





### **Session Agenda**

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
- Rule-based classification
- Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
- Model selection
- Summary

Icons / Meta	phors	
Z	Information	
	Common Realization	
	Knowledge/Competency Pattern	
	Governance	
and the second s	Alignment	
in the second	Solution Approach	
	5	i

Agend	a		
	1	Session Overview	and the second se
	2	Classification and Prediction	and the second se
	3	Summary and Conclusion	
			6

### **Classification and Prediction – Sub-Topics**

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
- Rule-based classification
- Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
- Model selection











## **Classification and Prediction – Sub-Topics**

- What is classification? What is prediction?
- Issues regarding classification and prediction
  - Classification by decision tree induction
  - Bayesian classification
  - Rule-based classification
  - Classification by back propagation
  - Support Vector Machines (SVM)
  - Lazy learners (or learning from your neighbors)
  - Frequent-pattern-based classification
  - Other classification methods
  - Prediction
  - Accuracy and error measures
  - Ensemble methods
  - Model selection

## Issues: 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: 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

### **Classification and Prediction – Sub-Topics**

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
  - Bayesian classification
  - Rule-based classification
  - Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
- Model selection

## Decision Tree Induction: Training Dataset

	age	income	student	credit_rating	buys_computer
	<=30	high	no	fair	no
	<=30	high	no	excellent	no
This follows	3140	high	no	fair	yes
n overnle	>40	medium	no	fair	yes
an example	>40	low	yes	fair	yes
of Quinian's	>40	low	yes	excellent	no
ID3 (Playing	3140	low	yes	excellent	yes
Tennis)	<=30	medium	no	fair	no
	<=30	low	yes	fair	yes
	>40	medium	yes	fair	yes
	<=30	medium	yes	excellent	yes
	3140	medium	no	excellent	yes
	3140	high	yes	fair	yes
	>40	medium	no	excellent	no



### **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



### **Attribute Selection: Information Gain**





### 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)

$$SplitInfo_A(D) = -\sum_{j=1}^{\nu} \frac{|D_j|}{|D|} \times \log_2(\frac{|D_j|}{|D|})$$

» GainRatio(A) = Gain(A)/SplitInfo(A)

• Ex. SplitInfo<sub>A</sub>(D) = 
$$-\frac{4}{14} \times \log_2(\frac{4}{14}) - \frac{6}{14} \times \log_2(\frac{6}{14}) - \frac{4}{14} \times \log_2(\frac{4}{14}) = 0.926$$
  
» gain\_ratio(income) = 0.029/0.926 = 0.031

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

### Gini index (CART, IBM IntelligentMiner)

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

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

where  $p_i$  is the relative frequency  $\overline{of}^1$  class *j* in *D* 

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

gini <sub>A</sub>(D) = 
$$\frac{|D_1|}{|D|}$$
gini (D<sub>1</sub>) +  $\frac{|D_2|}{|D|}$ gini (D<sub>2</sub>)

Reduction in Impurity: 

 $\Delta gini(A) = gini(D) - gini_A(D)$ • The attribute provides the smallest gini<sub>split</sub>(D) (or the largest reduction in impurity) is chosen to split the node (need to enumerate all the possible splitting points for each attribute)

### Gini index (CART, IBM IntelligentMiner)

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<sub>1</sub>: {low, medium} and 4 in D<sub>2</sub>

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

but gini{medium,high} is 0.30 and thus the best since it is the lowest

- 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



## Ťı

### **Other Attribute Selection Measures**

- <u>CHAID</u>: a popular decision tree algorithm, measure based on χ<sup>2</sup> test for independence
- <u>C-SEP</u>: performs better than info. gain and gini index in certain cases
- <u>G-statistic</u>: has a close approximation to χ<sup>2</sup> distribution
- <u>MDL (Minimal Description Length) principle</u> (i.e., the simplest solution is preferred):
  - The best tree as the one that requires the fewest # of bits to both (1) encode the tree, and (2) encode the exceptions to the tree
- Multivariate splits (partition based on multiple variable combinations)
  - » <u>CART</u>: finds multivariate splits based on a linear comb. of attrs.
- Which attribute selection measure is the best?
  - » Most give good results, none is significantly superior than others

# 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"



### Classification in Large Databases







Rainforest: Training Set and Its AVC Sets													
		ŀ	AVC-set on Age			AVC-set on <i>income</i>							
Training Examples						Age	Buy_0	Computer		income	Buy_	Computer	]
age	income	student	redit rating	com			yes	no			yes	no	
<=30	high	no	fair	no		<=30	3	2		high	2	2	
<=30	high	no	excellent	no		314			medium	4	2		
3140	high	no	fair	yes		0	4	0	0	low	3	1	1
>40	medium	no	fair	yes		>40	3	2		1011	Ŭ	•	J
>40	low	yes	fair	yes									
>40	low	yes	excellent	no	AVC-set on <i>Student</i> AVC-set or						et on	n	
3140	low	yes	excellent	yes		AVC-SEL OIT SLUDER					edit rating		
<=30	medium	no	fair	no	aturiant Dury Computer			_					
<=30	low	yes	fair	yes	l s	student Buy_C		Buy_Computer		Credit	Bu	ıy_Compu	ter
>40	medium	yes	fair	yes			yes	no		rating	yes	no	
<=30	medium	yes	excellent	yes		yes	6	1		fair	6	2	
3140	medium	no	excellent	yes		no	3	4		excellent	3	3	
3140	high	yes	fair	yes									
>40	medium	no	excellent	no									
												3	33

### 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

34









### **Classification and Prediction – Sub-Topics**

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
  - Rule-based classification
  - Classification by back propagation
  - Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
- Model selection

## **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 40









## Naïve Bayesian Classifier: Training Dataset



	age	income	student	credit_rating	_com
	<=30	high	no	fair	no
	<=30	high	no	excellent	no
Class:	3140	high	no	fair	yes
C1:buys_computer = 'yes'	>40	medium	no	fair	yes
C2:buys_computer = `no'	>40	low	yes	fair	yes
Data assessa	>40	low	yes	excellent	no
Data sample $X = (200 c - 20)$	3140	low	yes	excellent	yes
$\Lambda = (age < -50,$ Income = medium	<=30	medium	no	fair	no
Student = ves	<=30	low	yes	fair	yes
Credit rating = Fair)	>40	medium	yes	fair	yes
_ 5 ,	<=30	medium	yes	excellent	yes
	3140	medium	no	excellent	yes
	3140	high	yes	fair	yes
	>40	medium	no	excellent	no
					45

Naïve Bayesian Classifier: An Example	
<ul> <li>P(C<sub>i</sub>): P(buys_computer = "yes") = 9/14 = 0.643</li> <li>P(buys_computer = "no") = 5/14= 0.357</li> </ul>	
<ul> <li>Compute P(X C<sub>i</sub>) for each class         P(age = "&lt;=30"   buys_computer = "yes") = 2/9 = 0.222         P(age = "&lt;= 30"   buys_computer = "no") = 3/5 = 0.6         P(income = "medium"   buys_computer = "yes") = 4/9 = 0.444         P(income = "medium"   buys_computer = "no") = 2/5 = 0.4         P(student = "yes"   buys_computer = "yes) = 6/9 = 0.667         P(student = "yes"   buys_computer = "no") = 1/5 = 0.2         P(credit_rating = "fair"   buys_computer = "yes") = 6/9 = 0.667         P(credit_rating = "fair"   buys_computer = "no") = 2/5 = 0.4</li></ul>	
<ul> <li>X = (age &lt;= 30, income = medium, student = yes, credit_rating = fair)</li> </ul>	
P(X C <sub>i</sub> ) : P(X buys_computer = "yes") = 0.222 x 0.444 x 0.667 x 0.667 = 0.044 P(X buys_computer = "no") = 0.6 x 0.4 x 0.2 x 0.4 = 0.019 P(X C <sub>i</sub> )*P(C <sub>i</sub> ) : P(X buys_computer = "yes") * P(buys_computer = "yes") = 0.028 P(X buys_computer = "no") * P(buys_computer = "no") = 0.007	
Therefore, X belongs to class ("buys_computer = yes")	46













### **Using IF-THEN Rules for Classification**







### **Rule Induction: Sequential Covering Method**



- Sequential covering algorithm: Extracts rules directly from training data
- Typical sequential covering algorithms: FOIL, AQ, CN2, RIPPER
- Rules are learned *sequentially*, each for a given class C<sub>i</sub> will cover many tuples of C<sub>i</sub> but none (or few) of the tuples of other classes
- Steps:
  - » Rules are learned one at a time
  - » Each time a rule is learned, the tuples covered by the rules are removed
  - The process repeats on the remaining tuples unless *termination condition*, e.g., when no more training examples or when the quality of a rule returned is below a user-specified threshold
- Comp. w. decision-tree induction: learning a set of rules *simultaneously*



### How to Learn-One-Rule?

- Start with the most general rule possible: condition = empty
- Adding new attributes by adopting a greedy depth-first strategy
   » Picks the one that most improves the rule quality
- Rule-Quality measures: consider both coverage and accuracy
  - » Foil-gain (in FOIL & RIPPER): assesses info\_gain by extending condition

$$FOIL\_Gain=pos \times (\log_2 \frac{pos'}{pos+neg'} - \log_2 \frac{pos}{pos+neg})$$

It favors rules that have high accuracy and cover many positive tuples

Rule pruning based on an independent set of test tuples

$$FOIL\_Prune(R) = \frac{pos - neg}{pos + neg}$$

Pos/neg are # of positive/negative tuples covered by R. If *FOIL\_Prune* is higher for the pruned version of R, prune R



T,



- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
- Rule-based classification
- Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
- Model selection









# Classification by Backpropagation Backpropagation: A neural network learning algorithm Started by psychologists and neurobiologists to develop and test computational analogues of neurons A neural network: A set of connected input/output units where each connection has a weight associated with it During the learning phase, the network learns by adjusting the weights so as to be able to predict the correct class label of the input tuples Also referred to as connectionist learning due to the connections between units

### **Neural Network as a Classifier**

### Weakness

- » Long training time
- » Require a number of parameters typically best determined empirically, e.g., the network topology or "structure."
- » Poor interpretability: Difficult to interpret the symbolic meaning behind the learned weights and of "hidden units" in the network

### Strength

- » High tolerance to noisy data
- » Ability to classify untrained patterns
- » Well-suited for continuous-valued inputs and outputs
- » Successful on a wide array of real-world data
- » Algorithms are inherently parallel
- » Techniques have recently been developed for the extraction of rules from trained neural networks





# How A Multi-Layer Neural Network Works? The inputs to the network correspond to the attributes measured for each training tuple Inputs are fed simultaneously into the units making up the input layer They are then weighted and fed simultaneously to a hidden layer The number of hidden layers is arbitrary, although usually only one The weighted outputs of the last hidden layer are input to units making up the output layer, which emits the network's prediction The network is feed-forward in that none of the weights cycles back to an input unit or to an output unit of a previous layer From a statistical point of view, networks perform nonlinear regression: Given enough hidden units and enough training samples, they can closely approximate any function

### Defining a Network Topology

- First decide the network topology: # of units in the input layer, # of hidden layers (if > 1), # of units in each hidden layer, and # of units in the output layer
- Normalizing the input values for each attribute measured in the training tuples to [0.0—1.0]
- One input unit per domain value, each initialized to 0
- **Output**, if for classification and more than two classes, one output unit per class is used
- Once a network has been trained and its accuracy is unacceptable, repeat the training process with a different network topology or a different set of initial weights

### Backpropagation

- Iteratively process a set of training tuples & compare the network's prediction with the actual known target value
- For each training tuple, the weights are modified to minimize the mean squared error between the network's prediction and the actual target value
- Modifications are made in the "backwards" direction: from the output layer, through each hidden layer down to the first hidden layer, hence "backpropagation"
- Steps
  - » Initialize weights (to small random #s) and biases in the network
  - » Propagate the inputs forward (by applying activation function)
  - » Backpropagate the error (by updating weights and biases)
  - » Terminating condition (when error is very small, etc.)

### **Backpropagation and Interpretability**

- Efficiency of backpropagation: Each epoch (one interation through the training set) takes O(|D| \* w), with |D| tuples and w weights, but # of epochs can be exponential to n, the number of inputs, in the worst case
- Rule extraction from networks: network pruning
  - Simplify the network structure by removing weighted links that have the least effect on the trained network
  - » Then perform link, unit, or activation value clustering
  - The set of input and activation values are studied to derive rules describing the relationship between the input and hidden unit layers
- Sensitivity analysis: assess the impact that a given input variable has on a network output. The knowledge gained from this analysis can be represented in rules

## **Classification and Prediction – Sub-Topics** What is classification? What is prediction? Issues regarding classification and prediction Classification by decision tree induction Bayesian classification Rule-based classification Classification by back propagation Support Vector Machines (SVM) . Lazy learners (or learning from your neighbors) Frequent-pattern-based classification Other classification methods Prediction Accuracy and error measures Ensemble methods Model selection
## SVM—Support Vector Machines

- A new classification method for both <u>linear and nonlinear</u> data
- It uses a <u>nonlinear mapping</u> to transform the original training data into a higher dimension
- With the new dimension, it searches for the linear optimal separating hyperplane (i.e., "decision boundary")
- With an appropriate nonlinear mapping to a sufficiently high dimension, data from two classes can always be separated by a hyperplane
- SVM finds this hyperplane using support vectors ("essential" training tuples) and margins (defined by the support vectors)

## **SVM**—History and Applications

- Vapnik and colleagues (1992)—groundwork from Vapnik
   & Chervonenkis' statistical learning theory in 1960s
- Features: training can be slow but accuracy is high owing to their ability to model complex nonlinear decision boundaries (margin maximization)
- Used both for classification and prediction
- Applications:
  - » handwritten digit recognition, object recognition, speaker identification, benchmarking time-series prediction tests















- The support vectors are the essential or critical training examples they lie closest to the decision boundary (MMH)
- If all other training examples are removed and the training is repeated, the same separating hyperplane would be found
- The number of support vectors found can be used to compute an (upper) bound on the expected error rate of the SVM classifier, which is independent of the data dimensionality
- Thus, an SVM with a small number of support vectors can have good generalization, even when the dimensionality of the data is high



sider the following example. A 3D input vector  $\mathbf{X} = (x_1, x_2, x_3)$  is mapped into a 6D space Z using the mappings  $\phi_1(\mathbf{X}) = x_1, \phi_2(\mathbf{X}) = x_2, \phi_3(\mathbf{X}) = x_3, \phi_4(\mathbf{X}) = (x_1)^2, \phi_5(\mathbf{X}) = x_1x_2$ , and  $\phi_6(\mathbf{X}) = x_1x_3$ . A decision hyperplane in the new space is  $d(\mathbf{Z}) = \mathbf{WZ} + b$ , where W and Z are vectors. This is linear. We solve for W and b and then substitute back so that we see that the linear decision hyperplane in the new (Z) space corresponds to a nonlinear second order polynomial in the original 3-D input space,

$$\begin{aligned} d(Z) &= w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 (x_1)^2 + w_5 x_1 x_2 + w_6 x_1 x_3 + b \\ &= w_1 z_1 + w_2 z_2 + w_3 z_3 + w_4 z_4 + w_5 z_5 + w_6 z_6 + b \end{aligned}$$

Search for a linear separating hyperplane in the new space





## **CB-SVM: Clustering-Based SVM**

- Training data sets may not even fit in memory
- Read the data set once (minimizing disk access)
  - » Construct a statistical summary of the data (i.e., hierarchical clusters) given a limited amount of memory
  - » The statistical summary maximizes the benefit of learning SVM
- The summary plays a role in indexing SVMs
- Essence of Micro-clustering (Hierarchical indexing structure)
  - » Use micro-cluster hierarchical indexing structure
    - provide finer samples closer to the boundary and coarser samples farther from the boundary
  - » Selective de-clustering to ensure high accuracy









riment on a L	arge Data Se	et		
S-Rate	# of data	# of errors	T-Time	S-Time
0.0001%	23	6425	0.000114	822.97
0.001%	226	2413	0.000972	825.40
0.01%	2333	1132	0.03	828.61
0.1%	23273	1012	6.287	835.87
1%	230380	1015	1192.793	838.92
5%	1151714	1020	20705.4	842.92
ASVM	307	865	5487	2.213
CB-SVM	2893	876	1.639	2528.213

Table 4: Performance results on the very large data set (# of training data = 23066169, # of testing data = 233890). S-Rate: sampling rate; T-Time: training time; S-Time: sampling time; ASVM: selective sampling

#### SVM vs. Neural Network

# SVM

- » Relatively new concept
- » Deterministic algorithm
- » Nice Generalization properties
- Hard to learn learned in batch mode using quadratic programming techniques
- » Using kernels can learn very complex functions

# Neural Network

- » Relatively old
- » Nondeterministic algorithm
- Generalizes well but doesn't have strong mathematical foundation
- » Can easily be learned in incremental fashion
- To learn complex functions—use multilayer perceptron (not that trivial)



### Notes about SVM - Introductory Literature

- "Statistical Learning Theory" by Vapnik: difficult to understand, containing many errors.
- C. J. C. Burges. <u>A Tutorial on Support Vector Machines for Pattern</u> <u>Recognition</u>. *Knowledge Discovery and Data Mining*, 2(2), 1998.
  - » Easier than Vapnik's book, but still not introductory level; the examples are not so intuitive
- The book <u>An Introduction to Support Vector Machines</u> by Cristianini and Shawe-Taylor
  - » Not introductory level, but the explanation about Mercer's Theorem is better than above literatures

92

- <u>Neural Networks and Learning Machines</u> by Haykin
  - » Contains a nice chapter on SVM introduction

### **Classification and Prediction – Sub-Topics**

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
- Rule-based classification
- Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
- Model selection

# Lazy vs. Eager Learning

- Lazy vs. eager learning
  - » Lazy learning (e.g., instance-based learning): Simply stores training data (or only minor processing) and waits until it is given a test tuple
  - » Eager learning (the above discussed methods): Given a set of training set, constructs a classification model before receiving new (e.g., test) data to classify
- Lazy: less time in training but more time in predicting
- Accuracy
  - » 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
  - » Eager: must commit to a single hypothesis that covers the entire instance space

Lazy Learner: Instance-Based Methods	Ż
<ul> <li>Instance-based learning:         <ul> <li>Store training examples and delay the processing ("lazy evaluation") until a new instance must be classified</li> </ul> </li> <li>Typical approaches         <ul> <li><u>k-nearest neighbor approach</u></li> <li>Instances represented as points in a Euclidean space.</li> <li><u>Locally weighted regression</u></li> <li>Constructs local approximation</li> <li><u>Case-based reasoning</u></li> <li>Uses symbolic representations and knowledgebased inference</li> </ul> </li> </ul>	





### **Case-Based Reasoning (CBR)**

- **X**
- CBR: Uses a database of problem solutions to solve new problems
- Store <u>symbolic description</u> (tuples or cases)—not points in a Euclidean space
- <u>Applications:</u> Customer-service (product-related diagnosis), legal ruling
- Methodology
  - » Instances represented by rich symbolic descriptions (e.g., function graphs)
  - » Search for similar cases, multiple retrieved cases may be combined
  - » Tight coupling between case retrieval, knowledge-based reasoning, and problem solving
- Challenges
  - » Find a good similarity metric
  - Indexing based on syntactic similarity measure, and when failure, backtracking, and adapting to additional cases

### **Classification and Prediction – Sub-Topics**

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
- Rule-based classification
- Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
- Model selection

### **Associative Classification**

- Associative classification
  - » Association rules are generated and analyzed for use in classification
  - » Search for strong associations between frequent patterns (conjunctions of attribute-value pairs) and class labels
  - » Classification: Based on evaluating a set of rules in the form of

 $P_1 \wedge p_2 \dots \wedge p_l \rightarrow "A_{class} = C" (conf, sup)$ 

- Why effective?
  - » It explores highly confident associations among multiple attributes and may overcome some constraints introduced by decision-tree induction, which considers only one attribute at a time
  - In many studies, associative classification has been found to be more accurate than some traditional classification methods, such as C4.5



### A Closer Look at CMAR



- CMAR (Classification based on Multiple Association Rules: Li, Han, Pei, ICDM'01)
- Efficiency: Uses an enhanced FP-tree that maintains the distribution of class labels among tuples satisfying each frequent itemset
- Rule pruning whenever a rule is inserted into the tree
  - » Given two rules, R<sub>1</sub> and R<sub>2</sub>, if the antecedent of R<sub>1</sub> is more general than that of R<sub>2</sub> and conf(R<sub>1</sub>) ≥ conf(R<sub>2</sub>), then R<sub>2</sub> is pruned
  - » Prunes rules for which the rule antecedent and class are not positively correlated, based on a  $\chi^2$  test of statistical significance
- Classification based on generated/pruned rules
  - » If only one rule satisfies tuple X, assign the class label of the rule
  - » If a rule set S satisfies X, CMAR
    - · divides S into groups according to class labels
    - uses a weighted  $\chi^2$  measure to find the strongest group of rules, based on the statistical correlation of rules within a group
    - · assigns X the class label of the strongest group



# Associative Classification May Achieve High Accuracy and Efficiency (Cong et al. SIGMOD05)

Dataset	RCBT	CBA	IRG Classifier	C4.5 family			SVM
				single tree	bagging	boosting	
AML/ALL (ALL)	91.18%	91.18%	64.71%	91.18%	91.18%	91.18%	97.06%
Lung Cancer(LC)	97.99%	81.88%	89.93%	81.88%	96.64%	81.88%	96.64%
Ovarian Cancer(OC)	97.67%	93.02%	-	97.67%	97.67%	97.67%	97.67%
Prostate Cancer(PC)	97.06%	82.35%	88.24%	26.47%	26.47%	26.47%	79.41%
Average Accuracy	95.98%	87.11%	80.96%	74.3%	77.99%	74.3%	92.70%











# **Experimental Results**



107

Table 1. Accuracy by SVM on Frequent Combined Features vs. Single Features

Data	Si	ngle Fea	ture	Freq. 1	Pattern
	Item_All	$Item\_FS$	Item_RBF	Pat_Ali	$Pat_FS$
anneal	99.78	99.78	99.11	99.33	99.67
austral	85.01	85.50	85.01	81.79	91.14
auto	83.25	84.21	78.80	74.97	90.79
breast	97.46	97.46	96.98	96.83	97.78
cleve	84.81	84.81	85.80	78.55	95.04
diabetes	74.41	74.41	74.55	77.73	78.31
glass	75.19	75.19	74.78	79.91	81.32
heart	84.81	84.81	84.07	82.22	88.15
hepatic	84.50	89.04	85.83	81.29	96.83
horse	83.70	84.79	82.36	82.35	92.39
iono	93.15	94.30	92.61	89.17	95.44
iris	94.00	96.00	94.00	95.33	96.00
labor	89.99	91.67	91.67	94.99	95.00
lymph	81.00	81.62	84.29	83.67	96.67
pima	74.56	74.56	76.15	76.43	77.16
sonar	82.71	86.55	82.71	84.60	90.86
vehicle	70.43	72.93	72.14	73.33	76.34
wine	98.33	99.44	98.33	98.30	100
Z00	97.09	97.09	95.09	94.18	99.00

Dataset	Single 1	Features	Freque	nt Patterns
	Item_All	$Item\_FS$	$Pat\_All$	$Pat\_FS$
anneal	98.33	98.33	97.22	98.44
austral	84.53	84.53	84.21	88.24
auto	71.70	77.63	71.14	78.77
breast	95.56	95.56	95.40	96.35
cleve	80.87	80.87	80.84	91.42
diabetes	77.02	77.02	76.00	76.58
glass	75.24	75.24	76.62	79.89
heart	81.85	81.85	80.00	86.30
hepatic	78.79	85.21	80.71	93.04
horse	83.71	83.71	84.50	87.77
iono	92.30	92.30	92.89	94.87
iris	94.00	94.00	93.33	93.33
labor	86.67	86.67	95.00	91.67
lymph	76.95	77.62	74.90	83.67
pima	75.86	75.86	76.28	76.72
sonar	80.83	81.19	83.67	83.67
vehicle	70.70	71.49	74.24	73.06
wine	95.52	93.82	96.63	99.44
Z00	91.18	91.18	95.09	97.09

Table 2. Accuracy by C4.5 on Frequent Combined Features vs. Single Features

Scalability	Tests				
	Tabl	e 3. Accura	cy & Time	on Chess	Data
	$min\_sup$	#Patterns	Time (s)	SVM (%)	C4.5~(%)
	1	N/A	N/A	N/A	N/A
	2000	68,967	44.703	92.52	97.59
	2200	$28,\!358$	19.938	91.68	97.84
	2500	$6,\!837$	2.906	91.68	97.62
	2800	$1,\!031$	0.469	91.84	97.37
	3000	136	0.063	91.90	97.06
	Table	4. Accuracy	y & Time c	on Wavefor SVM (%)	m Data C4.5 (%)
	1	9,468,109	N/A	N/A	N/A
	80	26,576	176.485	92.40	88.35
	100	$15,\!316$	90.406	92.19	87.29
	150	5,408	23.610	91.53	88.80
	200	2.481	8.234	91.22	87.32





# Classification and Prediction – Sub-Topics

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
- Rule-based classification
- Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
- Model selection

### **Genetic Algorithms (GA)**

- Genetic Algorithm: based on an analogy to biological evolution
- An initial population is created consisting of randomly generated rules
  - » Each rule is represented by a string of bits
  - » E.g., if  $A_1$  and  $\neg A_2$  then  $C_2$  can be encoded as 100
  - » If an attribute has k > 2 values, k bits can be used
- Based on the notion of survival of the fittest, a new population is formed to consist 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
- The process continues until a population P evolves when each rule in P satisfies a prespecified threshold
- Slow but easily parallelizable

### **Rough Set Approach**



- 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)
- Finding the minimal subsets (reducts) of attributes for feature reduction is NP-hard but a discernibility matrix (which stores the differences between attribute values for each pair of data tuples) is used to reduce the computation intensity





# **Classification and Prediction – Sub-Topics**

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
- Rule-based classification
- Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
- Model selection



### **Linear Regression**

 <u>Linear regression</u>: involves a response variable y and a single predictor variable x

 $y = w_0 + w_1 x$ 

where  $w_0$  (y-intercept) and  $w_1$  (slope) are regression coefficients

Method of least squares: estimates the best-fitting straight line

$$w_{1} = \frac{\sum_{i=1}^{|D|} (x_{i} - \bar{x})(y_{i} - \bar{y})}{\sum_{i=1}^{|D|} (x_{i} - \bar{x})^{2}} \qquad w_{0} = \bar{y} - w_{1}\bar{x}$$

- <u>Multiple linear regression</u>: involves more than one predictor variable
  - » Training data is of the form  $(\mathbf{X}_1, \mathbf{y}_1), (\mathbf{X}_2, \mathbf{y}_2), ..., (\mathbf{X}_{|\mathsf{D}|}, \mathbf{y}_{|\mathsf{D}|})$
  - » Ex. For 2-D data, we may have:  $y = w_0 + w_1 x_1 + w_2 x_2$
  - » Solvable by extension of least square method or using SAS, S-Plus
  - » Many nonlinear functions can be transformed into the above



### **Other Regression-Based Models**

- Generalized linear model:
  - » Foundation on which linear regression can be applied to modeling categorical response variables
  - » Variance of y is a function of the mean value of y, not a constant
  - » Logistic regression: models the prob. of some event occurring as a linear function of a set of predictor variables
  - » <u>Poisson regression</u>: models the data that exhibit a Poisson distribution
- Log-linear models: (for categorical data)
  - » Approximate discrete multidimensional prob. distributions
  - » Also useful for data compression and smoothing
- Regression trees and model trees
  - » Trees to predict continuous values rather than class labels







# **Prediction: Categorical Data**



124

# Classification and Prediction – Sub-Topics

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
- Rule-based classification
- Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
  - Ensemble methods
  - Model selection

Classifier Accuracy Measures 🦸								
			Real c	lass\Predicted class	C <sub>1</sub>		~C <sub>1</sub>	
				C <sub>1</sub>	True positive		False negative	٦
		ſ		~C <sub>1</sub>	Fals	e positive	True negative	
ſ	Real class\Predicted class	buy_computer	= yes	buy_computer =	no	total	recognition(%)	٦
Ī	buy_computer = yes	6954		46		7000	99.34	
	buy_computer = no	412		2588		3000	86.27	1
Ī	total	7366		2634		10000	95.52	
•	<ul> <li>Accuracy of a cla correctly classifie</li> <li>Error rate (mise</li> <li>Given <i>m</i> classe in class <i>i</i> that a</li> <li>Alternative accur sensitivity = t-post specificity = t-neg, precision = t-post accuracy = sensit</li> <li>This model car</li> </ul>	ISSIFIER IM, ac ad by the modula of the module of the mo	c(M): p del M ate) of htry in a the cla es (e.g. true po true po true no s + neg for cos	M = 1 – acc(M) confusion massifier as class , for cancer d sitive recogniti egative recogniti egative recogni ) + specificity *	atrix, j iagno on rat tion ra neg/( sis	indicates osis) te */ ate */ (pos + ne	es that are s # of tuples eg)	5







![](_page_64_Figure_0.jpeg)

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
- Rule-based classification
- Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
  - Model selection

![](_page_64_Figure_15.jpeg)

# **Bagging: Boostrap Aggregation**

- Analogy: Diagnosis based on multiple doctors' majority vote
- Training
  - Siven a set D of d tuples, at each iteration i, a training set D<sub>i</sub> of d tuples is sampled with replacement from D (i.e., boostrap)
  - » A classifier model M<sub>i</sub> is learned for each training set D<sub>i</sub>
- Classification: classify an unknown sample X
  - » Each classifier M<sub>i</sub> returns its class prediction
  - $\,$  > The bagged classifier M\* counts the votes and assigns the class with the most votes to  ${\rm X}$

- Prediction: can be applied to the prediction of continuous values by taking the average value of each prediction for a given test tuple
- Accuracy
  - » Often significant better than a single classifier derived from D
  - » For noise data: not considerably worse, more robust
  - » Proved improved accuracy in prediction

![](_page_65_Figure_13.jpeg)

### Adaboost (Freund and Schapire, 1997)

- Given a set of *d* class-labeled tuples, (X<sub>1</sub>, y<sub>1</sub>), ..., (X<sub>d</sub>, y<sub>d</sub>)
- Initially, all the weights of tuples are set the same (1/d)
- Generate k classifiers in k rounds. At round i,
  - Tuples from D are sampled (with replacement) to form a training set D<sub>i</sub> of the same size
  - » Each tuple's chance of being selected is based on its weight
  - A classification model M<sub>i</sub> is derived from D<sub>i</sub>
  - » Its error rate is calculated using D<sub>i</sub> as a test set
  - » If a tuple is misclssified, its weight is increased, o.w. it is decreased
- Error rate: err(X<sub>j</sub>) is the misclassification error of tuple X<sub>j</sub>. Classifier M<sub>i</sub> error rate is the sum of the weights of the misclassified tuples:

$$error(M_i) = \sum_{j}^{d} w_j \times err(\mathbf{X}_j)$$

The weight of classifier M<sub>i</sub>'s vote is

 $\log \frac{1 - error(M_i)}{error(M_i)}$ 

### **Classification and Prediction – Sub-Topics**

- What is classification? What is prediction?
- Issues regarding classification and prediction
- Classification by decision tree induction
- Bayesian classification
- Rule-based classification
- Classification by back propagation
- Support Vector Machines (SVM)
- Lazy learners (or learning from your neighbors)
- Frequent-pattern-based classification
- Other classification methods
- Prediction
- Accuracy and error measures
- Ensemble methods
- Model selection

134

![](_page_67_Figure_0.jpeg)

C <sub>1</sub>	~C <sub>1</sub>
True Positives (TP)	False Negatives (FN
False Positives (FP)	True Negatives (TN)
correctly	classified,
TP + TN	
= $TP + TN + FP$	+FN
	$C_{1}$ True Positives (TP) False Positives (FP) cognition rate: percentage correctly $= \frac{TP + TN}{TP + TN + FP}$

# Classifier Evaluation Metrics: Example - Confusion Matrix

Actual class\Predicted class	buy_computer = yes	buy_computer = no	Total	Recognition(%)
buy_computer = yes	6954	46	7000	99.34
buy_computer = no	412	2588	3000	86.27
Total	7366	2634	1000 0	95.42

- Siven *m* classes, an entry, *CMi,j* in a confusion matrix indicates # of tuples in class *i* that were labeled by the classifier as class *j*.
- » May be extra rows/columns to provide totals or recognition rate per class.

![](_page_68_Figure_4.jpeg)

![](_page_69_Figure_0.jpeg)

Cla	lassifier Evaluation Metrics: Example									
						_				
	Actual class\Predicted class	cancer = yes	cancer = no	Total	Recognition( %)					
	cancer = yes	90	210	300	30.00 sensitivity					
	cancer = no	140	9560	9700	98.56 specificity					
	Total	230	9770	1000 0	96.40 accuracy					

Precision = 90/230 = 39.13%; Recall = 90/300 = 30.00%

![](_page_70_Figure_0.jpeg)

![](_page_70_Figure_1.jpeg)

![](_page_71_Figure_0.jpeg)

![](_page_71_Figure_1.jpeg)

- » Works well with small data sets
- » Samples the given training tuples uniformly with replacement
  - i.e., each time a tuple is selected, it is equally likely to be selected again and re-added to the training set
- Several bootstrap methods, and a common one is .632 boostrap
  - ≫ A data set with *d* tuples is sampled *d* times, with replacement, resulting in a training set of *d* samples. The data tuples that did not make it into the training set end up forming the test set. About 63.2% of the original data end up in the bootstrap, and the remaining 36.8% form the test set (since  $(1 1/d)^d \approx e^{-1} = 0.368$ )
  - Repeat the sampling procedure k times, overall accuracy of the model:

$$acc(M) = \sum_{i=1}^{k} (0.632 \times acc(M_i)_{test\_set} + 0.368 \times acc(M_i)_{train\_set})$$

![](_page_71_Figure_9.jpeg)








## Model Selection: ROC Curves

- ROC (Receiver Operating Characteristics) curves: for visual comparison of classification models
- Originated from signal detection theory
- Shows the trade-off between the true positive rate and the false positive rate
- The area under the ROC curve is a measure of the accuracy of the model
- Rank the test tuples in decreasing order: the one that is most likely to belong to the positive class appears at the top of the list
- The closer to the diagonal line (i.e., the closer the area is to 0.5), the less accurate is the model



- Horizontal axis rep. the false positive rate
- The plot also shows a diagonal line
- A model with perfect accuracy will have an area of 1.0

## **Issues Affecting Model Selection**

# Accuracy

» classifier accuracy: predicting class label

## 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

150

Agenda					
	1	Session Overview	15-10 P		
	2	Classification and Prediction	part the second s		
	3	Summary and Conclusion	and the second se		
			151		







## Summary (III)

#### Backpropagation:

- » <u>Neural network</u> algorithm that uses gradient descent
- Searches for a set of weights that model the data so as to minimize the error between the network's class prediction and the actual class label of data tuples
- » Rules can be extracted for improved interpretability

#### Support Vector Machine (SVM):

- » For classification of both linear and nonlinear data
- Transforms original data into a higher dimension, from where it finds a <u>hyperplane</u> for separation of the data using essential training tuples called **support vectors**.

#### Pattern-Based Classification:

- » Uses <u>association mining techniques</u> that search for frequently occurring patterns in large databases.
- The patterns may generate rules, which can be analyzed for use in classification.

## Summary (IV)

- Lazy Learners:
  - store all training tuples and wait until presented with a test tuple before performing generalization.
  - » k-nearest neighbor and case-based reasoning
- Genetic Algorithms: populations of rules "evolve" via operations of crossover and mutation until all rules within a population satisfy specified threshold.
- Rough Set Approach: approximately define classes that are not distinguishable based on the available attributes
- Fuzzy Set Approaches: replace ``brittle" threshold cutoffs for continuous-valued attributes with degree of membership functions.

155

### 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, <u>Discriminative Frequent Pattern Analysis</u> for Effective Classification, ICDE'07.
- H. Cheng, X. Yan, J. Han, and P. S. Yu, <u>Direct Discriminative Pattern Mining for</u> <u>Effective Classification</u>, 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.

157

References (2)			
<ul> <li>A. J. Dobson. An Introduction to Generalized Linear Models. Chapman &amp; Hall, 1990.</li> </ul>			
<ul> <li>G. Dong and J. Li. Efficient mining of emerging patterns: Discovering trends and differences. KDD'99.</li> </ul>			
R. O. Duda, P. E. Hart, and D. G. Stork. Pattern Classification, 2ed. John Wiley, 2001			
<ul> <li>U. M. Fayyad. Branching on attribute values in decision tree generation. AAAI'94.</li> </ul>			
<ul> <li>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.</li> </ul>			
<ul> <li>J. Gehrke, R. Ramakrishnan, and V. Ganti. Rainforest: A framework for fast decision tree construction of large datasets. VLDB'98.</li> </ul>			
<ul> <li>J. Gehrke, V. Gant, R. Ramakrishnan, and WY. Loh, BOAT Optimistic Decision Tree Construction. SIGMOD'99.</li> </ul>			
<ul> <li>T. Hastie, R. Tibshirani, and J. Friedman. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer-Verlag, 2001.</li> </ul>			
<ul> <li>D. Heckerman, D. Geiger, and D. M. Chickering. Learning Bayesian networks: The combination of knowledge and statistical data. Machine Learning, 1995.</li> </ul>			
<ul> <li>W. Li, J. Han, and J. Pei, CMAR: Accurate and Efficient Classification Based on Multiple Class-Association Rules, ICDM'01.</li> </ul>			
158			

## **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.

159

• J. R. Quinlan. Bagging, boosting, and c4.5. AAAI'96.

## 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. Indurkhya. 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. 160



Next Session: Cluster Analysis				
	162			