



# Chapter 12: Data Warehousing and Mining





# Outline

- Decision Support Systems
- Data Warehousing
- Data Mining
- Classification
- Association Rules
- Clustering



# Decision Support Systems

- **Decision-support systems** are used to make business decisions, often based on data collected by on-line transaction-processing systems.
- Examples of business decisions:
  - What items to stock?
  - What insurance premium to change?
  - To whom to send advertisements?
- Examples of data used for making decisions
  - Retail sales transaction details
  - Customer profiles (income, age, gender, etc.)



# Decision-Support Systems: Overview

- **Data analysis** tasks are simplified by specialized tools and SQL extensions
  - Example tasks
    - ▶ For each product category and each region, what were the total sales in the last quarter and how do they compare with the same quarter last year
    - ▶ As above, for each product category and each customer category
- **Statistical analysis** packages can be interfaced with databases
  - Statistical analysis is a large field, but not covered here
- **Data mining** seeks to discover knowledge automatically in the form of statistical rules and patterns from large databases.
- A **data warehouse** archives information gathered from multiple sources, and stores it under a unified schema, at a single site.
  - Important for large businesses that generate data from multiple divisions, possibly at multiple sites
  - Data may also be purchased externally



# RDBMS used for OLTP

- Database Systems have been used traditionally for **OLTP**  
**(online transaction processing)**
  - clerical data processing tasks
  - detailed, up to date data
  - structured repetitive tasks
  - read/update a few records
  - isolation, recovery and integrity are critical

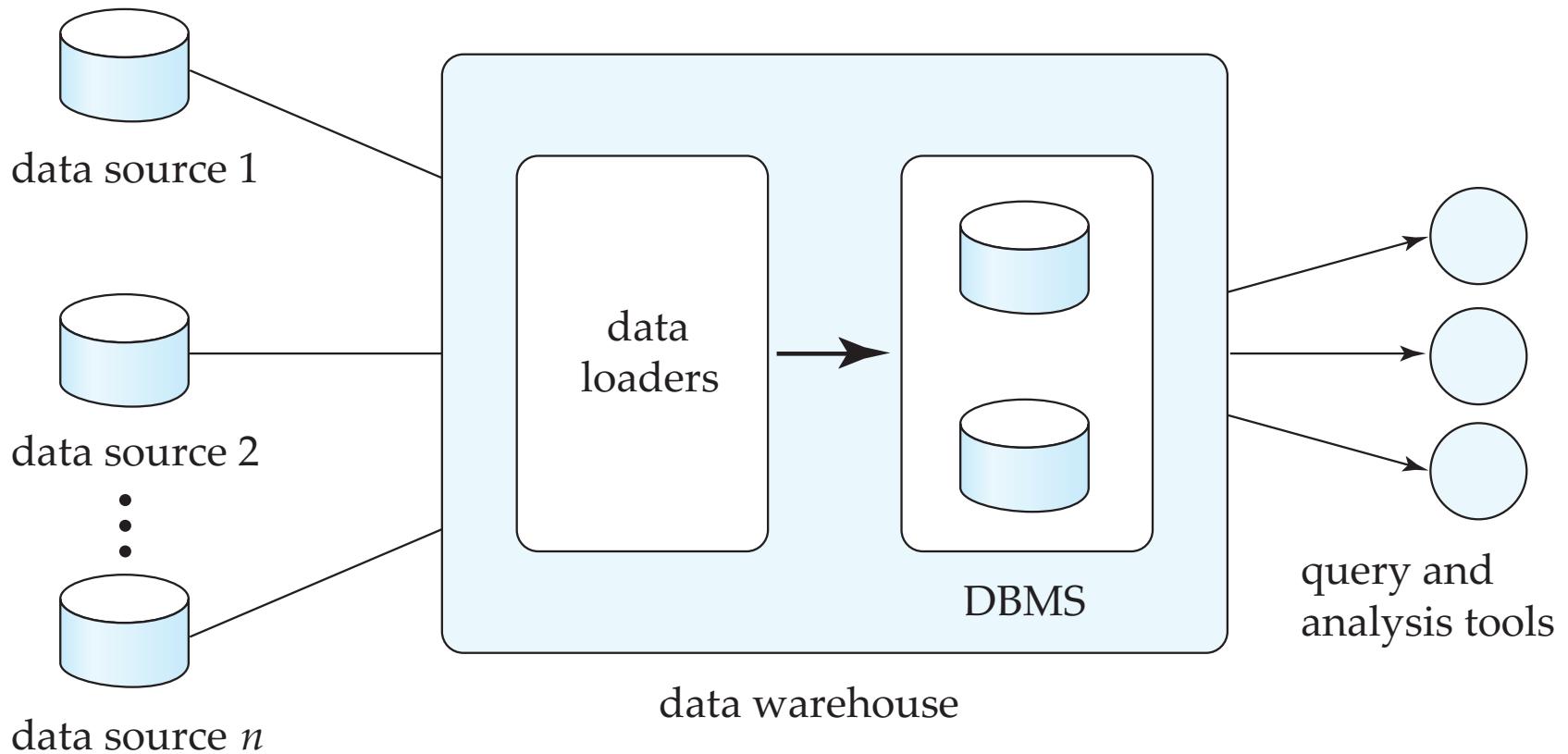


# Data Warehousing

- Data sources often store only current data, not historical data
- Corporate decision making requires a unified view of all organizational data, including **historical data**
- A **data warehouse** is a repository (archive) of information gathered from multiple sources, stored under a unified schema, at a single site
  - Greatly simplifies querying, permits study of historical trends
  - Shifts decision support query load away from transaction processing systems



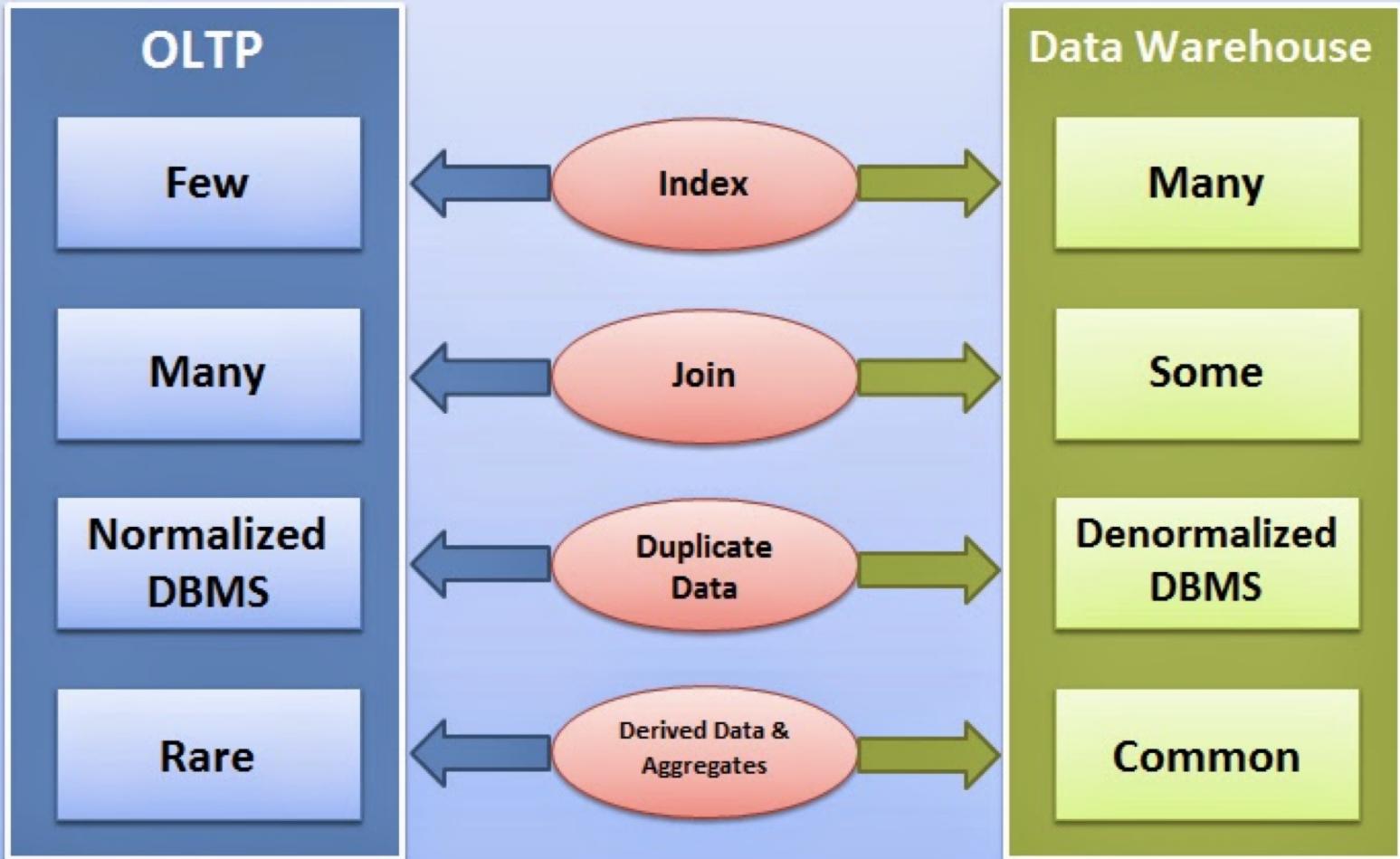
# Data Warehousing





# OLTP (DB) vs Data Warehouse

## OLTP Vs Data Warehouse





# Warehouse Schemas

- Dimension values are usually encoded using small integers and mapped to full values via dimension tables
- Resultant schema is called a **star schema**
  - More complicated schema structures
    - ▶ **Snowflake schema**: multiple levels of dimension tables
    - ▶ **Constellation**: multiple fact tables



# Schema Design

- Database organization
  - Must look like business
  - Must be recognizable by business user
  - Approachable by business user
  - Must be simple
- Schema types
  - Star schema
  - Fact constellation schema
  - Snowflake schema



# Dimension Tables

## ■ Dimension tables

- Define business in terms already familiar to users
- **Wide rows** with lots of descriptive text
- **Small tables** (about a million rows)
- Joined to fact table by a foreign key
- Heavily indexed
- Typical dimensions
  - ▶ time periods, geographic region (markets, cities), products, customers, salesperson, etc.



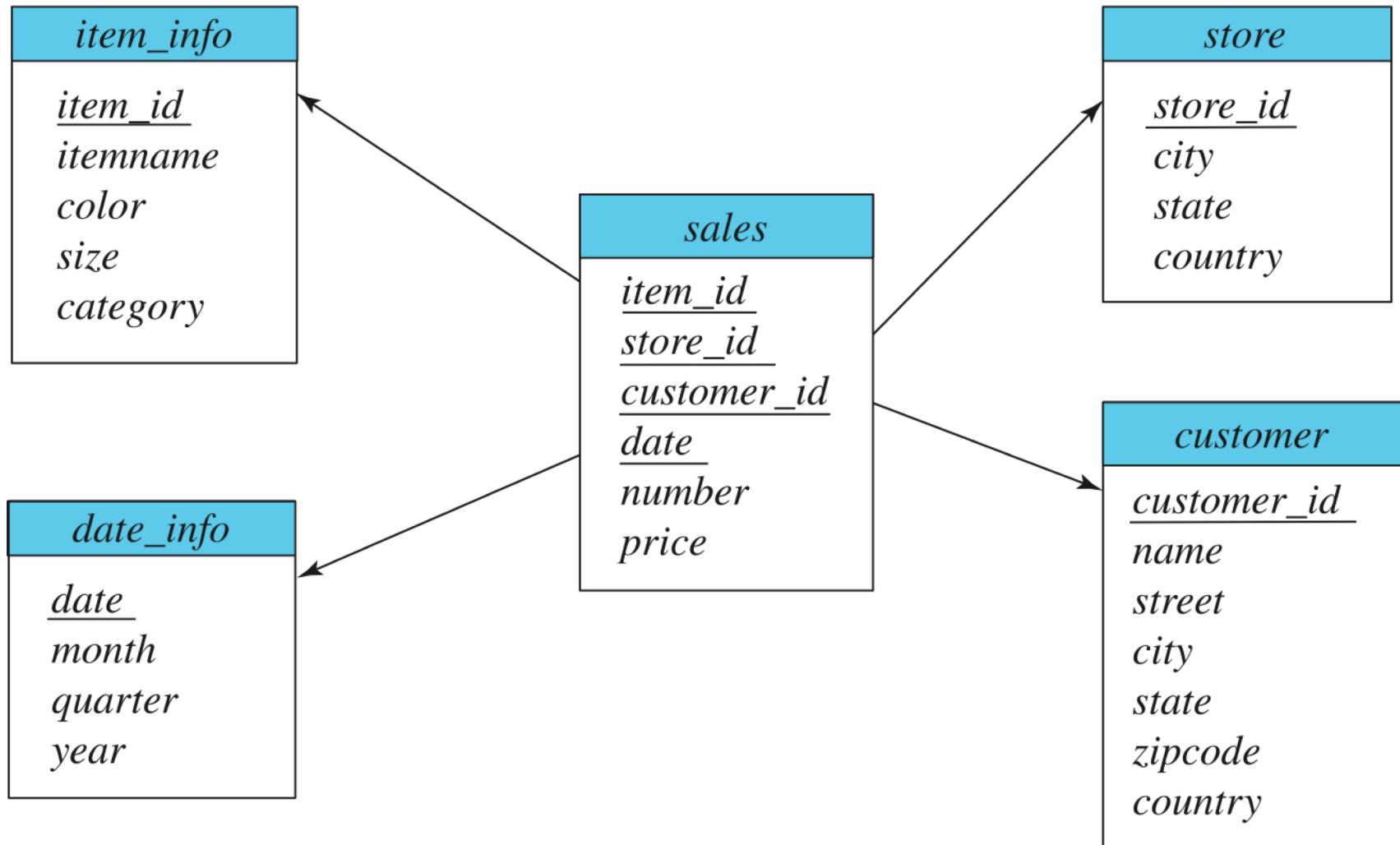
# Fact Table

## ■ Central table

- Mostly raw numeric items
- **Narrow rows**, a few columns at most
- **Large number of rows** (millions to a billion)
- Access via dimensions



# Data Warehouse Schema

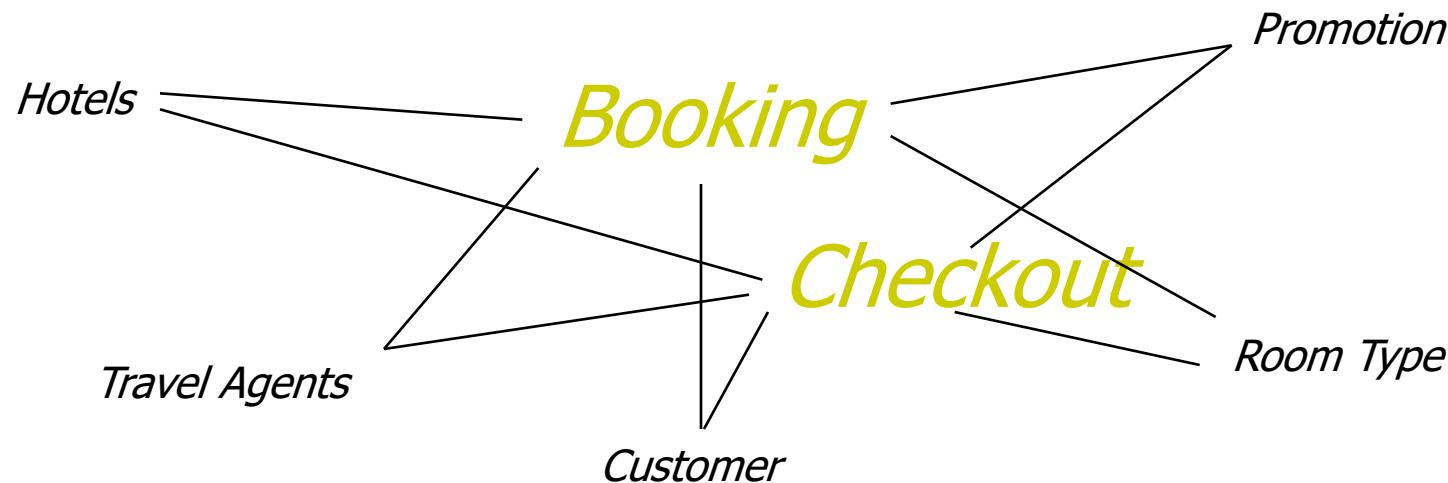




# Fact Constellation

## Fact Constellation

- Multiple fact tables that share many dimension tables
- Booking** and **Checkout** may share many dimension tables in the hotel industry





# De-normalization

- Normalization in a data warehouse may lead to lots of small tables
- Can lead to excessive I/O's since many tables have to be accessed
- De-normalization is the answer especially since **updates are rare**



# Data Mining

- Data mining is the process of semi-automatically analyzing large databases to find useful patterns
- **Prediction** based on past history
  - Predict if a credit card applicant poses a good credit risk, based on some attributes (income, job type, age, ...) and past history
  - Predict if a pattern of phone calling card usage is likely to be fraudulent
- Some examples of prediction mechanisms:
  - **Classification**
    - ▶ Given a new item whose class is unknown, predict to which class it belongs
  - **Regression**
    - ▶ Given a set of mappings for an unknown function, predict the function result for a new parameter value



# Data Mining (Cont.)

## ■ Descriptive Patterns

- **Associations**

- ▶ Find books that are often bought by “similar” customers. If a new such customer buys one such book, suggest the others too.
- Associations may be used as a first step in detecting **causation**
  - ▶ E.g., association between exposure to chemical X and cancer,

- **Clusters**

- ▶ E.g., typhoid cases were clustered in an area surrounding a contaminated well
- ▶ Detection of clusters remains important in detecting epidemics



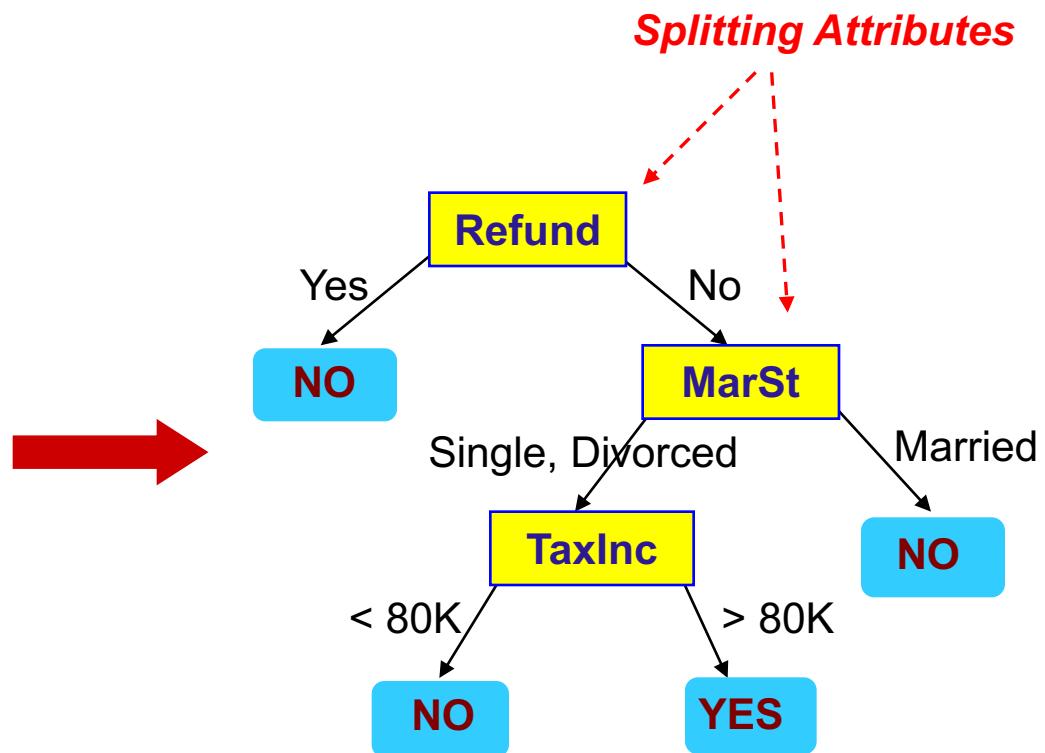
# Classification Rules

- Classification rules help assign new objects to classes.
  - E.g., given a new automobile insurance applicant, should he or she be classified as low risk, medium risk or high risk?
- Classification rules for above example could use a variety of data, such as educational level, salary, age, etc.
  - $\forall$  person P, P.degree = masters **and** P.income > 75,000  
 $\Rightarrow$  P.credit = excellent
  - $\forall$  person P, P.degree = bachelors **and**  
 $(P.\text{income} \geq 25,000 \text{ and } P.\text{income} \leq 75,000)$   
 $\Rightarrow$  P.credit = good
- Rules are not necessarily exact: there may be some misclassifications
- Classification rules can be shown compactly as a decision tree.



# Example of a Decision Tree

Tid	Categorical				Continuous class
	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	



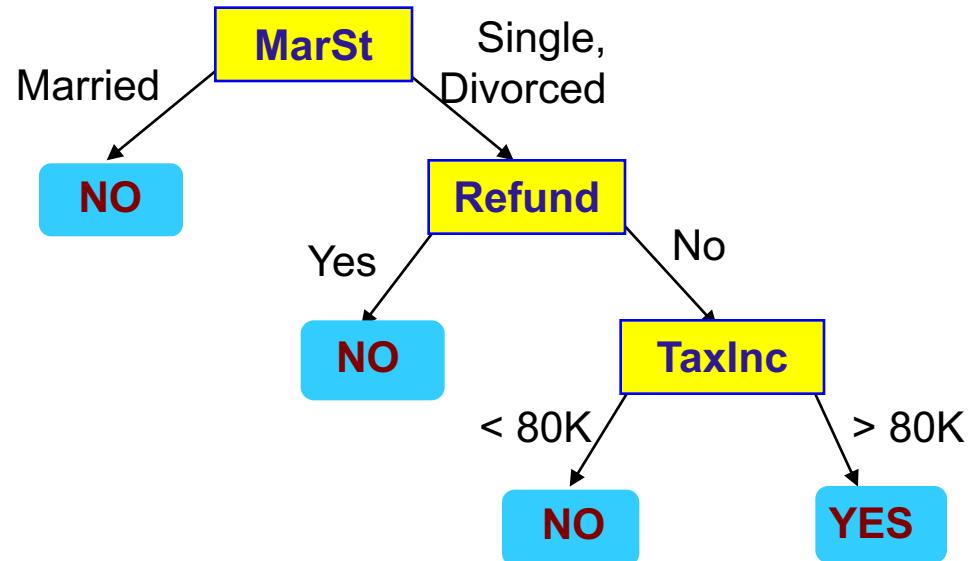
Training Data

Model: Decision Tree



# Another Example of Decision Tree

Tid	Refund	Marital Status	Taxable Income	Cheat	categorical categorical continuous class
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	



There could be more than one tree that fits the same data!



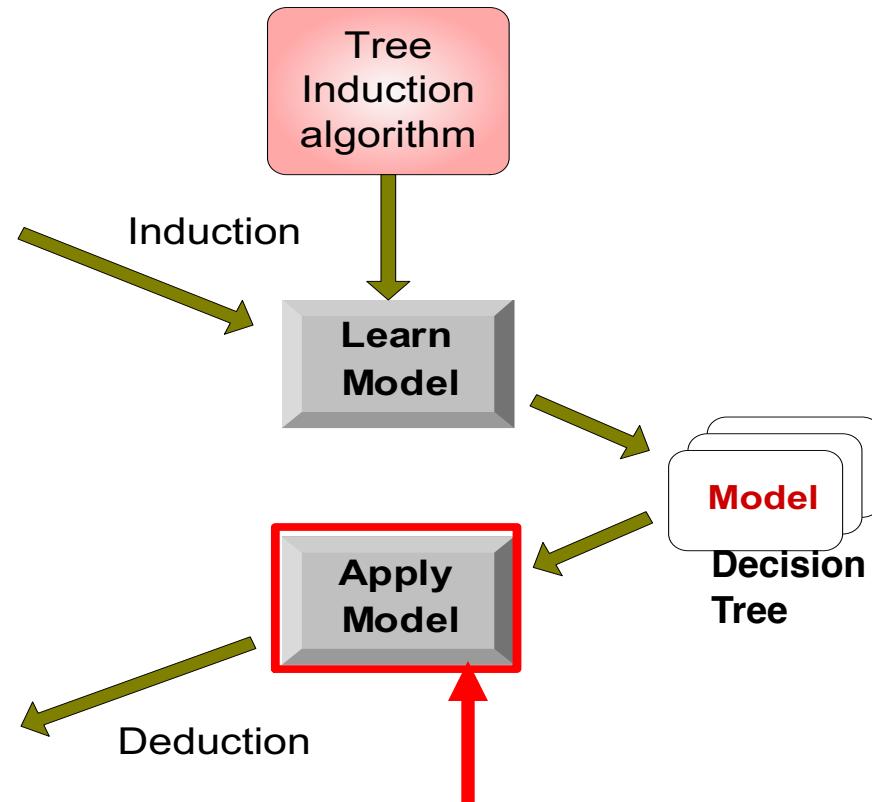
# Decision Tree Classification Task

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

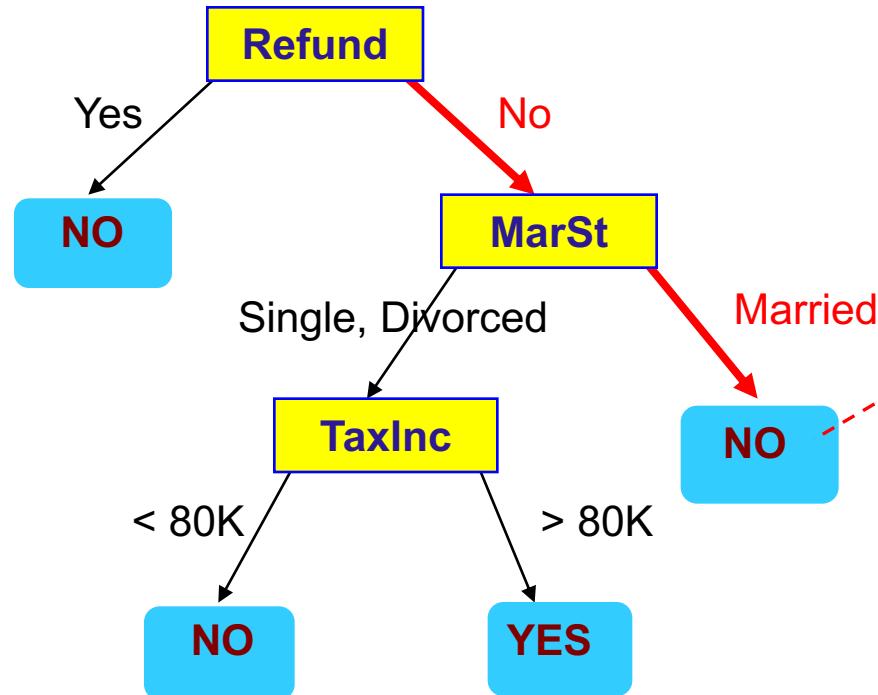
Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set





# Apply Model to Test Data



## Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

Assign Cheat to “No”



# Example

- Example: Play Tennis

*PlayTennis: training examples*

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



# Association Rules

- Retail shops are often interested in associations between different items that people buy.
  - Someone who buys **bread** is quite likely also to buy **milk**
  - A person who bought the book *Database System Concepts* is quite likely also to buy the book *Operating System Concepts*.
- Associations information can be used in several ways.
  - E.g., when a customer buys a particular book, an online shop may suggest associated books.
- **Association rules:**

$bread \Rightarrow milk$        $DB\text{-}Concepts, OS\text{-}Concepts \Rightarrow Networks$

- Left hand side: **antecedent**, right hand side: **consequent**
- An association rule must have an associated **population**; the population consists of a set of **instances**
  - ▶ E.g., each transaction (sale) at a shop is an instance, and the set of all transactions is the population



# Other Types of Associations

- Basic association rules have several limitations
- Deviations from the expected probability are more interesting
  - E.g., if many people purchase bread, and many people purchase cereal, quite a few would be expected to purchase both
  - We are interested in **positive** as well as **negative correlations** between sets of items
    - ▶ Positive correlation: co-occurrence is higher than predicted
    - ▶ Negative correlation: co-occurrence is lower than predicted
- Sequence associations / correlations
  - E.g., whenever bonds go up, stock prices go down in 2 days
- Deviations from temporal patterns
  - E.g., deviation from a steady growth
  - E.g., sales of winter wear go down in summer
    - ▶ not surprising, part of a known pattern.
    - ▶ look for deviation from value predicted using past patterns

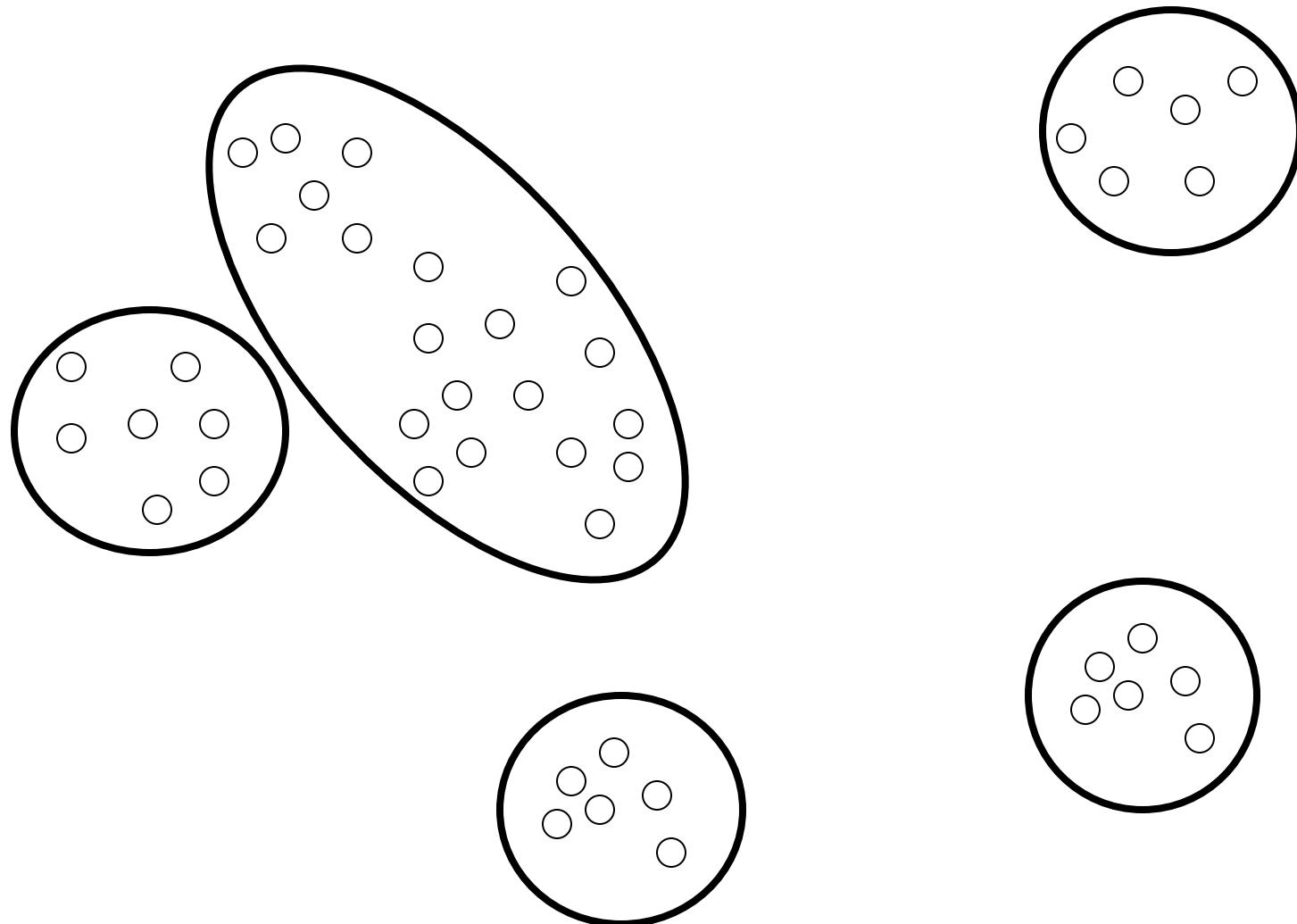


# Clustering

- Clustering: Intuitively, finding clusters of points in the given data such that **similar points** lie in the same cluster
- Can be formalized using distance metrics in several ways
  - Group points into  $k$  sets (for a given  $k$ ) such that the average distance of points from the centroid of their assigned group is minimized
    - ▶ Centroid: point defined by taking average of coordinates in each dimension.
  - Another metric: minimize average distance between every pair of points in a cluster
- Has been studied extensively in statistics, but on small data sets
  - Data mining systems aim at clustering techniques that can handle **very large data sets**



# Clustering





# Clustering: Navigation of search keywords

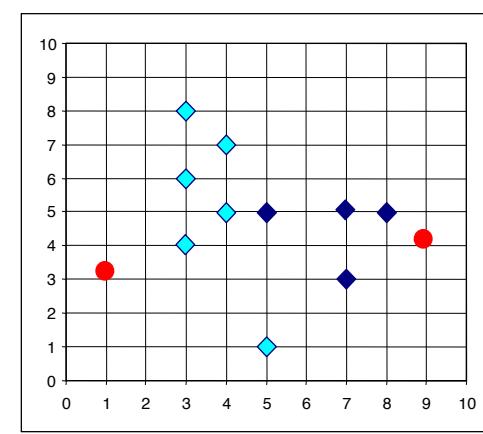
- For grouping search keywords for query suggestion





# The K-Means Clustering

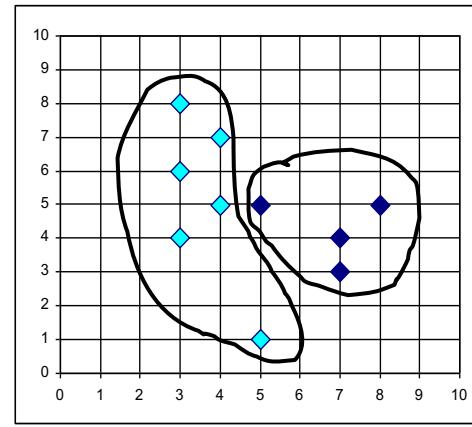
## ■ Example



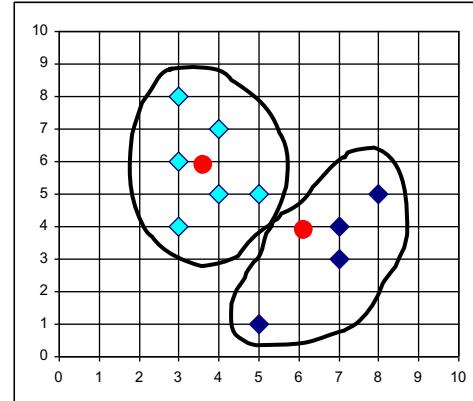
K=2

Arbitrarily choose K object as initial cluster center

Assign each objects to most similar center

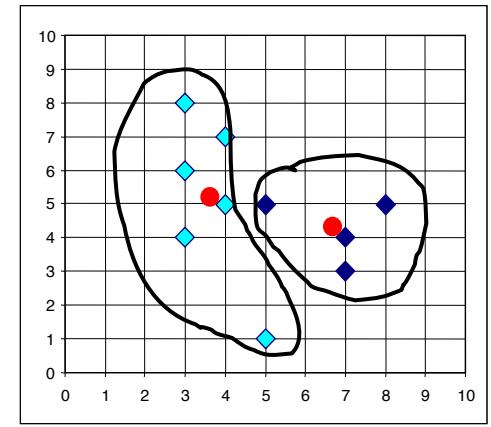


reassign

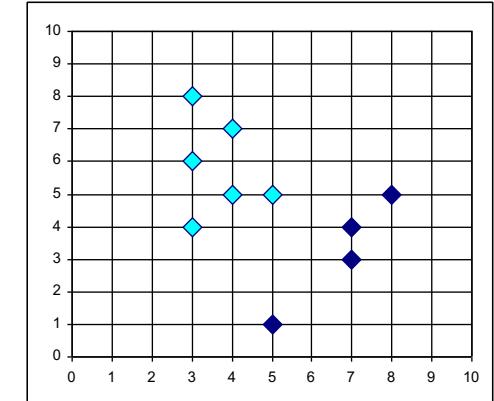


reassign

Update the cluster means



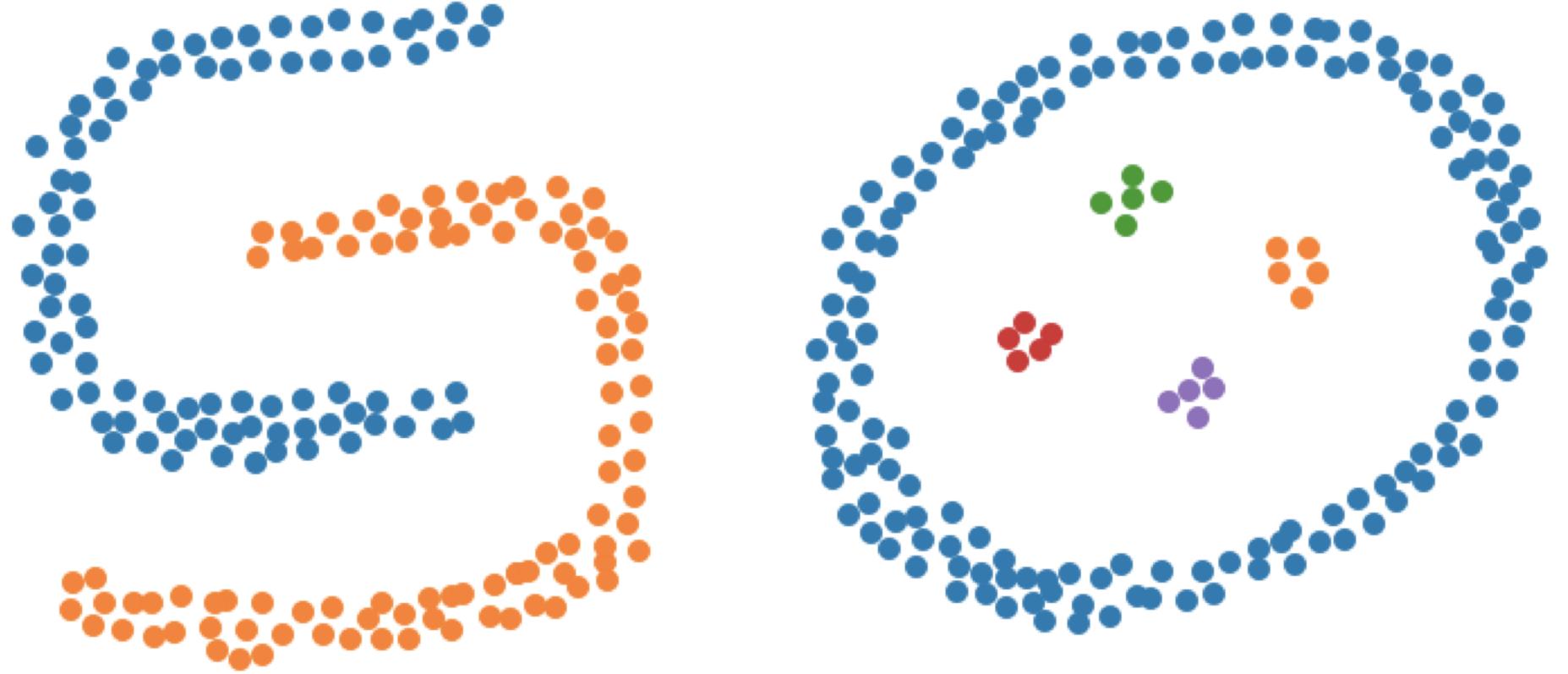
Update the cluster means



reassign



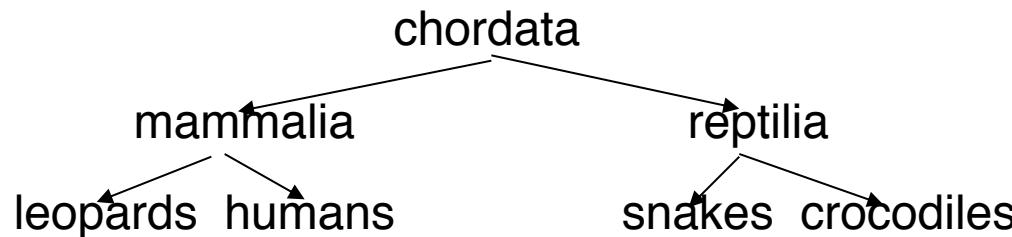
# DBSCAN Algorithm





# Hierarchical Clustering

- Example from biological classification
  - (the word classification here does not mean a prediction mechanism)

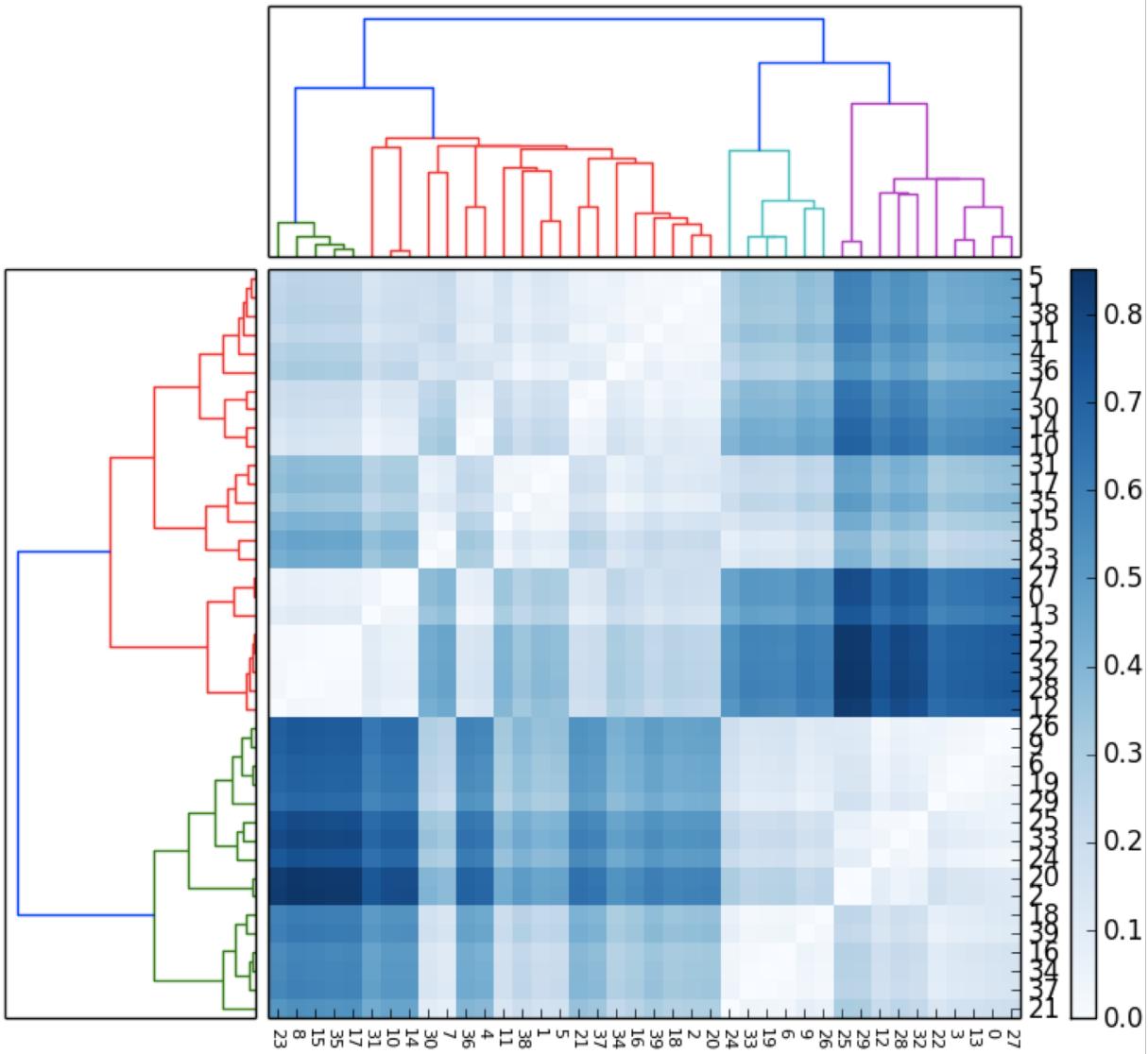


- Other examples: Internet directory systems (e.g., Yahoo, more on this later)
- **Agglomerative clustering algorithms**

- Build small clusters, then cluster small clusters into bigger clusters, and so on

## ■ **Divisive clustering algorithms**

- Start with all items in a single cluster, repeatedly refine (break) clusters into smaller ones





# End of Chapter