



# **Data Mining: Concepts and Techniques**

# Introduction

---

- Why Data Mining? 
- What Is Data Mining?
- A Multi-Dimensional View of Data Mining
- What Kind of Data Can Be Mined?
- What Kinds of Patterns Can Be Mined?
- What Technology Are Used?
- What Kind of Applications Are Targeted?
- Major Issues in Data Mining
- A Brief History of Data Mining and Data Mining Society
- Summary

# 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
- We are drowning in data, but starving for knowledge!
- “Necessity is the mother of invention”—Data mining—Automated analysis of massive data sets

# Introduction

---

- Why Data Mining?
- What Is Data Mining? 
- A Multi-Dimensional View of Data Mining
- What Kind of Data Can Be Mined?
- What Kinds of Patterns Can Be Mined?
- What Technology Are Used?
- What Kind of Applications Are Targeted?
- Major Issues in Data Mining
- A Brief History of Data Mining and Data Mining Society
- Summary

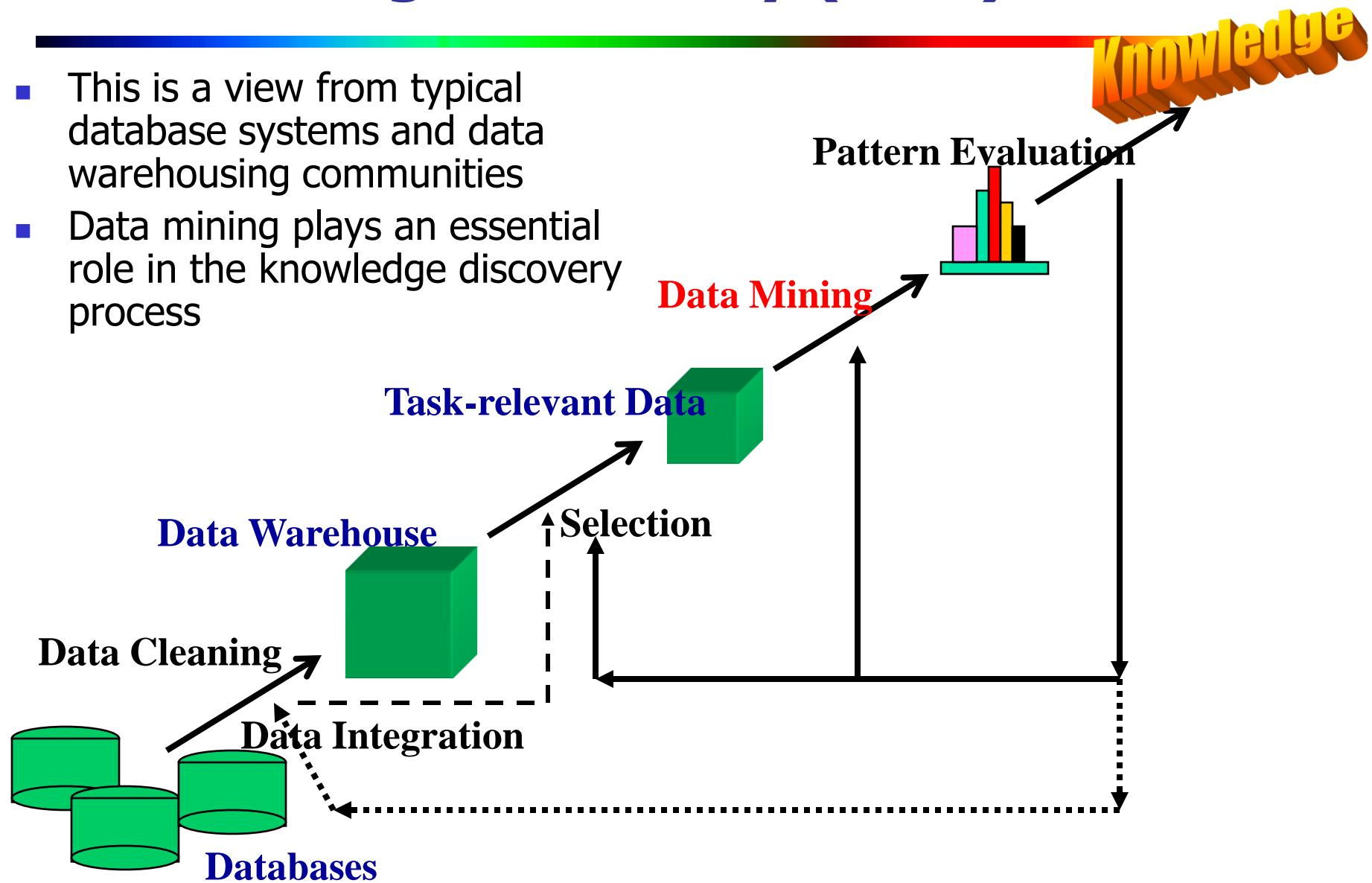
# What Is Data Mining?

---

- Data mining (knowledge discovery from data)
  - Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data

# Knowledge Discovery (KDD) Process

- This is a view from typical database systems and data warehousing communities
- Data mining plays an essential role in the knowledge discovery process



# Introduction

---

- Why Data Mining?
- What Is Data Mining?
- A Multi-Dimensional View of Data Mining
- What Kind of Data Can Be Mined?
- What Kinds of Patterns Can Be Mined?
- What Technology Are Used?
- What Kind of Applications Are Targeted?
- Major Issues in Data Mining
- A Brief History of Data Mining and Data Mining Society
- Summary



# **Multi-Dimensional View of Data Mining**

---

- **Data to be mined**
  - Database data (extended-relational, object-oriented, heterogeneous, legacy), data warehouse, transactional data, stream, spatiotemporal, time-series, sequence, text and web, multi-media, graphs & social and information networks
- **Knowledge to be mined (or: Data mining functions)**
  - association, classification, clustering, trend/deviation, outlier analysis, etc.
  - Multiple/integrated functions and mining at multiple levels
- **Techniques utilized**
  - Data-intensive, data warehouse (OLAP), machine learning, statistics, pattern recognition, visualization, high-performance, etc.
- **Applications adapted**
  - Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, text mining, Web mining, etc.

# Introduction

---

- Why Data Mining?
- What Is Data Mining?
- A Multi-Dimensional View of Data Mining
- What Kind of Data Can Be Mined? 
- What Kinds of Patterns Can Be Mined?
- What Technology Are Used?
- What Kind of Applications Are Targeted?
- Major Issues in Data Mining
- A Brief History of Data Mining and Data Mining Society
- Summary

# Data Mining: On What Kinds of Data?

---

- Database-oriented data sets and applications
  - Relational database, data warehouse, transactional database
- Advanced data sets and advanced applications
  - Data streams and sensor data
  - Time-series data, sequence data (incl. bio-sequences)
  - graphs, social networks and multi-linked data
  - Heterogeneous databases
  - Spatial data
  - Multimedia database
  - Text databases
  - The World-Wide Web

# Introduction

---

- Why Data Mining?
- What Is Data Mining?
- A Multi-Dimensional View of Data Mining
- What Kind of Data Can Be Mined?
- What Kinds of Patterns Can Be Mined? 
- What Technology Are Used?
- What Kind of Applications Are Targeted?
- Major Issues in Data Mining
- A Brief History of Data Mining and Data Mining Society
- Summary

# Data Mining Function: (1) Generalization

---

- Information integration and data warehouse construction
  - Data cleaning, transformation, integration, and multidimensional data model
- Data cube technology
  - Scalable methods for computing (i.e., materializing) multidimensional aggregates
  - OLAP (online analytical processing)
- Multidimensional concept description: Characterization and discrimination
  - Generalize, summarize, and contrast data characteristics, e.g., dry vs. wet region

# Data Mining Function: (2) Association and Correlation Analysis

---

- Frequent patterns (or frequent itemsets)
  - What items are frequently purchased together in your Walmart?
- Association, correlation vs. causality
  - A typical association rule
    - Diaper → Beer [0.5%, 75%] (support, confidence)
  - Are strongly associated items also strongly correlated?
- How to mine such patterns and rules efficiently in large datasets?
- How to use such patterns for classification, clustering, and other applications?

# Data Mining Function: (3) Classification

---

- Classification and label prediction
  - Construct models (functions) based on some training examples
  - Describe and distinguish classes or concepts for future prediction
    - E.g., classify countries based on (climate), or classify cars based on (gas mileage)
  - Predict some unknown class labels
- Typical methods
  - Decision trees, naïve Bayesian classification, support vector machines, neural networks, rule-based classification, pattern-based classification, logistic regression, ...
- Typical applications:
  - Credit card fraud detection, direct marketing, classifying diseases, web-pages, ...

# Data Mining Function: (4) Cluster Analysis

---

- Unsupervised learning (i.e., Class label is unknown)
- Group data to form new categories (i.e., clusters), e.g., cluster houses to find distribution patterns
- Principle: Maximizing intra-class similarity & minimizing interclass similarity
- Many methods and applications

# Data Mining Function: (5) Outlier Analysis

---

- Outlier analysis
  - Outlier: A data object that does not comply with the general behavior of the data
  - Noise or exception?
  - Methods: by product of clustering or regression analysis, ...
  - Useful in fraud detection, rare events analysis

# Time and Ordering: Sequential Pattern, Trend and Evolution Analysis

---

- Sequence, trend and evolution analysis
  - Trend, time-series, and deviation analysis: e.g., regression and value prediction
  - Sequential pattern mining
    - e.g., first buy digital camera, then buy large SD memory cards
  - Motifs and biological sequence analysis
    - Approximate and consecutive motifs
  - Similarity-based analysis
- Mining data streams
  - Ordered, time-varying, potentially infinite, data streams

# Structure and Network Analysis

---

- Graph mining
  - Finding frequent subgraphs (e.g., chemical compounds), trees (XML), substructures (web fragments)
- Information network analysis
  - Social networks: actors (objects, nodes) and relationships (edges)
    - e.g., author networks in CS, terrorist networks
  - Multiple heterogeneous networks
    - A person could be multiple information networks: friends, family, classmates, ...
  - Links carry a lot of semantic information: Link mining
- Web mining
  - Web is a big information network: from PageRank to Google
  - Analysis of Web information networks
    - Web community discovery, opinion mining, usage mining, ...

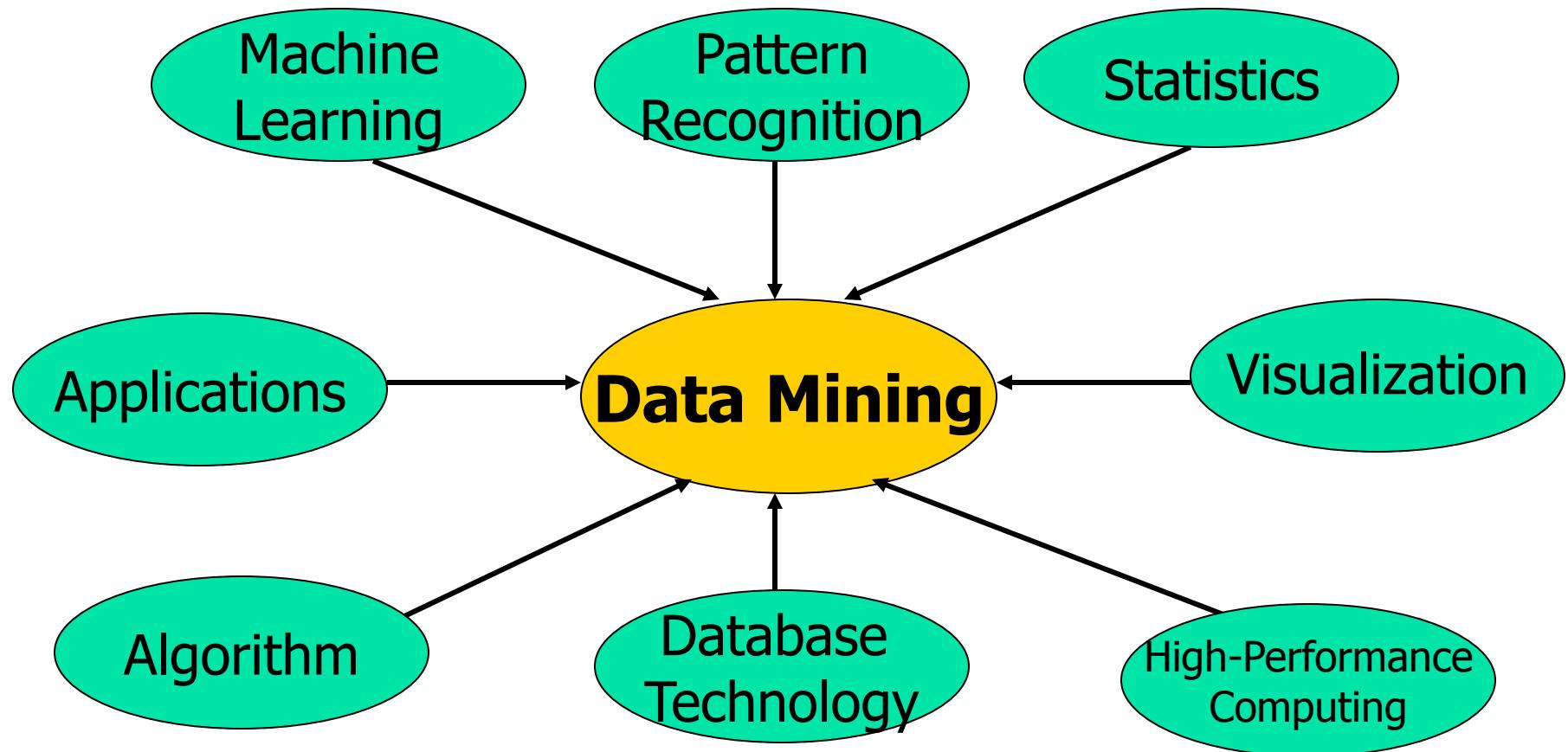
# Introduction

---

- Why Data Mining?
- What Is Data Mining?
- A Multi-Dimensional View of Data Mining
- What Kind of Data Can Be Mined?
- What Kinds of Patterns Can Be Mined?
- What Technology Are Used? 
- What Kind of Applications Are Targeted?
- Major Issues in Data Mining
- A Brief History of Data Mining and Data Mining Society
- Summary

# Data Mining: Confluence of Multiple Disciplines

---



# Why Confluence of Multiple Disciplines?

---

- 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
  - Spatial, multimedia, text and Web data
- New and sophisticated applications

# Introduction

---

- Why Data Mining?
- What Is Data Mining?
- A Multi-Dimensional View of Data Mining
- What Kind of Data Can Be Mined?
- What Kinds of Patterns Can Be Mined?
- What Technology Are Used?
- What Kind of Applications Are Targeted? 
- Major Issues in Data Mining
- A Brief History of Data Mining and Data Mining Society
- Summary

# Applications of Data Mining

---

- Web page analysis: from web page classification, clustering to PageRank algorithms
- Collaborative analysis & recommender systems
- Basket data analysis to targeted marketing
- Biological and medical data analysis: classification, cluster analysis (microarray data analysis), biological sequence analysis, biological network analysis

# Introduction

---

- Why Data Mining?
- What Is Data Mining?
- A Multi-Dimensional View of Data Mining
- What Kind of Data Can Be Mined?
- What Kinds of Patterns Can Be Mined?
- What Technology Are Used?
- What Kind of Applications Are Targeted?
- Major Issues in Data Mining
- A Brief History of Data Mining and Data Mining Society
- Summary



# Major Issues in Data Mining (1)

---

- Mining Methodology
  - Mining various and new kinds of knowledge
  - Mining knowledge in multi-dimensional space
  - Data mining: An interdisciplinary effort
  - Boosting the power of discovery in a networked environment
  - Handling noise, uncertainty, and incompleteness of data
  - Pattern evaluation and pattern- or constraint-guided mining
- User Interaction
  - Interactive mining
  - Incorporation of background knowledge
  - Presentation and visualization of data mining results

# Major Issues in Data Mining (2)

---

- Efficiency and Scalability
  - Efficiency and scalability of data mining algorithms
  - Parallel, distributed, stream, and incremental mining methods
- Diversity of data types
  - Handling complex types of data
  - Mining dynamic, networked, and global data repositories
- Data mining and society
  - Social impacts of data mining

# Introduction

---

- Why Data Mining?
- What Is Data Mining?
- A Multi-Dimensional View of Data Mining
- What Kind of Data Can Be Mined?
- What Kinds of Patterns Can Be Mined?
- What Technology Are Used?
- What Kind of Applications Are Targeted?
- Major Issues in Data Mining
- A Brief History of Data Mining and Data Mining Society
- Summary



# 1.7 Data Mining Task Primitives

---

- How to construct a data mining query?
  - The primitives allow the user to interactively communicate with the data mining system during discovery to direct the mining process, or examine the findings

# 1.7 Data Mining Task Primitives

---

- The primitives specify:
  - (1) The set of task-relevant data – which portion of the database to be used
    - Database or data warehouse name
    - Database tables or data warehouse cubes
    - Condition for data selection
    - Relevant attributes or dimensions
    - Data grouping criteria

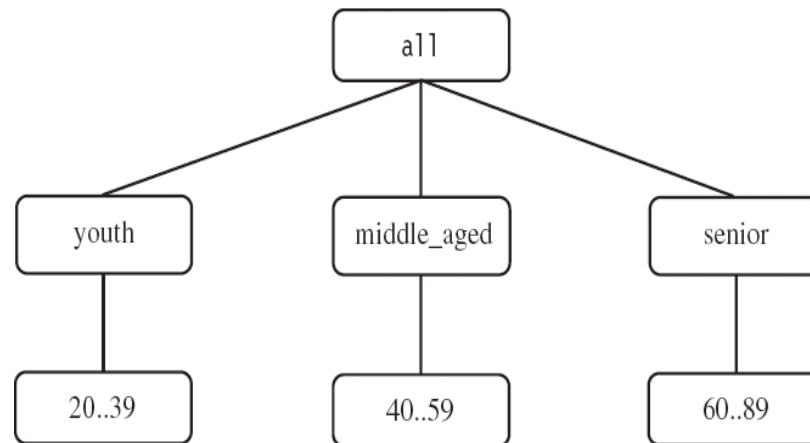
# 1.7 Data Mining Task Primitives

---

- The primitives specify:
  - (2) The kind of knowledge to be mined – what DB functions to be performed
    - Characterization
    - Discrimination
    - Association
    - Classification/prediction
    - Clustering
    - Outlier analysis
    - Other data mining tasks

# 1.7 Data Mining Task Primitives

(3) The background knowledge to be used – what domain knowledge, concept hierarchies, etc.



- (4) Interestingness measures and thresholds – support, confidence, etc.
- (5) Visualization methods – what form to display the result, e.g. rules, tables, charts, graphs, ...

# 1.7 Data Mining Task Primitives

---

- DMQL – Data Mining Query Language
  - Designed to incorporate these primitives
  - Allow user to interact with DM systems
  - Providing a **standardized language** like SQL

# A Brief History of Data Mining Society

---

- 1989 IJCAI Workshop on Knowledge Discovery in Databases
  - Knowledge Discovery in Databases (G. Piatetsky-Shapiro and W. Frawley, 1991)
- 1991-1994 Workshops on Knowledge Discovery in Databases
  - Advances in Knowledge Discovery and Data Mining (U. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, 1996)
- 1995-1998 International Conferences on Knowledge Discovery in Databases and Data Mining (KDD'95-98)
  - Journal of Data Mining and Knowledge Discovery (1997)
- ACM SIGKDD conferences since 1998 and SIGKDD Explorations
- More conferences on data mining
  - PAKDD (1997), PKDD (1997), SIAM-Data Mining (2001), (IEEE) ICDM (2001), etc.
- ACM Transactions on KDD starting in 2007

# Conferences and Journals on Data Mining

---

- KDD Conferences
  - ACM SIGKDD Int. Conf. on Knowledge Discovery in Databases and Data Mining (**KDD**)
  - SIAM Data Mining Conf. (**SDM**)
  - (IEEE) Int. Conf. on Data Mining (**ICDM**)
  - European Conf. on Machine Learning and Principles and practices of Knowledge Discovery and Data Mining (**ECML-PKDD**)
  - Pacific-Asia Conf. on Knowledge Discovery and Data Mining (**PAKDD**)
  - Int. Conf. on Web Search and Data Mining (**WSDM**)
- Other related conferences
  - DB conferences: ACM SIGMOD, VLDB, ICDE, EDBT, ICDT, ...
  - Web and IR conferences: WWW, SIGIR, WSDM
  - ML conferences: ICML, NIPS
  - PR conferences: CVPR,
- Journals
  - Data Mining and Knowledge Discovery (DAMI or DMKD)
  - IEEE Trans. On Knowledge and Data Eng. (TKDE)
  - KDD Explorations
  - ACM Trans. on KDD

# Where to Find References? DBLP, CiteSeer, Google

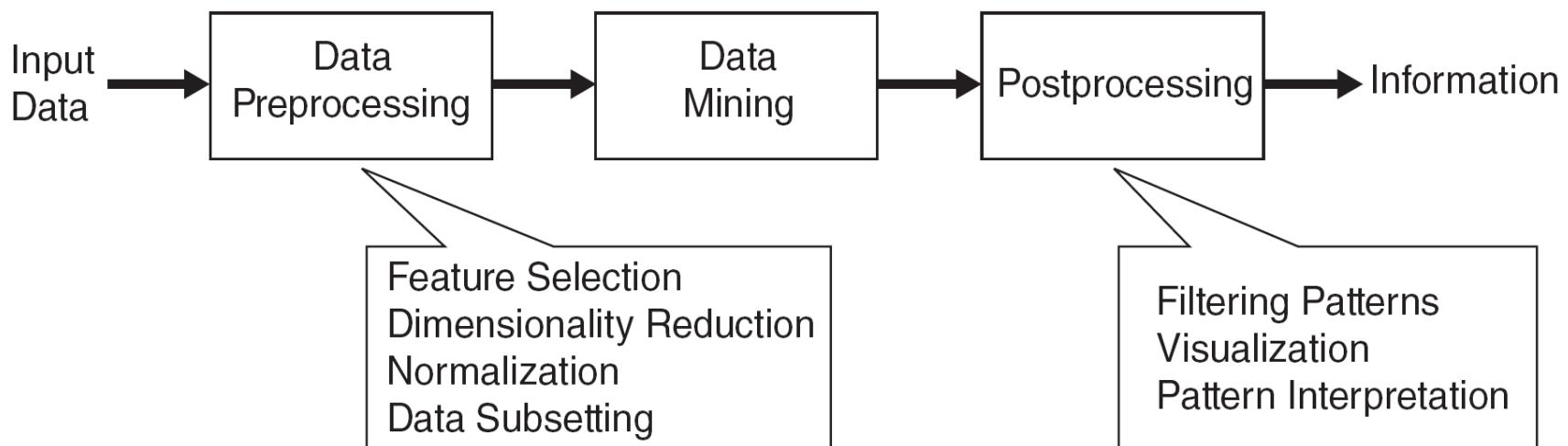
---

- Data mining and KDD (SIGKDD: CDROM)
  - Conferences: ACM-SIGKDD, IEEE-ICDM, SIAM-DM, PKDD, PAKDD, etc.
  - Journal: Data Mining and Knowledge Discovery, KDD Explorations, ACM TKDD
- Database systems (SIGMOD: ACM SIGMOD Anthology—CD ROM)
  - Conferences: ACM-SIGMOD, ACM-PODS, VLDB, IEEE-ICDE, EDBT, ICDT, DASFAA
  - Journals: IEEE-TKDE, ACM-TODS/TOIS, JIIS, J. ACM, VLDB J., Info. Sys., etc.
- AI & Machine Learning
  - Conferences: Machine learning (ML), AAAI, IJCAI, COLT (Learning Theory), CVPR, NIPS, etc.
  - Journals: Machine Learning, Artificial Intelligence, Knowledge and Information Systems, IEEE-PAMI, etc.
- Web and IR
  - Conferences: SIGIR, WWW, CIKM, etc.
  - Journals: WWW: Internet and Web Information Systems,
- Statistics
  - Conferences: Joint Stat. Meeting, etc.
  - Journals: Annals of statistics, etc.
- Visualization
  - Conference proceedings: CHI, ACM-SIGGraph, etc.
  - Journals: IEEE Trans. visualization and computer graphics, etc.

# What is Data Mining?

## • Many Definitions

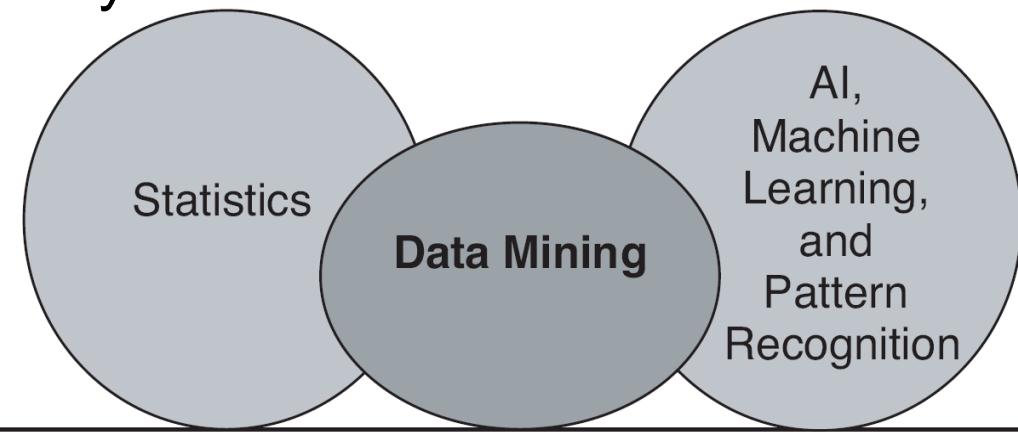
- Non-trivial extraction of implicit, previously unknown and potentially useful information from data
- Exploration & analysis, by automatic or semi-automatic means, of large quantities of data in order to discover meaningful patterns



# Origins of Data Mining

---

- Draws ideas from machine learning/AI, pattern recognition, statistics, and database systems
- Traditional techniques may be unsuitable due to data that is
  - Large-scale
  - High dimensional
  - Heterogeneous
  - Complex
  - Distributed
- A key component of the emerging field of data science and data-driven discovery



Database Technology, Parallel Computing, Distributed Computing

# Data Mining Tasks

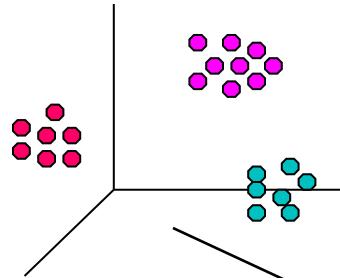
---

---

- Prediction Methods
  - Use some variables to predict unknown or future values of other variables.
- Description Methods
  - Find human-interpretable patterns that describe the data.

From [Fayyad, et.al.] Advances in Knowledge Discovery and Data Mining, 1996

# Data Mining Tasks ...



*Clustering*

**Data**

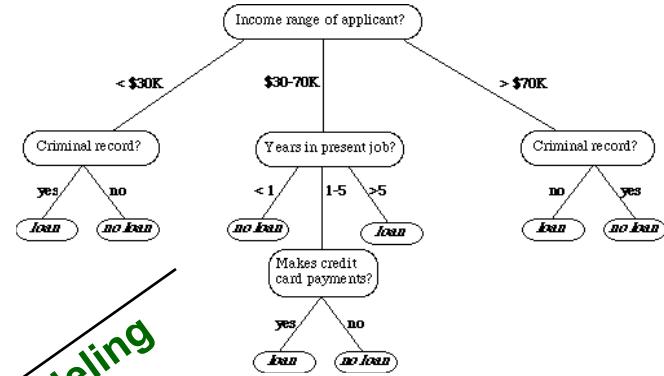
Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes
11	No	Married	60K	No
12	Yes	Divorced	220K	No
13	No	Single	85K	Yes
14	No	Married	75K	No
15	No	Single	90K	Yes

*Association Rules*



09/09/2020

Introduction to Data Mining, 2nd Edition  
Tan, Steinbach, Karpatne, Kumar



*Predictive Modeling*

*Anomaly Detection*

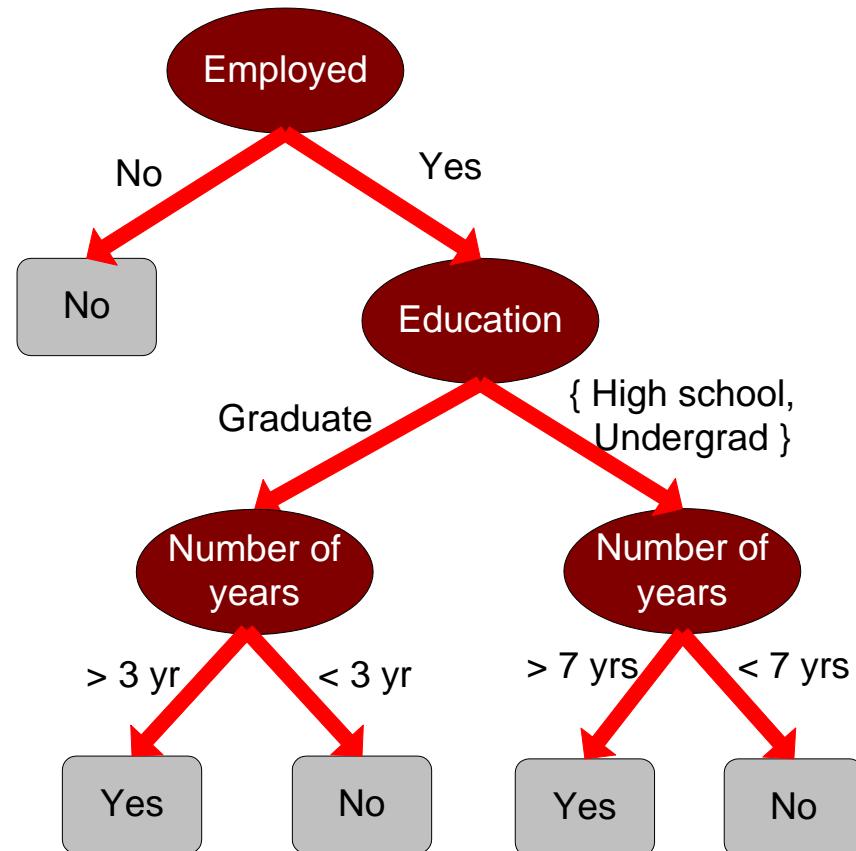


# Predictive Modeling: Classification

- Find a model for class attribute as a function of the values of other attributes

Tid	Employed	Level of Education	# years at present address	Credit Worthy
1	Yes	Graduate	5	Yes
2	Yes	High School	2	No
3	No	Undergrad	1	No
4	Yes	High School	10	Yes
...	...	...	...	...

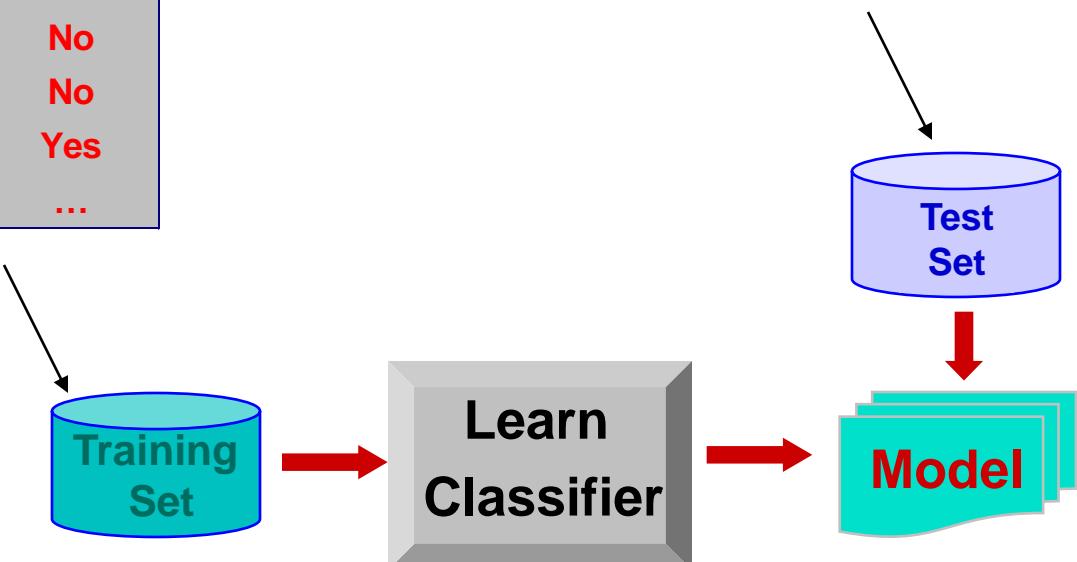
Model for predicting credit worthiness



# Classification Example

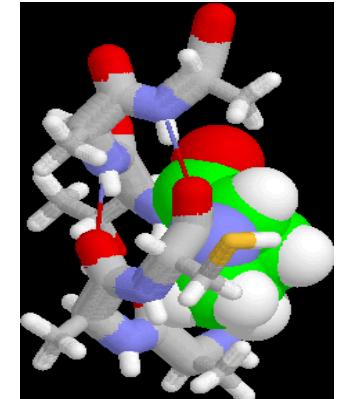
categorical categorical quantitative class				
Tid	Employed	Level of Education	# years at present address	Credit Worthy
1	Yes	Graduate	5	Yes
2	Yes	High School	2	No
3	No	Undergrad	1	No
4	Yes	High School	10	Yes
...	...	...	...	...

Tid	Employed	Level of Education	# years at present address	Credit Worthy
1	Yes	Undergrad	7	?
2	No	Graduate	3	?
3	Yes	High School	2	?
...	...	...	...	...



# Examples of Classification Task

- Classifying credit card transactions as legitimate or fraudulent
- Classifying land covers (water bodies, urban areas, forests, etc.) using satellite data
- Categorizing news stories as finance, weather, entertainment, sports, etc
- Predicting tumor cells as benign or malignant
- Classifying secondary structures of protein as alpha-helix, beta-sheet, or random coil



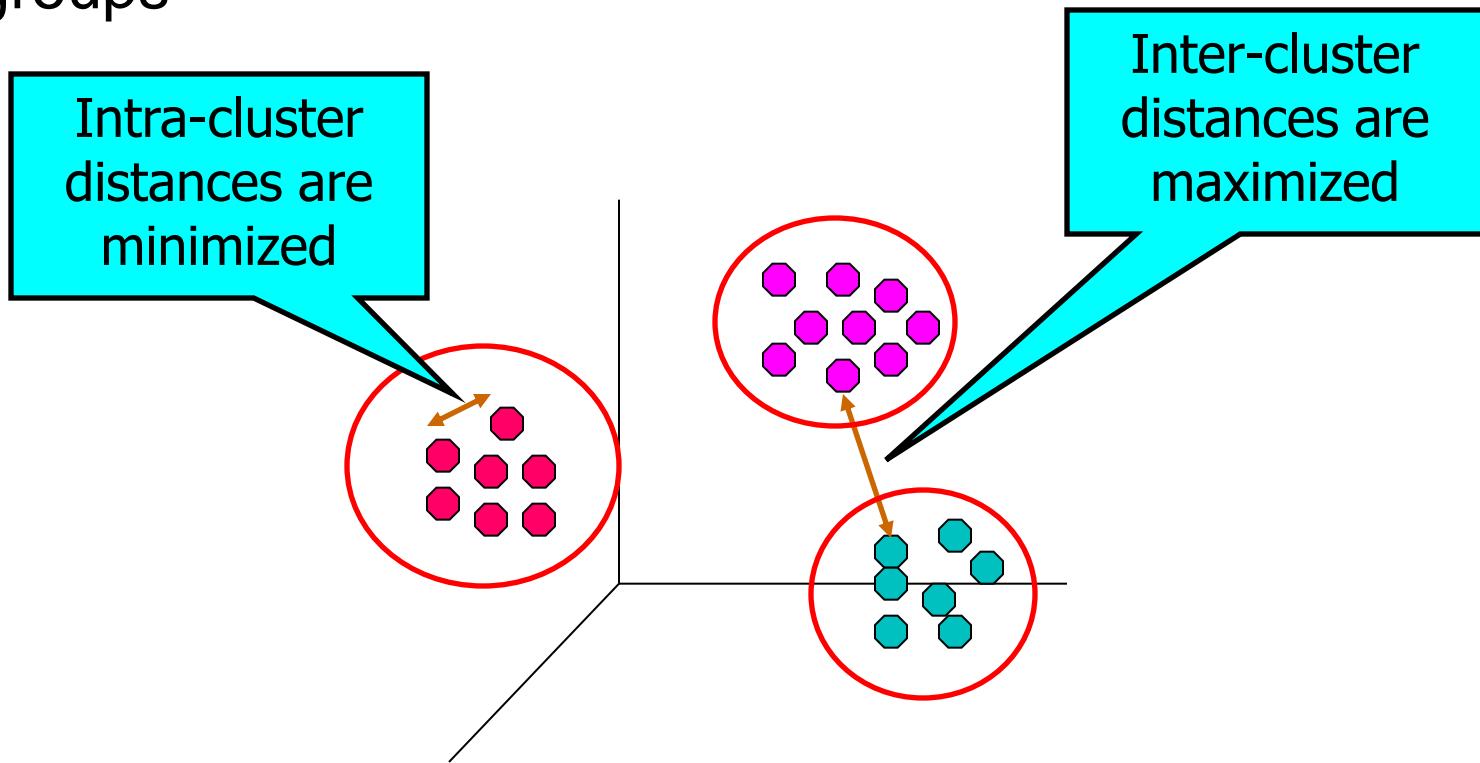
# Classification: Application

---

- Fraud Detection
  - **Goal:** Predict fraudulent cases in credit card transactions.
  - **Approach:**
    - ◆ Use credit card transactions and the information on its account-holder as attributes.
      - When does a customer buy, what does he buy, how often he pays on time, etc
    - ◆ Label past transactions as fraud or fair transactions. This forms the class attribute.
    - ◆ Learn a model for the class of the transactions.
    - ◆ Use this model to detect fraud by observing credit card transactions on an account.

# Clustering

- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups



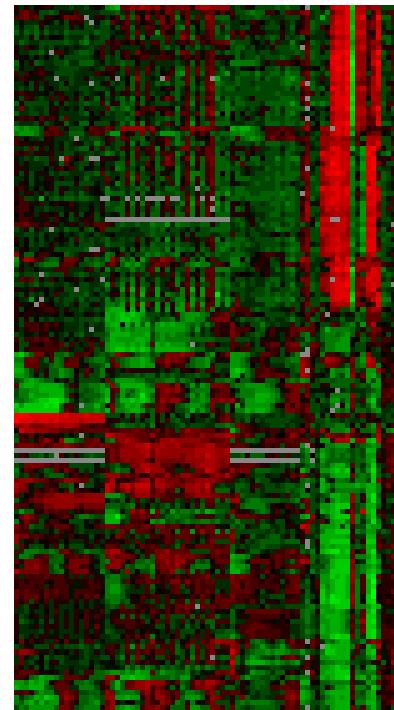
# Applications of Cluster Analysis

## • Understanding

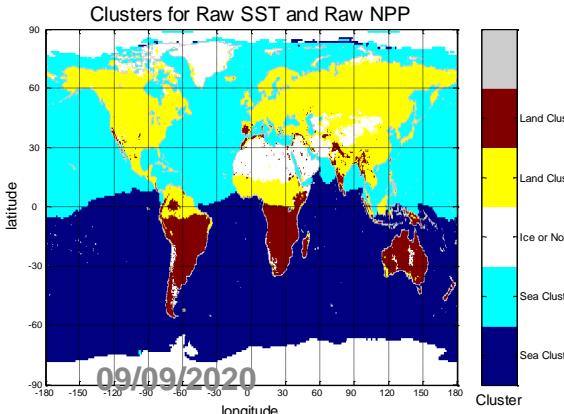
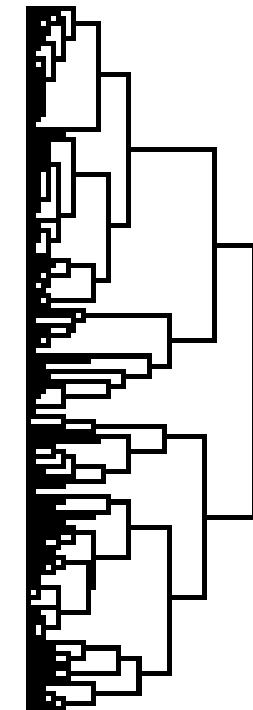
- Custom profiling for targeted marketing
- Group related documents for browsing
- Group genes and proteins that have similar functionality
- Group stocks with similar price fluctuations

## • Summarization

- Reduce the size of large data sets

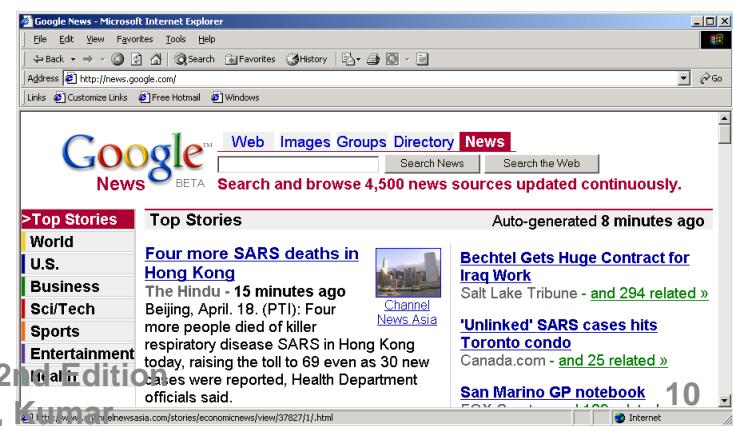


Courtesy: Michael Eisen



**Use of K-means to partition Sea Surface Temperature (SST) and Net Primary Production (NPP) into clusters that reflect the Northern and Southern Hemispheres.**

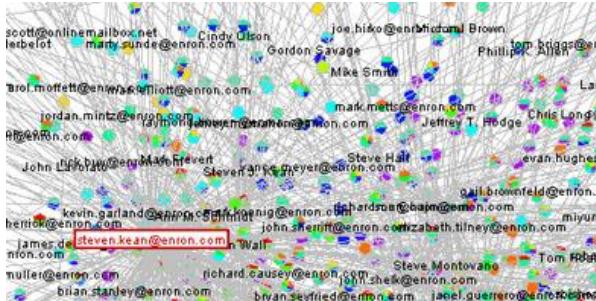
Introduction to Data Mining, 2nd Edition  
Tan, Steinbach, Karpatne, Kumar



# Clustering: Application

- Document Clustering:
  - **Goal:** To find groups of documents that are similar to each other based on the important terms appearing in them.
  - **Approach:** To identify frequently occurring terms in each document. Form a similarity measure based on the frequencies of different terms. Use it to cluster.

Enron email dataset



# Association Rule Discovery: Definition

- Given a set of records each of which contain some number of items from a given collection
  - Produce dependency rules which will predict occurrence of an item based on occurrences of other items.

TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

Rules Discovered:

$\{\text{Milk}\} \rightarrow \{\text{Coke}\}$

$\{\text{Diaper}, \text{Milk}\} \rightarrow \{\text{Beer}\}$

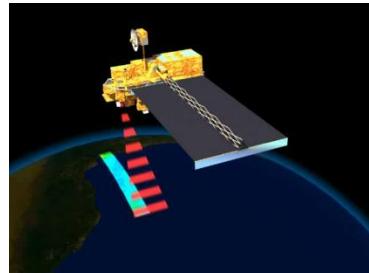
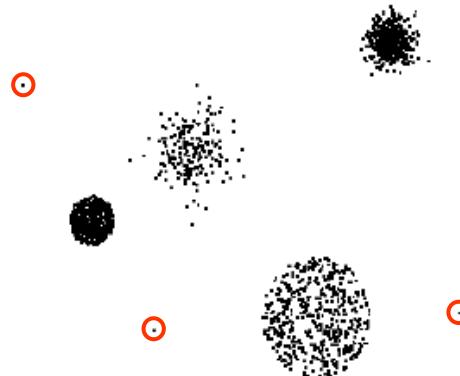
# Association Analysis: Applications

---

- Market-basket analysis
  - Rules are used for sales promotion, shelf management, and inventory management
- Telecommunication alarm diagnosis
  - Rules are used to find combination of alarms that occur together frequently in the same time period
- Medical Informatics
  - Rules are used to find combination of patient symptoms and test results associated with certain diseases

# Deviation/Anomaly/Change Detection

- Detect significant deviations from normal behavior
- Applications:
  - Credit Card Fraud Detection
  - Network Intrusion Detection
  - Identify anomalous behavior from sensor networks for monitoring and surveillance.



# Motivating Challenges

---

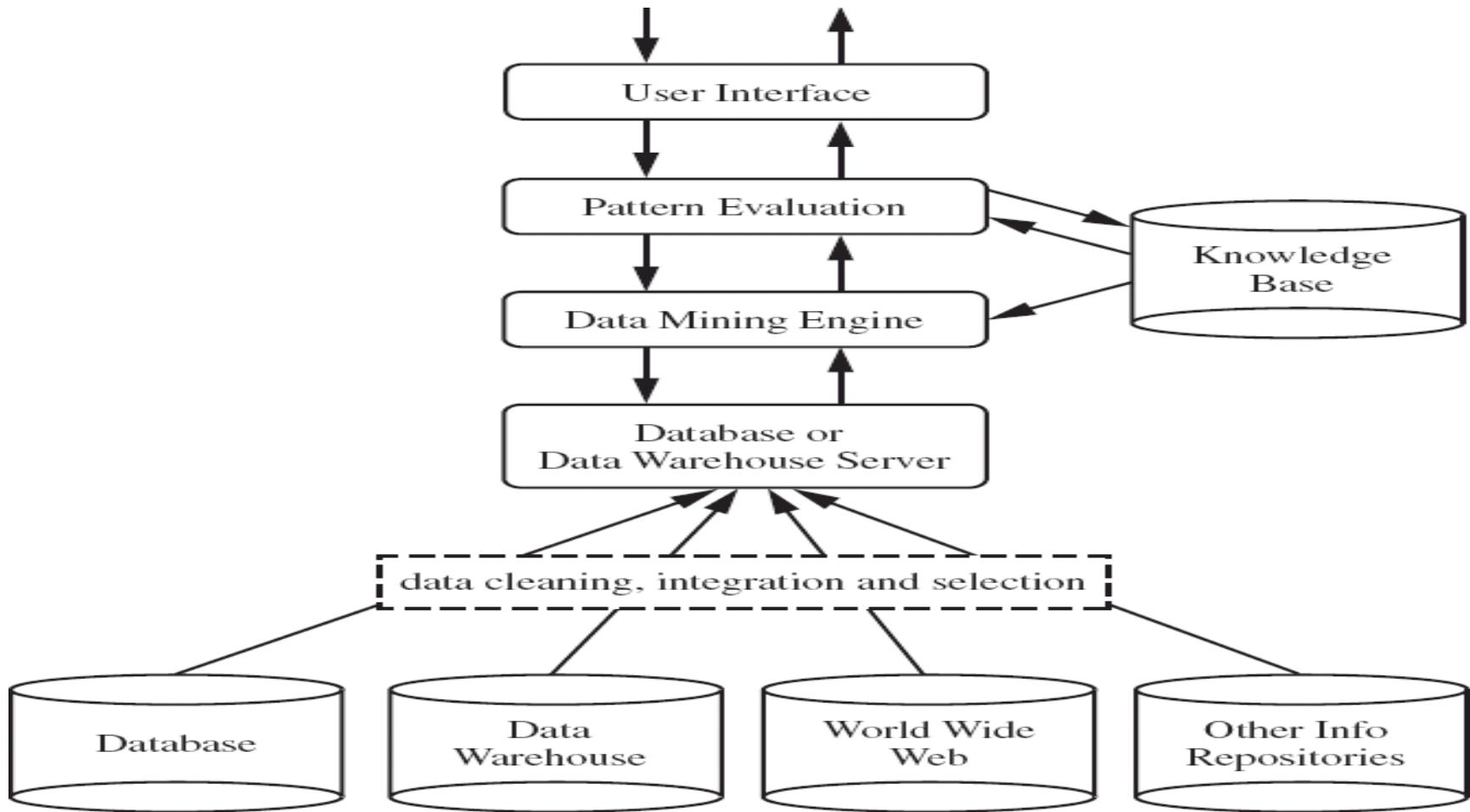
---

- Scalability
- High Dimensionality
- Heterogeneous and Complex Data
- Data Ownership and Distribution
- Non-traditional Analysis

# A typical DM System Architecture

- Database, data warehouse, WWW or other information repository (**store data**)
- Database or data warehouse server (**fetch and combine data**)
- Knowledge base (**turn data into meaningful groups according to domain knowledge**)
- Data mining engine (**perform mining tasks**)
- Pattern evaluation module (**find interesting patterns**)
- User interface (**interact with the user**)

# A typical DM System Architecture



## DEPARTMENT OF INFORMATION TECHNOLOGY, NITK, SURATHKAL

<b>Course Code:</b>	<b>IT-414</b>	<b>Course Name:</b>	<b>Data Warehousing and Data Mining</b>
<b>Core/Elective/MLC:</b>	Elective	<b>L-T-P:</b>	(3-0-2) 4
<b>Pre-requisites:</b>	Database System	<b>Contact Hours:</b>	8 AM to 5 PM
<b>Type of course:</b> <b>(Lecture/Tutorial/ Seminar/Project)</b>	Lecture/Project	<b>Course Assessment Methods:</b> <b>(Both Continuous and Semester-End Assessment)</b>	Mid Sem Exam=20% Lab Assignments=15% End Sem Exam=35% Course Project=30%

**Course Description:** The Course Covers the following main aspects

- What is data warehouse and data Mining, importance and challenges of Data Mining
- Data Warehouse Architecture and schemas
- Data Mining Techniques and its real world application

**Course Objectives:** The aim of the course is

- To understand the data warehouse architecture, data warehouse schemas, Multidimensional model
- To understand and apply different data preprocessing techniques
- To learn different feature selection methods
- To study different Association Rule mining methods to identify important association between frequent itemsets and to apply frequent itemset mining on the data
- To understand and apply different classification and clustering methods

**Course Outcomes:** After the completion of this course, the student will be able to:

<b>CO1</b>	Understand the differences of OLTP and OLAP, understand the data warehousing techniques, data warehouse architecture and process of constructing data warehouses with different models
<b>CO2</b>	Learn data mining techniques such as association rule mining, various classification and clustering techniques and apply for knowledge representation
<b>CO3</b>	Understand the different data preprocessing techniques along with feature selection methods
<b>CO4</b>	To apply DM techniques to real world dataset and Analyze results

#### **Detailed Course Plan:**

Weeks	Topics
Week 1 & 2	Introduction to data mining: Motivation and significance of data mining, data mining functionalities, interestingness measures, classification of data mining system, major issues in data mining, Data Mining task primitives, types of data for data mining
Week 3 & 4:	Association Analysis: Basic concepts: Frequent itemsets, Association rules, Market basket analysis example. Frequent ItemSet Mining methods: Apriori algorithm, Rule Generation, FP-growth Algorithm; Constraint based Association analysis.
Week 5 & 6	Compact Representation of Frequent ItemSets, Colossal itemset Mining methods, multilevel and multidimensional association rule mining
Week 7 & 8:	Data preprocessing: Why preprocess the data?, Data cleaning: Missing value, Noisy data; Data integration and transformation: Correlation analysis, Min Max, Z-score and decimal scaling normalization, ; Attribute subset selection. Filter based and Wrapper based feature selection methods, Chi square method, Heuristic methods for feature selection
Week 9 & 10:	PCA, OLTP/OLAP differences; Data Warehouse: What is Data Warehouse and data warehouse schemas, Data Cubes and multidimensional data model, data warehouse architecture
Week 11 & 12:	Classification and clustering methods: Density based clusters and comparison with partition and hierarchical methods, Association based classification, ensemble based classifiers, Boosted Decision tree classifier,
Week 13	Data mining on complex data and applications, trends , Research paper discussion.

**Reference Books:**

- 1 Han, J. and Kamber, M., "Data Mining - Concepts and Techniques", 3rd Ed., Morgan Kaufmann Series, (Elsevier), 2008.
- 2 Alex Berson , S. J. Smith, "Data Warehousing, Data Mining & OLAP" , McGraw Hill
- 3 Tan, P.N., Steinbach, M. and Kumar, V., "Introduction to Data Mining", Addison Wesley – Pearson, 2006
- 4 Pujari, A. K., "Data Mining Techniques", 4<sup>th</sup> Ed., Sangam Books.
- 5 Oded Maimon, Lior Rokach, The Data Mining and Knowledge Discovery Handbook, Springer, 2005.
- 6 S. Weiss and N. Indurkhy, Predictive Data-Mining: A Practical Guide, Morgan Kaufmann, 1998
- 7 S. Weiss, N. Indurkhy, T. Zhang and F. Damerau, Text Mining: Predictive Methods for Analyzing Unstructured Information, Springer, 2004.

**Course Instructors:**

**Dr. Nagamma Patil**

# Outline

- Basics
- Market Basket Analysis: A Motivating Example
- Preliminaries
- Frequent Itemset Mining
- Apriori Algorithm
- Frequent Pattern growth (FP-growth) Algorithm
- Frequent Closed Itemset Mining
- Frequent Closed Itemset Mining from High Dimensional Data

# Basics

- Frequent patterns are patterns (such as itemsets, subsequences, or substructures) that appear in a data set frequently.
- A set of items, such as milk and bread, that appear frequently together in a transaction data set is a frequent itemset.
- A subsequence, such as buying first a PC, then a digital camera, and then a memory card, if it occurs frequently in a shopping history database, is a frequent sequential pattern.
- A substructure can refer to different structural forms, such as subgraphs, subtrees, or sublattices. If a substructure occurs frequently, it is called a frequent structured pattern.

# Market Basket Analysis: A Motivating Example

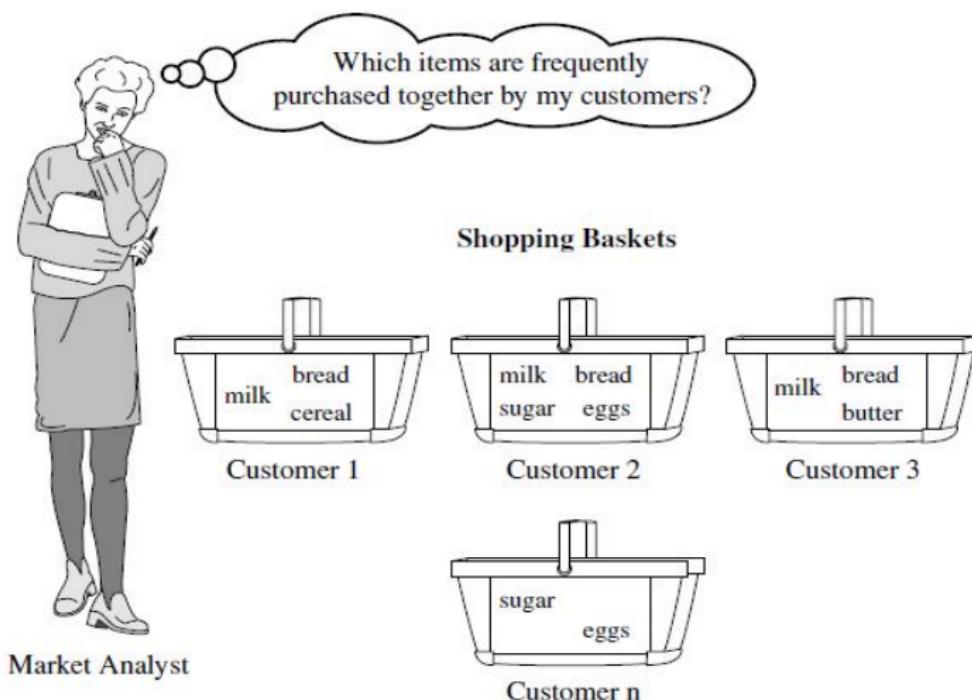


Figure 1. Apriori

# Market Basket Analysis

- Frequent itemset mining leads to the discovery of associations and correlations among items in large transactional or relational data sets.
- A typical example of frequent itemset mining is market basket analysis.
- Data: collection of transactions of customers.
- Goal: find sets of products frequently occurring together.
- The discovery of associations helps in many business decision making processes, such as catalog design and customer shopping behavior analysis.

# Applications

- Market basket analysis.
- Catalog design.
- Customer shopping behavior analysis.
- Web log analysis.
- DNA sequence analysis.
- Sale campaign analysis.
- Software bug detection.
- Chemical Compound Prediction.
- Text analysis.

# Preliminaries

Let the Dataset  $D$  consist of  $m$  number of transactions (rows) and  $n$  of attributes or products (features)

- $R = \{r_1, r_2, \dots, r_m\}$
- $F = \{f_1, f_2, \dots, f_n\}$
- Each row  $r_i$  has unique row identifier,  $rid$  and consist of set of products (features).
- A non-empty subset of features  $X \subseteq F$  is defined as an itemset.
- Let  $r(f_j)$  signify the rows in which  $j^{th}$  feature of the dataset is present.
- A non-empty subset of rids  $Y \subseteq R$  is defined as rowset.
- Let  $f(r_i)$  signify the features present in the  $i^{th}$  row of the dataset.

# Preliminaries

## Example 1

Table 1 shows an example of a Dataset  $D$  consisting of 8 rows, where each row is described with unique row identifier ( $rid$ ),  $R = \{1, 2, 3, 4, 5, 6, 7, 8\}$  and 11 features,  $F = \{f_1, f_2, f_3, f_4, f_5, f_6, f_7, f_8, f_9, f_{10}, f_{11}\}$ .

Table 1  
Dataset  $D$

row id ( $rid$ )	features
1	$f_1, f_2, f_4, f_6, f_{10}$
2	$f_1, f_2, f_4, f_7, f_8$
3	$f_2, f_4, f_7, f_8$
4	$f_1, f_2, f_6, f_8, f_9, f_{10}$
5	$f_1, f_3, f_4, f_7, f_8, f_{10}$
6	$f_2, f_4, f_9$
7	$f_5, f_7$
8	$f_5, f_{11}$

# Preliminaries

## Definition 1 (Support)

The number of rows in which an itemset  $X$  occurs is called the support of an itemset, denoted by  $\text{sup}(X)$ .

## Example 2

In Table 1, the support of an itemset  $X = \{f_2, f_4, f_7, f_8\}$ ,  $\text{sup}(X)$  is 2.

## Definition 2 (Support Set)

The rows in which an itemset  $X$  occurs is called support set of an itemset, denoted by  $\text{supset}(X)$ .

## Example 3

In Table 1, the support set of an itemset  $X = \{f_2, f_4, f_7, f_8\}$ ,  $\text{supset}(X)$  is 23.

# Preliminaries

## Definition 3 (Cardinality)

The number of items in an itemset  $X$  is called as the cardinality of an itemset, denoted by  $card(X)$ .

## Example 4

In Table 1, the cardinality of an itemset  $X = \{f_2, f_4, f_7, f_8\}$ ,  $card(X)$  is 4.

## Definition 4 (Frequent Itemset)

An itemset  $X$  is called frequent itemset if and only if  $sup(X) \geq minsup$ , where  $minsup$  is user specified least support threshold.

## Example 5

In Table 1, the itemset  $X = \{f_2, f_8\}$  is frequent itemset with minimum support threshold set to 2,  $sup(X) \geq 2$ .

# Preliminaries

## Definition 5 (Association Rule)

Let  $A$  and  $B$  be the set of items. An association rule is an implication of the form  $A \Rightarrow B$ , where  $A \subset F$ ,  $B \subset F$  and  $A \cap B = \emptyset$ . The association rule  $A \Rightarrow B$  holds in the dataset with **support**  $s$  and has **confidence**  $c$ .

**Support**  $s$ , is the percentage of transactions in  $D$  that contain  $A \cup B$  (i.e., the union of sets  $A$  and  $B$ , or say, both  $A$  and  $B$ ).

**Confidence**  $c$ , is the percentage of transactions in  $D$  containing  $A$  that also contain  $B$ . This is taken to be the conditional probability,  $P(B|A)$ .

$$\text{support}(A \Rightarrow B) = P(A \cup B) \quad (1)$$

$$\text{confidence}(A \Rightarrow B) = P(B|A) = \frac{\text{support}(A \cup B)}{\text{support}(A)} \quad (2)$$

# Frequent Itemset Mining

Table 2  
Dataset  $D$

TID	Items Bought
1	Beer, Nuts, Chips
2	Beer, Coffee, Chips
3	Beer, Chips, Eggs
4	Nuts, Eggs, Milk
5	Nuts, Coffee, Chips, Eggs, Milk

- Problem: To Mine the Frequent Itemsets with minimum support threshold (*minsup*) set to 50% and minimum confidence threshold (*minconf*) set to 50%.
- Frequent Itemsets are: Beer:3, Nuts:3, Chips:4, Eggs:3, {Beer, Chips}:3.
- Example of association rules  
 $\text{Beer} \rightarrow \text{Chips}$  (60%, 100%).  
 $\text{Chips} \rightarrow \text{Beer}$  (60%, 75%).

# Frequent Itemset Mining

- Frequent Itemset Mining Algorithms
  - Apriori Algorithm
  - Frequent Pattern growth (FP-growth) algorithm
- Frequent Closed Itemset Mining Algorithms
- Frequent Maximal Itemset Mining Algorithms
- Frequent Colossal Itemset Mining Algorithms
- Frequent Colossal Closed Itemset Mining Algorithms

# Apriori Algorithm

- Apriori is an algorithm proposed by R. Agrawal and R. Srikant in 1994 for mining frequent itemsets from transactional datasets for generating association rules.
- The name of the algorithm is based on the fact that the algorithm uses prior knowledge of frequent itemset properties.
- Apriori employs an iterative approach known as a level-wise search, where k-itemsets are used to explore  $(k+1)$ -itemsets.
- First, the set of frequent 1-itemsets is found by scanning the database to accumulate the count for each item, and collecting those items that satisfy minimum support.

# Apriori Algorithm

- The resulting set is denoted  $L_1$ . Next,  $L_1$  is used to find  $L_2$ , the set of frequent 2-itemsets, which is used to find  $L_3$ , and so on, until no more frequent k-itemsets can be found.
- The finding of each  $L_k$  requires one full scan of the database.
- To improve the efficiency of the level-wise generation of frequent itemsets, an important property called the Apriori property.
- Apriori property: All nonempty subsets of a frequent itemset must also be frequent.
- The property belongs to a special category of properties called antimonotone in the sense that if a set cannot pass a test, all of its supersets will fail the same test as well.

# Apriori Algorithm

- Apriori Algorithm has two steps

- The Join step
- The Prune step

- **The Join step:**

- To find  $L_k$ , a set of candidate k-itemsets is generated by joining  $L_{k-1}$  with itself. This set of candidates is denoted  $C_k$ .
- Apriori assumes that items within a transaction or itemset are sorted in lexicographic order.
- The join,  $L_{k-1} \bowtie L_{k-1}$ , is performed, where members of  $L_{k-1}$  are joinable if their first  $(k-2)$  items are in common.

# Apriori Algorithm

- **The Prune step:**

- $C_k$  is a superset of  $L_k$ , that is, its members may or may not be frequent, but all of the frequent k-itemsets are included in  $C_k$ .
- A scan of the database to determine the count of each candidate in  $C_k$  would result in the determination of  $L_k$ .
- Any (k-1)-itemset that is not frequent cannot be a subset of a frequent k-itemset.
- If any (k-1)-subset of a candidate k-itemset is not in  $L_{k-1}$ , then the candidate cannot be frequent either and so can be removed from  $C_k$ .

# Apriori Algorithm

Table 3  
Dataset  $D$

TID	List of items
1	I1, I2, I5
2	I2, I4
3	I2, I3
4	I1, I2, I4
5	I1, I3
6	I2, I3
7	I1, I3
8	I1, I2, I3, I5
9	I1, I2, I3

# Apriori Algorithm

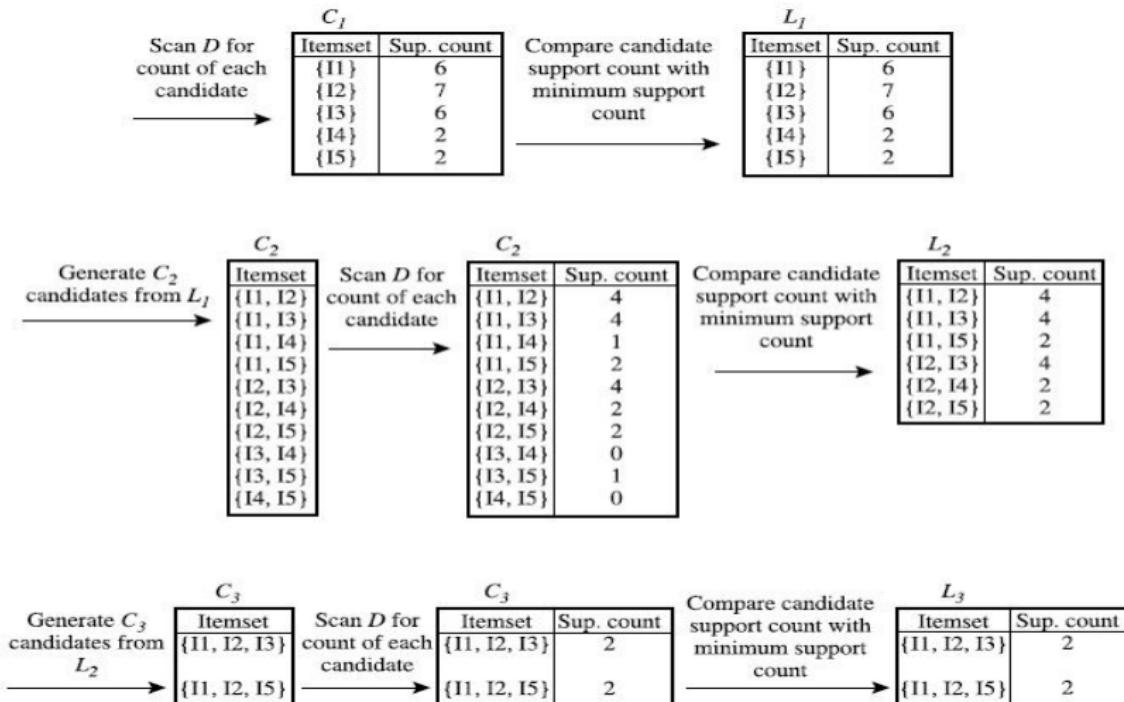


Figure 2. Steps Apriori Algorithm

# Apriori Algorithm

$$\begin{aligned}\text{Join: } C_3 &= L_2 \bowtie L_2 = \{\{I1, I2\}, \{I1, I3\}, \{I1, I5\}, \{I2, I3\}, \{I2, I4\}, \{I2, I5\}\} \bowtie \\ &\quad \{\{I1, I2\}, \{I1, I3\}, \{I1, I5\}, \{I2, I3\}, \{I2, I4\}, \{I2, I5\}\} \\ &= \{\{I1, I2, I3\}, \{I1, I2, I5\}, \{I1, I3, I5\}, \{I2, I3, I4\}, \{I2, I3, I5\}, \{I2, I4, I5\}\}.\end{aligned}$$

Prune using the Apriori property: All nonempty subsets of a frequent itemset must also be frequent. Do any of the candidates have a subset that is not frequent?

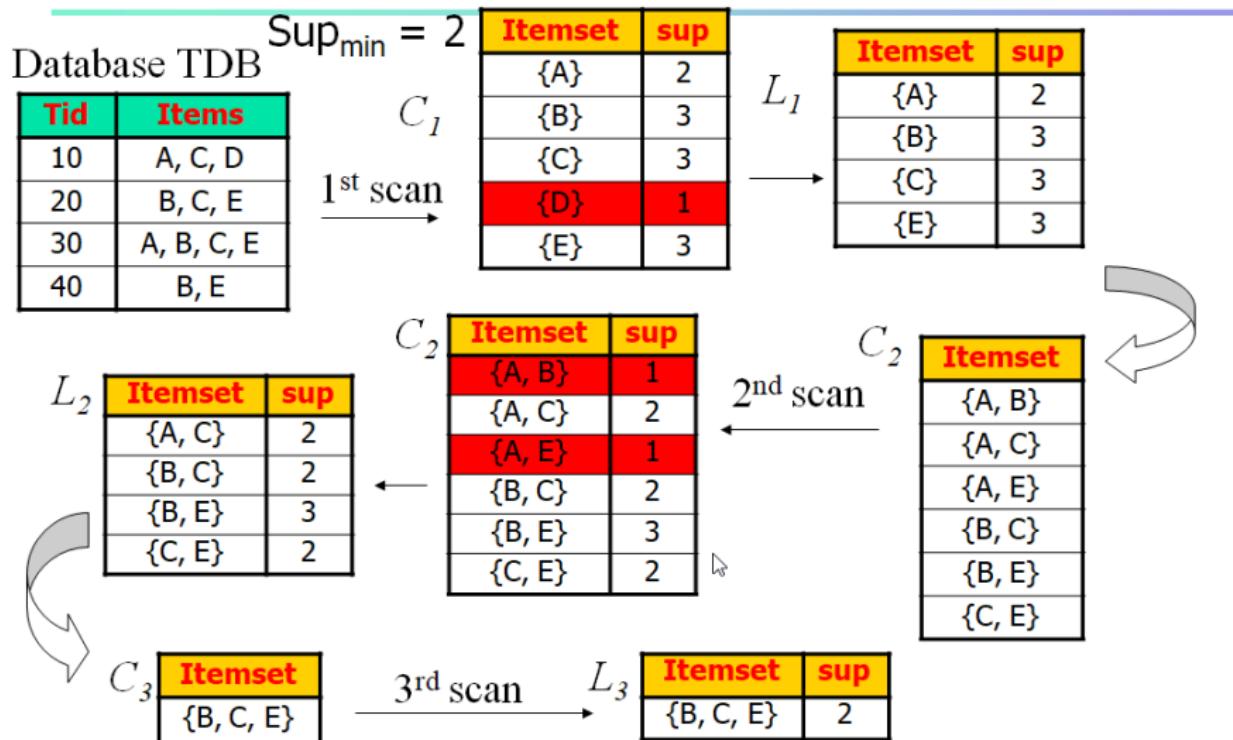
- The 2-item subsets of  $\{I1, I2, I3\}$  are  $\{I1, I2\}$ ,  $\{I1, I3\}$ , and  $\{I2, I3\}$ . All 2-item subsets of  $\{I1, I2, I3\}$  are members of  $L_2$ . Therefore, keep  $\{I1, I2, I3\}$  in  $C_3$ .
- The 2-item subsets of  $\{I1, I2, I5\}$  are  $\{I1, I2\}$ ,  $\{I1, I5\}$ , and  $\{I2, I5\}$ . All 2-item subsets of  $\{I1, I2, I5\}$  are members of  $L_2$ . Therefore, keep  $\{I1, I2, I5\}$  in  $C_3$ .
- The 2-item subsets of  $\{I1, I3, I5\}$  are  $\{I1, I3\}$ ,  $\{I1, I5\}$ , and  $\{I3, I5\}$ .  $\{I3, I5\}$  is not a member of  $L_2$ , and so it is not frequent. Therefore, remove  $\{I1, I3, I5\}$  from  $C_3$ .
- The 2-item subsets of  $\{I2, I3, I4\}$  are  $\{I2, I3\}$ ,  $\{I2, I4\}$ , and  $\{I3, I4\}$ .  $\{I3, I4\}$  is not a member of  $L_2$ , and so it is not frequent. Therefore, remove  $\{I2, I3, I4\}$  from  $C_3$ .
- The 2-item subsets of  $\{I2, I3, I5\}$  are  $\{I2, I3\}$ ,  $\{I2, I5\}$ , and  $\{I3, I5\}$ .  $\{I3, I5\}$  is not a member of  $L_2$ , and so it is not frequent. Therefore, remove  $\{I2, I3, I5\}$  from  $C_3$ .
- The 2-item subsets of  $\{I2, I4, I5\}$  are  $\{I2, I4\}$ ,  $\{I2, I5\}$ , and  $\{I4, I5\}$ .  $\{I4, I5\}$  is not a member of  $L_2$ , and so it is not frequent. Therefore, remove  $\{I2, I4, I5\}$  from  $C_3$ .

Therefore,  $C_3 = \{\{I1, I2, I3\}, \{I1, I2, I5\}\}$  after pruning.

# Apriori Algorithm

- Generating association rules. The frequent itemset considered is {I1, I2, I5}
- The nonempty subsets of frequent itemset are {I1, I2}, {I1, I5}, {I2, I5}, {I1}, {I2}, and {I5}.
- The resulting association rules are as shown below, each listed with its confidence:
  - $I1 \wedge I2 \Rightarrow I5$ , confidence =  $2/4 = 50\%$
  - $I1 \wedge I5 \Rightarrow I2$ , confidence =  $2/2 = 100\%$
  - $I2 \wedge I5 \Rightarrow I1$ , confidence =  $2/2 = 100\%$
  - $I1 \Rightarrow I2 \wedge I5$ , confidence =  $2/6 = 33\%$
  - $I2 \Rightarrow I1 \wedge I5$ , confidence =  $2/7 = 29\%$
  - $I5 \Rightarrow I1 \wedge I2$ , confidence =  $2/2 = 50\%$

# Apriori Algorithm



# FP-Tree/FP-Growth Algorithm

- Use a compressed representation of the database using an FP-tree
- Once an FP-tree has been constructed, it uses a recursive divide-and-conquer approach to mine the frequent itemsets.

## Building the FP-Tree

1. Scan data to determine the support count of each item.  
Infrequent items are discarded, while the frequent items are sorted in decreasing support counts.
2. Make a second pass over the data to construct the FPtree.  
As the transactions are read, before being processed, their items are sorted according to the above order.

# FP-tree Example: step 1

Step 1: Scan DB for the first time to generate L  
(minimum support=3)

TID	Items bought
100	{f, a, c, d, g, i, m, p}
200	{a, b, c, f, l, m, o}
300	{b, f, h, j, o}
400	{b, c, k, s, p}
500	{a, f, c, e, l, p, m, n}



L

Item	frequency
f	4
c	4
a	3
b	3
m	3
p	3

By-Product of First Scan  
of Database

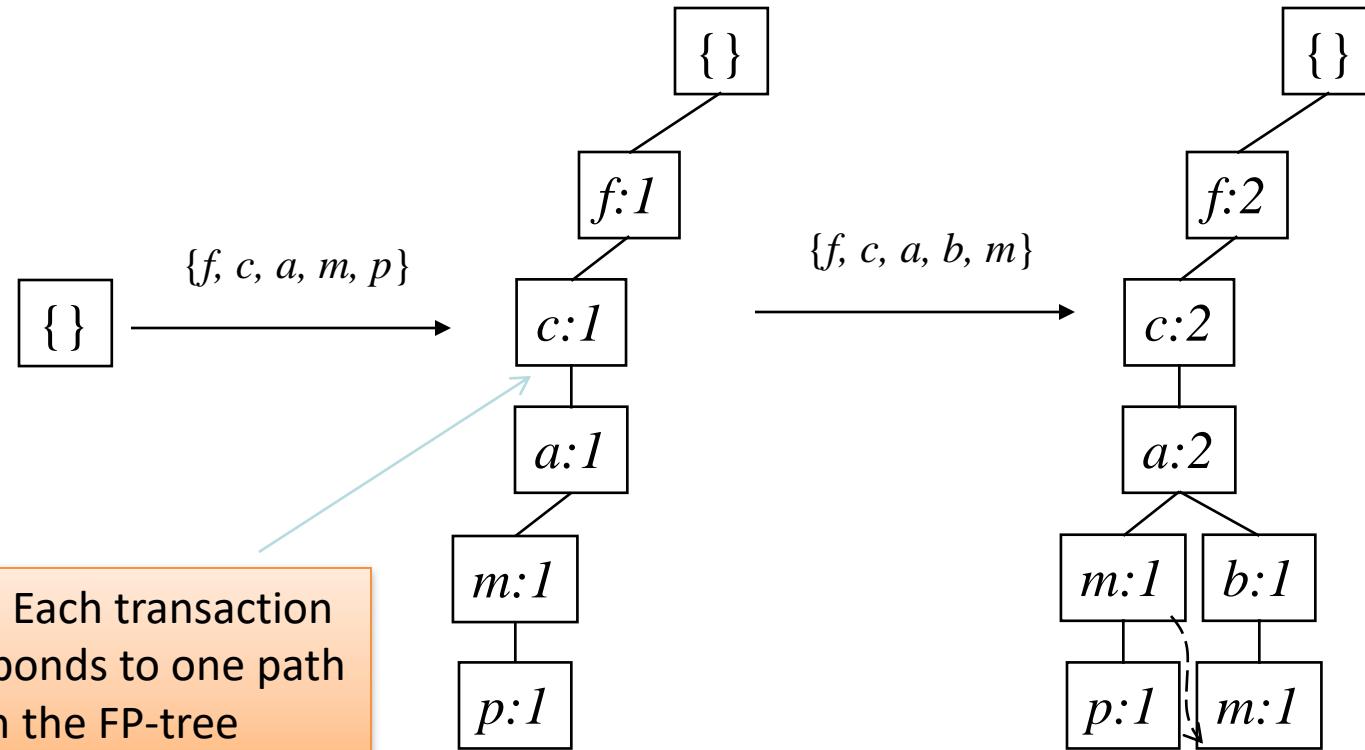
# FP-tree Example: step 2

**Step 2: scan the DB for the second time, order frequent items in each transaction**

<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

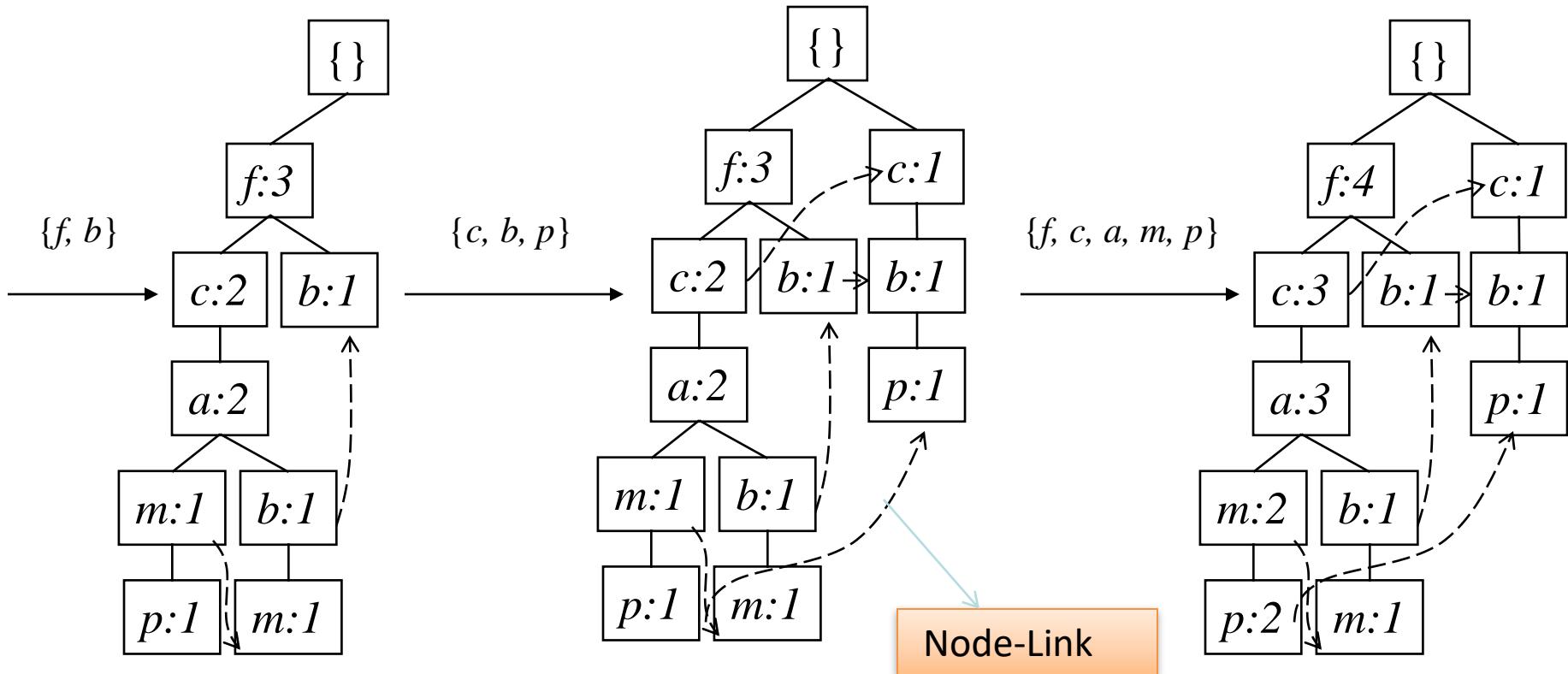
# FP-tree Example: step 2

## Step 2: construct FP-tree



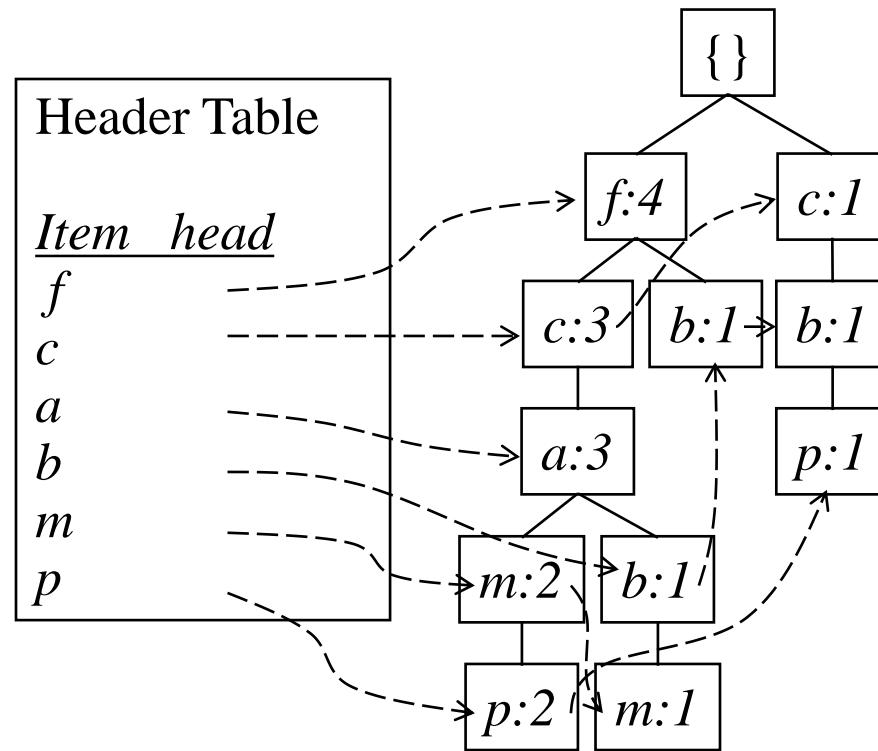
# FP-tree Example: step 2

## Step 2: construct FP-tree



# Construction Example

## Final FP-tree



# FP-growth: Mining Frequent Patterns Using FP-tree

# 3 Major Steps

---

Starting the processing from the end of list L:

Step 1:

Construct **conditional pattern base** for each item in the header table

Step 2

Construct **conditional FP-tree** from each conditional pattern base

Step 3

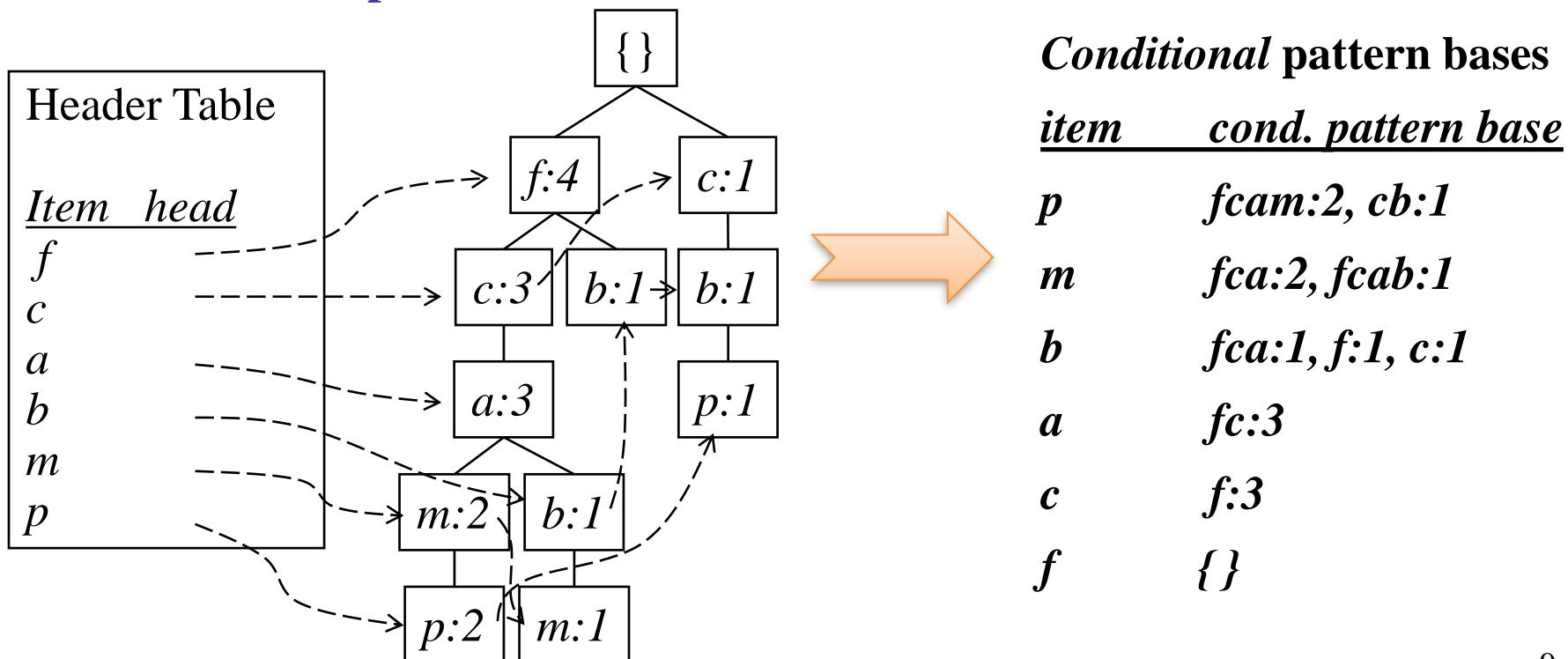
**Recursively mine** conditional FP-trees and grow frequent patterns obtained so far. If the conditional FP-tree contains a **single path**, simply enumerate all the patterns

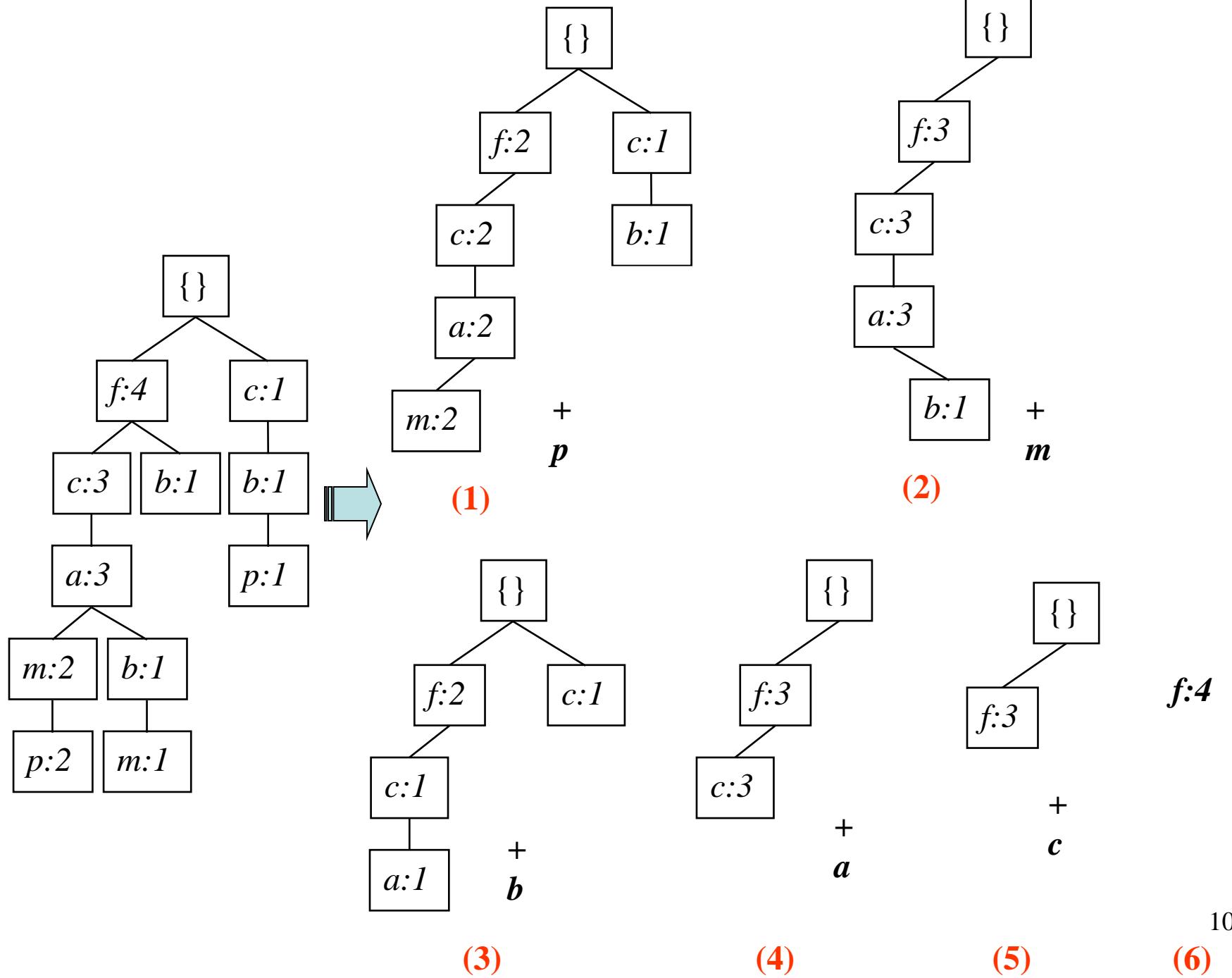
---

# Step 1: Construct Conditional Pattern Base

---

- Starting at the bottom of frequent-item header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item
- Accumulate all of **transformed prefix paths** of that item to form a **conditional pattern base**





# Conditional Pattern Bases and Conditional FP-Tree

Item	Conditional pattern base	Conditional FP-tree
p	$\{(fcam:2), (cb:1)\}$	$\{(c:3)\} p$
m	$\{(fca:2), (fcab:1)\}$	$\{(f:3, c:3, a:3)\} m$
b	$\{(fca:1), (f:1), (c:1)\}$	Empty
a	$\{(fc:3)\}$	$\{(f:3, c:3)\} a$
c	$\{(f:3)\}$	$\{(f:3)\} c$
f	Empty	Empty

order of L

min\_sup = 3

1	<b>f, c, a, m, p</b>
2	<b>f, c, a, b, m</b>
3	<b>f, b</b>
4	<b>c, b, p</b>
5	<b>f, c, a, m, p</b>

1	<b>f, c, a, m</b>
4	<b>c, b</b>
5	<b>f, c, a, m</b>

+ p

1	<b>c</b>
4	<b>c</b>
5	<b>c</b>

+ p

p: 3  
cp: 3

1	<b>f, c, a</b>
2	<b>f, c, a, b</b>
5	<b>f, c, a</b>

+ m

1	<b>f, c, a</b>
2	<b>f, c, a</b>
5	<b>f, c, a</b>

+ m

m: 3  
fm: 3  
cm: 3  
am: 3  
fcm: 3  
fam: 3  
cam: 3  
fcam: 3

2	<b>f, c, a</b>
3	<b>f</b>
4	<b>c</b>

+ b

b: 3

1	<b>f, c</b>
2	<b>f, c</b>
5	<b>f, c</b>

+ a

a: 3  
fa: 3  
ca: 3  
fca: 3

1	<b>f</b>
2	<b>f</b>
4	
5	<b>f</b>

+ c

c: 4  
fc: 3

**f: 1,2,3,5**

f: 4

# Discussion

- Advantages of FP-Growth
  - only 2 passes over data-set
  - no candidate generation
  - much faster than Apriori
- Disadvantages of FP-Growth
  - FP-Tree may not fit in memory!!
  - FP-Tree is expensive to build

# Challenges

- A major challenge in mining frequent itemsets from a large data set is the fact that such mining often generates a huge number of itemsets satisfying the minimum support threshold, especially when *minsup* is set low.
- This is because if an itemset is frequent, each of its subsets is frequent as well.
- A long itemset will contain a combinatorial number of shorter, frequent sub-itemsets.
- For example, a frequent itemset of length 100, such as  $\{a_1, a_2, \dots, a_{100}\}$ , contains  $\binom{100}{1} = 100$  frequent 1-itemsets:  $a_1, a_2, \dots, a_{100}$ ,  $\binom{100}{2}$  frequent 2-itemsets:  $(a_1, a_2), (a_1, a_3), \dots, (a_{99}, a_{100})$ , and so on. The total number of frequent itemsets that it contains is thus,

$$\binom{100}{1} + \binom{100}{2} + \dots + \binom{100}{100} = 2^{100} - 1 \approx 1.27 \times 10^{30} \quad (3)$$

# Frequent Closed Itemset Mining

## Definition 6 (Frequent Closed Itemset)

An itemset  $X$  is called frequent closed itemset if and only if it is frequent and there exists no proper superset  $X''$ , ( $X \subset X''$ ) such that support of  $X$  is same as the support of  $X''$ ,  $\text{sup}(X) = \text{sup}(X'')$ .

## Definition 7 (Frequent Maximal Itemset)

An itemset  $X$  is called frequent maximal itemset if and only if it is frequent and there exists no proper superset  $X''$ , ( $X \subset X''$ ).

- Suppose that a transaction database has only two transactions:  $\{\langle a_1, a_2, \dots, a_{100} \rangle; \langle a_1, a_2, \dots, a_{50} \rangle\}$ . Let the minimum support count threshold be  $\text{minsup} = 1$ .
- Two closed frequent itemsets and their support counts, that is,  $C = \{\{a_1, a_2, \dots, a_{100}\} : 1; \{a_1, a_2, \dots, a_{50}\} : 2\}$ .

# Frequent Closed Itemset Mining

- Suppose that a transaction database has only two transactions:  $\{\langle a_1, a_2, \dots, a_{100} \rangle; \langle a_1, a_2, \dots, a_{50} \rangle\}$ . Let the minimum support count threshold be  $\text{minsup} = 1$ .
- Two closed frequent itemsets and their support counts, that is,  $C = \{\{a_1, a_2, \dots, a_{100}\} : 1; \{a_1, a_2, \dots, a_{50}\} : 2\}$ .
- There is one maximal frequent itemset  $M = \{\{a_1, a_2, \dots, a_{100}\} : 1\}$ . (We cannot include  $\{a_1, a_2, \dots, a_{50}\}$  as a maximal frequent itemset because it has a frequent super-set,  $\{a_1, a_2, \dots, a_{100}\}$ .)

# Frequent Closed Itemset Mining from High Dimensional Dataset

- The conventional algorithms mine frequent itemsets, frequent closed itemset and frequent maximal itemset from the transactional datasets.
- In the modern era, the abundant data across variety of domains, including bioinformatics has led to the new form of dataset known as a high dimensional dataset, whose data characteristics are different from that of transactional datasets.
- The high dimensional datasets consist of less number of rows and considerably large number of features.
- The amount of information that can be extracted from high dimensional datasets is potentially huge, but extraction of information and knowledge from these datasets is a non-trivial task.

# Frequent Closed Itemset Mining from High Dimensional Dataset

- The conventional algorithms adopt feature enumeration based approach for mine frequent closed itemsets.
- The conventional algorithms face an uphill task in mining frequent closed itemsets from the high dimensional datasets.
- To overcome the inefficiency and uphill task of these algorithms, sequential row enumerated algorithms were proposed to mine FCI from high dimensional datasets.
- This problem of inefficiency can be solved to the greater extent by parallel row enumerated algorithms.

# Frequent Colossal Itemset Mining

- The result of frequent closed itemset mining algorithms includes small and mid-sized itemsets, which does not enclose valuable and complete information in many applications.
- In application dealing with high dimensional datasets such as bioinformatics, association rule mining gives greater importance to the large-sized itemsets called as colossal itemsets.

## Definition 8 (Frequent Colossal Itemset)

An itemset  $X$  is called frequent colossal itemset if and only if it is frequent and  $\text{card}(X) \geq \text{mincard}$ , where  $\text{mincard}$  is user specified least cardinality threshold.

## Example 6

In Table 1, the itemset  $X = \{f_1, f_2, f_6, f_{10}\}$  is frequent colossal itemset with minimum support threshold set to 2 and minimum cardinality threshold set to 4,  $\text{sup}(X) \geq 2$  and  $\text{card}(X) \geq 4$ .

# Frequent Colossal Closed Itemset Mining

## Definition 9 (Frequent Colossal Closed Itemset)

An itemset  $X$  is called frequent colossal closed itemset if and only if it is frequent closed and  $\text{card}(X) \geq \text{mincard}$ , where  $\text{mincard}$  is user specified least cardinality threshold.

## Example 7

In Table 1, the itemset  $X = \{f_2, f_4, f_7, f_8\}$ , is frequent colossal closed itemset with minimum support threshold set to 2 and minimum cardinality threshold set to 4,  $\text{sup}(X) \geq 2$  and  $\text{card}(X) \geq 4$ .

# Apriori Algorithm

- Apriori: A Candidate Generation-and-Test Approach
- Improving the Efficiency of Apriori
- FP-Growth: A Frequent Pattern-Growth Approach
- ECLAT: Frequent Pattern Mining with Vertical Data Format
- Mining Close Frequent Patterns and Maxpatterns
- The disadvantage of Apriori algorithm are:
  - It may need to generate a huge number of candidate sets.
  - It may need to repeatedly scan the database and check a large set of candidates by pattern matching.
  - Breadth-first search approach

# FP-growth Algorithm

- Frequent Pattern growth (FP-growth) Algorithm was proposed by J. Han, J. Pei, and Y. Yin.
- Mining frequent itemsets without candidate generation. FP-growth adopts depth-first search approach.
- The first scan of the database is the same as Apriori, which derives the set of frequent items (1-itemsets) and their support counts (frequencies).
- Let the minimum support count be 2. The set of frequent items is sorted in the order of descending support count.

# FP-growth Algorithm

- An FP-tree is then constructed as follows. First, create the root of the tree, labeled with “null.” Scan database D a second time.
- The items in each transaction are processed according to descending support count order and a branch is created for each transaction.

Table 4  
Dataset  $D$

TID	List of items
1	I1, I2, I5
2	I2, I4
3	I2, I3
4	I1, I2, I4
5	I1, I3
6	I2, I3
7	I1, I3
8	I1, I2, I3, I5
9	I1, I2, I3

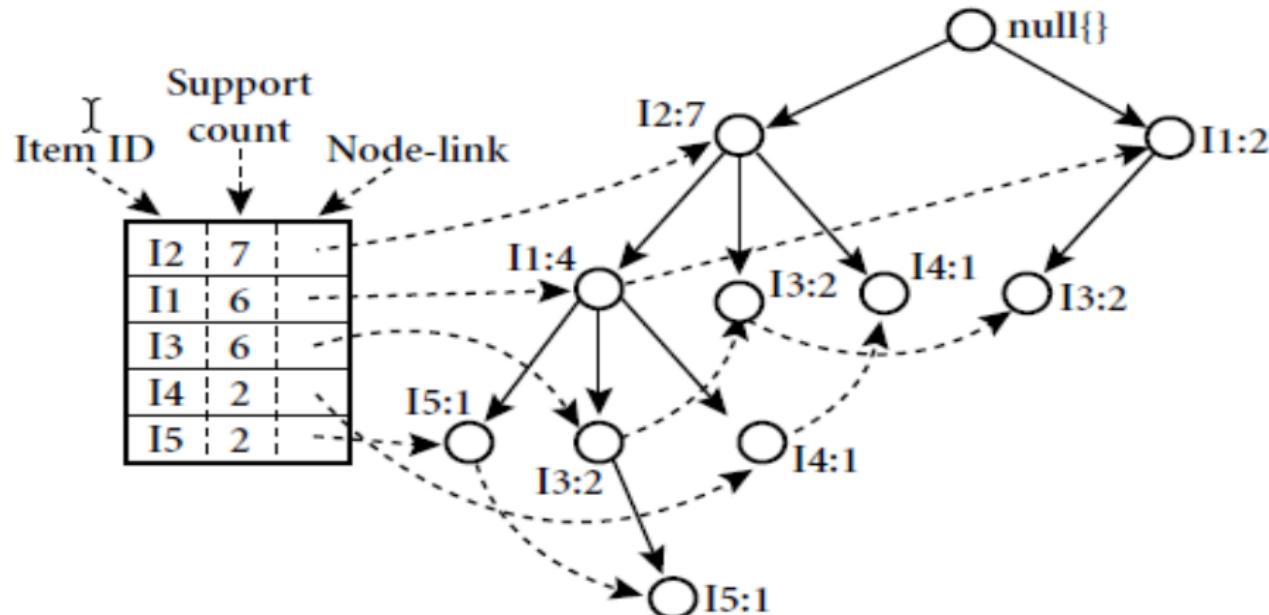
# FP-growth Algorithm

- The set of frequent items is sorted in the order of descending support count. Thus, we have  $\{\{I2: 7\}, \{I1: 6\}, \{I3: 6\}, \{I4: 2\}, \{I5: 2\}\}$ .
- The items in each transaction are processed according to descending support count order and a branch is created for each transaction.

Table 5  
Dataset  $D$

TID	List of items
1	I2, I1, I5
2	I2, I4
3	I2, I3
4	I2, I1, I4
5	I1, I3
6	I2, I3
7	I1, I3
8	I2, I1, I3, I5
9	I2, I1, I3

# FP-growth Algorithm



An FP-tree registers compressed, frequent pattern information.

# FP-growth Algorithm

- The FP-tree is mined as follows.
- Start from each frequent length-1 pattern (as an initial suffix pattern).
- Construct its conditional pattern base (a “subdatabase,” which consists of the set of prefix paths in the FP-tree co-occurring with the suffix pattern).
- Then construct its (conditional) FP-tree, and perform mining recursively on such a tree.
- The pattern growth is achieved by the concatenation of the suffix pattern with the frequent patterns generated from a conditional FP-tree.

# FP-growth Algorithm

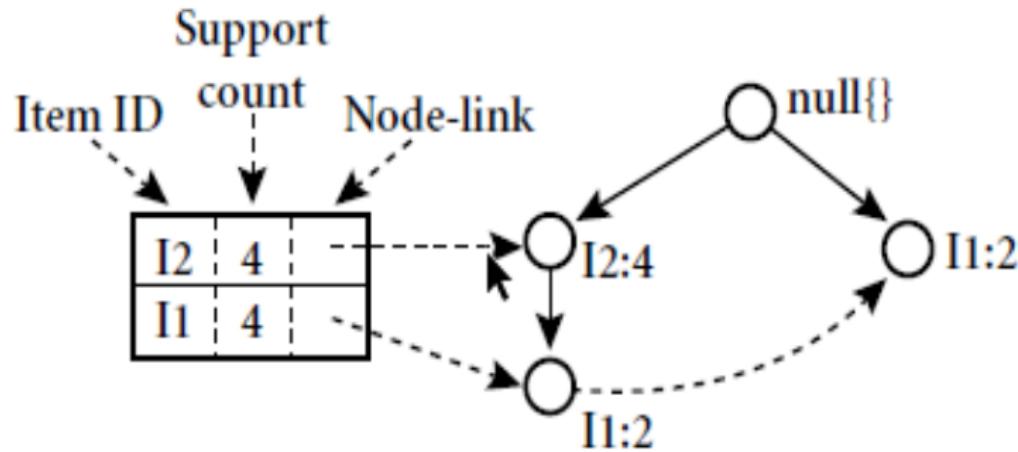
Mining the FP-tree by creating conditional (sub-)pattern bases.

Item	Conditional Pattern Base	Conditional FP-tree	Frequent Patterns Generated
I5	$\{\{I2, I1: 1\}, \{I2, I1, I3: 1\}\}$	$\langle I2: 2, I1: 2 \rangle$	$\{I2, I5: 2\}, \{I1, I5: 2\}, \{I2, I1, I5: 2\}$
I4	$\{\{I2, I1: 1\}, \{I2: 1\}\}$	$\langle I2: 2 \rangle$	$\{I2, I4: 2\}$
I3	$\{\{I2, I1: 2\}, \{I2: 2\}, \{I1: 2\}\}$	$\langle I2: 4, I1: 2 \rangle, \langle I1: 2 \rangle$	$\{I2, I3: 4\}, \{I1, I3: 4\}, \{I2, I1, I3: 2\}$
I1	$\{\{I2: 4\}\}$	$\langle I2: 4 \rangle$	$\{I2, I1: 4\}$

# FP-growth Algorithm

- Lets first consider I5. I5 occurs in two branches of the FP-tree.
- The paths formed by these branches are  $\langle I2, I1, I5: 1 \rangle$  and  $\langle I2, I1, I3, I5: 1 \rangle$ .
- Considering I5 as a suffix, its corresponding two prefix paths are  $\langle I2, I1: 1 \rangle$  and  $\langle I2, I1, I3: 1 \rangle$ , which form its conditional pattern base.
- Its conditional FP-tree contains only a single path,  $\langle I2: 2, I1: 2 \rangle$ ; I3 is not included because its support count of 1 is less than the minimum support count.
- The single path generates all the combinations of frequent patterns:  $\{I2, I5: 2\}$ ,  $\{I1, I5: 2\}$ ,  $\{I2, I1, I5: 2\}$ .

# FP-growth Algorithm



---

The conditional FP-tree associated with the conditional node I3.

# FP-growth Algorithm

- I3's conditional pattern base is  $\{\{I2, I1: 2\}, \{I2: 2\}, \{I1: 2\}\}$ .
- Its conditional FP-tree has two branches,  $\langle I2: 4, I1: 2 \rangle$  and  $\langle I1: 2 \rangle$
- Generates the set of patterns,  $\{\{I2, I3: 4\}, \{I1, I3: 4\}, \{I2, I1, I3: 2\}\}$ .

# **Data Warehousing and On-line Analytical Processing**

---

- Data Warehouse: Basic Concepts
- Data Warehouse Modeling: Data Cube and OLAP
- Data Warehouse Design and Usage
- Data Warehouse Implementation

# What is a Data Warehouse?

---

- Defined in many different ways, but not rigorously.
  - A decision support database that is maintained separately from the organization's operational database
  - Support information processing by providing a solid platform of consolidated, historical data for analysis.
- "A data warehouse is a subject-oriented, integrated, time-variant, and nonvolatile collection of data in support of management's decision-making process."—W. H. Inmon
- Data warehousing:
  - The process of constructing and using data warehouses

# Data Warehouse—Subject-Oriented

---

- Organized around major subjects, such as customer, product, sales
- Focusing on the modeling and analysis of data for decision makers, not on daily operations or transaction processing
- Provides a simple and concise view around particular subject issues by excluding data that are not useful in the decision support process

# Data Warehouse—Integrated

---

- Constructed by integrating multiple, heterogeneous data sources
  - relational databases, flat files, on-line transaction records
- Data cleaning and data integration techniques are applied.
  - Ensure consistency in naming conventions, encoding structures, attribute measures, etc. among different data sources
    - E.g., Hotel price: currency, tax, breakfast covered, etc.
  - When data is moved to the warehouse, it is converted.

# Data Warehouse—Nonvolatile

---

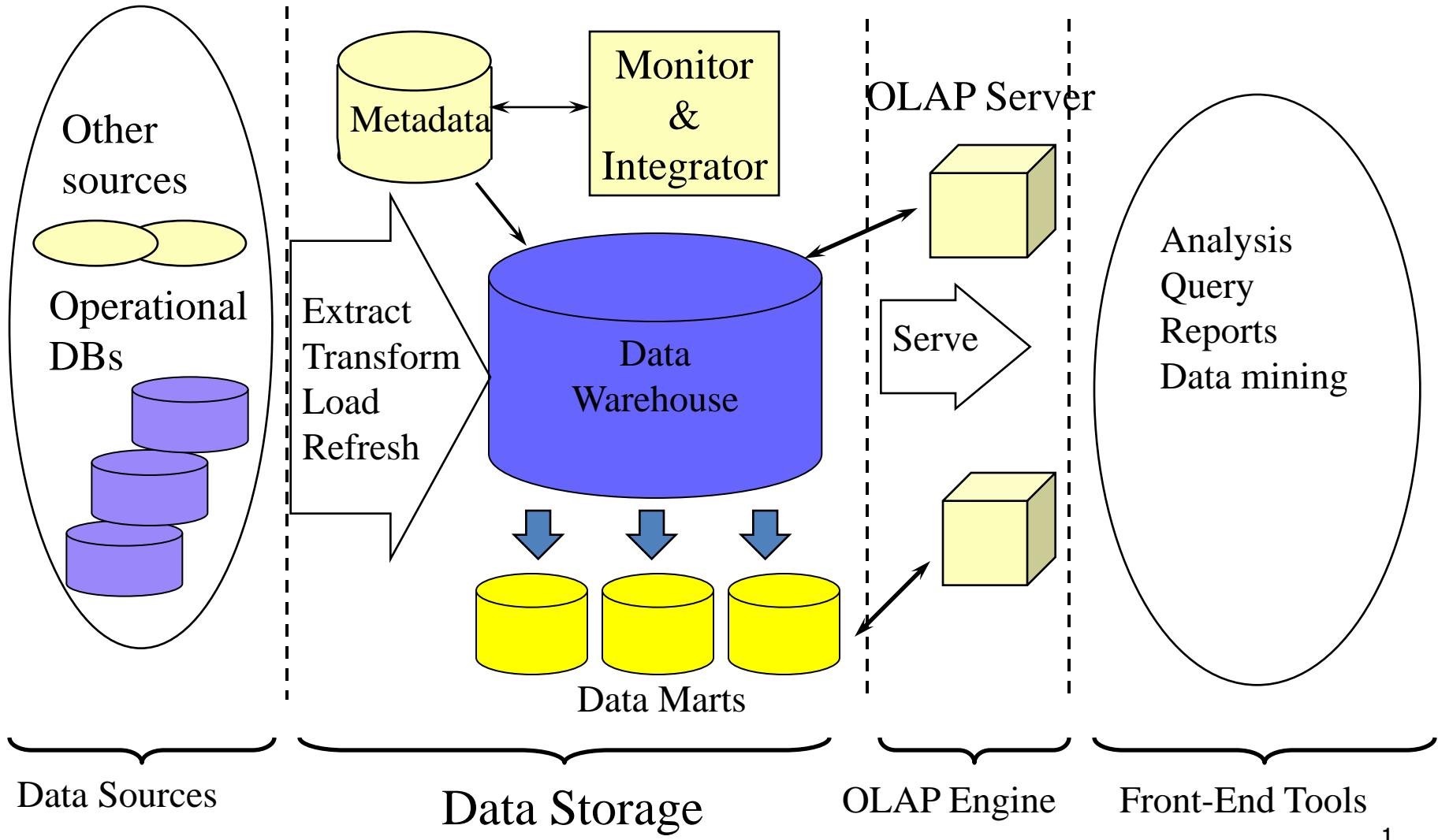
- A physically separate store of data transformed from the operational environment
- Operational update of data does not occur in the data warehouse environment
  - Does not require transaction processing, recovery, and concurrency control mechanisms
  - Requires only two operations in data accessing:
    - *initial loading of data* and *access of data*

# OLTP vs. OLAP

---

	<b>OLTP</b>	<b>OLAP</b>
<b>users</b>	clerk, IT professional	knowledge worker
<b>function</b>	day to day operations	decision support
<b>DB design</b>	application-oriented	subject-oriented
<b>data</b>	current, up-to-date detailed, flat relational isolated	historical, summarized, multidimensional integrated, consolidated
<b>usage</b>	repetitive	ad-hoc
<b>access</b>	read/write index/hash on prim. key	lots of scans
<b>unit of work</b>	short, simple transaction	complex query
<b># records accessed</b>	tens	millions
<b>#users</b>	thousands	hundreds
<b>DB size</b>	100MB-GB	100GB-TB
<b>metric</b>	transaction throughput	query throughput, response

# Data Warehouse: A Multi-Tiered Architecture



# Three Data Warehouse Models

- Enterprise warehouse
  - collects all of the information about subjects spanning the entire organization
- Data Mart
  - a subset of corporate-wide data that is of value to a specific groups of users. Its scope is confined to specific, selected groups, such as marketing data mart

- Virtual warehouse
  - A set of views over operational databases
  - Only some of the possible summary views may be materialized

# Extraction, Transformation, and Loading (ETL)

- **Data extraction**
  - get data from multiple, heterogeneous, and external sources
- **Data cleaning**
  - detect errors in the data and rectify them when possible
- **Data transformation**
  - convert data from legacy or host format to warehouse format
- **Load**
  - sort, summarize, consolidate, compute views, check integrity, and build indices and partitions
- **Refresh**
  - propagate the updates from the data sources to the warehouse

# Metadata Repository

- **Meta data** is the data defining warehouse objects. It stores:
- Description of the structure of the data warehouse
  - schema, view, dimensions, hierarchies, derived data defn, data mart locations and contents
- Operational meta-data
  - data lineage (history of migrated data and transformation path), currency of data (active, archived, or purged), monitoring information (warehouse usage statistics, error reports, audit trails)

- The algorithms used for summarization
- The mapping from operational environment to the data warehouse
- Data related to system performance
  - warehouse schema, view and derived data definitions
- Business data
  - business terms and definitions, ownership of data, charging policies

# Data Warehousing and On-line Analytical Processing

- Data Warehouse: Basic Concepts
- Data Warehouse Modeling: Data Cube and OLAP
- Data Warehouse Design and Usage

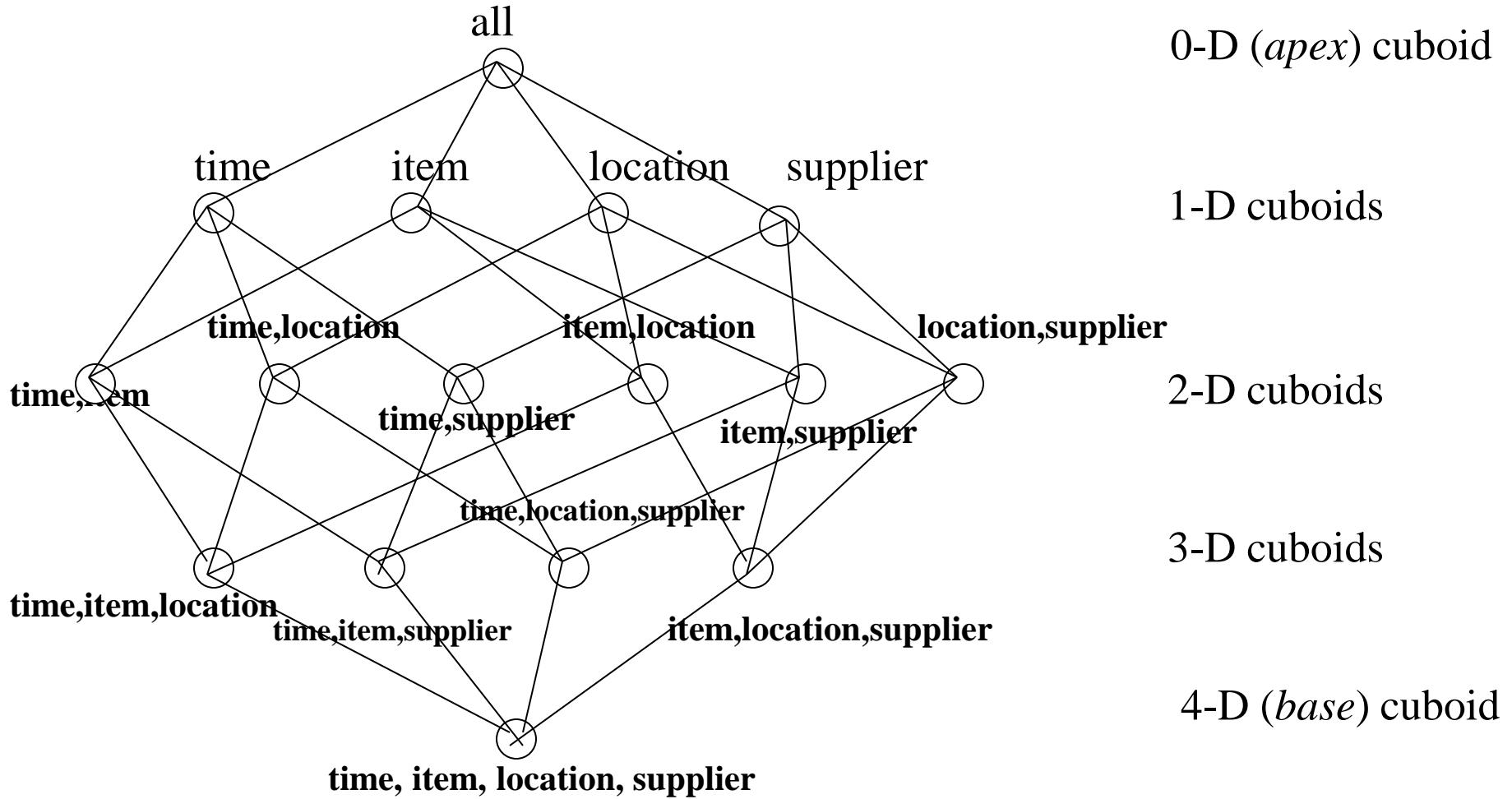


# From Tables and Spreadsheets to Data Cubes

- A **data warehouse** is based on a multidimensional data model which views data in the form of a data cube
- A data cube, such as sales, allows data to be modeled and viewed in multiple dimensions
  - **Dimension tables**, such as item (item\_name, brand, type), or time(day, week, month, quarter, year)

- **Fact table** contains **measures** (such as dollars\_sold) and keys to each of the related dimension tables
- In data warehousing literature, an n-D base cube is called a base cuboid. The top most 0-D cuboid, which holds the highest-level of summarization, is called the apex cuboid. The lattice of cuboids forms a data cube.

# Cube: A Lattice of Cuboids

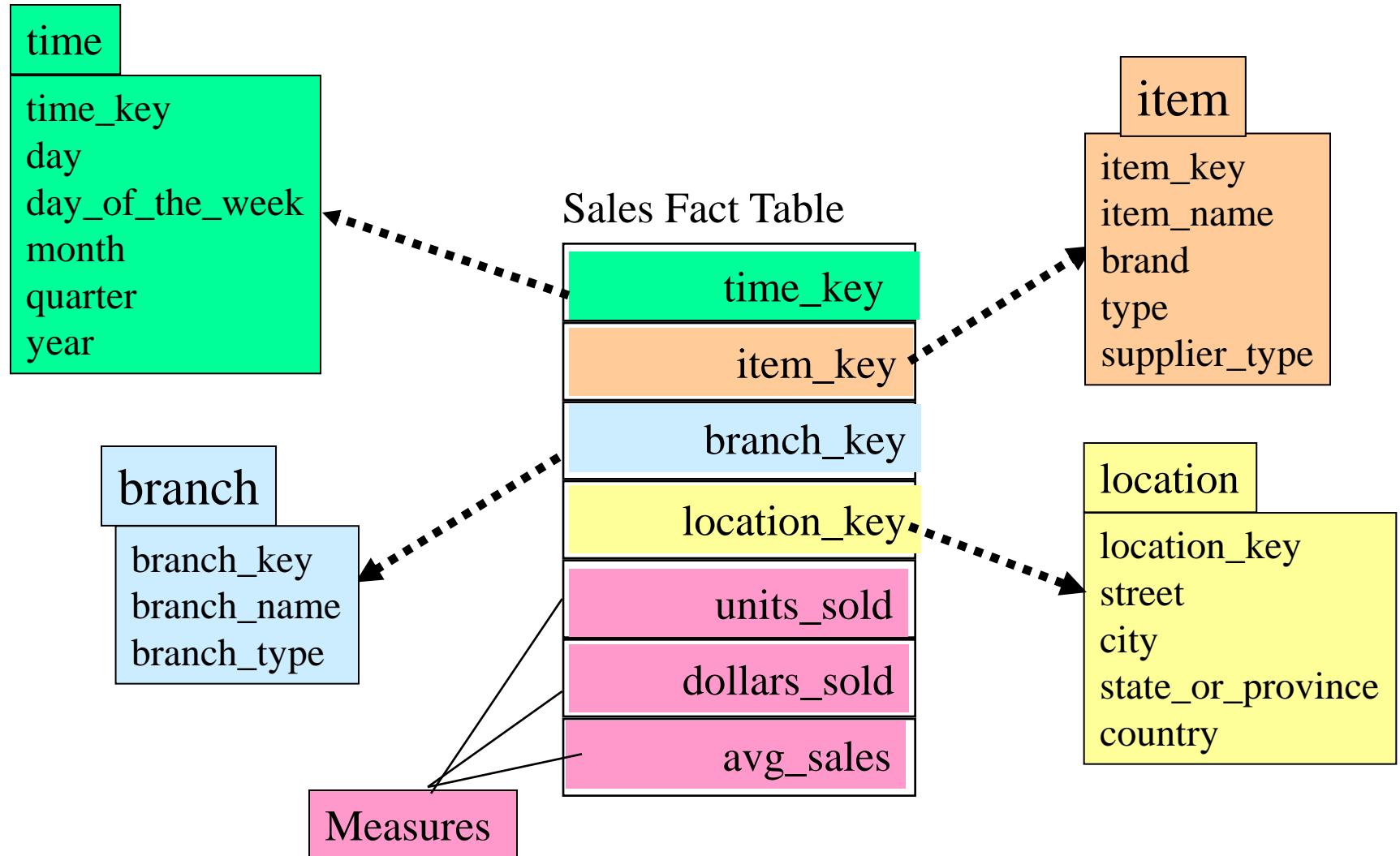


# Conceptual Modeling of Data Warehouses

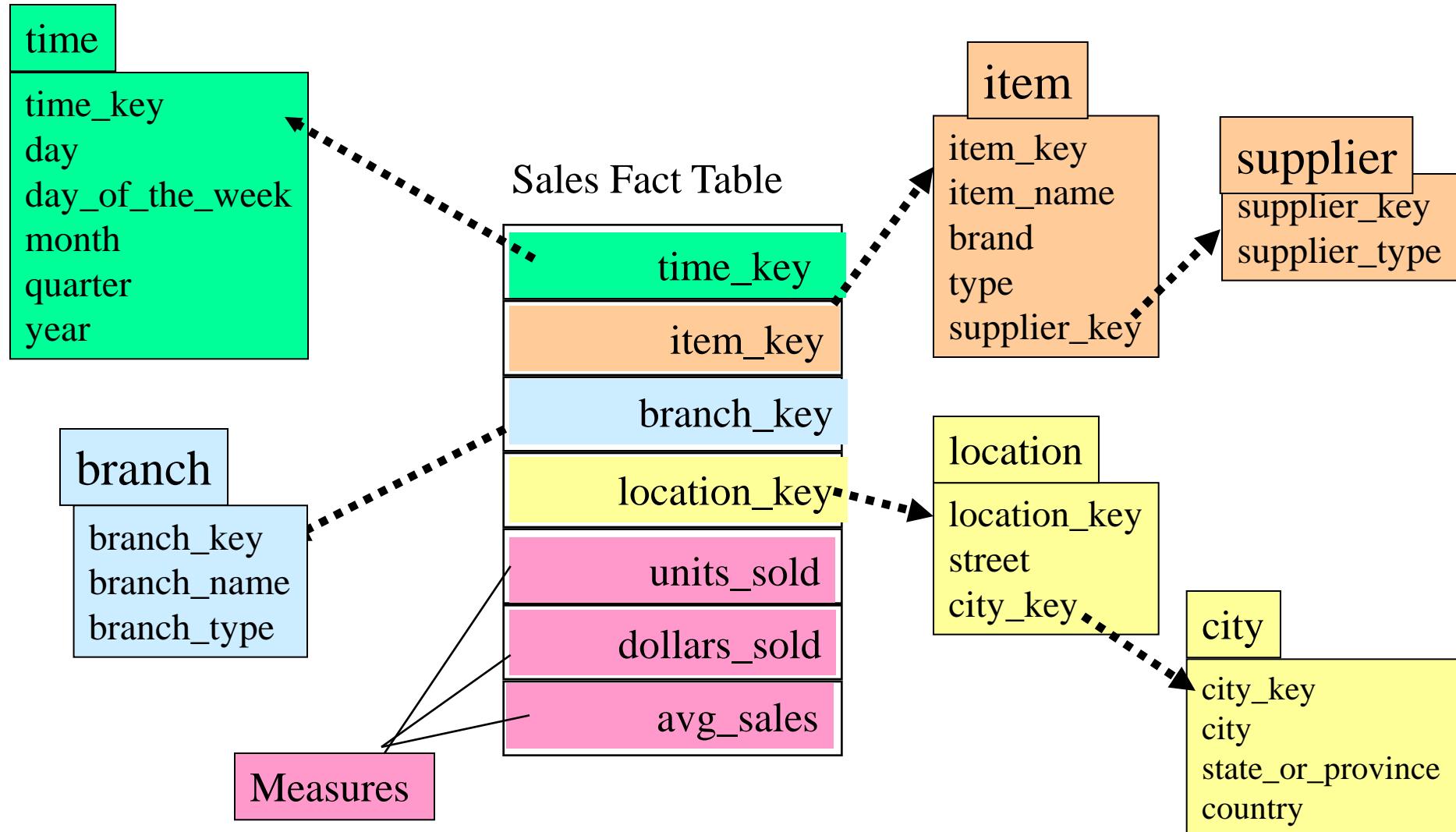
- Modeling data warehouses: dimensions & measures
  - Star schema: A fact table in the middle connected to a set of dimension tables
  - Snowflake schema: A refinement of star schema where some dimensional hierarchy is normalized into a set of smaller dimension tables, forming a shape similar to snowflake

- Fact constellations: Multiple fact tables share dimension tables, viewed as a collection of stars, therefore called galaxy schema or fact constellation

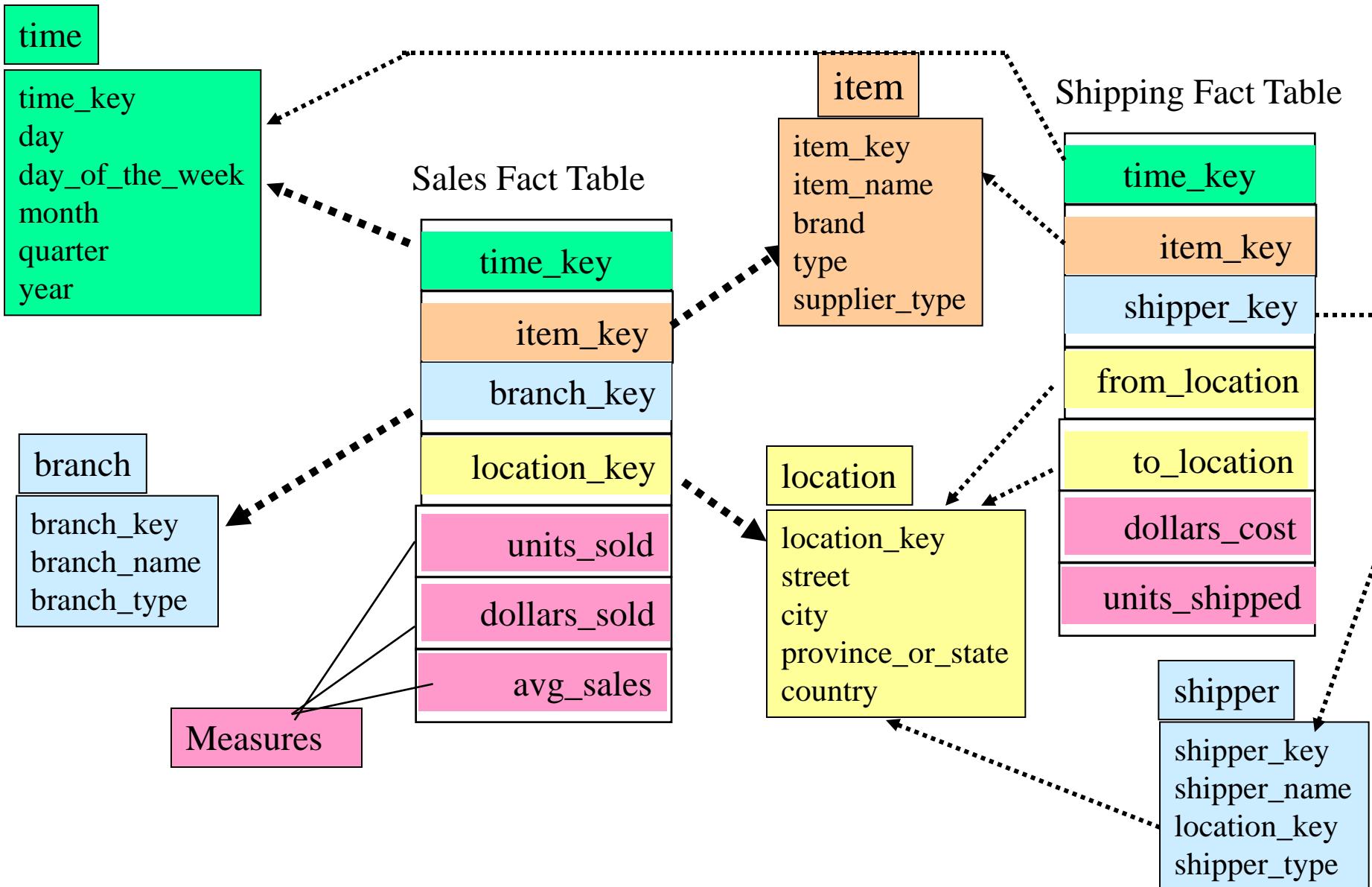
# Example of Star Schema



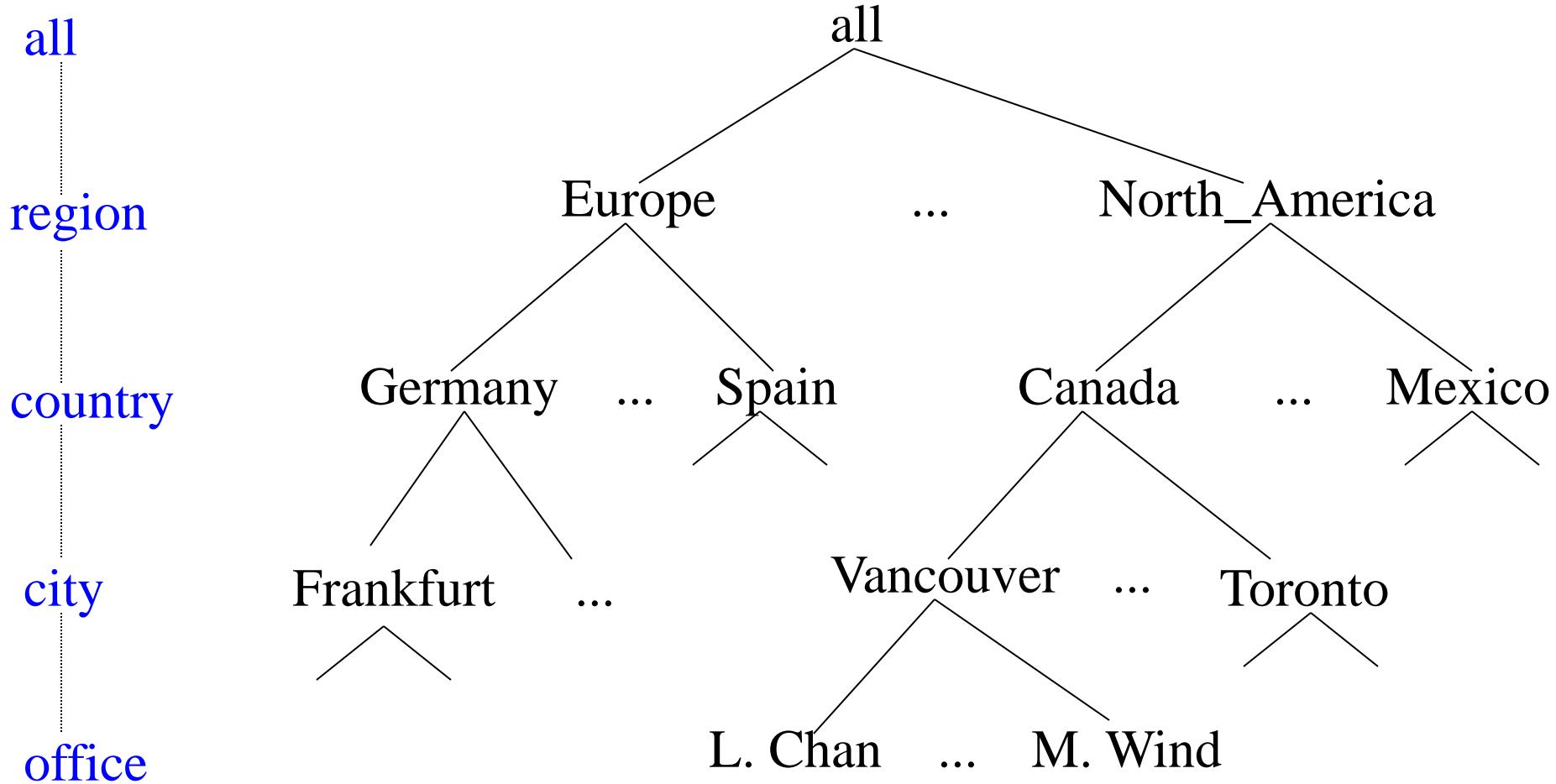
# Example of Snowflake Schema



# Example of Fact Constellation



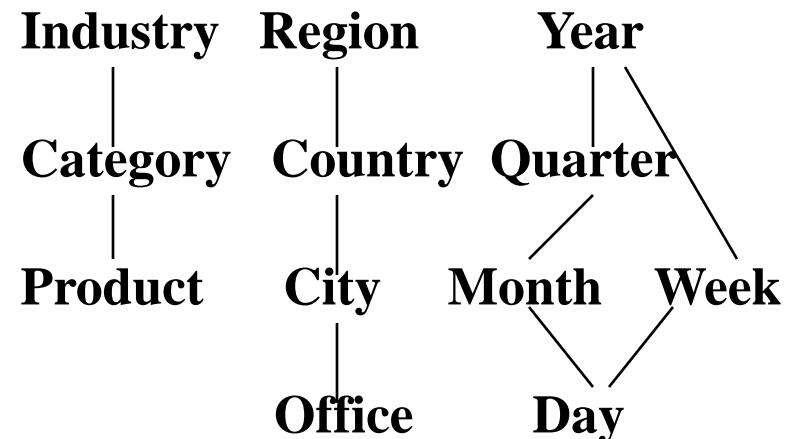
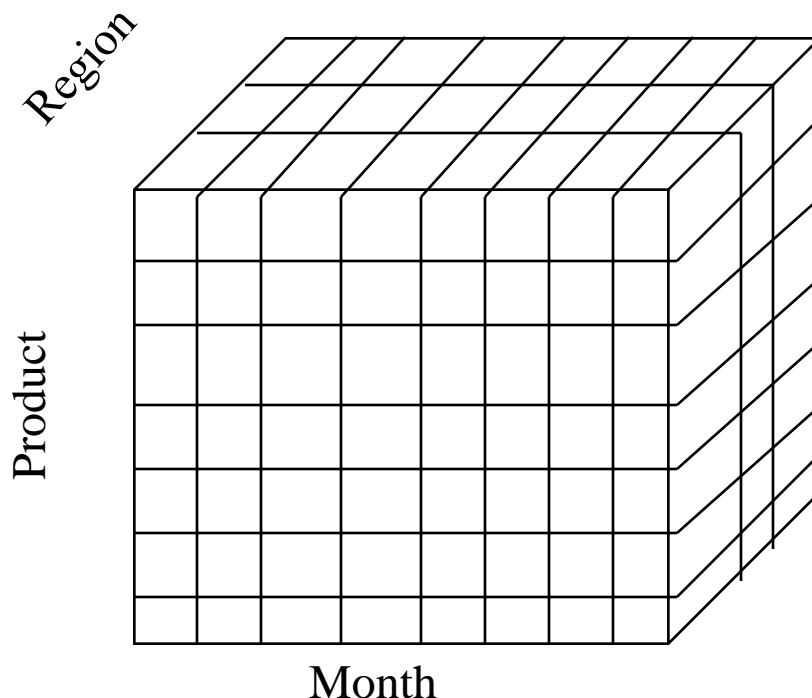
# A Concept Hierarchy: **Dimension** (location)



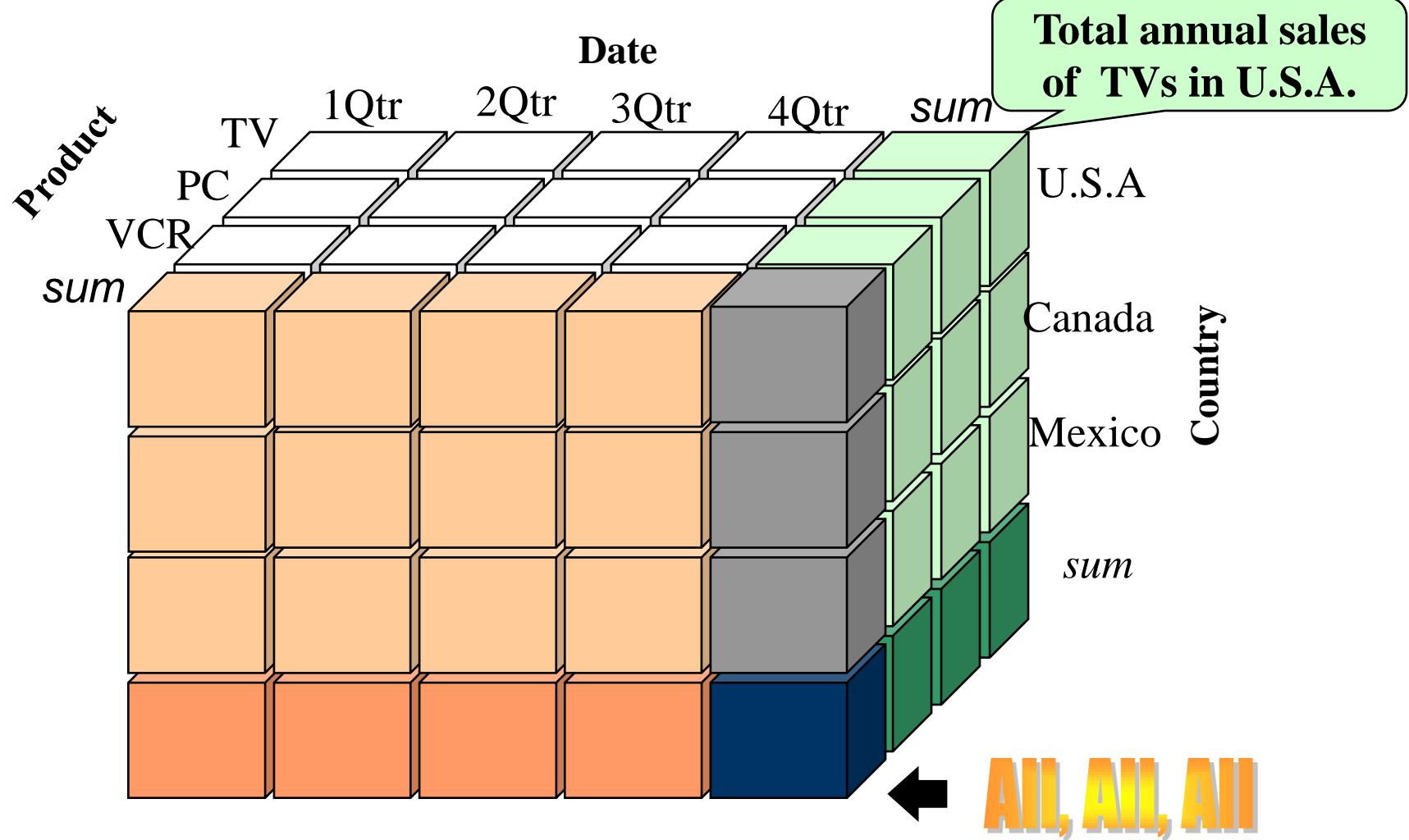
# Multidimensional Data

- Sales volume as a function of product, month, and region

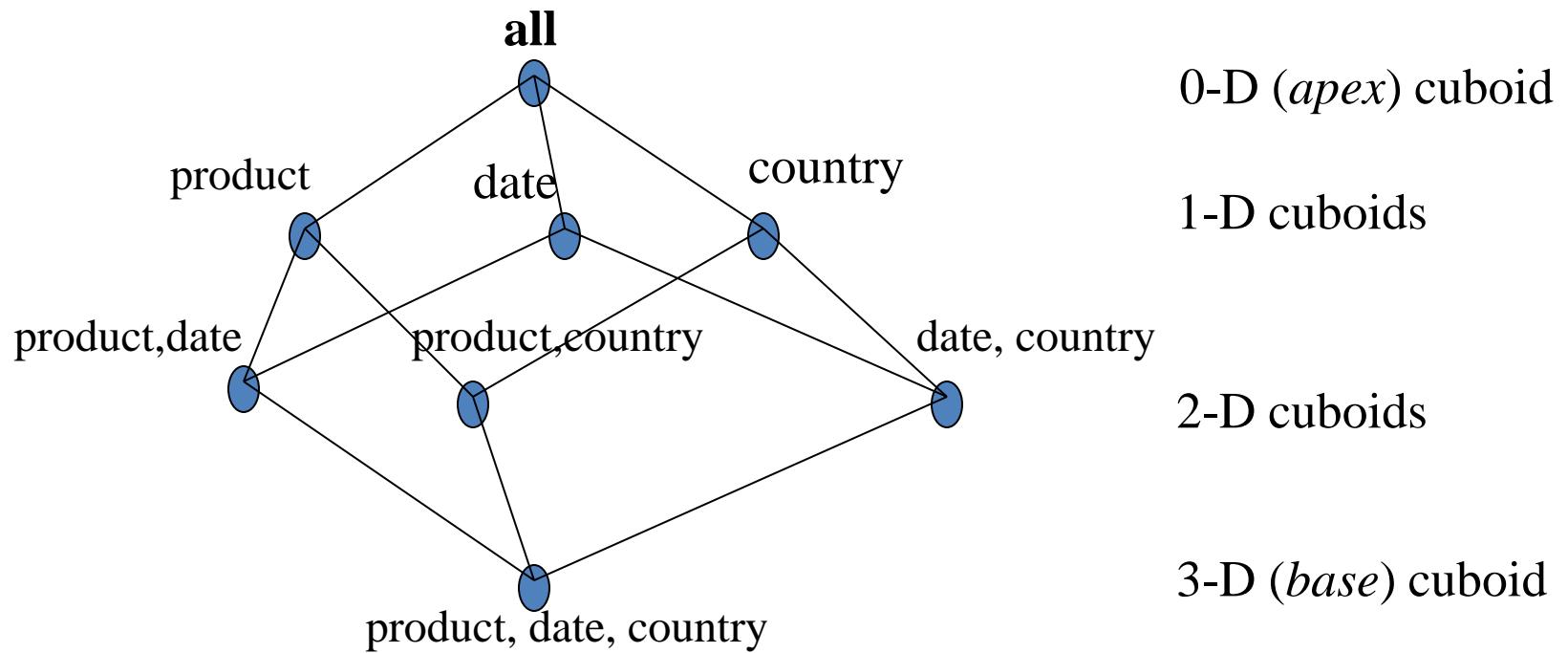
Dimensions: *Product, Location, Time*  
Hierarchical summarization paths



# A Sample Data Cube



# Cuboids Corresponding to the Cube



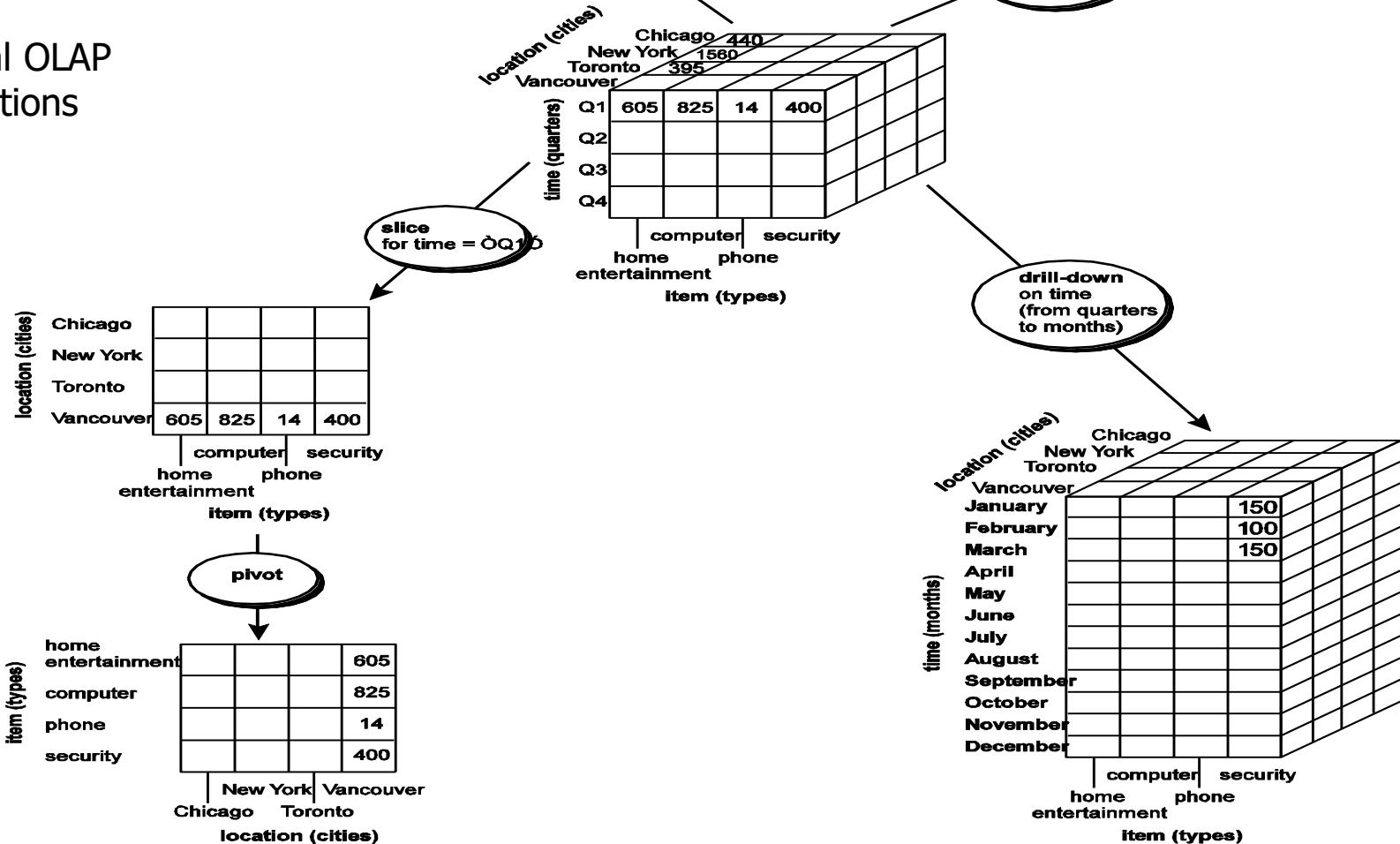
# Typical OLAP Operations

- Roll up (drill-up): summarize data
  - *by climbing up hierarchy or by dimension reduction*
- Drill down (roll down): reverse of roll-up
  - *from higher level summary to lower level summary or detailed data, or introducing new dimensions*
- Slice and dice: *project and select*

- Pivot (rotate):
  - *reorient the cube, visualization, 3D to series of 2D planes*

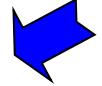


## Typical OLAP Operations



# Data Warehousing and On-line Analytical Processing

---

- Data Warehouse: Basic Concepts
- Data Warehouse Modeling: Data Cube and OLAP
- Data Warehouse Design and Usage 
- Data Warehouse Implementation
- Data Generalization by Attribute-Oriented Induction
- Summary

# **Design of Data Warehouse: A Business Analysis Framework**

---

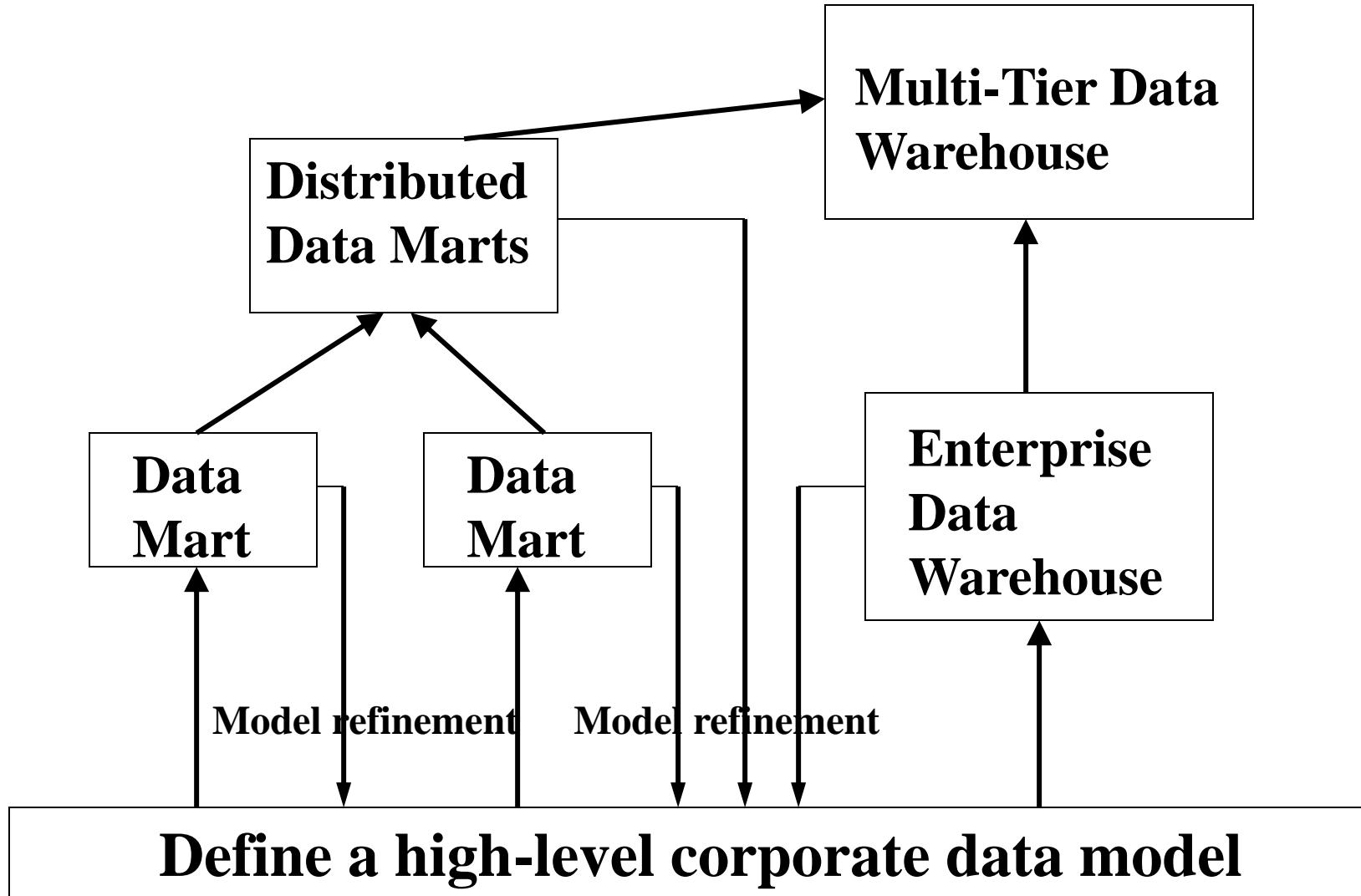
- Four views regarding the design of a data warehouse
  - **Top-down view**
    - allows selection of the relevant information necessary for the data warehouse
  - **Data source view**
    - exposes the information being captured, stored, and managed by operational systems
  - **Data warehouse view**
    - consists of fact tables and dimension tables
  - **Business query view**
    - sees the perspectives of data in the warehouse from the view of end-user

# Data Warehouse Design Process

---

- **Top-down, bottom-up approaches or a combination** of both
  - Top-down: Starts with overall design and planning (mature)
  - Bottom-up: Starts with experiments and prototypes (rapid)
- **Typical data warehouse design process**
  - Choose a **business process** to model, e.g., orders, invoices, etc.
  - Choose the *grain (atomic level of data)* of the business process
  - Choose the **dimensions** that will apply to each fact table record
  - Choose the **measure** that will populate each fact table record

# Data Warehouse Development: A Recommended Approach



# Data Warehouse Usage

---

- Three kinds of data warehouse applications
  - **Information processing**
    - supports querying, basic statistical analysis, and reporting using crosstabs, tables, charts and graphs
  - **Analytical processing**
    - multidimensional analysis of data warehouse data
    - supports basic OLAP operations, slice-dice, drilling, pivoting
  - **Data mining**
    - knowledge discovery from hidden patterns
    - supports associations, constructing analytical models, performing classification and prediction, and presenting the mining results using visualization tools

# **From On-Line Analytical Processing (OLAP) to On Line Analytical Mining (OLAM)**

---

- Why online analytical mining?
  - High quality of data in data warehouses
    - DW contains integrated, consistent, cleaned data
  - Available information processing structure surrounding data warehouses
    - Web accessing, service facilities, reporting and OLAP tools
  - OLAP-based exploratory data analysis
    - Mining with drilling, dicing, pivoting, etc.
  - On-line selection of data mining functions
    - Integration and swapping of multiple mining functions, algorithms, and tasks

# Data Warehousing and On-line Analytical Processing

---

- Data Warehouse: Basic Concepts
- Data Warehouse Modeling: Data Cube and OLAP
- Data Warehouse Design and Usage
- Data Warehouse Implementation 
- Data Generalization by Attribute-Oriented Induction
- Summary

# Efficient Data Cube Computation

---

- Data cube can be viewed as a lattice of cuboids
  - The bottom-most cuboid is the base cuboid
  - The top-most cuboid (apex) contains only one cell
  - How many cuboids in an n-dimensional cube with L levels?
- Materialization of data cube
  - Materialize every (cuboid) (**full materialization**), none (**no materialization**), or some (**partial materialization**)
  - Selection of which cuboids to materialize
    - Based on size, sharing, access frequency, etc.

# The “Compute Cube” Operator

- Cube definition and computation in DMQL

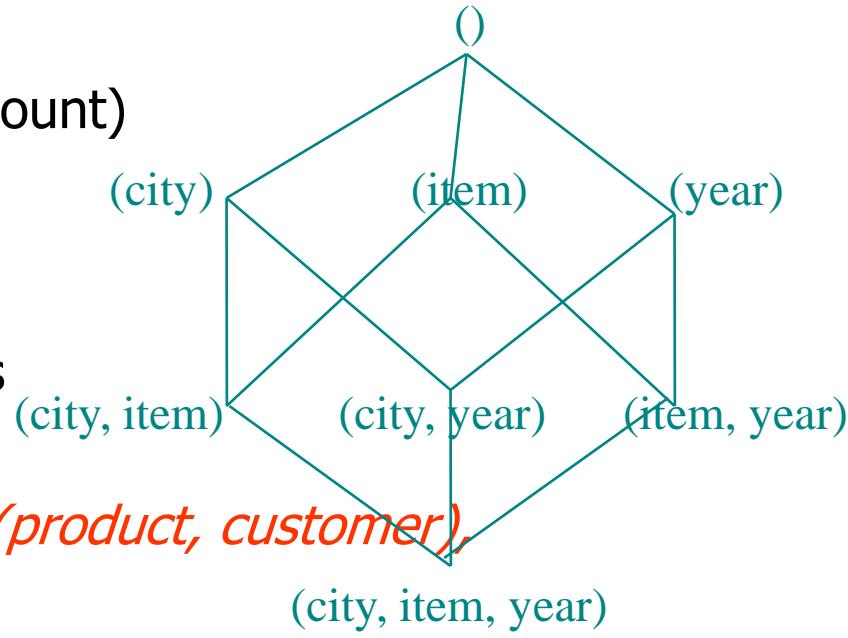
```
define cube sales [item, city, year]: sum (sales_in_dollars)  
compute cube sales
```

- Transform it into a SQL-like language (with a new operator **cube by**, introduced by Gray et al.'96)

```
SELECT item, city, year, SUM (amount)  
FROM SALES  
CUBE BY item, city, year
```

- Need compute the following Group-Bys

*(date, product, customer),  
(date,product),(date, customer), (product, customer),  
(date), (product), (customer)  
()*



# OLAP Server Architectures

---

- Relational OLAP (ROLAP)
  - Use relational or extended-relational DBMS to store and manage warehouse data and OLAP middle ware
  - Include optimization of DBMS backend, implementation of aggregation navigation logic, and additional tools and services
  - Greater scalability
- Multidimensional OLAP (MOLAP)
  - Sparse array-based multidimensional storage engine
  - Fast indexing to pre-computed summarized data
- Hybrid OLAP (HOLAP) (e.g., Microsoft SQLServer)
  - Flexibility, e.g., low level: relational, high-level: array
- Specialized SQL servers (e.g., Redbricks)
  - Specialized support for SQL queries over star/snowflake schemas

# Mining Multi-Dimensional Association

---

- Single-dimensional rules:

buys(X, "milk")  $\Rightarrow$  buys(X, "bread")

- Multi-dimensional rules:  $\geq 2$  dimensions or predicates

- Inter-dimension assoc. rules (*no repeated predicates*)

age(X,"19-25")  $\wedge$  occupation(X,"student")  $\Rightarrow$  buys(X, "coke")

- hybrid-dimension assoc. rules (*repeated predicates*)

age(X,"19-25")  $\wedge$  buys(X, "popcorn")  $\Rightarrow$  buys(X, "coke")

# Constraint-based (Query-Directed) Mining

---

- Finding **all** the patterns in a database **autonomously**? — unrealistic!
  - The patterns could be too many but not focused!
- Data mining should be an **interactive** process
  - User directs what to be mined using a **data mining query language** (or a graphical user interface)
- Constraint-based mining
  - User flexibility: provides **constraints** on what to be mined

# Constraints in Data Mining

---

- Knowledge type constraint:
  - classification, association, etc.
- Data constraint
  - find product pairs sold together in stores in Chicago this year
- Dimension/level constraint
  - in relevance to region, price, brand, customer category
- Rule (or pattern) constraint
  - small sales (price < \$10) triggers big sales (sum > \$200)
- Interestingness constraint
  - strong rules:  $\text{min\_support} \geq 3\%$ ,  $\text{min\_confidence} \geq 60\%$

# Meta-Rule Guided Mining

---

- Meta-rule can be in the rule form with partially instantiated predicates and constants

$$P_1(X, Y) \wedge P_2(X, W) \Rightarrow \text{buys}(X, \text{"iPad"})$$

- The resulting rule derived can be

$$\text{age}(X, \text{"15-25"}) \wedge \text{profession}(X, \text{"student"}) \Rightarrow \text{buys}(X, \text{"iPad"})$$

- In general, it can be in the form of

$$P_1 \wedge P_2 \wedge \dots \wedge P_l \Rightarrow Q_1 \wedge Q_2 \wedge \dots \wedge Q_r$$

# Challenges

---

- A major challenge in mining frequent itemsets from a large data set is the fact that such mining often generates a huge number of itemsets satisfying the minimum support threshold, especially when *minsup* is set low.
- This is because if an itemset is frequent, each of its subsets is frequent as well.
- A long itemset will contain a combinatorial number of shorter, frequent sub-itemsets.

# Closed Patterns and Max-Patterns

---

- Example : A long pattern contains a combinatorial number of sub-patterns, e.g.,  $\{a_1, \dots, a_{100}\}$  contains  $\binom{100}{1} + \binom{100}{2} + \dots + \binom{100^0}{1^00^0} = 2^{100} - 1 = 1.27*10^{30}$  sub-patterns!
- Solution: Mine ***closed patterns*** and ***max-patterns instead***
- An itemset X is ***closed*** if X is *frequent* and there exists *no super-pattern Y ⊃ X, with the same support as X*
- An itemset X is a ***max-pattern*** if X is frequent and there exists no frequent super-pattern Y ⊃ X
- Closed pattern is a lossless compression of freq. patterns

# Closed Patterns and Max-Patterns

---

- Exercise. DB = { $\langle a_1, \dots, a_{100} \rangle$ ,  $\langle a_1, \dots, a_{50} \rangle$ }
  - Min\_sup = 1.
- What is the set of **closed itemset**?
  - $\langle a_1, \dots, a_{100} \rangle$ : 1
  - $\langle a_1, \dots, a_{50} \rangle$ : 2
- What is the set of **max-pattern**?
  - $\langle a_1, \dots, a_{100} \rangle$ : 1

# Colossal itemset

---

- The result of frequent closed itemset mining algorithms includes small and mid-sized itemsets, which does not enclose valuable and complete information in many applications.
- In application dealing with high dimensional datasets such as bioinformatics (Micro array analysis, biological sequence analysis) , association rule mining gives greater importance to the large sized itemsets called as colossal Itemsets
- An itemset X is called frequent colossal closed itemset if and only if it is frequent closed and  $\text{card}(X) \geq \text{mincard}$ , where mincard is user specified least cardinality threshold

Table 1

---

Tid	features
1	$f_1, f_2, f_4, f_6, f_{10}$
2	$f_1, f_2, f_4, f_7, f_8$
3	$f_2, f_4, f_7, f_8$
4	$f_1, f_2, f_6, f_8, f_9, f_{10}$
5	$f_1, f_3, f_4, f_7, f_8, f_{10}$
6	$f_2, f_4, f_9$
7	$f_5, f_7$
8	$f_5, f_{11}$

- 
- In Table 1, the itemset  $X = \{f_2, f_4, f_7, f_8\}$ , is frequent colossal closed itemset with minimum support threshold set to 2 and minimum cardinality threshold set to 4,  $sup(X) \geq 2$  and  $card(X) \geq 4$ .

# Colossal Patterns: A Motivating Example

---

Let's make a set of 40 transactions

**T<sub>1</sub> = 1 2 3 4 ..... 39 40**

**T<sub>2</sub> = 1 2 3 4 ..... 39 40**

: .

: .

: .

: .

**T<sub>40</sub>=1 2 3 4 ..... 39 40**

Then delete the items on the diagonal

**T<sub>1</sub> = 2 3 4 ..... 39 40**

**T<sub>2</sub> = 1 3 4 ..... 39 40**

: .

: .

: .

: .

**T<sub>40</sub>=1 2 3 4 ..... 39**

# A Show of Colossal Pattern Mining!

---

**T<sub>1</sub> = 2 3 4 ..... 39 40**

**T<sub>2</sub> = 1 3 4 ..... 39 40**

: .

: .

: .

: .

**T<sub>40</sub>=1 2 3 4 ..... 39**

**T<sub>41</sub>= 41 42 43 ..... 79**

**T<sub>42</sub>= 41 42 43 ..... 79**

: .

: .

**T<sub>60</sub>= 41 42 43 ... 79**

Let the min-support threshold  $\sigma= 20$

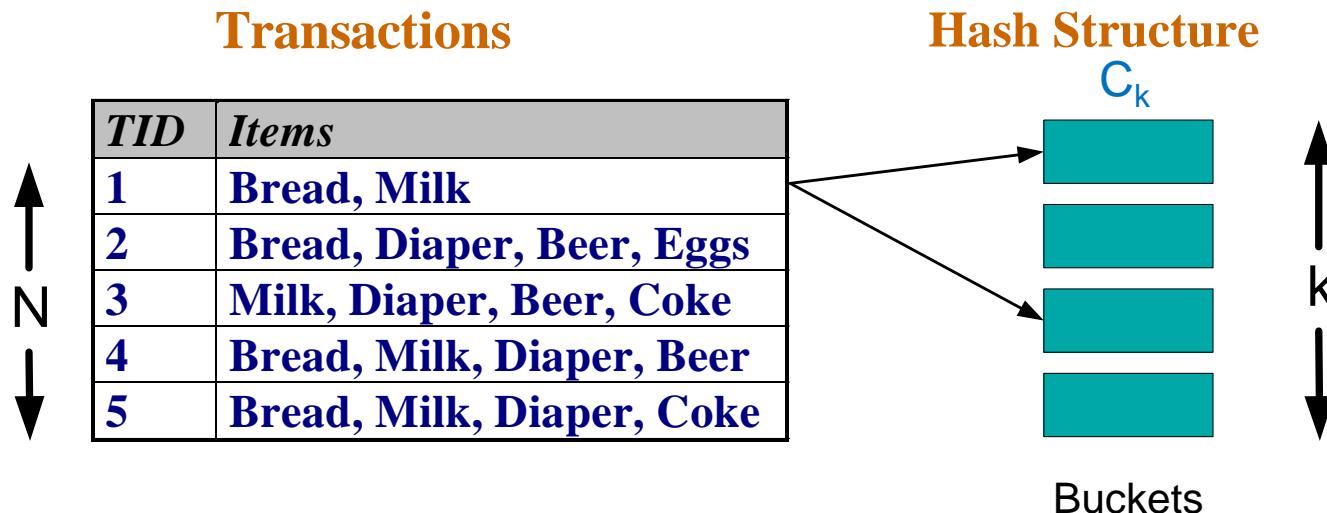
Then there are  $\binom{40}{20}$  closed/maximal frequent patterns of size 20

However, there is only one colossal pattern with size greater than 20,

$\alpha= \{41,42,\dots,79\}$  of size 39

# Computing Frequent Itemsets

- Given the set of **candidate** itemsets  $C_k$ , we need to compute the support and find the **frequent** itemsets  $L_k$ .
- Scan the data, and use a **hash structure** to keep a counter for each candidate itemset that appears in the data



# A simple hash structure

- Create a dictionary (hash table) that stores the candidate itemsets as keys, and the number of appearances as the value.
- Increment the counter for each itemset that you see in the

# Example

Suppose you have 15 candidate itemsets of length 3:

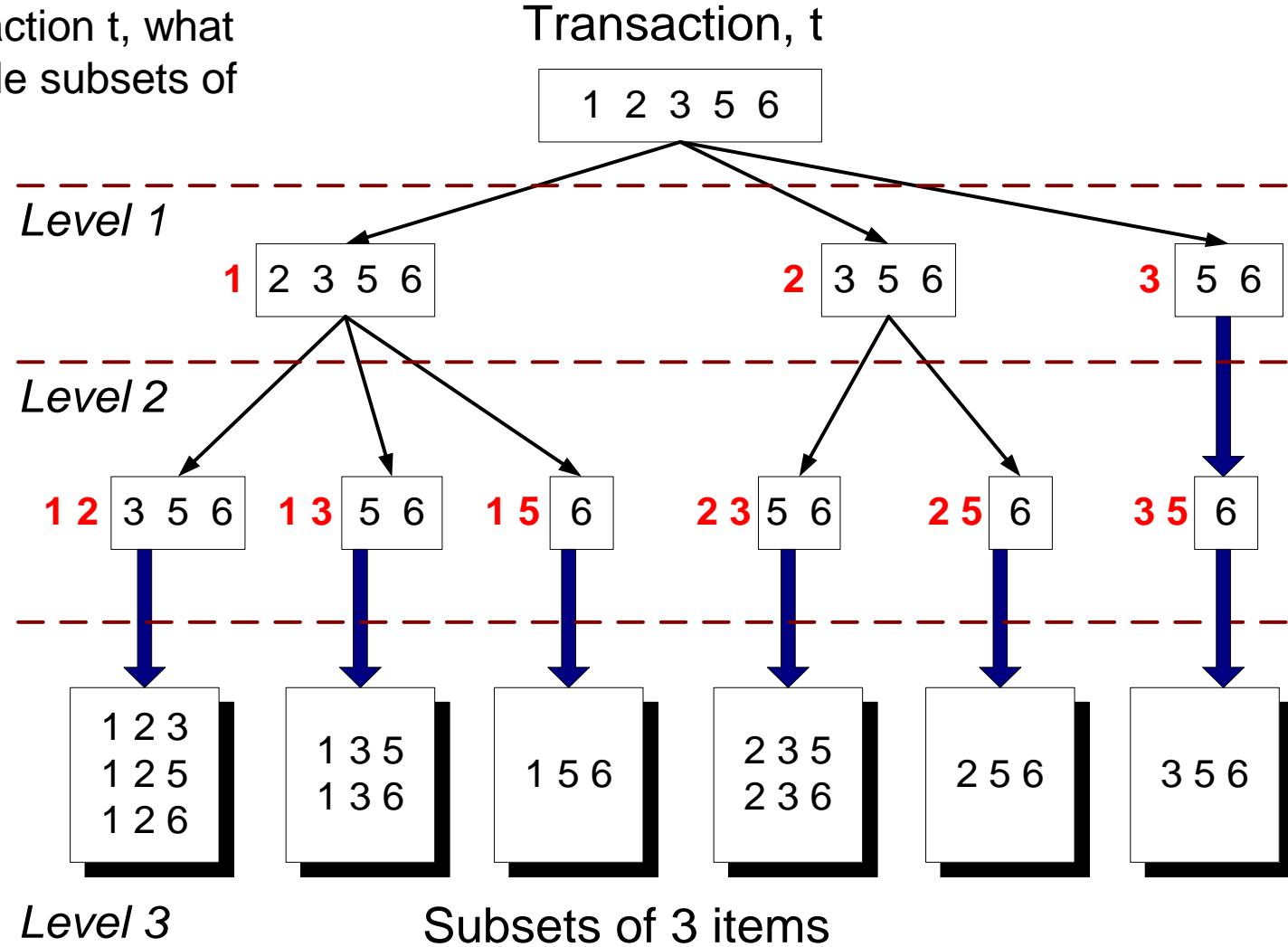
$\{1\ 4\ 5\}$ ,  $\{1\ 2\ 4\}$ ,  $\{4\ 5\ 7\}$ ,  $\{1\ 2\ 5\}$ ,  $\{4\ 5\ 8\}$ ,  
 $\{1\ 5\ 9\}$ ,  $\{1\ 3\ 6\}$ ,  $\{2\ 3\ 4\}$ ,  $\{5\ 6\ 7\}$ ,  $\{3\ 4\ 5\}$ ,  
 $\{3\ 5\ 6\}$ ,  $\{3\ 5\ 7\}$ ,  $\{6\ 8\ 9\}$ ,  $\{3\ 6\ 7\}$ ,  $\{3\ 6\ 8\}$

Hash table stores the counts of the candidate itemsets as they have been computed so far

Key	Value
$\{3\ 6\ 7\}$	0
$\{3\ 4\ 5\}$	1
$\{1\ 3\ 6\}$	3
$\{1\ 4\ 5\}$	5
$\{2\ 3\ 4\}$	2
$\{1\ 5\ 9\}$	1
$\{3\ 6\ 8\}$	0
$\{4\ 5\ 7\}$	2
$\{6\ 8\ 9\}$	0
$\{5\ 6\ 7\}$	3
$\{1\ 2\ 4\}$	8
$\{3\ 5\ 7\}$	1
$\{1\ 2\ 5\}$	0
$\{3\ 5\ 6\}$	1
$\{4\ 5\ 8\}$	0

# Subset Generation

Given a transaction t, what are the possible subsets of size 3?



Recursion!

# Example

Tuple {1,2,3,5,6} generates the following itemsets of length 3:

{1 2 3}, {1 2 5}, {1 2 6}, {1 3 5}, {1 3 6},  
 {1 5 6}, {2 3 5}, {2 3 6}, {3 5 6},

Increment the counters for the itemsets in the dictionary

Key	Value
{3 6 7}	0
{3 4 5}	1
{1 3 6}	3
{1 4 5}	5
{2 3 4}	2
{1 5 9}	1
{3 6 8}	0
{4 5 7}	2
{6 8 9}	0
{5 6 7}	3
{1 2 4}	8
{3 5 7}	1
{1 2 5}	0
{3 5 6}	1
{4 5 8}	0

# Example

Tuple {1,2,3,5,6} generates the following itemsets of length 3:

{1 2 3}, {1 2 5}, {1 2 6}, {1 3 5}, {1 3 6},  
{1 5 6}, {2 3 5}, {2 3 6}, {3 5 6},

Increment the counters for the itemsets in the dictionary

Key	Value
{3 6 7}	0
{3 4 5}	1
{1 3 6}	4
{1 4 5}	5
{2 3 4}	2
{1 5 9}	1
{3 6 8}	0
{4 5 7}	2
{6 8 9}	0
{5 6 7}	3
{1 2 4}	8
{3 5 7}	1
{1 2 5}	1
{3 5 6}	2
{4 5 8}	0

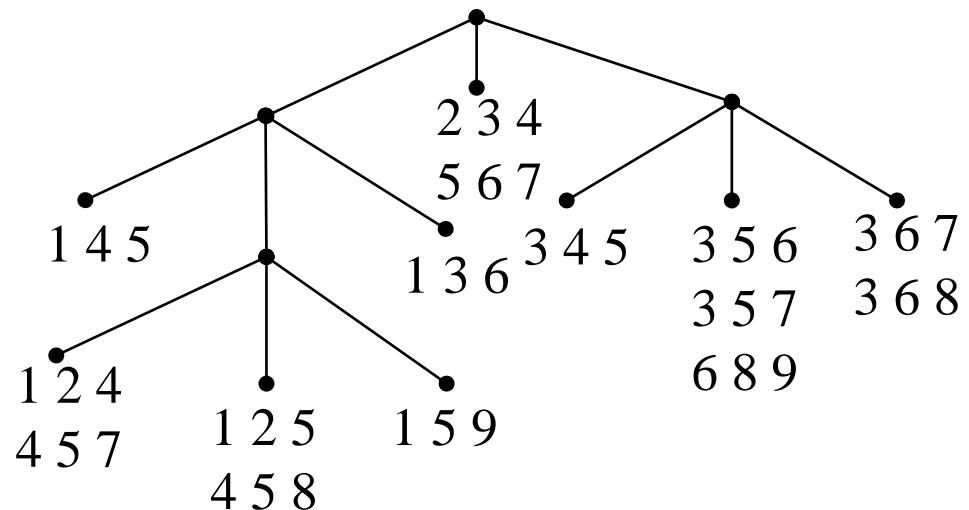
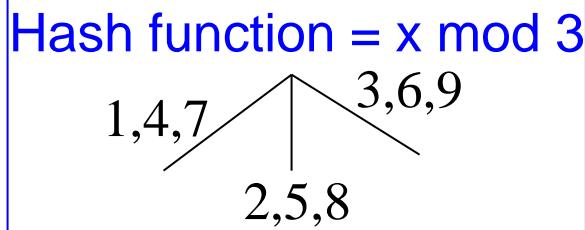
# The Hash Tree Structure

Suppose you have the same 15 candidate itemsets of length 3:

{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4},  
{5 6 7}, {3 4 5}, {3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}

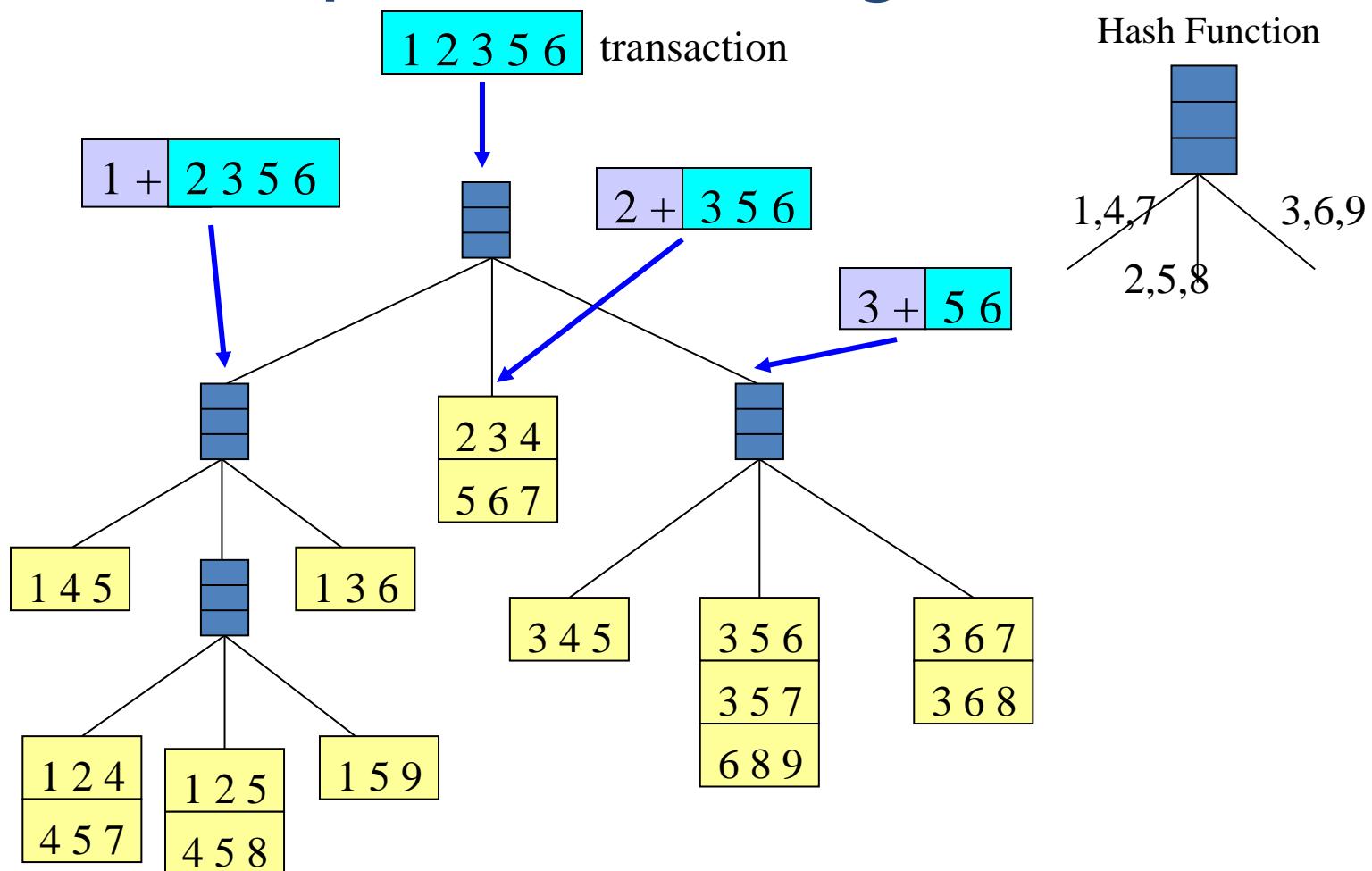
You need:

- Hash function
- Leafs: Store the itemsets

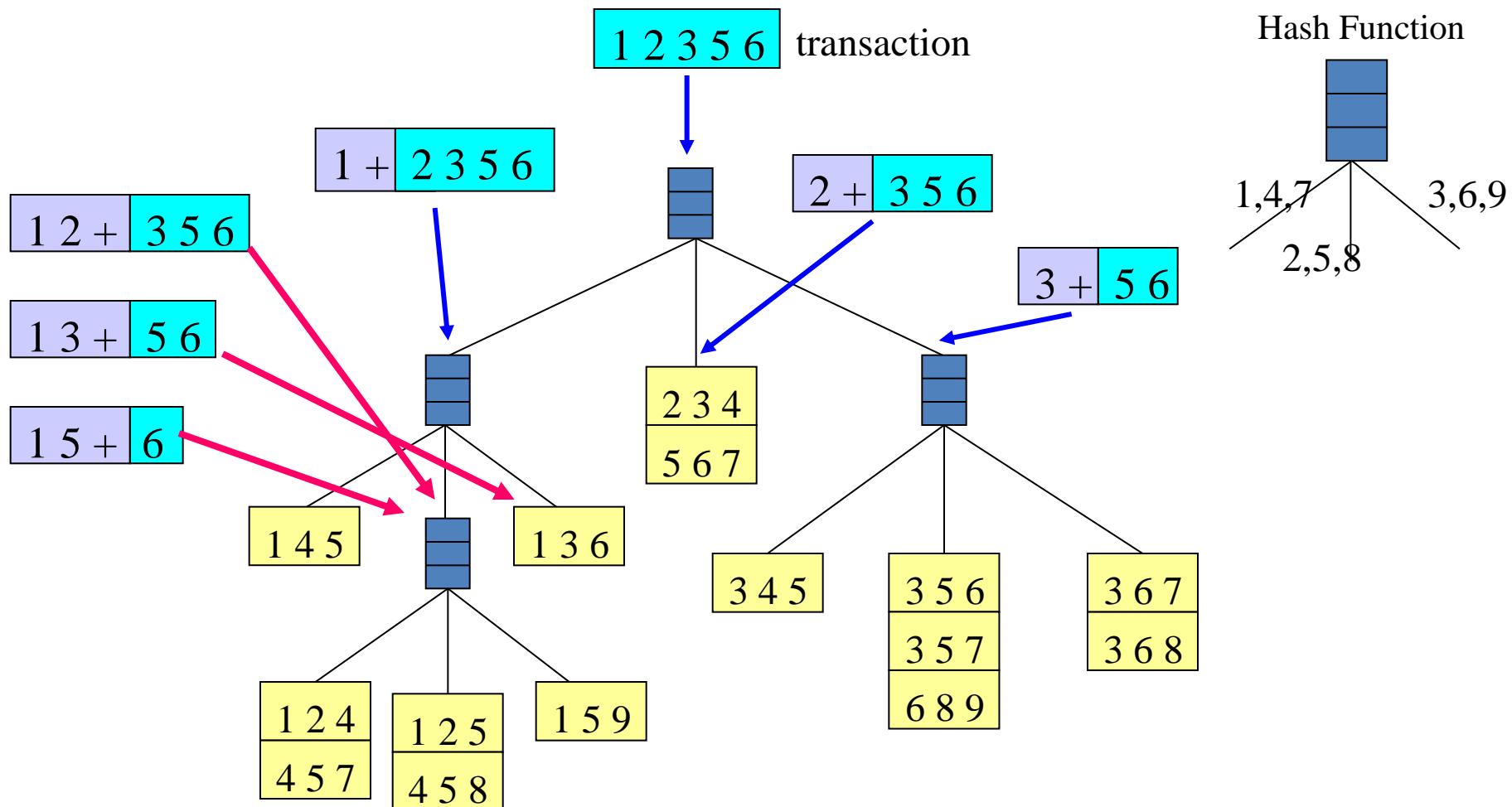


At the i-th level we hash on the i-th item

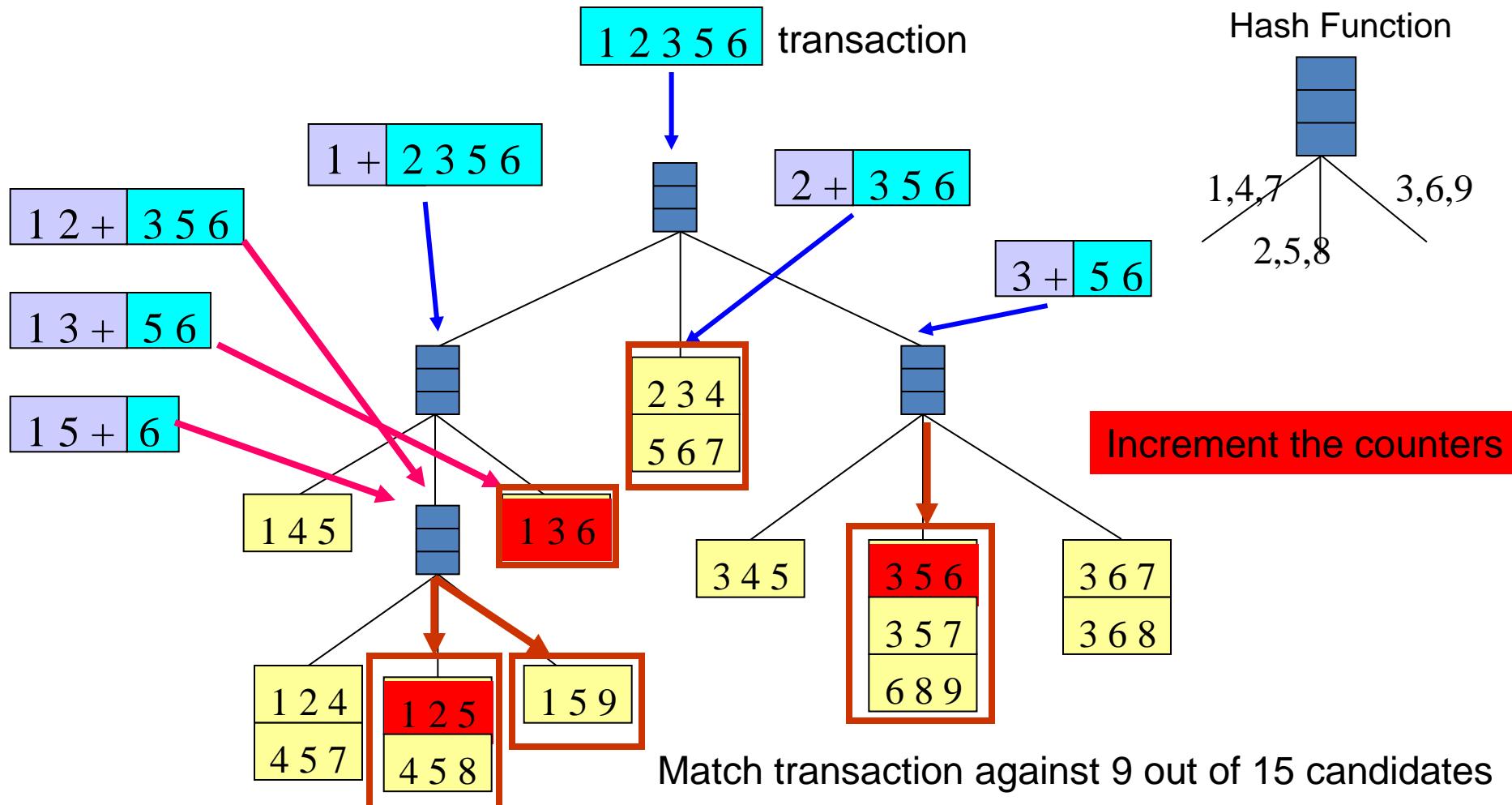
# Subset Operation Using Hash Tree



# Subset Operation Using Hash Tree



# Subset Operation Using Hash Tree



# Data Objects

---

- Data sets are made up of data objects.
- A **data object** represents an entity.
- Examples:
  - sales database: customers, store items, sales
  - medical database: patients, treatments
  - university database: students, professors, courses
- Also called *samples*, *examples*, *instances*, *data points*, *objects*, *tuples*.
- Data objects are described by **attributes**.
- Database rows -> data objects; columns -> attributes.

# Attributes

---

- **Attribute** (or **dimensions, features, variables**): a data field, representing a characteristic or feature of a data object.
  - *E.g., customer\_ID, name, address*
- Types:
  - Nominal
  - Binary
  - Numeric: quantitative
    - Interval-scaled
    - Ratio-scaled

# Attribute Types

- **Nominal:** categories, states, or “names of things”
  - $Hair\_color = \{auburn, black, blond, brown, grey, red, white\}$
  - marital status, occupation, ID numbers, zip codes
- **Binary**
  - Nominal attribute with only 2 states (0 and 1)
  - Symmetric binary: both outcomes equally important
    - e.g., gender
  - Asymmetric binary: outcomes not equally important.
    - e.g., medical test (positive vs. negative)
    - Convention: assign 1 to most important outcome (e.g., HIV positive)
- **Ordinal**
  - Values have a meaningful order (ranking) but magnitude between successive values is not known.
  - $Size = \{small, medium, large\}$ , grades, army rankings

# Numeric Attribute Types

---

- Quantity (integer or real-valued)
- **Interval**
  - Measured on a scale of **equal-sized units**
  - Values have order
    - E.g., *temperature in C° or F°, calendar dates*
- **Ratio**
  - Inherent **zero-point**
  - We can speak of values as being an order of magnitude larger than the unit of measurement (10 K° is twice as high as 5 K°).
    - e.g., *temperature in Kelvin, length, counts*

# Discrete vs. Continuous Attributes

---

## ■ Discrete Attribute

- Has only a finite or countably infinite set of values
  - E.g., zip codes, profession, or the set of words in a collection of documents
- Sometimes, represented as integer variables
- Note: Binary attributes are a special case of discrete attributes

## ■ Continuous Attribute

- Has real numbers as attribute values
  - E.g., temperature, height, or weight
- Practically, real values can only be measured and represented using a finite number of digits
- Continuous attributes are typically represented as floating-point variables

# **Similarity and Dissimilarity**

---

- **Similarity**
  - Numerical measure of how alike two data objects are
  - Value is higher when objects are more alike
  - Often falls in the range [0,1]
- **Dissimilarity** (e.g., distance)
  - Numerical measure of how different two data objects are
  - Lower when objects are more alike
  - Minimum dissimilarity is often 0
  - Upper limit varies
- **Proximity** refers to a similarity or dissimilarity

# Data Matrix and Dissimilarity Matrix

## ■ Data matrix

- n data points with p dimensions

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{np} \end{bmatrix}$$

## ■ Dissimilarity matrix

- n data points, but registers only the distance
- A triangular matrix

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

# Proximity Measure for Nominal Attributes

---

- Can take 2 or more states, e.g., red, yellow, blue, green (generalization of a binary attribute)
- Method 1: Simple matching
  - $m$ : # of matches,  $p$ : total # of variables

$$d(i, j) = \frac{p - m}{p}$$

- Method 2: Use a large number of binary attributes
  - creating a new binary attribute for each of the  $M$  nominal states

# Proximity Measure for Binary Attributes

- A contingency table for binary data
- Distance measure for symmetric binary variables:
- Distance measure for asymmetric binary variables:
- Jaccard coefficient (*similarity* measure for *asymmetric* binary variables):

		Object <i>j</i>		
		1	0	sum
Object <i>i</i>	1	<i>q</i>	<i>r</i>	<i>q + r</i>
	0	<i>s</i>	<i>t</i>	<i>s + t</i>
sum		<i>q + s</i>	<i>r + t</i>	<i>p</i>

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

$$d(i, j) = \frac{r + s}{q + r + s}$$

$$\text{sim}_{\text{Jaccard}}(i, j) = \frac{q}{q + r + s}$$

# Dissimilarity between Binary Variables

## ■ Example

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

- Gender is a symmetric attribute
- The remaining attributes are asymmetric binary
- Let the values Y and P be 1, and the value N 0

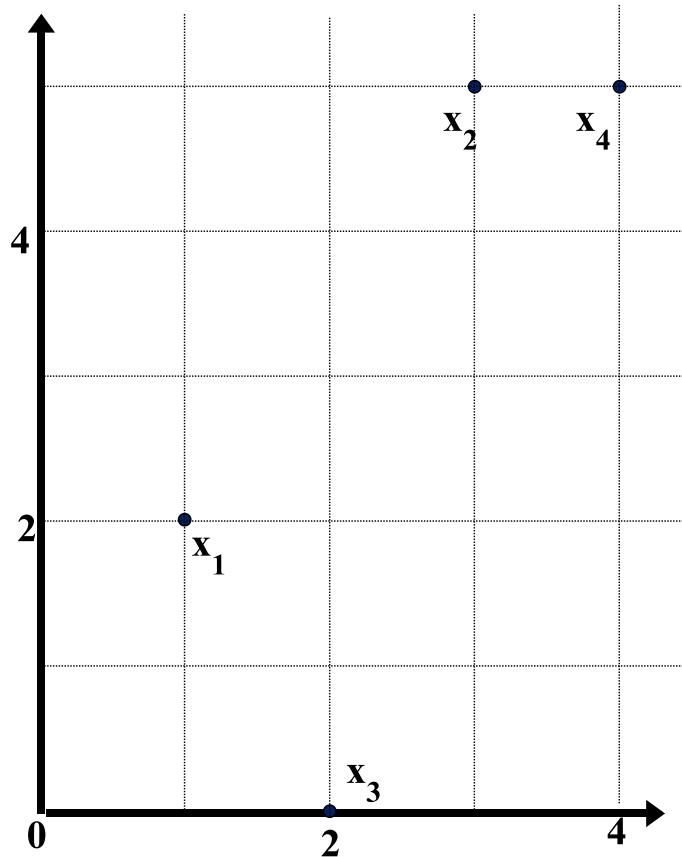
$$d(jack, mary) = \frac{0+1}{2+0+1} = 0.33$$

$$d(jack, jim) = \frac{1+1}{1+1+1} = 0.67$$

$$d(jim, mary) = \frac{1+2}{1+1+2} = 0.75$$

# Example:

## Data Matrix and Dissimilarity Matrix



**Data Matrix**

point	attribute1	attribute2
$x_1$	1	2
$x_2$	3	5
$x_3$	2	0
$x_4$	4	5

**Dissimilarity Matrix**

(with Euclidean Distance)

L2	$x_1$	$x_2$	$x_3$	$x_4$
$x_1$	0			
$x_2$	3.61	0		
$x_3$	2.24	5.1	0	
$x_4$	4.24	1	5.39	0

# Distance on Numeric Data: Minkowski Distance

- *Minkowski distance*: A popular distance measure

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \cdots + |x_{ip} - x_{jp}|^h}$$

where  $i = (x_{i1}, x_{i2}, \dots, x_{ip})$  and  $j = (x_{j1}, x_{j2}, \dots, x_{jp})$  are two  $p$ -dimensional data objects, and  $h$  is the order (the distance so defined is also called L- $h$  norm)

# Special Cases of Minkowski Distance

- $h = 1$ : Manhattan (city block,  $L_1$  norm) distance
  - E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

- $h = 2$ : ( $L_2$  norm) Euclidean distance

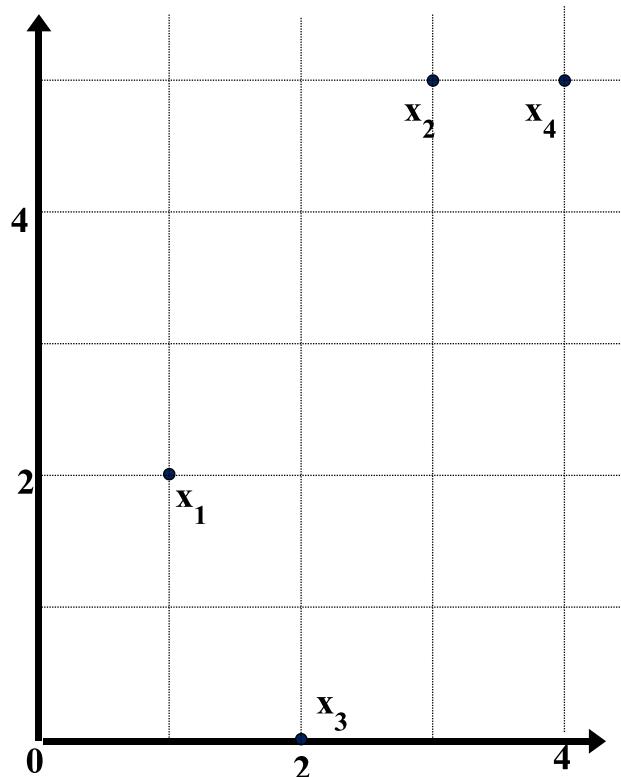
$$d(i, j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2)}$$

- $h \rightarrow \infty$ . “supremum” ( $L_{\max}$  norm,  $L_\infty$  norm) distance.
  - This is the maximum difference between any component (attribute) of the vectors

$$d(i, j) = \lim_{h \rightarrow \infty} \left( \sum_{f=1}^p |x_{if} - x_{jf}|^h \right)^{\frac{1}{h}} = \max_f^p |x_{if} - x_{jf}|$$

# Example: Minkowski Distance

point	attribute 1	attribute 2
x1	1	2
x2	3	5
x3	2	0
x4	4	5



## Dissimilarity Matrices

### Manhattan ( $L_1$ )

$L$	x1	x2	x3	x4
x1	0			
x2	5	0		
x3	3	6	0	
x4	6	1	7	0

### Euclidean ( $L_2$ )

$L_2$	x1	x2	x3	x4
x1	0			
x2	3.61	0		
x3	2.24	5.1	0	
x4	4.24	1	5.39	0

### Supremum

$L_\infty$	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0

# Ordinal Variables

---

- An ordinal variable can be discrete or continuous
- Order is important, e.g., rank
- Can be treated like interval-scaled
  - replace  $x_{if}$  by their rank  $r_{if} \in \{1, \dots, M_f\}$
  - map the range of each variable onto [0, 1] by replacing  $i$ -th object in the  $f$ -th variable by

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

- compute the dissimilarity using methods for interval-scaled variables

# Cosine Similarity

- A **document** can be represented by thousands of attributes, each recording the *frequency* of a particular word (such as keywords) or phrase in the document.

Document	team	coach	hockey	baseball	soccer	penalty	score	win	loss	season
Document1	5	0	3	0	2	0	0	2	0	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	0
Document4	0	1	0	0	1	2	2	0	3	0

- Other vector objects: gene features in micro-arrays, ...
- Applications: information retrieval, biologic taxonomy, gene feature mapping, ...
- Cosine measure: If  $d_1$  and  $d_2$  are two vectors (e.g., term-frequency vectors), then

$$\cos(d_1, d_2) = (d_1 \bullet d_2) / \|d_1\| \|d_2\|,$$

where  $\bullet$  indicates vector dot product,  $\|d\|$ : the length of vector  $d$

# Example: Cosine Similarity

---

- $\cos(d_1, d_2) = (d_1 \bullet d_2) / \|d_1\| \|d_2\|$ ,  
where  $\bullet$  indicates vector dot product,  $\|d\|$ : the length of vector  $d$
- Ex: Find the **similarity** between documents 1 and 2.

$$d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$d_2 = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$$

$$d_1 \bullet d_2 = 5*3 + 0*0 + 3*2 + 0*0 + 2*1 + 0*1 + 0*1 + 2*1 + 0*0 + 0*1 = 25$$

$$\|d_1\| = \sqrt{(5^2 + 0^2 + 3^2 + 0^2 + 2^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2)} = \sqrt{42} \approx 6.481$$

$$\|d_2\| = \sqrt{(3^2 + 0^2 + 2^2 + 0^2 + 1^2 + 1^2 + 0^2 + 1^2 + 0^2 + 1^2)} = \sqrt{17} \approx 4.12$$

$$\cos(d_1, d_2) = 0.94$$

# Basic Statistical Descriptions of Data

---

- Motivation
  - To better understand the data: central tendency, variation and spread
- Data dispersion characteristics
  - median, max, min, quantiles, outliers, variance, etc.
- Numerical dimensions correspond to sorted intervals
  - Data dispersion: analyzed with multiple granularities of precision
  - Boxplot or quantile analysis on sorted intervals

# Measuring the Central Tendency

- Mean (algebraic measure) (sample vs. population):

Note:  $n$  is sample size and  $N$  is population size.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \mu = \frac{\sum x}{N}$$

- Weighted arithmetic mean:

- Median:

- Middle value if odd number of values, or average of the middle two values otherwise

- Estimated by interpolation (for *grouped data*):

$$\text{median} = L_1 + \left( \frac{n/2 - (\sum freq)l}{freq_{median}} \right) width$$

- Mode

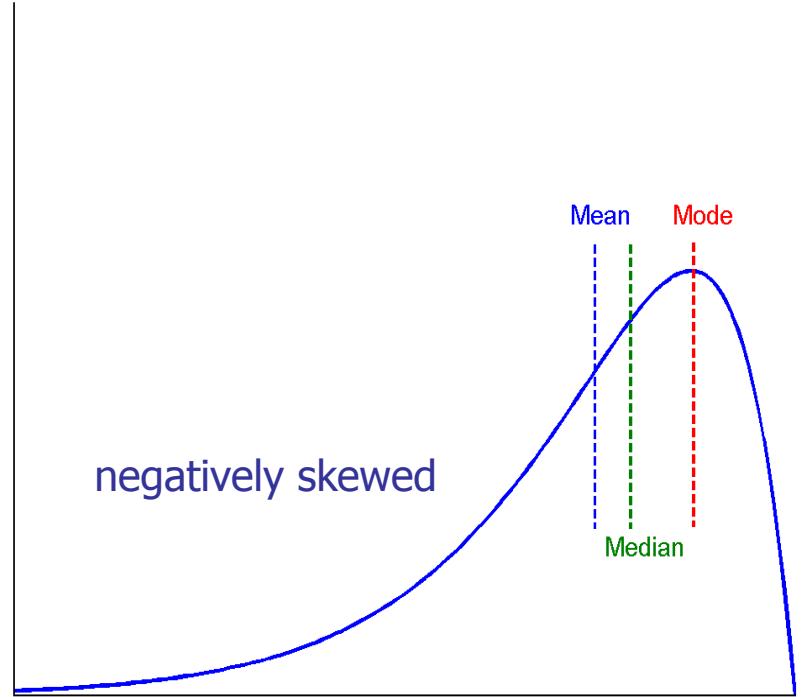
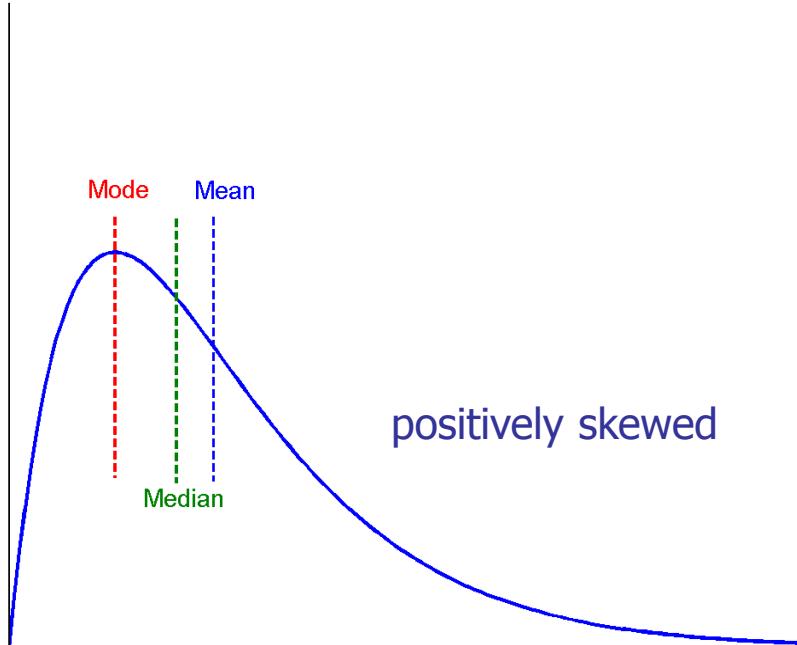
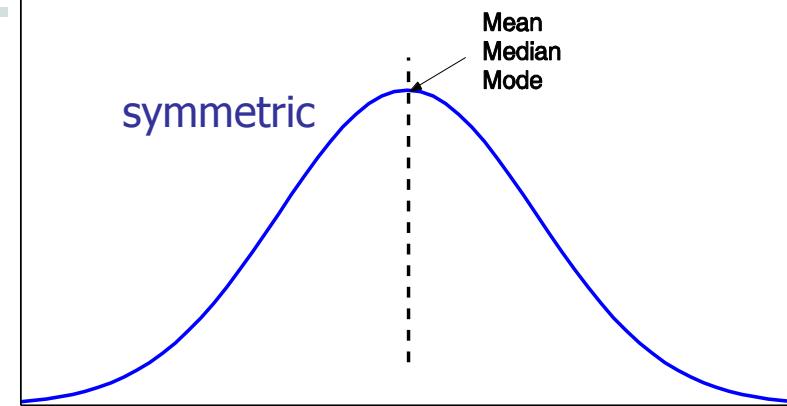
- Value that occurs most frequently in the data
  - Unimodal, bimodal, trimodal
  - Empirical formula:

$$mean-mode = 3 \times (mean - median)$$

age	frequency
1–5	200
6–15	450
16–20	300
21–50	1500
51–80	700
81–110	44

# Symmetric vs. Skewed Data

- Median, mean and mode of symmetric, positively and negatively skewed data



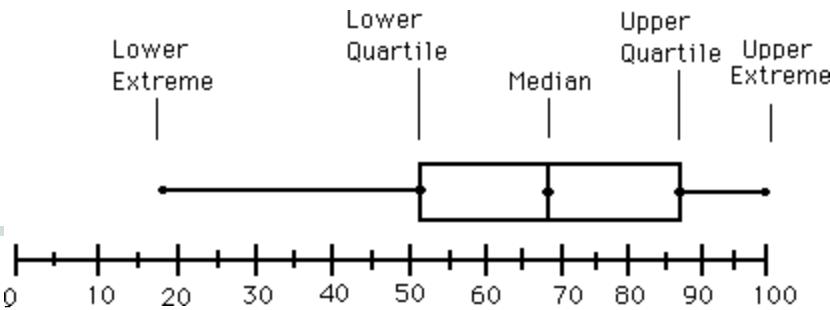
# Measuring the Dispersion of Data

- Quartiles, outliers and boxplots
  - **Quartiles:**  $Q_1$  ( $25^{\text{th}}$  percentile),  $Q_3$  ( $75^{\text{th}}$  percentile)
  - **Inter-quartile range:**  $\text{IQR} = Q_3 - Q_1$
  - **Five number summary:** min,  $Q_1$ , median,  $Q_3$ , max
  - **Boxplot:** ends of the box are the quartiles; median is marked; add whiskers, and plot outliers individually
  - **Outlier:** usually, a value higher/lower than  $1.5 \times \text{IQR}$
- Variance and standard deviation (*sample: s, population: σ*)
  - **Variance:** (algebraic, scalable computation)

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[ \sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2 \right] \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{N} \sum_{i=1}^n x_i^2 - \mu^2$$

- **Standard deviation**  $s$  (or  $\sigma$ ) is the square root of variance  $s^2$  (or  $\sigma^2$ )

# Boxplot Analysis

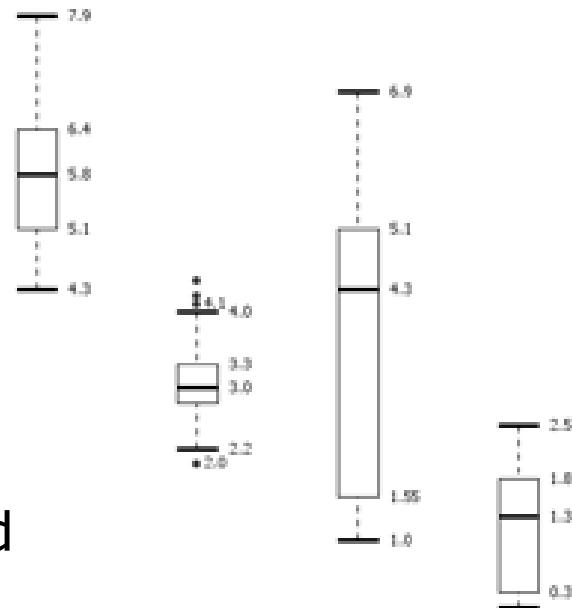


- **Five-number summary** of a distribution

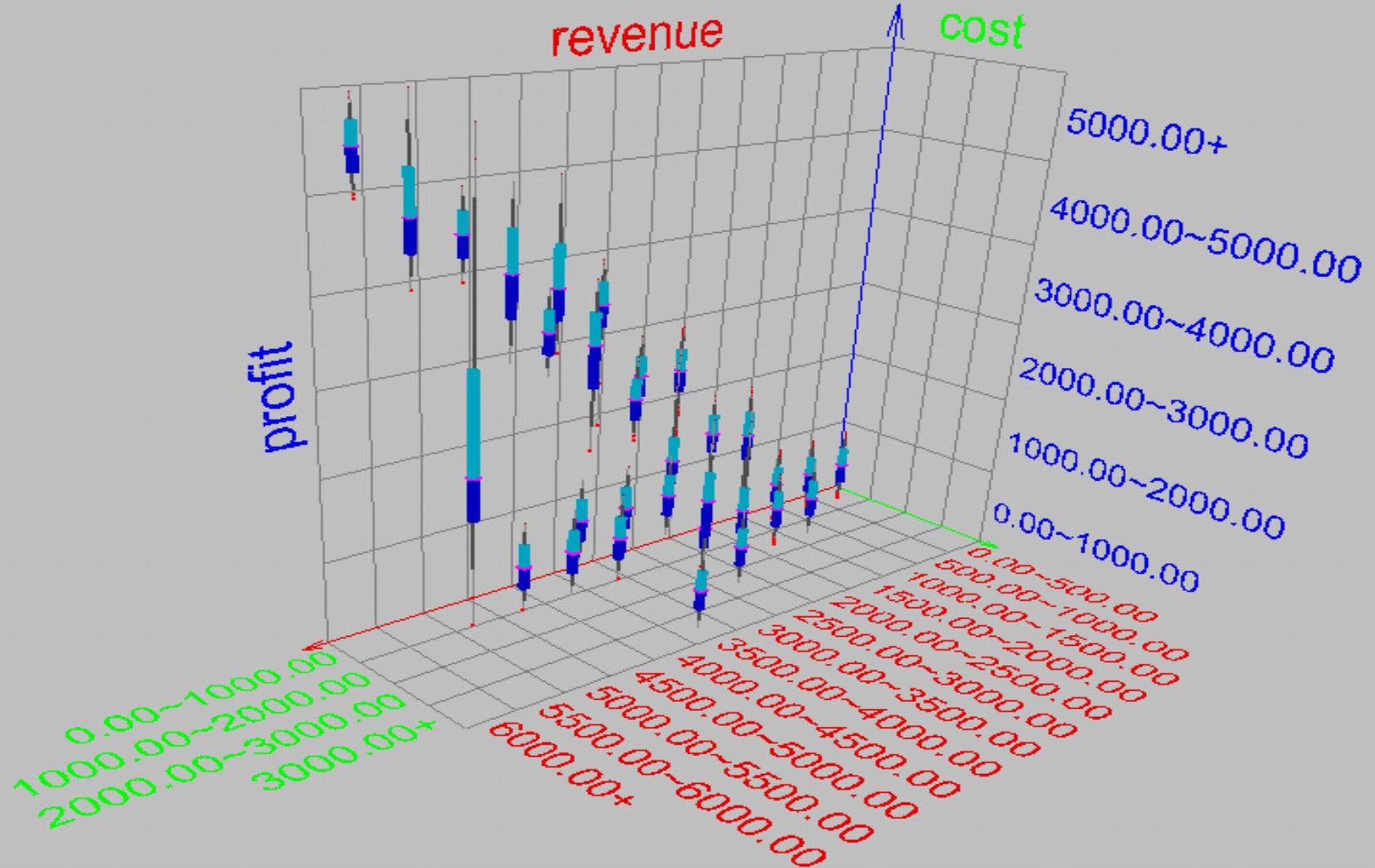
- Minimum, Q1, Median, Q3, Maximum

- **Boxplot**

- Data is represented with a box
  - The ends of the box are at the first and third quartiles, i.e., the height of the box is IQR
  - The median is marked by a line within the box
  - Whiskers: two lines outside the box extended to Minimum and Maximum
  - Outliers: points beyond a specified outlier threshold, plotted individually

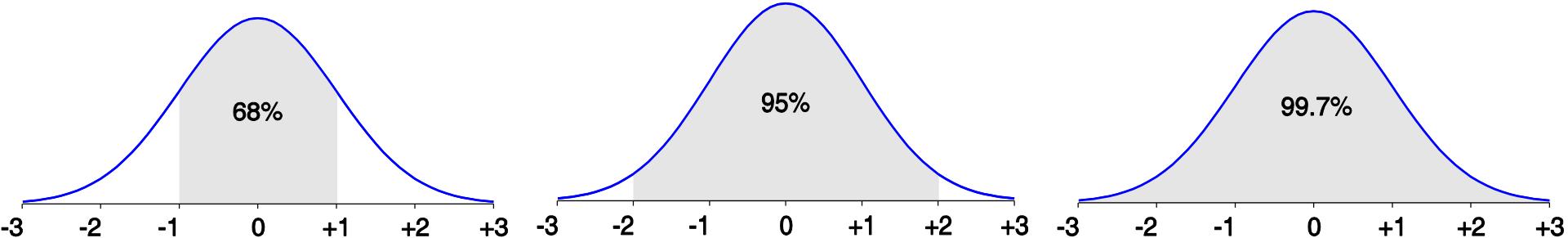


# Visualization of Data Dispersion: 3-D Boxplots



# Properties of Normal Distribution Curve

- The normal (distribution) curve
  - From  $\mu-\sigma$  to  $\mu+\sigma$ : contains about 68% of the measurements ( $\mu$ : mean,  $\sigma$ : standard deviation)
  - From  $\mu-2\sigma$  to  $\mu+2\sigma$ : contains about 95% of it
  - From  $\mu-3\sigma$  to  $\mu+3\sigma$ : contains about 99.7% of it



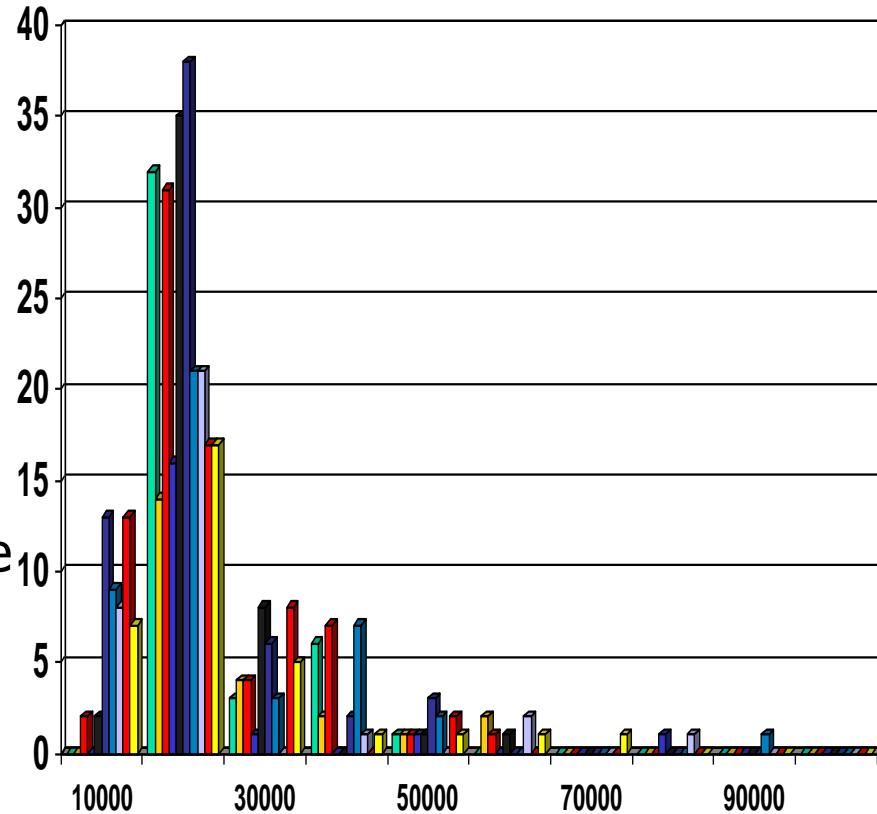
# Graphic Displays of Basic Statistical Descriptions

---

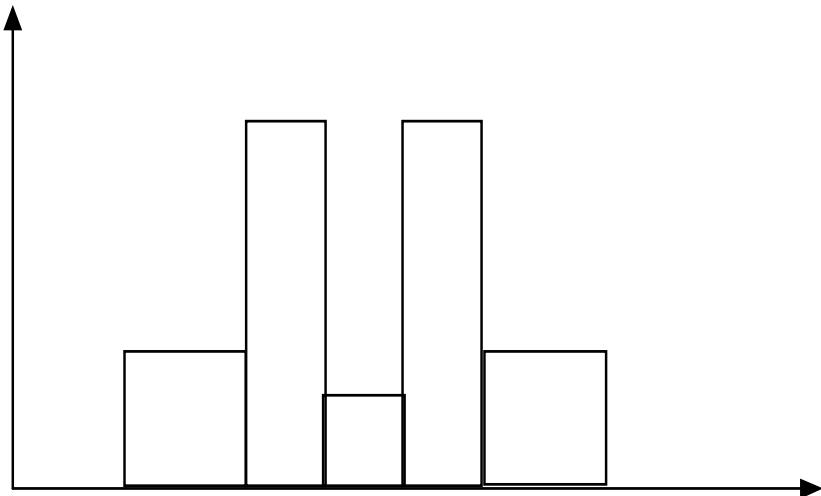
- **Boxplot:** graphic display of five-number summary
- **Histogram:** x-axis are values, y-axis repres. frequencies
- **Quantile plot:** each value  $x_i$  is paired with  $f_i$  indicating that approximately  $100 f_i\%$  of data are  $\leq x_i$
- **Quantile-quantile (q-q) plot:** graphs the quantiles of one univariant distribution against the corresponding quantiles of another
- **Scatter plot:** each pair of values is a pair of coordinates and plotted as points in the plane

# Histogram Analysis

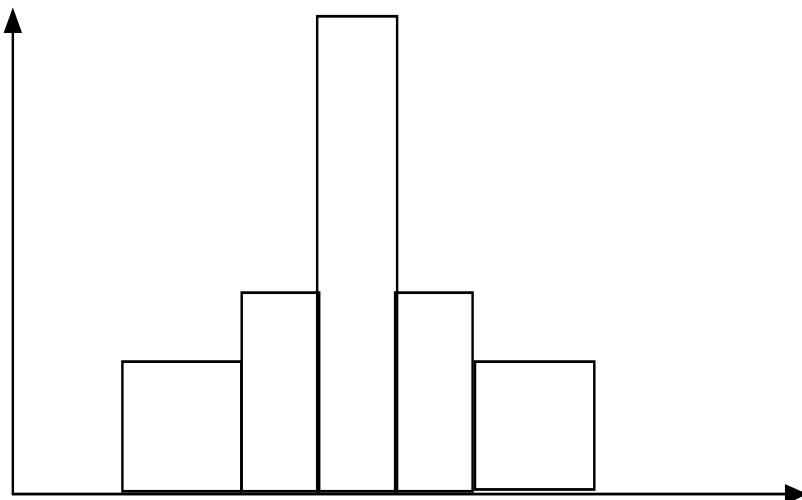
- Histogram: Graph display of tabulated frequencies, shown as bars
- It shows what proportion of cases fall into each of several categories
- Differs from a bar chart in that it is the *area* of the bar that denotes the value, not the height as in bar charts, a crucial distinction when the categories are not of uniform width
- The categories are usually specified as non-overlapping intervals of some variable. The categories (bars) must be adjacent



# Histograms Often Tell More than Boxplots

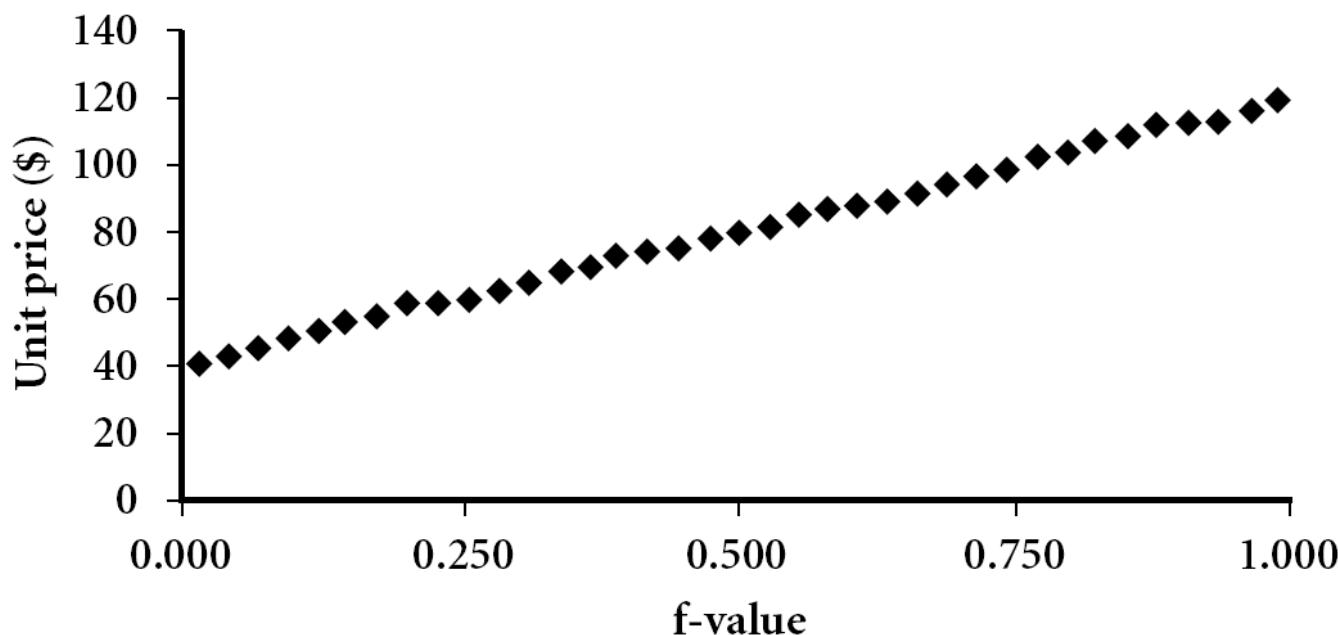


- The two histograms shown in the left may have the same boxplot representation
  - The same values for: min, Q1, median, Q3, max
- But they have rather different data distributions



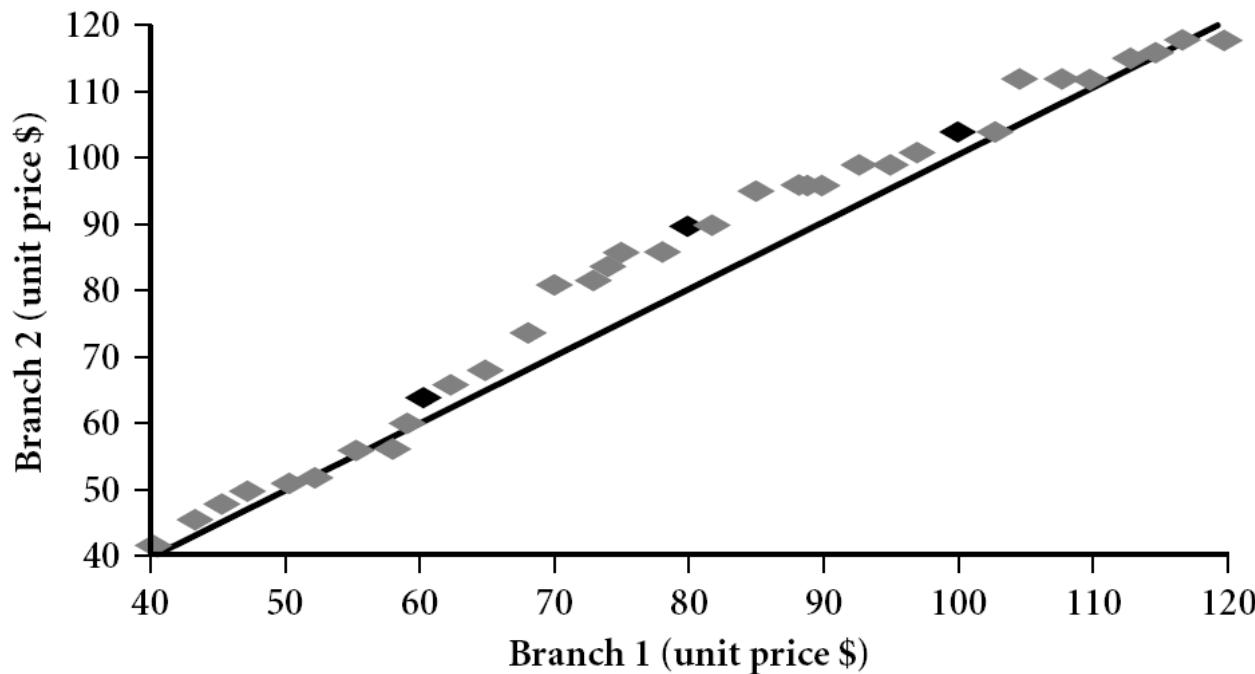
# Quantile Plot

- Displays all of the data (allowing the user to assess both the overall behavior and unusual occurrences)
- Plots **quantile** information
  - For a data  $x_i$ , data sorted in increasing order,  $f_i$  indicates that approximately  $100 f_i\%$  of the data are below or equal to the value  $x_i$



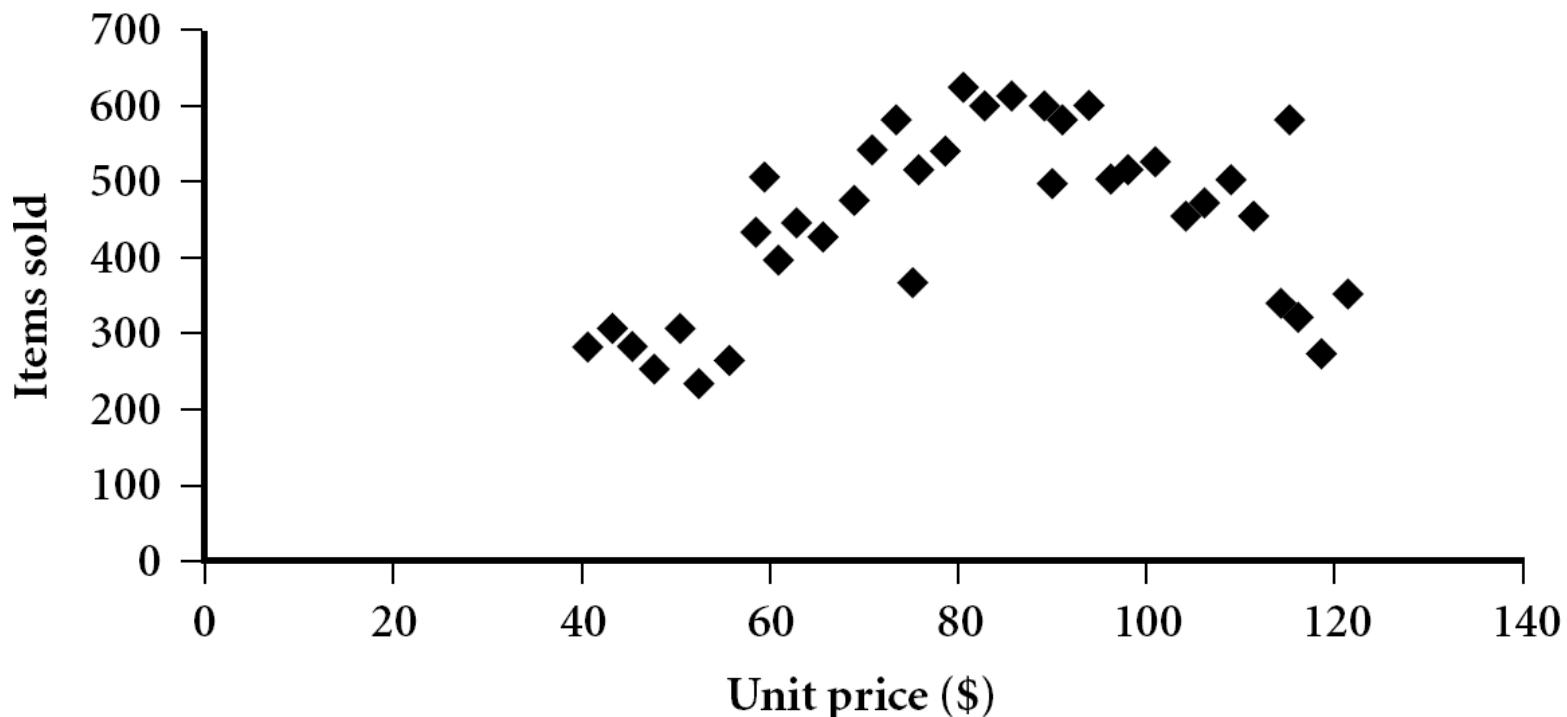
# Quantile-Quantile (Q-Q) Plot

- Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- View: Is there is a shift in going from one distribution to another?
- Example shows unit price of items sold at Branch 1 vs. Branch 2 for each quantile. Unit prices of items sold at Branch 1 tend to be lower than those at Branch 2.

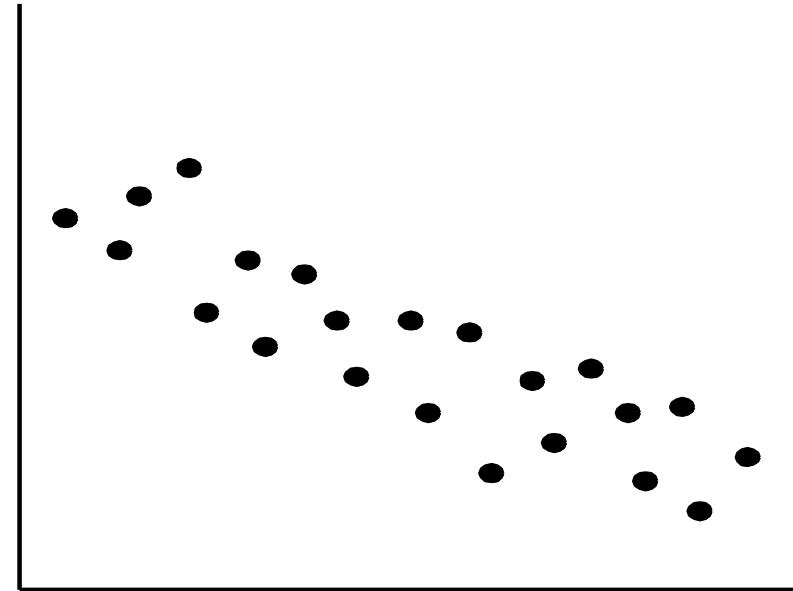
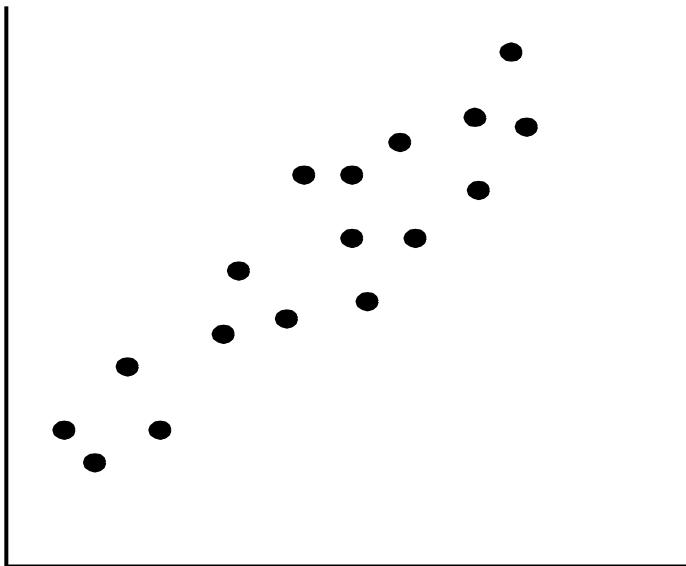


# Scatter plot

- Provides a first look at bivariate data to see clusters of points, outliers, etc
- Each pair of values is treated as a pair of coordinates and plotted as points in the plane



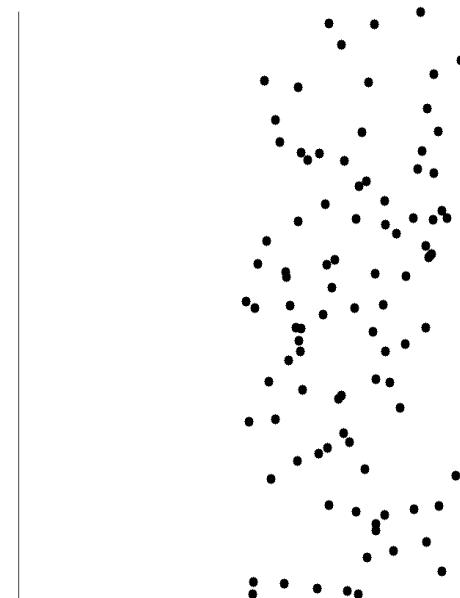
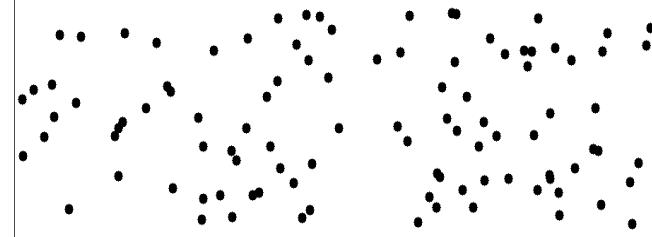
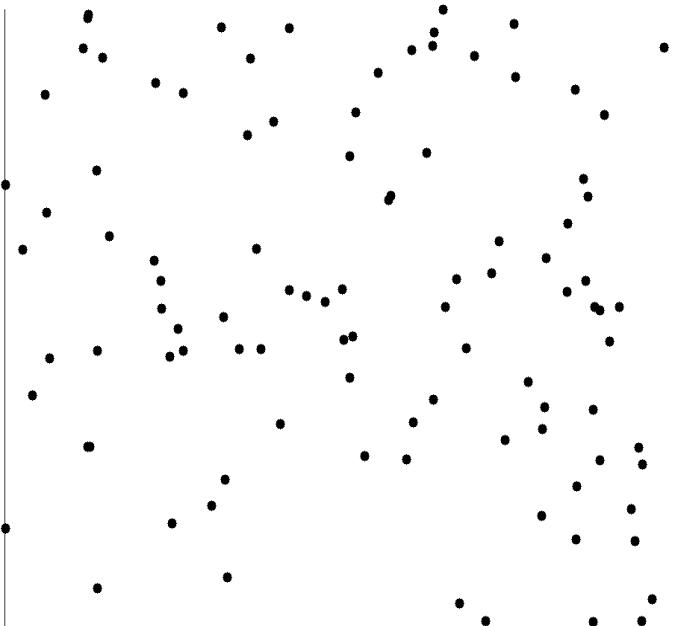
# Positively and Negatively Correlated Data



- The left half fragment is positively correlated
- The right half is negative correlated

# Uncorrelated Data

---



# DWDM

Data integration: Redundancy issues: An attr. might be redundant if it is derived from another attribute

$r_{AB} \rightarrow$  Correlation measure b/w attr. A & B

$$r_{AB} = \frac{\sum_{i=1}^n (A_i - \bar{A})(B_i - \bar{B})}{n \sigma_A \sigma_B} \quad (\text{for numeric data})$$

$\bar{A} \rightarrow$  Mean.  $\sigma_A \rightarrow$  Std. dev.  $n \rightarrow$  No. of tuples

If  $r_{AB} > 0 \rightarrow$  A, B are +vely dependent

$r_{AB} = 0 \rightarrow$  No correlation

$r_{AB} < 0 \rightarrow$  A, B are -vely correlated

For categorical data:

Chi Sq. ( $\chi^2$ )

(Pearson correlation)

Eg:

	M	F
Magazine	250	200
Newspaper	50	1000

Let A have C distinct values.  $a_1, a_2, \dots, a_c$   
B have r distinct values  $b_1, b_2, \dots, b_r$

For  $(A_i, B_j)$

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$O_{ij} \rightarrow$  Actual count (Observed freq.)

$E_{ij} \rightarrow$  Expected freq.

$$E_{ij} = \text{Count}(A=a_i) \times \text{Count}(B=b_j)$$

N

$$\chi^2 = \sum_{i=1}^r \frac{(O_i - E_i)^2}{E_i}$$

M      F

	M	F	
Mag.	250	200	450
News	50	100	150
	300	1200	1500

$$E_{ij} = \frac{C_{Male} C_{Mag.}}{(1200 + 300)} = \frac{300 \times 450}{1500} = 90$$

$$E_{11, News} = \frac{300 \times 150}{1500} = 210$$

$$E_{F, Mag.} = \frac{450 \times 200}{1500} = 360$$

$$E_{F, N} = \frac{210 \times 360}{1500} = 840$$

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(200-360)^2}{360} + \frac{(100-840)^2}{840} = 507.93$$

$$= 4 \times \frac{(160)^2}{90} = 4 \times \frac{25600}{90}$$

From stat. table, for  $2 \times 2$  matrix,

$$\chi^2 = 10.828(r-1)(c-1) \text{ deg. of freedom}$$

$$= 10.828$$

Calculated value ( $\chi^2$ )  $\geq$  Table value ( $\chi^2$ )

$\hookrightarrow$  Correlated values

$\therefore$  Gender & Reading materials are correlated.

## Feature selection

- Choose optimal subset of features
  - Remove irrelevant features.
  - Inc. pred. accuracy
  - Improve learning efficiency → Reduce storage space  
" computation cost

Searching best subset of features

Criteria for evaluating different subsets.

Selection of best subset:

Heuristic methods:

→ Stepwise forward selection

{A<sub>1</sub>, A<sub>2</sub>, ..., A<sub>n</sub>}

Start with {} add best attr. {A<sub>best</sub>}, add next best {A<sub>best+1</sub>} - - -

→ Stepwise backward elimination:

Start with all, keep removing worst attr.

→ Combination of forward selection & backward elimination: At each step, select a best attr. & eliminate a worst attr.

Criteria for best attr. split:

Gain measure, gini measure

Attributes that do not occur on generating Decision tree are considered irrelevant.

(Wrapper based)

Genetic algorithm: → Maximise accuracy measure

Init. chromosome.

1011010

3 absent

feature

present

fitness fn

Evaluate chromosome with fitness function



Selection, crossover, mutation

Genetic operations  
Check stopping criteria.

→ Max iter,

→ Converged test fn  
subset of

Criteria to evaluate features:

$$\textcircled{2} \text{ Entropy measure} = - \sum_{i=1}^c P(C_i) \log_2(C_i)$$

→ no. of class label

Information gain

$$Inf(A) = - \sum_{j=1}^{|D|} \frac{|D_j|}{|D|} \pi(D_j)$$

$|D| \rightarrow$  no. of types partition j.

$|D_j| \rightarrow$  No. of tuple in  $D_j$

Distance measure:

$$\left( \sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$$

(Depend on p,  
can be Manhattan dist,  
euclidean dist, etc)

Dependency measure / Correlation measure:

$$r(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{n} \sigma_x \sigma_y}$$

Consistency measure:

If input feature values are same

Same class label ( $O_1$ )

Accuracy measurement

Evaluate all possible subsets from accuracy method. Subset with highest acc. is considered.

### Models of FS:

Filter based model (Faster) No

Original  
FS

Gen.  
subsets

Measures

General  
characteristics  
(D data)

Stop?

Yes

Eg:- dependency  
measure,

Gini measure,  
consistency mea.  
etc

Validation

Accuracy  
Analysis

Wrapper based model (Slower)

Og  
features

Generate  
subsets

Evaluate by  
using learning  
algorithm/  
classifier

No

Stop?

Yes

Analyse

Validate

Usually better results, longer  
computation energy & time.

## Data transformation:

A fn that maps the entire data set of values of a given attribute to a new set.

### Methods

- Smoothing Remove noise from data
- Feature selection construction
- Aggregation → Summarisation
- Normalisation
- Generalisation

### Normalisation:

Min max normalisation

$$V' = \frac{V - \text{min}}{\text{Max} - \text{min}} \quad (\text{new max-new min})$$

+ new min

$V'$  → normalised value

$V$  → Attr. value

new max new min → scales to fit data to (Range).

### Z score normalisation

$$V' = \frac{V - \mu}{\sigma} \quad \begin{cases} \mu \rightarrow \text{mean} \\ \sigma \rightarrow \text{SD} \end{cases}$$

### Decimal scaling based normalisation:

$$V' = \frac{V}{10^j} \quad j \rightarrow \text{smallest int. such that } \text{Max}(V') < 1$$

Ex:- 986 → 9.17.  $j$  is 1000 → 3  
 $0.917 \rightarrow 0.986$

# Principal Component Analysis

Dimensionality reduction:

Steps:

- ① Collect data
- ② Subtract mean from data
- ③ Calc covariance matrix
- ④ Calc eigenvalues & eigenvectors of the covariance matrix
- ⑤ Order eigenvectors by eigenvalues (high to low)
- ⑥ Select 'p' eigenvectors
- ⑦ Derive new data.

X	Y	$\bar{x}$	$\bar{y}$
2.5	2.4	2.1	1.91
0.5	0.7	2.1	1.91
2.2	2.9	2.1	1.91
1.9	2.2	$\bar{x} - \bar{x} = 0.69, -1.31, 0.39, 0.09,$	
3.1	3.0		$1.29, 0.49, 0.19, -0.81,$
2.3	2.7		$-0.31, -0.71$
2	1.6	$y - \bar{y} = 0.49, -1.21, 0.99, 0.29,$	
1	1.1		$1.09, 0.79, -0.31, -0.81, -0.31$
1.5	1.6		$-1.01$
1.1	0.9		
$\bar{x} - \bar{x}$	$\bar{y} - \bar{y}$		

(3) Covariance matrix ( $2 \times 2$ )

$$\text{cov}(x_i, y_j) = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$= \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix}^{(n-1)}$$

$$= \begin{bmatrix} 0.616155 & 0.61544 \\ 0.61544 & 0.71655 \end{bmatrix}$$

$$\text{Eigenvalues} = \begin{pmatrix} 0.0490833 \\ 1.28462771 \end{pmatrix}$$

$$\text{Eigen vectors} = \begin{pmatrix} -0.6778 & -0.7351 \\ 0.6778 & -0.7351 \end{pmatrix}$$

Order the vectors:

$$\begin{bmatrix} -0.6778 & -0.7351 \\ -0.7351 & 0.6778 \end{bmatrix}$$

↓  
Select

Derive new data

$$(-0.6778 \times 0.69) + (-0.7351 \times 0.49) \text{ so on.}$$

Modified  $X$  :-  $-0.827, 1.7, -0.992, -0.274,$   
 $-1.6, 0.91, 0.099, 1.14, 0.43, 0.122$

Clustering :-

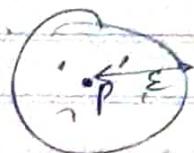
- Simple but sensitive to noise (Every obj.  $\in 1$  cluster)
- Partition based (No need to know  $k$ )
- Hierarchical based ( $k$ ;  $O(n^2)$ )  
 (Agglomerative & divisive) Dendrogram
- Clusters
- $O(n^k t)$  → Iterations
- no. of data points

FCM :- Fuzzy C means; Object & more than 1 cluster  
 with membership matrix

Density based clustering

↳ No. of points in given no. of radius

DBScan (Density based) :-



P is a core point if with  $\epsilon$  radius, there are given min. number of points around it.

# INTRODUCTION

A point is border point if a point is not core point but lies within the neighbourhood of a core point.

A noise point is a point that is neither a core nor a border point. These points are to be discarded.

Adv : Not sensitive to noise ; Any shape of cluster is possible

Disadv :- Picking parameter values ( $\epsilon$  & min val.)

## OPTICS: Ordering Points To Identify Clustering

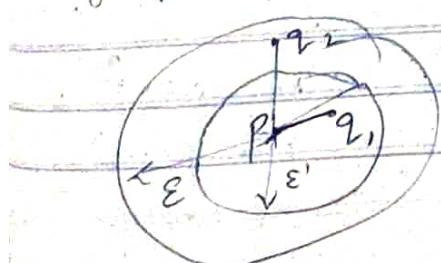
Need to identify core dist, reachability dist along with core, border point



In fig 2,  $\epsilon' < \epsilon$  & point is core pt even with  $\epsilon'$ .

Core dist. of an objective 'p' is the smallest  $\epsilon'$  value that makes 'p' as a core point

Reachability dist of an object 'q' w.r.t. another object 'p' is the greater value of the core dist of 'p' & euclidian dist of  $(p \rightarrow q)$ .



$$\text{Reachability dist}(p, q_1) = \epsilon'$$

$$\begin{aligned} \text{Reachability}(p, q_2) &= \max(\epsilon', p_{q_2}) \\ &= \text{Euclidean dist}(p, q_2) \end{aligned}$$



# SRIRAM

Options → It stores core dist & reachability dist. for each obj.

→ Maintains a list called seedlist → objects are sorted by reach. dist from their resp. closest core points.

? ? ?

## CLIQUE (Clustering in Quest)

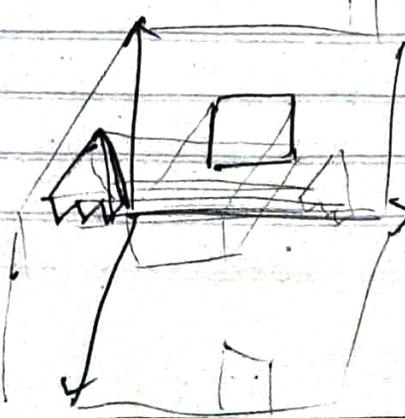
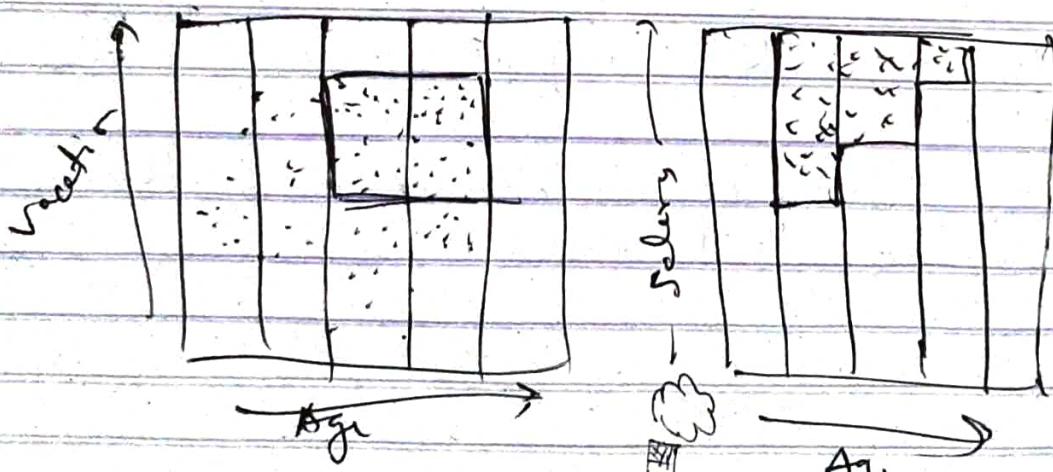
A dimension growth subspace clustering  
(first of its kind)

Starts with a single dimension, divides this dim into grids.

Identifies dense units in grids (no. of minpts)  
The subspaces representing these dense units are intersected to form candidate search for higher dim.

Get maximal & minimal region

Highly  $\downarrow$  dense      Low  $\downarrow$  dense  
units                          region



Grid size??  
& min points  
to be chosen

# FREERAM GHEERAM

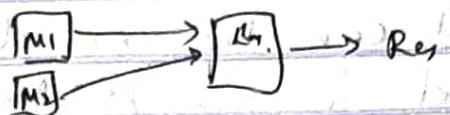
## Classification

- ↳ Naive Bayes
- ↳ Decision tree
- ↳ KNN
- ↳ Artif. Neural Net
  - ↳ Backprop.
  - ↳ Perception

## Ensemble classifiers

- ↳ Set of classifiers

(Learning is based on comb. of multiple classifiers)



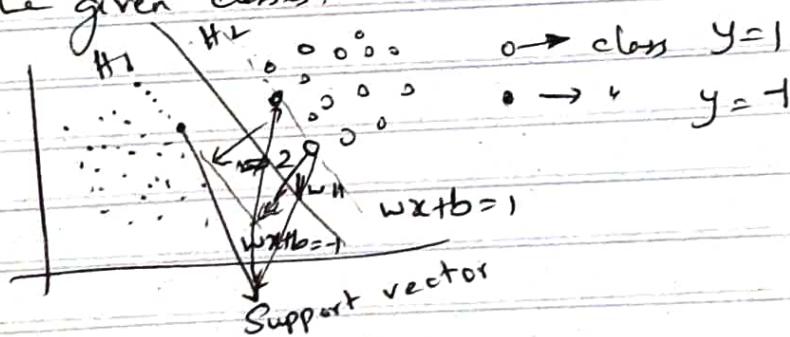
Bagging → Each classifier returns class pred.  
for  $x$  which counts as a vote. Democracy

Boosting → Assign wt to each classifier based on how well it performs.

Sum weights of each classifier that assigned class c. Class with highest wt. wins.

## SVM (Support Vector Machine)

→ Max. margin hyperplane (MMH) that separates the given classes.



Any hyperplane  $wx+b=0$   
 $w$  is normal vector to the hyperplane  $w_1, w_2, \dots, w_n$   
 $n$  is no. of attrs.

$x$  is the training tuple

$b$  is the bias

Need to maximise dist. b/w  $H_1, H_2$   
 $\text{dist}(H_1, H_2) = \frac{2\|w\|}{\sqrt{w^2}}$

$$d(x_t) = \sum_{i=1}^l y_i \alpha_i x_i \cdot x_t + b_0$$

$y_i \rightarrow$  Class Label for sup vector  $x_i$

$l \rightarrow$  No. of support vector

$x_t \rightarrow$  test tuple

$\alpha_i, b_0 \rightarrow$  Numeric parameters defined by user

$d(x_i) > 0 \rightarrow$  +ve class

$d(x_i) < 0 \rightarrow$  -ve class

SVM can be non-linearly separable data  
Kernel Functions → Maps data to higher dim feature space where the data points are linearly separable

~~7b/2d~~

## Data Preprocessing :

→ Y?

real-world data might be incomplete, noisy,  
inconsistent.

e.g.: DOB  
year / birth

missing values

outliers  
errors

① Data cleaning → To fill the missing values  
& to smooth the data.

② Data integration .

data from diff sources are integrated .



③ Data Transformation .

\* -2, 32, 100, 50, 48 → -0.02, 0.32, 0.100, 0.5, 0.48  
Normalisation .

\* Aggregation

monthly sales data → yearly sales data .

\* generalization → Mapping <sup>from</sup> lower level  
concepts to high level

street → city .

④ Data Reduction .

	$A_1$	$A_2$	$A_3$	$(A_4)$	$\dots$	$A_{70}$
$T_1$						
$T_2$						
$\vdots$						
$T_{70}$						

tuples  
remain  
same  
→

$A_1$	$A_2$	$A_3 \dots A_{70}$
$T_1$		
$T_2$		
$\vdots$		
$T_{70}$		

→ FS, PCA  
feature principal  
selection comp Analysis

### ① Data Cleaning

- i) → - fill the missing values
- ii) → - smooth the data

→ \* ① ignore the tuple.

if % of missing value is high, ~~too less~~ data will ~~not~~ affect

if % of missing value is less (eg. < 1%)  
~~we might lose imp information~~

② Fill the missing values manually.

for smaller datasets & for fewer values  
its feasible.

③ Use some global constant to fill missing values.  
eg: unknown,  $\infty$

#### Drawbacks

sometimes mining alg may consider this pattern as interesting patterns, but will not be useful for users at all.

→ not a good method

✓ ④ Use attribute mean value

- better method compared to others

✓ ⑤ Use most probable value to fill the missing value.  
LDT, Bayesian inference

Based on highest gain

because only few values present,  
so including them also wont have much effect

ignore

cannot ignore

ii) Smooth the data .

→ Binning — method are going to smooth the sorted data .

e.g: price .

4, 8, 15, ~~21~~, 21, 24, 25, 28, 34

→ \* use equal frequency binning  
same no of samples in each bin .

4 8 15    ~~21~~ 21 24    25 28 34  
Bin 1              Bin 2              Bin 3

\* smoothing by bin mean .

Bin 1 → 4, 8, 15 → 9, 9, 9 .

Bin 2 → ~~21~~, 21, 24 → 22, 22, 22 .

Bin 3 → 25, 28, 34 → 29, 29, 29 .

\* smooth by bin boundary .

Bin 1 → 4, 8, 15 → 4, 4, 15

Bin 2 → 21, 21, 24 → 21, 21, 24 .

Bin 3 → 25, 28, 34 → 25, 25, 34 .

increase such non boundary values to the closest boundary value .

\* smooth by median .

change every value to median value .

→ equal width binning .

bin size -  $\frac{\max - \min}{\text{no of elements}}$

- divide the range into  $N$  intervals of width  $w$

- interval =  $\frac{\text{Max} - \text{Min}}{N}$

→ Boundaries -

$$\text{min} + w, \text{min} + 2w, \dots, \text{min} + nw$$

Data

eg : 5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204  
222, 230

Width :  $\frac{230 - 5}{3} = 75$   
( $N=3$ )

Bin 1 = 5, 10, 11, 13, 15, 35, 50, 55, 72 [5 \$ 75]

Bin 2 = 92 (75 \$ 150)

Bin 3 = ~~204, 222,~~ 230 (150 \$ 225)

③ Data Preprocessing

② Data Integration

→ check for redundancy in the attributes

→ An attribute may be redundant if it is derived from another attribute.

$A, B \rightarrow$  dependent attribute

\* time of training  $T$

\* model won't be efficient

\* remove redundant attributes, keep only

independent ones

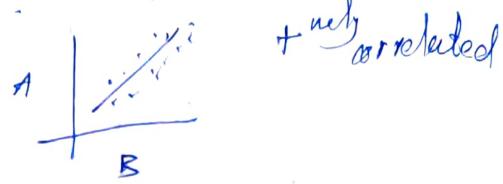
Correlation  $r_{A,B} = \frac{\sum_{i=1}^n (A_i - \bar{A})(B_i - \bar{B})}{\sqrt{n} \sqrt{\sum (B_i - \bar{B})^2}}$

$\bar{A}, \bar{B} \rightarrow$  mean values of  $A$  &  $B$  respectively.

$\sigma_A \sigma_B \rightarrow$  std deviation of A & B

$n \rightarrow$  no of tuples

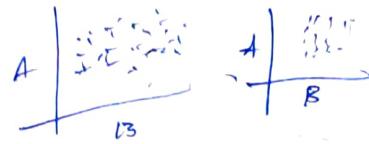
$r_{A,B} \rightarrow$  the  
Dependent



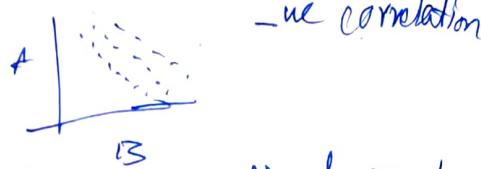
choose  
these  
attributes  $r_{A,B} \rightarrow 0$

no correlation

independent



$r_{A,B} \rightarrow -ve$   
dependent



$\Rightarrow$  only for numerical data this method can be used

$\Rightarrow$  Categorical Attributes

Chi-square Measure ( $\chi^2$ )

(Pearson correlation measure)

gender  $\rightarrow$  A has c distinct values  $a_1, a_2, \dots, a_c$   
 reading M F  $\rightarrow$  B has r " "  $b_1, b_2, \dots, b_r$   
 magazine also 100  
 newspaper 50 1000  
 let  $(A_i, B_j)$  denote the event that

A take value  $a_i$  & B take value  $b_j$ .

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$O_{ij} \rightarrow$  actual count  $\rightarrow$  observed freq

$E_{ij} \rightarrow$  expected frequency

$$E_{ij} = \frac{\text{count}(A=a_i) \times \text{count}(B=b_j)}{N}$$

$N \rightarrow$  no of tuples

	M	F	T
M	250	200	450
N	50	1000	1050
T	300	1800	

$$C_{11} = \frac{\text{count (male)} \times \text{count (magazine)}}{N}$$

$$= \frac{300 \times 450}{1800} = \underline{\underline{90}} \rightarrow c_{11}$$

C<sub>12</sub> = male

$$C_{(male, news)} = \frac{250 \times 1050}{1800} = \underline{\underline{210}}$$

$$C_{(F, mag)} = \frac{200 \times 450}{1800} = \underline{\underline{360}}$$

$$C_{(F, news P)} = \frac{1200 \times 1050}{1800} = \underline{\underline{840}}$$

$$\begin{aligned} \chi^2 &= \frac{(250 - 90)^2}{90} + \frac{(\cancel{50} - 210)^2}{210} \\ &\quad + \frac{(\cancel{200} - 360)^2}{360} + \frac{(1000 - 840)^2}{840} \\ &= \frac{160^2}{90} + \frac{170^2}{210} + \frac{160^2}{360} + \frac{160^2}{840} \end{aligned}$$

$$= \underline{\underline{567.93}}$$

statistics  $2 \times 2 \text{ matrix} = 10.828$

$(r-1)(c-1)$  degrees of freedom  
 $1 \times 1 = 1$

→ calculated ( $\chi^2$ ) ≥ Actual  $\chi^2$  value  
 From statistics, → both are correlated

$$507.93 \geq 10.828$$

Hence,  $\sigma$  & reading mat are correlated.

→ if there is no deviation between observed & expected freq,  $\chi^2$  value will be less.

15) 3/22

### Feature Selection

→ FS is a process that chooses an optimal subset of features from the original set.

Why FS?

→ To remove irrelevant features

→ To increase the predictive accuracy

→ To improve the learning efficiency — reduce the storage space — reduce the computational cost.

### Criterias for feature Selection

① searching for best subset of features

② Criteria for evaluating diff subsets

①  $\Rightarrow$  Heuristic Methods -

① Stepwise forward Selection

$$\{A_1, A_2, A_3, A_4, A_5, A_6\}$$

(all attributes given in D)

{ } 3 start with empty set .

{ $A_1$ } — best attribute

{ $A_1, A_2$ } → add see best attribute

: continue the process

$\{A_1, A_4, A_5\}$  — reduced set

→ every step add the best attribute from the given set.

(i) Stepwise backward elimination.

$\{A_1, A_2, A_3, A_4, A_5, A_6\}$  — original set.



every step  
remove one worst  
attribute.

$\{A_1, A_3, A_5, A_6\}$

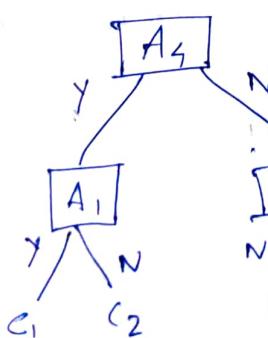
$\{A_1, A_5, A_6\}$  ← reduced set.

(ii) Combination of forward selection and Backward elimination.

→ each step select one best attribute and eliminate worst attribute.

(iii) Decision Tree.

→ gain measure or gini measure to find the best split attribute.



splitting attributes —  
best attributes.

← class label.

→ The attributes that does not occur in the tree are considered as irrelevant attributes.

DT → classification prob as well as feature selection

### v) GA (Genetic Alg)

→ ① Initialise the chromosomes

10110110 — length of chromosome  
| |  
feature absence of equal to the no of  
present feature features in the DS.

② → evaluate the chromosome by using fitness function.

Accuracy measure → fitness function maximization

③ → selection, crossover, mutation

select best exchange flip one bit.  
chromosomes genes

④ → stop check the stopping criteria —

max no of iterations or convergence of fitness function

⑤ ⇒ Criteria to evaluate the Subsets of feature

i) Entropy  $= \sum_{i=1}^c P(C_i) \log_2(C_i)$

$c$  — no. of class labels

$p(c_i) \rightarrow$  prob of class  $c_i$

i) Gain

$$\text{Info}(A) = - \sum_{j=1}^c \frac{|D_j|}{|D|} I(D_j)$$

$|D_j|$  — no. of tuples in the partition  $j$

$|D|$  — total no. of tuples in DS

$I(D_j)$  —

— If a particular feat gives the highest gain, we can keep the feature else eliminate it.

Attribute should give  
best gain?

ii) Distance Measure

$$\sum_{i=1}^n \left( |x_i - y_i|^p \right)^{\frac{1}{p}}$$

(not for supremum dist.)

$p=1 \rightarrow$  manhattan distance

$p=2 \rightarrow$  euclidean dist

similar features  $\rightarrow$  same class label

Keep the features  
(relevant)

else no consistency of class label — irrelevant features

③ Dependency Measures or Correlation measure

$$\gamma(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n \sigma_x \sigma_y}$$

$\gamma(x, y) = 0 \leftarrow$  independent.

eliminate  $\begin{cases} \gamma(x, y) = +ve & \leftarrow x \& y \text{ are +vely correlated.} \\ \text{the features} & \\ \gamma(x, y) = -ve & \leftarrow x \& y \text{ are -vely correlated.} \end{cases}$

(4) consistency measure

I/P feature values are same  $\Rightarrow$  same class label (o/p).

(with dist measure check if features are similar)

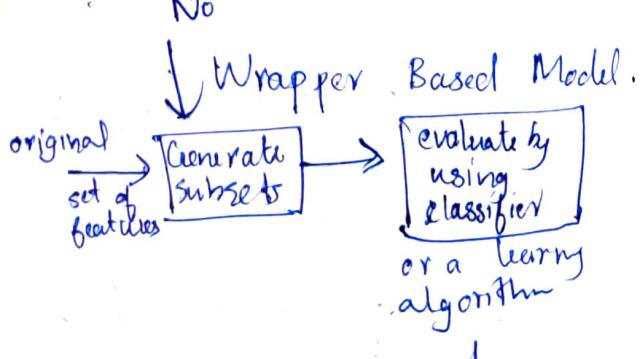
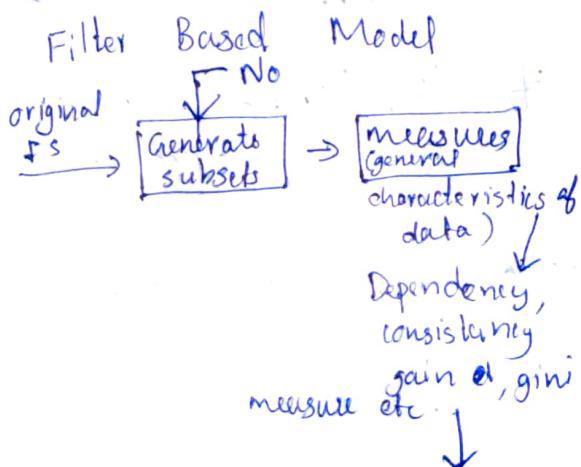
$\rightarrow$  keep the features if they are consistent

(same similar feature values  $\rightarrow$  same class labels).

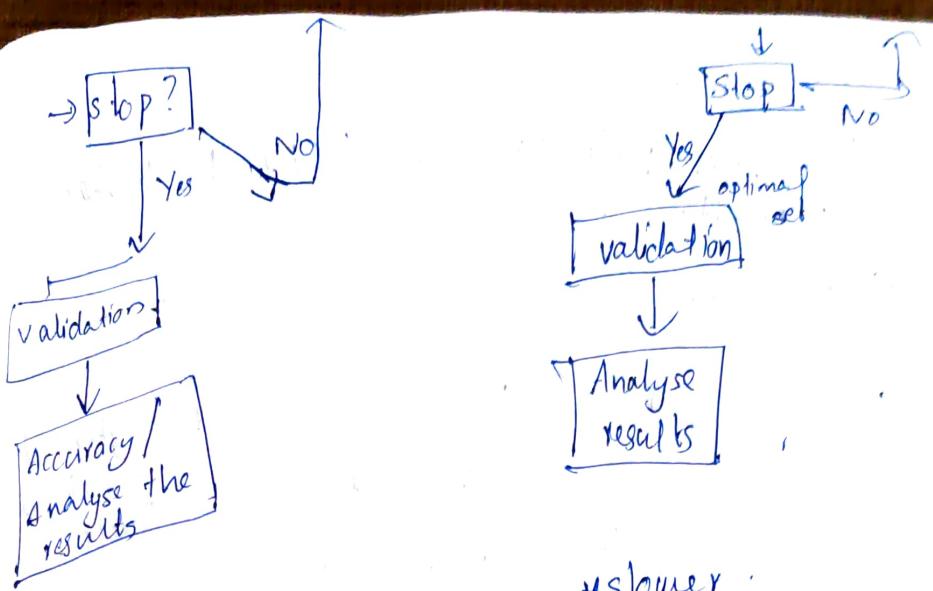
(5) accuracy measure

- eval ~~all~~ <sup>all</sup> diff subsets using accuracy measure.  
keep the <sup>subset</sup> ~~one~~ with highest accuracy.

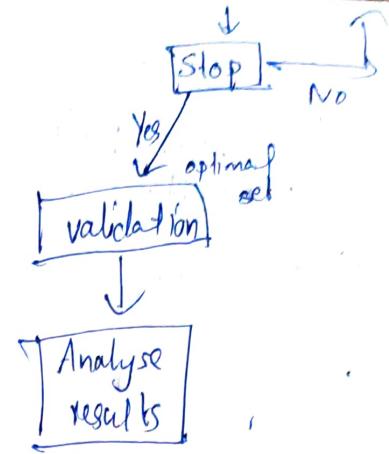
16/3/22  
Models of F.S



evaluate by using classifier or a learning algorithm



\* Faster  
no learning algorithm present



- \* slower training time is required
- \* computationally expensive learning algorithm
- \* higher accuracy / good results, because already all features are evaluated.

→ GA - FS :- wrapper based model.

every chromosome is evaluated by classifier.  
Fitness fn — is the accuracy got by using a learning algorithm.

Data Transformation :-

→ A fn that maps entire set of values of a given attribute to a new set.

Methods :-

└ smoothing - remove noise from data.

└ feature construction - new features constructed from the given ones.

└ aggregation - summarization

eg: sales data aggregation

- Normalization (0 to 1)

- concept hierarchy
- generalisation
    - lower level primitives are mapped to higher level primitives
  - street → city → country
  - concept hierarchy (lower to higher / higher to lower)
    - we can mine at diff. levels
    - lower / higher / mixed

### Normalization

└ min max normalization

$$v' = \frac{v - \min}{\max - \min} (\text{new Max} - \text{new min}) + \text{new min}$$

$v'$  → normalised value

$v$  → attribute value

max, min — of the attribute

new max, new min — scale to which we want to convert the original value

Eg: Income \$12,000 to \$98,000

normalise in range (0-1)

value — \$73,000

$$v' = \frac{73,000 - 12,000}{98,000 - 12,000} (1-0) + 0 = 0.716$$

└ 3-score Normalisation

$$v' = \frac{v - \mu}{\sigma} (\frac{\mu - \text{mean}}{\sigma - \text{s.d}})$$

$$\mu = 54,000 \quad \sigma = 16,000$$

$$v' = \frac{770 - 54,000}{16000}$$

$$= \underline{\underline{1.03}}$$

↳ Decimal Scaling Based Normalization

$$v' = \frac{v}{10^j}$$

$j$  — smallest integer such that  $\max |v'| < 1$

Attribute range from -986 to 917

Max abs value is 986

$$j = 3$$

norm value in the rage of :

-0.986 to 0.917 → normalized value.

21/3/22  
↳ Principal Comp Analysis

↳ dimension reduction

↳ n of dim, select only, p no of dimensions  
 $p < n$

without losing any info from the given data

steps : ① Collect Data

② Subtract mean from data

③ Calculate covariance matrix

④ calculate eigen values & eigen vectors

of the covariance matrix

⑤ Order the eigen vectors by eigen values (highest to lowest) least significant component

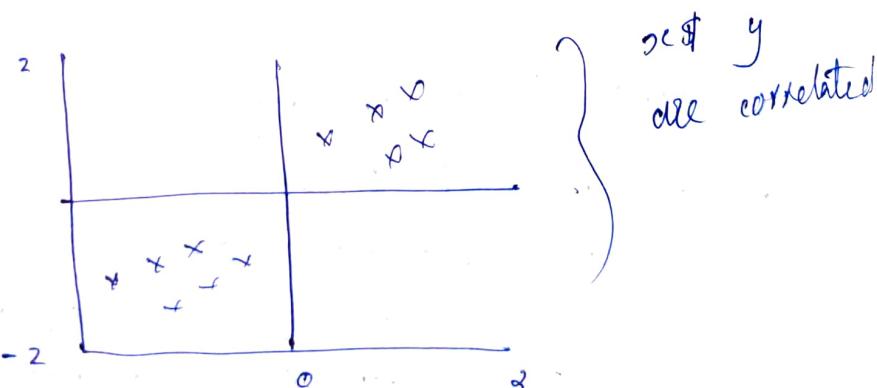
failure vector : see principal component  
 eigen vec 1, eigVec 2, eigVec 3 ... eigVec n  
 PC<sup>1</sup> choose only 'p' no of eigen vectors, where  $p \leq n$ .  
 ('n' is given no of dimensions).  
 'p'  $\rightarrow$  reduced dimension.

### ⑤ Define new data -

[If initial data is correlated, and data will be uncorrelated]

eg:

	x	y	Subtract mean	
	2.5	2.4	$x - \bar{x}$	$y - \bar{y}$
	0.5	0.7	0.69	0.49
	2.2	2.9	-1.31	-1.21
	1.9	2.2	0.39	0.99
	3.1	3.0	0.09	0.29
	2.3	2.7	0.49	0.79
	2	1.6	0.19	-0.31
	1	1.1	-0.81	-0.81
	1.5	1.6	-0.31	0.31
	1.1	0.9	-0.71	-1.01
$\bar{x} =$	1.81	$\bar{y} = 1.91$		



### ⑥ Find covariance matrix ( $2 \times 2$ )

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(\bar{y}_i - \bar{y})}{(n-1)}$$

$$C = \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix}$$

$$\text{cov}(x, x) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n-1)}$$

$$C = \begin{bmatrix} 0.616155 & 0.61544 \\ 0.61544 & 0.71655 \end{bmatrix}$$

Find eigen vectors & eigen values for C:  
 by:

$$\begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

$$\begin{bmatrix} 12 \\ 8 \end{bmatrix} \approx 4 \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

$$\mu \lambda = \sqrt{9}$$

$\mu \rightarrow n \times n$  matrix.

$\lambda \rightarrow$  non zero  $n \times 1$  vector.

$\nu \rightarrow$  eigen value.

$\vec{v} \rightarrow$  eigen vector

$$\textcircled{1} \text{ eigen values} = \begin{pmatrix} 0.0490833 \\ 1.28402771 \end{pmatrix}$$

( $n$  eigen vectors &  $n$  eigen values)  
 $\hookrightarrow$  dimension of data.

$$\text{eigen vectors} = \begin{bmatrix} -0.7351 & -0.6778 \\ 0.6778 & -0.7351 \end{bmatrix}$$

\textcircled{2} arrange eigen vectors based on eigen values

$$\begin{bmatrix} -0.6778 & -0.7351 \\ -0.7351 & 0.6778 \end{bmatrix}$$

(PC1)

(PC2)

⑥ Derive new Data (P=1)

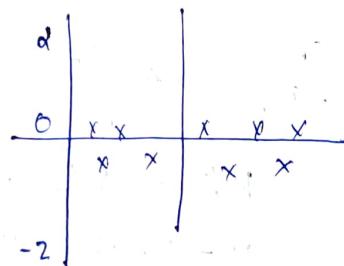
$$PC_1 = -0.6778, \quad -0.7351$$

$$-0.6778 \times 0.69 + 0 \cdot -0.7351 \times 0.49 = -0.827$$

$$-0.6778 \times (-1.31) + (-0.7351) \times (-1.21) = 1.7$$

$$\text{modified } x = -0.827, 1.7, -0.992, -0.274, \\ -1.6, 0.91, 0.099, 1.14, 0.43, 1.29$$

$$\text{Modified } X \\ (PC_2) = -0.17, 0.14, 0.38, 0.13, 0.20, \\ 0.17, -0.34, -0.04, 0.01, \\ -0.16$$



→ now the dimensions are independent.

Best sig component of data, without losing any information.

### Clustering

- Group the similar objects
- unsupervised (class label 'x')

Partition Based

Hierarchical based

- k-mean, k medoid

→ agglomerative and divisive.

k - no of clusters

\* every obj in a single cluster.

- k should be known

\* merge

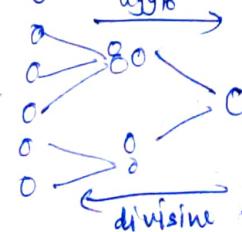
- every obj belongs to only one cluster.

\* till every obj in a single cluster.

- choose k no of cluster center (x)

aggl

- Find the dist of every pt with given cluster center.



divisive

- assign the pt to nearest centroid
- repeat until no change in the centroids

Adv  
→ easy / simple

Dis  
→ very sensitive to noise

→  $O(n k t)$  no of iterations  
no of data pts no of clusters

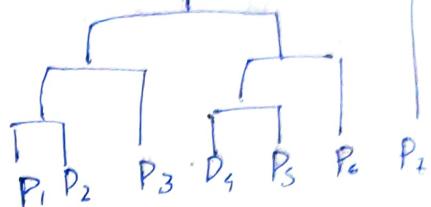
→ hard clustering

\* FCM (soft clustering).  
→ obj can belong to more than one cluster with membership matrix.

→ K value not required  
→  $O(n^2)$

↳ no of data pts

→ Dendograms



Disadv

for high dimensional data, height of dendrogram is going to increase.

→ cut the tree at some optimal level to get optimum no of clusters (using some optimisation techniques).

### Density Based

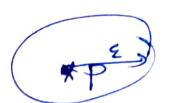
Density → no of pts with given radius ( $\epsilon$ ) of cluster

→ identify the dense and sparse region

→ DBSCAN and OPTICS

→ min pts &  $\epsilon$

\* core pts at  $\epsilon$  rad, has min pts



P is core pt, if it covers min pts within  $\epsilon$  radius

→ core pt, border pt & noise pt } entire dataset will be divided into 3

its not core pt, but lies within the neighbourhood of core pt. not core pt nor border pt. DISCARD pt.

→ if a pt is a core pt, make a cluster  
if a border pt, add it to the nearest core pt.  
discards the noise pt.

Adv

- not sensitive to noise
- any shape of cluster is possible

Dis

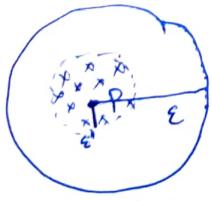
- have to choose  $\epsilon$  & min pts - diff for high dimensional data  
(optimisation techniques should be used)

## OPTICS (Ordering Points To Identify Clustering Strategy)

- need to identify
  - core distance
  - reachability distance

$$-\min p \in S$$

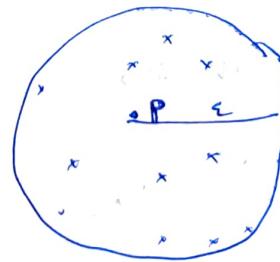
$$\epsilon = 6 \text{ mm}$$



$$\epsilon' < \epsilon$$

within ' $\epsilon'$  dist, 'P' is a core pt.

$$\epsilon' \rightarrow \text{core distance}$$

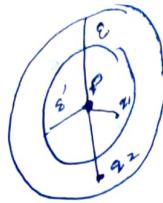


$\epsilon$  is used by 'P' to become core pt.

→ core distance of an obj 'P' is the smallest ' $\epsilon$ ' value that makes 'P' as core pt

→ The reachability dist of an obj 'q' wrt another object 'p' is the greater of the value of the core dist of

$P$  and euclidean dist between  $P \& q$



$P \rightarrow$  core pt with core dist  $\leq r$

$\text{Reach}(P, q_1) = r_1$

$\text{Reach}(P, q_2) = \text{euclidean dist}(P, q_2)$

Optics — It stores core dist & reachability dist for each obj.

— optics method maintains a list called seed list in which objects are sorted by reachability dist from their respective closest core objects

## 28/3/22 Optics

### core point

→ core distance

→ reachability distance

→ optics computes an ordering of all objects in given db

→ It stores core dist and reachability dist of each object

→ It maintains a seed list

— objects are sorted by reachability distance from the respective closest core objects

— points are processed from smallest to highest reachability dist

### Steps

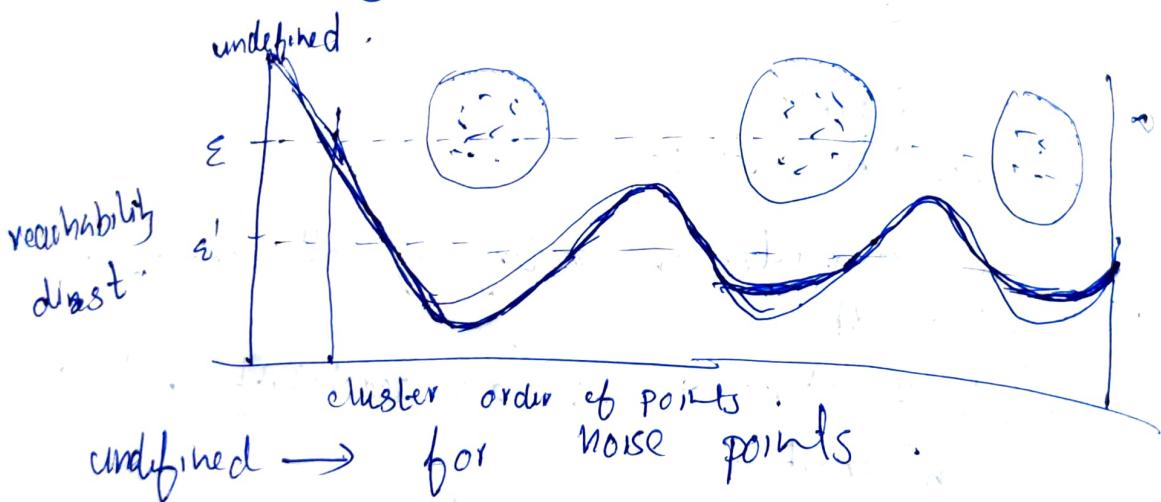
① Begin with point  $P$

② within  $\epsilon$  radius  $P$  retrieves all objs, determines core dist & reachability dist

⑥ O/P P

- \* If p is core pt,
    - then For each obj q in the ε dist. of p, update its reachability dist.
    - insert q into list if its not processed
- repeat until the I/P is fully consumed and list is empty.

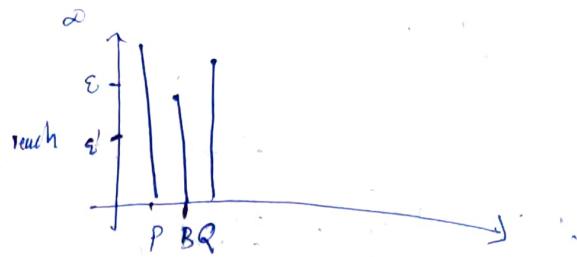
- \* If p is not core pt
  - move on to the next object in the list or move on to the I/P ~~data.s. db~~ if list is empty



→ lower bumps, the pts will have very less reachability dist, & will be connected to a dense cluster

→ pts are processed on a specific order & this order selects an object q that is densely reachable with lowest ε' value so that clusters with high density will be processed first

ag:



list: [ B 30 C 40 ]

B  
[ P Q R S ]

① list starting from P

[ B 30 ] C 40

② let densely reachable pts of B be q, r, s

[ S 41 ] P 39 [ R 42 ] C 40

→ only

[ q 39 ] C 40 [ S 41 ] R 42

output pt: q

③ q is. reachable to C & R

[ C 40 ] S 41 [ R 42 ] C 39 [ Z 50 ]

[ C 39 ] S 41 [ R 42 ] Z 50

continued until list is empty or db is completely processed

→ pts processed in a manner that highly clustered clusters processed first

→ clusters can be visualized after processing the pts

30/3/22

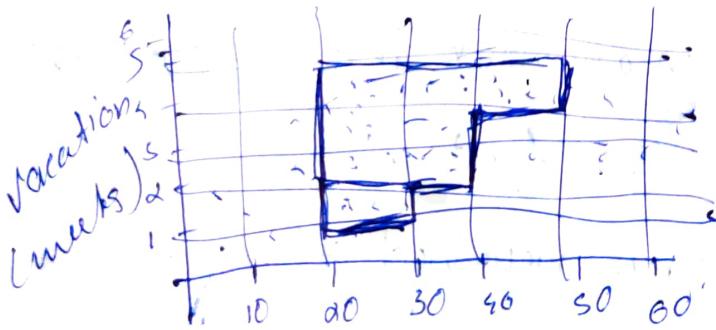
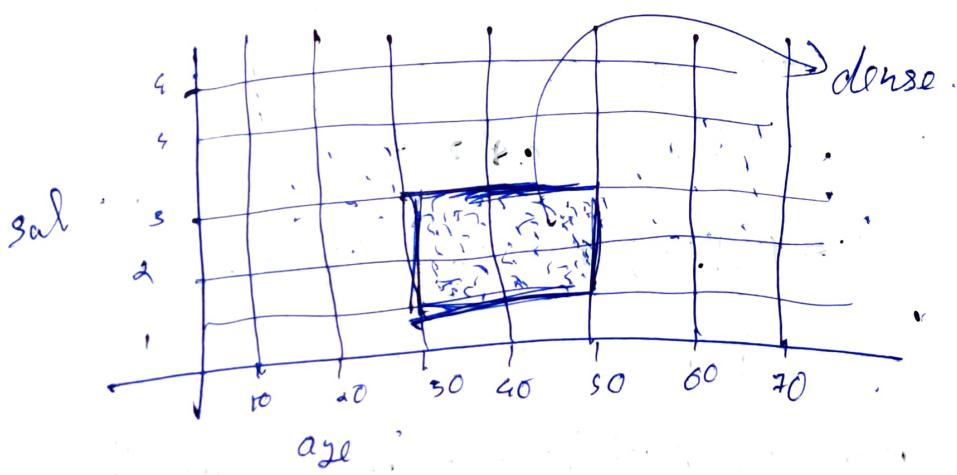
## CLIQUE

subspace clustering

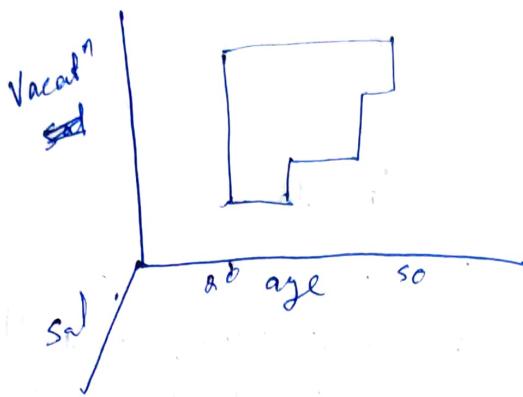
- First alg method for dimension growth of subspace clustering.
- It starts with a single dimension, dives
- It identifies dense units (no of minpts).
- Subspace representing these dense units are intersected to form candidate search for higher dimension.

Maximal region - highly dense units.

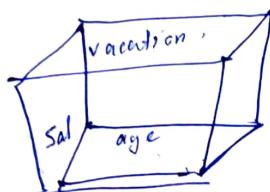
Minimal region - low dense region.



- intersecting (dimensions and dense regions)



- cube (we can view the data in the form of cube after intersection) ; at the end



- density based clustering

Drawback -

- deciding grid size
  - Min pts .
- } vary the parameters and find the optimal ones .

Classifiers -

- NB, DT, KNN, ANN  
| gain, gini  
probability based classifier

| backpropagation

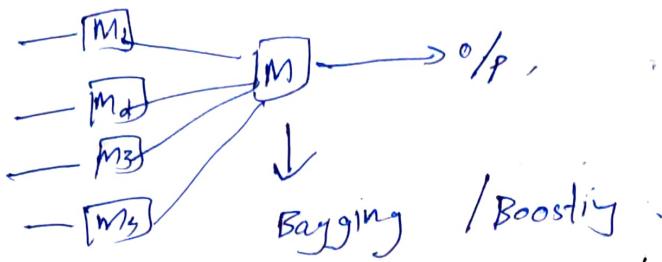
perceptron model

Ensemble Classifiers

→ set of classifiers

→ the learning is based on multiple classifiers ;  
not on a single classifier .

single  
classifiers -



→ Bagging → each classifier returns ~~and~~ the class prediction for 'x' which counts as a vote.

x is assigned to the class with most votes -

$M_1 \xrightarrow{\text{X}} \text{yes } ①$

$M_2 \longrightarrow \text{yes } ②$

$M_3 \longrightarrow \text{no } ③$

$M_4 \longrightarrow \text{yes } ④$

$x \longrightarrow \text{Yes}$

→ Boosting → assign wt to each classifier not based on how well the ~~stare~~ classifier performed

→ sum the wts of each classifier that assigned class 'c' to 'x'

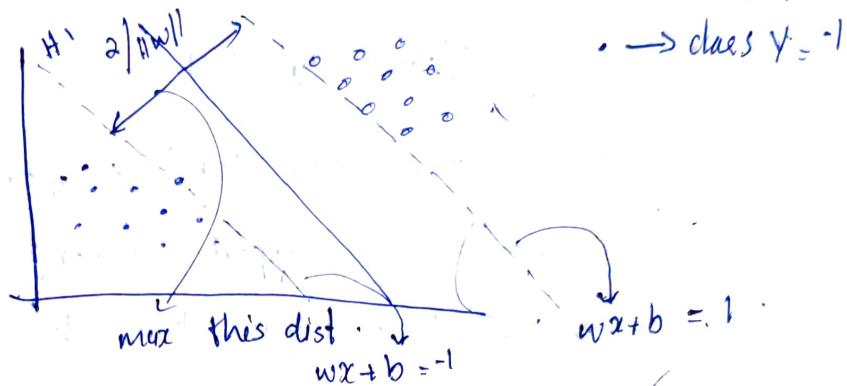
→ The class with highest sum is winner

→ prior knowledge of classifiers are req.  
as weights should be assigned

[Eg: Naive Bayes — should have independent attributes. If you have dependent attributes, assign less weight to Naive Bayes]

## SVM (Support Vector Machine)

→ Find the maximum margin hyperplane (MMH) that separates the given classes:



Any hyperplane can be written as  $wx + b = 0$   
w is normal vector to the hyperplane,

$w_1, w_2, \dots, w_n$ , n is the no of attributes in the dataset.

x — is the training tuple

b — bias

→ Any pt that falls above hyperplane  
 $wx + b > 0$

→ any pts that falls below hyperplane  
 $wx + b < 0$

→ By geometry,  
dist betwn  $H_1 \& H_2 = 2/\|w\|$ .

Maximise the distance between  $H_1 \& H_2$ .

— this will well separate the classes.

→ Minimise  $w$ .

## Support Vectors

→ The pts that lie on  $H_1$  &  $H_2$  are called as support vectors.

## Finding Class Labels

$$d(x_i) = \sum_{i=1}^k y_i \cdot l_i \cdot \text{dot}(x_i + b)$$

$y_i$  → is the class label of support vector  $x_i$ .

$k$  → no of support vectors.

$x_i$  → best tuple to which class label has to be fit.

$l_i, b$  → are constant values

numeric parameters ~~are~~ defined

by user.

$d(x_i) > 0 \rightarrow +ve$  class.

$d(x_i) < 0 \rightarrow -ve$  class.

SVM - non linearly separable data (can be used for)

- Kernel functions are used for non linear separation, data vectors are mapped to high dimensional feature space, where those can be separated in a linear fashion.

↳ polynomial, sigmoid, radial -