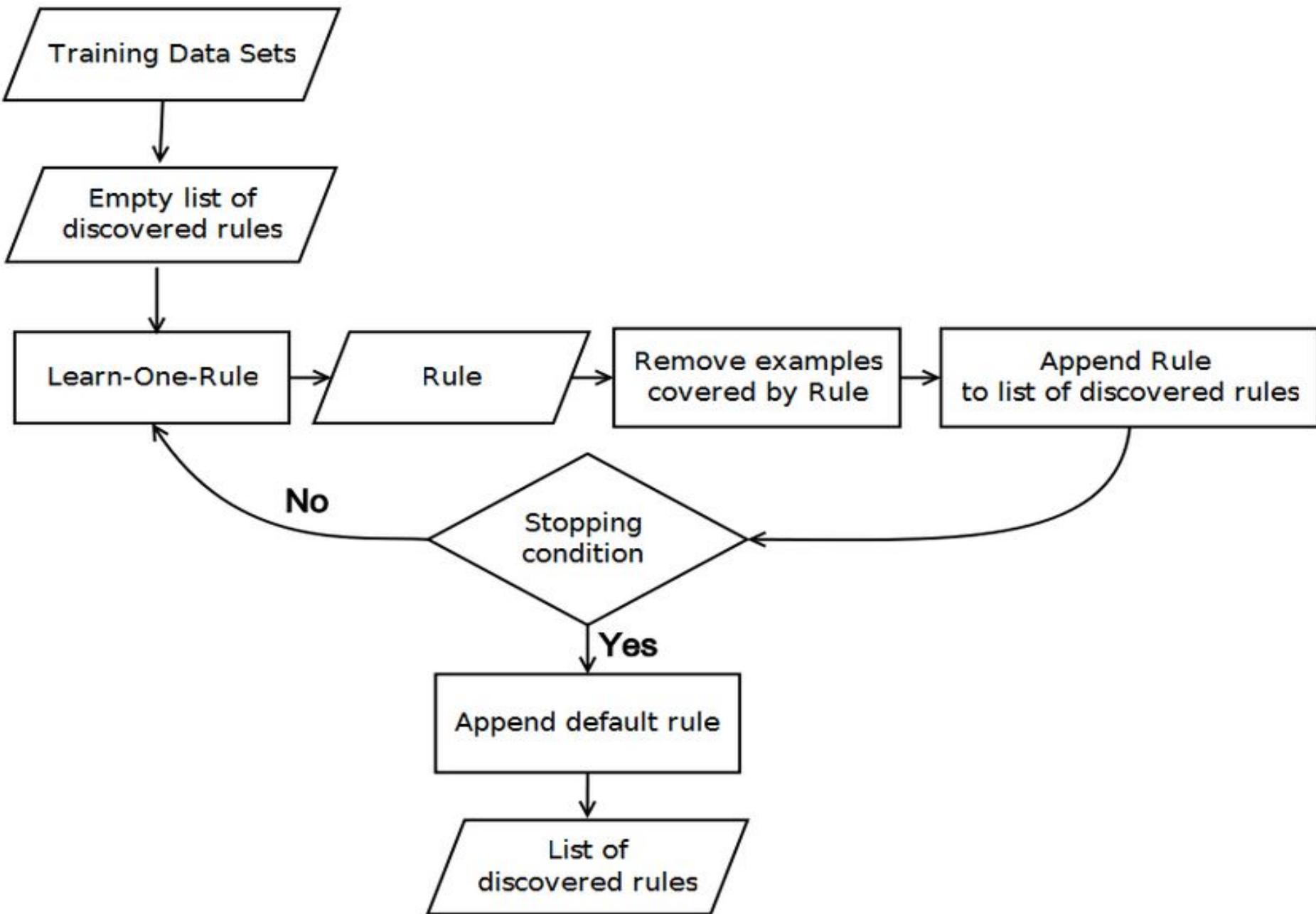
Input:

- D , a data set of class-labeled tuples;
- Att_vals , the set of all attributes and their possible values.

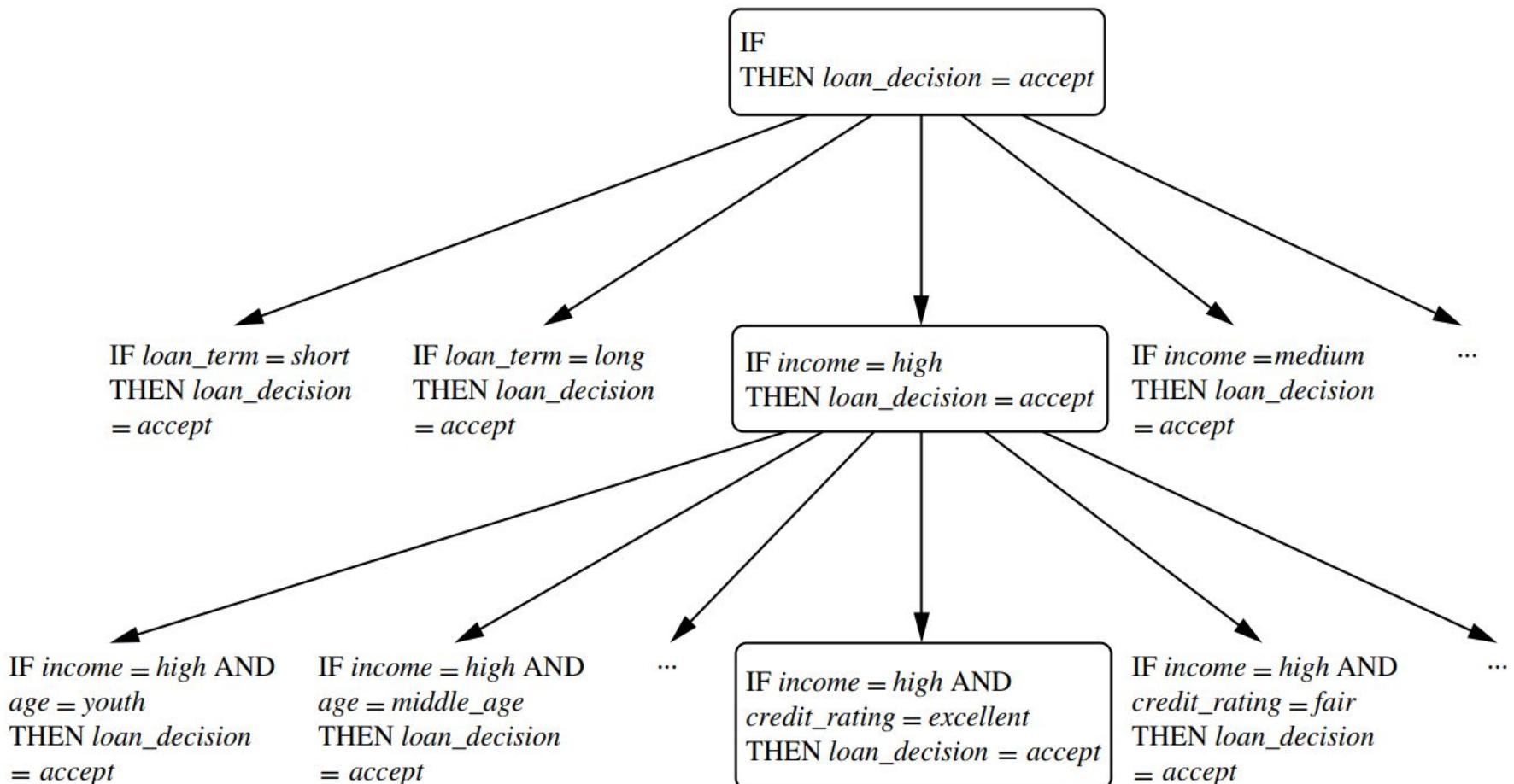
Output: A set of IF-THEN rules.

Method:

- (1) $Rule_set = \{\}$; // initial set of rules learned is empty
- (2) **for each** class c **do**
- (3) **repeat**
- (4) $Rule = \text{Learn_One_Rule}(D, Att_vals, c);$
- (5) remove tuples covered by $Rule$ from D ;
- (6) $Rule_set = Rule_set + Rule;$ // add new rule to rule set
- (7) **until** terminating condition;
- (8) **endfor**
- (9) return $Rule_Set;$



Rule Growing(How are rule learn?)



- Learn_One_Rule() needs a measure of rule quality.
- Accuracy may seem like an obvious choice at first, but consider Example
- Thus, R2 has greater accuracy than R1, but it is not the better rule because of its small coverage.

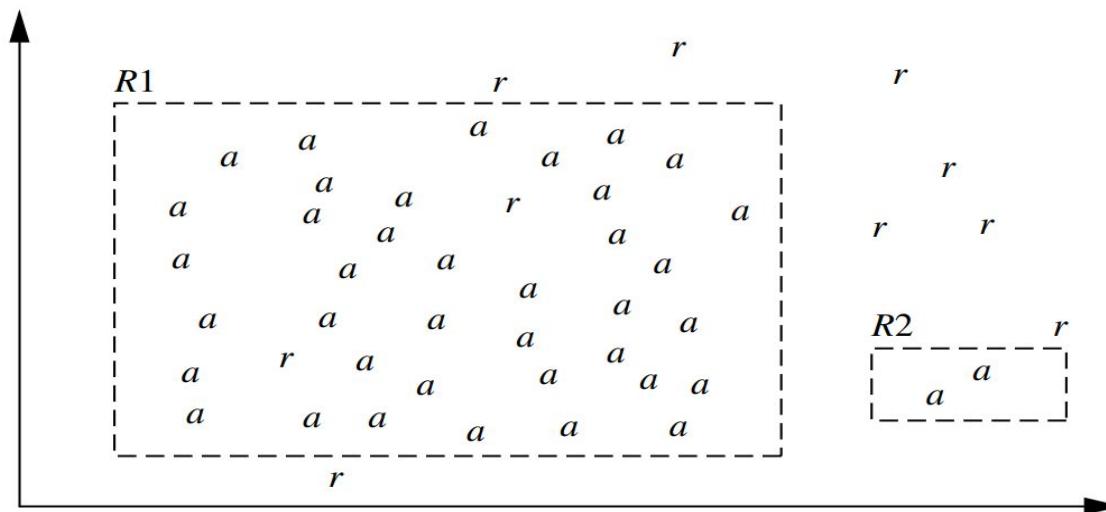
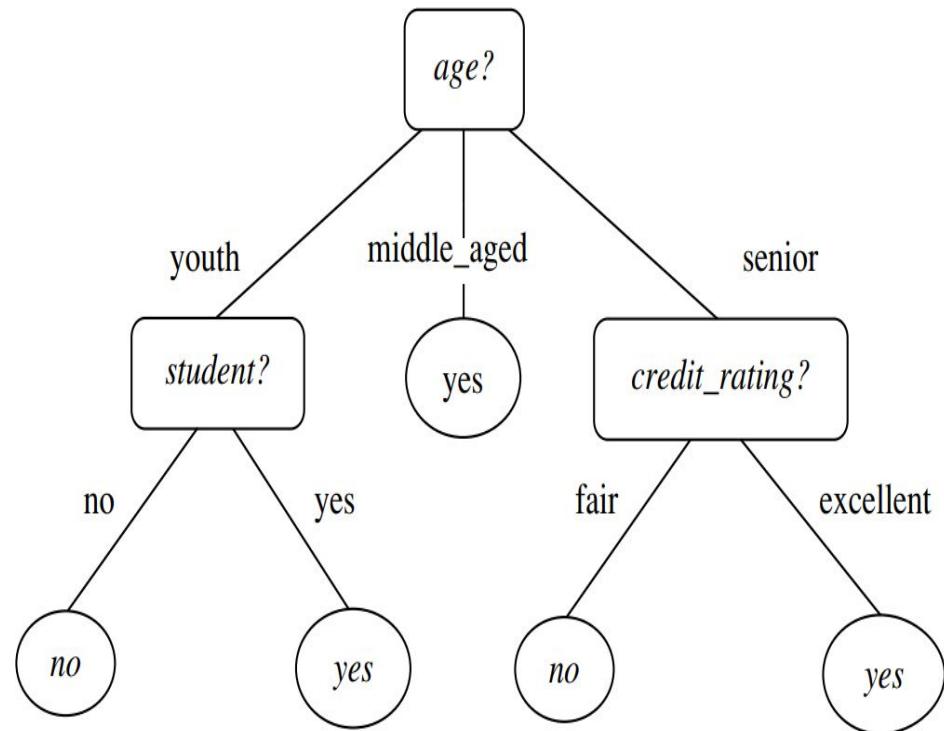


Figure 8.12 Rules for the class *loan_decision = accept*, showing *accept* (*a*) and *reject* (*r*) tuples.

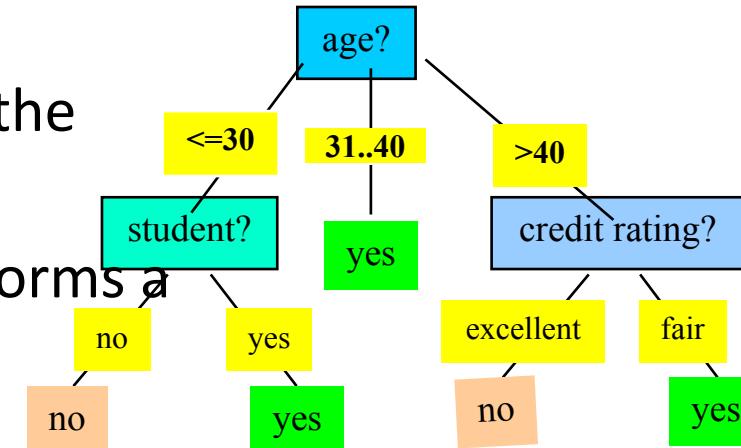
Indirect Method



- | | | |
|---------------------------------|--------------------------------------|---------------------------------|
| <i>R1: IF age = youth</i> | <i>AND student = no</i> | <i>THEN buys_computer = no</i> |
| <i>R2: IF age = youth</i> | <i>AND student = yes</i> | <i>THEN buys_computer = yes</i> |
| <i>R3: IF age = middle_aged</i> | | <i>THEN buys_computer = yes</i> |
| <i>R4: IF age = senior</i> | <i>AND credit_rating = excellent</i> | <i>THEN buys_computer = yes</i> |
| <i>R5: IF age = senior</i> | <i>AND credit_rating = fair</i> | <i>THEN buys_computer = no</i> |

Rule Extraction from a Decision Tree

- Rules are *easier to understand* than large trees
- One rule is created *for each path* from the root to a leaf
- Each attribute-value pair along a path forms a conjunction: the leaf holds the class prediction
- Rules are mutually exclusive and exhaustive
- Example: Rule extraction from our *buys_computer* decision-tree



IF <i>age</i> = young AND <i>student</i> = no	THEN <i>buys_computer</i> = no
IF <i>age</i> = young AND <i>student</i> = yes	THEN <i>buys_computer</i> = yes
IF <i>age</i> = mid-age	THEN <i>buys_computer</i> = yes
IF <i>age</i> = old AND <i>credit_rating</i> = excellent	THEN <i>buys_computer</i> = no
IF <i>age</i> = old AND <i>credit_rating</i> = fair	THEN <i>buys_computer</i> = yes

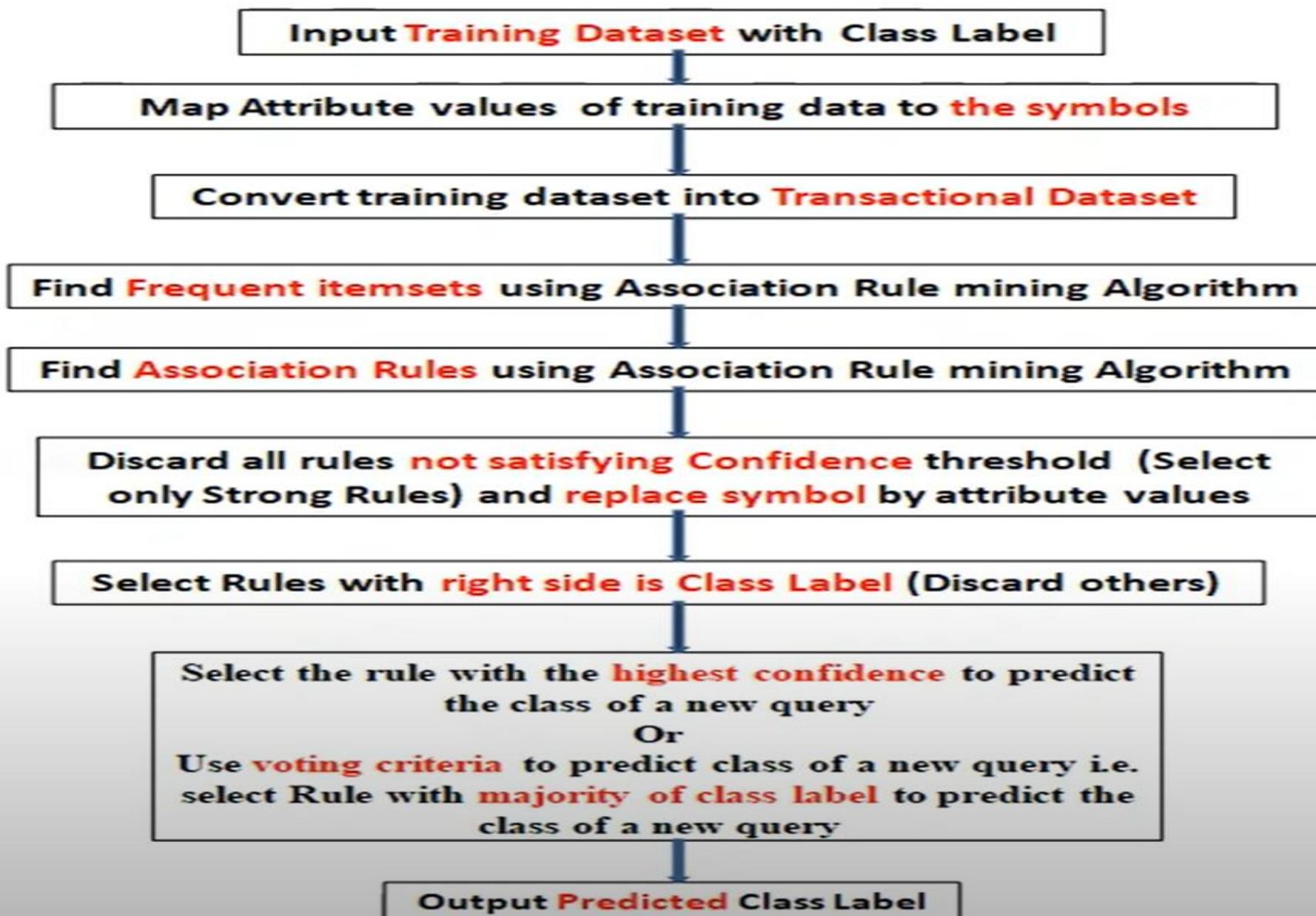
Associative Classification(ARM+Classification)

- Association rules are mined in a two-step process
 - The first step searches for patterns of attribute–value pairs that occur repeatedly in a data set(also referred to as frequent patterns).
 - The second step analyzes the frequent itemset to generate association rules.
 - All association rules must satisfy certain criteria regarding their “accuracy” (or confidence) and the proportion of the data set that they actually represent (referred to as support).
 - For example, the following is an association rule mined from a data set, D, shown with its confidence and support:

-
- From a classification point of view, this is akin to rule accuracy.
 - **For example**
 - A confidence of 93% for Rule means that 93% of the customers in D who are young and have an OK credit rating belong to the class buys computer = yes.
 - The percentage of tuples in D satisfying the rule antecedent and having class label C is called the support of R.
 - A support of 20% for Rule above means that 20% of the customers in D are young, have an OK credit rating

$age = youth \wedge credit = OK \Rightarrow buys_computer$
= yes [support = 20%, confidence = 93%],

**Steps for Associative Classifier
(Association Rule mining + Classification)**



Weekend	Weather	Parents	Money	Decision
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay In
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

Time	Weather	Company	Goes
Morning	Sunny	Yes	Yes
Evening	Rainy	No	No
Morning	Sunny	Yes	Yes
Evening	Sunny	Yes	Yes

Instance	a2	a3	Classification
1	Hot	High	No
2	Hot	High	No
6	Cool	High	No
7	Hot	High	No
8	Hot	Normal	Yes

Classifier Evaluation Metrics: Confusion Matrix

Confusion Matrix:

Actual class\Predicted class	C_1	$\neg C_1$
C_1	True Positives (TP)	False Negatives (FN)
$\neg C_1$	False Positives (FP)	True Negatives (TN)

Example of Confusion Matrix:

Actual class\Predicted class	buy_computer = yes	buy_computer = no	Total
buy_computer = yes	6954	46	7000
buy_computer = no	412	2588	3000
Total	7366	2634	10000

- Given m classes, an entry, $CM_{i,j}$, in a **confusion matrix** indicates # of tuples in class i that were labeled by the classifier as class j
- May have extra rows/columns to provide totals

Classifier Evaluation Metrics: Accuracy, Error Rate, Sensitivity and Specificity

A\P	C	$\neg C$	
C	TP	FN	P
$\neg C$	FP	TN	N
	P'	N'	All

- **Classifier Accuracy**, or recognition rate: percentage of test set tuples that are correctly classified

$$\text{Accuracy} = (\text{TP} + \text{TN})/\text{All}$$

- **Error rate**: $1 - \text{accuracy}$, or
$$\text{Error rate} = (\text{FP} + \text{FN})/\text{All}$$

- **Class Imbalance Problem**:
 - One class may be *rare*, e.g. fraud, or HIV-positive
 - Significant *majority of the negative class* and minority of the positive class
- **Sensitivity**: True Positive recognition rate
 - $$\text{Sensitivity} = \text{TP}/\text{P}$$
- **Specificity**: True Negative recognition rate
 - $$\text{Specificity} = \text{TN}/\text{N}$$

Classifier Evaluation Metrics: Precision and Recall, and F-measures

- **Precision:** exactness – what % of tuples that the classifier labeled as positive are actually positive

$$precision = \frac{TP}{TP + FP}$$

- **Recall:** completeness – what % of positive tuples did the classifier label as positive?

$$recall = \frac{TP}{TP + FN}$$

- Perfect score is 1.0
- Inverse relationship between precision & recall
- **F measure (F_1 or F-score):** harmonic mean of precision and recall,

$$F = \frac{2 \times precision \times recall}{precision + recall}$$

- **F_β :** weighted measure of precision and recall
 - assigns β times as much weight to recall as to precision

$$F_\beta = \frac{(1 + \beta^2) \times precision \times recall}{\beta^2 \times precision + recall}$$

Classifier Evaluation Metrics: Example

Actual Class\Predicted class	cancer = yes	cancer = no	Total	Recognition(%)
cancer = yes	90	210	300	30.00 (<i>sensitivity</i>)
cancer = no	140	9560	9700	98.56 (<i>specificity</i>)
Total	230	9770	10000	96.40 (<i>accuracy</i>)

- $Precision = 90/230 = 39.13\%$ $Recall = 90/300 = 30.00\%$