

---

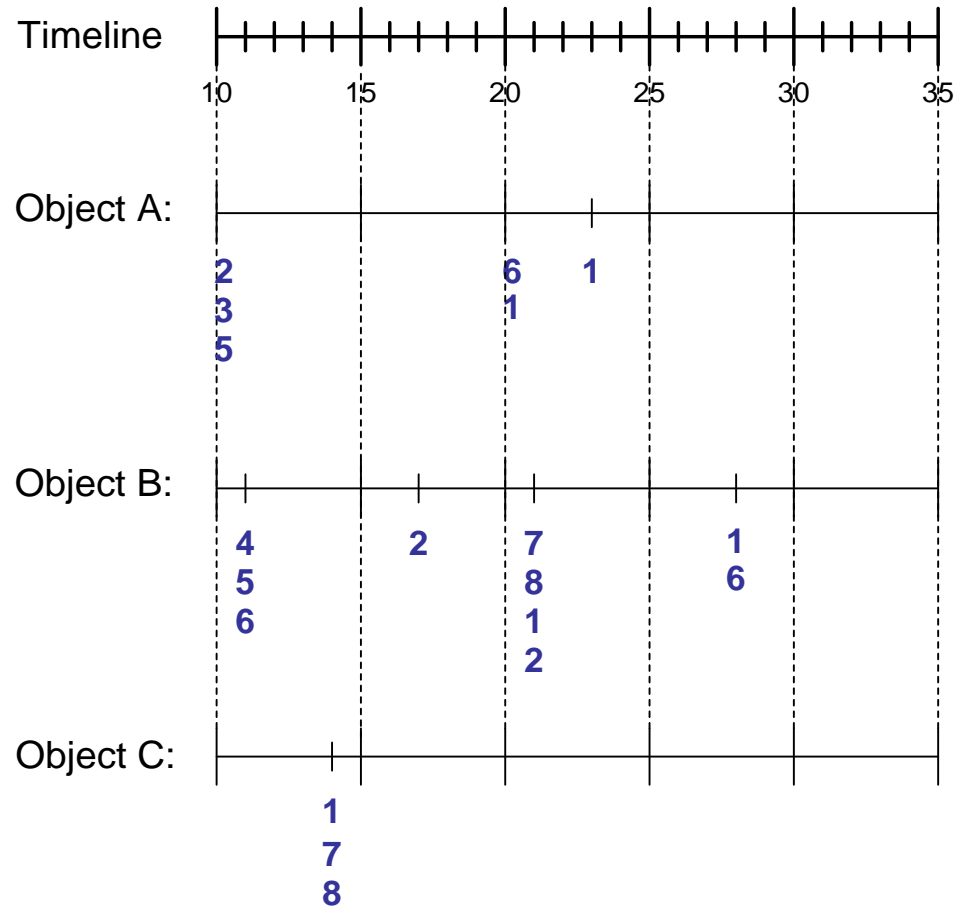
# **Association Rule Mining**

Dr. Faisal Kamiran

# Sequence Data

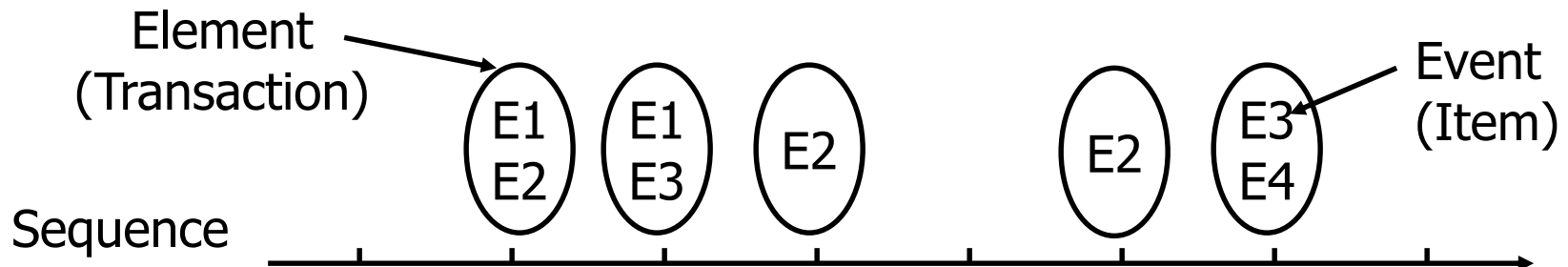
## Sequence Database:

Object	Timestamp	Events
A	10	2, 3, 5
A	20	6, 1
A	23	1
B	11	4, 5, 6
B	17	2
B	21	7, 8, 1, 2
B	28	1, 6
C	14	1, 8, 7



# Sequence Data

Sequence Database	Sequence	Element (Transaction)	Event (Item)
Customer	Purchase history of a given customer	A set of items bought by a customer at time $t$	Books, diary products, CDs, etc
Web Data	Browsing activity of a particular Web visitor	A collection of files viewed by a Web visitor after a single mouse click	Home page, index page, contact info, etc
Event data	History of events generated by a given sensor	Events triggered by a sensor at time $t$	Types of alarms generated by sensors
Genome sequences	DNA sequence of a particular species	An element of the DNA sequence	Bases A,T,G,C



# Formal Definition of a Sequence

---

- A sequence is an ordered list of elements (transactions)

$$S = \langle e_1 e_2 e_3 \dots \rangle$$

- Each element contains a collection of events (items)

$$e_i = \{i_1, i_2, \dots, i_k\}$$

- Each element is attributed to a specific time or location
- Length of a sequence,  $|s|$ , is given by the number of elements of the sequence
- A  $k$ -sequence is a sequence that contains  $k$  events (items)

# Examples of Sequence

---

## □ Web sequence:

< {Homepage} {Electronics} {Digital Cameras} {Canon Digital Camera}  
{Shopping Cart} {Order Confirmation} {Return to Shopping} >

## □ Sequence of books checked out at a library:

<{Fellowship of the Ring} {The Two Towers} {Return of the King}>

# Formal Definition of a Subsequence

- A sequence  $\langle a_1 a_2 \dots a_n \rangle$  is contained in another sequence  $\langle b_1 b_2 \dots b_m \rangle$  ( $m \geq n$ ) if there exist integers  $i_1 < i_2 < \dots < i_n$  such that  $a_1 \subseteq b_{i_1}$ ,  $a_2 \subseteq b_{i_2}$ , ...,  $a_n \subseteq b_{i_n}$

Data sequence	Subsequence	Contain?
$\langle \{2,4\} \{3,5,6\} \{8\} \rangle$	$\langle \{2\} \{3,5\} \rangle$	Yes
$\langle \{1,2\} \{3,4\} \rangle$	$\langle \{1\} \{2\} \rangle$	No
$\langle \{2,4\} \{2,4\} \{2,5\} \rangle$	$\langle \{2\} \{4\} \rangle$	Yes

- The support of a subsequence  $w$  is defined as the fraction of data sequences that contain  $w$
- A *sequential pattern* is a frequent subsequence (i.e., a subsequence whose support is  $\geq \text{minsup}$ )

# Sequential Pattern Mining: Definition

---

## □ Given:

- a database of sequences
- a user-specified minimum support threshold, *minsup*

## □ Task:

- Find all subsequences with support  $\geq \textit{minsup}$

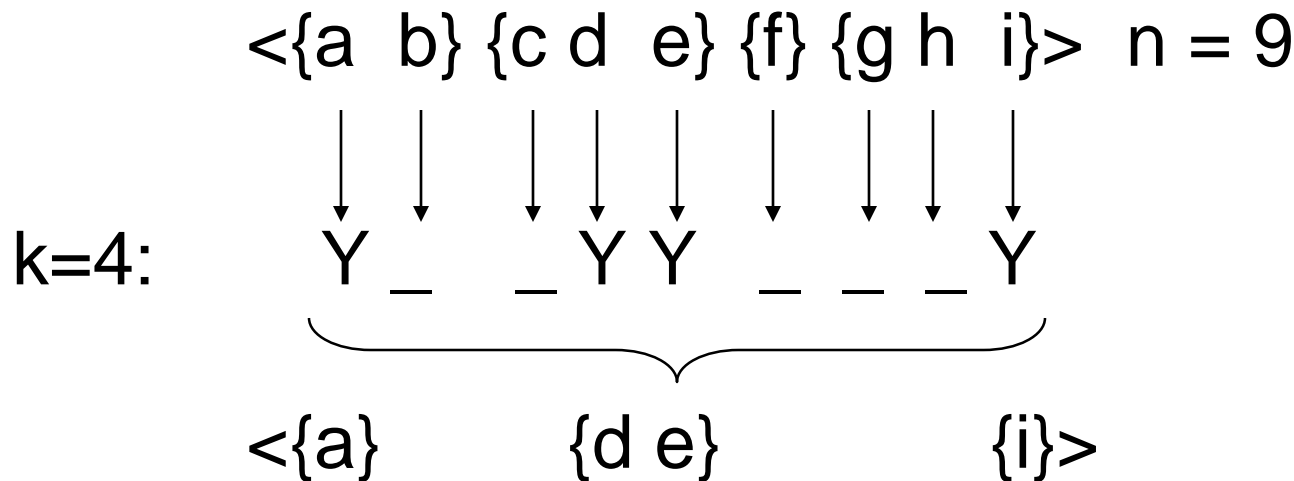
# Sequential Pattern Mining: Challenge

□ Given a sequence:  $\langle \{a\} \{b\} \{c\} \{d\} \{e\} \{f\} \{g\} \{h\} \{i\} \rangle$

— Examples of subsequences:

$\langle \{a\} \{c\} \{d\} \{f\} \{g\} \rangle$ ,  $\langle \{c\} \{d\} \{e\} \rangle$ ,  $\langle \{b\} \{g\} \rangle$ , etc.

□ How many k-subsequences can be extracted from a given n-sequence?



Answer :

$$\binom{n}{k} = \binom{9}{4} = 126$$



# Sequential Pattern Mining: Example

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

*Minsup* = 50%

**Examples of Frequent Subsequences:**

$\langle \{1,2\} \rangle$   $s=60\%$   
 $\langle \{2,3\} \rangle$   
 $\langle \{2,4\} \rangle$   
 $\langle \{3\} \{5\} \rangle$   
 $\langle \{1\} \{2\} \rangle$   
 $\langle \{2\} \{2\} \rangle$   
 $\langle \{1\} \{2,3\} \rangle$   
 $\langle \{2\} \{2,3\} \rangle$   
 $\langle \{1,2\} \{2,3\} \rangle$

# Sequential Pattern Mining: Example

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

*Minsup* = 50%

**Examples of Frequent Subsequences:**

< {1,2} >                      s=60%

< {2,3} >                      s=60%

< {2,4} >

< {3} {5} >

< {1} {2} >

< {2} {2} >

< {1} {2,3} >

< {2} {2,3} >

< {1,2} {2,3} >

# Sequential Pattern Mining: Example

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

*Minsup* = 50%

**Examples of Frequent Subsequences:**

< {1,2} >                      s=60%

< {2,3} >                      s=60%

< {2,4}>                      s=80%

< {3} {5}>

< {1} {2} >

< {2} {2} >

< {1} {2,3} >

< {2} {2,3} >

< {1,2} {2,3} >

# Sequential Pattern Mining: Example

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

*Minsup* = 50%

**Examples of Frequent Subsequences:**

< {1,2} >                      s=60%

< {2,3} >                      s=60%

< {2,4} >                      s=80%

< {3} {5} >                    s=80%

< {1} {2} >

< {2} {2} >

< {1} {2,3} >

< {2} {2,3} >

< {1,2} {2,3} >

# Sequential Pattern Mining: Example

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

*Minsup* = 50%

## Examples of Frequent Subsequences:

$\langle \{1,2\} \rangle$   $s=60\%$

$\langle \{2,3\} \rangle$   $s=60\%$

$\langle \{2,4\} \rangle$   $s=80\%$

$\langle \{3\} \{5\} \rangle$   $s=80\%$

$\langle \{1\} \{2\} \rangle$   $s=80\%$

$\langle \{2\} \{2\} \rangle$

$\langle \{1\} \{2,3\} \rangle$

$\langle \{2\} \{2,3\} \rangle$

$\langle \{1,2\} \{2,3\} \rangle$

# Sequential Pattern Mining: Example

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

*Minsup* = 50%

## Examples of Frequent Subsequences:

$\langle \{1,2\} \rangle$   $s=60\%$

$\langle \{2,3\} \rangle$   $s=60\%$

$\langle \{2,4\} \rangle$   $s=80\%$

$\langle \{3\} \{5\} \rangle$   $s=80\%$

$\langle \{1\} \{2\} \rangle$   $s=80\%$

$\langle \{2\} \{2\} \rangle$   $s=60\%$

$\langle \{1\} \{2,3\} \rangle$

$\langle \{2\} \{2,3\} \rangle$

$\langle \{1,2\} \{2,3\} \rangle$

# Sequential Pattern Mining: Example

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

*Minsup* = 50%

## Examples of Frequent Subsequences:

< {1,2} >	s=60%
< {2,3} >	s=60%
< {2,4}>	s=80%
< {3} {5}>	s=80%
< {1} {2} >	s=80%
< {2} {2} >	s=60%
< {1} {2,3} >	s=60%
< {2} {2,3} >	
< {1,2} {2,3} >	

# Sequential Pattern Mining: Example

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

*Minsup* = 50%

## Examples of Frequent Subsequences:

< {1,2} >	s=60%
< {2,3} >	s=60%
< {2,4}>	s=80%
< {3} {5}>	s=80%
< {1} {2} >	s=80%
< {2} {2} >	s=60%
< {1} {2,3} >	s=60%
< {2} {2,3} >	s=60%
< {1,2} {2,3} >	



# Sequential Pattern Mining: Example

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

*Minsup* = 50%

## Examples of Frequent Subsequences:

< {1,2} >	s=60%
< {2,3} >	s=60%
< {2,4}>	s=80%
< {3} {5}>	s=80%
< {1} {2} >	s=80%
< {2} {2} >	s=60%
< {1} {2,3} >	s=60%
< {2} {2,3} >	s=60%
< {1,2} {2,3} >	s=60%

# Extracting Sequential Patterns

□ Given  $n$  events:  $i_1, i_2, i_3, \dots, i_n$

□ Candidate 1-subsequences:

$\langle \{i_1\} \rangle, \langle \{i_2\} \rangle, \langle \{i_3\} \rangle, \dots, \langle \{i_n\} \rangle$

□ Candidate 2-subsequences:

$\langle \{i_1, i_2\} \rangle, \langle \{i_1, i_3\} \rangle, \dots, \langle \{i_1\} \{i_1\} \rangle, \langle \{i_1\} \{i_2\} \rangle, \dots, \langle \{i_{n-1}\} \{i_n\} \rangle$

□ Candidate 3-subsequences:

$\langle \{i_1, i_2, i_3\} \rangle, \langle \{i_1, i_2, i_4\} \rangle, \dots, \langle \{i_1, i_2\} \{i_1\} \rangle, \langle \{i_1, i_2\} \{i_2\} \rangle, \dots,$   
 $\langle \{i_1\} \{i_1, i_2\} \rangle, \langle \{i_1\} \{i_1, i_3\} \rangle, \dots, \langle \{i_1\} \{i_1\} \{i_1\} \rangle, \langle \{i_1\} \{i_1\} \{i_2\} \rangle, \dots$

# Generalized Sequential Pattern (GSP)

---

## □ Step 1:

- Make the first pass over the sequence database  $D$  to yield all the 1-element frequent sequences

## □ Step 2:

Repeat until no new frequent sequences are found

### – **Candidate Generation:**

- ◆ Merge pairs of frequent subsequences found in the  $(k-1)th$  pass to generate candidate sequences that contain  $k$  items

### – **Candidate Pruning:**

- ◆ Prune candidate  $k$ -sequences that contain infrequent  $(k-1)$ -subsequences

### – **Support Counting:**

- ◆ Make a new pass over the sequence database  $D$  to find the support for these candidate sequences

### – **Candidate Elimination:**

- ◆ Eliminate candidate  $k$ -sequences whose actual support is less than *minsup*

# Candidate Generation

---

## □ Base case ( $k=2$ ):

- Merging two frequent 1-sequences  $\langle\{i_1\}\rangle$  and  $\langle\{i_2\}\rangle$  will produce two candidate 2-sequences:  $\langle\{i_1\} \{i_2\}\rangle$  and  $\langle\{i_1 i_2\}\rangle$

## □ General case ( $k>2$ ):

- A frequent  $(k-1)$ -sequence  $w_1$  is merged with another frequent  $(k-1)$ -sequence  $w_2$  to produce a candidate  $k$ -sequence if the subsequence obtained by removing the first event in  $w_1$  is the same as the subsequence obtained by removing the last event in  $w_2$ 
  - ◆ The resulting candidate after merging is given by the sequence  $w_1$  extended with the last event of  $w_2$ .
    - If the last two events in  $w_2$  belong to the same element, then the last event in  $w_2$  becomes part of the last element in  $w_1$
    - Otherwise, the last event in  $w_2$  becomes a separate element appended to the end of  $w_1$

# Candidate Generation Examples

---

- Merging the sequences

$w_1 = \langle \{1\} \{2\ 3\} \{4\} \rangle$  and  $w_2 = \langle \{2\ 3\} \{4\ 5\} \rangle$

will produce the candidate sequence  $\langle \{1\} \{2\ 3\} \{4\ 5\} \rangle$  because the last two events in  $w_2$  (4 and 5) belong to the same element

- Merging the sequences

$w_1 = \langle \{1\} \{2\ 3\} \{4\} \rangle$  and  $w_2 = \langle \{2\ 3\} \{4\} \{5\} \rangle$

will produce the candidate sequence  $\langle \{1\} \{2\ 3\} \{4\} \{5\} \rangle$  because the last two events in  $w_2$  (4 and 5) do not belong to the same element

- We do not have to merge the sequences

$w_1 = \langle \{1\} \{2\} \{3\} \rangle$  and  $w_2 = \langle \{1\} \{2,5\} \rangle$

# GSP Example

Frequent  
3-sequences

< {1} {2} {3} >  
< {1} {2 5} >  
< {1} {5} {3} >  
< {2} {3} {4} >  
< {2 5} {3} >  
< {3} {4} {5} >  
< {5} {3 4} >

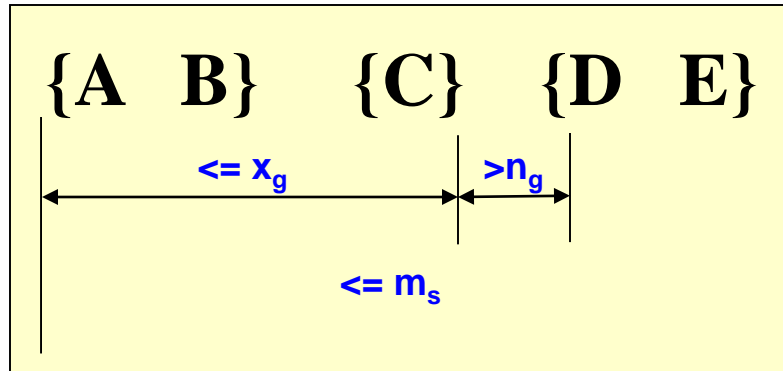
Candidate  
Generation

< {1} {2} {3} {4} >  
< {1} {2 5} {3} >  
< {1} {5} {3 4} >  
< {2} {3} {4} {5} >  
< {2 5} {3 4} >

Candidate  
Pruning

< {1} {2 5} {3} >

# Timing Constraints (I)



$x_g$ : max-gap

$n_g$ : min-gap

$$x_g = 2, n_g = 0$$

Data sequence	Subsequence	Contain?
$\langle \{2,4\} \{3,5,6\} \{4,7\} \{4,5\} \{8\} \rangle$	$\langle \{6\} \{5\} \rangle$	Yes
$\langle \{1\} \{2\} \{3\} \{4\} \{5\} \rangle$	$\langle \{1\} \{4\} \rangle$	No
$\langle \{1\} \{2,3\} \{3,4\} \{4,5\} \rangle$	$\langle \{2\} \{3\} \{5\} \rangle$	Yes
$\langle \{1,2\} \{3\} \{2,3\} \{3,4\} \{2,4\} \{4,5\} \rangle$	$\langle \{1,2\} \{5\} \rangle$	No

# Timing Constraints (I)

- Maxgap = 3
- Mingap = 1

Data Sequence, $s$	Sequential Pattern, $t$	$maxgap$	$mingap$
$\langle \{1,3\} \{3,4\} \{4\} \{5\} \{6,7\} \{8\} \rangle$	$\langle \{3\} \{6\} \rangle$		
$\langle \{1,3\} \{3,4\} \{4\} \{5\} \{6,7\} \{8\} \rangle$	$\langle \{6\} \{8\} \rangle$		
$\langle \{1,3\} \{3,4\} \{4\} \{5\} \{6,7\} \{8\} \rangle$	$\langle \{1,3\} \{6\} \rangle$		
$\langle \{1,3\} \{3,4\} \{4\} \{5\} \{6,7\} \{8\} \rangle$	$\langle \{1\} \{3\} \{8\} \rangle$		



# Timing Constraints (I)

- Maxgap = 3
- Mingap = 1

Data Sequence, $s$	Sequential Pattern, $t$	$maxgap$	$mingap$
$\langle \{1,3\} \{3,4\} \{4\} \{5\} \{6,7\} \{8\} \rangle$	$\langle \{3\} \{6\} \rangle$	Pass	Pass
$\langle \{1,3\} \{3,4\} \{4\} \{5\} \{6,7\} \{8\} \rangle$	$\langle \{6\} \{8\} \rangle$	Pass	Fail
$\langle \{1,3\} \{3,4\} \{4\} \{5\} \{6,7\} \{8\} \rangle$	$\langle \{1,3\} \{6\} \rangle$	Fail	Pass
$\langle \{1,3\} \{3,4\} \{4\} \{5\} \{6,7\} \{8\} \rangle$	$\langle \{1\} \{3\} \{8\} \rangle$	Fail	Fail

# Mining Sequential Patterns with Timing Constraints

---

## □ Approach 1:

- Mine sequential patterns without timing constraints
- Postprocess the discovered patterns

## □ Approach 2:

- Modify GSP to directly prune candidates that violate timing constraints
- Question:
  - ◆ Does Apriori principle still hold?

# Apriori Principle for Sequence Data

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

Suppose:

$$x_g = 1 \text{ (max-gap)}$$

$$n_g = 0 \text{ (min-gap)}$$

$$\text{minsup} = 60\%$$

$$\langle \{2\} \{5\} \rangle \text{ support} = 40\%$$

but

$$\langle \{2\} \{3\} \{5\} \rangle \text{ support} = 60\%$$

Problem exists because of max-gap constraint

No such problem if max-gap is infinite

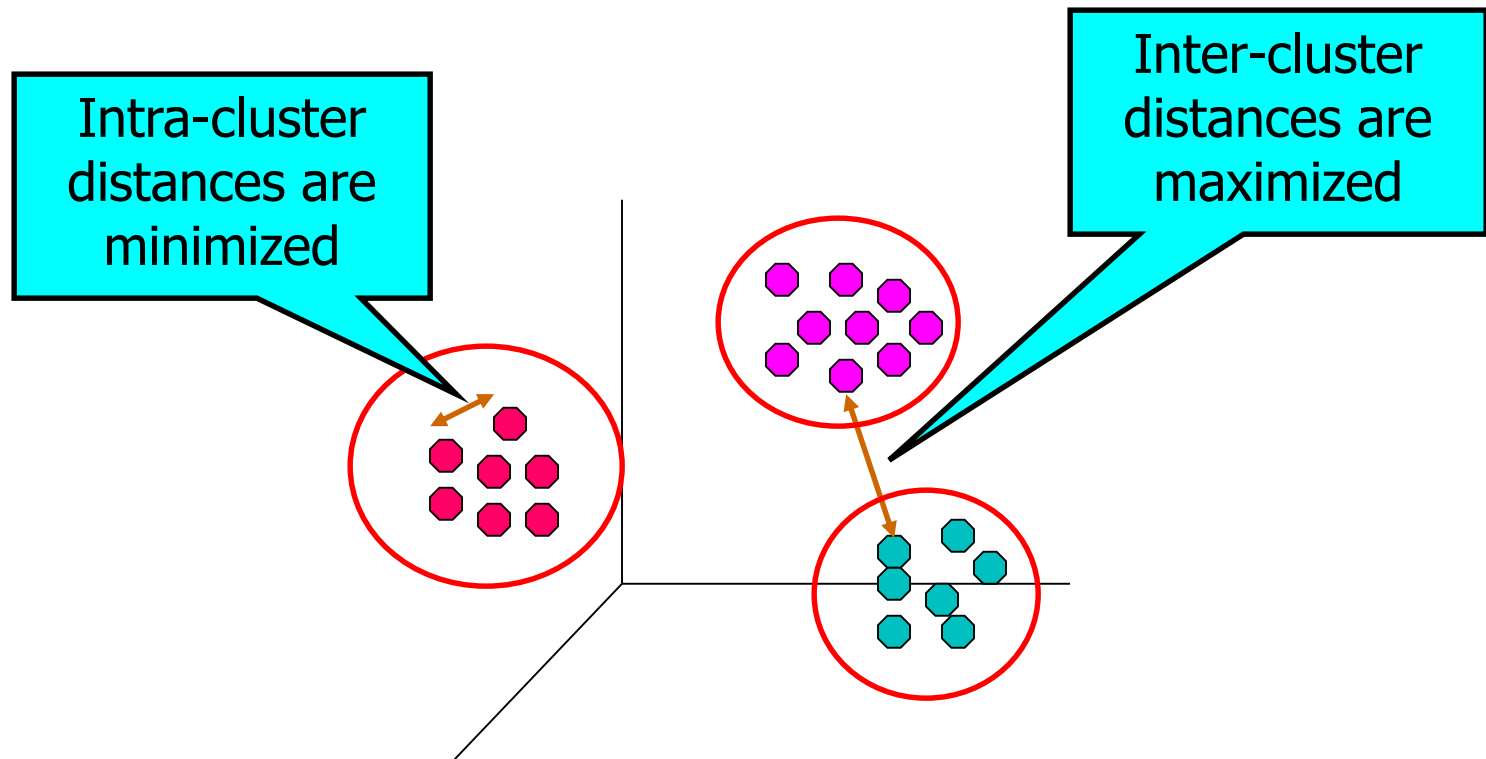
---

# Clustering

Dr. Faisal Kamiran

# What is 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



# Applications of Cluster Analysis

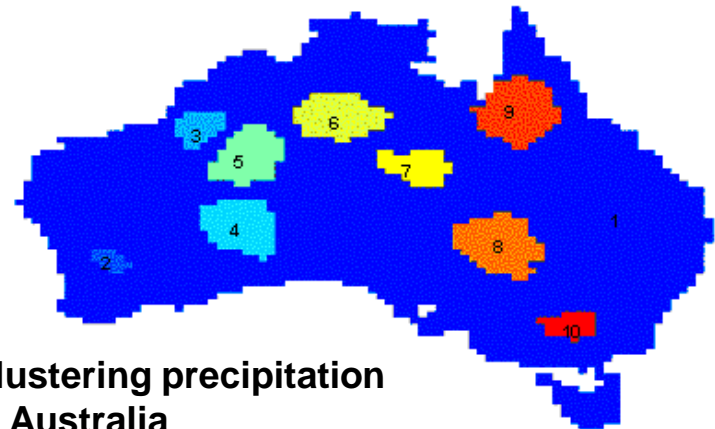
## □ Understanding

- Group related documents for browsing, group genes and proteins that have similar functionality, or group stocks with similar price fluctuations

	<i>Discovered Clusters</i>	<i>Industry Group</i>
<b>1</b>	Applied-Matl-DOWN,Bay-Network-Down,3-COM-DOWN, Cabletron-Sys-DOWN,CISCO-DOWN,HP-DOWN, DSC-Comm-DOWN,INTEL-DOWN,LSI-Logic-DOWN, Micron-Tech-DOWN,Texas-Inst-Down,Tellabs-Inc-Down, Natl-Semiconduct-DOWN,Oracl-DOWN,SGI-DOWN, Sun-DOWN	Technology1-DOWN
<b>2</b>	Apple-Comp-DOWN,Autodesk-DOWN,DEC-DOWN, ADV-Micro-Device-DOWN,Andrew-Corp-DOWN, Computer-Assoc-DOWN,Circuit-City-DOWN, Compaq-DOWN, EMC-Corp-DOWN, Gen-Inst-DOWN, Motorola-DOWN,Microsoft-DOWN,Scientific-Atl-DOWN	Technology2-DOWN
<b>3</b>	Fannie-Mae-DOWN,Fed-Home-Loan-DOWN, MBNA-Corp-DOWN,Morgan-Stanley-DOWN	Financial-DOWN
<b>4</b>	Baker-Hughes-UP,Dresser-Inds-UP,Halliburton-HLD-UP, Louisiana-Land-UP,Phillips-Petro-UP,Unocal-UP, Schlumberger-UP	Oil-UP

## □ Summarization

- Reduce the size of large data sets



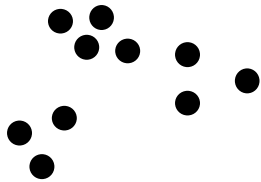
Clustering precipitation  
in Australia

# What is not Cluster Analysis?

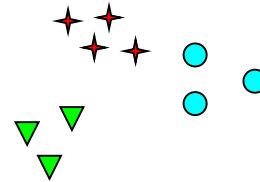
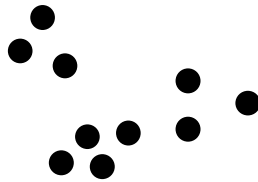
---

- Supervised classification
  - Have class label information
- Simple segmentation
  - Dividing students into different registration groups alphabetically, by last name
- Results of a query
  - Groupings are a result of an external specification

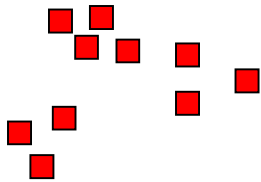
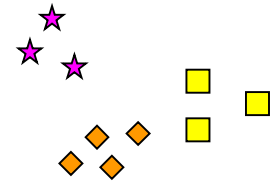
# Notion of a Cluster can be Ambiguous



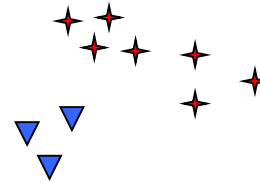
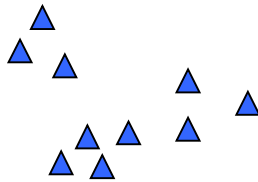
How many clusters?



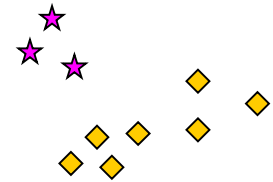
Six Clusters



Two Clusters



Four Clusters





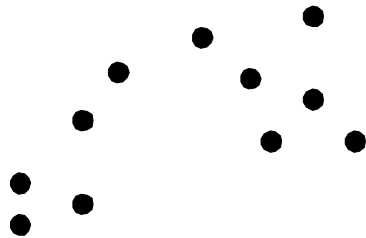
# Types of Clusterings

---

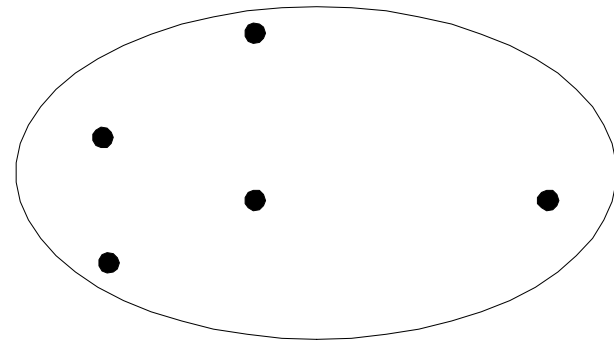
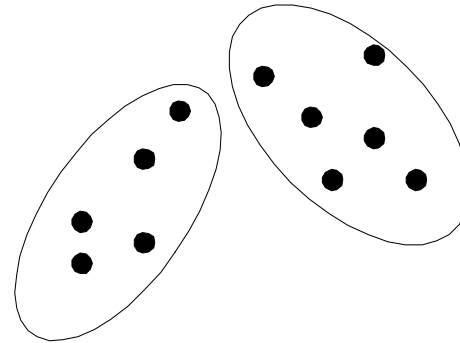
- A **clustering** is a set of clusters
- Important distinction between **hierarchical** and **partitional** sets of clusters
- Partitional Clustering
  - A division of data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset
- Hierarchical clustering
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

# Partitional Clustering

---



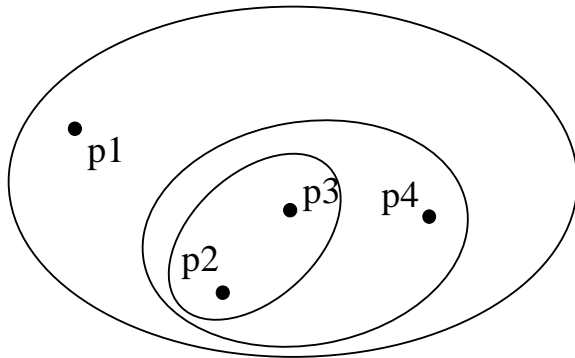
**Original Points**



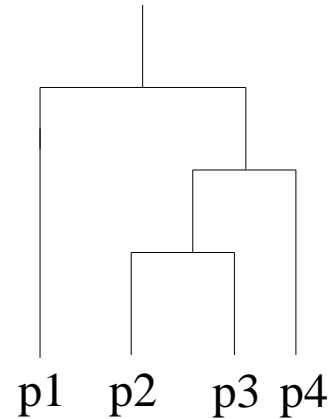
**A Partitional Clustering**

# Hierarchical Clustering

---



**Traditional Hierarchical Clustering**



**Traditional Dendrogram**

# Other Distinctions Between Sets of Clusters

---

## □ Exclusive versus non-exclusive

- In non-exclusive clusterings, points may belong to multiple clusters.
- Can represent multiple classes or ‘border’ points

## □ Fuzzy versus non-fuzzy

- In fuzzy clustering, a point belongs to every cluster with some weight between 0 and 1
- Weights must sum to 1
- Probabilistic clustering has similar characteristics

## □ Partial versus complete

- In some cases, we only want to cluster some of the data

## □ Heterogeneous versus homogeneous

- Cluster of widely different sizes, shapes, and densities

# Types of Clusters

---

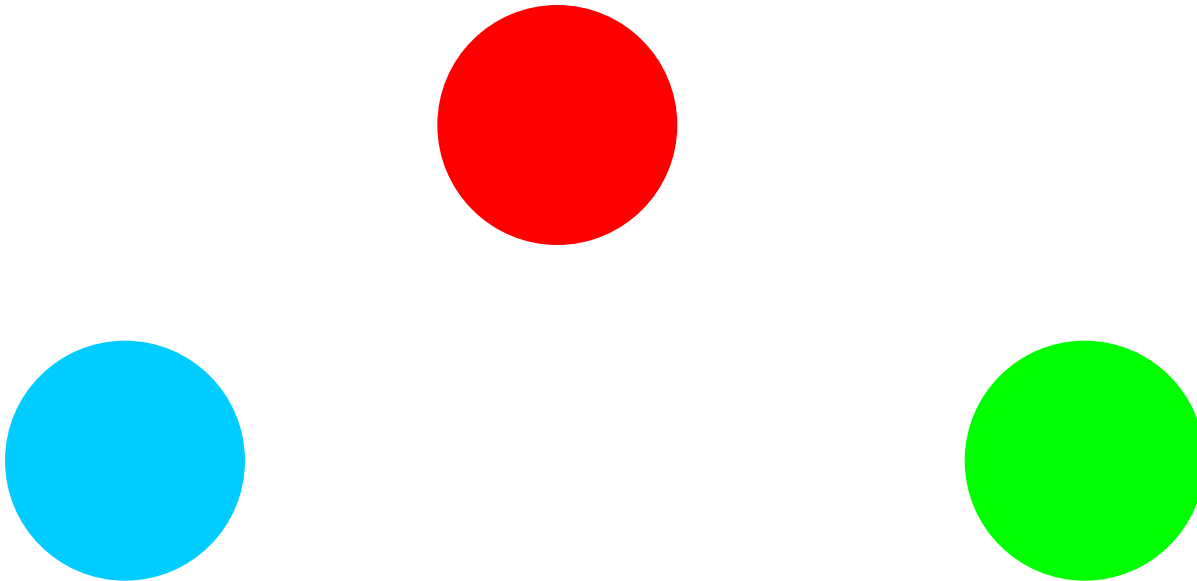
- Well-separated clusters
- Center-based clusters
- Contiguous clusters
- Density-based clusters
- Property or Conceptual
- Described by an Objective Function

# Types of Clusters: Well-Separated

---

## □ Well-Separated Clusters:

- A cluster is a set of points such that any point in a cluster is closer (or more similar) to every other point in the cluster than to any point not in the cluster.



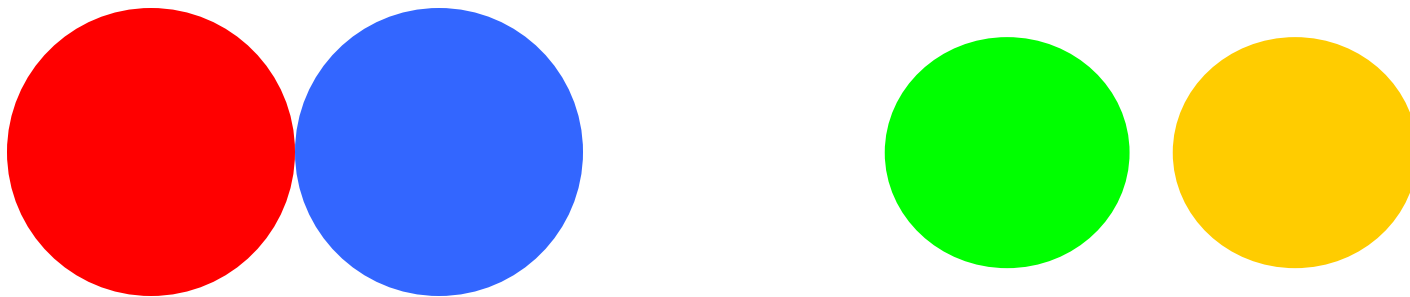
**3 well-separated clusters**

# Types of Clusters: Center-Based

---

## □ Center-based

- A cluster is a set of objects such that an object in a cluster is closer (more similar) to the “center” of a cluster, than to the center of any other cluster
- The center of a cluster is often a **centroid**, the average of all the points in the cluster, or a **medoid**, the most “representative” point of a cluster

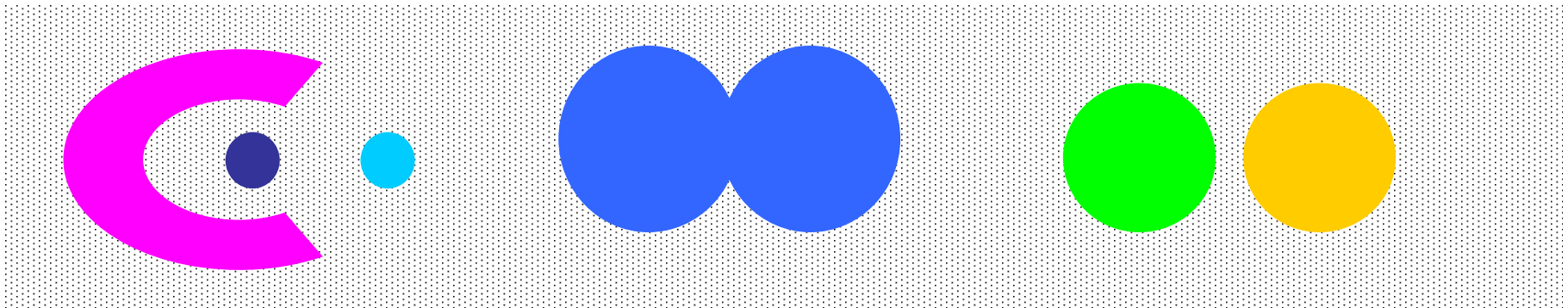


**4 center-based clusters**

# Types of Clusters: Density-Based

## □ Density-based

- A cluster is a dense region of points, which is separated by low-density regions, from other regions of high density.
- Used when the clusters are irregular or intertwined, and when noise and outliers are present.



**6 density-based clusters**



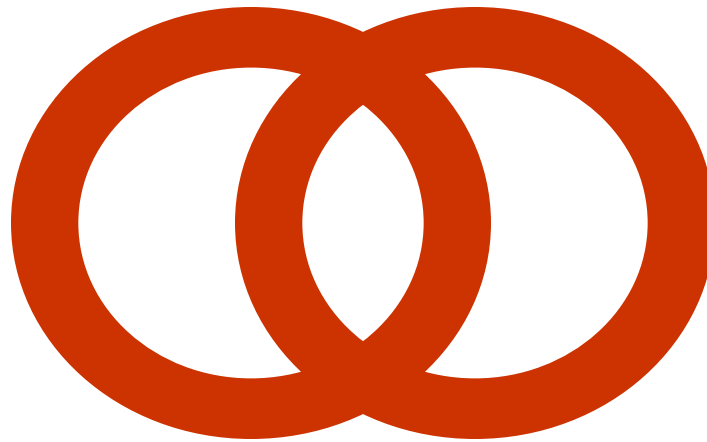
# Types of Clusters: Conceptual Clusters

---

## □ Shared Property or Conceptual Clusters

- Finds clusters that share some common property or represent a particular concept.

.



**2 Overlapping Circles**

# Types of Clusters: Objective Function

---

## □ Clusters Defined by an Objective Function

- Finds clusters that minimize or maximize an objective function.
- Enumerate all possible ways of dividing the points into clusters and evaluate the 'goodness' of each potential set of clusters by using the given objective function. (NP Hard)
- Can have global or local objectives.
  - ◆ Hierarchical clustering algorithms typically have local objectives
  - ◆ Partitional algorithms typically have global objectives

# Characteristics of the Input Data Are Important

---

- Type of proximity or density measure
  - This is a derived measure, but central to clustering
- Sparseness
  - Dictates type of similarity
  - Adds to efficiency
- Attribute type
  - Dictates type of similarity
- Type of Data
  - Dictates type of similarity
  - Other characteristics, e.g., autocorrelation
- Dimensionality
- Noise and Outliers
- Type of Distribution

# Clustering Algorithms

---

- K-means and its variants
- Hierarchical clustering
- Density-based clustering