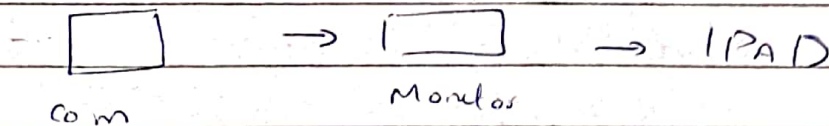## GSP

Data Mining : extract useful knowledge from data

data. discrete sequence

A sequence is an ordered list of symbols

□ → □ → IPAD
com      Monitor

sequences

1 → go → back → home

The goal is to find all subsequences that appear frequently in a set of discrete sequences

4 Items

I = {a, b, c, d, e}

a → apple  b = bread  , c = cake  d = dodes
                                    e = eggs

An itemset is set of items that is a subset of
                                    IMP    T

{a, b, c}    {d, e}    | An itemset cannot
                       | have doublicates
                       | items
                       | a same item two

A discrete sequence is an ordered
list of itemsets $S = \langle x_1, x_2, \ldots, x_n \rangle$
where $x_j \subset I$ for any $j \in \{1, 2, n\}$

Example
$\langle \{a, b\}, \{c\} \rangle$     $a, b \to c$

It means customer purchased $a, b$
at same time     and then purchased
cake

$\langle \{a\}, \{a\}, \{c\} \rangle$

$a \to a \to c$

subsequence
$\langle \{a, c\} \rangle \subseteq \langle \{a, b, c\} \rangle$
$\langle \{a, c\} \rangle \not\subseteq \langle \{a\}, \{c\} \rangle$

$\langle \{a\}, \{c\} \rangle \subseteq \langle \{a, b\}, \{d\}, \{b, c\} \rangle$

$\varnothing$

$$D = \{S_1, S_2, S_3, S_4\}$$

Input      Sequence Database

$$S_1 = \langle \{a, b\}, \{c\}, \{a\} \rangle$$
$$S_2 = \langle \{a\}, \{b\}, \{c\} \rangle$$
$$S_3 = \langle \{b\}, \{c\}, \{d\} \rangle$$
$$S_4 = \langle \{b\}, \{a, b\}, \{c\} \rangle$$

$$\text{Support} (\langle \{a\} \rangle) = 3$$
$$\text{sup} (\langle \{b\} \rangle) = 4$$
$$\text{sup} (\langle \{a\}, \{b\} \rangle) = 1$$
$$\text{sup} (\langle \{a, b\} \rangle) = 2$$

Output

$$\text{minsup} = 3$$

$$\langle \{a\} \rangle \qquad\qquad \text{sup} = 3$$
$$\langle \{b\} \rangle \qquad\qquad \text{sup} = 4$$
$$\langle \{c\} \rangle \qquad\qquad \text{sup} = 4$$
$$\langle \{a\}, \{c\} \rangle \qquad\qquad \text{sup} = 3$$
$$\langle \{b\}, \{c\} \rangle \qquad\qquad \text{sup} = 4$$
$$\langle \{a, b\}, \{c\} \rangle \qquad\qquad \text{sup} = 3$$

$\frac{di}{tb} \times 5$

$30$

1sequence

{1} {2} {3}   {4} {5}

2sequence

a. 25   <{1} {2}> , {

| 3-sequence | Merge for candidate seven |
|---|---|
| < {1} {2} {3} > | |
| < {1} {2 5} > | < {1} {2} {3} > |
| < {1} {5} {3} > | < {2} {3} {4} > |
| < {2} {3} {4} > | < {1} {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 Generation

Candidate Pruning

# Candidate Generation



**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

$$\leftarrow \{1\} \{2\} \{3\} \{4\} \rightarrow$$

find support

$$< \{1\} \{2\} [3] \{45\} >$$

$$< \{1\} \{2\} \{3\} \{4\} > \qquad < \{2\} \{3\} \{4\} >$$

$$< \{1\} \{2\} \{3\} \{4\} > \qquad < \{1\} \{3\} \{3\} \{4\} >$$

$$< \{12\} \{2\} \{3\} \{4\} > \qquad < \{1\} \{2\} \{4\} >$$

$$< \{1\} \{2\} \{3\} \{4\} > \qquad \{\{1\} \{2\} \{3\}\}$$

Pruning Imppoint

$$\{1, 2\} \{4\}$$

$$\{2\} \{4\} \{6\}$$

$$\{1, 2\} \{4\} \{6\} \checkmark \qquad \left( \times \{1\} \{2\} \{4\} \{6\} \right)$$

# Candidate Pruning

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

**Candidate Generation**

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

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

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

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

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

# Candidate Pruning

<{1} {2} {3} {4}>    <{1} {2} {3} {4}>      <{2} {3} {4}> ✓
                           <{1} {2} {3} {4}>      <{1} {3} {4}> ✗
                           <{1} {2} {3} {4}>      <{1} {2} {4}> ✗
                           <{1} {2} {3} {4}>      <{1} {2} {3}> ✓

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

## Candidate Generation

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

<{1} {5} {3 4}>    <{1} {5} {3 4}>      <{5} {3 4}> ✓
                           <{1} {5} {3 4}>      <{1} {3 4}> ✗
                           <{1} {5} {3 4}>      <{1} {5} {4}> ✗
                           <{1} {5} {3 4}>      <{1} {5} {3}> ✓

<{2} {3} {4} {5}>    <{2} {3} {4} {5}>      <{3} {4} {5}> ✓
                           <{2} {3} {4} {5}>      <{2} {4} {5}> ✗
                           <{2} {3} {4} {5}>      <{2} {3} {5}> ✗
                           <{2} {3} {4} {5}>      <{2} {3} {4}> ✓

<{2 5} {3 4}>    <{2 5} {3 4}>      <{5} {3 4}> ✓
                           <{2 5} {3 4}>      <{2} {3 4}> ✗
                           <{2 5} {3 4}>      <{2 5} {4}> ✗
                           <{2 5} {3 4}>      <{2 5} {3}> ✓

## 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 Pruning

**Candidate Generation**

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

---

< {1} {2} {3} {4} >  < {1} {2} {3} {4} >      < {2} {3} {4} > ✓
                      < {1} {2} {3} {4} >      < {1} {3} {4} > ✗
                      < {1} {2} {3} {4} >      < {1} {2} {4} > ✗
                      < {1} {2} {3} {4} >      < {1} {2} {3} > ✓

---

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

---

< {1} {5} {3 4} >  < {1} {5} {3 4} >      < {5} {3 4} > ✓
                   < {1} {5} {3 4} >      < {1} {3 4} > ✗
                   < {1} {5} {3 4} >      < {1} {5} {4} > ✗
                   < {1} {5} {3 4} >      < {1} {5} {3} > ✓

---

< {2} {3} {4} {5} >  < {2} {3} {4} {5} >      < {3} {4} {5} > ✓
                     < {2} {3} {4} {5} >      < {2} {4} {5} > ✗
                     < {2} {3} {4} {5} >      < {2} {3} {5} > ✗
                     < {2} {3} {4} {5} >      < {2} {3} {4} > ✓

---

< {2 5} {3 4} >  < {2 5} {3 4} >      < {5} {3 4} > ✓
                 < {2 5} {3 4} >      < {2} {3 4} > ✗
                 < {2 5} {3 4} >      < {2 5} {4} > ✗
                 < {2 5} {3 4} >      < {2 5} {3} > ✓

---

**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 Pruning**

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

---