# Prediction

By

Dr. T. Sree Sharmila,
Associate Professor, Dept. of IT,
SSN College of Engineering

SSN

# What Is Prediction?

- (Numerical) prediction is similar to classification
  - construct a model
  - use model to predict continuous or ordered value for a given input
- Prediction is different from classification
  - Classification refers to predict categorical class label
  - 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 $(\mathbf{X_1}, y_1), (\mathbf{X_2}, y_2),..., (\mathbf{X_{|D|}}, y_{|D|})$
  - 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

- <u>Generalized linear model</u>:
  - 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
  - <u>Logistic regression</u>: 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
- <u>Log-linear models</u>: (for categorical data)
  - Approximate discrete multidimensional prob. distributions
  - Also useful for data compression and smoothing
- <u>Regression trees and model trees</u>
  - 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