



Classification

Dr. N. Abdolvand

Data Mining

Data Mining: Concepts and Techniques

(3rd ed.)

— Chapter 8 —

Jiawei Han, Micheline Kamber, and Jian Pei
University of Illinois at Urbana-Champaign &
Simon Fraser University

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



Classification

Dr. N. Abdolvand

Data Mining

Data Mining: Concepts and Techniques

(3rd ed.)

— Chapter 8 —

Jiawei Han, Micheline Kamber, and Jian Pei
University of Illinois at Urbana-Champaign &
Simon Fraser University

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

Chapter 8. Classification: Basic Concepts

- ◆ Classification: Basic Concepts
- ◆ Decision Tree Induction
- ◆ Bayes Classification Methods
- ◆ Rule-Based Classification
- ◆ Model Evaluation and Selection
- ◆ Techniques to Improve Classification Accuracy: Ensemble Methods
- ◆ Summary

3

Case Study

- ◆ Gill runs a sports academy designed to help high school aged athletes achieve their maximum athletic potential. On the boys side of his academy, he focuses on four major sports: Football, Basketball, Baseball and Hockey. He has found that while many high school athletes enjoy participating in a number of sports in high school, as they begin to consider playing a sport at the college level, they would prefer to specialize in one sport. As he's worked with athletes over the years, Gill has developed an extensive data set, and he now is wondering if he can use past performance from some of his previous clients to predict prime sports for up-and-coming high school athletes. Ultimately, he hopes he can make a recommendation to each athlete as to the sport in which they should most likely choose to specialize. By evaluating each athlete's performance across a battery of test, Gill hopes we can help him figure out for which sport each athlete has the highest aptitude.



Dr. N. Abdolvand

Chapter 8. Classification: Basic Concepts

- ◆ Classification: Basic Concepts
- ◆ Decision Tree Induction
- ◆ Bayes Classification Methods
- ◆ Rule-Based Classification
- ◆ Model Evaluation and Selection
- ◆ Techniques to Improve Classification Accuracy: Ensemble Methods
- ◆ Summary

3

Case Study

- ◆ Gill runs a sports academy designed to help high school aged athletes achieve their maximum athletic potential. On the boys side of his academy, he focuses on four major sports: Football, Basketball, Baseball and Hockey. He has found that while many high school athletes enjoy participating in a number of sports in high school, as they begin to consider playing a sport at the college level, they would prefer to specialize in one sport. As he's worked with athletes over the years, Gill has developed an extensive data set, and he now is wondering if he can use past performance from some of his previous clients to predict prime sports for up-and-coming high school athletes. Ultimately, he hopes he can make a recommendation to each athlete as to the sport in which they should most likely choose to specialize. By evaluating each athlete's performance across a battery of test, Gill hopes we can help him figure out for which sport each athlete has the highest aptitude.



Dr. N. Abdolvand

ORGANIZATIONAL UNDERSTANDING

- ◆ Gill's objective is to examine young athletes and, based upon their performance across a number of metrics, help them decide which sport is the most prime for their specialized success. Gill recognizes that all of his clients possess some measure of athleticism, and that they enjoy participating in a number of sports. Being young, athletic, and adaptive, most of his clients are quite good at a number of sports, and he has seen over the years that some people are so naturally gifted that they would excel in any sport they choose for specialization. Thus, he recognizes, as a limitation of this data mining exercise, that he may not be able to use data to determine an athlete's "best" sport. Still, he has seen metrics and evaluations work in the past, and has seen that some of his previous athletes really were pre-disposed to a certain sport, and that they were successful as they went on to specialize in that sport. Based on his industry experience, he has decided to go ahead with an experiment in mining data for athletic aptitude, and has enlisted our help.



Dr. N. Abdolvand

Attributes

- ◆ **Age:** This is the age in years (one decimal precision for the part of the year since the client's last birthday) at the time that the athletic and personality trait battery test was administered. Participants ranged in age from 13-19 years old at the time they took the battery.
- ◆ **Strength:** This is the participant's strength measured through a series of weight lifting exercises and recorded on a scale of 0-10, with 0 being limited strength and 10 being sufficient strength to perform all lifts without any difficulty. No participant scored 8, 9 or 10, but some participants did score 0.
- ◆ **Quickness:** This is the participant's performance on a series of responsiveness tests. Participants were timed on how quickly they were able to press buttons when they were illuminated or to jump when a buzzer sounded. Their response times were tabulated on a scale of 0-6, with 6 being extremely quick response and 0 being very slow. Participants scored all along the spectrum for this attribute.
- ◆ **Injury:** This is a simple yes (1) / no (0) column indicating whether or not the young athlete had already suffered an athletic-related injury that was severe enough to require surgery or other major medical intervention. Common injuries treated with ice, rest, stretching, etc. were entered as 0. Injuries that took more than three week to heal, that required physical therapy or surgery were flagged as 1.



Dr. N. Abdolvand

ORGANIZATIONAL UNDERSTANDING

- ◆ Gill's objective is to examine young athletes and, based upon their performance across a number of metrics, help them decide which sport is the most prime for their specialized success. Gill recognizes that all of his clients possess some measure of athleticism, and that they enjoy participating in a number of sports. Being young, athletic, and adaptive, most of his clients are quite good at a number of sports, and he has seen over the years that some people are so naturally gifted that they would excel in any sport they choose for specialization. Thus, he recognizes, as a limitation of this data mining exercise, that he may not be able to use data to determine an athlete's "best" sport. Still, he has seen metrics and evaluations work in the past, and has seen that some of his previous athletes really were pre-disposed to a certain sport, and that they were successful as they went on to specialize in that sport. Based on his industry experience, he has decided to go ahead with an experiment in mining data for athletic aptitude, and has enlisted our help.



Dr. N. Abdolvand

Attributes

- ◆ **Age:** This is the age in years (one decimal precision for the part of the year since the client's last birthday) at the time that the athletic and personality trait battery test was administered. Participants ranged in age from 13-19 years old at the time they took the battery.
- ◆ **Strength:** This is the participant's strength measured through a series of weight lifting exercises and recorded on a scale of 0-10, with 0 being limited strength and 10 being sufficient strength to perform all lifts without any difficulty. No participant scored 8, 9 or 10, but some participants did score 0.
- ◆ **Quickness:** This is the participant's performance on a series of responsiveness tests. Participants were timed on how quickly they were able to press buttons when they were illuminated or to jump when a buzzer sounded. Their response times were tabulated on a scale of 0-6, with 6 being extremely quick response and 0 being very slow. Participants scored all along the spectrum for this attribute.
- ◆ **Injury:** This is a simple yes (1) / no (0) column indicating whether or not the young athlete had already suffered an athletic-related injury that was severe enough to require surgery or other major medical intervention. Common injuries treated with ice, rest, stretching, etc. were entered as 0. Injuries that took more than three week to heal, that required physical therapy or surgery were flagged as 1.



Dr. N. Abdolvand

Attributes

- ◆ **Vision:** Athletes were not only tested on the usual 20/20 vision scale using an eye chart, but were also tested using eye-tracking technology to see how well they were able to pick up objects visually. This test challenged participants to identify items that moved quickly across their field of vision, and to estimate speed and direction of moving objects. Their scores were recorded on a 0 to 4 scale with 4 being perfect vision and identification of moving objects. No participant scored a perfect 4, but the scores did range from 0 to 3.
- ◆ **Endurance:** Participants were subjected to an array of physical fitness tests including running, calisthenics, aerobic and cardiovascular exercise, and distance swimming. Their performance was rated on a scale of 0-10, with 10 representing the ability to perform all tasks without fatigue of any kind. Scores ranged from 0 to 6 on this attribute. Gill has acknowledged to us that even finely tuned professional athletes would not be able to score a 10 on this portion of the battery, as it is specifically designed to test the limits of human endurance.
- ◆ **Agility:** This is the participant's score on a series of tests of their ability to move, twist, turn, jump, change direction, etc. The test checked the athlete's ability to move nimbly, precisely, and powerfully in a full range of directions. This metric is comprehensive in nature, and is influenced by some of the other metrics, as agility is often dictated by one's strength, quickness, etc. Participants were scored between 0 and 100 on this attribute, and in our data set from Gill, we have found performance between 13 and 80.



Dr. N. Abdolvand

Attributes

- ◆ **Decision_Making:** This portion of the battery tests the athlete's process of deciding what to do in athletic situations. Athlete's participated in simulations that tested their choices of whether or not to swing a bat, pass a ball, move to a potentially advantageous location of a playing surface, etc. Their scores were to have been recorded on a scale of 0 to 100, though Gill has indicated that no one who completed the test should have been able to score lower than a 3, as three points are awarded simply for successfully entering and exiting the decision making part of the battery. Gill knows that all 493 of his former athletes represented in this data set successfully entered and exited this portion, but there are a few scores lower than 3, and also a few over 100 in the data set, so we know we have some data preparation in our future.
- ◆ **Prime_Sport:** This attribute is the sport each of the 453 athletes went on to specialize in after they left Gill's academy. This is the attribute Gill is hoping to be able to predict for his current clients. For the boys in this study, this attribute will be one of four sports: football (American, not soccer; sorry soccer fans), Basketball, Baseball, or Hockey.



Dr. N. Abdolvand

Attributes

- ◆ **Vision:** Athletes were not only tested on the usual 20/20 vision scale using an eye chart, but were also tested using eye-tracking technology to see how well they were able to pick up objects visually. This test challenged participants to identify items that moved quickly across their field of vision, and to estimate speed and direction of moving objects. Their scores were recorded on a 0 to 4 scale with 4 being perfect vision and identification of moving objects. No participant scored a perfect 4, but the scores did range from 0 to 3.
- ◆ **Endurance:** Participants were subjected to an array of physical fitness tests including running, calisthenics, aerobic and cardiovascular exercise, and distance swimming. Their performance was rated on a scale of 0-10, with 10 representing the ability to perform all tasks without fatigue of any kind. Scores ranged from 0 to 6 on this attribute. Gill has acknowledged to us that even finely tuned professional athletes would not be able to score a 10 on this portion of the battery, as it is specifically designed to test the limits of human endurance.
- ◆ **Agility:** This is the participant's score on a series of tests of their ability to move, twist, turn, jump, change direction, etc. The test checked the athlete's ability to move nimbly, precisely, and powerfully in a full range of directions. This metric is comprehensive in nature, and is influenced by some of the other metrics, as agility is often dictated by one's strength, quickness, etc. Participants were scored between 0 and 100 on this attribute, and in our data set from Gill, we have found performance between 13 and 80.



Dr. N. Abdolvand

Attributes

- ◆ **Decision_Making:** This portion of the battery tests the athlete's process of deciding what to do in athletic situations. Athlete's participated in simulations that tested their choices of whether or not to swing a bat, pass a ball, move to a potentially advantageous location of a playing surface, etc. Their scores were to have been recorded on a scale of 0 to 100, though Gill has indicated that no one who completed the test should have been able to score lower than a 3, as three points are awarded simply for successfully entering and exiting the decision making part of the battery. Gill knows that all 493 of his former athletes represented in this data set successfully entered and exited this portion, but there are a few scores lower than 3, and also a few over 100 in the data set, so we know we have some data preparation in our future.
- ◆ **Prime_Sport:** This attribute is the sport each of the 453 athletes went on to specialize in after they left Gill's academy. This is the attribute Gill is hoping to be able to predict for his current clients. For the boys in this study, this attribute will be one of four sports: football (American, not soccer; sorry soccer fans), Basketball, Baseball, or Hockey.



Dr. N. Abdolvand

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

9

Prediction Problems: Classification vs. Numeric Prediction

- ◆ Classification
 - predicts categorical class labels (discrete or nominal)
 - classifies data (constructs a model) based on the training set and the values (class labels) in a classifying attribute and uses it in classifying new data
- ◆ Numeric Prediction
 - models continuous-valued functions, i.e., predicts unknown or missing values
- ◆ Typical applications
 - Credit/loan approval:
 - Medical diagnosis: if a tumor is cancerous or benign
 - Fraud detection: if a transaction is fraudulent
 - Web page categorization: which category it is

10

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

9

Prediction Problems: Classification vs. Numeric Prediction

- ◆ Classification
 - predicts categorical class labels (discrete or nominal)
 - classifies data (constructs a model) based on the training set and the values (class labels) in a classifying attribute and uses it in classifying new data
- ◆ Numeric Prediction
 - models continuous-valued functions, i.e., predicts unknown or missing values
- ◆ Typical applications
 - Credit/loan approval:
 - Medical diagnosis: if a tumor is cancerous or benign
 - Fraud detection: if a transaction is fraudulent
 - Web page categorization: which category it is

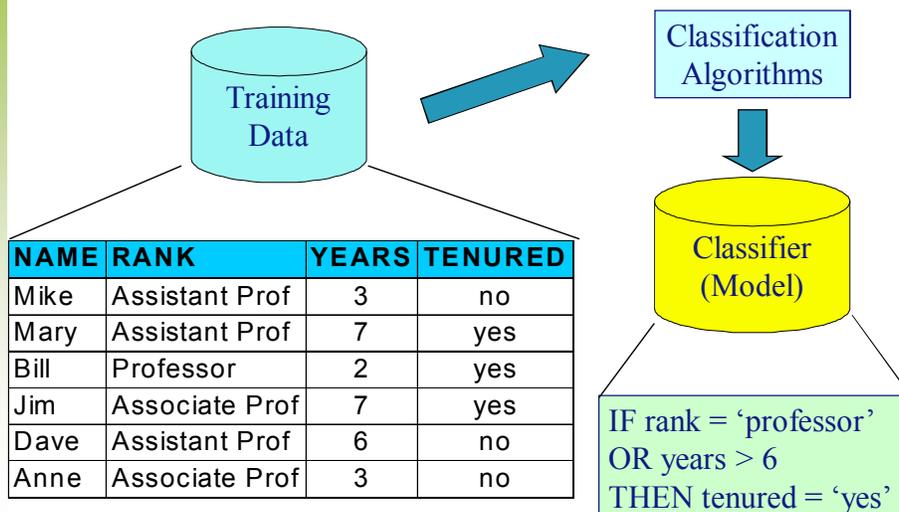
10

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

11

Process (1): Model Construction



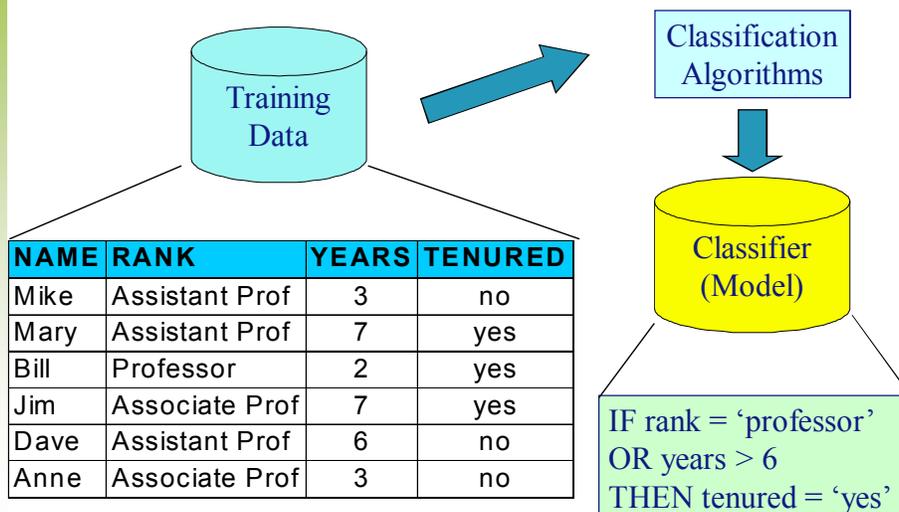
12

Classification—A Two-Step Process

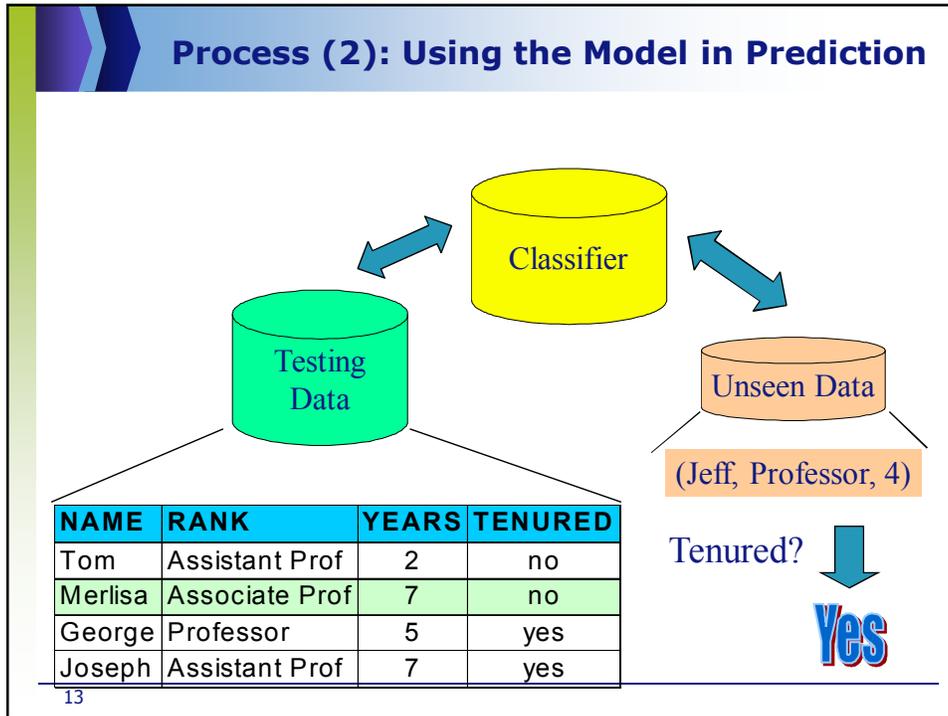
- ◆ Model construction: describing a set of predetermined classes
 - Each tuple/sample is assumed to belong to a predefined class, as determined by the class label attribute
 - The set of tuples used for model construction is training set
 - The model is represented as classification rules, decision 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

11

Process (1): Model Construction



12

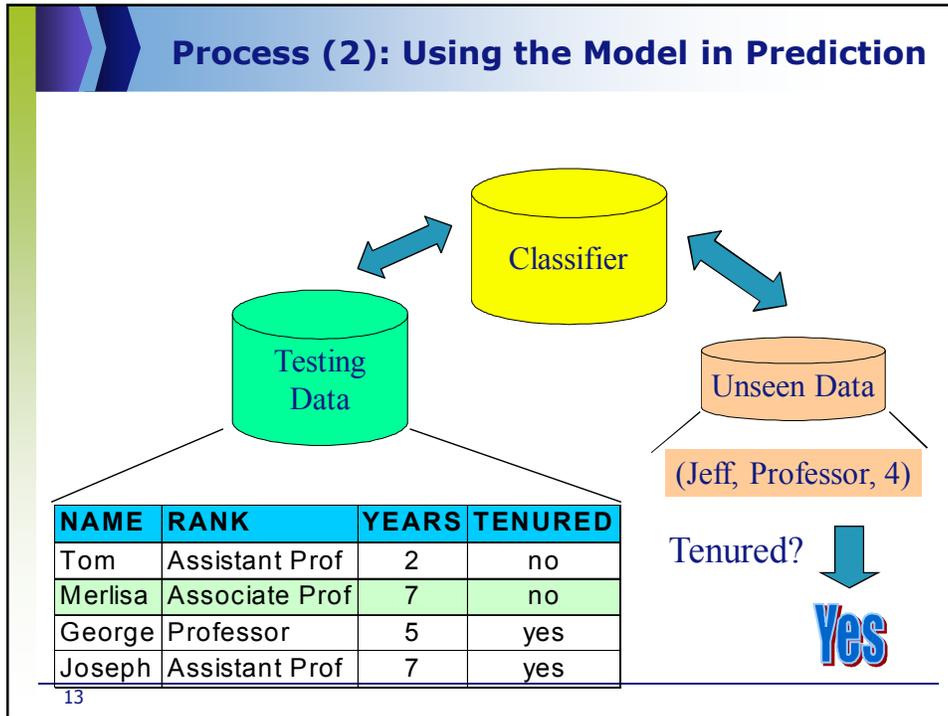


which of the attributes would be best to segment these people into groups, in a way that will distinguish write-offs from non-write-offs?

The figure shows a simple segmentation problem: twelve people represented as stick figures. There are two types of heads: square and circular; and two types of bodies: rectangular and oval; and two of the people have gray bodies while the rest are white. These are the attributes we will use to describe the people. Above each person is the binary target label, Yes or No, indicating (for example) whether the person becomes a loan write-off. We could describe the data on these people as:

- Attributes:
 - head-shape: square, circular
 - body-shape: rectangular, oval
 - body-color: gray, white
- Target variable:
 - write-off: Yes, No

Dr. N. Abdolvand

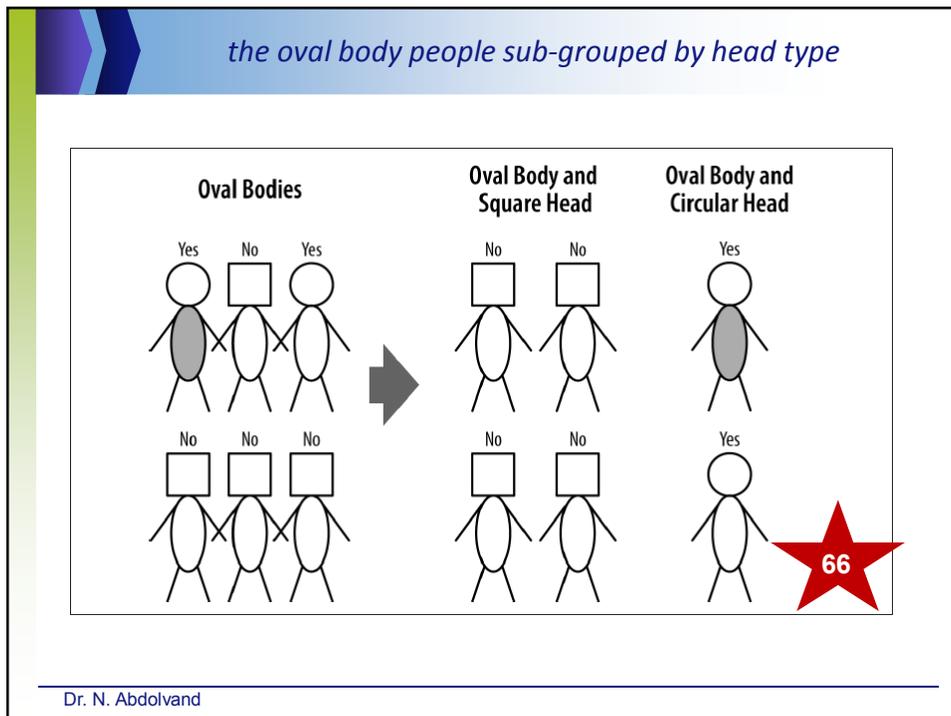
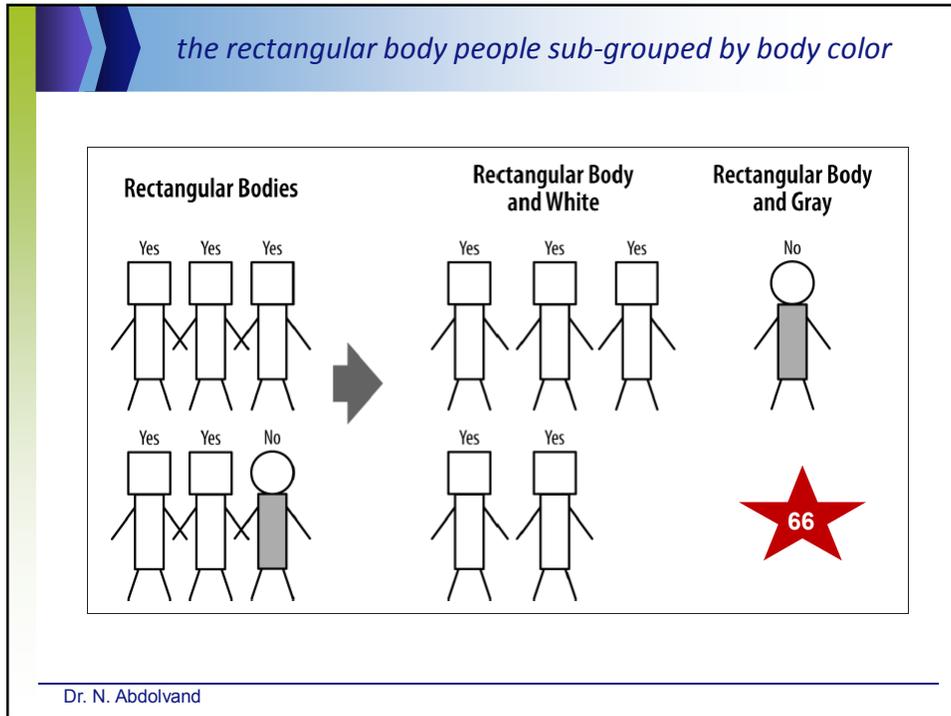


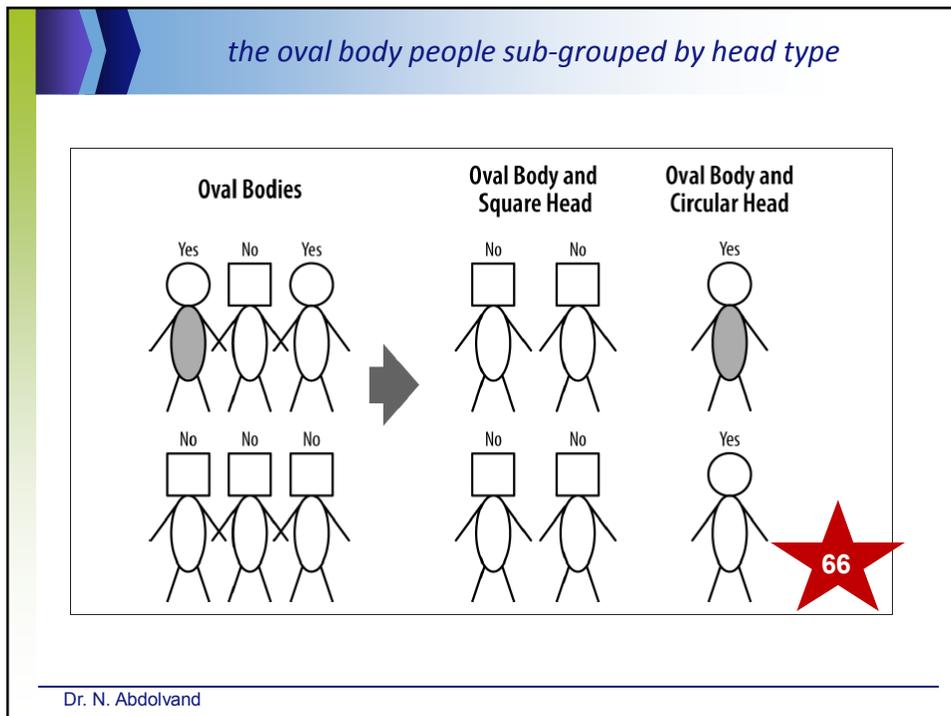
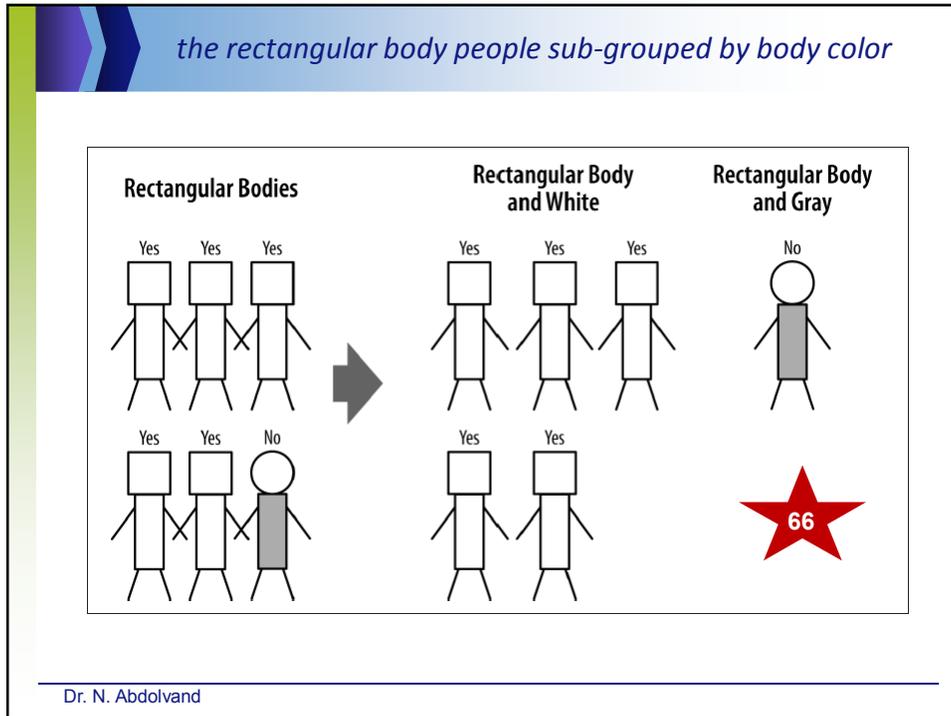
which of the attributes would be best to segment these people into groups, in a way that will distinguish write-offs from non-write-offs?

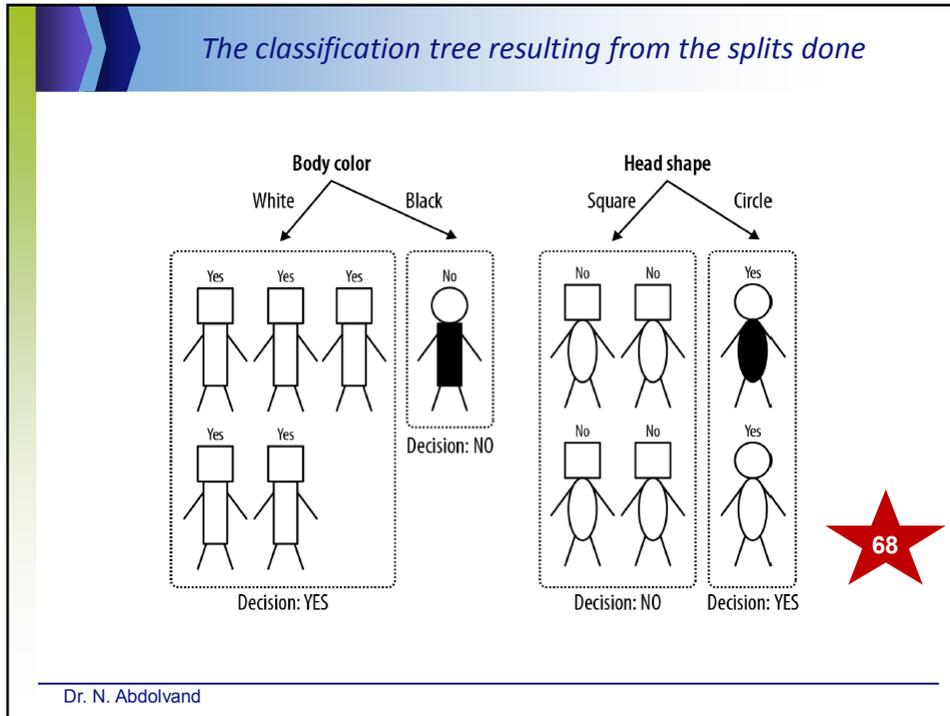
The figure shows a simple segmentation problem: twelve people represented as stick figures. There are two types of heads: square and circular; and two types of bodies: rectangular and oval; and two of the people have gray bodies while the rest are white. These are the attributes we will use to describe the people. Above each person is the binary target label, Yes or No, indicating (for example) whether the person becomes a loan write-off. We could describe the data on these people as:

- Attributes:
 - head-shape: square, circular
 - body-shape: rectangular, oval
 - body-color: gray, white
- Target variable:
 - write-off: Yes, No

Dr. N. Abdolvand







Decision Tree Induction: An Example

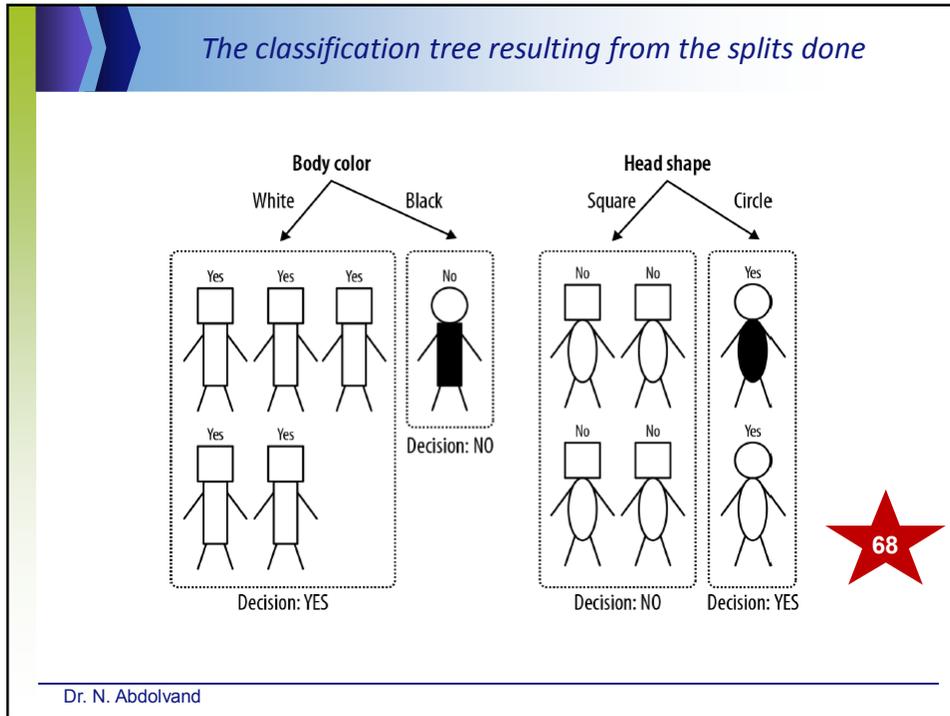
- Training data set: Buys_computer
- The data set follows an example of Quinlan's ID3 (Playing Tennis)
- Resulting tree:

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

```

    graph TD
      A[age?] --> B["<=30"]
      A --> C["31..40"]
      A --> D[">40"]
      C --> E[yes]
      B --> F[student?]
      F --> G[no]
      F --> H[yes]
      G --> I[no]
      H --> J[yes]
      D --> K[credit rating?]
      K --> L[excellent]
      K --> M[fair]
      L --> N[no]
      M --> O[yes]
  
```

18



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

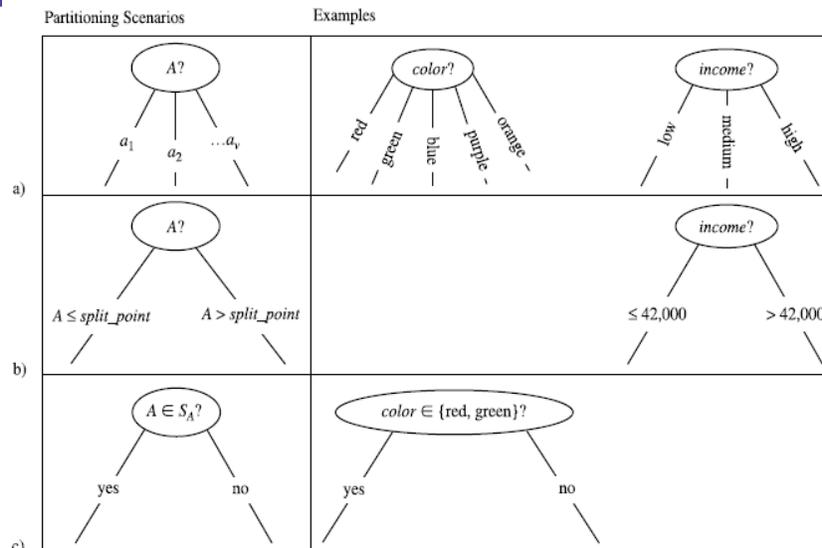
18

Algorithm for Decision Tree Induction

- ◆ Basic algorithm (a greedy algorithm)
 - Tree is constructed in a **top-down recursive divide-and-conquer manner**
 - At start, all the training examples are at the root
 - Attributes are categorical (if continuous-valued, they are discretized in advance)
 - Examples are partitioned recursively based on selected attributes
 - Test attributes are selected on the basis of a heuristic or statistical measure (e.g., **information gain**)
- ◆ Conditions for stopping partitioning
 - All samples for a given node belong to the same class
 - There are no remaining attributes for further partitioning – **majority voting** is employed for classifying the leaf
 - There are no samples left

19

Three possibilities for partitioning tuples based on the splitting criterion



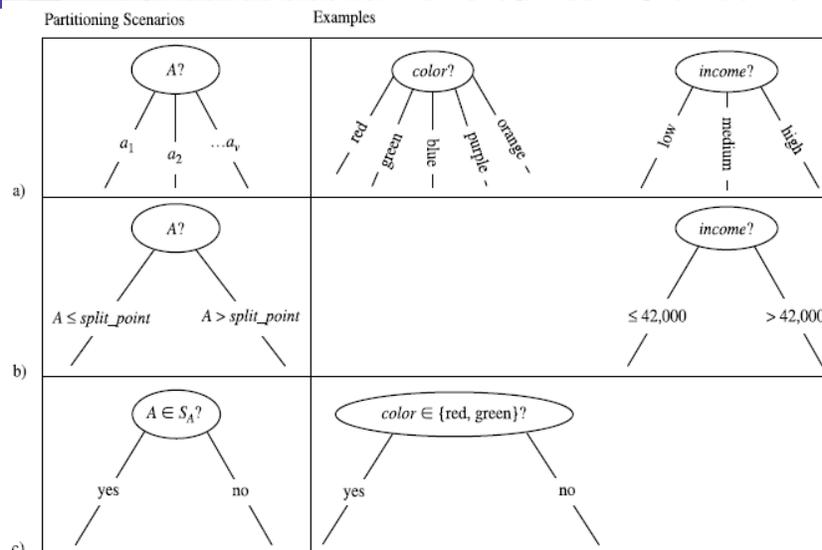
Dr. N. Abdolvand

Algorithm for Decision Tree Induction

- ◆ Basic algorithm (a greedy algorithm)
 - Tree is constructed in a **top-down recursive divide-and-conquer manner**
 - At start, all the training examples are at the root
 - Attributes are categorical (if continuous-valued, they are discretized in advance)
 - Examples are partitioned recursively based on selected attributes
 - Test attributes are selected on the basis of a heuristic or statistical measure (e.g., **information gain**)
- ◆ Conditions for stopping partitioning
 - All samples for a given node belong to the same class
 - There are no remaining attributes for further partitioning – **majority voting** is employed for classifying the leaf
 - There are no samples left

19

Three possibilities for partitioning tuples based on the splitting criterion



Dr. N. Abdolvand

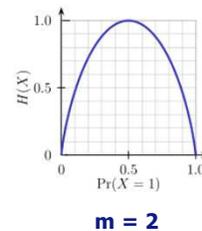
Attribute Selection Measures

- ◆ The attribute selection measure provides a ranking for each attribute describing the given training tuples.
- ◆ The attribute having the best score for the measure is chosen as the *splitting attribute for the given tuples*.
- ◆ If the *splitting attribute* is continuous-valued or if we are restricted to binary trees then, respectively, either a *split point* or a *splitting subset* must also be determined as part of the *splitting criterion*.
- ◆ Three popular attribute selection measures
 - Information gain,
 - gain ratio, and
 - gini index

Dr. N. Abdolvand

Brief Review of Entropy

- ◆ Entropy (Information Theory)
 - A measure of uncertainty associated with a random variable
 - Calculation: For a discrete random variable Y taking m distinct values $\{y_1, \dots, y_m\}$,
 - $H(Y) = -\sum_{i=1}^m p_i \log(p_i)$, where $p_i = P(Y = y_i)$
 - Interpretation:
 - Higher entropy \Rightarrow higher uncertainty
 - Lower entropy \Rightarrow lower uncertainty
- ◆ Conditional Entropy
 - $H(Y|X) = \sum_x p(x)H(Y|X = x)$



22

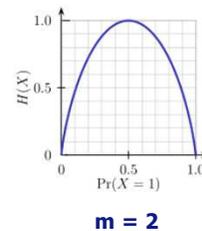
Attribute Selection Measures

- ◆ The attribute selection measure provides a ranking for each attribute describing the given training tuples.
- ◆ The attribute having the best score for the measure is chosen as the *splitting attribute for the given tuples*.
- ◆ If the *splitting attribute* is continuous-valued or if we are restricted to binary trees then, respectively, either a *split point* or a *splitting subset* must also be determined as part of the *splitting criterion*.
- ◆ Three popular attribute selection measures
 - Information gain,
 - gain ratio, and
 - gini index

Dr. N. Abdolvand

Brief Review of Entropy

- ◆ Entropy (Information Theory)
 - A measure of uncertainty associated with a random variable
 - Calculation: For a discrete random variable Y taking m distinct values $\{y_1, \dots, y_m\}$,
 - $H(Y) = -\sum_{i=1}^m p_i \log(p_i)$, where $p_i = P(Y = y_i)$
 - Interpretation:
 - Higher entropy => higher uncertainty
 - Lower entropy => lower uncertainty
- Conditional Entropy
 - $H(Y|X) = \sum_x p(x)H(Y|X = x)$



22

Attribute Selection Measure: Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- Let p_i be the probability that an arbitrary tuple in D belongs to class C_i , estimated by $|C_{i,D}|/|D|$
- Expected information (entropy) needed to classify a tuple in D :

$$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$$
- Information needed (after using A to split D into v partitions) to classify D :

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$
- Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

23

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

$Info(D) = I(9,5) = -\frac{9}{14} \log_2(\frac{9}{14}) - \frac{5}{14} \log_2(\frac{5}{14}) = 0.940$

age	p_i	n_i	$I(p_i, n_i)$
<=30	2	3	0.971
31...40	4	0	0
>40	3	2	0.971

$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$

$\frac{5}{14} I(2,3)$ means "age <=30" has 5 out of 14 samples, with 2 yes'es and 3 no's. Hence

$Gain(age) = Info(D) - Info_{age}(D) = 0.246$

Similarly,

$Gain(income) = 0.029$
 $Gain(student) = 0.151$
 $Gain(credit_rating) = 0.048$

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

Attribute Selection Measure: Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- Let p_i be the probability that an arbitrary tuple in D belongs to class C_i , estimated by $|C_{i,D}|/|D|$
- Expected information (entropy) needed to classify a tuple in D :

$$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$$
- Information needed (after using A to split D into v partitions) to classify D :

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$
- Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

23

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

$Info(D) = I(9,5) = -\frac{9}{14} \log_2(\frac{9}{14}) - \frac{5}{14} \log_2(\frac{5}{14}) = 0.940$

age	p_i	n_i	$I(p_i, n_i)$
<=30	2	3	0.971
31...40	4	0	0
>40	3	2	0.971

$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$

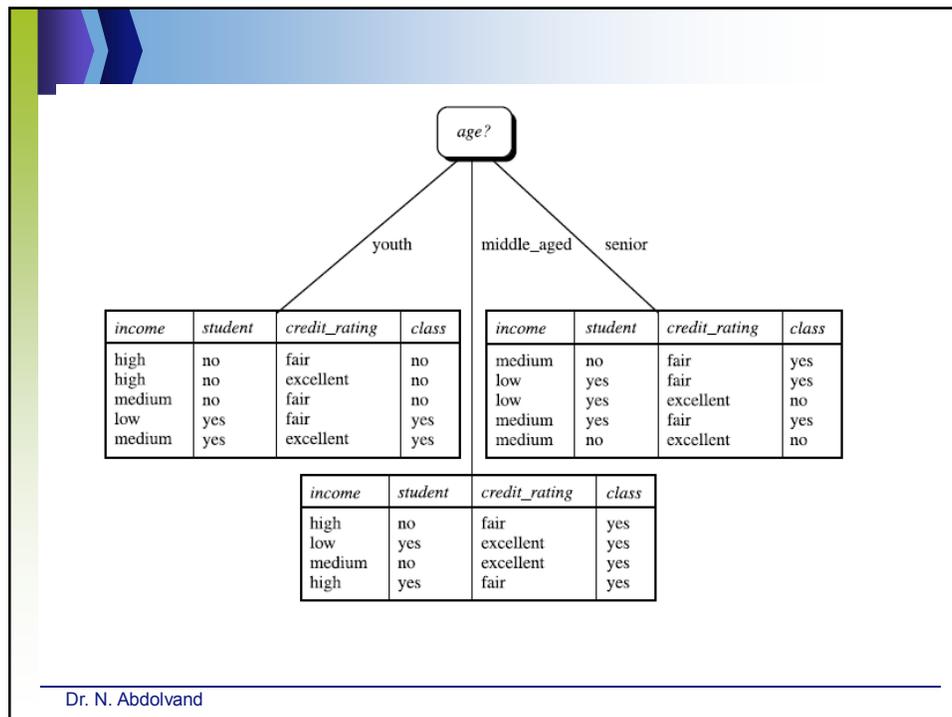
$\frac{5}{14} I(2,3)$ means "age <=30" has 5 out of 14 samples, with 2 yes'es and 3 no's. Hence

$Gain(age) = Info(D) - Info_{age}(D) = 0.246$

Similarly,

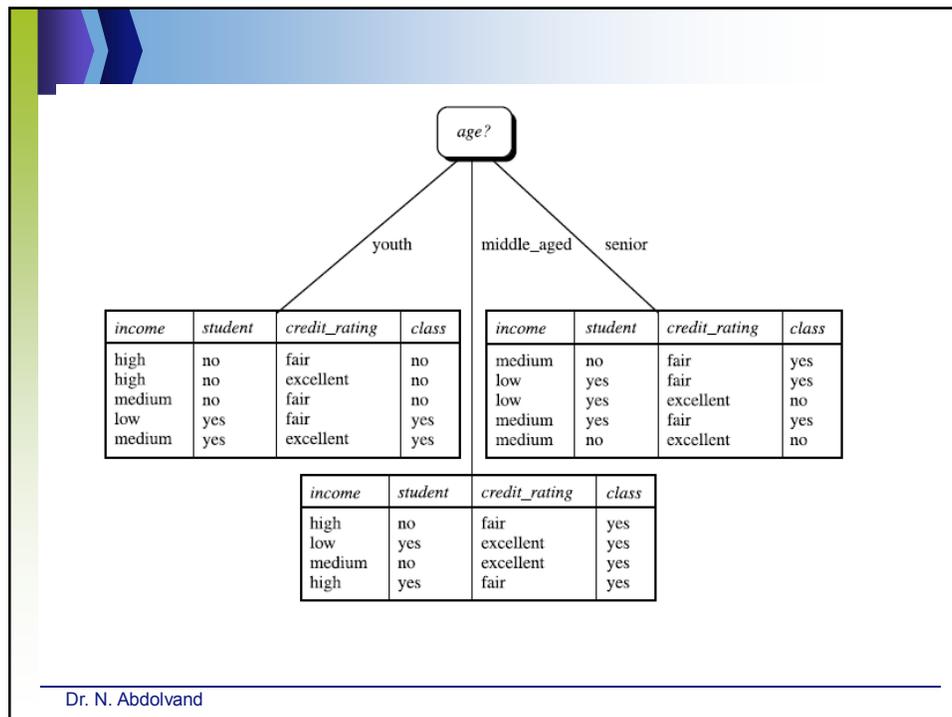
$Gain(income) = 0.029$
 $Gain(student) = 0.151$
 $Gain(credit_rating) = 0.048$

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



Computing Information-Gain for Continuous-Valued Attributes

- ◆ Let attribute A be a continuous-valued attribute
- ◆ Must determine the *best split point* for A
 - Sort the value A in increasing order
 - Typically, the midpoint between each pair of adjacent values is considered as a possible *split point*
 - $(a_i + a_{i+1})/2$ is the midpoint between the values of a_i and a_{i+1}
 - The point with the *minimum expected information requirement* for A is selected as the split-point for A
- ◆ Split:
 - D1 is the set of tuples in D satisfying $A \leq \text{split-point}$, and D2 is the set of tuples in D satisfying $A > \text{split-point}$



Computing Information-Gain for Continuous-Valued Attributes

- ◆ Let attribute A be a continuous-valued attribute
- ◆ Must determine the *best split point* for A
 - Sort the value A in increasing order
 - Typically, the midpoint between each pair of adjacent values is considered as a possible *split point*
 - $(a_i + a_{i+1})/2$ is the midpoint between the values of a_i and a_{i+1}
 - The point with the *minimum expected information requirement* for A is selected as the split-point for A
- ◆ Split:
 - D1 is the set of tuples in D satisfying $A \leq \text{split-point}$, and D2 is the set of tuples in D satisfying $A > \text{split-point}$

Gain Ratio for Attribute Selection (C4.5)

- ◆ Information gain measure is biased towards attributes with a large number of values
- ◆ C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

$$\text{SplitInfo}_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

- $\text{GainRatio}(A) = \text{Gain}(A)/\text{SplitInfo}(A)$

- ◆ Ex

$$\text{SplitInfo}_{\text{income}}(D) = -\frac{4}{14} \times \log_2\left(\frac{4}{14}\right) - \frac{6}{14} \times \log_2\left(\frac{6}{14}\right) - \frac{4}{14} \times \log_2\left(\frac{4}{14}\right) = 1.557$$

- $\text{gain_ratio}(\text{income}) = 0.029/1.557 = 0.019$

- ◆ The attribute with the maximum gain ratio is selected as the splitting attribute

27

Gini Index (CART, IBM IntelligentMiner)

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

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

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 $\text{gini}_A(D)$ is defined as

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

- ◆ Reduction in Impurity:

$$\Delta \text{gini}(A) = \text{gini}(D) - \text{gini}_A(D)$$

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

28

Gain Ratio for Attribute Selection (C4.5)

- ◆ Information gain measure is biased towards attributes with a large number of values
- ◆ C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

$$\text{SplitInfo}_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

- $\text{GainRatio}(A) = \text{Gain}(A)/\text{SplitInfo}(A)$

- ◆ Ex

$$\text{SplitInfo}_{\text{income}}(D) = -\frac{4}{14} \times \log_2\left(\frac{4}{14}\right) - \frac{6}{14} \times \log_2\left(\frac{6}{14}\right) - \frac{4}{14} \times \log_2\left(\frac{4}{14}\right) = 1.557$$

- $\text{gain_ratio}(\text{income}) = 0.029/1.557 = 0.019$

- ◆ The attribute with the maximum gain ratio is selected as the splitting attribute

27

Gini Index (CART, IBM IntelligentMiner)

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

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

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 $\text{gini}_A(D)$ is defined as

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

- ◆ Reduction in Impurity:

$$\Delta \text{gini}(A) = \text{gini}(D) - \text{gini}_A(D)$$

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

28

Computation of Gini Index

- ◆ Ex. D has 9 tuples in buys_computer = "yes" and 5 in "no"

$$gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

- ◆ Suppose the attribute income partitions D into 10 in D_1 : {low, medium} and 4 in D_2

$$gini_{income \in \{low, medium\}}(D) = \left(\frac{10}{14}\right)Gini(D_1) + \left(\frac{4}{14}\right)Gini(D_2)$$

$$\begin{aligned} &Gini_{income \in \{low, medium\}}(D) \\ &= \frac{10}{14}Gini(D_1) + \frac{4}{14}Gini(D_2) \\ &= \frac{10}{14} \left(1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2\right) \\ &= 0.450 \\ &= Gini_{income \in \{high\}}(D). \end{aligned}$$

$Gini_{\{low, high\}}$ is 0.458; $Gini_{\{medium, high\}}$ is 0.450. Thus, split on the {low, medium} (and {high}) since it has the lowest Gini index

- ◆ All attributes are assumed continuous-valued
- ◆ May need other tools, e.g., clustering, to get the possible split values
- ◆ Can be modified for categorical attributes

Comparing Attribute Selection Measures

- ◆ The three measures, in general, return good results but
 - **Information gain:**
 - biased towards multivalued attributes
 - **Gain ratio:**
 - tends to prefer unbalanced splits in which one partition is much smaller than the others
 - **Gini index:**
 - biased to multivalued attributes
 - has difficulty when # of classes is large
 - tends to favor tests that result in equal-sized partitions and purity in both partitions

Computation of Gini Index

- ◆ Ex. D has 9 tuples in buys_computer = "yes" and 5 in "no"

$$gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

- ◆ Suppose the attribute income partitions D into 10 in D_1 : {low, medium} and 4 in D_2

$$gini_{income \in \{low, medium\}}(D) = \left(\frac{10}{14}\right)Gini(D_1) + \left(\frac{4}{14}\right)Gini(D_2)$$

$$\begin{aligned} &Gini_{income \in \{low, medium\}}(D) \\ &= \frac{10}{14}Gini(D_1) + \frac{4}{14}Gini(D_2) \\ &= \frac{10}{14} \left(1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2\right) \\ &= 0.450 \\ &= Gini_{income \in \{high\}}(D). \end{aligned}$$

$Gini_{\{low, high\}}$ is 0.458; $Gini_{\{medium, high\}}$ is 0.450. Thus, split on the {low, medium} (and {high}) since it has the lowest Gini index

- ◆ All attributes are assumed continuous-valued
- ◆ May need other tools, e.g., clustering, to get the possible split values
- ◆ Can be modified for categorical attributes

Comparing Attribute Selection Measures

- ◆ The three measures, in general, return good results but
 - **Information gain:**
 - biased towards multivalued attributes
 - **Gain ratio:**
 - tends to prefer unbalanced splits in which one partition is much smaller than the others
 - **Gini index:**
 - biased to multivalued attributes
 - has difficulty when # of classes is large
 - tends to favor tests that result in equal-sized partitions and purity in both partitions

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”

32

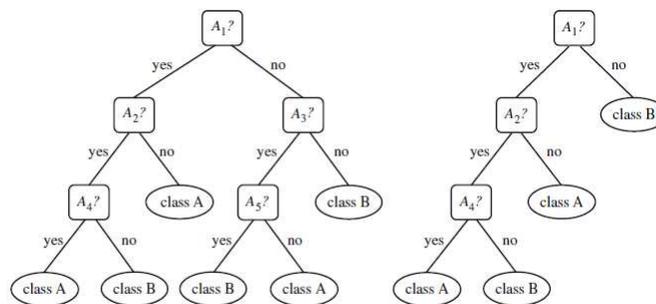


Figure 6.6 An unpruned decision tree and a pruned version of it.

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”

32

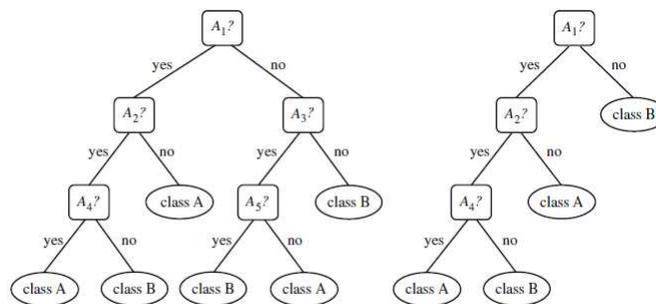


Figure 6.6 An unpruned decision tree and a pruned version of it.

Enhancements to Basic Decision Tree Induction

- ◆ Allow for **continuous-valued attributes**
 - Dynamically define new discrete-valued attributes that partition the continuous attribute value into a discrete set of intervals
- ◆ Handle **missing attribute values**
 - Assign the most common value of the attribute
 - Assign probability to each of the possible values
- ◆ **Attribute construction**
 - Create new attributes based on existing ones that are sparsely represented
 - This reduces fragmentation, repetition, and replication

34

Classification in Large Databases

- ◆ Classification—a classical problem extensively studied by statisticians and machine learning researchers
- ◆ Scalability: Classifying data sets with millions of examples and hundreds of attributes with reasonable speed
- ◆ Why is decision tree induction popular?
 - relatively faster learning speed (than other classification methods)
 - convertible to simple and easy to understand classification rules
 - can use SQL queries for accessing databases
 - comparable classification accuracy with other methods
- ◆ RainForest (VLDB'98 — Gehrke, Ramakrishnan & Ganti)
 - Builds an AVC-list (attribute, value, class label)

35

Enhancements to Basic Decision Tree Induction

- ◆ Allow for **continuous-valued attributes**
 - Dynamically define new discrete-valued attributes that partition the continuous attribute value into a discrete set of intervals
- ◆ Handle **missing attribute values**
 - Assign the most common value of the attribute
 - Assign probability to each of the possible values
- ◆ **Attribute construction**
 - Create new attributes based on existing ones that are sparsely represented
 - This reduces fragmentation, repetition, and replication

34

Classification in Large Databases

- ◆ Classification—a classical problem extensively studied by statisticians and machine learning researchers
- ◆ Scalability: Classifying data sets with millions of examples and hundreds of attributes with reasonable speed
- ◆ Why is decision tree induction popular?
 - relatively faster learning speed (than other classification methods)
 - convertible to simple and easy to understand classification rules
 - can use SQL queries for accessing databases
 - comparable classification accuracy with other methods
- ◆ RainForest (VLDB'98 — Gehrke, Ramakrishnan & Ganti)
 - Builds an AVC-list (attribute, value, class label)

35

Rainforest: Training Set and Its AVC Sets

Training Examples

age	income	student	credit_rating	comp
<=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

37

Presentation of Classification Results

The screenshot shows the dbminer interface with a classification tree for the 'cost' dimension. The tree structure is as follows:

- Root Node: cost (0.00-2000.00)
 - Left Child: cost (0.00-1000.00)
 - Left Child: region (Europe)
 - Left Child: region (Far East)
 - Left Child: region (North America)
 - Left Child: cost (1000.00-2000.00)
 - Left Child: revenue (2000.00-4000.00)
 - Left Child: revenue (4000.00-6000.00)
 - Left Child: revenue (6000.00+)
 - Right Child: revenue (Not Specified)

Classification attribute: product

- Environmental Line (Red)
- GO Sport Line (Green)
- Outdoor Products (Blue)

Rainforest: Training Set and Its AVC Sets

Training Examples

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

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

37

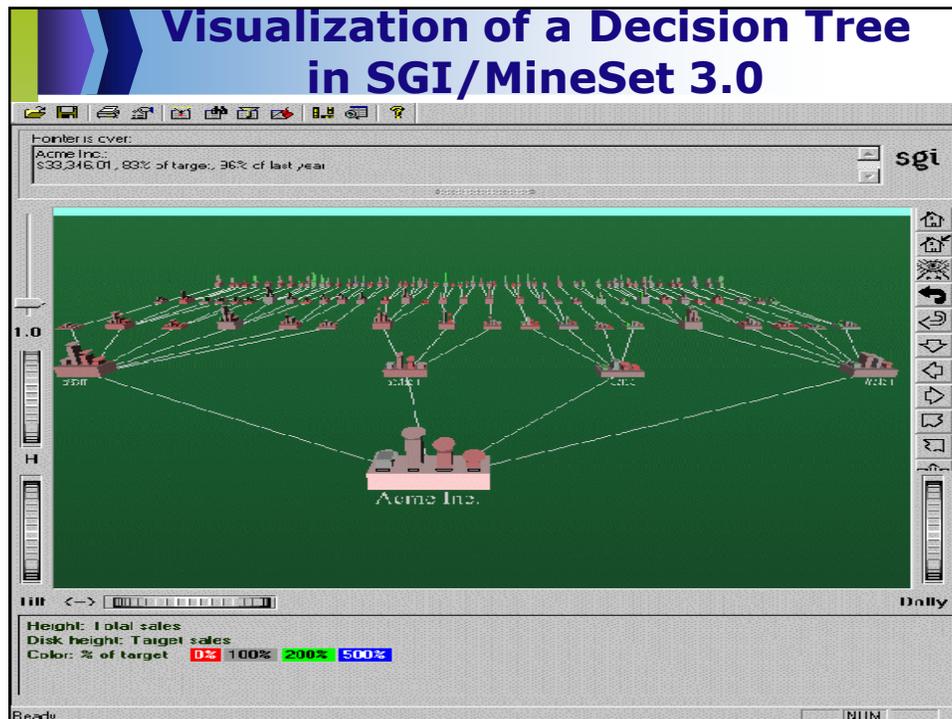
Presentation of Classification Results

The screenshot shows the dbminer interface with a classification tree for the 'cost' dimension. The tree structure is as follows:

- Root node: cost (0.00-2000.00)
 - Left child: cost (0.00-1000.00)
 - Left child: region (Europe)
 - Left child: region (Far East)
 - Left child: region (North America)
 - Left child: cost (1000.00-2000.00)
 - Left child: revenue (2000.00-4000.00)
 - Left child: revenue (4000.00-6000.00)
 - Left child: revenue (6000.00+)
 - Right child: revenue (Not Specified)
- Right child: revenue (0.00-2000.00)

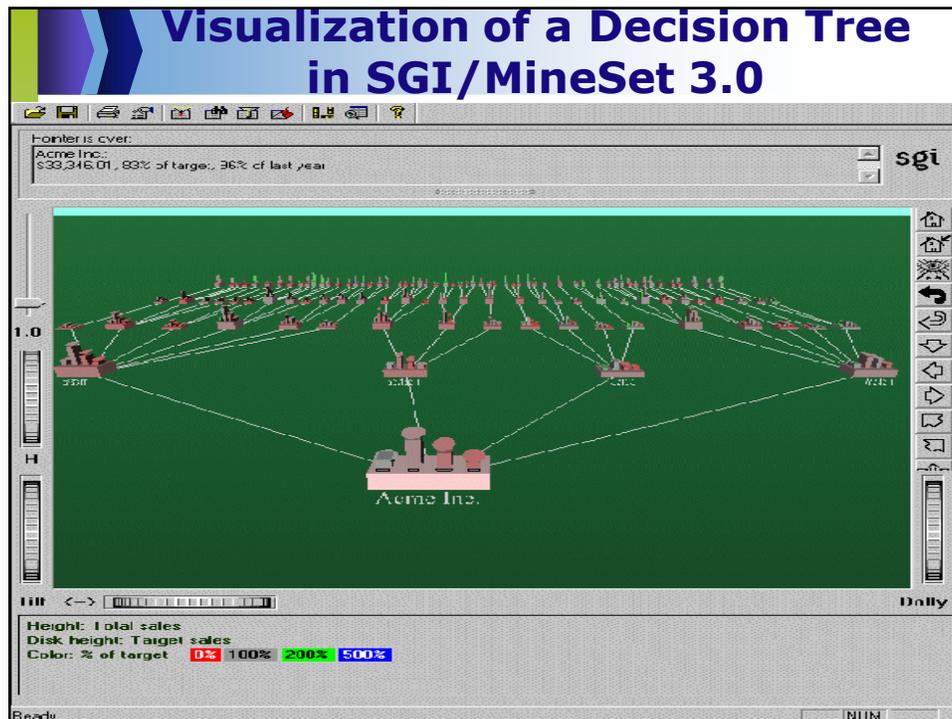
The legend indicates the classification attribute 'product' with the following categories:

- Environmental Line (Red)
- GO Sport Line (Green)
- Outdoor Products (Blue)



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



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 \mathbf{X} is 31.40, medium income

41

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

42

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 \mathbf{X} is 31.40, medium income

41

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

42

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

43

Naïve Bayes Classifier

- ◆ A simplified assumption: attributes are conditionally independent (i.e., no dependence relation between attributes):
- $$P(\mathbf{X} | C_i) = \prod_{k=1}^n P(x_k | C_i) = P(x_1 | C_i) \times P(x_2 | C_i) \times \dots \times P(x_n | C_i)$$
- ◆ This greatly reduces the computation cost: Only counts the class distribution
 - ◆ If A_k is categorical, $P(x_k | C_i)$ is the # of tuples in C_i having value x_k for A_k divided by $|C_{i,D}|$ (# of tuples of C_i in D)
 - ◆ If A_k is continuous-valued, $P(x_k | C_i)$ is usually computed based on Gaussian distribution with a mean μ and standard deviation σ

and $P(x_k | C_i)$ is

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$P(\mathbf{X} | C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i})$$

44

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

43

Naïve Bayes Classifier

- ◆ A simplified assumption: attributes are conditionally independent (i.e., no dependence relation between attributes):
- $$P(\mathbf{X} | C_i) = \prod_{k=1}^n P(x_k | C_i) = P(x_1 | C_i) \times P(x_2 | C_i) \times \dots \times P(x_n | C_i)$$
- ◆ This greatly reduces the computation cost: Only counts the class distribution
 - ◆ If A_k is categorical, $P(x_k | C_i)$ is the # of tuples in C_i having value x_k for A_k divided by $|C_{i,D}|$ (# of tuples of C_i in D)
 - ◆ If A_k is continuous-valued, $P(x_k | C_i)$ is usually computed based on Gaussian distribution with a mean μ and standard deviation σ

and $P(x_k | C_i)$ is

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$P(\mathbf{X} | C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i})$$

44

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

45

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

46

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

45

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 = (age <= 30 , income = medium, student = yes, credit_rating = 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

46

Avoiding the Zero-Probability Problem

- ◆ Naïve Bayesian prediction requires each conditional prob. be **non-zero**. Otherwise, the predicted prob. will be zero

$$P(X | C_i) = \prod_{k=1}^n P(x_k | C_i)$$

- ◆ Ex. Suppose a dataset with 1000 tuples, income=low (0), income= medium (990), and income = high (10)
- ◆ Use **Laplacian correction** (or Laplacian estimator)
 - *Adding 1 to each case*
 - Prob(income = low) = 1/1003
 - Prob(income = medium) = 991/1003
 - Prob(income = high) = 11/1003
 - The “corrected” prob. estimates are close to their “uncorrected” counterparts

47

Naïve Bayes Classifier: Comments

- ◆ Advantages
 - Easy to implement
 - Good results obtained in most of the cases
- ◆ Disadvantages
 - Assumption: class conditional independence, therefore loss of accuracy
 - Practically, dependencies exist among variables
 - E.g., hospitals: patients: Profile: age, family history, etc.
Symptoms: fever, cough etc., Disease: lung cancer, diabetes, etc.
 - Dependencies among these cannot be modeled by Naïve Bayes Classifier
- ◆ How to deal with these dependencies? Bayesian Belief Networks (Chapter 9)

48

Avoiding the Zero-Probability Problem

- ◆ Naïve Bayesian prediction requires each conditional prob. be **non-zero**. Otherwise, the predicted prob. will be zero

$$P(X | C_i) = \prod_{k=1}^n P(x_k | C_i)$$

- ◆ Ex. Suppose a dataset with 1000 tuples, income=low (0), income= medium (990), and income = high (10)
- ◆ Use **Laplacian correction** (or Laplacian estimator)
 - *Adding 1 to each case*
 - Prob(income = low) = 1/1003
 - Prob(income = medium) = 991/1003
 - Prob(income = high) = 11/1003
 - The “corrected” prob. estimates are close to their “uncorrected” counterparts

47

Naïve Bayes Classifier: Comments

- ◆ Advantages
 - Easy to implement
 - Good results obtained in most of the cases
- ◆ Disadvantages
 - Assumption: class conditional independence, therefore loss of accuracy
 - Practically, dependencies exist among variables
 - E.g., hospitals: patients: Profile: age, family history, etc.
Symptoms: fever, cough etc., Disease: lung cancer, diabetes, etc.
 - Dependencies among these cannot be modeled by Naïve Bayes Classifier
- ◆ How to deal with these dependencies? Bayesian Belief Networks (Chapter 9)

48

Using IF-THEN Rules for Classification

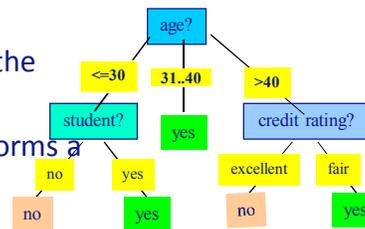
- ◆ Represent the knowledge in the form of **IF-THEN** rules
 - R: IF *age* = youth AND *student* = yes THEN *buys_computer* = yes
 - Rule antecedent/precondition vs. rule consequent
- ◆ Assessment of a rule: *coverage* and *accuracy*
 - n_{covers} = # of tuples covered by R
 - n_{correct} = # of tuples correctly classified by R
 - coverage(R) = $n_{\text{covers}} / |D|$ /* D: training data set */
 - accuracy(R) = $n_{\text{correct}} / n_{\text{covers}}$
- ◆ If more than one rule are triggered, need **conflict resolution**
 - Size ordering: assign the highest priority to the triggering rules that has the “toughest” requirement (i.e., with the *most attribute tests*)
 - Class-based ordering: decreasing order of *prevalence* or *misclassification cost per class*
 - Rule-based ordering (**decision list**): rules are organized into one long priority list, according to some measure of rule quality or by experts

49

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



50

Using IF-THEN Rules for Classification

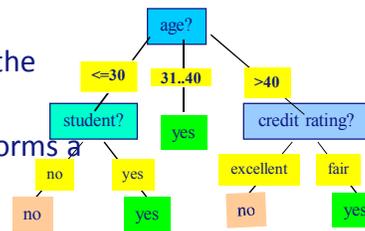
- ◆ Represent the knowledge in the form of **IF-THEN** rules
 - R: IF *age* = youth AND *student* = yes THEN *buys_computer* = yes
 - Rule antecedent/precondition vs. rule consequent
- ◆ Assessment of a rule: *coverage* and *accuracy*
 - n_{covers} = # of tuples covered by R
 - n_{correct} = # of tuples correctly classified by R
 - coverage(R) = $n_{\text{covers}}/|D|$ /* D: training data set */
 - accuracy(R) = $n_{\text{correct}}/n_{\text{covers}}$
- ◆ If more than one rule are triggered, need **conflict resolution**
 - Size ordering: assign the highest priority to the triggering rules that has the “toughest” requirement (i.e., with the *most attribute tests*)
 - Class-based ordering: decreasing order of *prevalence* or *misclassification cost per class*
 - Rule-based ordering (**decision list**): rules are organized into one long priority list, according to some measure of rule quality or by experts

49

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



50

Model Evaluation and Selection

- ◆ Evaluation metrics: How can we measure accuracy? Other metrics to consider?
- ◆ Use **validation test set** of class-labeled tuples instead of training set when assessing accuracy
- ◆ Methods for estimating a classifier's accuracy:
 - Holdout method, random subsampling
 - Cross-validation
 - Bootstrap
- ◆ Comparing classifiers:
 - Confidence intervals
 - Cost-benefit analysis and ROC Curves

55

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

56

Model Evaluation and Selection

- ◆ Evaluation metrics: How can we measure accuracy? Other metrics to consider?
- ◆ Use **validation test set** of class-labeled tuples instead of training set when assessing accuracy
- ◆ Methods for estimating a classifier's accuracy:
 - Holdout method, random subsampling
 - Cross-validation
 - Bootstrap
- ◆ Comparing classifiers:
 - Confidence intervals
 - Cost-benefit analysis and ROC Curves

55

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

56

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

A \ P	C	-C	
C	TP	FN	P
-C	FP	TN	N
	P'	N'	All

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

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

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

- **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
 - $\text{Sensitivity} = TP / P$
- **Specificity**: True Negative recognition rate
 - $\text{Specificity} = TN / N$

57

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

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

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

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

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

- ◆ Perfect score is 1.0

- ◆ Inverse relationship between precision & recall

- ◆ **F measure (F_1 or F-score)**: harmonic mean of precision and recall,

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

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

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

58

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

A \ P	C	-C	
C	TP	FN	P
-C	FP	TN	N
	P'	N'	All

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

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

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

- **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
 - $\text{Sensitivity} = TP / P$
- **Specificity**: True Negative recognition rate
 - $\text{Specificity} = TN / N$

57

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

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

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

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

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

- ◆ Perfect score is 1.0

- ◆ Inverse relationship between precision & recall

- ◆ **F measure (F_1 or F-score)**: harmonic mean of precision and recall,

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

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

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

58

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>)

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

59

Evaluating Classifier Accuracy: Holdout & Cross-Validation Methods

◆ Holdout method

- Given data is randomly partitioned into two independent sets
 - Training set (e.g., 2/3) for model construction
 - Test set (e.g., 1/3) for accuracy estimation
- Random sampling:** a variation of holdout
 - Repeat holdout k times, accuracy = avg. of the accuracies obtained

◆ Cross-validation (*k*-fold, where *k* = 10 is most popular)

- Randomly partition the data into *k* *mutually exclusive* subsets, each approximately equal size
- At *i*-th iteration, use D_i as test set and others as training set
- Leave-one-out:** *k* folds where *k* = # of tuples, for small sized data
- *Stratified cross-validation*:** folds are stratified so that class dist. in each fold is approx. the same as that in the initial data

60

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>)

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

59

Evaluating Classifier Accuracy: Holdout & Cross-Validation Methods

◆ Holdout method

- Given data is randomly partitioned into two independent sets
 - Training set (e.g., 2/3) for model construction
 - Test set (e.g., 1/3) for accuracy estimation
- Random sampling:** a variation of holdout
 - Repeat holdout k times, accuracy = avg. of the accuracies obtained

◆ Cross-validation (*k*-fold, where *k* = 10 is most popular)

- Randomly partition the data into *k* *mutually exclusive* subsets, each approximately equal size
- At *i*-th iteration, use D_i as test set and others as training set
- Leave-one-out:** *k* folds where *k* = # of tuples, for small sized data
- *Stratified cross-validation*:** folds are stratified so that class dist. in each fold is approx. the same as that in the initial data

60

Summary (I)

- ◆ **Classification** is a form of data analysis that extracts **models** describing important data classes.
- ◆ Effective and scalable methods have been developed for **decision tree induction**, **Naive Bayesian classification**, **rule-based classification**, and many other classification methods.
- ◆ **Evaluation metrics** include: accuracy, sensitivity, specificity, precision, recall, F measure, and F_β measure.
- ◆ **Stratified k-fold cross-validation** is recommended for accuracy estimation. **Bagging** and **boosting** can be used to increase overall accuracy by learning and combining a series of individual models.

75

Summary (II)

- ◆ **Significance tests** and **ROC curves** are useful for model selection.
- ◆ There have been numerous **comparisons of the different classification** methods; the matter remains a research topic
- ◆ No single method has been found to be superior over all others for all data sets
- ◆ Issues such as accuracy, training time, robustness, scalability, and interpretability must be considered and can involve trade-offs, further complicating the quest for an overall superior method

76

Summary (I)

- ◆ **Classification** is a form of data analysis that extracts **models** describing important data classes.
- ◆ Effective and scalable methods have been developed for **decision tree induction**, **Naive Bayesian classification**, **rule-based classification**, and many other classification methods.
- ◆ **Evaluation metrics** include: accuracy, sensitivity, specificity, precision, recall, F measure, and F_β measure.
- ◆ **Stratified k-fold cross-validation** is recommended for accuracy estimation. **Bagging** and **boosting** can be used to increase overall accuracy by learning and combining a series of individual models.

75

Summary (II)

- ◆ **Significance tests** and **ROC curves** are useful for model selection.
- ◆ There have been numerous **comparisons of the different classification** methods; the matter remains a research topic
- ◆ No single method has been found to be superior over all others for all data sets
- ◆ Issues such as accuracy, training time, robustness, scalability, and interpretability must be considered and can involve trade-offs, further complicating the quest for an overall superior method

76

References (1)

- ◆ C. Apte and S. Weiss. **Data mining with decision trees and decision rules.** Future Generation Computer Systems, 13, 1997
- ◆ C. M. Bishop, **Neural Networks for Pattern Recognition.** Oxford University Press, 1995
- ◆ L. Breiman, J. Friedman, R. Olshen, and C. Stone. **Classification and Regression Trees.** Wadsworth International Group, 1984
- ◆ C. J. C. Burges. **A Tutorial on Support Vector Machines for Pattern Recognition.** *Data Mining and Knowledge Discovery*, 2(2): 121-168, 1998
- ◆ P. K. Chan and S. J. Stolfo. **Learning arbiter and combiner trees from partitioned data for scaling machine learning.** KDD'95
- ◆ H. Cheng, X. Yan, J. Han, and C.-W. Hsu, [Discriminative Frequent Pattern Analysis for Effective Classification](#), ICDE'07
- ◆ H. Cheng, X. Yan, J. Han, and P. S. Yu, [Direct Discriminative Pattern Mining for Effective Classification](#), ICDE'08
- ◆ W. Cohen. **Fast effective rule induction.** ICML'95
- ◆ G. Cong, K.-L. Tan, A. K. H. Tung, and X. Xu. **Mining top-k covering rule groups for gene expression data.** SIGMOD'05

77

References (2)

- ◆ A. J. Dobson. **An Introduction to Generalized Linear Models.** Chapman & Hall, 1990.
- ◆ G. Dong and J. Li. **Efficient mining of emerging patterns: Discovering trends and differences.** KDD'99.
- ◆ R. O. Duda, P. E. Hart, and D. G. Stork. **Pattern Classification**, 2ed. John Wiley, 2001
- ◆ U. M. Fayyad. **Branching on attribute values in decision tree generation.** AAAI'94.
- ◆ Y. Freund and R. E. Schapire. **A decision-theoretic generalization of on-line learning and an application to boosting.** J. Computer and System Sciences, 1997.
- ◆ J. Gehrke, R. Ramakrishnan, and V. Ganti. **Rainforest: A framework for fast decision tree construction of large datasets.** VLDB'98.
- ◆ J. Gehrke, V. Gant, R. Ramakrishnan, and W.-Y. Loh, **BOAT -- Optimistic Decision Tree Construction.** SIGMOD'99.
- ◆ T. Hastie, R. Tibshirani, and J. Friedman. **The Elements of Statistical Learning: Data Mining, Inference, and Prediction.** Springer-Verlag, 2001.
- ◆ D. Heckerman, D. Geiger, and D. M. Chickering. **Learning Bayesian networks: The combination of knowledge and statistical data.** Machine Learning, 1995.
- ◆ W. Li, J. Han, and J. Pei, **CMAR: Accurate and Efficient Classification Based on Multiple Class-Association Rules**, ICDM'01.

78

References (1)

- ◆ C. Apte and S. Weiss. **Data mining with decision trees and decision rules.** Future Generation Computer Systems, 13, 1997
- ◆ C. M. Bishop, **Neural Networks for Pattern Recognition.** Oxford University Press, 1995
- ◆ L. Breiman, J. Friedman, R. Olshen, and C. Stone. **Classification and Regression Trees.** Wadsworth International Group, 1984
- ◆ C. J. C. Burges. **A Tutorial on Support Vector Machines for Pattern Recognition.** *Data Mining and Knowledge Discovery*, 2(2): 121-168, 1998
- ◆ P. K. Chan and S. J. Stolfo. **Learning arbiter and combiner trees from partitioned data for scaling machine learning.** KDD'95
- ◆ H. Cheng, X. Yan, J. Han, and C.-W. Hsu, [Discriminative Frequent Pattern Analysis for Effective Classification](#), ICDE'07
- ◆ H. Cheng, X. Yan, J. Han, and P. S. Yu, [Direct Discriminative Pattern Mining for Effective Classification](#), ICDE'08
- ◆ W. Cohen. **Fast effective rule induction.** ICML'95
- ◆ G. Cong, K.-L. Tan, A. K. H. Tung, and X. Xu. **Mining top-k covering rule groups for gene expression data.** SIGMOD'05

77

References (2)

- ◆ A. J. Dobson. **An Introduction to Generalized Linear Models.** Chapman & Hall, 1990.
- ◆ G. Dong and J. Li. **Efficient mining of emerging patterns: Discovering trends and differences.** KDD'99.
- ◆ R. O. Duda, P. E. Hart, and D. G. Stork. **Pattern Classification**, 2ed. John Wiley, 2001
- ◆ U. M. Fayyad. **Branching on attribute values in decision tree generation.** AAAI'94.
- ◆ Y. Freund and R. E. Schapire. **A decision-theoretic generalization of on-line learning and an application to boosting.** J. Computer and System Sciences, 1997.
- ◆ J. Gehrke, R. Ramakrishnan, and V. Ganti. **Rainforest: A framework for fast decision tree construction of large datasets.** VLDB'98.
- ◆ J. Gehrke, V. Gant, R. Ramakrishnan, and W.-Y. Loh, **BOAT -- Optimistic Decision Tree Construction.** SIGMOD'99.
- ◆ T. Hastie, R. Tibshirani, and J. Friedman. **The Elements of Statistical Learning: Data Mining, Inference, and Prediction.** Springer-Verlag, 2001.
- ◆ D. Heckerman, D. Geiger, and D. M. Chickering. **Learning Bayesian networks: The combination of knowledge and statistical data.** Machine Learning, 1995.
- ◆ W. Li, J. Han, and J. Pei, **CMAR: Accurate and Efficient Classification Based on Multiple Class-Association Rules**, ICDM'01.

78

References (3)

- ◆ T.-S. Lim, W.-Y. Loh, and Y.-S. Shih. **A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms.** Machine Learning, 2000.
- ◆ J. Magidson. **The Chaid approach to segmentation modeling: Chi-squared automatic interaction detection.** In R. P. Bagozzi, editor, Advanced Methods of Marketing Research, Blackwell Business, 1994.
- ◆ M. Mehta, R. Agrawal, and J. Rissanen. **SLIQ : A fast scalable classifier for data mining.** EDBT'96.
- ◆ T. M. Mitchell. **Machine Learning.** McGraw Hill, 1997.
- ◆ S. K. Murthy, **Automatic Construction of Decision Trees from Data: A Multi-Disciplinary Survey,** Data Mining and Knowledge Discovery 2(4): 345-389, 1998
- ◆ J. R. Quinlan. **Induction of decision trees.** *Machine Learning*, 1:81-106, 1986.
- ◆ J. R. Quinlan and R. M. Cameron-Jones. **FOIL: A midterm report.** ECML'93.
- ◆ J. R. Quinlan. **C4.5: Programs for Machine Learning.** Morgan Kaufmann, 1993.
- ◆ J. R. Quinlan. **Bagging, boosting, and c4.5.** AAAI'96.

79

References (4)

- ◆ R. Rastogi and K. Shim. **Public: A decision tree classifier that integrates building and pruning.** VLDB'98.
- ◆ J. Shafer, R. Agrawal, and M. Mehta. **SPRINT : A scalable parallel classifier for data mining.** VLDB'96.
- ◆ J. W. Shavlik and T. G. Dietterich. **Readings in Machine Learning.** Morgan Kaufmann, 1990.
- ◆ P. Tan, M. Steinbach, and V. Kumar. **Introduction to Data Mining.** Addison Wesley, 2005.
- ◆ S. M. Weiss and C. A. Kulikowski. **Computer Systems that Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning, and Expert Systems.** Morgan Kaufman, 1991.
- ◆ S. M. Weiss and N. Indurkha. **Predictive Data Mining.** Morgan Kaufmann, 1997.
- ◆ I. H. Witten and E. Frank. **Data Mining: Practical Machine Learning Tools and Techniques,** 2ed. Morgan Kaufmann, 2005.
- ◆ X. Yin and J. Han. **CPAR: Classification based on predictive association rules.** SDM'03
- ◆ H. Yu, J. Yang, and J. Han. **Classifying large data sets using SVM with hierarchical clusters.** KDD'03.

80

References (3)

- ◆ T.-S. Lim, W.-Y. Loh, and Y.-S. Shih. **A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms.** Machine Learning, 2000.
- ◆ J. Magidson. **The Chaid approach to segmentation modeling: Chi-squared automatic interaction detection.** In R. P. Bagozzi, editor, Advanced Methods of Marketing Research, Blackwell Business, 1994.
- ◆ M. Mehta, R. Agrawal, and J. Rissanen. **SLIQ : A fast scalable classifier for data mining.** EDBT'96.
- ◆ T. M. Mitchell. **Machine Learning.** McGraw Hill, 1997.
- ◆ S. K. Murthy, **Automatic Construction of Decision Trees from Data: A Multi-Disciplinary Survey,** Data Mining and Knowledge Discovery 2(4): 345-389, 1998
- ◆ J. R. Quinlan. **Induction of decision trees.** *Machine Learning*, 1:81-106, 1986.
- ◆ J. R. Quinlan and R. M. Cameron-Jones. **FOIL: A midterm report.** ECML'93.
- ◆ J. R. Quinlan. **C4.5: Programs for Machine Learning.** Morgan Kaufmann, 1993.
- ◆ J. R. Quinlan. **Bagging, boosting, and c4.5.** AAAI'96.

79

References (4)

- ◆ R. Rastogi and K. Shim. **Public: A decision tree classifier that integrates building and pruning.** VLDB'98.
- ◆ J. Shafer, R. Agrawal, and M. Mehta. **SPRINT : A scalable parallel classifier for data mining.** VLDB'96.
- ◆ J. W. Shavlik and T. G. Dietterich. **Readings in Machine Learning.** Morgan Kaufmann, 1990.
- ◆ P. Tan, M. Steinbach, and V. Kumar. **Introduction to Data Mining.** Addison Wesley, 2005.
- ◆ S. M. Weiss and C. A. Kulikowski. **Computer Systems that Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning, and Expert Systems.** Morgan Kaufman, 1991.
- ◆ S. M. Weiss and N. Indurkha. **Predictive Data Mining.** Morgan Kaufmann, 1997.
- ◆ I. H. Witten and E. Frank. **Data Mining: Practical Machine Learning Tools and Techniques,** 2ed. Morgan Kaufmann, 2005.
- ◆ X. Yin and J. Han. **CPAR: Classification based on predictive association rules.** SDM'03
- ◆ H. Yu, J. Yang, and J. Han. **Classifying large data sets using SVM with hierarchical clusters.** KDD'03.

80

Issues: Evaluating Classification Methods

- ◆ Accuracy
 - classifier accuracy: predicting class label
 - predictor accuracy: guessing value of predicted attributes
- ◆ Speed
 - time to construct the model (training time)
 - time to use the model (classification/prediction time)
- ◆ Robustness: handling noise and missing values
- ◆ Scalability: efficiency in disk-resident databases
- ◆ Interpretability
 - understanding and insight provided by the model
- ◆ Other measures, e.g., goodness of rules, such as decision tree size or compactness of classification rules

81

Predictor Error Measures

- ◆ Measure predictor accuracy: measure how far off the predicted value is from the actual known value
- ◆ **Loss function:** measures the error betw. y_i and the predicted value y_i'
 - Absolute error: $|y_i - y_i'|$
 - Squared error: $(y_i - y_i')^2$
- ◆ Test error (generalization error): the average loss over the test set
 - Mean absolute error: $\frac{\sum_{i=1}^d |y_i - y_i'|}{d}$ Mean squared error: $\frac{\sum_{i=1}^d (y_i - y_i')^2}{d}$
 - Relative absolute error: $\frac{\sum_{i=1}^d |y_i - \bar{y}|}{\sum_{i=1}^d |y_i - \bar{y}|}$ Relative squared error: $\frac{\sum_{i=1}^d (y_i - y_i')^2}{\sum_{i=1}^d (y_i - \bar{y})^2}$

The mean squared-error exaggerates the presence of outliers

Popularly use (square) root mean-square error, similarly, root relative squared error

82

Issues: Evaluating Classification Methods

- ◆ Accuracy
 - classifier accuracy: predicting class label
 - predictor accuracy: guessing value of predicted attributes
- ◆ Speed
 - time to construct the model (training time)
 - time to use the model (classification/prediction time)
- ◆ Robustness: handling noise and missing values
- ◆ Scalability: efficiency in disk-resident databases
- ◆ Interpretability
 - understanding and insight provided by the model
- ◆ Other measures, e.g., goodness of rules, such as decision tree size or compactness of classification rules

81

Predictor Error Measures

- ◆ Measure predictor accuracy: measure how far off the predicted value is from the actual known value
- ◆ **Loss function**: measures the error betw. y_i and the predicted value y_i'
 - Absolute error: $|y_i - y_i'|$
 - Squared error: $(y_i - y_i')^2$
- ◆ Test error (generalization error): the average loss over the test set
 - Mean absolute error: $\frac{\sum_{i=1}^d |y_i - y_i'|}{d}$ Mean squared error: $\frac{\sum_{i=1}^d (y_i - y_i')^2}{d}$
 - Relative absolute error: $\frac{\sum_{i=1}^d |y_i - \bar{y}|}{\sum_{i=1}^d |y_i - \bar{y}|}$ Relative squared error: $\frac{\sum_{i=1}^d (y_i - y_i')^2}{\sum_{i=1}^d (y_i - \bar{y})^2}$

The mean squared-error exaggerates the presence of outliers

Popularly use (square) root mean-square error, similarly, root relative squared error

82

Data Cube-Based Decision-Tree Induction

- ◆ Integration of generalization with decision-tree induction (Kamber et al.'97)
- ◆ Classification at primitive concept levels
 - E.g., precise temperature, humidity, outlook, etc.
 - Low-level concepts, scattered classes, bushy classification-trees
 - Semantic interpretation problems
- ◆ Cube-based multi-level classification
 - Relevance analysis at multi-levels
 - Information-gain analysis with dimension + level

83



Questions?

Data Cube-Based Decision-Tree Induction

- ◆ Integration of generalization with decision-tree induction (Kamber et al.'97)
- ◆ Classification at primitive concept levels
 - E.g., precise temperature, humidity, outlook, etc.
 - Low-level concepts, scattered classes, bushy classification-trees
 - Semantic interpretation problems
- ◆ Cube-based multi-level classification
 - Relevance analysis at multi-levels
 - Information-gain analysis with dimension + level

83



Questions?