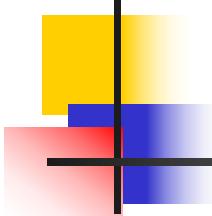


Data Mining: Concepts and Techniques

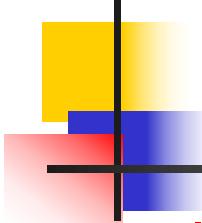
— Slides for Textbook —
— Chapter 7 —

©Jiawei Han and Micheline Kamber
Intelligent Database Systems Research Lab
School of Computing Science
Simon Fraser University, Canada
<http://www.cs.sfu.ca>



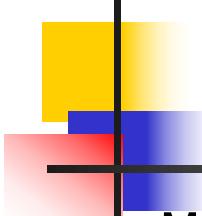
Chapter 7. Classification and Prediction

- n What is classification? What is prediction?
- n Issues regarding classification and prediction
- n Classification by decision tree induction
- n Bayesian Classification
- n Classification by backpropagation
- n Classification based on concepts from association rule mining
- n Other Classification Methods
- n Prediction
- n Classification accuracy
- n Summary



Classification vs. Prediction

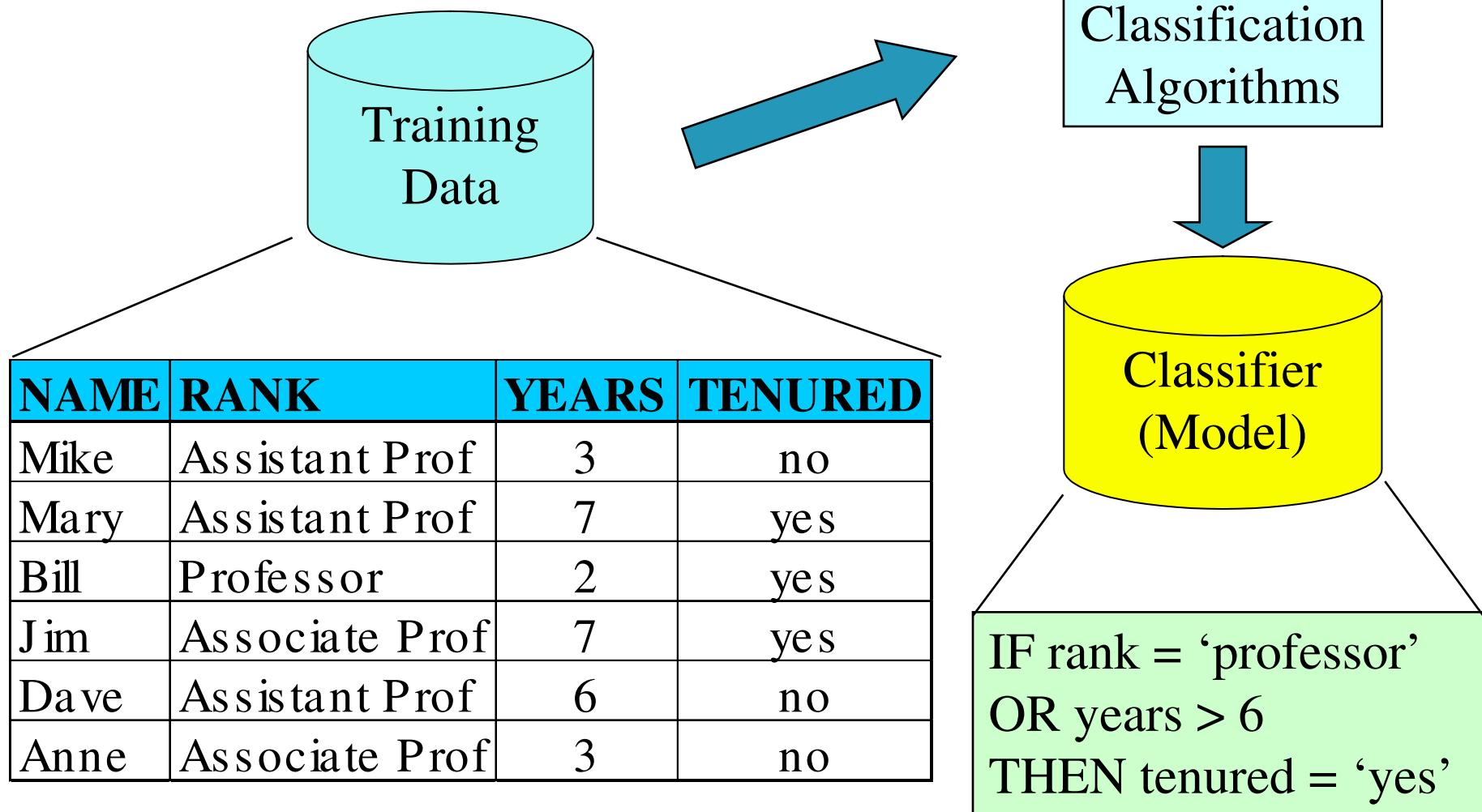
- n **Classification:**
 - n predicts categorical class labels
 - n classifies data (constructs a model) based on the training set and the values (**class labels**) in a classifying attribute and uses it in classifying new data
- n **Prediction:**
 - n models continuous-valued functions, i.e., predicts unknown or missing values
- n **Typical Applications**
 - n credit approval
 - n target marketing
 - n medical diagnosis
 - n treatment effectiveness analysis



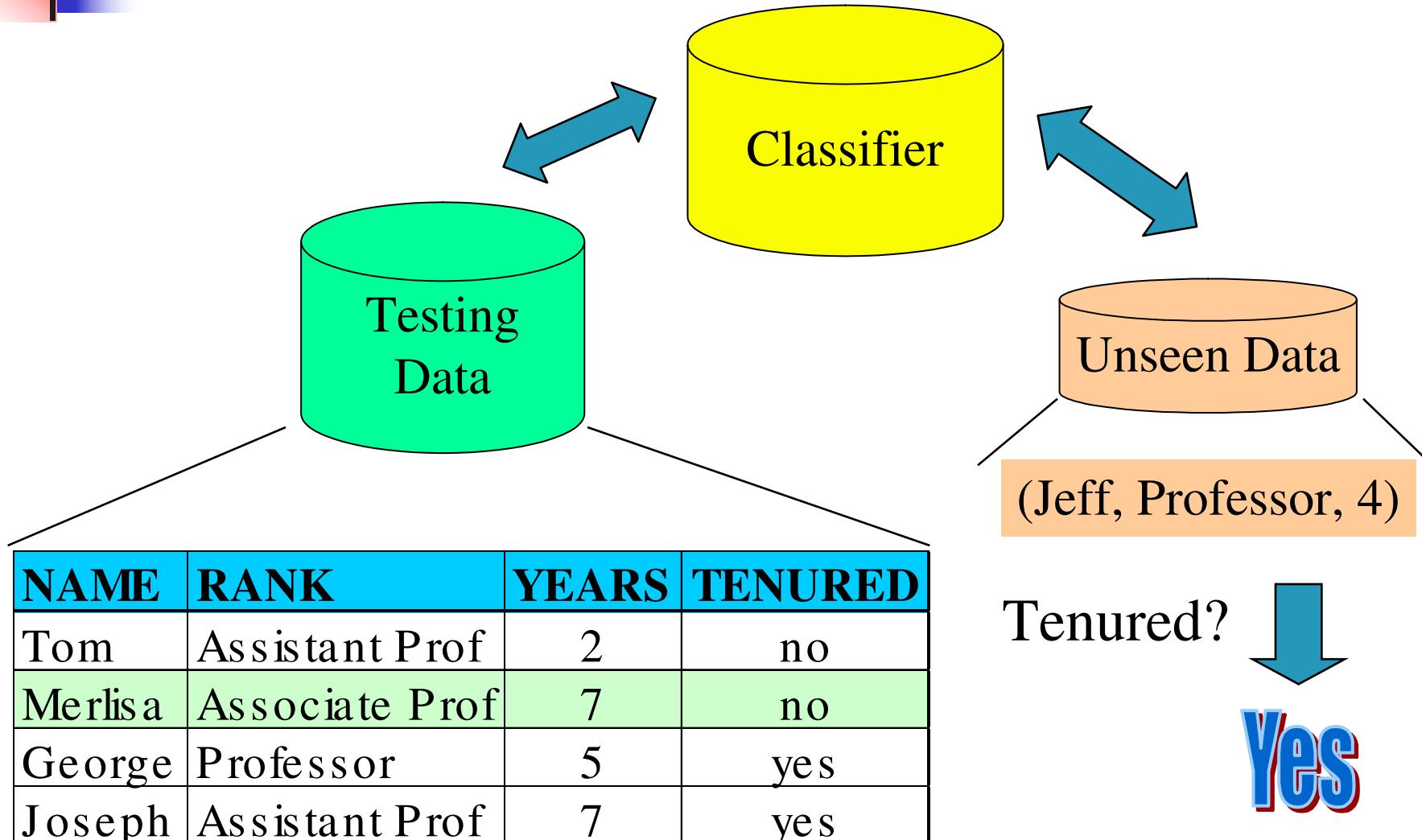
Classification—A Two-Step Process

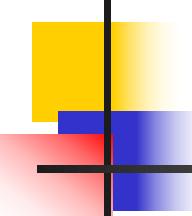
- n Model construction: describing a set of predetermined classes
 - n Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label attribute**
 - n The set of tuples used for model construction: **training set**
 - n The model is represented as classification rules,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

Classification Process (1): Model Construction



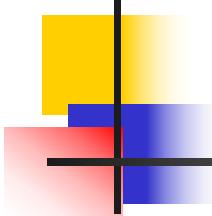
Classification Process (2): Use the Model in Prediction





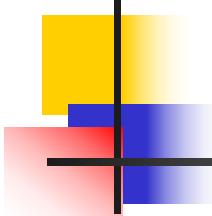
Supervised vs. Unsupervised Learning

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



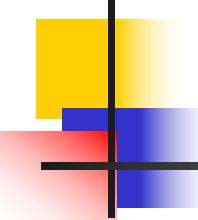
Chapter 7. Classification and Prediction

- n What is classification? What is prediction?
- n Issues regarding classification and prediction
- n Classification by decision tree induction
- n Bayesian Classification
- n Classification by backpropagation
- n Classification based on concepts from association rule mining
- n Other Classification Methods
- n Prediction
- n Classification accuracy
- n Summary



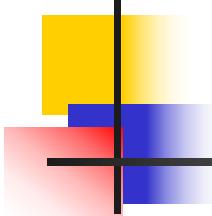
Issues regarding classification and prediction (1): Data Preparation

- „ Data cleaning
 - „ Preprocess data in order to reduce noise and handle missing values
- „ Relevance analysis (feature selection)
 - „ Remove the irrelevant or redundant attributes
- „ Data transformation
 - „ Generalize and/or normalize data



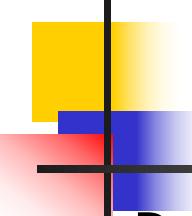
Issues regarding classification and prediction (2): Evaluating Classification Methods

- n Predictive accuracy
- n Speed and scalability
 - n time to construct the model
 - n time to use the model
- n Robustness
 - n handling noise and missing values
- n Scalability
 - n efficiency in disk-resident databases
- n Interpretability:
 - n understanding and insight provided by the model
- n Goodness of rules
 - n decision tree size
 - n compactness of classification rules



Chapter 7. Classification and Prediction

- „ What is classification? What is prediction?
- „ Issues regarding classification and prediction
- „ **Classification by decision tree induction**
- „ Bayesian Classification
- „ Classification by backpropagation
- „ Classification based on concepts from association rule mining
- „ Other Classification Methods
- „ Prediction
- „ Classification accuracy
- „ Summary



Classification by Decision Tree Induction

n Decision tree

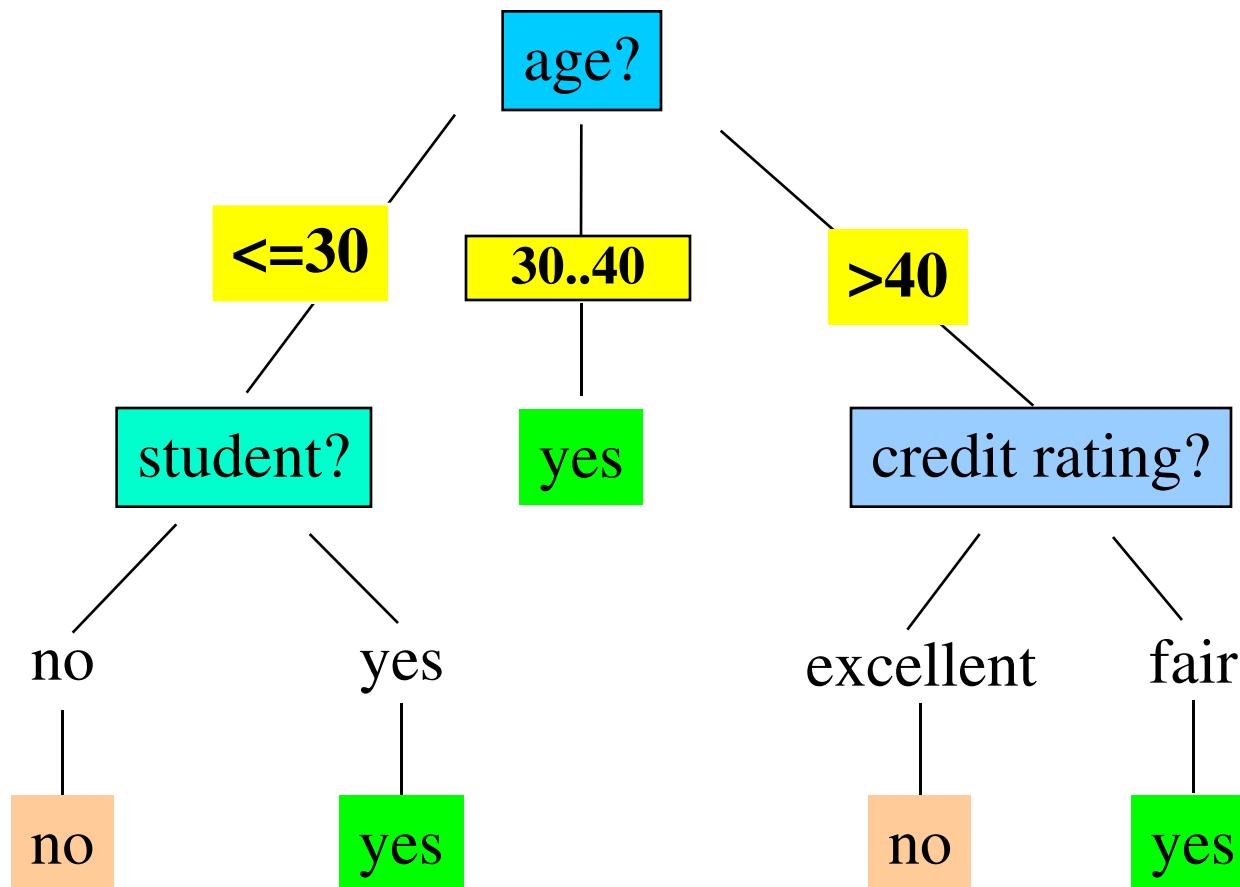
- n A flow-chart-like tree structure
- n Internal node denotes a test on an attribute
- n Branch represents an outcome of the test
- n Leaf nodes represent class labels or class distribution
- n Decision tree generation consists of two phases
 - n Tree construction
 - n At start, all the training examples are at the root
 - n Partition examples recursively based on selected attributes
 - n Tree pruning
 - n Identify and remove branches that reflect noise or outliers
- n Use of decision tree: Classifying an unknown sample
 - n Test the attribute values of the sample against the decision tree

Training Dataset

This follows an example from Quinlan's ID3

age	income	student	credit rating	buys computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
30...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

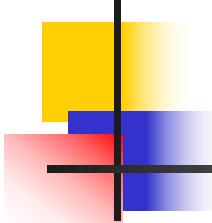
Output: A Decision Tree for “buys_computer”



Algorithm for Decision Tree Induction

n Basic algorithm (a greedy algorithm)

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



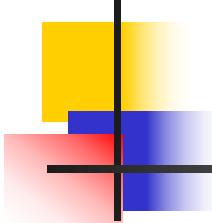
Attribute Selection Measure

- n **Information gain** (ID3/C4.5)
 - n All attributes are assumed to be categorical
 - n Can be modified for continuous-valued attributes
- n **Gini index** (IBM IntelligentMiner)
 - n All attributes are assumed continuous-valued
 - n Assume there exist several possible split values for each attribute
 - n May need other tools, such as clustering, to get the possible split values
 - n Can be modified for categorical attributes

Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- Assume there are two classes, P and N
 - Let the set of examples S contain p elements of class P and n elements of class N
 - The amount of information, needed to decide if an arbitrary example in S belongs to P or N is defined as

$$I(p, n) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$



Information Gain in Decision Tree Induction

- n Assume that using attribute A a set S will be partitioned into sets $\{S_1, S_2, \dots, S_v\}$
 - n If S_i contains p_i examples of P and n_i examples of N , the **entropy**, or the expected information needed to classify objects in all subtrees S_i , is

$$E(A) = \sum_{i=1}^v \frac{p_i + n_i}{p + n} I(p_i, n_i)$$

- n The encoding information that would be gained by branching on A

$$Gain(A) = I(p, n) - E(A)$$

Attribute Selection by Information Gain Computation

- g Class P: `buys_computer` = "yes"
- g Class N: `buys_computer` = "no"
- g $I(p, n) = I(9, 5) = 0.940$
- g Compute the entropy for `age`:

age	p_i	n_i	$I(p_i, n_i)$
≤ 30	2	3	0.971
$30 \dots 40$	4	0	0
>40	3	2	0.971

$$\begin{aligned}E(\text{age}) &= \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) \\&\quad + \frac{5}{14} I(3,2) = 0.971\end{aligned}$$

Hence

$$Gain(\text{age}) = I(p, n) - E(\text{age})$$

Similarly

$$Gain(\text{income}) = 0.029$$

$$Gain(\text{student}) = 0.151$$

$$Gain(\text{credit_rating}) = 0.048$$

Gini Index (IBM IntelligentMiner)

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

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

where p_j is the relative frequency of class j in T .

- If a data set T is split into two subsets T_1 and T_2 with sizes N_1 and N_2 respectively, the *gini* index of the split data contains examples from n classes, the *gini* index $gini(T)$ is defined as

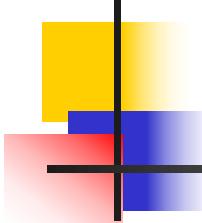
$$gini_{split}(T) = \frac{N_1}{N} gini(T_1) + \frac{N_2}{N} gini(T_2)$$

- The attribute provides the smallest $gini_{split}(T)$ is chosen to split the node (*need to enumerate all possible splitting points for each attribute*).

Extracting Classification Rules from Trees

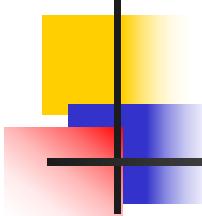
- Represent the knowledge in the form of **IF-THEN** rules
- One rule is created for each path from the root to a leaf
- Each attribute-value pair along a path forms a conjunction
- The leaf node holds the class prediction
- Rules are easier for humans to understand
- Example

IF *age* = “ ≤ 30 ” AND *student* = “no” THEN *buys_computer* = “no”
IF *age* = “ ≤ 30 ” AND *student* = “yes” THEN *buys_computer* = “yes”
IF *age* = “31...40” THEN *buys_computer* = “yes”
IF *age* = “ > 40 ” AND *credit_rating* = “excellent” THEN
 buys_computer = “yes”
IF *age* = “ ≤ 30 ” AND *credit_rating* = “fair” THEN *buys_computer* =
 “no”



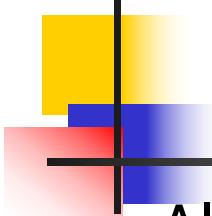
Avoid Overfitting in Classification

- The generated tree may overfit the training data
 - Too many branches, some may reflect anomalies due to noise or outliers
 - Result is in 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”



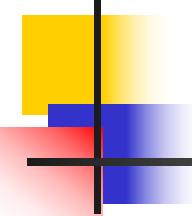
Approaches to Determine the Final Tree Size

- n Separate training (2/3) and testing (1/3) sets
- n Use cross validation, e.g., 10-fold cross validation
- n Use all the data for training
 - n but apply a **statistical test** (e.g., chi-square) to estimate whether expanding or pruning a node may improve the entire distribution
- n Use minimum description length (MDL) principle:



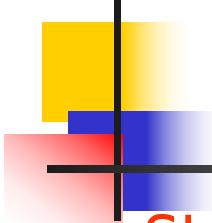
Enhancements to basic decision tree induction

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



Classification in Large Databases

- n Classification—a classical problem extensively studied by statisticians and machine learning researchers
- n Scalability: Classifying data sets with millions of examples and hundreds of attributes with reasonable speed
- n Why decision tree induction in data mining?
 - n relatively faster learning speed (than other classification methods)
 - n convertible to simple and easy to understand classification rules
 - n can use SQL queries for accessing databases
 - n comparable classification accuracy with other methods



Scalable Decision Tree Induction Methods in Data Mining Studies

n

SLIQ (EDBT'96 – Mehta et al.)

- n builds an index for each attribute and only class list and the current attribute list reside in memory

SPRINT (VLDB'96 – J. Shafer et al.)

- n constructs an attribute list data structure

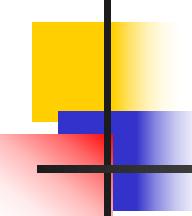
PUBLIC (VLDB'98 – Rastogi & Shim)

- n integrates tree splitting and tree pruning: stop growing the tree earlier

RainForest (VLDB'98 – Gehrke, Ramakrishnan & Ganti)

- n separates the scalability aspects from the criteria that determine the quality of the tree

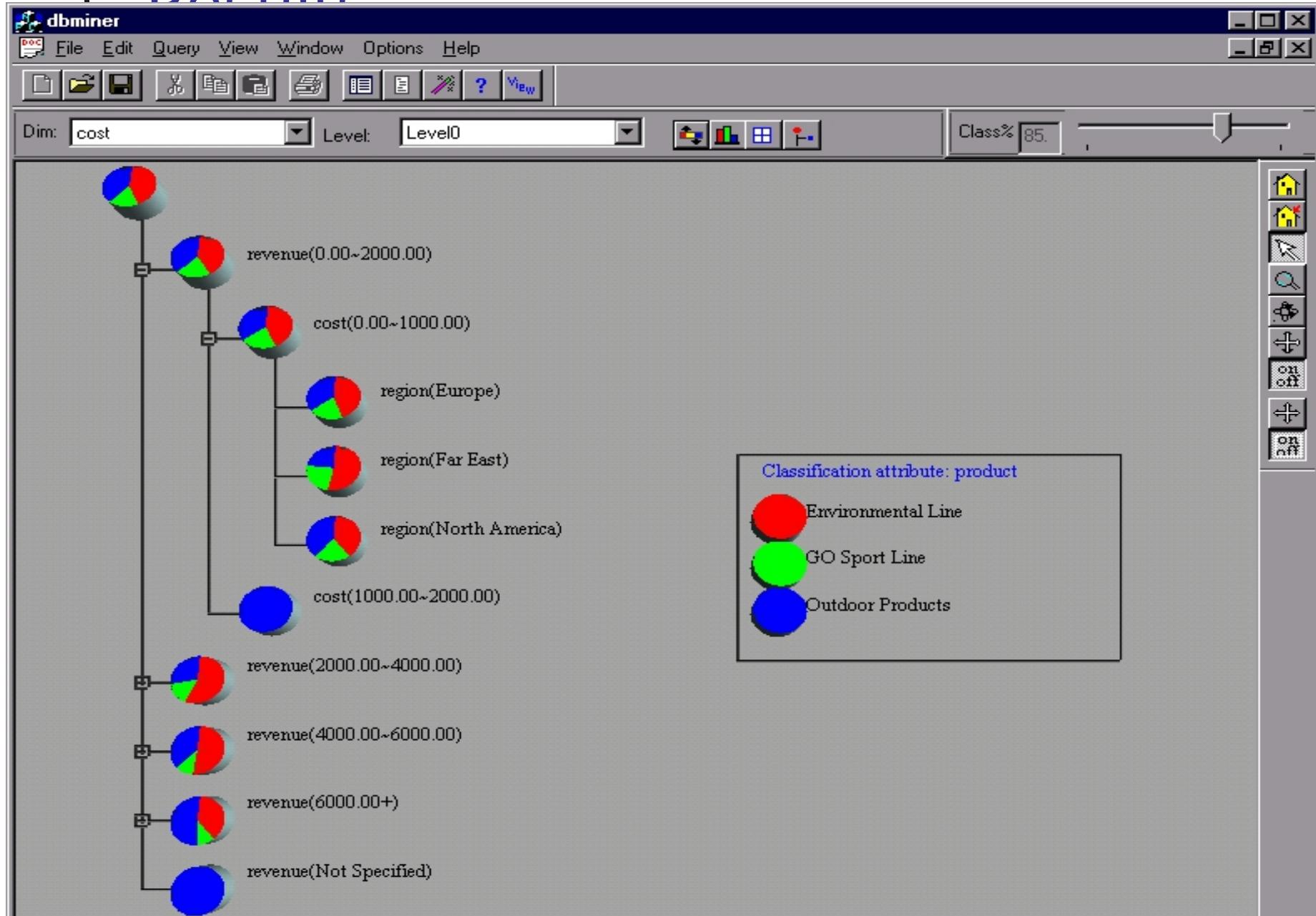
- n builds an AVC-list (attribute, value, class label)

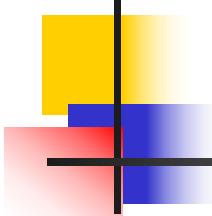


Data Cube-Based Decision-Tree Induction

- 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.

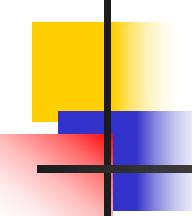
Presentation of Classification Results





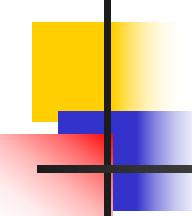
Chapter 7. Classification and Prediction

- n What is classification? What is prediction?
- n Issues regarding classification and prediction
- n Classification by decision tree induction
- n **Bayesian Classification**
- n Classification by backpropagation
- n Classification based on concepts from association rule mining
- n Other Classification Methods
- n Prediction
- n Classification accuracy
- n Summary



Bayesian Classification: Why?

- n Probabilistic learning: Calculate explicit probabilities for hypothesis, among the most practical approaches to certain types of learning problems
- n Incremental: Each training example can incrementally increase/decrease the probability that a hypothesis is correct. Prior knowledge can be combined with observed data.
- n Probabilistic prediction: Predict multiple hypotheses, weighted by their probabilities
- n 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 Theorem

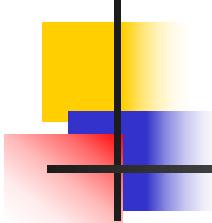
- n Given training data D , *posteriori probability of a hypothesis h* , $P(h|D)$ follows the Bayes theorem

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

- n MAP (maximum posteriori) hypothesis

$$h_{MAP} \equiv \arg \max_{h \in H} P(h|D) = \arg \max_{h \in H} P(D|h)P(h).$$

- n Practical difficulty: require initial knowledge of many probabilities, significant computational cost



Naïve Bayes Classifier (I)

- n A simplified assumption: attributes are conditionally independent:

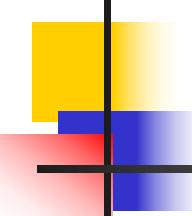
$$P(C_j|V) \propto P(C_j) \prod_{i=1}^n P(v_i|C_j)$$

- n Greatly reduces the computation cost, only count the class distribution.

Naive Bayesian Classifier (II)

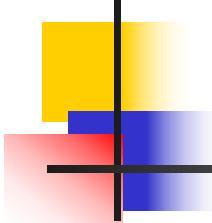
- Given a training set, we can compute the probabilities

Outlook	P	N	Humidity	P	N
sunny	2/9	3/5	high	3/9	4/5
overcast	4/9	0	normal	6/9	1/5
rain	3/9	2/5			
Tempreature			Windy		
hot	2/9	2/5	true	3/9	3/5
mild	4/9	2/5	false	6/9	2/5
cool	3/9	1/5			



Bayesian classification

- n The classification problem may be formalized using **a-posteriori probabilities**:
- n $P(C|X)$ = prob. that the sample tuple $X = \langle x_1, \dots, x_k \rangle$ is of class C.
- n E.g. $P(\text{class}=N \mid \text{outlook}=\text{sunny}, \text{windy}=\text{true}, \dots)$
- n Idea: assign to sample **X** the class label **C** such that **$P(C|X)$ is maximal**

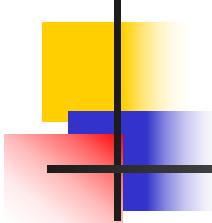


Estimating a-posteriori probabilities

n Bayes theorem:

$$P(C|X) = P(X|C) \cdot P(C) / P(X)$$

- n $P(X)$ is constant for all classes
- n $P(C)$ = relative freq of class C samples
- n C such that $P(C|X)$ is maximum =
C such that $P(X|C) \cdot P(C)$ is maximum
- n Problem: computing $P(X|C)$ is unfeasible!



Naïve Bayesian Classification

- n Naïve assumption: **attribute independence**
$$P(x_1, \dots, x_k | C) = P(x_1 | C) \cdot \dots \cdot P(x_k | C)$$
- n If i-th attribute is **categorical**:
 $P(x_i | C)$ is estimated as the relative freq of samples having value x_i as i-th attribute in class C
- n If i-th attribute is **continuous**:
 $P(x_i | C)$ is estimated thru a Gaussian density function
- n Computationally easy in both cases

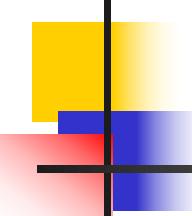
Play-tennis example: estimating $P(x_i|C)$

Outlook	Temperature	Humidity	Windy	Class
sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

$$P(p) = 9/14$$

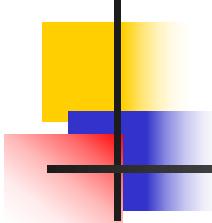
$$P(n) = 5/14$$

outlook		
$P(\text{sunny} p) = 2/9$	$P(\text{sunny} n) = 3/5$	
$P(\text{overcast} p) = 4/9$	$P(\text{overcast} n) = 0$	
$P(\text{rain} p) = 3/9$	$P(\text{rain} n) = 2/5$	
temperature		
$P(\text{hot} p) = 2/9$	$P(\text{hot} n) = 2/5$	
$P(\text{mild} p) = 4/9$	$P(\text{mild} n) = 2/5$	
$P(\text{cool} p) = 3/9$	$P(\text{cool} n) = 1/5$	
humidity		
$P(\text{high} p) = 3/9$	$P(\text{high} n) = 4/5$	
$P(\text{normal} p) = 6/9$	$P(\text{normal} n) = 2/5$	
windy		
$P(\text{true} p) = 3/9$	$P(\text{true} n) = 3/5$	



Play-tennis example: classifying X

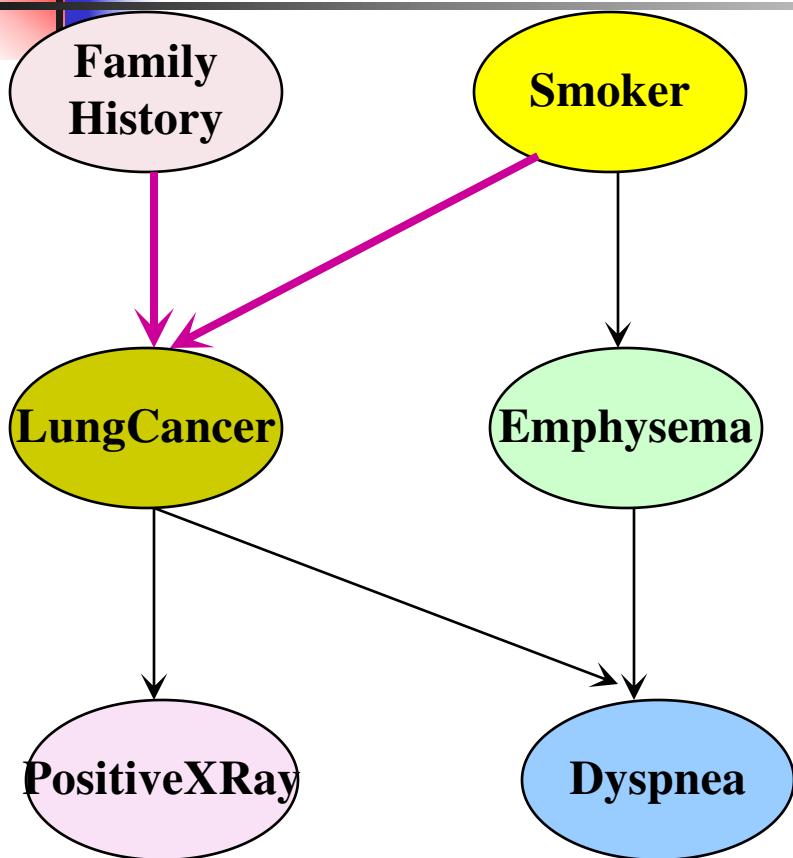
- n An unseen sample $X = \langle \text{rain}, \text{hot}, \text{high}, \text{false} \rangle$
- n $P(X|p) \cdot P(p) =$
 $P(\text{rain}|p) \cdot P(\text{hot}|p) \cdot P(\text{high}|p) \cdot P(\text{false}|p) \cdot P(p) =$
 $3/9 \cdot 2/9 \cdot 3/9 \cdot 6/9 \cdot 9/14 = 0.010582$
- n $P(X|n) \cdot P(n) =$
 $P(\text{rain}|n) \cdot P(\text{hot}|n) \cdot P(\text{high}|n) \cdot P(\text{false}|n) \cdot P(n) =$
 $2/5 \cdot 2/5 \cdot 4/5 \cdot 2/5 \cdot 5/14 = \textcolor{red}{0.018286}$
- n Sample X is classified in class n (don't play)



The independence hypothesis...

- n ... makes computation possible
- n ... yields optimal classifiers when satisfied
- n ... but is seldom satisfied in practice, as attributes (variables) are often correlated.
- n Attempts to overcome this limitation:
 - n **Bayesian networks**, that combine Bayesian reasoning with causal relationships between attributes
 - n **Decision trees**, that reason on one attribute at the time, considering most important attributes first

Bayesian Belief Networks (I)

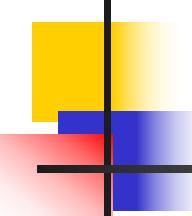


(FH, S) (FH, ~S)(~FH, S) (~FH, ~S)

LC	0.8	0.5	0.7	0.1
~LC	0.2	0.5	0.3	0.9

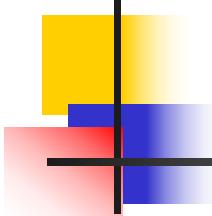
**The conditional probability table
for the variable LungCancer**

Bayesian Belief Networks



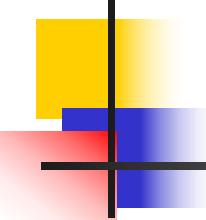
Bayesian Belief Networks (II)

- n Bayesian belief network allows a *subset* of the variables conditionally independent
- n A graphical model of causal relationships
- n Several cases of learning Bayesian belief networks
 - n Given both network structure and all the variables: easy
 - n Given network structure but only some variables
 - n When the network structure is not known in advance



Chapter 7. Classification and Prediction

- n What is classification? What is prediction?
- n Issues regarding classification and prediction
- n Classification by decision tree induction
- n Bayesian Classification
- n **Classification by backpropagation**
- n Classification based on concepts from association rule mining
- n Other Classification Methods
- n Prediction
- n Classification accuracy
- n Summary



Neural Networks

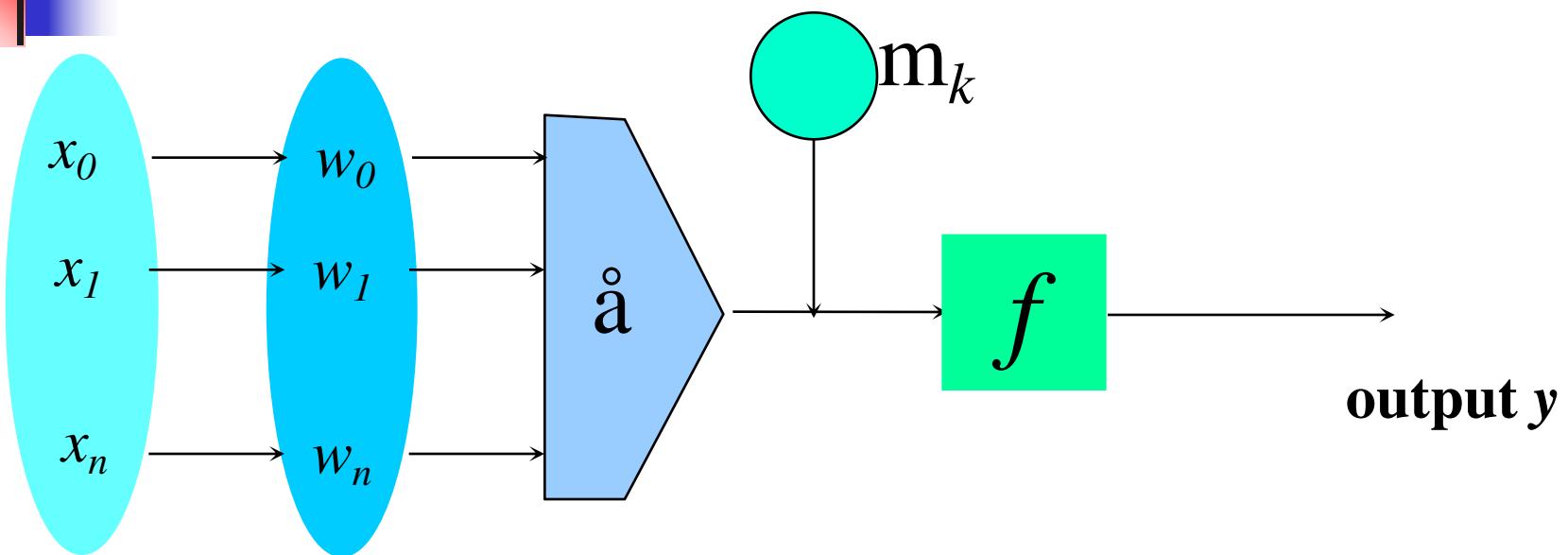
n Advantages

- n prediction accuracy is generally high
- n robust, works when training examples contain errors
- n output may be discrete, real-valued, or a vector of several discrete or real-valued attributes
- n fast evaluation of the learned target function

n Criticism

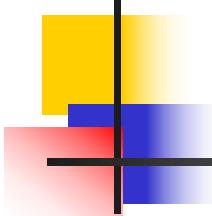
- n long training time
- n difficult to understand the learned function (weights)
- n not easy to incorporate domain knowledge

A Neuron



Input vector x **weight vector w** **weighted sum** **Activation function**

- n The n -dimensional input vector x is mapped into variable y by means of the scalar product and a nonlinear function mapping



Network Training

- n The ultimate objective of training
 - n obtain a set of weights that makes almost all the tuples in the training data classified correctly
- n Steps
 - n Initialize weights with random values
 - n Feed the input tuples into the network one by one
 - n For each unit
 - n Compute the net input to the unit as a linear combination of all the inputs to the unit
 - n Compute the output value using the activation function
 - n Compute the error
 - n Update the weights and the bias

Multi-Layer Perceptron

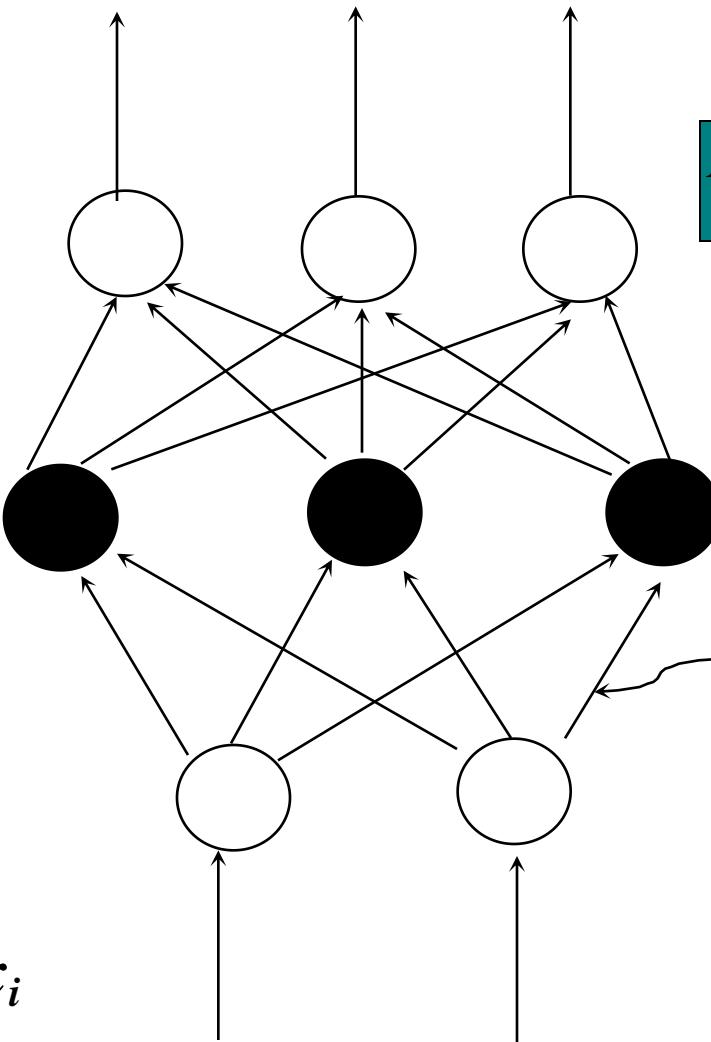
Output vector

Output nodes

Hidden nodes

Input nodes

Input vector: x_i



$$Err_j = O_j(1-O_j) \sum_k Err_k w_{jk}$$

$$\theta_j = \theta_j + (l) Err_j$$

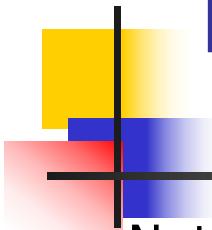
$$w_{ij} = w_{ij} + (l) Err_j O_i$$

$$Err_j = O_j(1-O_j)(T_j - O_j)$$

w_{ij}

$$O_j = \frac{1}{1+e^{-I_j}}$$

$$I_j = \sum_i w_{ij} O_i + \theta_j$$



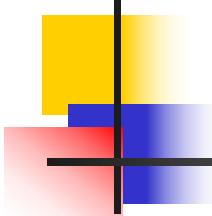
Network Pruning and Rule Extraction

n Network pruning

- n Fully connected network will be hard to articulate
- n N input nodes, h hidden nodes and m output nodes lead to $h(m+N)$ weights
- n Pruning: Remove some of the links without affecting classification accuracy of the network

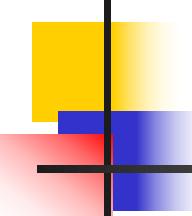
n Extracting rules from a trained network

- n Discretize activation values; replace individual activation value by the cluster average maintaining the network accuracy
- n Enumerate the output from the discretized activation values to find rules between activation value and output
- n Find the relationship between the input and activation value
- n Combine the above two to have rules relating the output to input



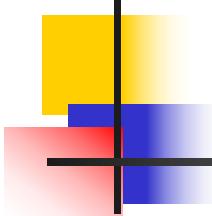
Chapter 7. Classification and Prediction

- n What is classification? What is prediction?
- n Issues regarding classification and prediction
- n Classification by decision tree induction
- n Bayesian Classification
- n Classification by backpropagation
- n **Classification based on concepts from association rule mining**
- n Other Classification Methods
- n Prediction
- n Classification accuracy
- n Summary



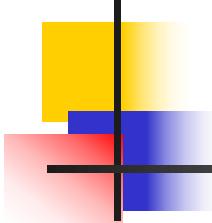
Association-Based Classification

- „ Several methods for association-based classification
 - „ ARCS: Quantitative association mining and clustering of association rules (Lent et al'97)
 - „ It beats C4.5 in (mainly) scalability and also accuracy
 - „ Associative classification: (Liu et al'98)
 - „ It mines high support and high confidence rules in the form of “cond_set => y”, where y is a class label
 - „ CAEP (Classification by aggregating emerging patterns) (Dong et al'99)
 - „ Emerging patterns (EPs): the itemsets whose support increases significantly from one class to another
 - „ Mine Eps based on minimum support and growth rate



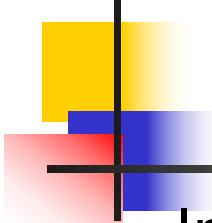
Chapter 7. Classification and Prediction

- n What is classification? What is prediction?
- n Issues regarding classification and prediction
- n Classification by decision tree induction
- n Bayesian Classification
- n Classification by backpropagation
- n Classification based on concepts from association rule mining
- n **Other Classification Methods**
- n Prediction
- n Classification accuracy
- n Summary



Other Classification Methods

- „ k-nearest neighbor classifier
- „ case-based reasoning
- „ Genetic algorithm
- „ Rough set approach
- „ Fuzzy set approaches



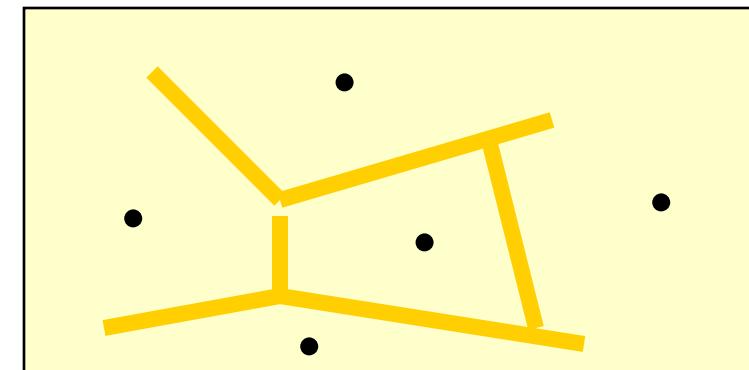
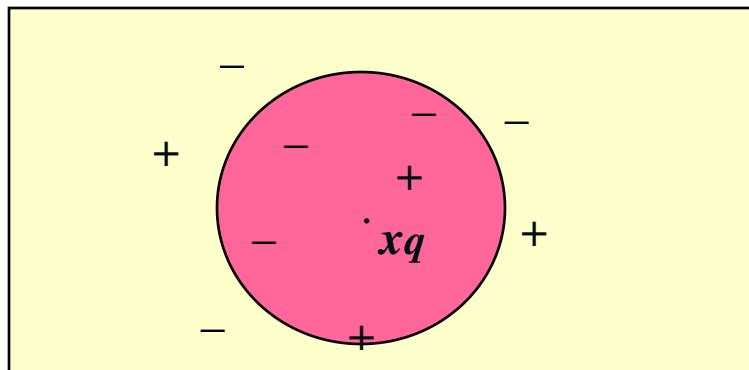
Instance-Based Methods

n Instance-based learning:

- n Store training examples and delay the processing (“lazy evaluation”) until a new instance must be classified
- n Typical approaches
 - n k -nearest neighbor approach
 - n Instances represented as points in a Euclidean space.
 - n Locally weighted regression
 - n Constructs local approximation
 - n Case-based reasoning
 - n Uses symbolic representations and knowledge-based inference

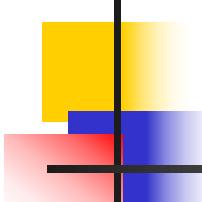
The k -Nearest Neighbor Algorithm

- All instances correspond to points in the n-D space.
- The nearest neighbor are defined in terms of Euclidean distance.
- The target function could be discrete- or real- valued.
- For discrete-valued, the k -NN returns the most common value among the k training examples nearest to x_q .
- Voronoi diagram: the decision surface induced by 1-NN for a typical set of training examples.



Discussion on the k -NN Algorithm

- n The k -NN algorithm for continuous-valued target functions
 - n Calculate the mean values of the k nearest neighbors
- n Distance-weighted nearest neighbor algorithm
 - n Weight the contribution of each of the k neighbors according to their distance to the query point x_q
 $w_i \equiv \frac{1}{d(x_q, x_i)^2}$
 - n giving greater weight to closer neighbors
 - n Similarly, for real-valued target functions
- n Robust to noisy data by averaging k -nearest neighbors
- n Curse of dimensionality: distance between neighbors could be dominated by irrelevant attributes.
 - n To overcome it, axes stretch or elimination of the least relevant attributes.

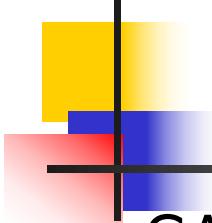


Case-Based Reasoning

- n Also uses: lazy evaluation + analyze similar instances
- n Difference: Instances are not “points in a Euclidean space”
- n Example: Water faucet problem in CADET (Sycara et al’92)
- n Methodology
 - n Instances represented by rich symbolic descriptions (e.g., function graphs)
 - n Multiple retrieved cases may be combined
 - n Tight coupling between case retrieval, knowledge-based reasoning, and problem solving
- n Research issues
 - n Indexing based on syntactic similarity measure, and when failure, backtracking, and adapting to additional cases

Remarks on Lazy vs. Eager Learning

- n Instance-based learning: lazy evaluation
- n Decision-tree and Bayesian classification: eager evaluation
- n Key differences
 - n Lazy method may consider query instance x_q when deciding how to generalize beyond the training data D
 - n Eager method cannot since they have already chosen global approximation when seeing the query
- n Efficiency: Lazy - less time training but more time predicting
- n Accuracy
 - n Lazy method effectively uses a richer hypothesis space since it uses many local linear functions to form its implicit global approximation to the target function
 - n Eager: must commit to a single hypothesis that covers the entire instance space

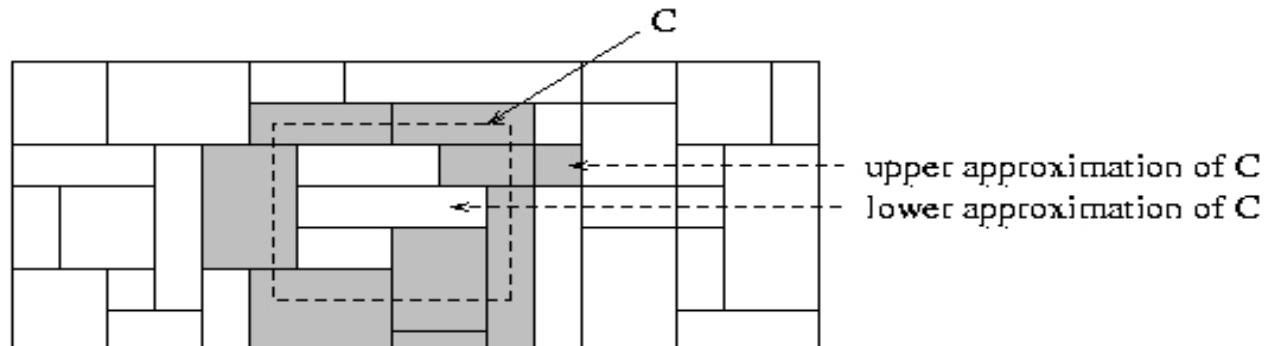


Genetic Algorithms

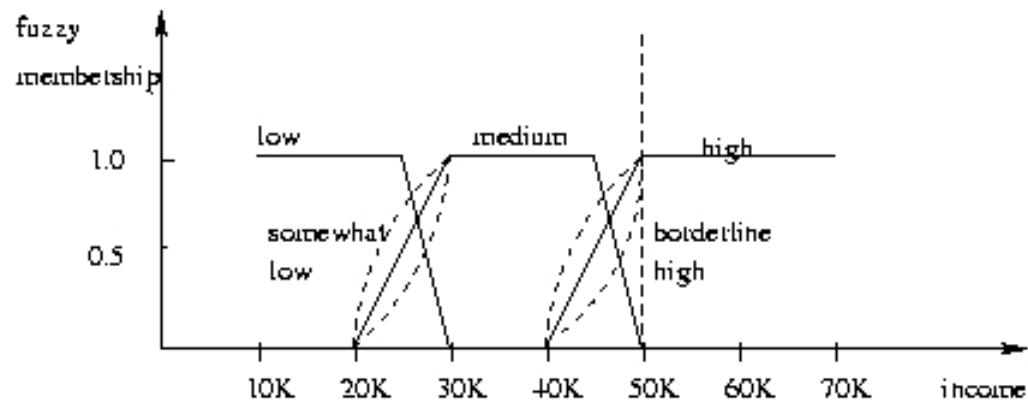
- GA: based on an analogy to biological evolution
- Each rule is represented by a string of bits
- An initial population is created consisting of randomly generated rules
 - e.g., IF A_1 and Not A_2 then C_2 can be encoded as 100
- Based on the notion of survival of the fittest, a new population is formed to consists of the fittest rules and their offsprings
- The fitness of a rule is represented by its classification accuracy on a set of training examples
- Offsprings are generated by crossover and mutation

Rough Set Approach

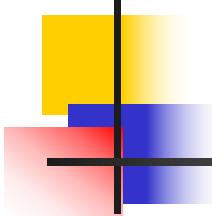
- n Rough sets are used to approximately or “roughly” define equivalent classes
- n A rough set for a given class C is approximated by two sets: a **lower approximation** (certain to be in C) and an **upper approximation** (cannot be described as not belonging to C)
- n Finding the minimal subsets (reducts) of attributes (for feature reduction) is NP-hard but a discernibility matrix is used to reduce the computation intensity



Fuzzy Set Approaches

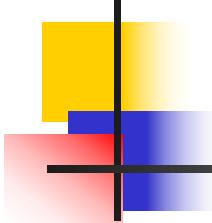


- n Fuzzy logic uses truth values between 0.0 and 1.0 to represent the degree of membership (such as using **fuzzy membership graph**)
- n Attribute values are converted to fuzzy values
 - n e.g., income is mapped into the discrete categories {low, medium, high} with fuzzy values calculated
- n For a given new sample, more than one fuzzy value may apply
- n Each applicable rule contributes a vote for membership in the categories
- n Typically, the truth values for each predicted category are summed



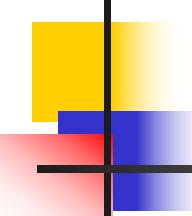
Chapter 7. Classification and Prediction

- „ What is classification? What is prediction?
- „ Issues regarding classification and prediction
- „ Classification by decision tree induction
- „ Bayesian Classification
- „ Classification by backpropagation
- „ Classification based on concepts from association rule mining
- „ Other Classification Methods
- „ **Prediction**
- „ Classification accuracy
- „ Summary



What Is Prediction?

- n Prediction is similar to classification
 - n First, construct a model
 - n Second, use model to predict unknown value
 - n Major method for prediction is regression
 - n Linear and multiple regression
 - n Non-linear regression
- n Prediction is different from classification
 - n Classification refers to predict categorical class label
 - n Prediction models continuous-valued functions



Predictive Modeling in Databases

- n Predictive modeling: Predict data values or construct generalized linear models based on the database data.
- n One can only predict value ranges or category distributions
- n Method outline:
 - n Minimal generalization
 - n Attribute relevance analysis
 - n Generalized linear model construction
 - n Prediction
- n Determine the major factors which influence the prediction
 - n Data relevance analysis: uncertainty measurement, entropy analysis, expert judgement, etc.

Regress Analysis and Log-Linear Models in Prediction

n Linear regression: $Y = \beta_0 + \beta_1 X$

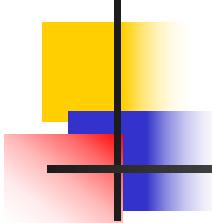
- n Two parameters , and specify the line and are to be estimated by using the data at hand.
- n using the least squares criterion to the known values of $Y_1, Y_2, \dots, X_1, X_2, \dots$

n Multiple regression: $Y = b_0 + b_1 X_1 + b_2 X_2$.

- n Many nonlinear functions can be transformed into the above.

n Log-linear models:

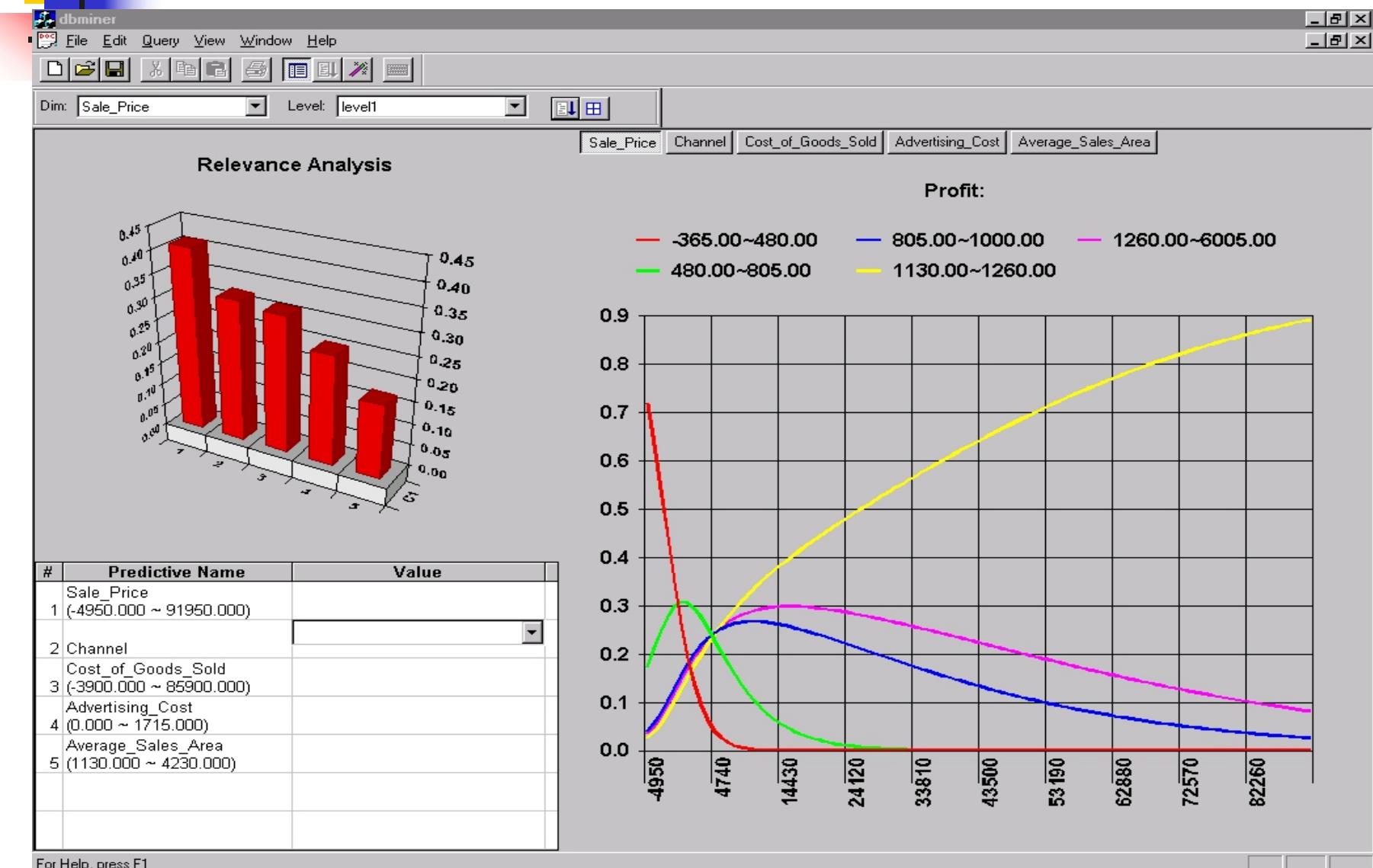
- n The multi-way table of joint probabilities is approximated by a product of lower-order tables.
- n Probability: $p(a, b, c, d) = ab ac ad bcd$



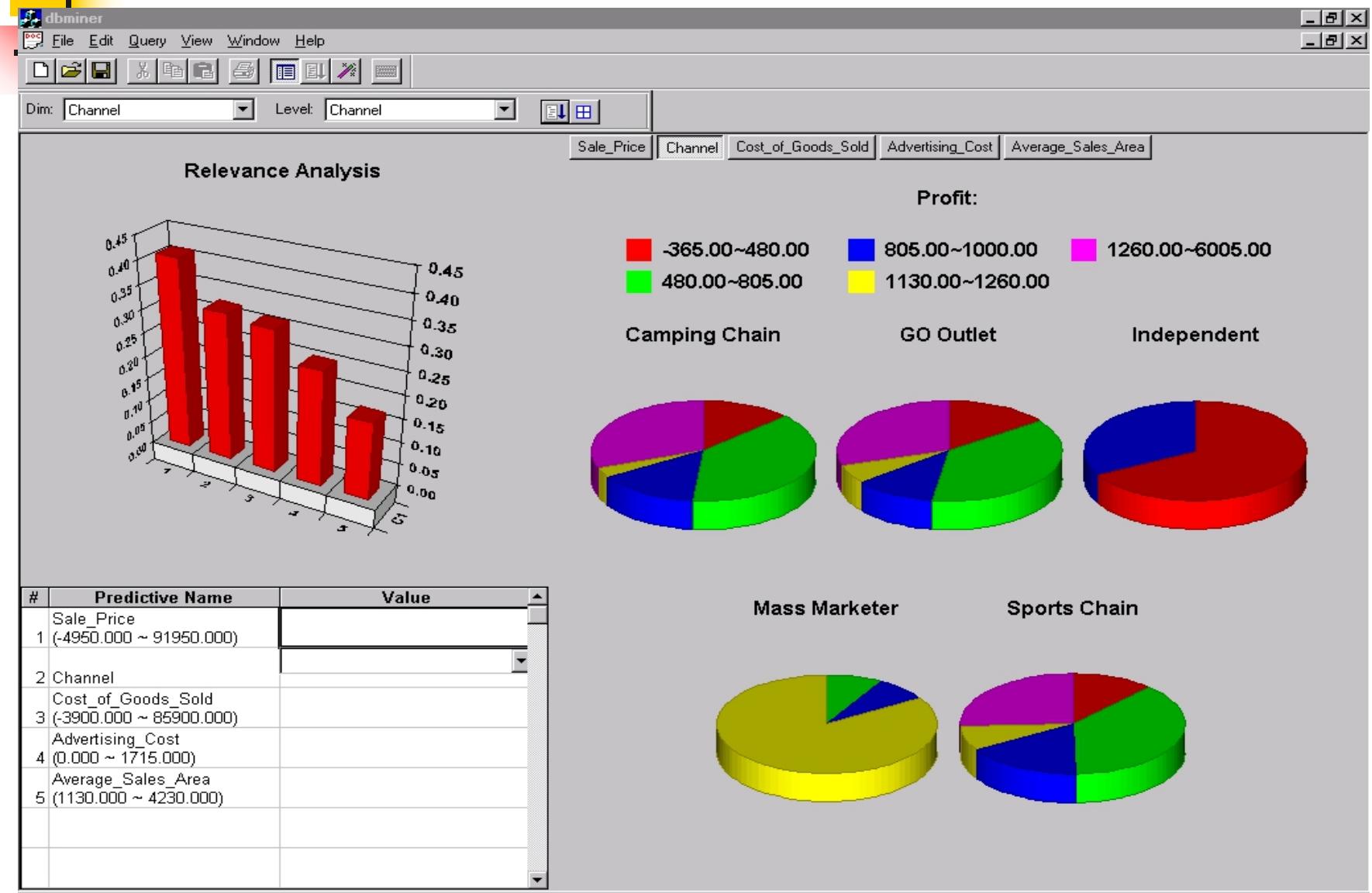
Locally Weighted Regression

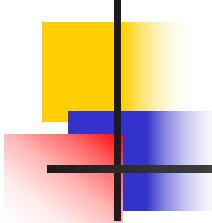
- n Construct an explicit approximation to f over a local region surrounding query instance x_q .
- n Locally weighted linear regression:
 - n The target function f is approximated near x_q using the linear function: $\hat{f}(x) = w_0 + w_1 a_1(x) + \dots + w_n a_n(x)$
 - n minimize the squared error: distance-decreasing weight K
$$E(x_q) \equiv \frac{1}{2} \sum_{x \in k_nearest_neighbors_of_x_q} K(d(x_q, x)) (f(x) - \hat{f}(x))^2$$
 - n the gradient descent training rule:
$$\Delta w_j \equiv \eta \sum_{x \in k_nearest_neighbors_of_x_q} K(d(x_q, x)) ((f(x) - \hat{f}(x)) a_j(x))$$
- n In most cases, the target function is approximated by a constant, linear, or quadratic function.

Prediction: Numerical Data



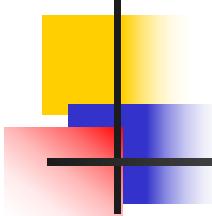
Prediction: Categorical Data





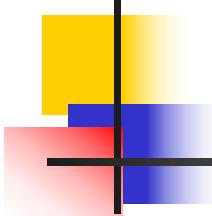
Chapter 7. Classification and Prediction

- „ What is classification? What is prediction?
- „ Issues regarding classification and prediction
- „ Classification by decision tree induction
- „ Bayesian Classification
- „ Classification by backpropagation
- „ Classification based on concepts from association rule mining
- „ Other Classification Methods
- „ Prediction
- „ **Classification accuracy**
- „ Summary



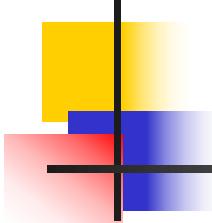
Classification Accuracy: Estimating Error Rates

- Partition: Training-and-testing
 - use two independent data sets, e.g., training set (2/3), test set(1/3)
 - used for data set with large number of samples
- Cross-validation
 - divide the data set into k subsamples
 - use $k-1$ subsamples as training data and one subsample as test data --- k -fold cross-validation
 - for data set with moderate size
- Bootstrapping (leave-one-out)
 - for small size data



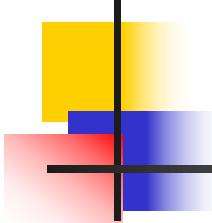
Boosting and Bagging

- Boosting increases classification accuracy
 - Applicable to decision trees or Bayesian classifier
- Learn a series of classifiers, where each classifier in the series pays more attention to the examples misclassified by its predecessor
- Boosting requires only linear time and constant space



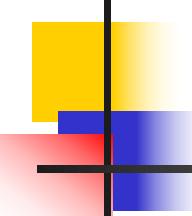
Boosting Technique (II) – Algorithm

- „ Assign every example an equal weight $1/N$
- „ For $t = 1, 2, \dots, T$ Do
 - „ Obtain a hypothesis (classifier) $h^{(t)}$ under $w^{(t)}$
 - „ Calculate the error of $h(t)$ and re-weight the examples based on the error
 - „ Normalize $w^{(t+1)}$ to sum to 1
- „ Output a weighted sum of all the hypothesis, with each hypothesis weighted according to its accuracy on the training set



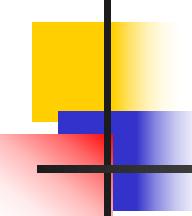
Chapter 7. Classification and Prediction

- n What is classification? What is prediction?
- n Issues regarding classification and prediction
- n Classification by decision tree induction
- n Bayesian Classification
- n Classification by backpropagation
- n Classification based on concepts from association rule mining
- n Other Classification Methods
- n Prediction
- n Classification accuracy
- n **Summary**



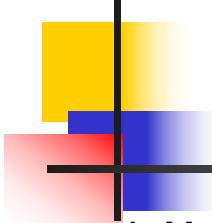
Summary

- n Classification is an **extensively studied** problem (mainly in statistics, machine learning & neural networks)
- n Classification is probably one of the most **widely used** data mining techniques with a lot of extensions
- n **Scalability** is still an important issue for database applications: thus combining classification **with database techniques** should be a promising topic
- n Research directions: classification of **non-relational data**, e.g., text, spatial, multimedia, etc..



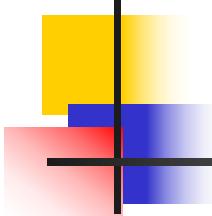
References (I)

- C. Apte and S. Weiss. Data mining with decision trees and decision rules. Future Generation Computer Systems, 13, 1997.
- L. Breiman, J. Friedman, R. Olshen, and C. Stone. Classification and Regression Trees. Wadsworth International Group, 1984.
- P. K. Chan and S. J. Stolfo. Learning arbiter and combiner trees from partitioned data for scaling machine learning. In Proc. 1st Int. Conf. Knowledge Discovery and Data Mining (KDD'95), pages 39-44, Montreal, Canada, August 1995.
- U. M. Fayyad. Branching on attribute values in decision tree generation. In Proc. 1994 AAAI Conf., pages 601-606, AAAI Press, 1994.
- J. Gehrke, R. Ramakrishnan, and V. Ganti. Rainforest: A framework for fast decision tree construction of large datasets. In Proc. 1998 Int. Conf. Very Large Data Bases, pages 416-427, New York, NY, August 1998.
- M. Kamber, L. Winstone, W. Gong, S. Cheng, and J. Han. Generalization and decision tree induction: Efficient classification in data mining. In Proc. 1997 Int. Workshop Research Issues on Data Engineering (RIDE'97), pages 111-120, Birmingham, England, April 1997.



References (II)

- J. Magidson. The Chaid approach to segmentation modeling: Chi-squared automatic interaction detection. In R. P. Bagozzi, editor, Advanced Methods of Marketing Research, pages 118-159. Blackwell Business, Cambridge Massachusetts, 1994.
- M. Mehta, R. Agrawal, and J. Rissanen. SLIQ : A fast scalable classifier for data mining. In Proc. 1996 Int. Conf. Extending Database Technology (EDBT'96), Avignon, France, March 1996.
- S. K. Murthy, Automatic Construction of Decision Trees from Data: A Multi-Diciplinary Survey, Data Mining and Knowledge Discovery 2(4): 345-389, 1998
- J. R. Quinlan. Bagging, boosting, and c4.5. In Proc. 13th Natl. Conf. on Artificial Intelligence (AAAI'96), 725-730, Portland, OR, Aug. 1996.
- R. Rastogi and K. Shim. Public: A decision tree classifier that integrates building and pruning. In Proc. 1998 Int. Conf. Very Large Data Bases, 404-415, New York, NY, August 1998.
- J. Shafer, R. Agrawal, and M. Mehta. SPRINT : A scalable parallel classifier for data mining. In Proc. 1996 Int. Conf. Very Large Data Bases, 544-555, Bombay, India, Sept. 1996.
- 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.



<http://www.cs.sfu.ca/~han>



Thank you !!!