



# Other Topics: Warehousing, Mining and Information Retrieval

**Database System Concepts, 6<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# 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 (e.g., : S++) 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

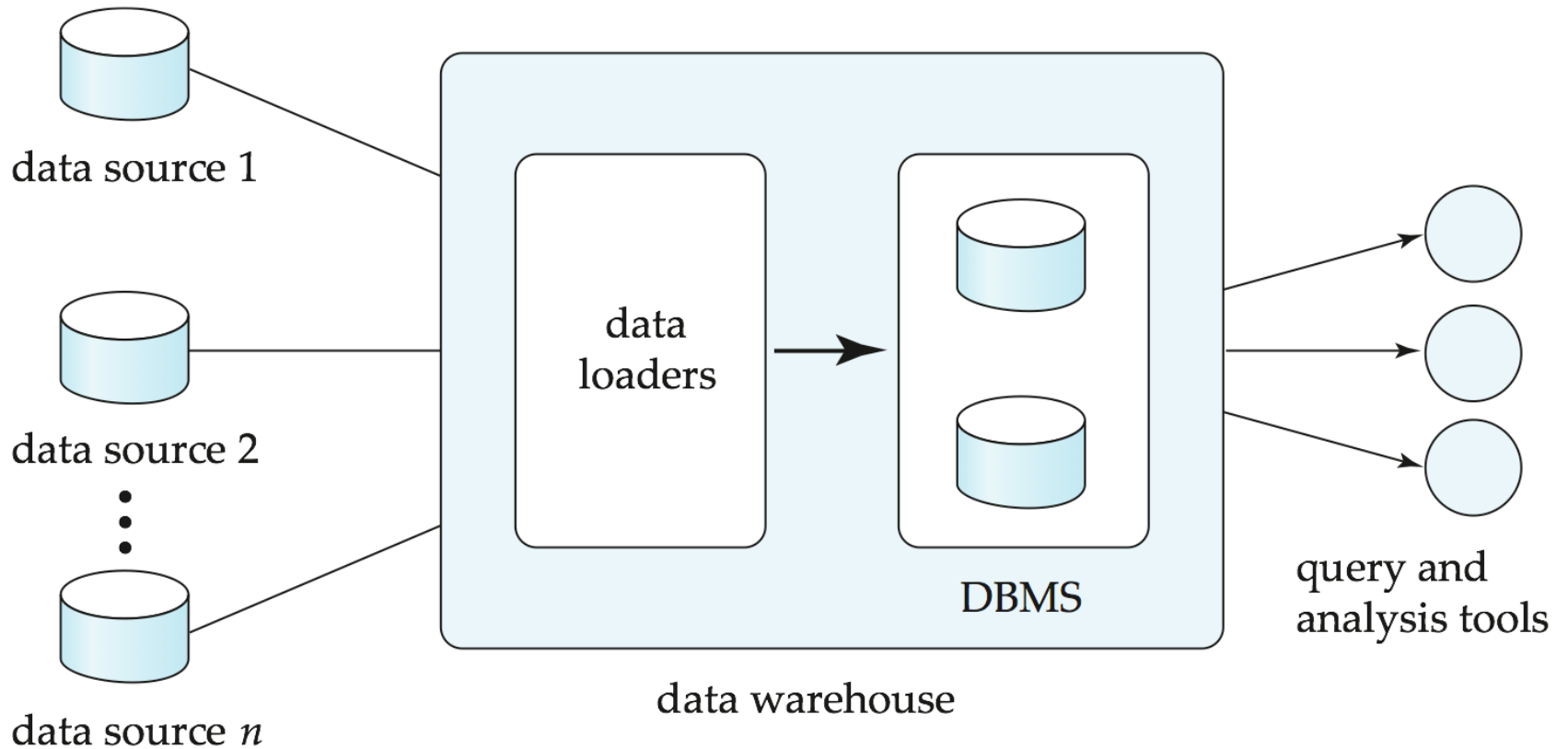


# Data Warehousing

- 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



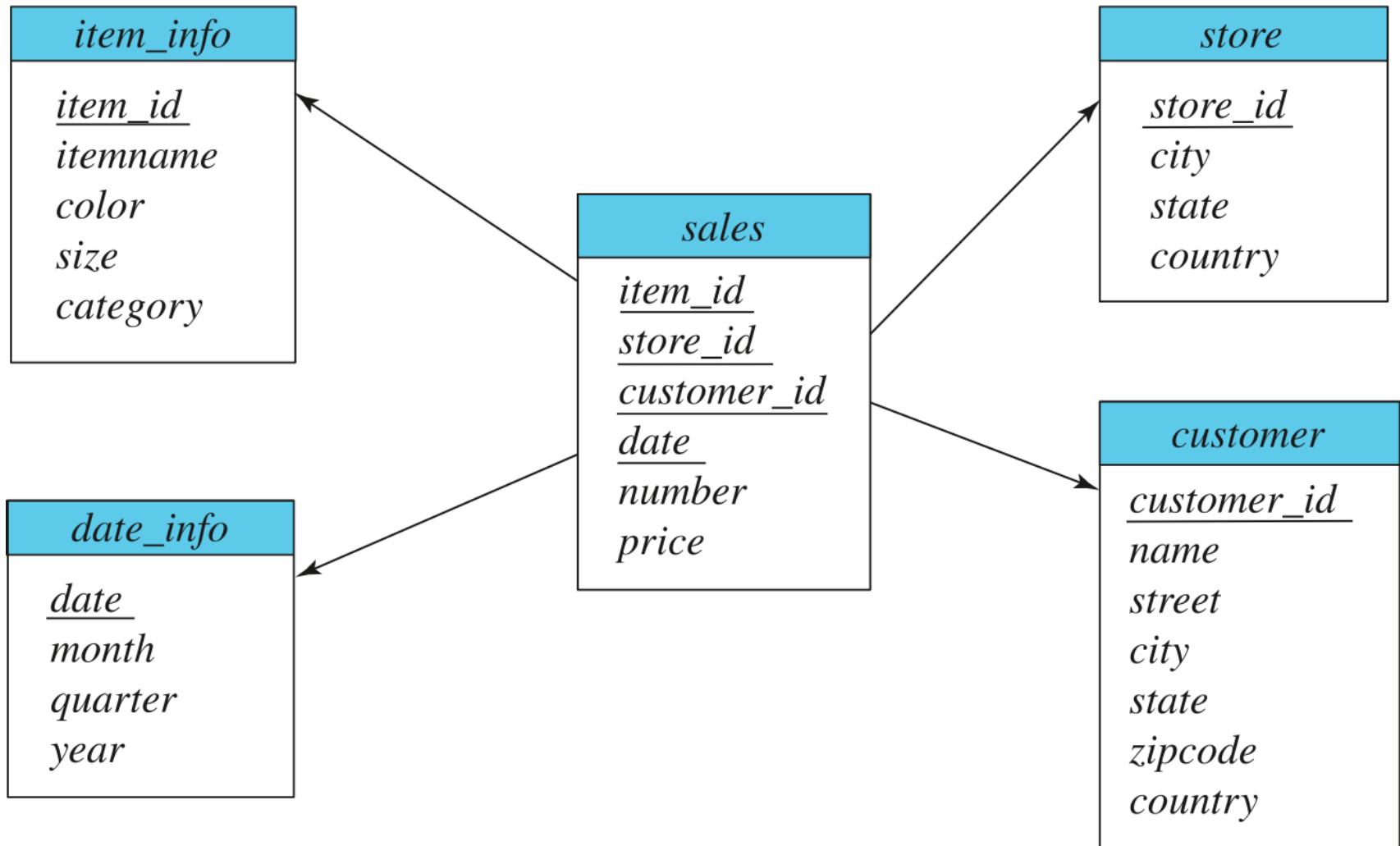


# Warehouse Schemas

- 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



# Data Warehouse Schema





# Data Mining





# 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** formulae
    - ▶ 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

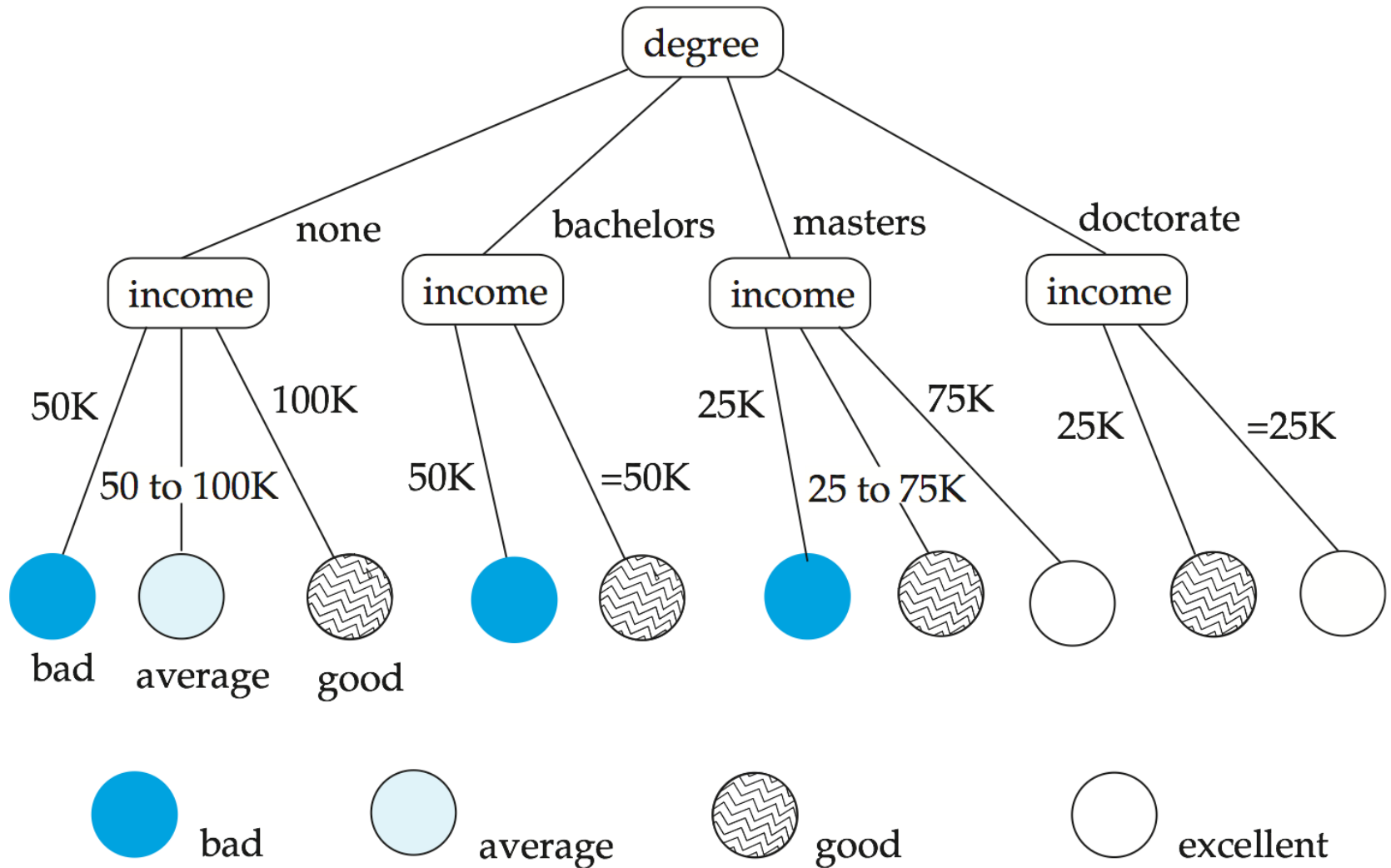


# Classification Rules

- 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  $P$ ,  $P.\text{degree} = \text{masters}$  **and**  $P.\text{income} > 75,000$   
 $\Rightarrow P.\text{credit} = \text{excellent}$
  - $\forall$  person  $P$ ,  $P.\text{degree} = \text{bachelors}$  **and**  
 $(P.\text{income} \geq 25,000 \text{ and } P.\text{income} \leq 75,000)$   
 $\Rightarrow P.\text{credit} = \text{good}$
- Rules are not necessarily exact: there may be some misclassifications
- Classification rules can be shown compactly as a decision tree.
- Several algorithms for constructing decision trees: see book for details



# Decision Tree





# Other Types of Classifiers

- Neural net classifiers are studied in artificial intelligence and are not covered here
- Bayesian classifiers (see book for details)
- Support Vector Machines (see book for details)



# 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



# Association Rules (Cont.)

- Rules have an associated support, as well as an associated confidence.
- **Support** is a measure of what fraction of the population satisfies both the antecedent and the consequent of the rule.
  - E.g. suppose only 0.001 percent of all purchases include milk and screwdrivers. The support for the rule is  $milk \Rightarrow screwdrivers$  is low.
- **Confidence** is a measure of how often the consequent is true when the antecedent is true.
  - E.g. the rule  $bread \Rightarrow milk$  has a confidence of 80 percent if 80 percent of the purchases that include bread also include milk.



# Clustering

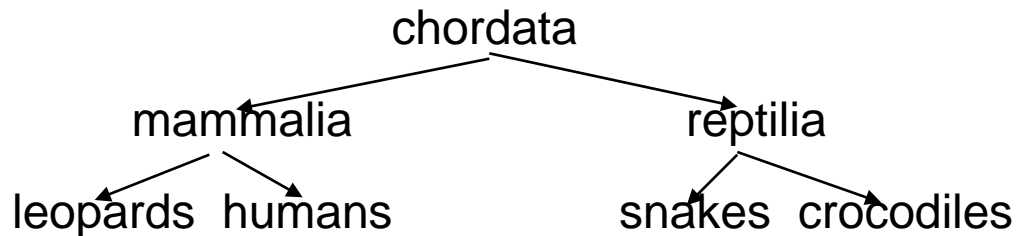
- Clustering: Intuitively, finding clusters of points in the given data such that similar points lie in the same cluster
- Can be formalized using distance metrics in several ways
  - Group points into  $k$  sets (for a given  $k$ ) such that the average distance of points from the centroid of their 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
  - E.g. the Birch clustering algorithm (more shortly)





# 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)



# Other Types of Mining

- **Text mining**: application of data mining to textual documents
  - cluster Web pages to find related pages
  - cluster pages a user has visited to organize their visit history
  - classify Web pages automatically into a Web directory



# Information Retrieval



# Information Retrieval Systems

- **Information retrieval (IR)** systems use a simpler data model than database systems
  - Information organized as a collection of documents
  - Documents are unstructured, no schema
- Information retrieval locates relevant documents, on the basis of user input such as keywords or example documents
  - e.g., find documents containing the words “database systems”
- Can be used even on textual descriptions provided with non-textual data such as images
- Web search engines are the most familiar example of IR systems



# Information Retrieval Systems (Cont.)

- Differences from database systems
  - IR systems don't deal with transactional updates (including concurrency control and recovery)
  - Database systems deal with structured data, with schemas that define the data organization
  - IR systems deal with some querying issues not generally addressed by database systems
    - ▶ Approximate searching by keywords
    - ▶ Ranking of retrieved answers by estimated degree of relevance



# Keyword Search

- In **full text** retrieval, all the words in each document are considered to be keywords.
  - We use the word **term** to refer to the words in a document
- Ranking of documents on the basis of estimated relevance to a keyword query is critical
  - Relevance ranking is based on factors such as
    - ▶ **Term frequency**
      - Frequency of occurrence of query keyword in document
    - ▶ **Inverse document frequency**
      - How many documents the query keyword occurs in
        - » Fewer → give more importance to keyword
    - ▶ **Hyperlinks to documents**
      - More links to a document → document is more important



# Relevance Using Hyperlinks

- Use number of hyperlinks to a site as a measure of the popularity or **prestige** of the site
  - Count only one hyperlink from each site (why? - see previous slide)
  - Popularity measure is for site, not for individual page
    - ▶ But, most hyperlinks are to root of site
    - ▶ Also, concept of “site” difficult to define since a URL prefix like `cs.yale.edu` contains many unrelated pages of varying popularity
- Refinements
  - When computing prestige based on links to a site, give more weight to links from sites that themselves have higher prestige
    - ▶ Definition is circular
    - ▶ Set up and solve system of simultaneous linear equations
  - Above idea is basis of the Google **PageRank** ranking mechanism



# Web Search Engines

- **Web crawlers** are programs that locate and gather information on the Web
  - Recursively follow hyperlinks present in known documents, to find other documents
    - ▶ Starting from a *seed* set of documents
  - Fetched documents
    - ▶ Handed over to an indexing system
    - ▶ Can be discarded after indexing, or store as a *cached* copy





# Information Retrieval and Structured Data

- Information retrieval systems originally treated documents as a collection of words
- **Information extraction systems** infer structure from documents, e.g.:
  - Extraction of house attributes (size, address, number of bedrooms, etc.) from a text advertisement
  - Extraction of topic and people named from a new article
- Relations or XML structures used to store extracted data
  - System seeks connections among data to answer queries
  - **Keyword querying on structured data**