# Unsupervised Learning

# What is Cluster Analysis?

- Cluster: A collection of data objects
  - similar (or related) to one another within the same group
  - dissimilar (or unrelated) to the objects in other groups
- Cluster analysis (or *clustering*, *data segmentation, …*)
  - Finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters
- Unsupervised learning: no predefined classes (i.e., *learning by observations* vs. learning by examples: supervised)
- Typical applications
  - As a stand-alone tool to get insight into data distribution
  - As a preprocessing step for other algorithms

# Clustering for Data Understanding and Applications

- Biology: taxonomy of living things: kingdom, phylum, class, order, family, genus and species

- Information retrieval: document clustering

- Land use: Identification of areas of similar land use in an earth observation database

- Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs

- City-planning: Identifying groups of houses according to their house type, value, and geographical location

- Earth-quake studies: Observed earth quake epicenters should be clustered along continent faults

- Climate: understanding earth climate, find patterns of atmospheric and ocean

- Economic Science: market resarch

# Clustering as a Preprocessing Tool (Utility)

- Summarization:

  - Preprocessing for regression, PCA, classification, and association analysis

- Compression:

  - Image processing: vector quantization

- Finding K-nearest Neighbors

  - Localizing search to one or a small number of clusters

- Outlier detection

  - Outliers are often viewed as those "far away" from any cluster

# Quality: What Is Good Clustering?

- A <u>good clustering</u> method will produce high quality clusters

  - high <u>intra-class</u> similarity: <span style="color:red">cohesive</span> within clusters

  - low <u>inter-class</u> similarity: <span style="color:red">distinctive</span> between clusters

- The <u>quality</u> of a clustering method depends on

  - the similarity measure used by the method

  - its implementation, and

  - Its ability to discover some or all of the <u>hidden</u> patterns

# Measure the Quality of Clustering

- <span style="color:red">Dissimilarity/Similarity metric</span>
    - Similarity is expressed in terms of a distance function, typically metric: $d(i, j)$
    - The definitions of <span style="color:red">distance functions</span> are usually rather different for interval-scaled, boolean, categorical, ordinal ratio, and vector variables
    - Weights should be associated with different variables based on applications and data semantics
- Quality of clustering:
    - There is usually a separate "quality" function that measures the "goodness" of a cluster.
    - It is hard to define "similar enough" or "good enough"
        - The answer is typically highly subjective

# Considerations for Cluster Analysis

- Partitioning criteria
  - Single level vs. hierarchical partitioning (often, multi-level hierarchical partitioning is desirable)

- Separation of clusters
  - Exclusive (e.g., one customer belongs to only one region) vs. non-exclusive (e.g., one document may belong to more than one class)

- Similarity measure
  - Distance-based (e.g., Euclidian, road network, vector)  vs. connectivity-based (e.g., density or contiguity)

- Clustering space
  - Full space (often when low dimensional) vs. subspaces (often in high-dimensional clustering)

# Requirements and Challenges

- Scalability
  - Clustering all the data instead of only on samples
- Ability to deal with different types of attributes
  - Numerical, binary, categorical, ordinal, linked, and mixture of these
- Constraint-based clustering
  - User may give inputs on constraints
  - Use domain knowledge to determine input parameters
- Interpretability and usability
- Others
  - Discovery of clusters with arbitrary shape
  - Ability to deal with noisy data
  - Incremental clustering and insensitivity to input order
  - High dimensionality

# Major Clustering Approaches

- Partitioning approach:
  - Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors
  - Typical methods: k-means, k-medoids, CLARANS, FCM
- Hierarchical approach:
  - Create a hierarchical decomposition of the set of data (or objects) using some criterion
  - Typical methods: Diana, Agnes, BIRCH, CAMELEON
- Density-based approach:
  - Based on connectivity and density functions
  - Typical methods: DBSCAN, OPTICS, DenClue
- Grid-based approach:
  - based on a multiple-level granularity structure
  - Typical methods: STING, WaveCluster, CLIQUE

# Partitioning Algorithms: Basic Concept

- Partitioning method: Partitioning a database **D** of **n** objects into a set of **k** clusters, such that the sum of squared distances is minimized (where $c_i$ is the centroid or medoid of cluster $C_i$)
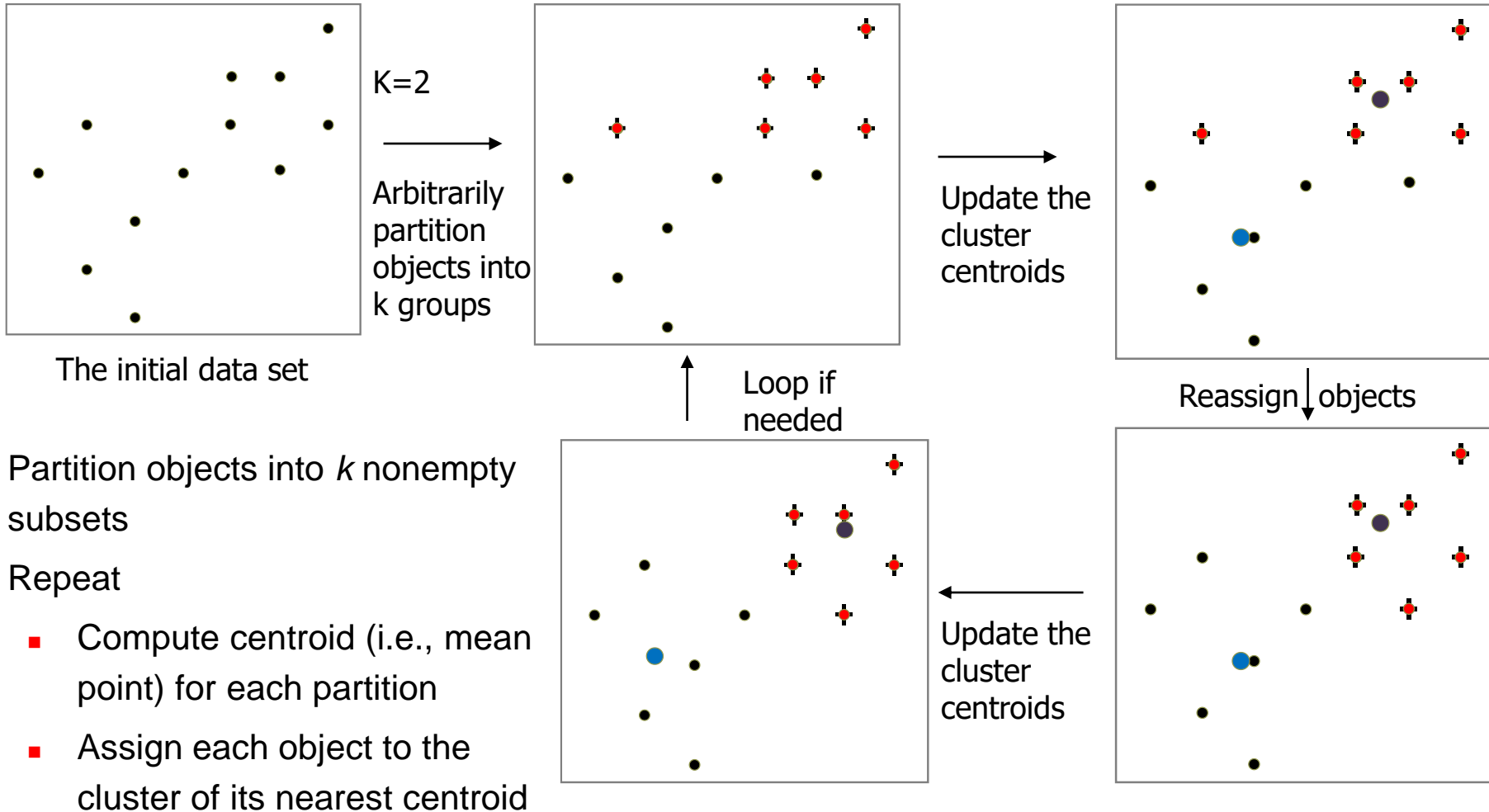
$$E = \Sigma_{i=1}^{k} \Sigma_{p \in C_i} (p - c_i)^2$$

- Given *k,* find a partition of *k clusters* that optimizes the chosen partitioning criterion
  - Global optimal: exhaustively enumerate all partitions
  - Heuristic methods: *k-means* and *k-medoids* algorithms
  - *k-means* : Each cluster is represented by the center of the cluster
  - *k-medoids* or PAM (Partition around medoids): Each cluster is represented by one of the objects in the cluster

# The *K-Means* Clustering Method

- Given *k*, the *k-means* algorithm is implemented in four steps:

  - Partition objects into *k* nonempty subsets

  - Compute seed points as the centroids of the clusters of the current partitioning (the centroid is the center, i.e., *mean point*, of the cluster)

  - Assign each object to the cluster with the nearest seed point

  - Go back to Step 2, stop when the assignment does not change

# An Example of *K-Means* Clustering



The initial data set

K=2

Arbitrarily partition objects into k groups

Update the cluster centroids

Loop if needed
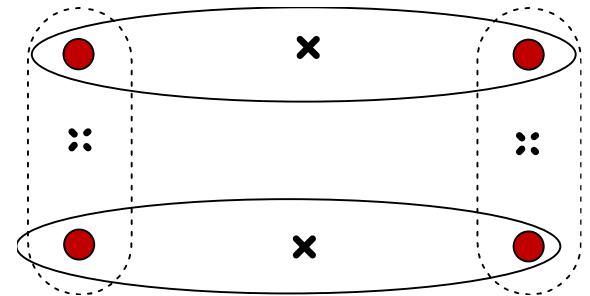
Reassign objects

Update the cluster centroids

- Partition objects into *k* nonempty subsets
- Repeat
  - Compute centroid (i.e., mean point) for each partition
  - Assign each object to the cluster of its nearest centroid
- Until no change

# Comments on the *K-Means* Method

- <u>Strength:</u> *Efficient*: O(*tkn*), where *n* is # objects, *k* is # clusters, and *t* is # iterations. Normally, *k, t* $<<$ *n*.

    - Comparing: PAM: $O(k(n-k)^2)$, CLARA: $O(ks^2 + k(n-k))$

- <u>Comment:</u> Often terminates at a *local optimal*.

- <u>Weakness</u>

    - Applicable only to objects in a continuous n-dimensional space

        - Using the k-modes method for categorical data

        - In comparison, k-medoids can be applied to a wide range of data

    - Need to specify *k,* the *number* of clusters, in advance

    - Sensitive to noisy data and *outliers*

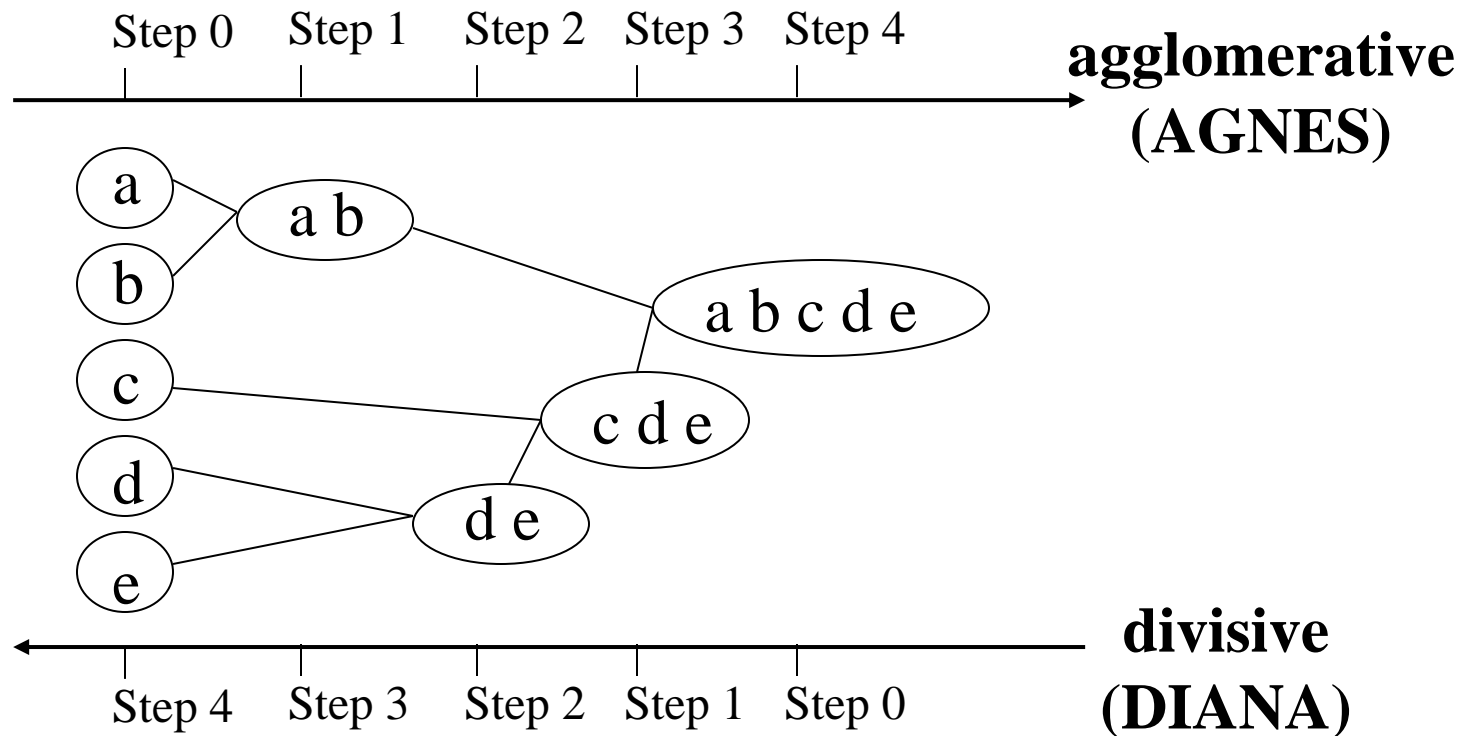    - Not suitable to discover clusters with *non-convex shapes*

# Variations of the *K-Means* Method

- Most of the variants of the *k-means* which differ in

  - Selection of the initial *k* means

  - Dissimilarity calculations

  - Strategies to calculate cluster means

- Handling categorical data: *k-modes*

  - Replacing means of clusters with <u>modes</u>

  - Using new dissimilarity measures to deal with categorical objects

  - Using a <u>frequency</u>-based method to update modes of clusters

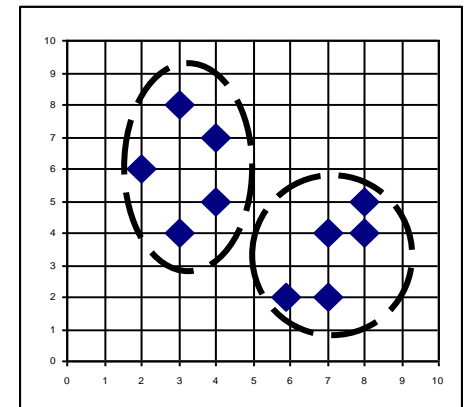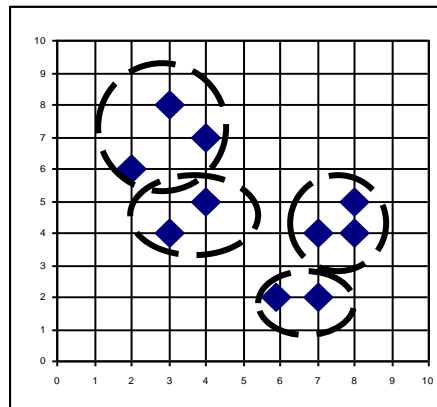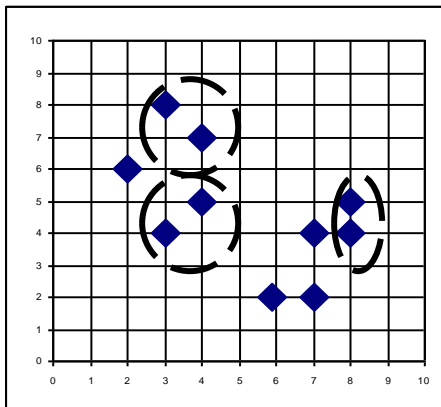  - A mixture of categorical and numerical data: *k-prototype* method

# Hierarchical Clustering

- Use distance matrix as clustering criteria.  This method does not require the number of clusters *k* as an input, but needs a termination condition
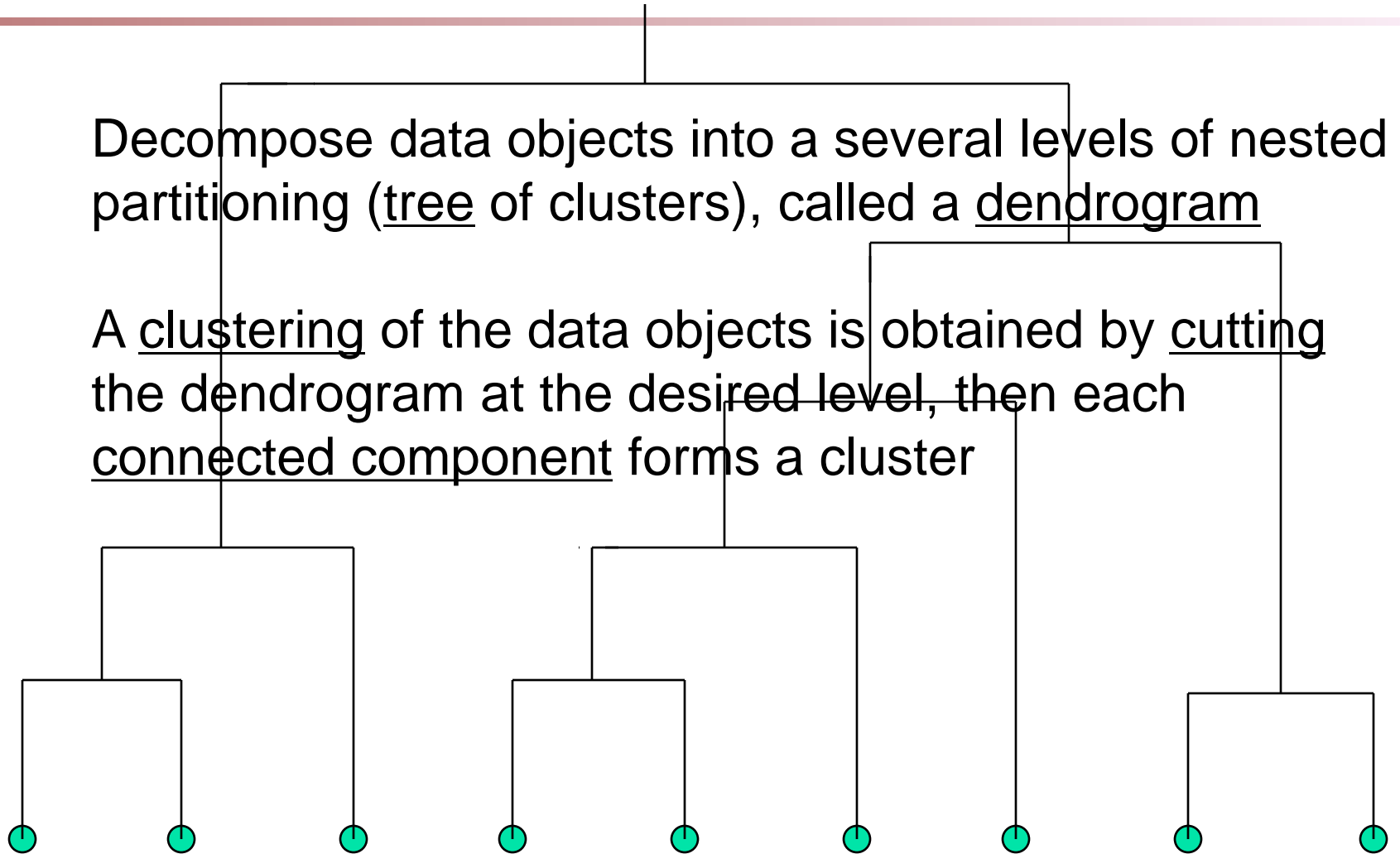
# AGNES (Agglomerative Nesting)

- Introduced in Kaufmann and Rousseeuw (1990)
- Implemented in statistical packages, e.g., Splus
- Use the **single-link** method and the dissimilarity matrix
- Merge nodes that have the least dissimilarity
- Go on in a non-descending fashion
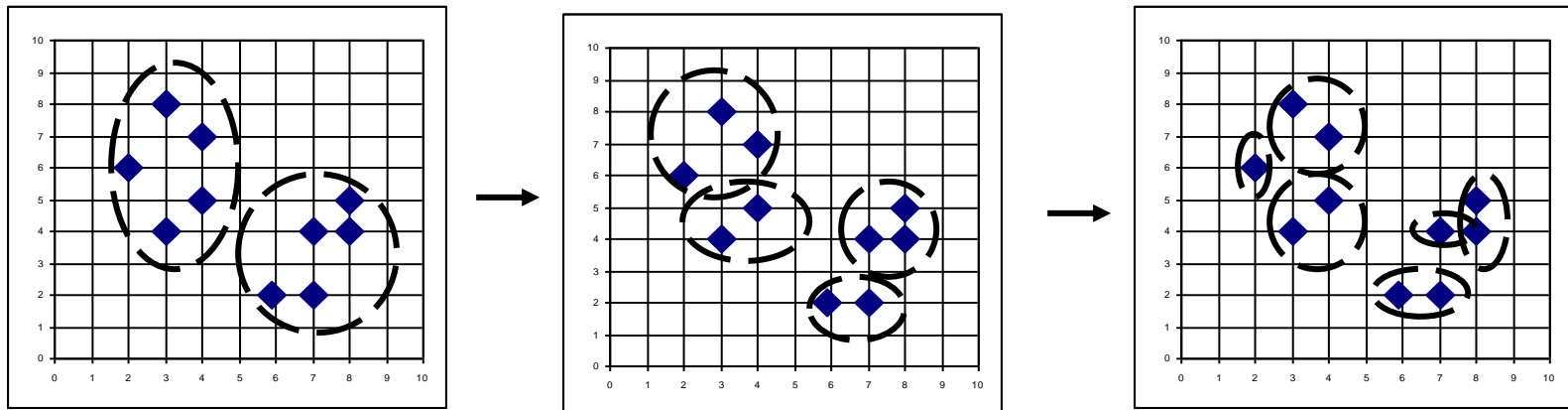- Eventually all nodes belong to the same cluster

# *Dendrogram:* Shows How Clusters are Merged

Decompose data objects into a several levels of nested partitioning (tree of clusters), called a dendrogram

A clustering of the data objects is obtained by cutting the dendrogram at the desired level, then each connected component forms a cluster
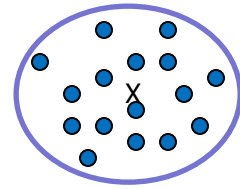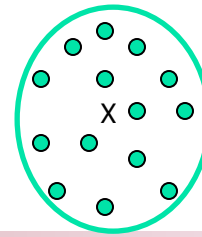
# DIANA (Divisive Analysis)

- Introduced in Kaufmann and Rousseeuw (1990)

- Implemented in statistical analysis packages, e.g., Splus

- Inverse order of AGNES

- Eventually each node forms a cluster on its own

# Distance between Clusters

- ■ Single link: smallest distance between an element in one cluster and an element in the other, i.e., dist($K_i$, $K_j$) = min($t_{ip}$, $t_{jq}$)

- ■ Complete link: largest distance between an element in one cluster and an element in the other, i.e., dist($K_i$, $K_j$) = max($t_{ip}$, $t_{jq}$)

- ■ Average: avg distance between an element in one cluster and an element in the other, i.e., dist($K_i$, $K_j$) = avg($t_{ip}$, $t_{jq}$)

- ■ Centroid: distance between the centroids of two clusters, i.e., dist($K_i$, $K_j$) = dist($C_i$, $C_j$)

- ■ Medoid: distance between the medoids of two clusters, i.e., dist($K_i$, $K_j$) = dist($M_i$, $M_j$)
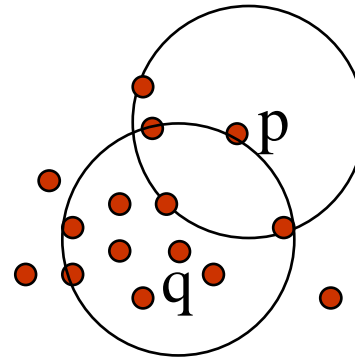  - ■ Medoid: a chosen, centrally located object in the cluster

# Density-Based Clustering Methods

- Clustering based on density (local cluster criterion), such as density-connected points
- Major features:
  - Discover clusters of arbitrary shape(S or oval shaped)
  - Handle noise
  - One scan
  - Need density parameters as termination condition
- Several interesting studies:
  - <u>DBSCAN:</u> Ester, et al. (KDD'96)
  - <u>OPTICS:</u> Ankerst, et al (SIGMOD'99).
  - <u>DENCLUE:</u> Hinneburg & D. Keim  (KDD'98)
  - <u>CLIQUE:</u> Agrawal, et al. (SIGMOD'98) (more grid-based)

# Density-Based Clustering: Basic Concepts

- Two parameters*:*

  - *Eps*: Maximum radius of the neighbourhood

  - *MinPts*: Minimum number of points in an Eps-neighbourhood of that point

- $N_{Eps}(p)$: {q belongs to D | dist(p,q) ≤ Eps}

- Directly density-reachable: A point *p* is directly density-reachable from a point *q* w.r.t. *Eps*, *MinPts* if

  - *p* belongs to $N_{Eps}(q)$

  - core point condition:
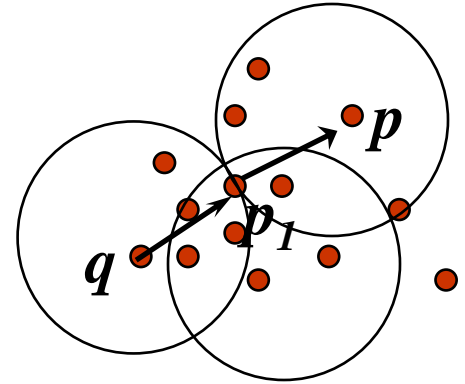
    $$|N_{Eps}(q)| \geq MinPts$$

MinPts = 5

Eps = 1 cm

# Density-Reachable and Density-Connected

- Density-reachable:

  - A point *p* is <span style="color:red">density-reachable</span> from a point *q* w.r.t. *Eps*, *MinPts* if there is a chain of points $p_1, \ldots, p_n$, $p_1 = q$, $p_n = p$ such that $p_{i+1}$ is directly density-reachable from $p_i$

- Density-connected

  - A point *p* is <span style="color:red">density-connected</span> to a point *q* w.r.t. *Eps*, *MinPts* if there is a point *o* such that both, *p* and *q* are density-reachable from *o* w.r.t. *Eps* and *MinPts*

# DBSCAN: Density-Based Spatial Clustering of Applications with Noise

- Relies on a *density-based* notion of cluster: A *cluster* is defined as a maximal set of density-connected points
- Discovers clusters of arbitrary shape in spatial databases with noise

Outlier

Border

Core

Eps = 1cm

MinPts = 5

**Algorithm: DBSCAN:** a density-based clustering algorithm.

**Input:**

- $D$: a data set containing $n$ objects,

- $\epsilon$: the radius parameter, and

- *MinPts*: the neighborhood density threshold.

**Output:** A set of density-based clusters.

**Method:**

```
(1)   mark all objects as unvisited;
(2)   do
(3)        randomly select an unvisited object p;
(4)        mark p as visited;
(5)        if the ε-neighborhood of p has at least MinPts objects
(6)            create a new cluster C, and add p to C;
(7)            let N be the set of objects in the ε-neighborhood of p;
(8)            for each point p' in N
(9)                if p' is unvisited
(10)                   mark p' as visited;
(11)                   if the ε-neighborhood of p' has at least MinPts points,
                       add those points to N;
(12)               if p' is not yet a member of any cluster, add p' to C;
(13)           end for
(14)           output C;
(15)       else mark p as noise;
(16)  until no object is unvisited;
```
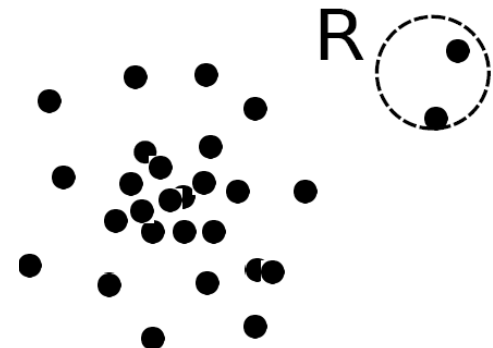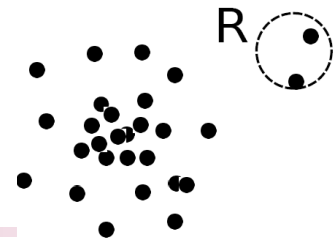
---

# DBSCAN: The Algorithm

- Arbitrary select a point $p$

- Retrieve all points density-reachable from $p$ w.r.t. *Eps* and *MinPts*

- If $p$ is a core point, a cluster is formed

- If $p$ is a border point, no points are density-reachable from $p$ and DBSCAN visits the next point of the database

- Continue the process until all of the points have been processed

# What Are Outliers?

- **Outlier**: A data object that **deviates significantly** from the normal objects as if it were **generated by a different mechanism**
  - Ex.:  Unusual credit card purchase, sports: Michael Jordon, Wayne Gretzky, ...
- Outliers are different from the noise data
  - Noise is random error or variance in a measured variable
  - Noise should be removed before outlier detection
- Outliers are interesting:  It violates the mechanism that generates the normal data
- Outlier detection vs. *novelty detection*: early stage, outlier; but later merged into the model
- Applications:
  - Credit card fraud detection
  - Telecom fraud detection
  - Customer segmentation
  - Medical analysis
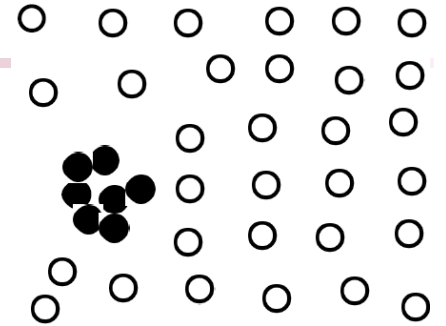
R

# Types of Outliers (I)



Global Outlier

- Three kinds: *global, contextual* and *collective* outliers
- **Global outlier** (or point anomaly)
  - Object is $O_g$ if it significantly deviates from the rest of the data set
  - Ex. Intrusion detection in computer networks
  - Issue: Find an appropriate measurement of deviation
- **Contextual outlier** (or *conditional outlier*)
  - Object is $O_c$ if it deviates significantly based on a selected context
  - Ex. $80^o$ F in Urbana: outlier? (depending on summer or winter?)
  - Attributes of data objects should be divided into two groups
    - Contextual attributes: defines the context, e.g., time & location
    - Behavioral attributes:  characteristics of the object, used in outlier evaluation, e.g., temperature
  - Can be viewed as a generalization of *local outliers*—whose density significantly deviates from its local area
  - Issue: How to define or formulate meaningful context?

# Types of Outliers (II)

- **Collective Outliers**


Collective Outlier

  - A subset of data objects *collectively* deviate significantly from the whole data set, even if the individual data objects may not be outliers

  - Applications: E.g., *intrusion detection*:
    - When a number of computers keep sending denial-of-service packages to each other

  - Detection of collective outliers
    - Consider not only behavior of individual objects, but also that of groups of objects
    - Need to have the background knowledge on the relationship among data objects, such as a distance or similarity measure on objects.

- A data set may have multiple types of outlier
- One object may belong to more than one type of outlier
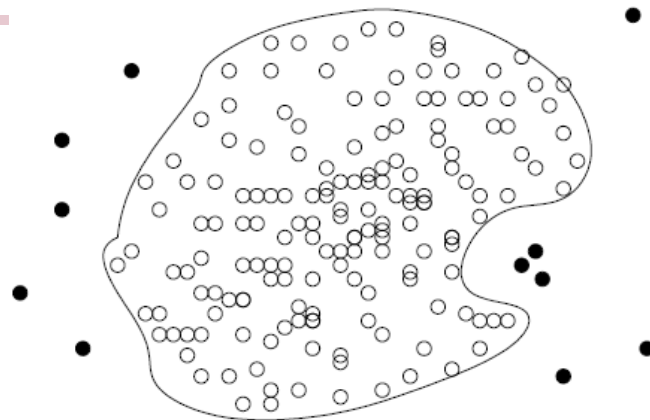
# Challenges of Outlier Detection

- Modeling normal objects and outliers properly
  - Hard to enumerate all possible normal behaviors in an application
  - The border between normal and outlier objects is often a gray area
- Application-specific outlier detection
  - Choice of distance measure among objects and the model of relationship among objects are often application-dependent
  - E.g., clinic data: a small deviation could be an outlier; while in marketing analysis, larger fluctuations
- Handling noise in outlier detection
  - Noise may distort the normal objects and blur the distinction between normal objects and outliers.  It may help hide outliers and reduce the effectiveness of outlier detection
- Understandability
  - Understand why these are outliers: Justification of the detection
  - Specify the degree of an outlier: the unlikelihood of the object being generated by a normal mechanism

# Supervised Methods

- Two ways to categorize outlier detection methods:
    - Based on <u>whether user-*labeled* examples of outliers can be obtained</u>:
        - Supervised, semi-supervised vs. unsupervised methods
    - Based on <u>*assumptions about normal data and outliers*</u>:
        - Statistical, proximity-based, and clustering-based methods
- **Outlier Detection I: Supervised Methods**
    - Modeling outlier detection as a classification problem
        - Samples examined by domain experts used for training & testing
    - Methods for Learning a classifier for outlier detection effectively:
        - Model normal objects & report those not matching the model as outliers, or
        - Model outliers and treat those not matching the model as normal
    - Challenges
        - Imbalanced classes, i.e., outliers are rare: Boost the outlier class and make up some artificial outliers
        - Catch as many outliers as possible, i.e., recall is more important than accuracy (i.e., not mislabeling normal objects as outliers)

# Classification-Based Method : One-Class Model

- Idea: Train a classification model that can distinguish "normal" data from outliers
- A brute-force approach: Consider a training set that contains samples labeled as "normal" and others labeled as "outlier"
  - But, the training set is typically heavily biased: # of "normal" samples likely far exceeds # of outlier samples
  - Cannot detect unseen anomaly
- One-class model: A classifier is built to describe only the normal class.
  - Learn the decision boundary of the normal class using classification methods such as SVM
  - Any samples that do not belong to the normal class (not within the decision boundary) are declared as outliers
  - Adv: can detect new outliers that may not appear close to any outlier objects in the training set
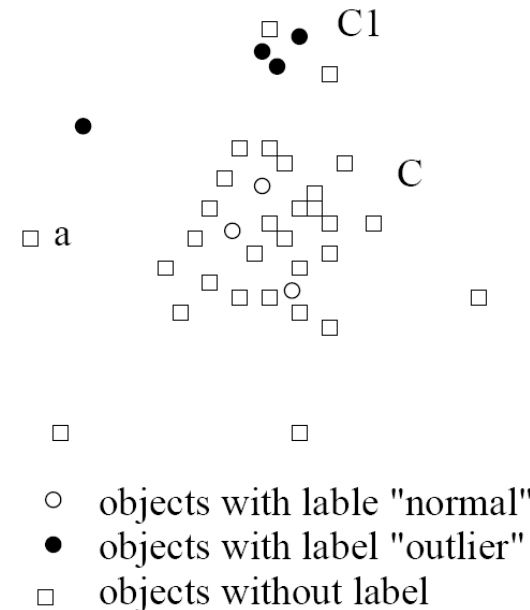  - Extension: Normal objects may belong to multiple classes

# Semi-Supervised Methods

- Semi-supervised learning: Combining classification-based and clustering-based methods
- Method
  - Using a clustering-based approach, find a large cluster, C, and a small cluster, $C_1$
  - Since some objects in C carry the label "normal", treat all objects in C as normal
  - Use the one-class model of this cluster to identify normal objects in outlier detection
  - Since some objects in cluster $C_1$ carry the label "outlier", declare all objects in $C_1$ as outliers
  - Any object that does not fall into the model for C (such as *a*) is considered an outlier as well
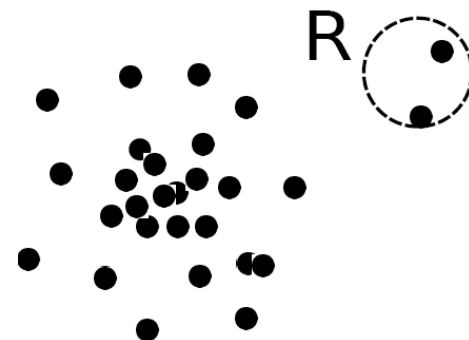- Comments on classification-based outlier detection methods
  - Strength: Outlier detection is fast
  - Bottleneck: Quality heavily depends on the availability and quality of the training set, but often difficult to obtain representative and high-quality training data

C1

C

a

○ objects with lable "normal"
● objects with label "outlier"
□ objects without label

# Unsupervised Methods

- Assume the normal objects are somewhat ``clustered'' into multiple groups, each having some distinct features
- An outlier is expected to be far away from any groups of normal objects
- Weakness: Cannot detect collective outlier effectively
    - Normal objects may not share any strong patterns, but the collective outliers may share high similarity in a small area
- Ex. In some intrusion or virus detection, normal activities are diverse
    - Unsupervised methods may have a high false positive rate but still miss many real outliers.
    - Supervised methods can be more effective, e.g., identify attacking some key resources
- Many clustering methods can be adapted for unsupervised methods
    - Find clusters, then outliers: not belonging to any cluster
    - Problem 1: Hard to distinguish noise from outliers
    - Problem 2: Costly since first clustering: but far less outliers than normal objects
        - Newer methods: tackle outliers directly

# Clustering-Based Methods

- Normal data belong to large and dense clusters, whereas outliers belong to small or sparse clusters, or do not belong to any clusters

- Example (right figure): two clusters
  - All points not in R form a large cluster
  - The two points in R form a tiny cluster, thus are outliers

- Since there are many clustering methods, there are many clustering-based outlier detection methods as well
- Clustering is expensive: straightforward adaption of a clustering method for outlier detection can be costly and does not scale up well for large data sets

# Clustering-Based Outlier Detection

- An object is an outlier if (1) it does not belong to any cluster, (2) there is a large distance between the object and its closest cluster , or (3) it belongs to a small or sparse cluster

- Case I: Not belong to any cluster
  - Identify animals not part of a flock:  Using a density-based clustering method such as DBSCAN

- Case 2:  Far from its closest cluster
  - Using k-means, partition data points of into clusters
  - For each object o, assign an outlier score based on its distance from its closest center
    - If $dist(o, c_o)/avg\_dist(c_o)$ is large, likely an outlier

- Ex. Intrusion detection: Consider the similarity between data points and the clusters in a training data set
  - Use a training set to find patterns of "normal" data, e.g., frequent itemsets in each segment, and cluster similar connections into groups
  - Compare new data points with the clusters mined—Outliers are possible attacks

# Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

- Frequent pattern: a pattern (a set of items, subsequences, substructures, etc.) that occurs frequently in a data set

- First proposed by Agrawal, Imielinski, and Swami [AIS93] in the context of frequent itemsets and association rule mining

- Motivation: Finding inherent regularities in data

  - What products were often purchased together?— Beer and diapers?!

  - What are the subsequent purchases after buying a PC?

  - What kinds of DNA are sensitive to this new drug?

  - Can we automatically classify web documents?

- Applications

  - Basket data analysis, cross-marketing, catalog design, sale campaign analysis, Web log (click stream) analysis, and DNA sequence analysis.

# Why Is Freq. Pattern Mining Important?

- Freq. pattern: An intrinsic and important property of datasets
- Foundation for many essential data mining tasks
  - Association, correlation, and causality analysis
  - Sequential, structural (e.g., sub-graph) patterns
  - Pattern analysis in spatiotemporal, multimedia, time-series, and stream data
  - Classification: discriminative, frequent pattern analysis
  - Cluster analysis: frequent pattern-based clustering
  - Data warehousing: iceberg cube and cube-gradient
  - Semantic data compression: fascicles
  - Broad applications

# Basic Concepts: Frequent Patterns

| Tid | Items bought |
|-----|--------------|
| 10 | Beer, Nuts, Diaper |
| 20 | Beer, Coffee, Diaper |
| 30 | Beer, Diaper, Eggs |
| 40 | Nuts, Eggs, Milk |
| 50 | Nuts, Coffee, Diaper, Eggs, Milk |



Customer buys both

Customer buys diaper

Customer buys beer

- **itemset**: A set of one or more items
- **k-itemset** $X = \{x_1, \ldots, x_k\}$
- *(absolute) support*, or, *support count* of X: Frequency or occurrence of an itemset X
- *(relative) support*, *s*, is the fraction of transactions that contains X (i.e., the probability that a transaction contains X)
- An itemset X is *frequent* if X's support is no less than a *minsup* threshold

# Basic Concepts: Association Rules

| Tid | Items bought |
|-----|--------------|
| 10 | Beer, Nuts, Diaper |
| 20 | Beer, Coffee, Diaper |
| 30 | Beer, Diaper, Eggs |
| 40 | Nuts, Eggs, Milk |
| 50 | Nuts, Coffee, Diaper, Eggs, Milk |

Customer buys both

Customer buys diaper

Customer buys beer

- Find all the rules $X \rightarrow Y$ with minimum support and confidence
  - support, $s$, probability that a transaction contains $X \cup Y$
  - confidence, $c$, conditional probability that a transaction having X also contains $Y$

Let  minsup = 50%, minconf = 50%

Freq. Pat.: Beer:3, Nuts:3, Diaper:4, Eggs:3, {Beer, Diaper}:3

- Association rules: (many more!)
  - *Beer → Diaper  (60%, 100%)*
  - *Diaper → Beer  (60%, 75%)*

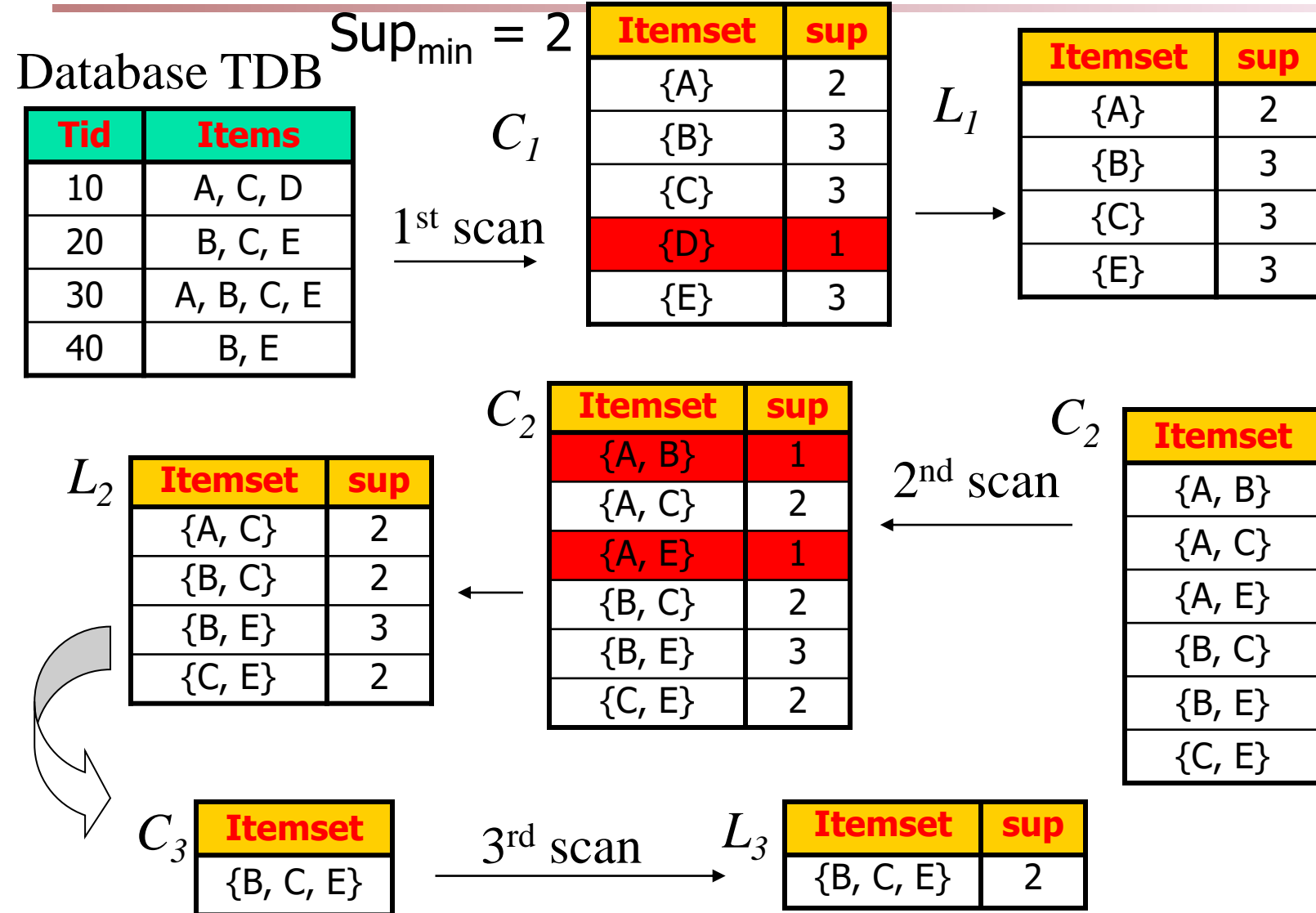# Apriori: A Candidate Generation & Test Approach

- Apriori pruning principle: If there is any itemset which is infrequent, its superset should not be generated/tested! Method:

  - Initially, scan DB once to get frequent 1-itemset

  - Generate length (k+1) candidate itemsets from length k frequent itemsets

  - Test the candidates against DB

  - Terminate when no frequent or candidate set can be generated

# The Apriori Algorithm—An Example

$Sup_{min} = 2$

Database TDB

| Tid | Items |
|-----|-------|
| 10 | A, C, D |
| 20 | B, C, E |
| 30 | A, B, C, E |
| 40 | B, E |

$1^{st}$ scan

$C_1$

| Itemset | sup |
|---------|-----|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {D} | 1 |
| {E} | 3 |

$L_1$

| Itemset | sup |
|---------|-----|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {E} | 3 |

$C_2$

| Itemset | sup |
|---------|-----|
| {A, B} | 1 |
| {A, C} | 2 |
| {A, E} | 1 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$2^{nd}$ scan

$C_2$

| Itemset |
|---------|
| {A, B} |
| {A, C} |
| {A, E} |
| {B, C} |
| {B, E} |
| {C, E} |

$L_2$

| Itemset | sup |
|---------|-----|
| {A, C} | 2 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$C_3$

| Itemset |
|---------|
| {B, C, E} |

$3^{rd}$ scan

$L_3$

| Itemset | sup |
|---------|-----|
| {B, C, E} | 2 |

41

# The Apriori Algorithm (Pseudo-Code)

$C_k$: Candidate itemset of size k

$L_k$ : frequent itemset of size k

$L_1$ = {frequent items};

**for** ($k$ = 1; $L_k$ !=$\varnothing$; $k$++) **do begin**

    $C_{k+1}$ = candidates generated from $L_k$;

    **for each** transaction $t$ in database do

       increment the count of all candidates in $C_{k+1}$ that are

        contained in $t$

    $L_{k+1}$  = candidates in $C_{k+1}$ with min_support

    **end**

**return** $\cup_k L_k$;

# Implementation of Apriori

- How to generate candidates?
    - Step 1: self-joining $L_k$
    - Step 2: pruning
- Example of Candidate-generation
    - $L_3=\{abc, abd, acd, ace, bcd\}$
    - Self-joining: $L_3*L_3$
        - *abcd* from *abc* and *abd*
        - *acde* from *acd* and *ace*
    - Pruning:
        - *acde* is removed because *ade* is not in $L_3$
    - $C_4 = \{abcd\}$