# Data warehousing and Mining

-Kiran Bhowmick

# Course structure

| Program: Third Year B.Tech. in Computer Engineering | | | | | | | | Semester : V | |
|---|---|---|---|---|---|---|---|---|---|
| Course : Data Mining and Warehouse | | | | | | | | Course Code:DJ19CEC501 | |
| Course :  Data Mining and Warehouse Laboratory | | | | | | | | Course Code: DJ19CEL501 | |

| Teaching Scheme (Hours / week) | | | | Evaluation Scheme | | | | | Total marks (A+ B) |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Semester End Examination Marks (A) | | Continuous Assessment  Marks (B) | | | |
| Lectures (hrs) | Practical (hrs) | Tutorial (hrs) | Total Credits | Theory | | Term Test 1 | Term Test 2 | Avg. | |
| | | | | 65 | | 20 | 15 | 35 | 100 |
| | | | | Laboratory Examination | | Term work | | | |
| 3 | 2 | 5 | 4 | Oral | Practical | Oral &Practical | Laboratory Work | Tutorial /  Mini project / presentation/ Journal | Total Term work | 50 |
| | | | | 25 | - | - | 15 | 10 | 25 | |

# Syllabus

| Program: Computer Engineering | | T.Y B.Tech. | Semester: V |
|---|---|---|---|
| Course: Data Warehousing and Mining (DJS22CEC501) | | | |
| Course: Data Warehousing and Mining Laboratory (DJS22CEL501) | | | |

**Pre-requisite:** Basic database concepts, Concepts of algorithm design and analysis

**Course Objectives:**
This course introduces data warehouse and data mining concepts.
1. To identify the need of and perform data modelling to provide strategic information for making business decisions.
2. To analyze data and identify and develop relevant mining models to discover knowledge from data in various applications.

**Outcomes:** On successful completion of course, learner will be able to:
1. Design data warehouse models using dimension-modeling techniques.
2. Analyse the data by applying Online Analytical Processing (OLAP) operations for strategic decisions.
3. Apply preprocessing techniques to the given raw data.
4. Apply appropriate data mining techniques on data sets to retrieve relevant information.

## Data Warehousing and Mining (DJS22CEC501)

| Unit | Description | Duration |
|------|-------------|----------|
| 1 | **Introduction to Data Warehouse and Dimensional modelling:** Introduction to Strategic Information, Need for Strategic Information, Features of Data Warehouse, Data warehouse versus Data Marts, Data warehouse versus Data Lake, Top-down versus Bottom-up approach. Data warehouse architecture, E-R modelling versus Dimensional Modelling, Information Package Diagram, STAR schema, STAR schema keys, Snowflake Schema, Fact Constellation Schema, Factless Fact tables, Update to the dimension tables, Aggregate fact tables. | 8 |
| 2 | **ETL Process and OLAP:** Major steps in ETL process, Data extraction: Techniques, Data transformation: Basic tasks, Major transformation types, Data Loading: Applying Data, OLTP Vs OLAP, OLAP definition, Dimensional Analysis, Hypercubes, OLAP operations: Drill down, Roll up, Slice, Dice and Rotation, OLAP models: MOLAP, ROLAP, HOLAP. | 6 |
| 3 | **Introduction to Data Mining, Data Exploration and Preprocessing:**<br><br>Data Mining Task and Techniques, KDD process, Issues in Data Mining, Applications of Data Mining, Data Exploration: Types of Attributes, Statistical | 6 |

| | | |
|---|---|---|
| | Description of Data, Data Visualization, Measuring data similarity and dissimilarity.<br><br>Data Preprocessing: Major tasks in preprocessing, Data Cleaning: Missing values, Noisy data; Data Integration: Entity Identification Problem, Redundancy and Correlation Analysis, Tuple Duplication, Data Value Conflict Detection and Resolution; Data Reduction: Attribute subset selection, Histograms, Clustering and Sampling; Data Transformation & Data Discretization: Data Transformation by Normalization, Discretization by Binning, Discretization by Histogram Analysis | |
| 4 | **Classification and Clustering:**<br><br>**Classification**<br><br>Basic Concepts of classification, Decision Tree Induction, Attribute Selection Measures using Information Gain, Tree pruning<br><br>Bayes Classification Methods: Bayes' Theorem, Naïve Bayesian Classification<br><br>Model Evaluation: Metrics for Evaluating Classifier Performance, Holdout Method and Random Subsampling, Cross Validation, Bootstrap<br><br>Improving Classification Accuracy: Ensemble classification, Bagging, Boosting and AdaBoost, Random Forests<br><br>**Clustering:**<br><br>Cluster Analysis and Requirements of Cluster Analysis<br><br>Partitioning Methods: k-Means, k-Medoids<br><br>Hierarchical Methods: Agglomerative, Divisive<br><br>Evaluation of Clustering: Assessing Clustering Tendency, Determining Number of Clusters and Measuring cluster quality: Intrinsic and Extrinsic methods | 8 |
| 5 | **Mining Frequent Patterns and Association Rules:**<br><br>Market Basket Analysis, Frequent Item sets, Closed Item sets, and Association Rule<br><br>Frequent Item set Mining Methods: Apriori Algorithm, Association Rule Generation,<br><br>FP growth | 5 |
| 6 | **Spatial and Web Mining:** Spatial Data, Spatial Vs. Classical Data Mining, Spatial | 6 |

| | | |
|---|---|---|
| Data Structures, Mining Spatial Association and Co-location Patterns, Spatial Clustering Techniques: CLARANS Extension<br><br>**Web Mining:** Web Content Mining, Web Structure Mining, Web Usage mining, Applications of Web Mining | |

**Books Recommended:**

1. Paulraj Ponniah, "Data Warehousing: Fundamentals for IT Professionals", 2nd Edition, Wiley India, 2013.
2. Theraja Reema, "Data Warehousing", 1st Edition, Oxford University Press, 2009.
3. Han, Kamber, "Data Mining Concepts and Techniques", 3rd Edition, Morgan Kaufmann, 2012.
4. P. N. Tan, M. Steinbach, Vipin Kumar, "Introduction to Data Mining", 2nd Edition, Pearson Education, 2018.
5. H. Dunham, "Data Mining: Introductory and Advanced Topics", 1st Edition, Pearson Education, 2006.

# Chapter 1. Introduction

- Motivation: Why data mining?

- What is data mining?

- Data Mining: On what kind of data?

- Data mining functionality

- Classification of data mining systems

- Top-10 most popular data mining algorithms

- Major issues in data mining

- Overview of the course

# Why Data Mining?

- The Explosive Growth of Data: from terabytes to petabytes
  - Data collection and data availability
    - Automated data collection tools, database systems, Web, computerized society
  - Major sources of abundant data
    - Business: Web, e-commerce, transactions, stocks, …
    - Science: Remote sensing, bioinformatics, scientific simulation, …
    - Society and everyone: news, digital cameras, YouTube
- <u>We are drowning in data, but starving for knowledge!</u>
- "Necessity is the mother of invention"—Data mining—Automated analysis of massive data sets

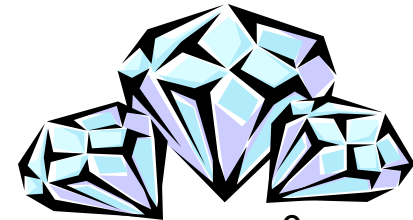# Why Not Traditional Data Analysis?

- Tremendous amount of data
  - Algorithms must be highly scalable to handle such as tera-bytes of data
- High-dimensionality of data
  - Micro-array may have tens of thousands of dimensions
- High complexity of data
  - Data streams and sensor data
  - Time-series data, temporal data, sequence data
  - Structure data, graphs, social networks and multi-linked data
  - Heterogeneous databases and legacy databases
  - Spatial, spatiotemporal, multimedia, text and Web data
  - Software programs, scientific simulations
- New and sophisticated applications

Dr. Kiran Bhowmick

# What Is Data Mining?

- Data mining (knowledge discovery from data)
  - Extraction of interesting (<u>non-trivial,</u> <u>implicit</u>, <u>previously unknown</u> and <u>potentially useful)</u> patterns or knowledge from huge amount of data
  - Data mining: a misnomer?
- Alternative names
  - Knowledge discovery (mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.
- Watch out: Is everything "data mining"?
  - Simple search and query processing
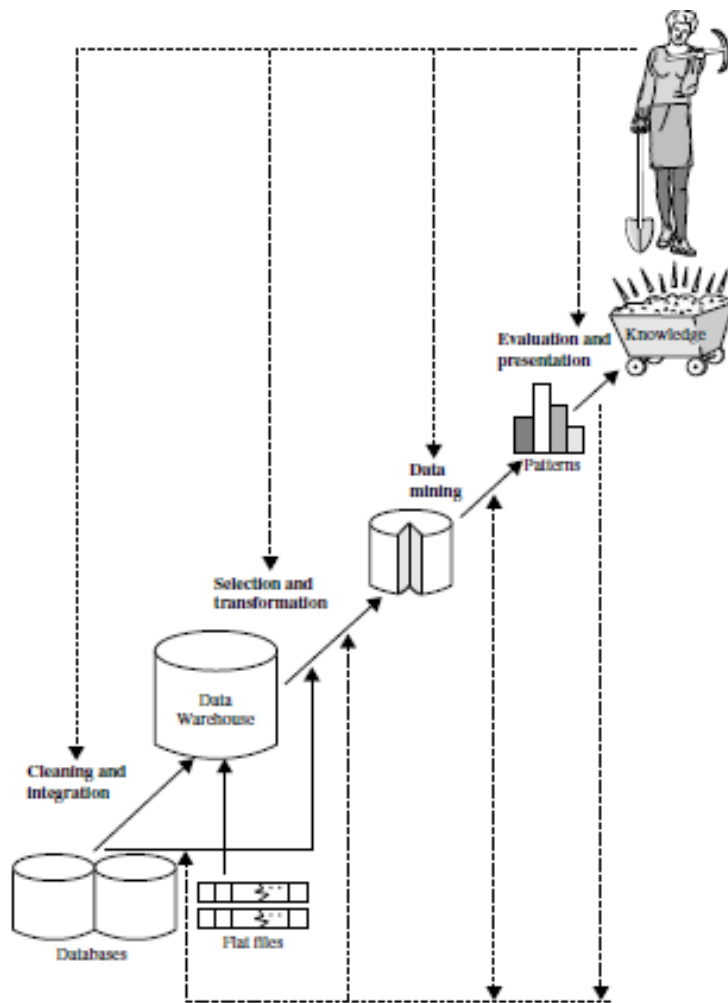  - (Deductive) expert systems

# Query Examples

- ## Database
  - Find all credit applicants with last name of Smith.
  - Identify customers who have purchased more than $10,000 in the last month.
  - Find all customers who have purchased IBM laptops
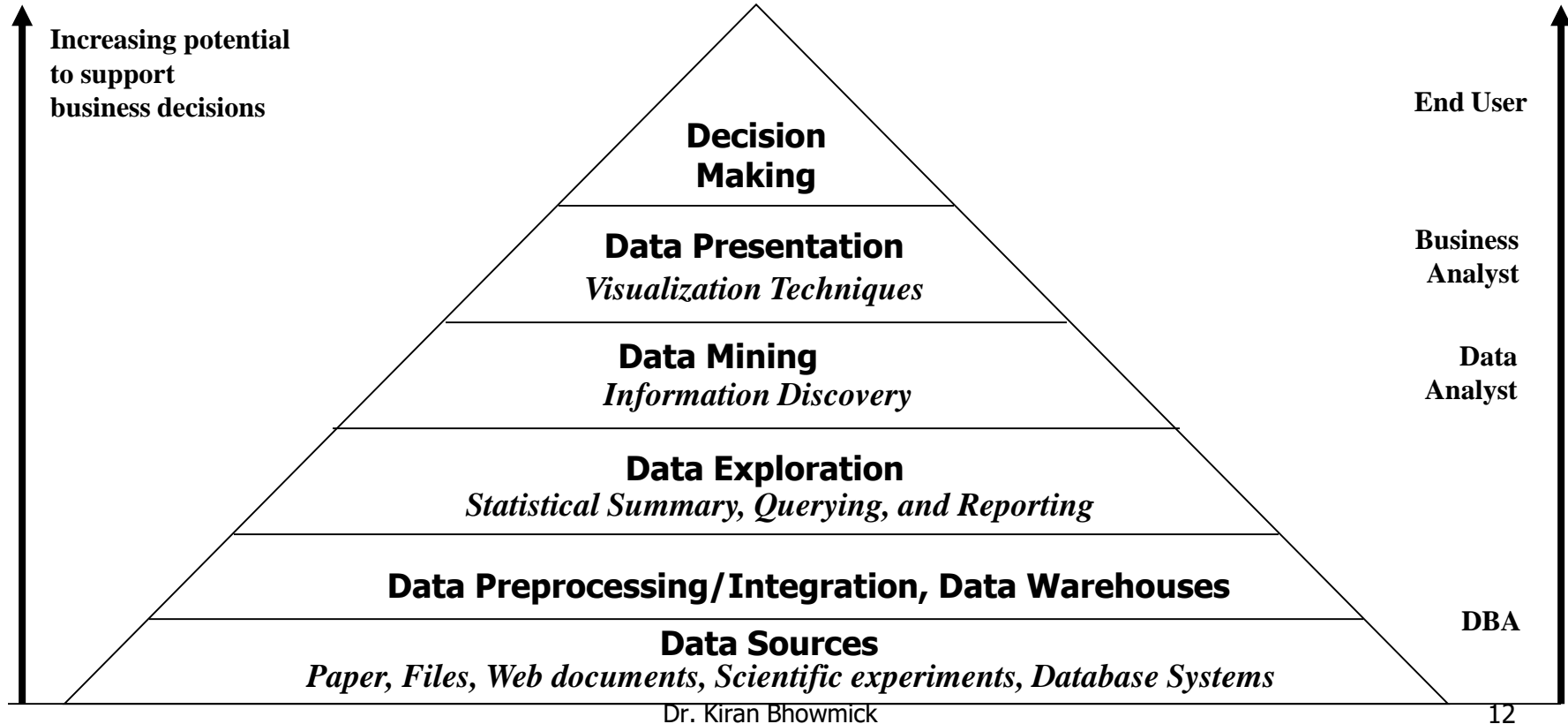
- ## Data Mining
  - Identify customers with similar buying habits. (Clustering)
  - Find all items which are frequently purchased with milk. (association rules)
  - Find all credit applicants who are poor credit risks. (classification)
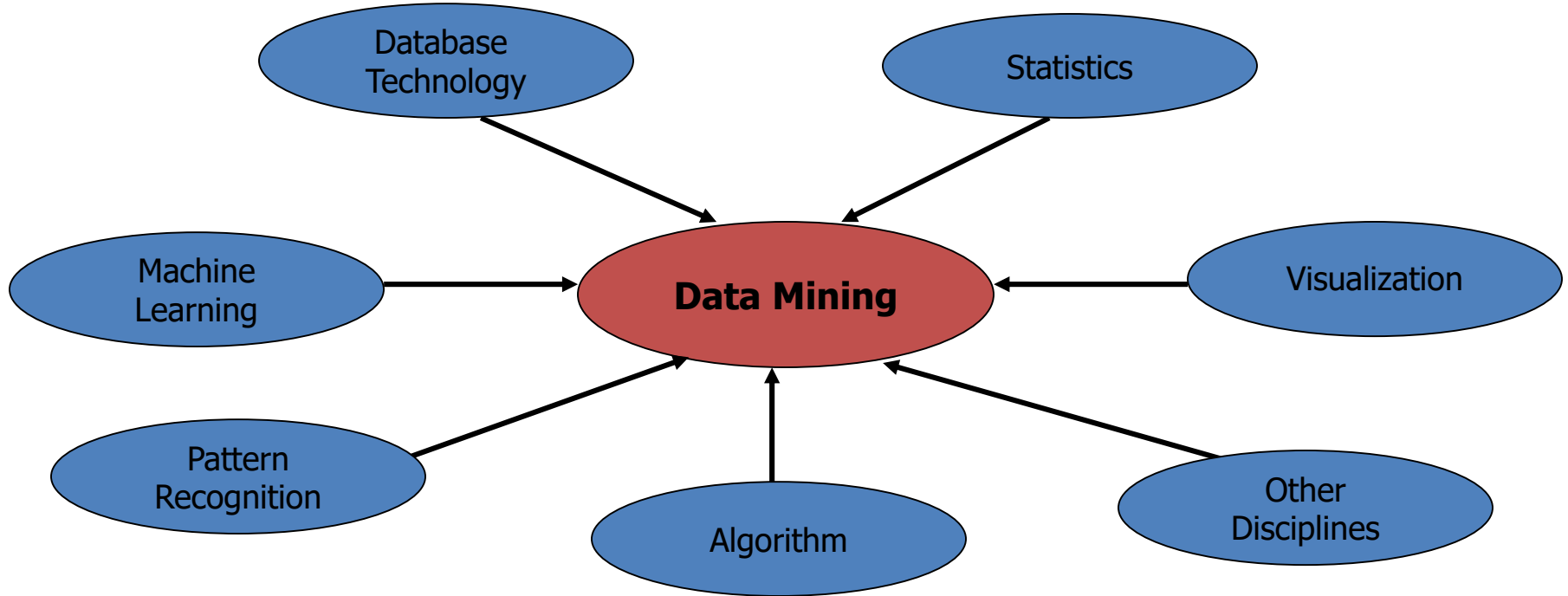
## Knowledge Discovery from Data (KDD)



**Figure 1.4** Data mining as a step in the process of knowledge discovery.

# Data Mining and Business Intelligence

Increasing potential
to support
business decisions

**End User**

**Decision Making**

**Business Analyst**

**Data Presentation**
*Visualization Techniques*

**Data Analyst**

**Data Mining**
*Information Discovery*

**Data Exploration**
*Statistical Summary, Querying, and Reporting*

**Data Preprocessing/Integration, Data Warehouses**

**DBA**

**Data Sources**
*Paper, Files, Web documents, Scientific experiments, Database Systems*

# Data Mining: Confluence of Multiple Disciplines

# Multi-Dimensional View of Data Mining

- **Data to be mined**
  - Relational, data warehouse, transactional, stream, object-oriented/relational, active, spatial, time-series, text, multi-media, heterogeneous, legacy, WWW

- **Knowledge to be mined**
  - Characterization, discrimination, association, classification, clustering, trend/deviation, outlier analysis, etc.
  - Multiple/integrated functions and mining at multiple levels

- **Techniques utilized**
  - Database-oriented, data warehouse (OLAP), machine learning, statistics, visualization, etc.

- **Applications adapted**
  - Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, text mining, Web mining, etc.
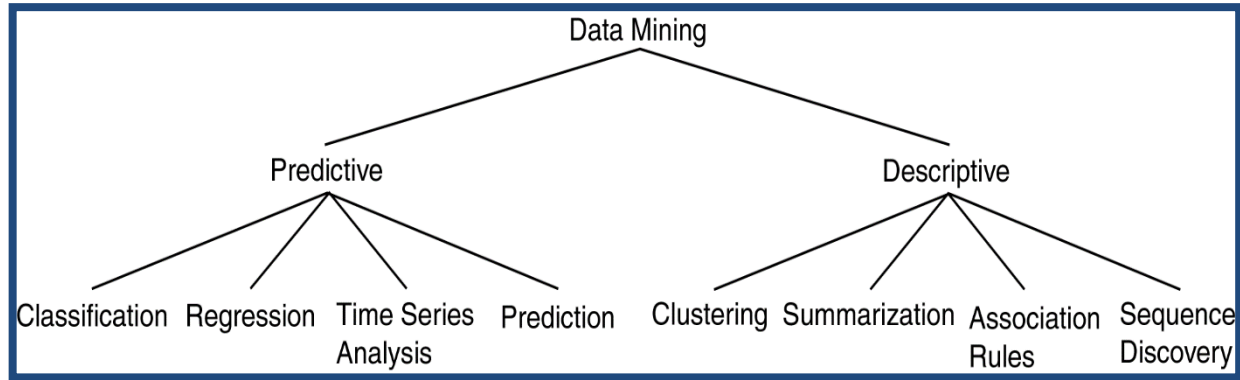
# Data Mining Definition

- Finding hidden information in a database

- Fit data to a model

- Similar terms

  - Exploratory data analysis

  - Data driven discovery

  - Deductive learning

# Data Mining Algorithm

- Objective:  Fit Data to a Model
  - Descriptive
  - Predictive
- Preference – Technique to choose the best model
- Search – Technique to search the data
  - "Query"

# Data Mining Models and Tasks



Predictive Data Model:

Makes prediction about values of data using known results found from different data.

Based on historical data. E.g: classification, regression, time series analysis, prediction

Descriptive data model:

Identifies patterns or relationships in data.

Offers a detailed description of the data, for example- it gives insight into what's going on inside the data without any prior idea.  This serves as a way to explore the properties of data.

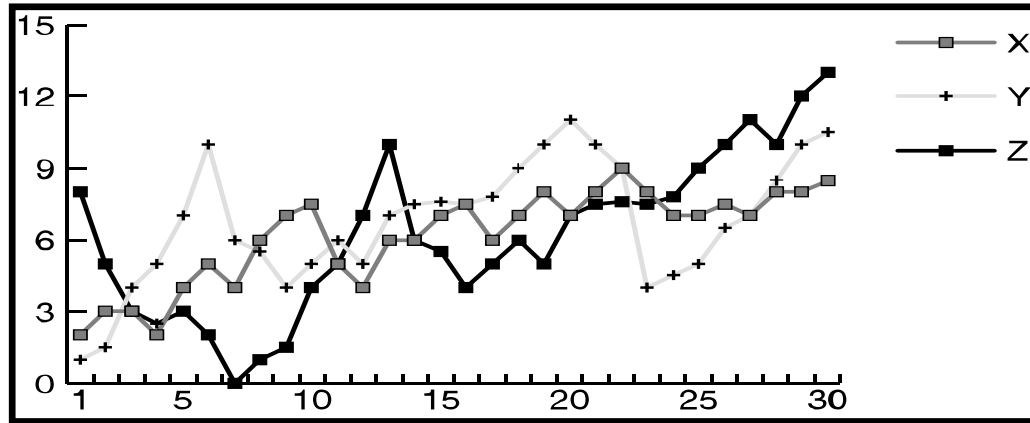Does not predict new properties.

E.g: Clustering, association rules etc

# Basic Data Mining Tasks

- ***Classification*** maps data into predefined groups or classes
  - Supervised learning
  - Pattern recognition
  - Prediction

- ***Regression*** is used to map a data item to a real valued prediction variable.
  - Regression involves predicting continuous, real-value quantities
  - regression involves the learning of the function that does this mapping.
  - E.g. Predicting house prices

- ***Time Series Analysis*** the value of an attribute is examined as it varies over time. Values obtained as evenly spaced time points (daily, weekly, hourly) Three basic functions:
  - Distance measures: determine similarity
  - Structure of line
  - Historical time series plot to predict future values
  - E.g. forecast monthly sales for a retail store

- ***Prediction*** predicting future states
  - Weather forecasts, earthquakes, floods etc
  - E.g. medical diagnosis, fraud detection etc.

# Ex: Time Series Analysis

- Example: Stock Market
- Predict future values
- Determine similar patterns over time



Time series plots

# Basic Data Mining Tasks

- ***Clustering*** groups similar data together into clusters.
  - Unsupervised learning
  - Segmentation
  - Partitioning
- ***Summarization*** maps data into subsets with associated simple descriptions.
  - Characterization or Generalization
  - It extracts or derives representative information about the database
  - E.g. Average CET score taking admission in an Engg college. Summarization will help estimate the type and intellect of student in the college
- ***Association rules*** uncovers relationships among data.
  - Also called link analysis, Affinity Analysis
  - An association rule is a model that identifies specific types of data associations.
  - E.g. Market-basket analysis

# Basic Data Mining Tasks (cont'd)

- *Sequence Discovery* to determine sequential patterns in data.
  - Patterns are based on a time sequence of actions.
  - Similar to association – data are found that are related.
  - Difference than association – relationship is based on time. For AR – items must be purchased at same time whereas for sequence discovery items can be purchased over a period of time.
  - E.g. Analyzing web logs of a website to understand what sequence of pages are frequently visited by users.
  - (A, B, C) or (A, D, B, C) or (A, E, B, C). Then add a link directly from page A to page C.