

# **Data Mining and Data Warehousing**

# Data Mining and Data Warehousing

- Decision Support Systems
- Data Warehousing
- Data Mining
- Classification
- 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.)

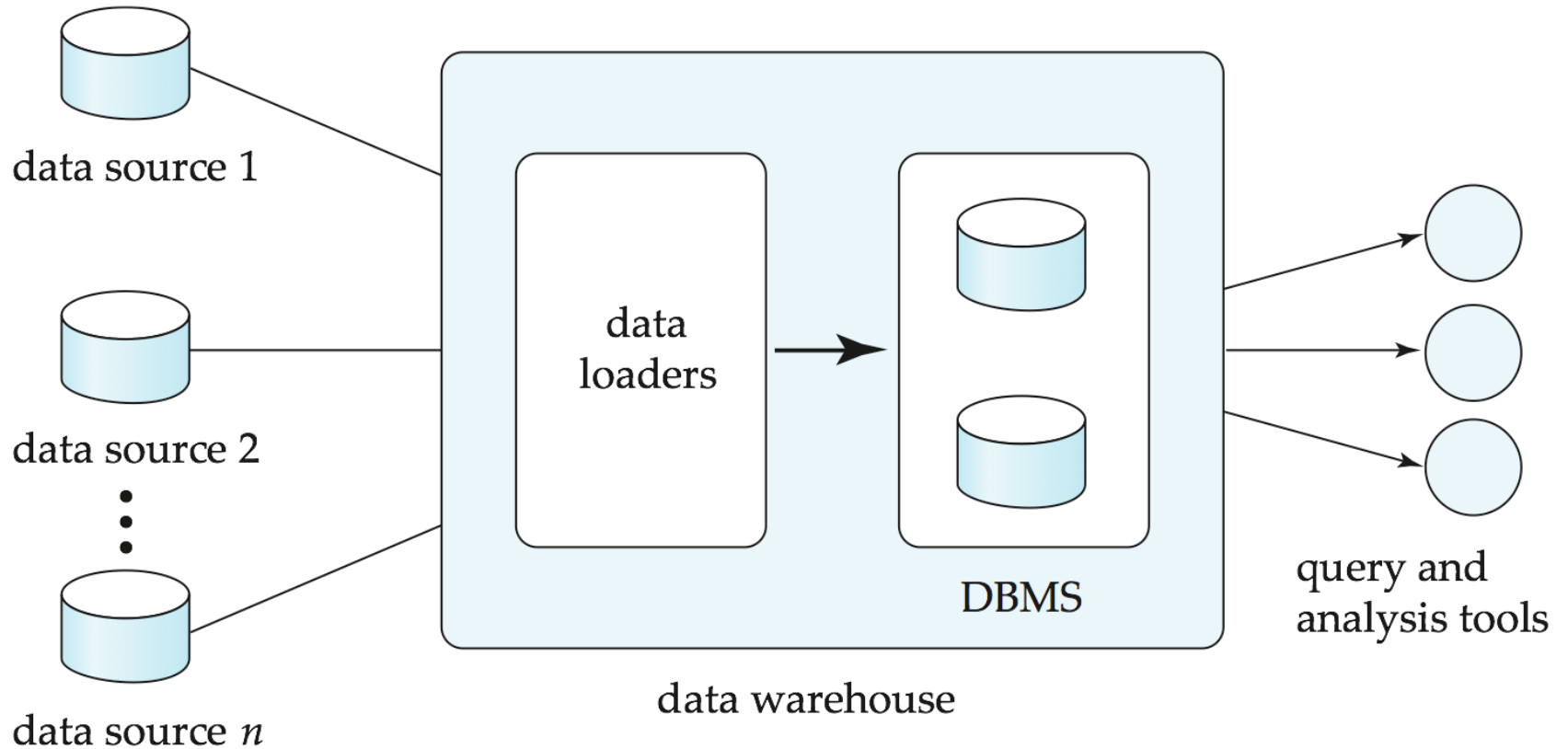
# Decision-Support Systems: Overview

- **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 (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



# Design Issues

- *When and how to gather data*
  - **Source driven architecture**: data sources transmit new information to warehouse, either continuously or periodically (e.g., at night)
  - **Destination driven architecture**: warehouse periodically requests new information from data sources
  - Keeping warehouse exactly synchronized with data sources (e.g., using two-phase commit) is too expensive
    - ▶ Usually OK to have slightly out-of-date data at warehouse
    - ▶ Data/updates are periodically downloaded from online transaction processing (OLTP) systems.
- *What schema to use*
  - Schema integration

# More Warehouse Design Issues

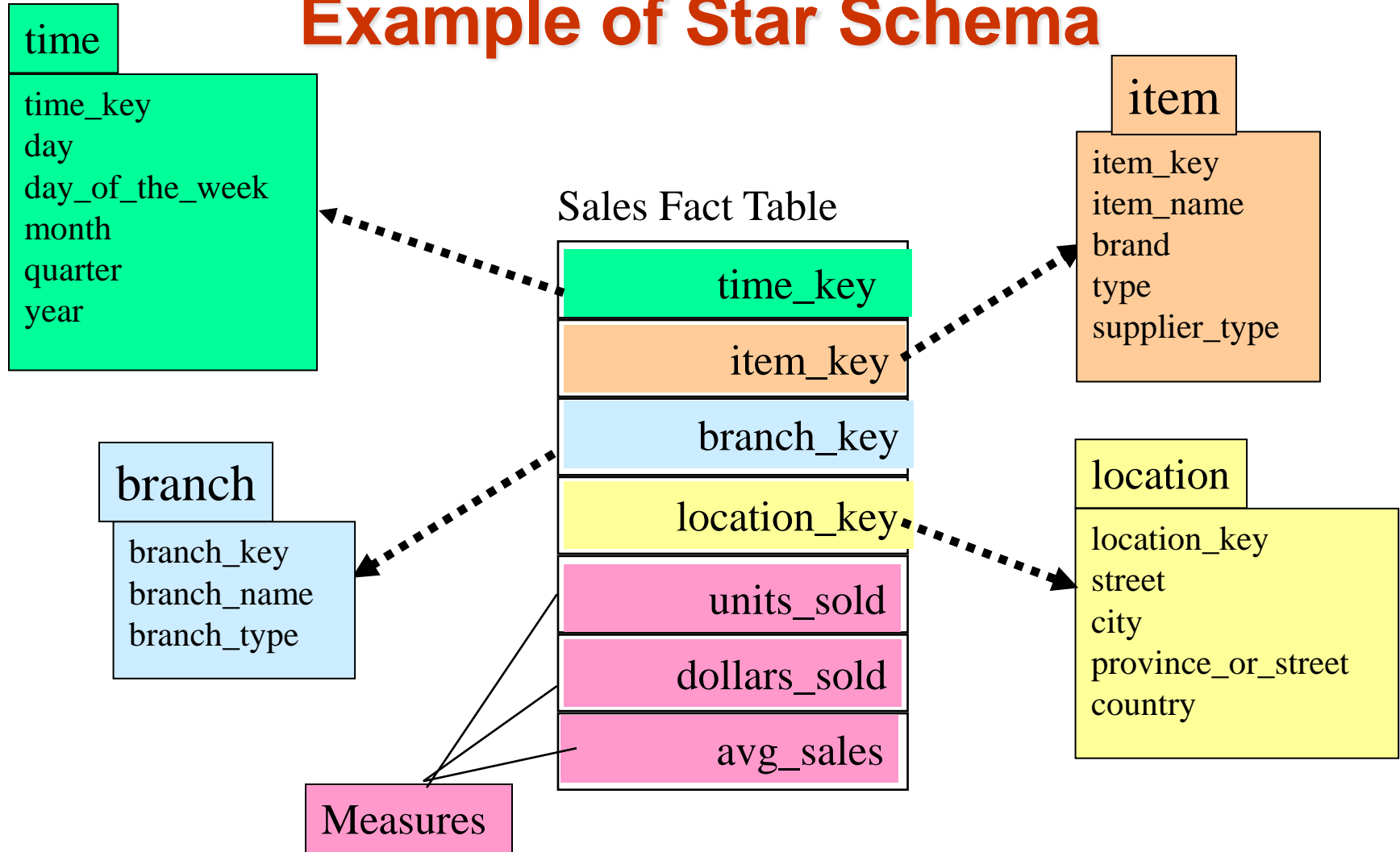
- *Data cleansing*
  - E.g., correct mistakes in addresses (misspellings, zip code errors)
  - **Merge** address lists from different sources and **purge** duplicates
- *How to propagate updates*
  - Warehouse schema may be a (materialized) view of schema from data sources
- *What data to summarize*
  - Raw data may be too large to store on-line
  - Aggregate values (totals/subtotals) often suffice
  - Queries on raw data can often be transformed by query optimizer to use aggregate values



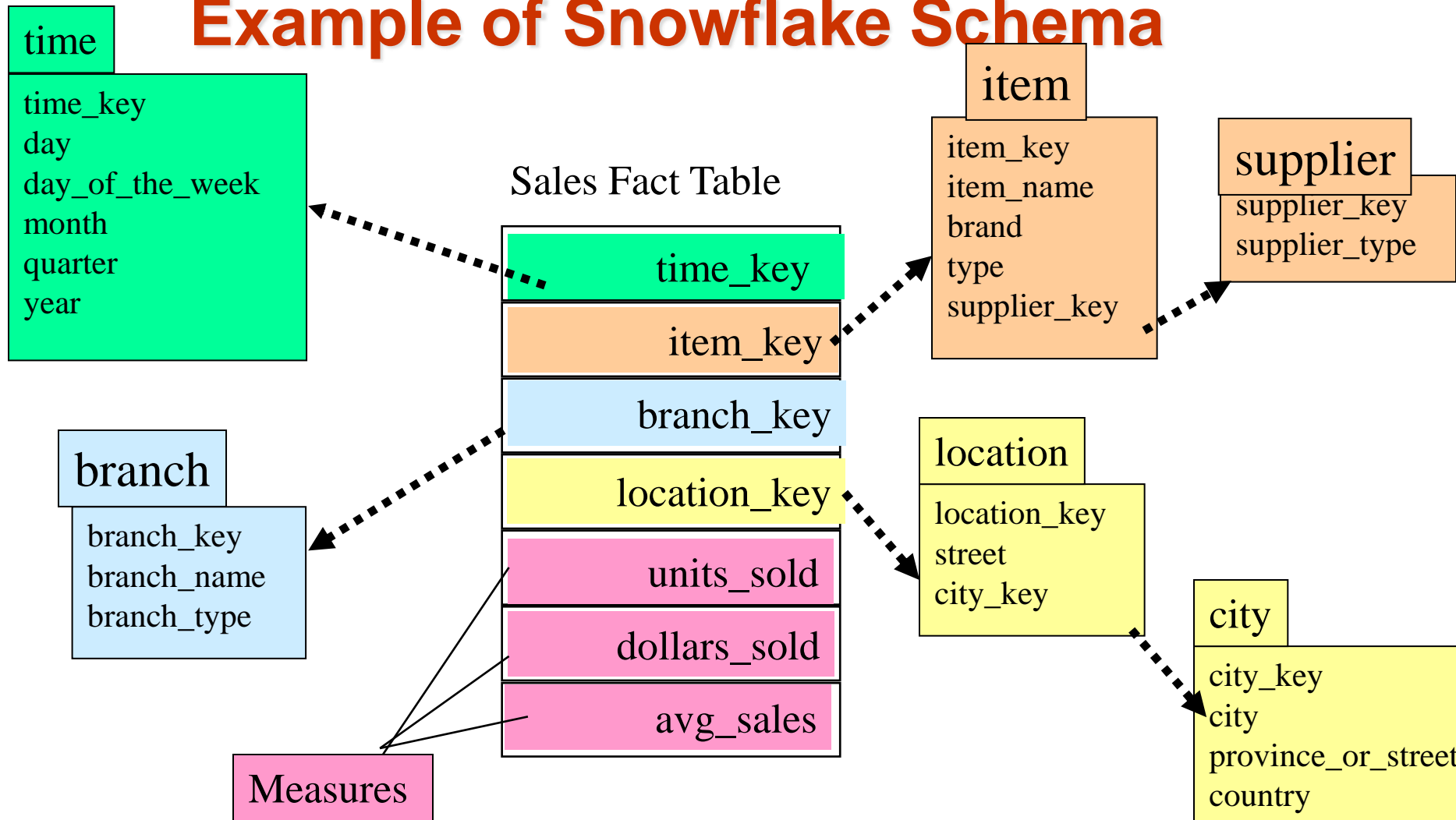
# Warehouse Schemas

- 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

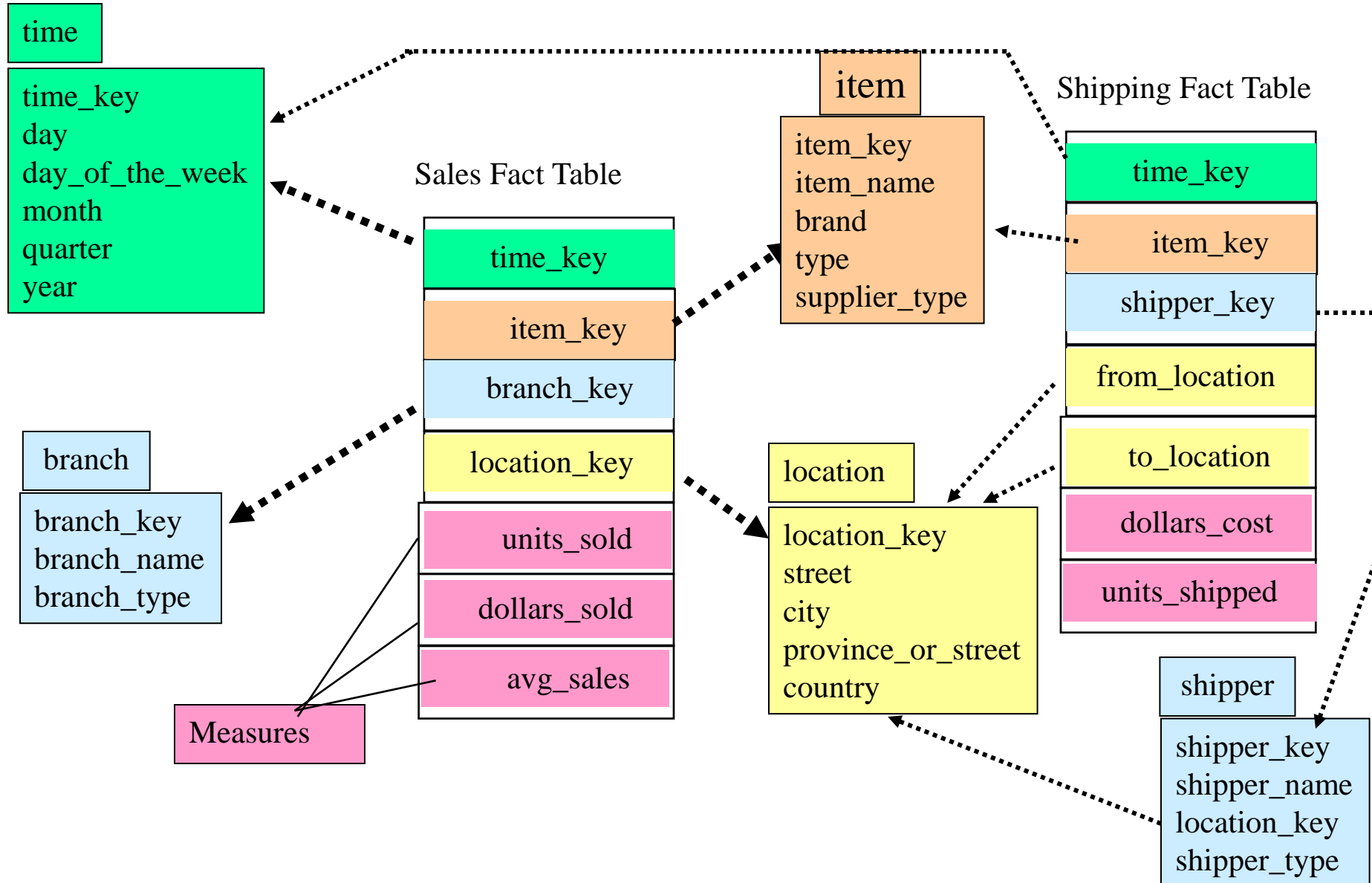
# Example of Star Schema



# Example of Snowflake Schema

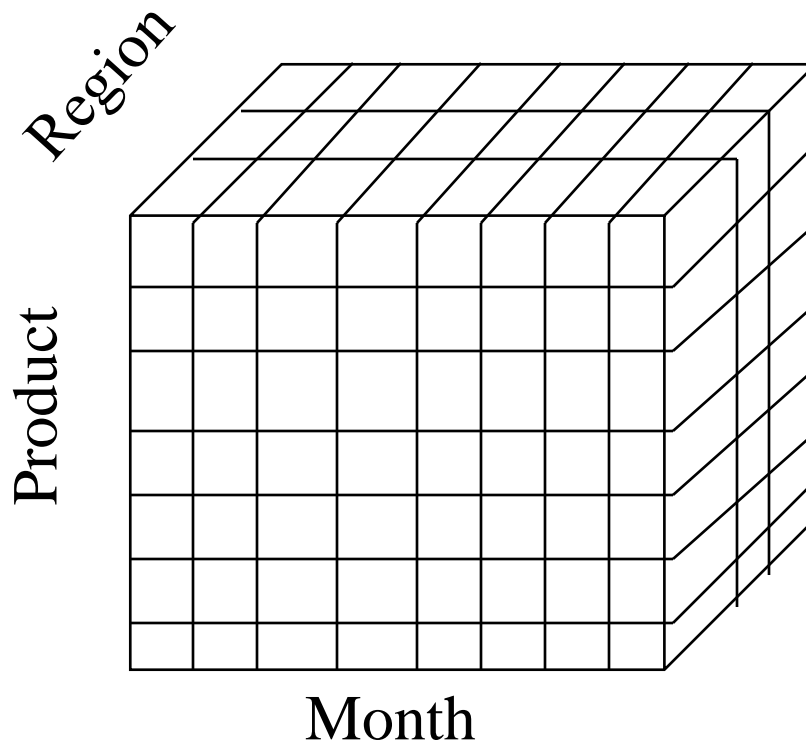


# Example of Fact Constellation

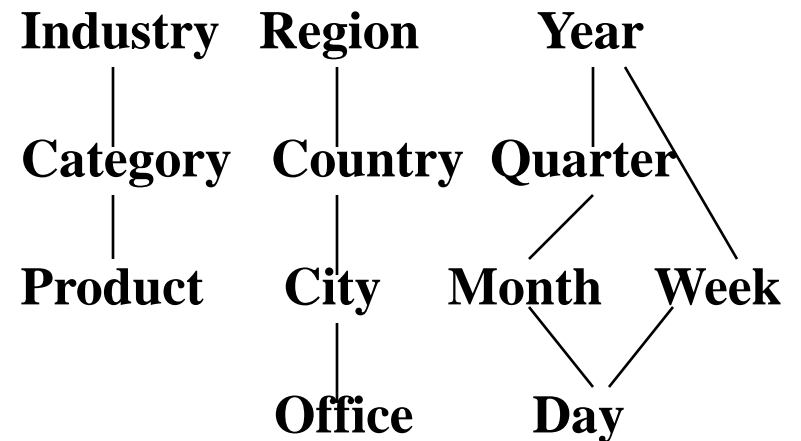


# Multidimensional Data

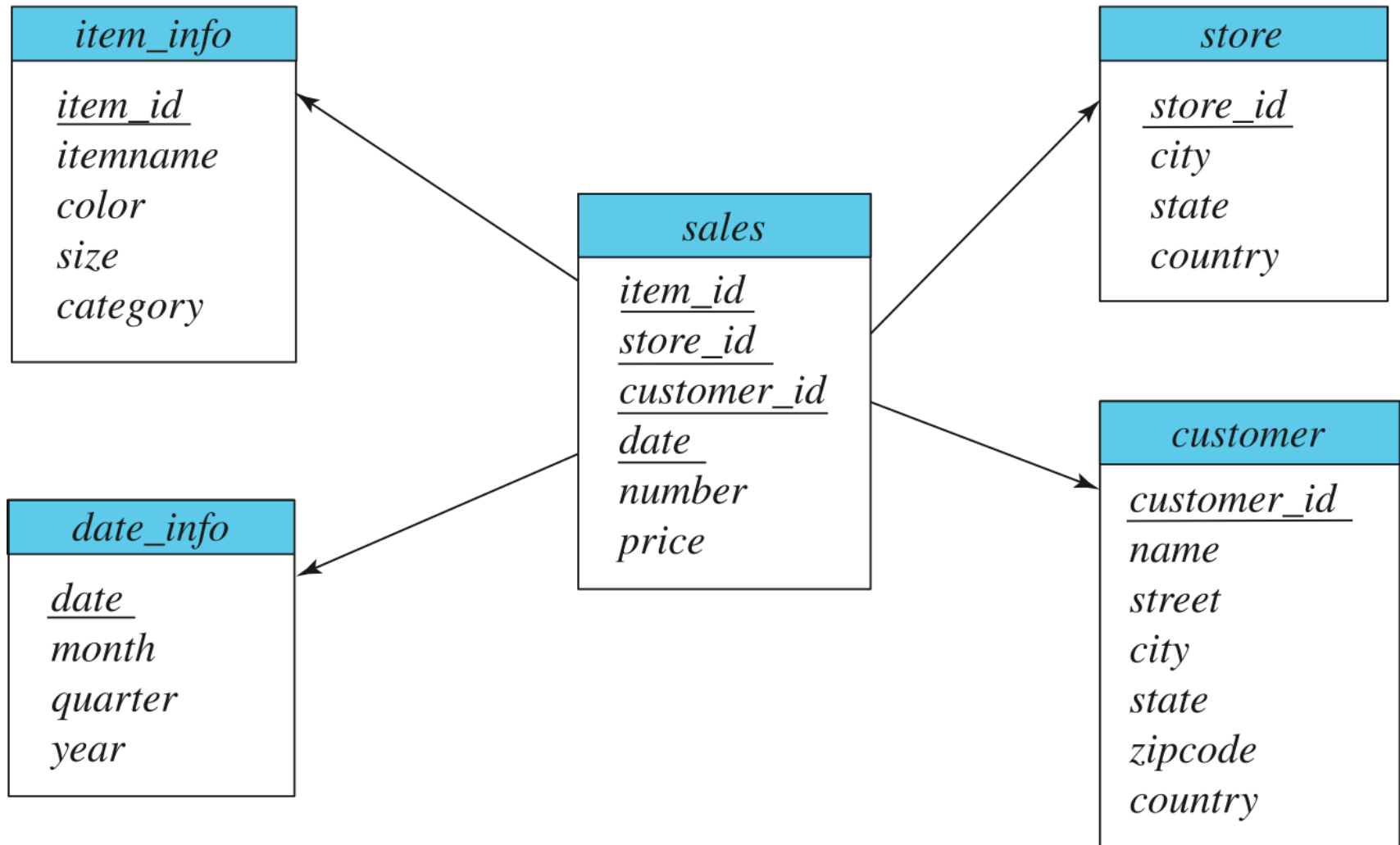
- Sales volume as a function of product, month, and region



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



# Data Warehouse Schema



# 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

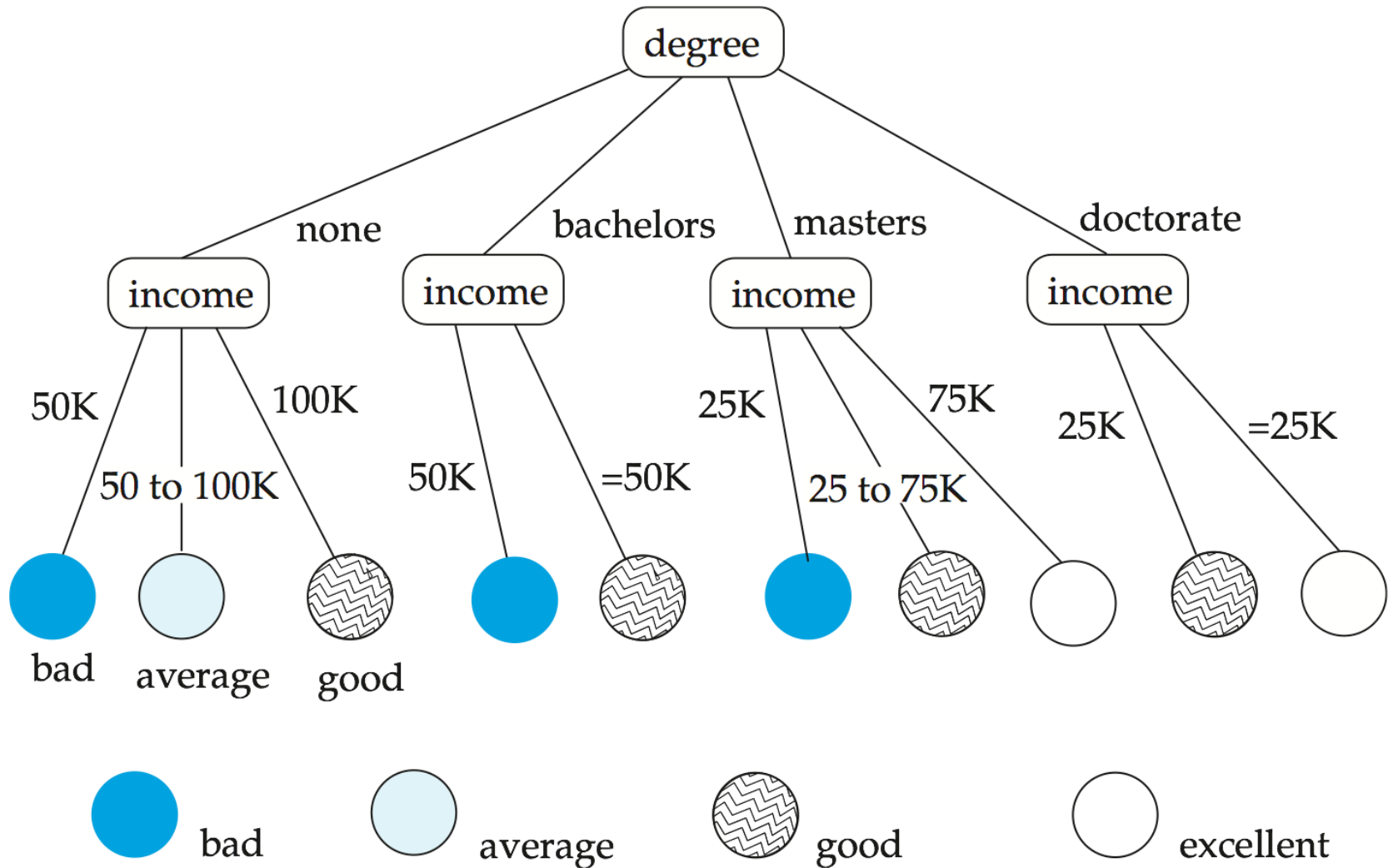
- ▶ E.g., typhoid cases were clustered in an area surrounding a contaminated well
- ▶ Detection of clusters remains important in detecting epidemics



# 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.

# Decision Tree



# Construction of Decision Trees

- **Training set:** a data sample in which the classification is already known.
- **Greedy** top down generation of decision trees.
  - Each internal node of the tree partitions the data into groups based on a **partitioning attribute**, and a **partitioning condition** for the node
  - **Leaf** node:
    - ▶ all (or most) of the items at the node belong to the same class, or
    - ▶ all attributes have been considered, and no further partitioning is possible.

# Best Splits

- Pick best attributes and conditions on which to partition
- The purity of a set  $S$  of training instances can be measured quantitatively in several ways.
  - Notation: number of classes =  $k$ , number of instances =  $|S|$ , fraction of instances in class  $i = p_i$ .
- The **Gini** measure of purity is defined as
  - [
$$\text{Gini}(S) = 1 - \sum_{i=1}^k p_i^2$$
]
  - When all instances are in a single class, the Gini value is 0
  - It reaches its maximum (of  $1 - 1/k$ ) if each class the same number of instances.

# Best Splits (Cont.)

- Another measure of purity is the **entropy** measure, which is defined as

$$\text{entropy}(S) = - \sum_{i=1}^k p_i \log_2 p_i$$

- When a set  $S$  is split into multiple sets  $S_i$ ,  $i=1, 2, \dots, r$ , we can measure the purity of the resultant set of sets as:

$$\text{purity}(S_1, S_2, \dots, S_r) = \sum_{i=1}^r \frac{|S_i|}{|S|} \text{purity}(S_i)$$

- The information gain due to particular split of  $S$  into  $S_i$ ,  $i = 1, 2, \dots, r$   
**Information-gain** ( $S, \{S_1, S_2, \dots, S_r\}$ ) =  $\text{purity}(S) - \text{purity}(S_1, S_2, \dots, S_r)$

# Best Splits (Cont.)

- Measure of “cost” of a split:

$$\text{Information-content } (S, \{S_1, S_2, \dots, S_r\}) = - \sum_{i=1}^r \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

- **Information-gain ratio** = 
$$\frac{\text{Information-gain } (S, \{S_1, S_2, \dots, S_r\})}{\text{Information-content } (S, \{S_1, S_2, \dots, S_r\})}$$
- The best split is the one that gives the maximum information gain ratio

# Finding Best Splits

- Categorical attributes (with no meaningful order):
  - Multi-way split, one child for each value
  - Binary split: try all possible breakup of values into two sets, and pick the best
- Continuous-valued attributes (can be sorted in a meaningful order)
  - Binary split:
    - ▶ Sort values, try each as a split point
      - E.g., if values are 1, 10, 15, 25, split at  $\leq 1$ ,  $\leq 10$ ,  $\leq 15$
    - ▶ Pick the value that gives best split
  - Multi-way split:
    - ▶ A series of binary splits on the same attribute has roughly equivalent effect

# Decision-Tree Construction Algorithm

**Procedure** *GrowTree* ( $S$ )

    Partition ( $S$ );

**Procedure** Partition ( $S$ )

**if** (  $\text{purity}(S) > \delta_p$  or  $|S| < \delta_s$  ) **then**  
        **return**;

**for each** attribute  $A$

        evaluate splits on attribute  $A$ ;

    Use best split found (across all attributes) to partition  
         $S$  into  $S_1, S_2, \dots, S_r$

**for**  $i = 1, 2, \dots, r$

        Partition ( $S_i$ );



# Other Types of Classifiers

- Neural net classifiers are studied in artificial intelligence and are not covered here
- Bayesian classifiers use **Bayes theorem**, which says

$$p(c_j | d) = \frac{p(d | c_j) p(c_j)}{p(d)}$$

where

$p(c_j | d)$  = probability of instance  $d$  being in class  $c_j$ ,

$p(d | c_j)$  = probability of generating instance  $d$  given class  $c_j$ ,

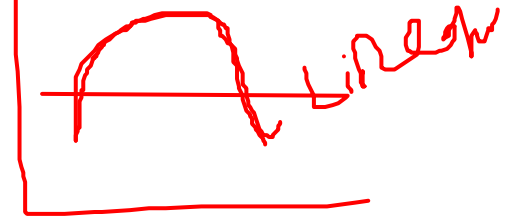
$p(c_j)$  = probability of occurrence of class  $c_j$ , and

$p(d)$  = probability of instance  $d$  occurring

# Naïve Bayesian Classifiers

- Bayesian classifiers require
  - computation of  $p(d | c_j)$
  - precomputation of  $p(c_j)$
  - $p(d)$  can be ignored since it is the same for all classes
- To simplify the task, **naïve Bayesian classifiers** assume attributes have independent distributions, and thereby estimate
$$p(d | c_j) = p(d_1 | c_j) * p(d_2 | c_j) * \dots * (p(d_n | c_j))$$
  - Each of the  $p(d_i | c_j)$  can be estimated from a histogram on  $d_i$  values for each class  $c_j$ 
    - ▶ the histogram is computed from the training instances
  - Histograms on multiple attributes are more expensive to compute and store

# Regression



□ Regression deals with the prediction of a value, rather than a class.

□ Given values for a set of variables,  $X_1, X_2, \dots, X_n$ , we wish to predict the value of a variable  $Y$ .

□ One way is to infer coefficients  $a_0, a_1, a_1, \dots, a_n$  such that

$$Y = a_0 + \underline{a_1 * X_1} + \underline{a_2 * X_2} + \dots + \underline{a_n * X_n}$$



□ Finding such a linear polynomial is called linear regression.

□ In general, the process of finding a curve that fits the data is also called **curve fitting**.

□ The fit may only be approximate

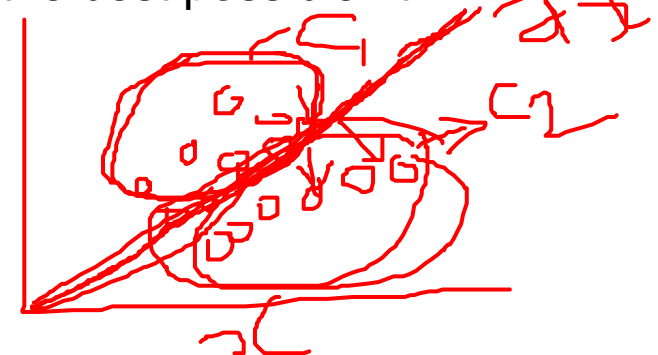
□ because of noise in the data, or

□ because the relationship is not exactly a polynomial

□ Regression aims to find coefficients that give the best possible fit.



$$y = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_n x_n$$

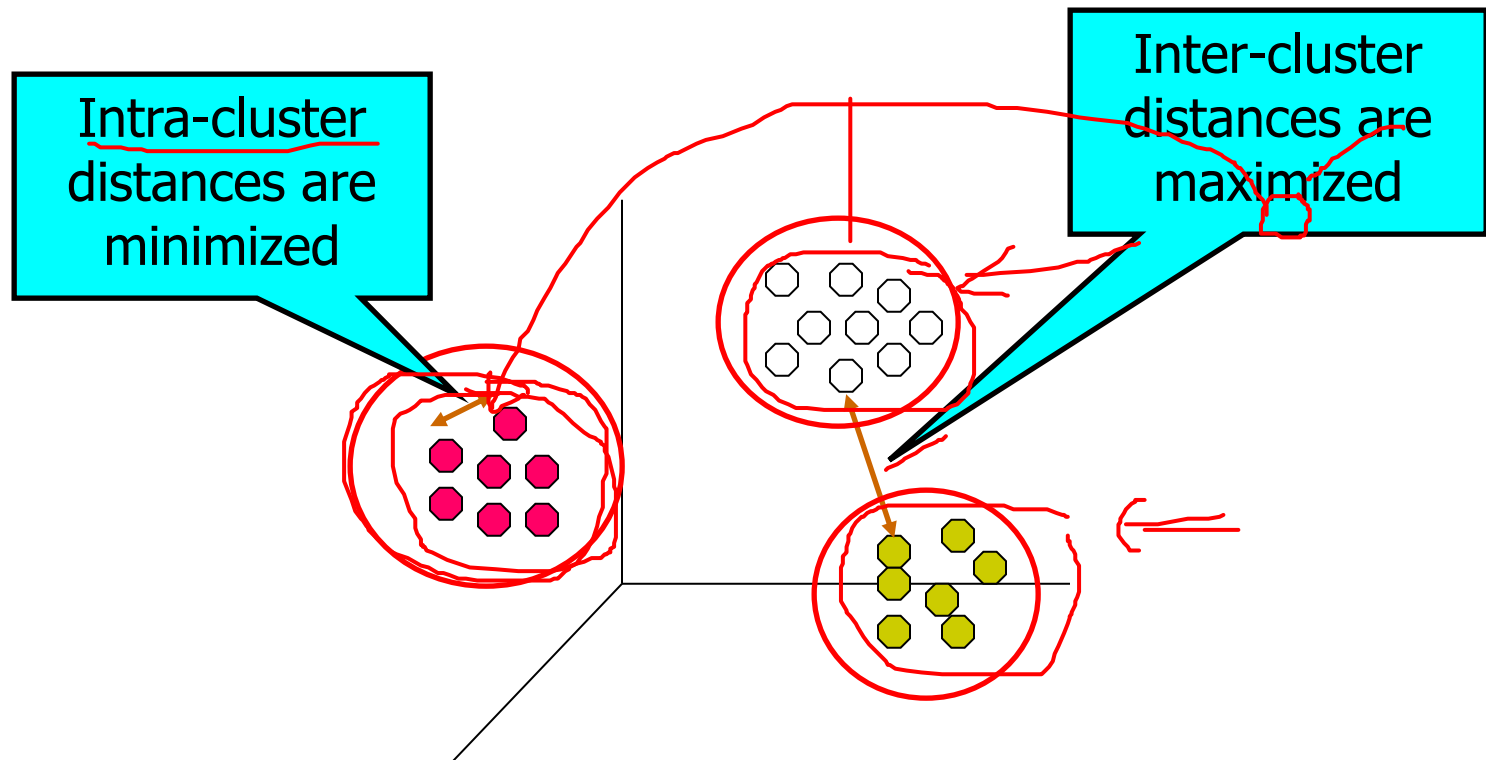


# What is clustering?

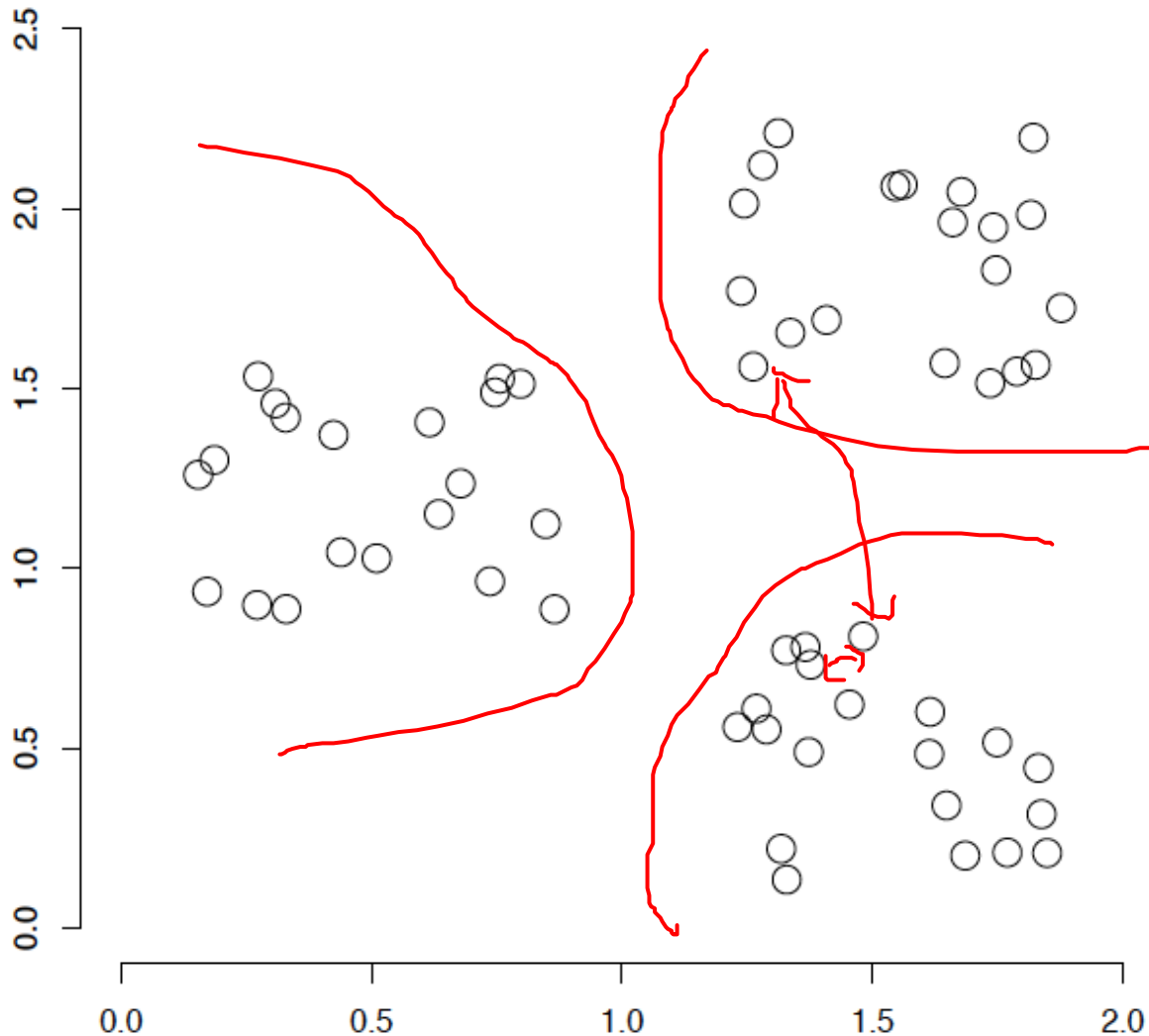
- **Clustering**: the process of grouping a set of objects into classes of similar objects
  - Documents within a cluster should be similar.
  - Documents from different clusters should be dissimilar.
- The commonest form of *unsupervised learning*
  - ▶ Unsupervised learning = learning from raw data, as opposed to supervised data where a classification of examples is given
  - A common and important task that finds many applications in IR and other places

# Cluster Analysis

- 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



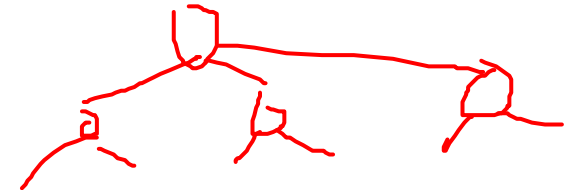
# A data set with clear cluster structure



- How would you design an algorithm for finding the three clusters in this case?

# Types of Clusterings

- A **clustering** is a set of clusters
- Important distinction between hierarchical and partitional sets of clusters
- Partitional Clustering
  - A division data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset
- Hierarchical clustering
  - A set of nested clusters organized as a hierarchical tree



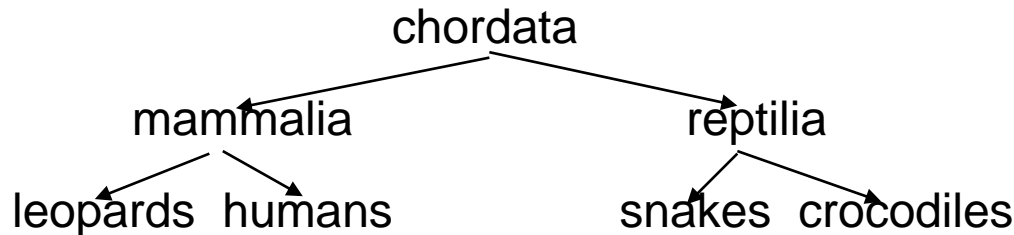
# Clustering

- 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)
- **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

# Clustering Algorithms

- Clustering algorithms have been designed to handle very large datasets
- E.g., the Birch algorithm
  - Main idea: use an in-memory R-tree to store points that are being clustered
  - Insert points one at a time into the R-tree, merging a new point with an existing cluster if it is less than some  $\delta$  distance away
  - If there are more leaf nodes than fit in memory, merge existing clusters that are close to each other
  - At the end of first pass we get a large number of clusters at the leaves of the R-tree
    - ▶ Merge clusters to reduce the number of clusters