# Data Mining:

## Concepts and Techniques

### (3rd ed.)

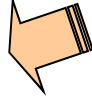### Classification: Advanced Methods

Jiawei Han, Micheline Kamber, and Jian Pei

University of Illinois at Urbana-Champaign &

Simon Fraser University
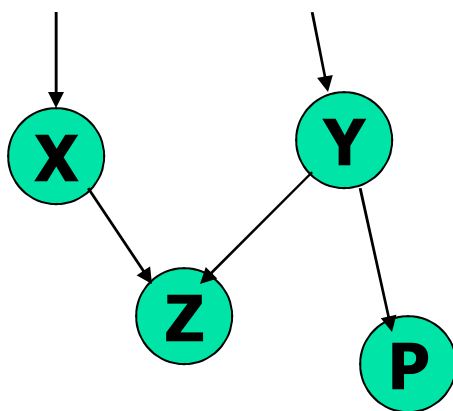
# Chapter 9. Classification: Advanced Methods

- Bayesian Belief Networks

- Support Vector Machines

- Lazy Learners (or Learning from Your Neighbors)

- Other Classification Methods

- Additional Topics Regarding Classification
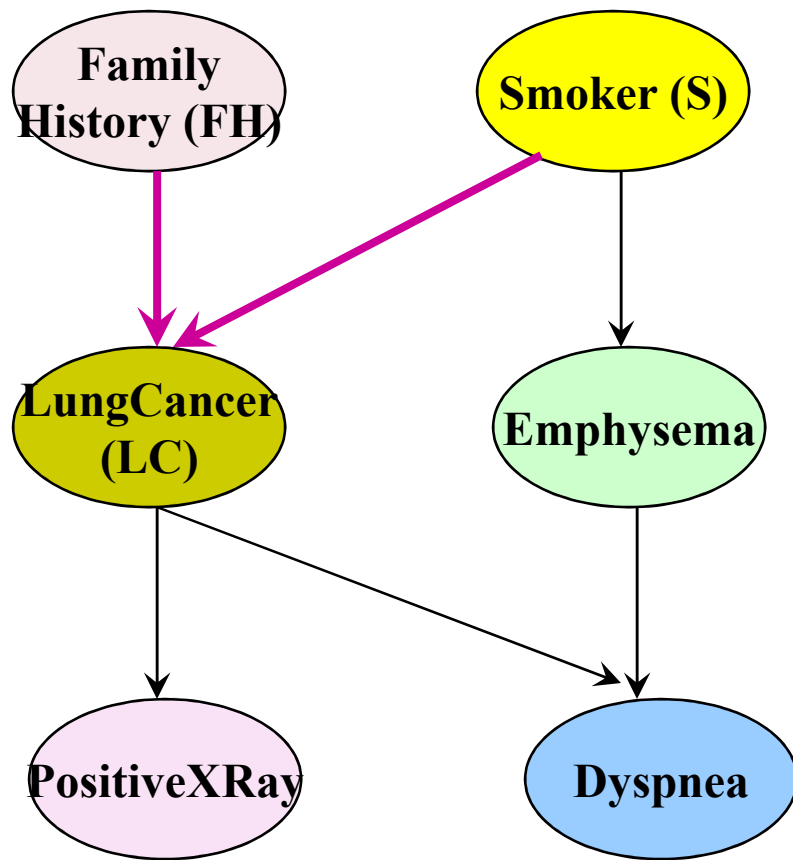
- Numerical Prediction

- Summary

# Bayesian Belief Networks

- **Bayesian belief networks** (also known as **Bayesian networks**, **probabilistic networks**): allow *class conditional independencies* between *subsets* of variables

- A (*directed acyclic*) graphical model of causal relationships
  - Represents <u>dependency</u> among the variables
  - Gives a specification of joint probability distribution



- ❑ Nodes: random variables
- ❑ Links: dependency
- ❑ X and Y are the parents of Z, and Y is the parent of P
- ❑ No dependency between Z and P
- ❑ Has no loops/cycles

# Bayesian Belief Network: An Example

**CPT**: **Conditional Probability Table** for variable LungCancer:

|       | (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       |

shows the conditional probability for each possible combination of its parents

Derivation of the probability of a particular combination of values of **X**, from CPT:

**Bayesian Belief Network**

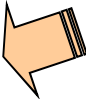$$P(x_1, ..., x_n) = \prod_{i=1}^{n} P(x_i \mid Parents(Y_i))$$

4

# Training Bayesian Networks: Several Scenarios

- Scenario 1:  Given both the network structure and all variables observable: *compute only the CPT entries*

- Scenario 2: Network structure known, some variables hidden: *gradient descent* (greedy hill-climbing) method, i.e., search for a solution along the steepest descent of a criterion function
  - Weights are initialized to random probability values
  - At each iteration, it moves towards what appears to be the best solution at the moment, w.o. backtracking
  - Weights are updated at each iteration & converge to local optimum

- Scenario 3: Network structure unknown, all variables observable: search through the model space to *reconstruct network topology*

- Scenario 4: Unknown structure, all hidden variables: No good algorithms known for this purpose

- D. Heckerman.  A Tutorial on Learning with Bayesian Networks.  In *Learning in Graphical Models,* M. Jordan, ed.. MIT Press, 1999.
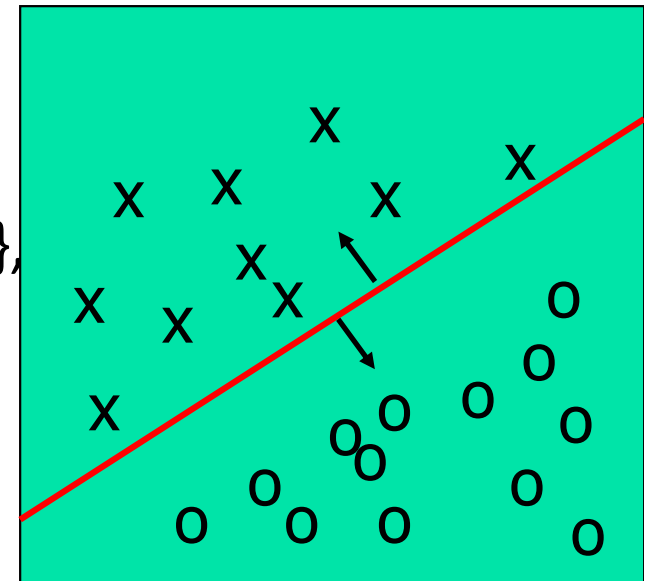
# Chapter 9. Classification: Advanced Methods

- Bayesian Belief Networks

- Support Vector Machines

- Lazy Learners (or Learning from Your Neighbors)

- Other Classification Methods

- Additional Topics Regarding Classification

- Numerical Prediction

- Summary

# Classification: A Mathematical Mapping

- Classification: predicts categorical class labels
  - E.g., Personal homepage classification
    - $x_i = (x_1, x_2, x_3, \dots)$, $y_i = +1$ or $-1$
    - $x_1$ : # of word "homepage"
    - $x_2$ : # of word "welcome"
- Mathematically, $x \in X = \Re^n$, $y \in Y = \{+1, -1\}$,
  - We want to derive a function f: X → Y
- Linear Classification
  - Binary Classification problem
  - Data above the red line belongs to class 'x'
  - Data below red line belongs to class 'o'
  - Examples: SVM, Perceptron, Probabilistic Classifiers

# Discriminative Classifiers

- Advantages
    - Prediction accuracy is generally high
        - As compared to Bayesian methods – in general
    - Robust, works when training examples contain errors
    - Fast evaluation of the learned target function
        - Bayesian networks are normally slow
- Criticism
    - Long training time
    - Difficult to understand the learned function (weights)
        - Bayesian networks can be used easily for pattern discovery
    - Not easy to incorporate domain knowledge
        - Easy in the form of priors on the data or distributions
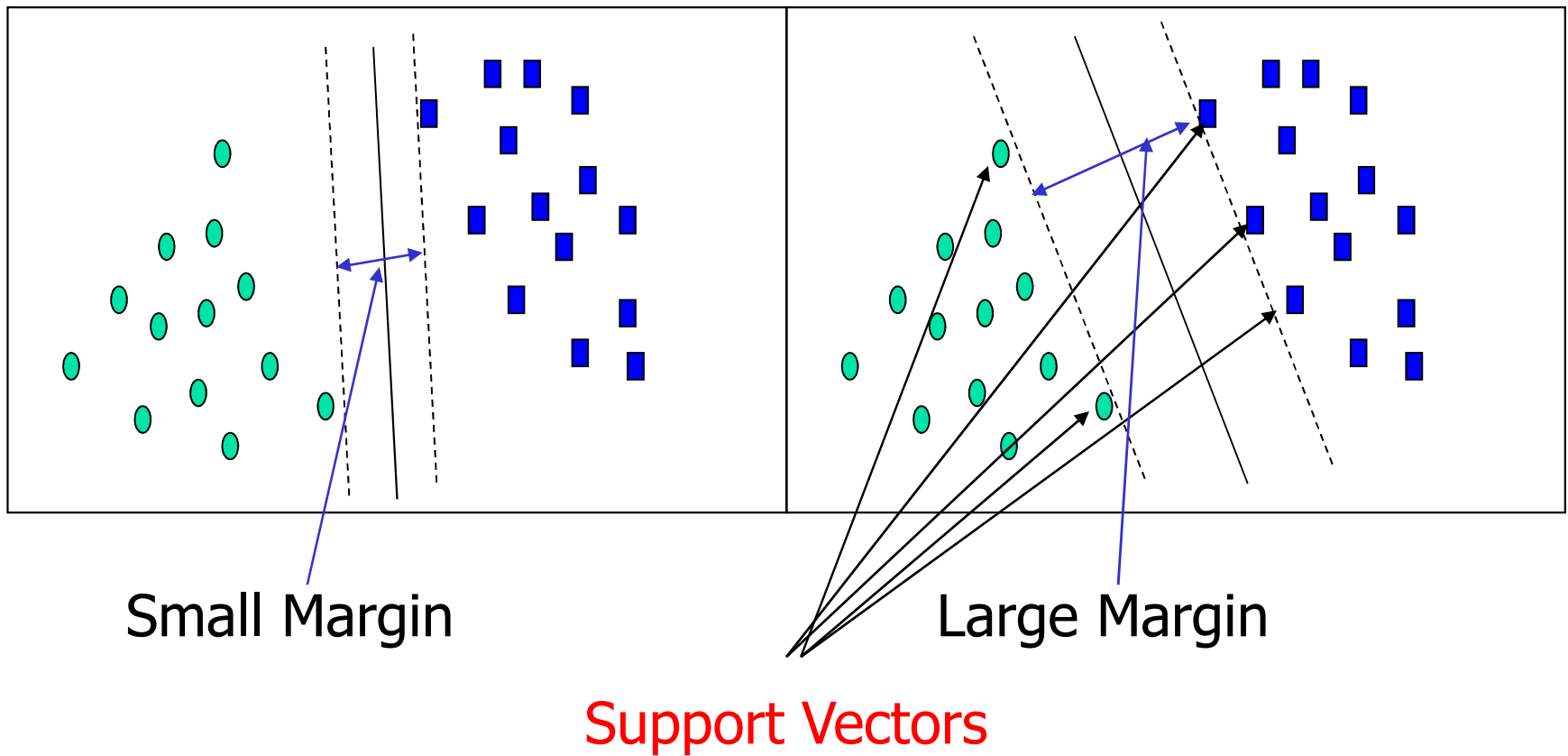
# SVM—Support Vector Machines

- A relatively 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)
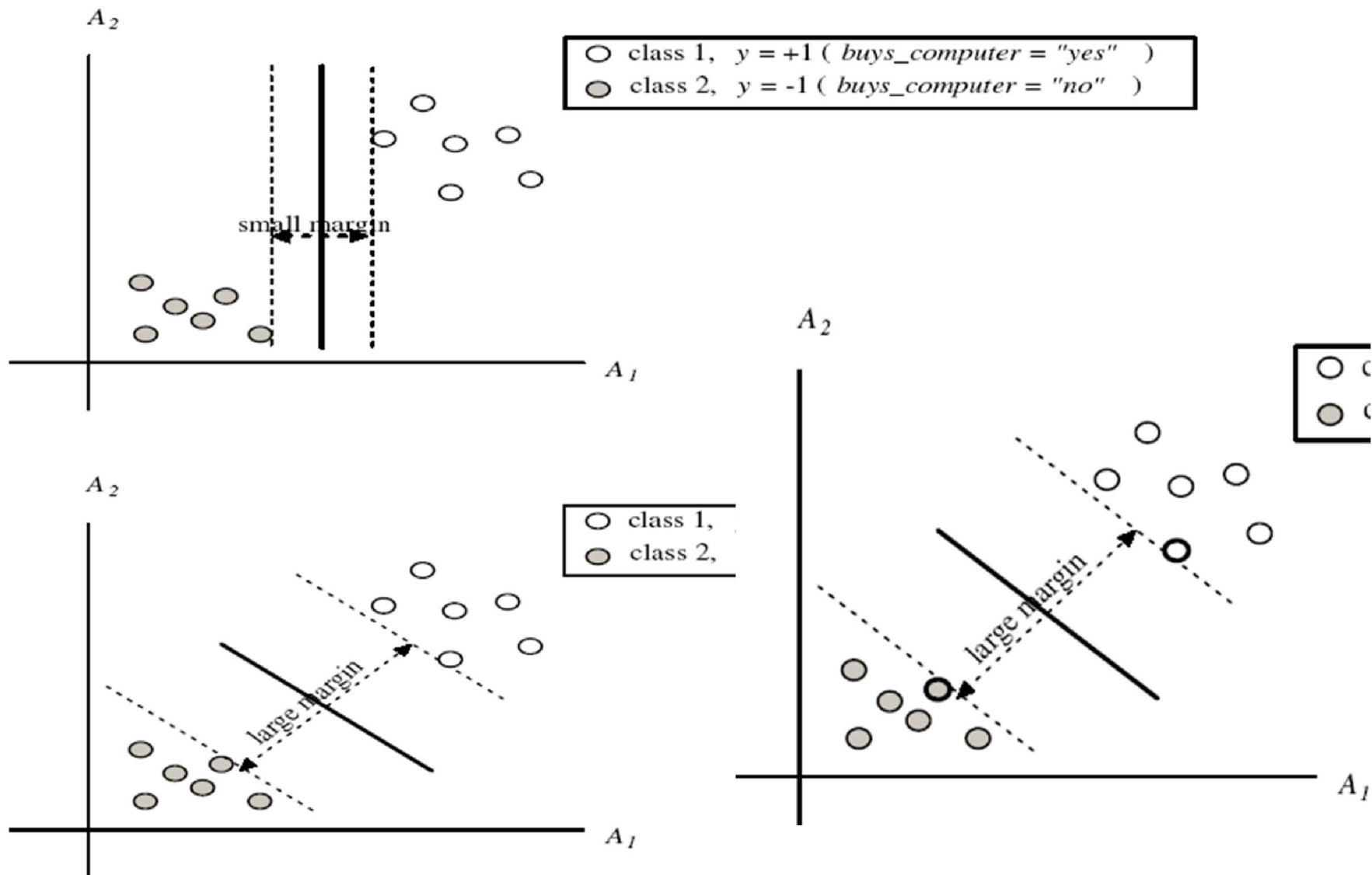
# SVM—History and Applications

- 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 for: classification and numeric prediction

- Applications:

    - handwritten digit recognition, object recognition, speaker identification, benchmarking time-series prediction tests

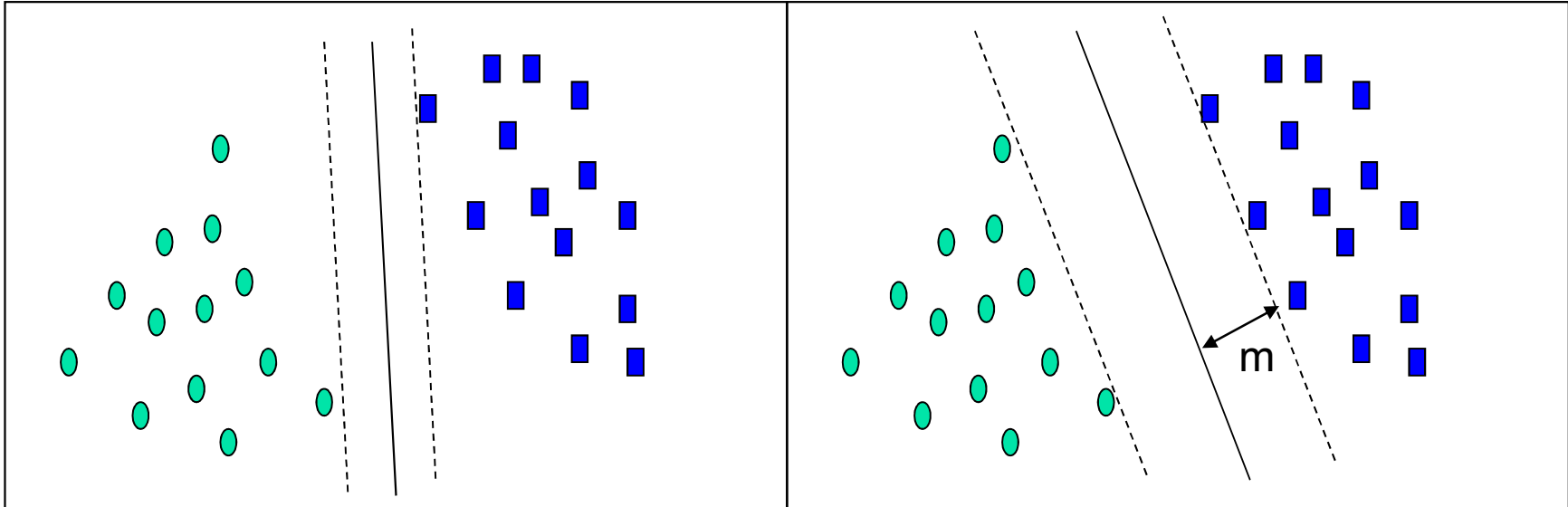# SVM—General Philosophy



Small Margin

Large Margin

Support Vectors

# SVM—Margins and Support Vectors



class 1, $y = +1$ ( *buys_computer = "yes"* )
class 2, $y = -1$ ( *buys_computer = "no"* )

small margin

large margin

class 1,
class 2,

# SVM—When Data Is Linearly Separable



Let data D be $(\mathbf{X}_1, y_1)$, ..., $(\mathbf{X}_{|D|}, y_{|D|})$, where $\mathbf{X}_i$ is the set of training tuples associated with the class labels $y_i$

There are infinite lines (<u>hyperplanes</u>) separating the two classes but we want to <u>find the best one</u> (the one that minimizes classification error on unseen data)

*SVM searches for the hyperplane with the largest margin*, i.e., **maximum marginal hyperplane** (MMH)

# SVM—Linearly Separable

- A separating hyperplane can be written as

  $$\mathbf{W} \bullet \mathbf{X} + b = 0$$

  where $\mathbf{W} = \{w_1, w_2, ..., w_n\}$ is a weight vector and $b$ a scalar (bias)

- For 2-D it can be written as

  $$w_0 + w_1 x_1 + w_2 x_2 = 0$$

- The hyperplane defining the sides of the margin:

  $$H_1: w_0 + w_1 x_1 + w_2 x_2 \geq 1 \quad \text{for } y_i = +1, \text{ and}$$
  $$H_2: w_0 + w_1 x_1 + w_2 x_2 \leq -1 \text{ for } y_i = -1$$
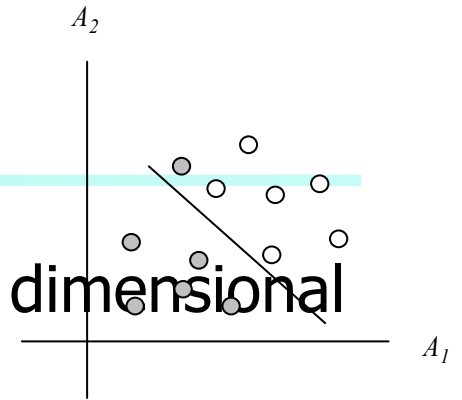
- Any training tuples that fall on hyperplanes $H_1$ or $H_2$ (i.e., the sides defining the margin) are **support vectors**

- This becomes a **constrained (convex) quadratic optimization** problem: Quadratic objective function and linear constraints → *Quadratic Programming (QP)* → Lagrangian multipliers

# Why Is SVM Effective on High Dimensional Data?

- The **complexity** of trained classifier is characterized by the # of support vectors rather than the dimensionality of the data

- 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

# SVM—Linearly Inseparable

- Transform the original input data into a higher dimensional space

$A_1$

Example 6.8 Nonlinear transformation of original input data into a higher dimensional space. Consider the following example. A 3D input vector $X = (x_1, x_2, x_3)$ is mapped into a 6D space $Z$ using the mappings $\phi_1(X) = x_1, \phi_2(X) = x_2, \phi_3(X) = x_3, \phi_4(X) = (x_1)^2, \phi_5(X) = x_1 x_2,$ and $\phi_6(X) = x_1 x_3$. A decision hyperplane in the new space is $d(Z) = WZ + b$, 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,

$$d(Z) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 (x_1)^2 + w_5 x_1 x_2 + w_6 x_1 x_3 + b$$
$$= w_1 z_1 + w_2 z_2 + w_3 z_3 + w_4 z_4 + w_5 z_5 + w_6 z_6 + b$$

∎

- Search for a linear separating hyperplane in the new space

# SVM: Different Kernel functions

- Instead of computing the dot product on the transformed data, it is math. equivalent to applying a kernel function $K(\mathbf{X_i}, \mathbf{X_j})$ to the original data, i.e., $K(\mathbf{X_i}, \mathbf{X_j}) = \Phi(\mathbf{X_i}).\Phi(\mathbf{X_j})$

- Typical Kernel Functions

$$\text{Polynomial kernel of degree } h: \quad K(X_i, X_j) = (X_i \cdot X_j + 1)^h$$

$$\text{Gaussian radial basis function kernel}: \quad K(X_i, X_j) = e^{-\|X_i - X_j\|^2 / 2\sigma^2}$$

$$\text{Sigmoid kernel}: \quad K(X_i, X_j) = \tanh(\kappa X_i \cdot X_j - \delta)$$

- SVM can also be used for classifying multiple (> 2) classes and for regression analysis (with additional parameters)

# SVM vs. Neural Network

- **SVM**

  - 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**

  - 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 (nontrivial)

# SVM Related Links

- SVM Website: http://www.kernel-machines.org/

- Representative implementations

  - **LIBSVM**: an efficient implementation of SVM, multi-class classifications, nu-SVM, one-class SVM, including also various interfaces with java, python, etc.

  - **SVM-light**: simpler but performance is not better than LIBSVM, support only binary classification and only in C

  - **SVM-torch**: another recent implementation also written in C

# Chapter 9. Classification: Advanced Methods

- Bayesian Belief Networks

- Support Vector Machines

- Lazy Learners (or Learning from Your Neighbors)

- Other Classification Methods

- Additional Topics Regarding Classification

- Numerical Prediction

- Summary

# Lazy vs. Eager Learning

- <u>Lazy vs. eager learning</u>
  - **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 tuples, 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 an implicit global approximation to the target function
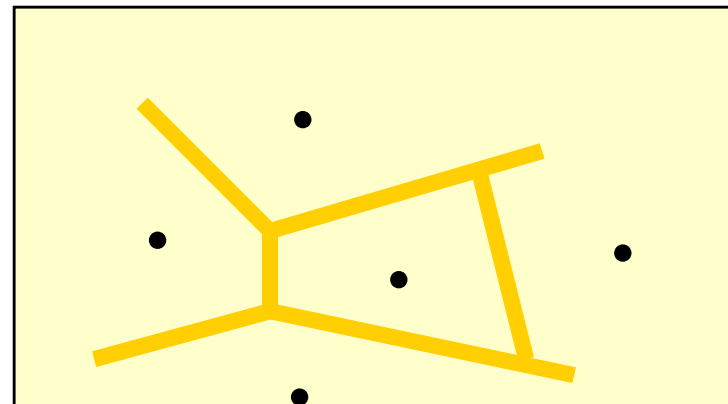  - Eager: must commit to a single hypothesis that covers the entire instance space
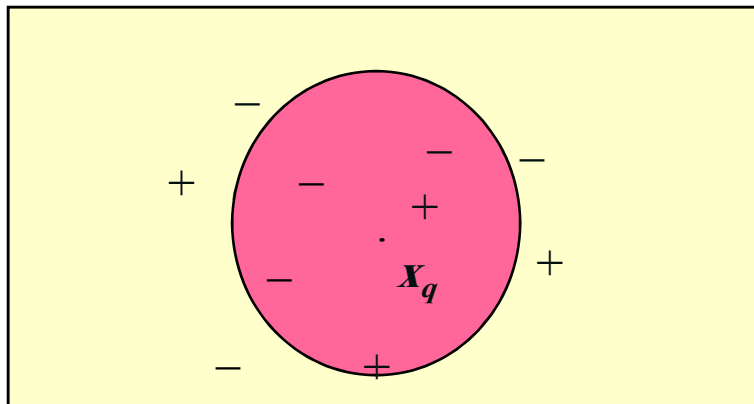
# Lazy Learner: Instance-Based Methods

- Instance-based learning:
  - Store training examples and delay the processing ("lazy evaluation") until a new instance must be classified
- Typical approaches
  - *k*-nearest neighbor approach
    - Instances represented as points in a Euclidean space.
  - Locally weighted regression
    - Constructs local approximation
  - Case-based reasoning
    - 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, dist($\mathbf{X_1}, \mathbf{X_2}$)
- Target function could be discrete- or real- valued
- For discrete-valued, *k*-NN returns the most common value among the *k* training examples nearest to $x_q$
- Vonoroi diagram: the decision surface induced by 1-NN for a typical set of training examples

# Discussion on the *k*-NN Algorithm

- *k*-NN for <u>real-valued prediction</u> for a given unknown tuple
  - Returns the mean values of the *k* nearest neighbors
- <u>Distance-weighted</u> nearest neighbor algorithm
  - Weight the contribution of each of the *k* neighbors according to their distance to the query $x_q$
    $$w \equiv \frac{1}{d(x_q, x_i)^2}$$
    - Give greater weight to closer neighbors
- <u>Robust</u> to noisy data by averaging *k*-nearest neighbors
- <u>Curse of dimensionality</u>: distance between neighbors could be dominated by irrelevant attributes
  - To overcome it, axes stretch or elimination of the least relevant attributes

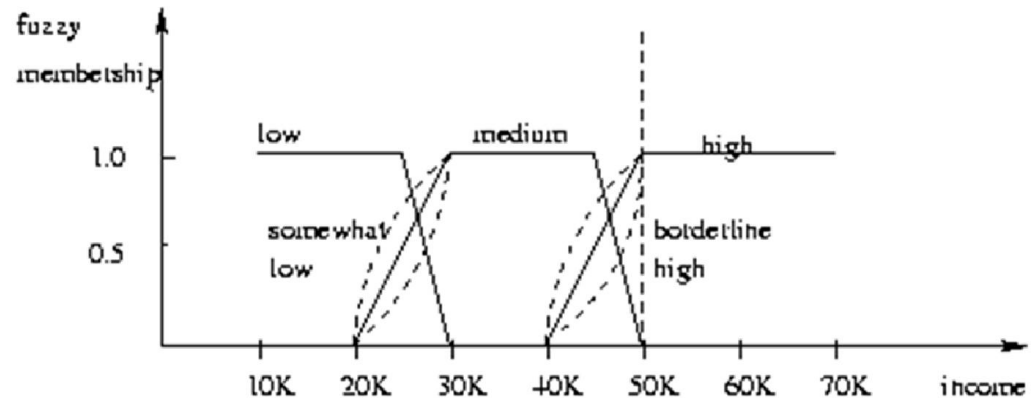# Chapter 9. Classification: Advanced Methods

- Bayesian Belief Networks

- Support Vector Machines

- Lazy Learners (or Learning from Your Neighbors)

- Other Classification Methods

- Additional Topics Regarding Classification

- Numerical Prediction

- Summary

# 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 offspring
- The *fitness of a rule* is represented by its classification accuracy on a set of training examples
- Offspring 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

# Fuzzy Set Approaches



- Fuzzy logic uses truth values between 0.0 and 1.0 to represent the degree of membership (such as in a *fuzzy membership graph*)
- Attribute values are converted to fuzzy values.  Ex.:
  - Income, $x$, is assigned a fuzzy membership value to each of the discrete categories {low, medium, high}, e.g. $49K belongs to "medium income" with fuzzy value 0.15 but belongs to "high income" with fuzzy value 0.96
  - Fuzzy membership values do not have to sum to 1.
- Each applicable rule contributes a vote for membership in the categories
- Typically, the truth values for each predicted category are summed, and these sums are combined

# Chapter 9. Classification: Advanced Methods

- Bayesian Belief Networks

- Support Vector Machines

- Lazy Learners (or Learning from Your Neighbors)

- Other Classification Methods

- Additional Topics Regarding Classification

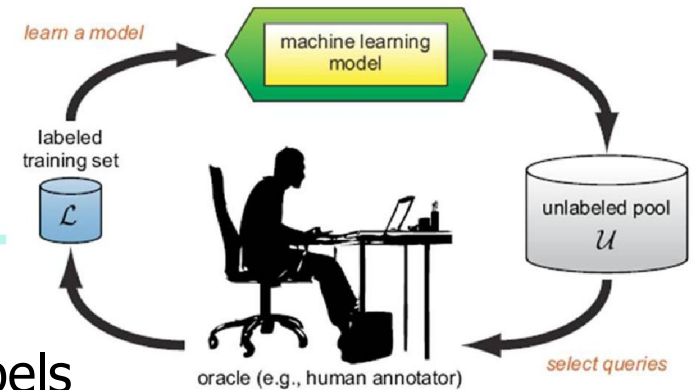- Numerical Prediction

- Summary

# Multiclass Classification

- Classification involving more than two classes (i.e., > 2 Classes)
- Method 1. **One-vs.-all** (OVA): Learn a classifier one at a time
  - Given m classes, train m classifiers: one for each class
  - Classifier j: treat tuples in class j as *positive* & all others as *negative*
  - To classify a tuple **X**, the set of classifiers vote as an ensemble
- Method 2. **All-vs.-all** (AVA): Learn a classifier for each pair of classes
  - Given m classes, construct m(m-1)/2 binary classifiers
  - A classifier is trained using tuples of the two classes
  - To classify a tuple **X**, each classifier votes.  X is assigned to the class with maximal vote
- Comparison
  - All-vs.-all tends to be superior to one-vs.-all
  - Problem: Binary classifier is sensitive to errors, and errors affect vote count

# Semi-Supervised Classification

- Semi-supervised: Uses labeled and unlabeled data to build a classifier
- Self-training:
    - Build a classifier using the labeled data
    - Use it to label the unlabeled data, and those with the most confident label prediction are added to the set of labeled data
    - Repeat the above process
    - Adv: easy to understand; disadv: may reinforce errors
- Co-training: Use two or more classifiers to teach each other
    - Each learner uses a mutually independent set of features of each tuple to train a good classifier, say $f_1$
    - Then $f_1$ and $f_2$ are used to predict the class label for unlabeled data X
    - Teach each other: The tuple having the most confident prediction from $f_1$ is added to the set of labeled data for $f_2$, & vice versa
- Other methods, e.g., joint probability distribution of features and labels
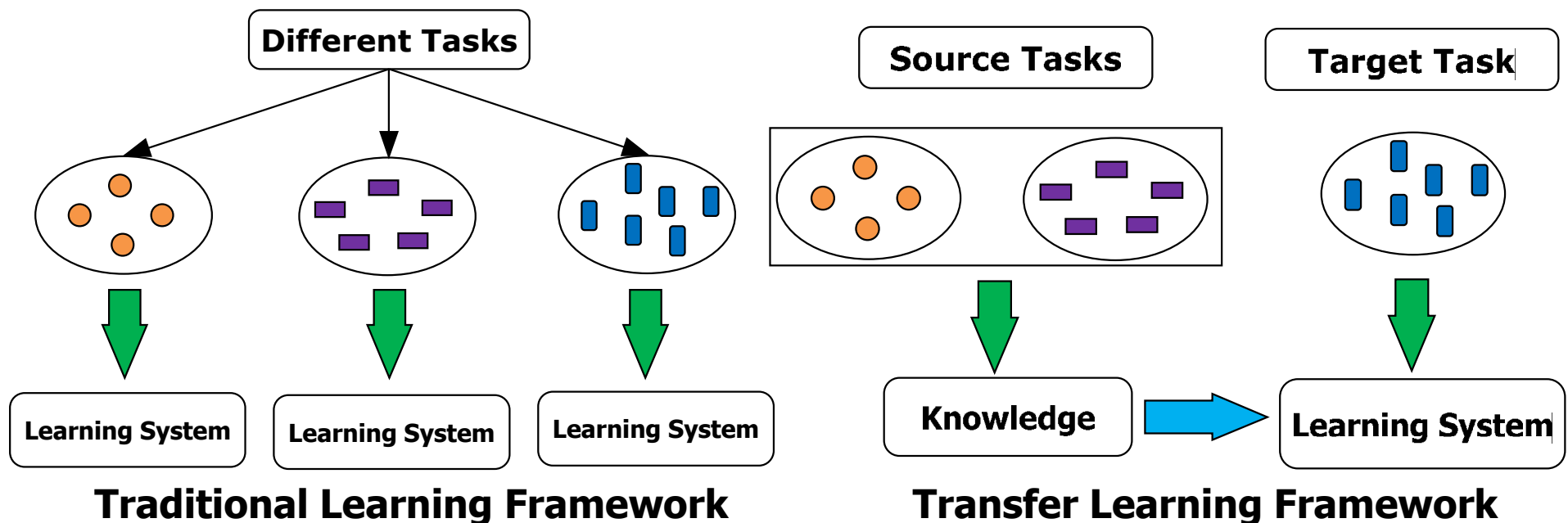
# Active Learning



- Class labels are expensive to obtain
- Active learner: query human (oracle) for labels
- Pool-based approach: Uses a pool of unlabeled data
  - L: a small subset of D is labeled, U: a pool of unlabeled data in D
  - Use a query function to carefully select one or more tuples from U and request labels from an oracle (a human annotator)
  - The newly labeled samples are added to L, and learn a model
  - Goal: Achieve high accuracy using as few labeled data as possible
- Evaluated using *learning curves*: Accuracy as a function of the number of instances queried (# of tuples to be queried should be small)
- Research issue: How to choose the data tuples to be queried?
  - Uncertainty sampling: choose the least certain ones
  - Reduce *version space*, the subset of hypotheses consistent w. the training data
  - Reduce expected entropy over U: Find the greatest reduction in the total number of incorrect predictions

# Transfer Learning: Conceptual Framework

- Transfer learning: Extract knowledge from one or more source tasks and apply the knowledge to a target task

- Traditional learning: Build a new classifier for each new task

- Transfer learning: Build new classifier by applying existing knowledge learned from source tasks



**Traditional Learning Framework**

**Transfer Learning Framework**

# Transfer Learning: Methods and Applications

- Applications: Especially useful when data is outdated or distribution changes, e.g., Web document classification, e-mail spam filtering
- *Instance-based transfer learning*:  Reweight some of the data from source tasks and use it to learn the target task
- TrAdaBoost (Transfer AdaBoost)
  - Assume source and target data each described by the same set of attributes (features) & class labels, but rather diff. distributions
  - Require only labeling a small amount of target data
  - Use source data in training: When a source tuple is misclassified, reduce the weight of such tupels so that they will have less effect on the subsequent classifier
- Research issues
  - Negative transfer: When it performs worse than no transfer at all
  - Heterogeneous transfer learning: Transfer knowledge from different feature space or multiple source domains
  - Large-scale transfer learning

# Chapter 9. Classification: Advanced Methods

- Bayesian Belief Networks

- Support Vector Machines

- Lazy Learners (or Learning from Your Neighbors)

- Other Classification Methods

- Additional Topics Regarding Classification

- Numerical Prediction 👉

- Summary

# Numerical Prediction

- (Numerical) prediction is similar to classification
  - construct a model
  - use model to predict continuous or ordered value for a given input
- Numerical prediction is different from classification
  - Classification refers to predict categorical class label
  - Numerical prediction models continuous-valued functions
- Major method for prediction: regression
  - model the relationship between one or more *independent* or **predictor** variables and a *dependent* or **response** variable
- Regression analysis
  - Linear and multiple regression
  - Non-linear regression
  - Other regression methods: generalized linear model, Poisson regression, log-linear models, regression trees

# 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

- <u>Method of least squares</u>: 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 $(\mathbf{X_1}, y_1), (\mathbf{X_2}, y_2), \ldots, (\mathbf{X_{|D|}}, y_{|D|})$
  - 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

# Nonlinear Regression

- Some nonlinear models can be modeled by a polynomial function

- A polynomial regression model can be transformed into linear regression model.  For example,

$$y = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

  convertible to linear with new variables: $x_2 = x^2$, $x_3 = x^3$

$$y = w_0 + w_1 x + w_2 x_2 + w_3 x_3$$

- Other functions, such as power function, can also be transformed to linear model

- Some models are intractable nonlinear (e.g., sum of exponential terms)

  - possible to obtain least square estimates through extensive calculation on more complex formulae

# 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
  - Poisson regression: 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

# Regression Trees and Model Trees

- Regression tree: proposed in CART system (Breiman et al. 1984)

  - CART: Classification And Regression Trees

  - Each leaf stores a *continuous-valued prediction*

  - It is the *average value of the predicted attribute* for the training tuples that reach the leaf

- Model tree: proposed by Quinlan (1992)

  - Each leaf holds a regression model—a multivariate linear equation for the predicted attribute

  - A more general case than regression tree

- Regression and model trees tend to be more accurate than linear regression when the data are not represented well by a simple linear model
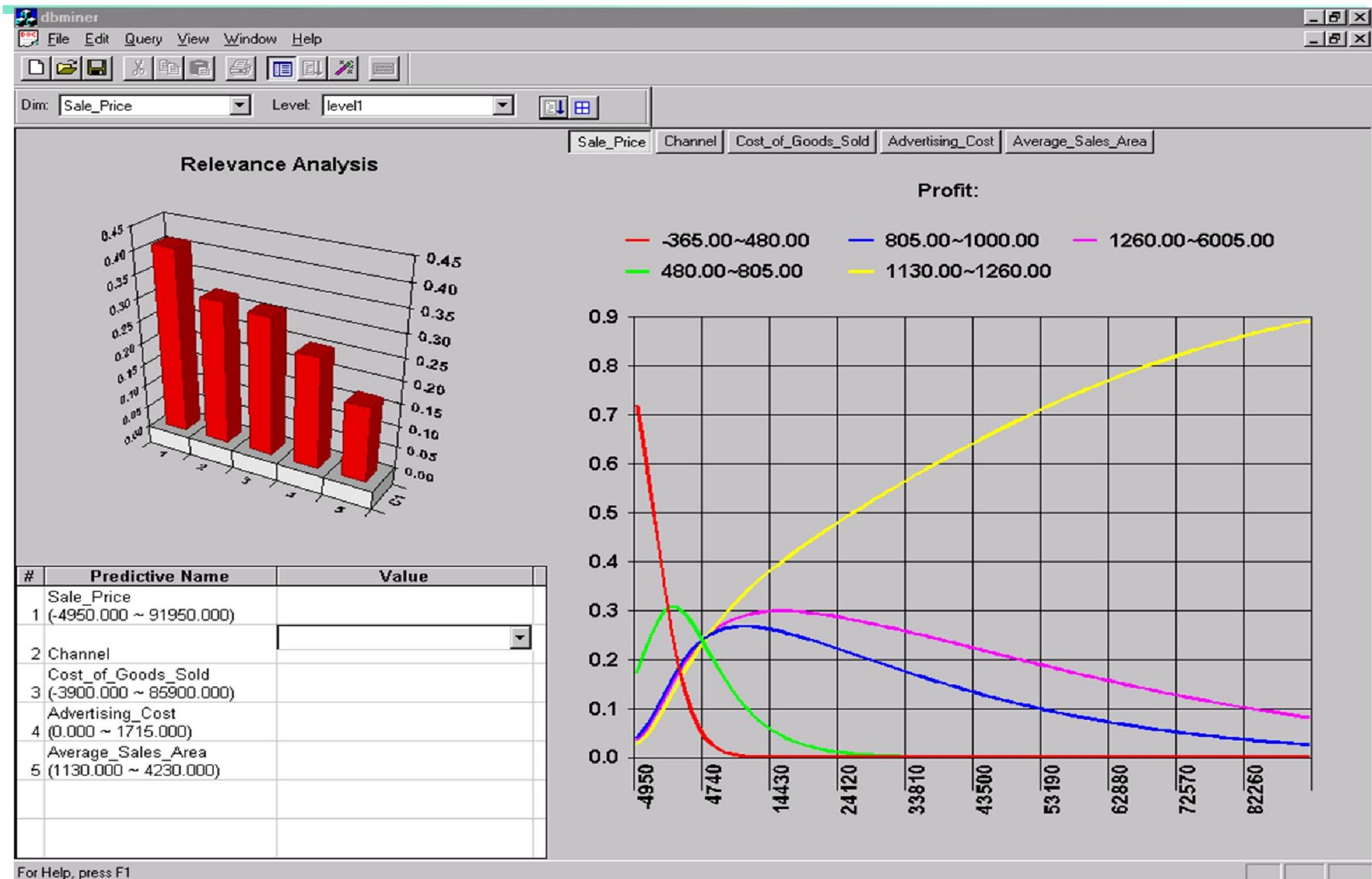
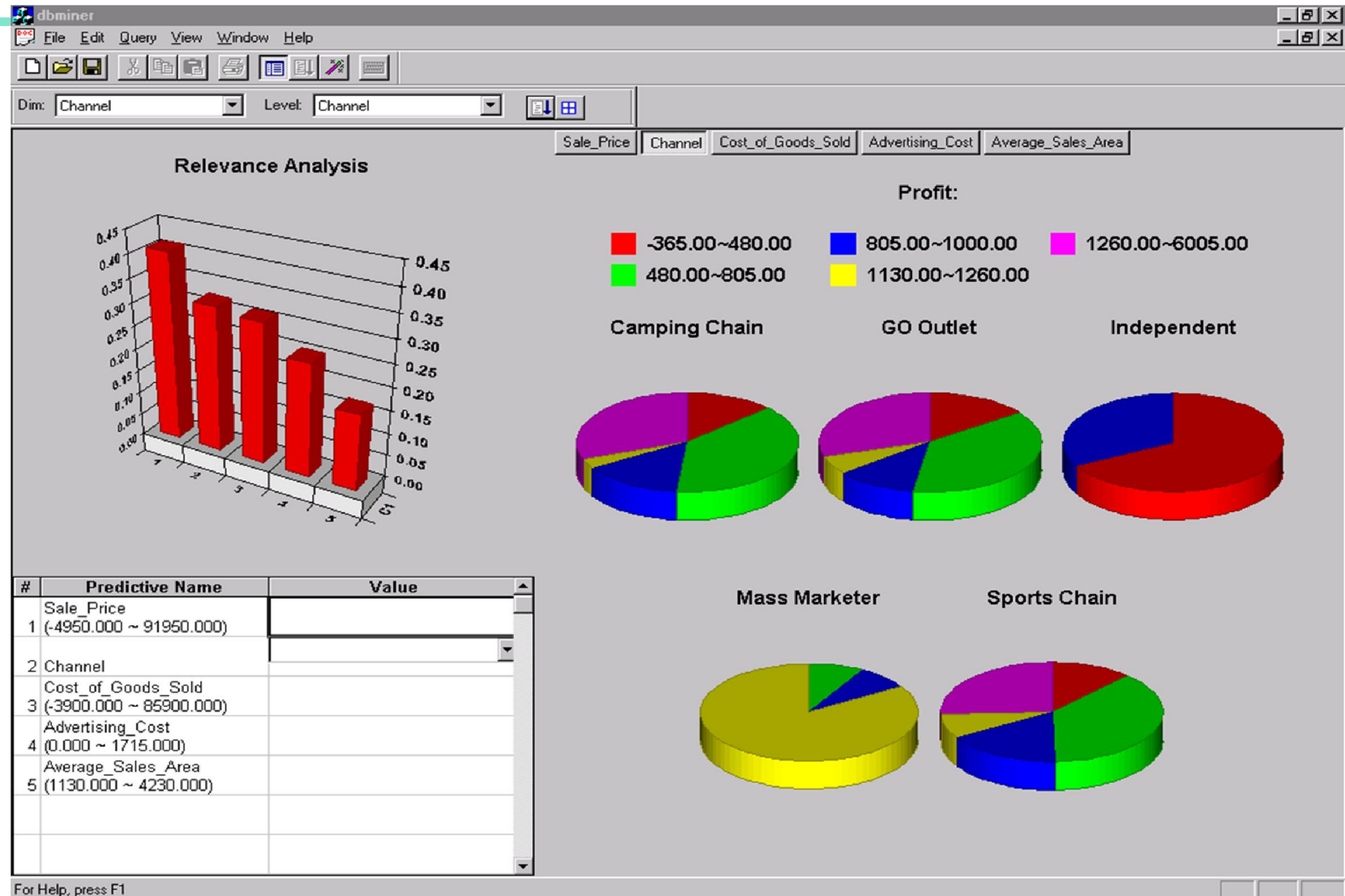# Predictive Modeling in Multidimensional Databases

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

# Prediction: Numerical Data

# Prediction: Categorical Data

# Chapter 9. Classification: Advanced Methods

- Bayesian Belief Networks

- Support Vector Machines

- Lazy Learners (or Learning from Your Neighbors)

- Other Classification Methods

- Additional Topics Regarding Classification

- Numerical Prediction

- Summary

# Summary

- Effective and advanced classification methods

    - Bayesian belief network (probabilistic networks)

    - Support Vector Machine (SVM)

    - Other classification methods: lazy learners (KNN, case-based reasoning), genetic algorithms, rough set and fuzzy set approaches

- Additional Topics on Classification

    - Multiclass classification

    - Semi-supervised classification

    - Active learning

    - Transfer learning

- Numerical prediction

    - (Linear) Regression