# Unit 4:
# Association Rules

# Section 2:
# Advanced Concepts and Algorithms

# Continuous and Categorical Attributes

How to apply association analysis formulation to non-asymmetric binary variables?

| Session Id | Country | Session Length (sec) | Number of Web Pages viewed | Gender | Browser Type | Buy |
|---|---|---|---|---|---|---|
| 1 | USA | 982 | 8 | Male | IE | No |
| 2 | China | 811 | 10 | Female | Netscape | No |
| 3 | USA | 2125 | 45 | Female | Mozilla | Yes |
| 4 | Germany | 596 | 4 | Male | IE | Yes |
| 5 | Australia | 123 | 9 | Male | Mozilla | No |
| … | … | … | … | … | … | … |

**Example of Association Rule:**

{Number of Pages $\in[5,10)$ ^ (Browser = Mozilla)} → {Buy = No}

# Handling Categorical Attributes

✗ TRANSFORM CATEGORICAL ATTRIBUTE INTO ASYMMETRIC BINARY VARIABLES

✗ INTRODUCE A NEW "ITEM" FOR EACH DISTINCT ATTRIBUTE-VALUE PAIR

  ✗ Example: replace Browser Type attribute with

    ✗ Browser Type = Internet Explorer
    ✗ Browser Type = Mozilla
    ✗ Browser Type = Mozilla

# Handling Categorical Attributes

- Potential Issues
    - What if attribute has many possible values
        - Example: attribute country has more than 200 possible values
        - Many of the attribute values may have very low support
            - Potential solution: Aggregate the low-support attribute values

    - What if distribution of attribute values is highly skewed
        - Example: 95% of the visitors have Buy = No
        - Most of the items will be associated with (Buy=No) item
            - Potential solution: drop the highly frequent items

# Handling Continuous Attributes

✗ DIFFERENT KINDS OF RULES:

  ✗ Age $\in$ [21,35) ^ Salary $\in$ [70k,120k) → Buy

  ✗ Salary $\in$ [70k,120k) ^ Buy → Age: $\mu$=28, $\sigma$=4

✗ DIFFERENT METHODS:

  ✗ Discretization-based

  ✗ Statistics-based

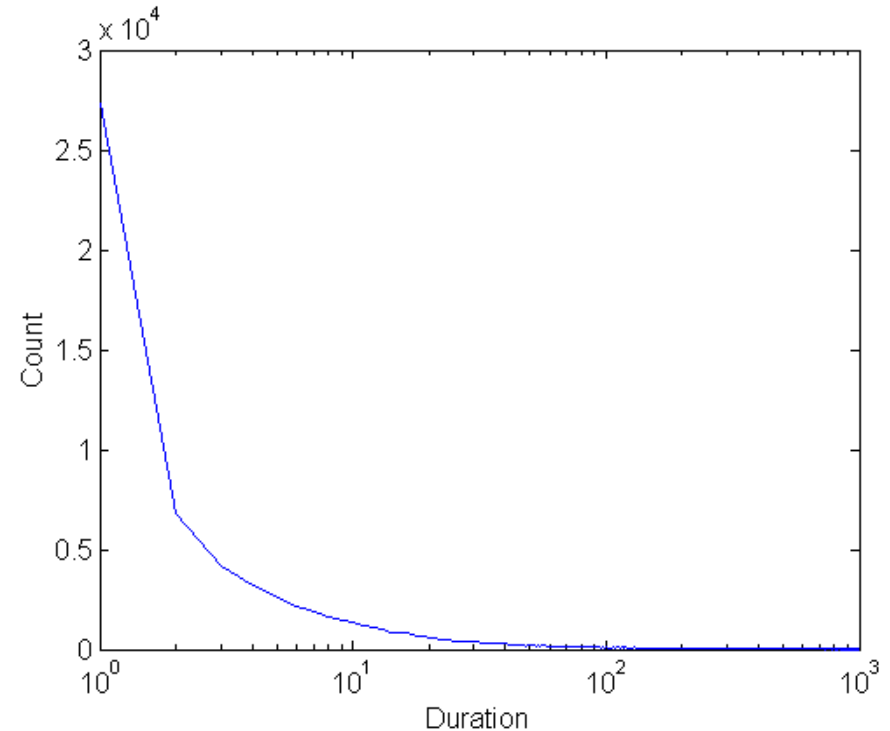  ✗ Non-discretization based

      ✗ minApriori

# Handling Continuous Attributes

✗ Use discretization
✗ Unsupervised:
  ✗ Equal-width binning
  ✗ Equal-depth binning
  ✗ Clustering

✗ Supervised:



Attribute values, v

| Class | $v_1$ | $v_2$ | $v_3$ | $v_4$ | $v_5$ | $v_6$ | $v_7$ | $v_8$ | $v_9$ |
|---|---|---|---|---|---|---|---|---|---|
| Anomalous | 0 | 0 | 20 | 10 | 20 | 0 | 0 | 0 | 0 |
| Normal | 150 | 100 | 0 | 0 | 0 | 100 | 100 | 150 | 100 |

$bin_1$　　　　$bin_2$　　　　$bin_3$

✘SIZE OF THE DISCRETIZED INTERVALS AFFECT SUPPORT & CONFIDENCE

$\{Refund = No, (Income = \$51,250)\} \rightarrow \{Cheat = No\}$

$\{Refund = No, (60K \leq Income \leq 80K)\} \rightarrow \{Cheat = No\}$

$\{Refund = No, (0K \leq Income \leq 1B)\} \rightarrow \{Cheat = No\}$

✘If intervals too small

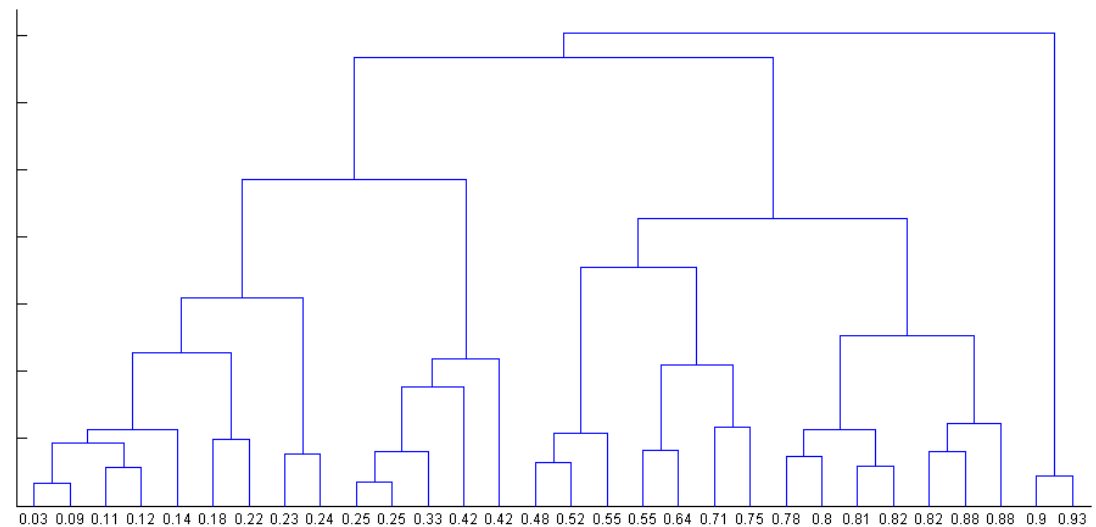  ✘ may not have enough support

✘If intervals too large

  ✘ may not have enough confidence

✘POTENTIAL SOLUTION: USE ALL POSSIBLE INTERVALS

✗EXECUTION TIME

✗If intervals contain n values, there are on average $O(n^2)$ possible ranges

✗TOO MANY RULES



0.03 0.09 0.11 0.12 0.14 0.18 0.22 0.23 0.24 0.25 0.25 0.33 0.42 0.42 0.48 0.52 0.55 0.55 0.64 0.71 0.75 0.78 0.8 0.81 0.82 0.82 0.88 0.88 0.9 0.93

{Refund = No, (Income = \$51,250)} → {Cheat = No}

{Refund = No, (51K ≤ Income ≤ 52K)} → {Cheat = No}

{Refund = No, (50K ≤ Income ≤ 60K)} → {Cheat = No}

# Approach by Srikant & Agrawal

✗ PREPROCESS THE DATA

    ✗ Discretize attribute using equi-depth partitioning

        ✗ Use *partial completeness measure* to determine number of partitions

        ✗ Merge adjacent intervals as long as support is less than max-support

✗ APPLY EXISTING ASSOCIATION RULE MINING ALGORITHMS

✗ DETERMINE INTERESTING RULES IN THE OUTPUT

✗ DISCRETIZATION WILL LOSE INFORMATION

**Approximated X**



X

✗ Use *partial completeness measure* to determine how much information is lost

C: FREQUENT ITEMSETS OBTAINED BY CONSIDERING ALL RANGES OF ATTRIBUTE VALUES

P: FREQUENT ITEMSETS OBTAINED BY CONSIDERING ALL RANGES OVER THE PARTITIONS

P IS *K-COMPLETE* W.R.T C IF $P \subseteq C$, AND $\forall X \in C, \exists X' \in P$ SUCH THAT:

1. X' IS A GENERALIZATION OF X AND SUPPORT $(X') \leq K \times$ SUPPORT$(X)$   $(K \geq 1)$
2. $\forall Y \subseteq X, \exists Y' \subseteq X'$ SUCH THAT SUPPORT $(Y') \leq K \times$ SUPPORT$(Y)$

GIVEN *K (PARTIAL COMPLETENESS LEVEL)*, CAN DETERMINE NUMBER OF INTERVALS (N)

# Interestingness Measure

{Refund = No, (Income = $51,250)} → {Cheat = No}

{Refund = No, (51K ≤ Income ≤ 52K)} → {Cheat = No}

{Refund = No, (50K ≤ Income ≤ 60K)} → {Cheat = No}

✗GIVEN AN ITEMSET: Z = {Z$_1$, Z$_2$, …, Z$_K$} AND ITS GENERALIZATION Z' = {Z$_1$', Z$_2$', …, Z$_K$'}

P(Z): SUPPORT OF Z
E$_{Z'}$(Z): EXPECTED SUPPORT OF Z BASED ON Z'

$$E_{z'}(Z) = \frac{P(z_1)}{P(z_1')} \times \frac{P(z_2)}{P(z_2')} \times \cdots \times \frac{P(z_k)}{P(z_k')} \times P(Z')$$

✗Z is R-interesting w.r.t. Z' if P(Z) ≥ R × E$_{Z'}$(Z)

✗For S: X → Y, and its generalization S': X' → Y'

      P(Y|X): confidence of X → Y

      P(Y'|X'): confidence of X' → Y'

      $E_{S'}$(Y|X): expected support of Z based on Z'

$$E(Y|X) = \frac{P(y_1)}{P(y_1{}')} \times \frac{P(y_2)}{P(y_2{}')} \times \cdots \times \frac{P(y_k)}{P(y_k{}')} \times P(Y'|X')$$

✗Rule S is R-interesting w.r.t its ancestor rule S' if

  ✗Support, P(S) ≥ R × $E_{S'}$(S) or

  ✗Confidence, P(Y|X) ≥ R × $E_{S'}$(Y|X)

# Statistics-based Methods

✗ EXAMPLE:

$$Browser=Mozilla \wedge Buy=Yes \rightarrow Age: \mu=23$$

✗ RULE CONSEQUENT CONSISTS OF A CONTINUOUS VARIABLE, CHARACTERIZED BY THEIR STATISTICS

- ✗ mean, median, standard deviation, etc.

✗ APPROACH:

- ✗ Withhold the target variable from the rest of the data
- ✗ Apply existing frequent itemset generation on the rest of the data
- ✗ For each frequent itemset, compute the descriptive statistics for the corresponding target variable
    - ✗ Frequent itemset becomes a rule by introducing the target variable as rule consequent
- ✗ Apply statistical test to determine interestingness of the rule

✗HOW TO DETERMINE WHETHER AN ASSOCIATION RULE INTERESTING?

   ✗Compare the statistics for segment of population covered by the rule vs segment of population not covered by the rule:

$$A \Rightarrow B: \mu \quad \text{versus} \quad A \Rightarrow B: \mu'$$

   ✗Statistical hypothesis testing:

      ✗ Null hypothesis:  H0: $\mu' = \mu + \Delta$

      ✗ Alternative hypothesis: H1: $\mu' > \mu + \Delta$

      ✗ Z has zero mean and variance 1 under null hypothesis

$$Z = \frac{\mu' - \mu - \Delta}{\sqrt{\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}}}$$

# Statistics-based Methods

✗Eхамрle:

r: Browser=Mozilla ^ Buy=Yes → Age: μ=23

✗Rule is interesting if difference between μ and μ' is greater than 5 years (i.e., Δ = 5)

✗For r, suppose          n1 = 50, s1 = 3.5

✗For r' (complement): n2 = 250, s2 = 6.5

$$Z = \frac{\mu' - \mu - \Delta}{\sqrt{\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}}} = \frac{30 - 23 - 5}{\sqrt{\dfrac{3.5^2}{50} + \dfrac{6.5^2}{250}}} = 3.11$$

✗For 1-sided test at 95% confidence level, critical Z-value for rejecting null hypothesis is 1.64.

✗Since Z is greater than 1.64, r is an interesting rule

# Min-Apriori (Han et al)

Document-term matrix:

| TID | W1 | W2 | W3 | W4 | W5 |
|-----|----|----|----|----|----|
| D1  | 2  | 2  | 0  | 0  | 1  |
| D2  | 0  | 0  | 1  | 2  | 2  |
| D3  | 2  | 3  | 0  | 0  | 0  |
| D4  | 0  | 0  | 1  | 0  | 1  |
| D5  | 1  | 1  | 1  | 0  | 2  |

Example:

W1 and W2 tends to appear together in the same document

# Min-Apriori

✗ DATA CONTAINS ONLY CONTINUOUS ATTRIBUTES OF THE SAME "TYPE"

✗ e.g., frequency of words in a document

| TID | W1 | W2 | W3 | W4 | W5 |
|-----|----|----|----|----|----|
| D1  | 2  | 2  | 0  | 0  | 1  |
| D2  | 0  | 0  | 1  | 2  | 2  |
| D3  | 2  | 3  | 0  | 0  | 0  |
| D4  | 0  | 0  | 1  | 0  | 1  |
| D5  | 1  | 1  | 1  | 0  | 2  |

✗ POTENTIAL SOLUTION:

✗ Convert into 0/1 matrix and then apply existing algorithms

✗ lose word frequency information

✗ Discretization does not apply as users want association among words not ranges of words

# Min-Apriori

✗ HOW TO DETERMINE THE SUPPORT OF A WORD?

  ✗ If we simply sum up its frequency, support count will be greater than total number of documents!

    ✗ Normalize the word vectors – e.g., using $L_1$ norm

    ✗ Each word has a support equals to 1.0

| TID | W1 | W2 | W3 | W4 | W5 |
|-----|----|----|----|----|----|
| D1 | 2 | 2 | 0 | 0 | 1 |
| D2 | 0 | 0 | 1 | 2 | 2 |
| D3 | 2 | 3 | 0 | 0 | 0 |
| D4 | 0 | 0 | 1 | 0 | 1 |
| D5 | 1 | 1 | 1 | 0 | 2 |

**Normalize** →

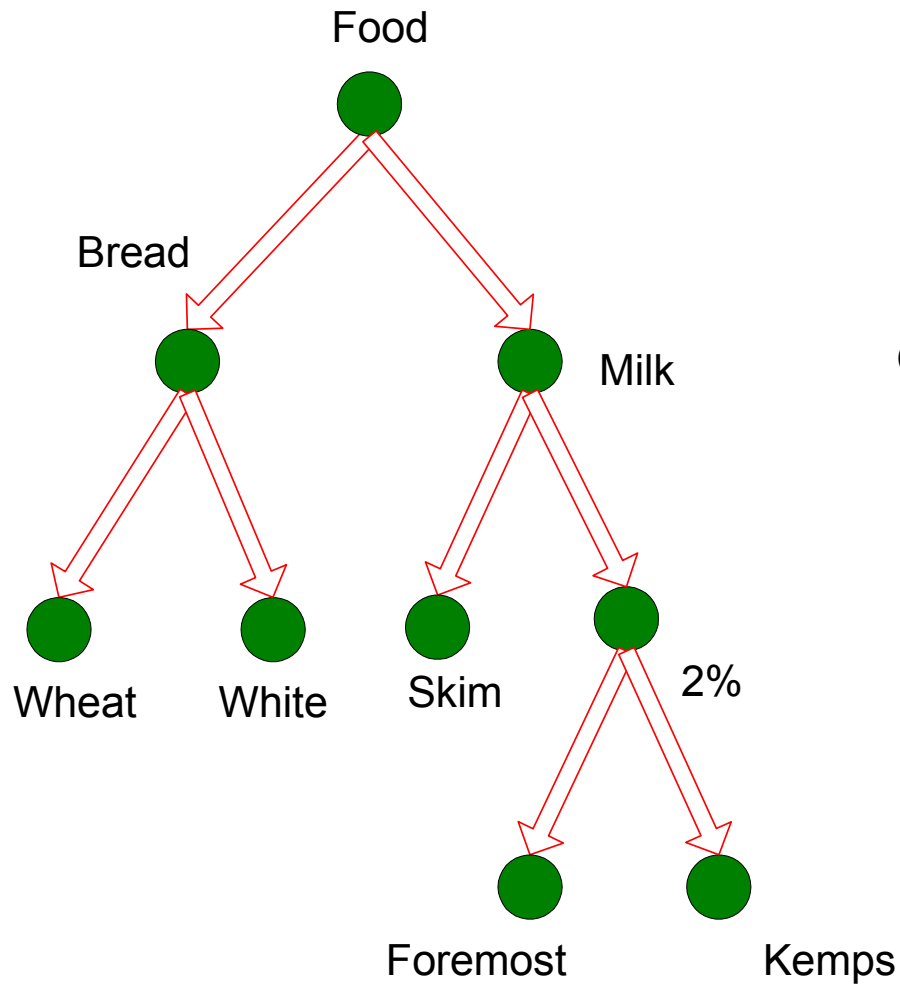| TID | W1 | W2 | W3 | W4 | W5 |
|-----|------|------|------|------|------|
| D1 | 0.40 | 0.33 | 0.00 | 0.00 | 0.17 |
| D2 | 0.00 | 0.00 | 0.33 | 1.00 | 0.33 |
| D3 | 0.40 | 0.50 | 0.00 | 0.00 | 0.00 |
| D4 | 0.00 | 0.00 | 0.33 | 0.00 | 0.17 |
| D5 | 0.20 | 0.17 | 0.33 | 0.00 | 0.33 |

✗ NEW DEFINITION OF SUPPORT:

$$\text{sup}(C) = \sum_{i \in T} \min_{j \in C} D(i, j)$$

| TID | W1 | W2 | W3 | W4 | W5 |
|-----|------|------|------|------|------|
| D1 | 0.40 | 0.33 | 0.00 | 0.00 | 0.17 |
| D2 | 0.00 | 0.00 | 0.33 | 1.00 | 0.33 |
| D3 | 0.40 | 0.50 | 0.00 | 0.00 | 0.00 |
| D4 | 0.00 | 0.00 | 0.33 | 0.00 | 0.17 |
| D5 | 0.20 | 0.17 | 0.33 | 0.00 | 0.33 |

**Example:**

**Sup(W1,W2,W3)**

**= 0 + 0 + 0 + 0 + 0.17**

**= 0.17**

# Anti-monotone property of Support

| TID | W1 | W2 | W3 | W4 | W5 |
|-----|------|------|------|------|------|
| D1 | 0.40 | 0.33 | 0.00 | 0.00 | 0.17 |
| D2 | 0.00 | 0.00 | 0.33 | 1.00 | 0.33 |
| D3 | 0.40 | 0.50 | 0.00 | 0.00 | 0.00 |
| D4 | 0.00 | 0.00 | 0.33 | 0.00 | 0.17 |
| D5 | 0.20 | 0.17 | 0.33 | 0.00 | 0.33 |

**Example:**

**Sup(W1) = 0.4 + 0 + 0.4 + 0 + 0.2 = 1**

**Sup(W1, W2) = 0.33 + 0 + 0.4 + 0 + 0.17 = 0.9**

**Sup(W1, W2, W3) = 0 + 0 + 0 + 0 + 0.17 = 0.17**

# Multi-level Association Rules

✗WHY SHOULD WE INCORPORATE CONCEPT HIERARCHY?

  ✗Rules at lower levels may not have enough support to appear in any frequent itemsets

  ✗Rules at lower levels of the hierarchy are overly specific

    ✗ e.g.,  skim milk → white bread, 2% milk → wheat bread, skim milk → wheat bread, etc.
  are indicative of association between milk and bread

✗How do support and confidence vary as we traverse the concept hierarchy?

✗If X is the parent item for both X1 and X2, then
$\sigma(X) \leq \sigma(X1) + \sigma(X2)$

✗If $\sigma(X1 \cup Y1) \geq$ minsup,
and X is parent of X1, Y is parent of Y1
then $\sigma(X \cup Y1) \geq$ minsup, $\sigma(X1 \cup Y) \geq$ minsup
$\sigma(X \cup Y) \geq$ minsup

✗If $conf(X1 \Rightarrow Y1) \geq$ minconf,
then $conf(X1 \Rightarrow Y) \geq$ minconf

# Multi-level Association Rules

✗ APPROACH 1:

  ✗ Extend current association rule formulation by augmenting each transaction with higher level items

  Original Transaction: {skim milk, wheat bread}

  Augmented Transaction:
  {skim milk, wheat bread, milk, bread, food}

✗ ISSUES:

  ✗ Items that reside at higher levels have much higher support counts
    ✗ if support threshold is low, too many frequent patterns involving items from the higher levels
  ✗ Increased dimensionality of the data

# Multi-level Association Rules

✗ APPROACH 2:

    ✗ Generate frequent patterns at highest level first

    ✗ Then, generate frequent patterns at the next highest level, and so on

✗ ISSUES:

    ✗ I/O requirements will increase dramatically because we need to perform more passes over the data

    ✗ May miss some potentially interesting cross-level association patterns

# 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 |

# Examples of 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 = < e1 e2 e3 ... >
  - Each element contains a collection of events (items)
    - ei = {i1, i2, ..., ik}
  - 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 INITIATING EVENTS CAUSING THE NUCLEAR ACCIDENT AT 3-MILE
ISLAND:
(HTTP://STELLAR-ONE.COM/NUCLEAR/STAFF_REPORTS/SUMMARY_SOE_THE_INITIATING_EVENT.HTM)

   <   {clogged resin} {outlet valve closure} {loss of feedwater}
       {condenser polisher outlet valve shut} {booster pumps trip}
       {main waterpump trips} {main turbine trips} {reactor pressure increases}>

✘ 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 $<A_1 A_2 \ldots A_N>$ IS CONTAINED IN ANOTHER SEQUENCE $<B_1 B_2 \ldots B_M>$ $(M \geq N)$ IF THERE EXIST INTEGERS

$I_1 < I_2 < \ldots < I_N$ SUCH THAT $A_1 \subseteq B_{I1}$, $A_2 \subseteq B_{I1}$, ..., $A_N \subseteq B_{IN}$

| Data sequence | Subsequence | Contain? |
|---|---|---|
| < {2,4} {3,5,6} {8} > | < {2} {3,5} > | Yes |
| < {1,2} {3,4} > | < {1} {2} > | No |
| < {2,4} {2,4} {2,5} > | < {2} {4} > | 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 $\geq$ *MINSUP*)

# Sequential Pattern Mining: Definition

✗Given:

  ✗a database of sequences

  ✗a user-specified minimum support threshold, *minsup*

✗Task:

  ✗Find all subsequences with support ≥ *minsup*

✗GIVEN A SEQUENCE:  <{A B} {C D E} {F} {G H I}>

  ✗Examples of subsequences:

  <{a} {c d} {f} {g} >, < {c d e} >, < {b} {g} >, etc.

✗HOW MANY K-SUBSEQUENCES CAN BE EXTRACTED FROM A GIVEN N-SEQUENCE?

  <{A  B} {C D  E} {F} {G H  I}>  N = 9

K=4:    Y _  _ Y Y _ _ _Y

  <{A}      {D E}    {I}>

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:**

| | |
|---|---|
| < {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:   I1, I2, I3, …, IN

✗CANDIDATE 1-SUBSEQUENCES:

<{i1}>, <{i2}>, <{i3}>, …, <{in}>

✗CANDIDATE 2-SUBSEQUENCES:

<{i1, i2}>, <{i1, i3}>, …, <{i1} {i1}>, <{i1} {i2}>, …, <{in-1} {in}>

✗CANDIDATE 3-SUBSEQUENCES:

<{i1, i2 , i3}>, <{i1, i2 , i4}>, …, <{i1, i2} {i1}>, <{i1, i2} {i2}>, …,

<{i1} {i1 , i2}>, <{i1} {i1 , i3}>, …, <{i1} {i1} {i1}>, <{i1} {i1} {i2}>, …

# 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*

- **BASE CASE (K=2):**
  - Merging two frequent 1-sequences $<\{i_1\}>$ and $<\{i_2\}>$ will produce two candidate 2-sequences: $<\{i_1\} \{i_2\}>$ and $<\{i_1 i_2\}>$

- **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

$W_1$ =<{1} {2 3} {4}> AND $W_2$ =<{2 3} {4 5}>

WILL PRODUCE THE CANDIDATE SEQUENCE < {1} {2 3} {4 5}> BECAUSE THE LAST TWO EVENTS IN $W_2$ (4 AND 5) BELONG TO THE SAME ELEMENT

✗ MERGING THE SEQUENCES

$W_1$ =<{1} {2 3} {4}> AND $W_2$ =<{2 3} {4} {5}>

WILL PRODUCE THE CANDIDATE SEQUENCE < {1} {2 3} {4} {5}> 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

$W_1$ =<{1} {2 6} {4}> AND $W_2$ =<{1} {2} {4 5}>

TO PRODUCE THE CANDIDATE < {1} {2 6} {4 5}> BECAUSE IF THE LATTER IS A VIABLE CANDIDATE, THEN IT CAN BE OBTAINED BY MERGING $W_1$ WITH

< {1} {2 6} {5}>

**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)

{A  B}    {C}    {D  E}

<= $x_g$    >$n_g$

<= $m_s$

$x_g$: max-gap

$n_g$: min-gap

$m_s$: maximum span

$x_g$ = 2, $n_g$ = 0, $m_s$ = 4

| Data sequence | Subsequence | Contain? |
|---|---|---|
| < {2,4} {3,5,6} {4,7} {4,5} {8} > | < {6} {5} > | Yes |
| < {1} {2} {3} {4} {5}> | < {1} {4} > | No |
| < {1} {2,3} {3,4} {4,5}> | < {2} {3} {5} > | Yes |
| < {1,2} {3} {2,3} {3,4} {2,4} {4,5}> | < {1,2} {5} > | No |

✗ 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:

$x_g = 1$ (max-gap)

$n_g = 0$ (min-gap)

$m_s = 5$ (maximum span)

*minsup* = 60%

<{2} {5}>   support = 40%

but

<{2} {3} {5}>   support = 60%

**Problem exists because of max-gap constraint**

**No such problem if max-gap is infinite**

# Contiguous Subsequences

- S IS A CONTIGUOUS SUBSEQUENCE OF

$$W = <E_1><E_2>...<E_K>$$

IF ANY OF THE FOLLOWING CONDITIONS HOLD:

1. s is obtained from w by deleting an item from either $e_1$ or $e_k$

2. s is obtained from w by deleting an item from any element $e_i$ that contains more than 2 items

3. s is a contiguous subsequence of s' and s' is a contiguous subsequence of w (recursive definition)

- EXAMPLES: S = < {1} {2} >

  – is a contiguous subsequence of
    < {1} {2 3}>, < {1 2} {2} {3}>, and < {3 4} {1 2} {2 3} {4} >

  – is not a contiguous subsequence of
    < {1} {3} {2}> and < {2} {1} {3} {2}>

# Modified Candidate Pruning Step

✗ WITHOUT MAXGAP CONSTRAINT:

  ✗ A candidate k-sequence is pruned if at least one of its (k–1)-subsequences is infrequent

✗ WITH MAXGAP CONSTRAINT:

  ✗ A candidate $k$-sequence is pruned if at least one of its **contiguous** $(k-1)$-subsequences is infrequent

# Timing Constraints (II)

{A  B}    {C}    {D  E}

<= $x_g$     >$n_g$     <= ws

<= $m_s$

$x_g$: max-gap

$n_g$: min-gap

**ws: window size**

$m_s$: maximum span

$x_g = 2$, $n_g = 0$, **ws = 1**, $m_s = 5$

| Data sequence | Subsequence | Contain? |
|---|---|---|
| < {2,4} {3,5,6} {4,7} {4,6} {8} > | < {3} {5} > | No |
| < {1} {2} {3} {4} {5}> | < {1,2} {3} > | Yes |
| < {1,2} {2,3} {3,4} {4,5}> | < {1,2} {3,4} > | Yes |

✗ GIVEN A CANDIDATE PATTERN: <{A, C}>

  ✗ Any data sequences that contain

    <... {a c} ... >,
    <... {a} ... {c}...>   ( where time({c}) − time({a}) ≤ ws)
    <...{c} ... {a} ...>   (where time({a}) − time({c}) ≤ ws)

    will contribute to the support count of candidate pattern

# Other Formulation

✗IN SOME DOMAINS, WE MAY HAVE ONLY ONE VERY LONG TIME SERIES
  ✗Example:
    ✗ monitoring network traffic events for attacks
    ✗ monitoring telecommunication alarm signals
✗GOAL IS TO FIND FREQUENT SEQUENCES OF EVENTS IN THE TIME SERIES
  ✗This problem is also known as frequent episode mining



Pattern: <E1> <E3>

# General Support Counting Schemes

Object's Timeline

Sequence: (p) (q)

| Method | Support Count |
|--------|---------------|
| COBJ | 1 |
| CWIN | 6 |
| CMINWIN | 4 |
| CDIST_O | 8 |
| CDIST | 5 |

Assume:

$x_g = 2$ (max-gap)

$n_g = 0$ (min-gap)

$ws = 0$ (window size)

$m_s = 2$ (maximum span)

# Frequent Subgraph Mining

✗Extend association rule mining to finding frequent subgraphs
✗Useful for Web Mining, computational chemistry, bioinformatics, spatial data sets, etc

# Graph Definitions



(a) Labeled Graph    (b) Subgraph    (c) Induced Subgraph

# Representing Transactions as Graphs

✗Each transaction is a clique of items

| Transaction Id | Items |
|---|---|
| 1 | {A,B,C,D} |
| 2 | {A,B,E} |
| 3 | {B,C} |
| 4 | {A,B,D,E} |
| 5 | {B,C,D} |

TID = 1:

# Representing Graphs as Transactions



G1

G2

G3

| | (a,b,p) | (a,b,q) | (a,b,r) | (b,c,p) | (b,c,q) | (b,c,r) | … | (d,e,r) |
|---|---|---|---|---|---|---|---|---|
| G1 | 1 | 0 | 0 | 0 | 0 | 1 | … | 0 |
| G2 | 1 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| G3 | 0 | 0 | 1 | 1 | 0 | 0 | … | 0 |
| G3 | … | … | … | … | … | … | … | … |

# Challenges

✗ NODE MAY CONTAIN DUPLICATE LABELS

✗ SUPPORT AND CONFIDENCE
- ✗ How to define them?

✗ ADDITIONAL CONSTRAINTS IMPOSED BY PATTERN STRUCTURE
- ✗ Support and confidence are not the only constraints
- ✗ Assumption: frequent subgraphs must be connected

✗ APRIORI-LIKE APPROACH:
- ✗ Use frequent k-subgraphs to generate frequent (k+1) subgraphs
  - ✗ What is k?

✗ SUPPORT:
- ✗ number of graphs that contain a particular subgraph

✗ APRIORI PRINCIPLE STILL HOLDS

✗ LEVEL-WISE (APRIORI-LIKE) APPROACH:
- ✗ Vertex growing:
    - ✗ k is the number of vertices
- ✗ Edge growing:
    - ✗ k is the number of edges

# Vertex Growing

G1       +       G2       G3 = join(G1,G2)

$$M_{G1} = \begin{pmatrix} 0 & p & p & q \\ p & 0 & r & 0 \\ p & r & 0 & 0 \\ q & 0 & 0 & 0 \end{pmatrix}$$

$$M_{G2} = \begin{pmatrix} 0 & p & p & 0 \\ p & 0 & r & 0 \\ p & r & 0 & r \\ 0 & 0 & r & 0 \end{pmatrix}$$

$$M_{G3} = \begin{pmatrix} 0 & p & p & 0 & q \\ p & 0 & r & 0 & 0 \\ p & r & 0 & r & 0 \\ 0 & 0 & r & 0 & 0 \\ q & 0 & 0 & 0 & 0 \end{pmatrix}$$
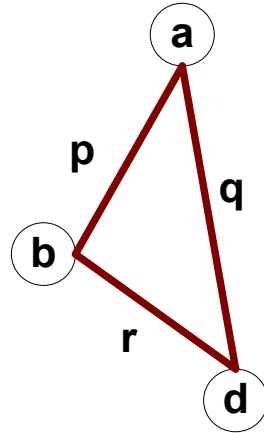
# Edge Growing



G1

+

G2

G3 = join(G1,G2)

# Apriori-like Algorithm

✗ FIND FREQUENT 1-SUBGRAPHS

✗ REPEAT

  ✗ Candidate generation

    ✗ Use frequent ($k$-$1$)-subgraphs to generate candidate $k$-subgraph

  ✗ Candidate pruning

    ✗ Prune candidate subgraphs that contain infrequent ($k$-$1$)-subgraphs

  ✗ Support counting

    ✗ Count the support of each remaining candidate

  ✗ Eliminate candidate $k$-subgraphs that are infrequent

**In practice, it is not as easy. There are many other issues**

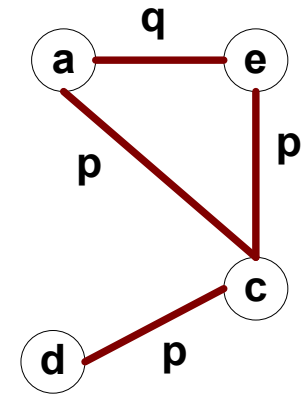G1                    G2                    G3                    G4

| | (a,b,p) | (a,b,q) | (a,b,r) | (b,c,p) | (b,c,q) | (b,c,r) | … | (d,e,r) |
|---|---|---|---|---|---|---|---|---|
| G1 | 1 | 0 | 0 | 0 | 0 | 1 | … | 0 |
| G2 | 1 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| G3 | 0 | 0 | 1 | 1 | 0 | 0 | … | 0 |
| G4 | 0 | 0 | 0 | 0 | 0 | 0 | … | 0 |

# Example

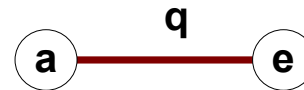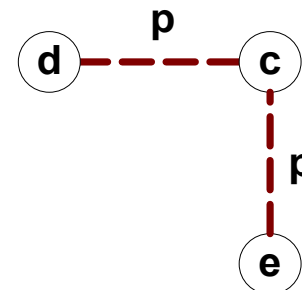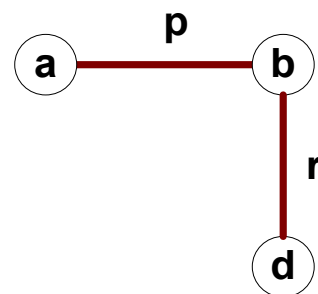Minimum support count = 2

**k=1**
**Frequent**
**Subgraphs**

(a)     (b)     (c)     (d)     (e)

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**k=2**
**Frequent**
**Subgraphs**

(a) —p— (b)     (a) —q— (e)     (b) —r— (d)

(c) —p— (d)     (c) —p— (e)

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**k=3**
**Candidate**
**Subgraphs**

(a) —p— (b)
         |
         r
         |
        (d)

(d) --p-- (c)
           |
           p
           |
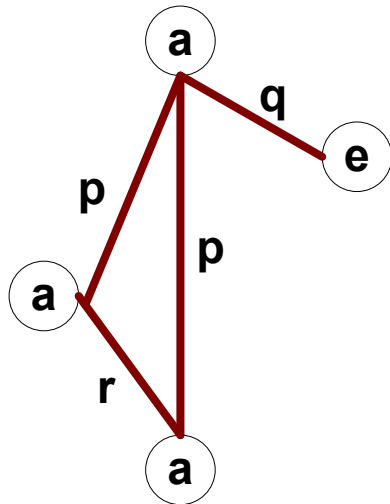          (e)

(Pruned candidate)

59

# Candidate Generation

✗ IN APRIORI:

  ✗ Merging two frequent *k*-itemsets will produce a candidate (*k+1*)-itemset

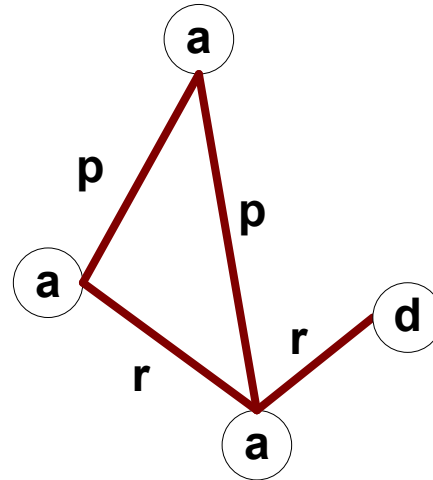✗ IN FREQUENT SUBGRAPH MINING (VERTEX/EDGE GROWING)

  ✗ Merging two frequent *k*-subgraphs may produce more than one candidate (*k+1*)-subgraph
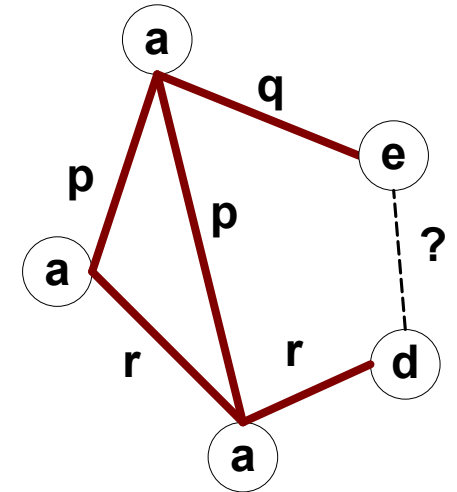
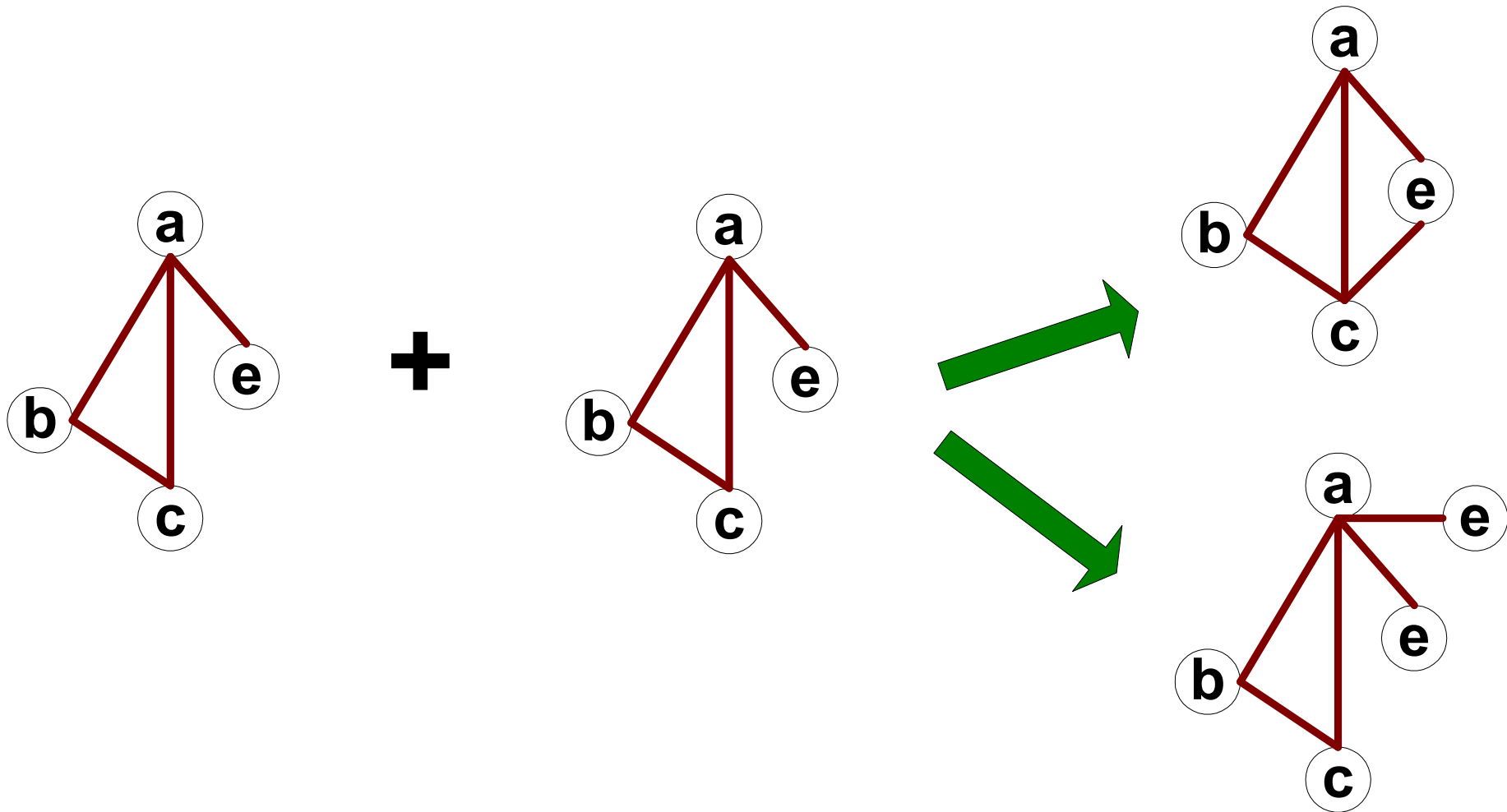# Multiplicity of Candidates (Vertex Growing)



G1

+

G2

G3 = join(G1,G2)

$$M_{G1} = \begin{pmatrix} 0 & p & p & q \\ p & 0 & r & 0 \\ p & r & 0 & 0 \\ q & 0 & 0 & 0 \end{pmatrix}$$

$$M_{G2} = \begin{pmatrix} 0 & p & p & 0 \\ p & 0 & r & 0 \\ p & r & 0 & r \\ 0 & 0 & r & 0 \end{pmatrix}$$

$$M_{G3} = \begin{pmatrix} 0 & p & p & 0 & q \\ p & 0 & r & 0 & 0 \\ p & r & 0 & r & 0 \\ 0 & 0 & r & 0 & ? \\ q & 0 & 0 & ? & 0 \end{pmatrix}$$
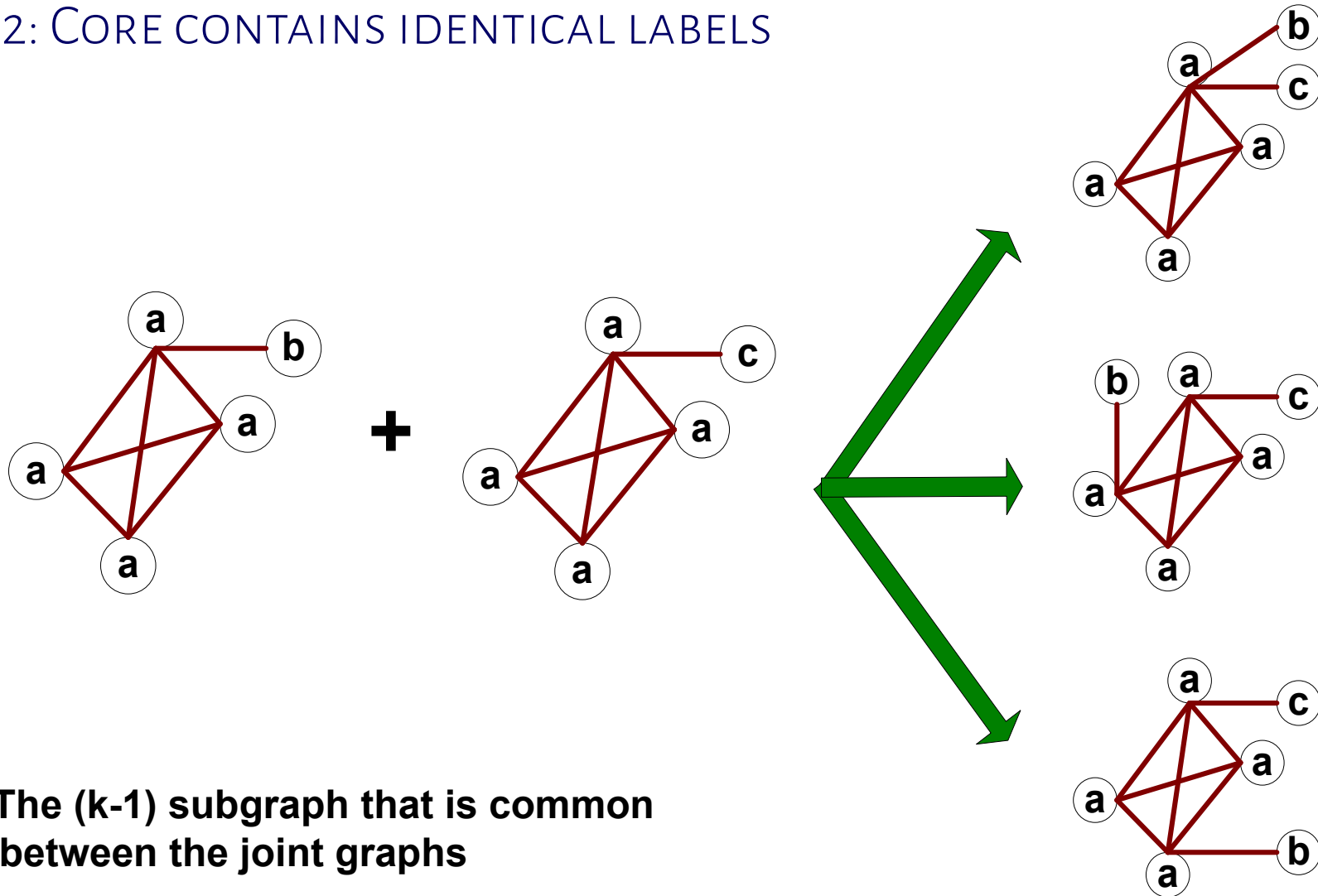
✗CASE 1: IDENTICAL VERTEX LABELS

✗ CASE 2: CORE CONTAINS IDENTICAL LABELS



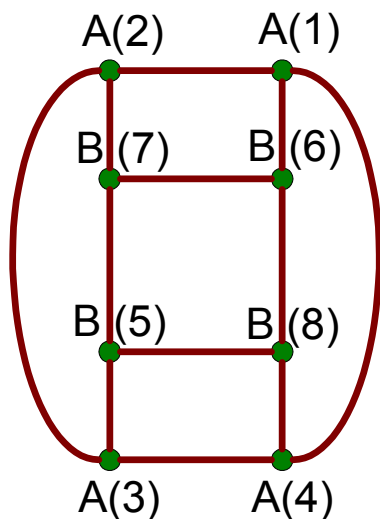**Core: The (k-1) subgraph that is common between the joint graphs**

✗ CASE 3: CORE MULTIPLICITY

# Adjacency Matrix Representation



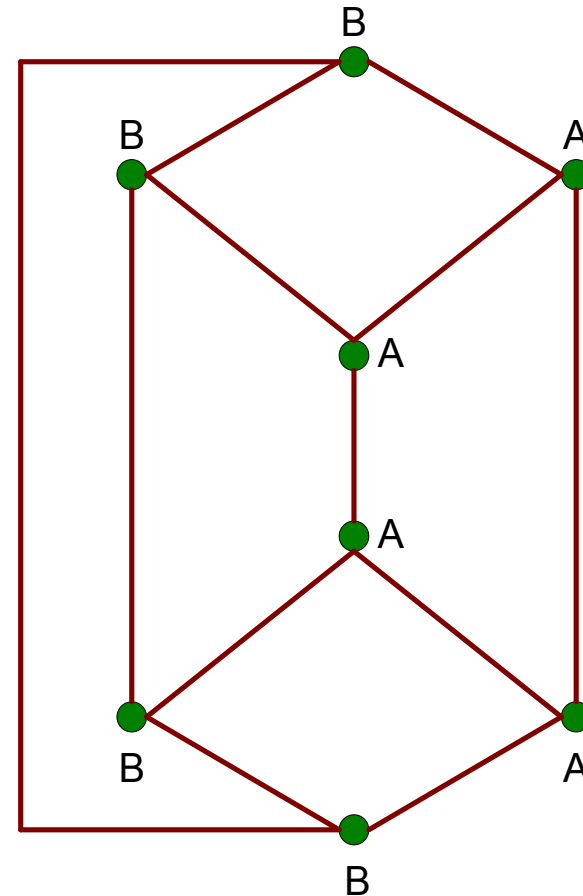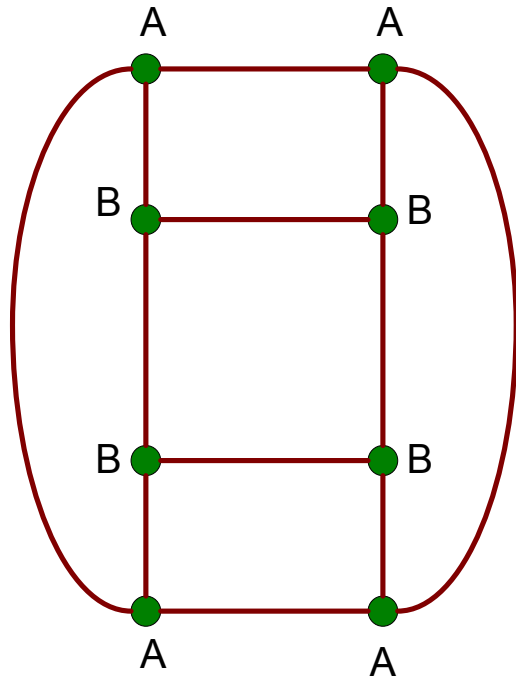|  | A(1) | A(2) | A(3) | A(4) | B(5) | B(6) | B(7) | B(8) |
|---|---|---|---|---|---|---|---|---|
| **A(1)** | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| **A(2)** | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| **A(3)** | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |
| **A(4)** | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 |
| **B(5)** | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 |
| **B(6)** | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 |
| **B(7)** | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 |
| **B(8)** | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 |

|  | A(1) | A(2) | A(3) | A(4) | B(5) | B(6) | B(7) | B(8) |
|---|---|---|---|---|---|---|---|---|
| **A(1)** | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| **A(2)** | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 |
| **A(3)** | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| **A(4)** | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 |
| **B(5)** | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 |
| **B(6)** | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| **B(7)** | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 |
| **B(8)** | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 |

• **The same graph can be represented in many ways**

# Graph Isomorphism

✗A GRAPH IS ISOMORPHIC IF IT IS TOPOLOGICALLY EQUIVALENT TO ANOTHER GRAPH

# Graph Isomorphism

✗ TEST FOR GRAPH ISOMORPHISM IS NEEDED:

   ✗ During candidate generation step, to determine whether a candidate has been generated

   ✗ During candidate pruning step, to check whether its (*k-1*)-subgraphs are frequent

   ✗ During candidate counting, to check whether a candidate is contained within another graph

✗ USE CANONICAL LABELING TO HANDLE ISOMORPHISM

  ✗ Map each graph into an ordered string representation (known as its code) such that two isomorphic graphs will be mapped to the same canonical encoding

  ✗ Example:

    ✗ Lexicographically largest adjacency matrix

$$\begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix} \qquad \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

**String: 0010001111010110**          **Canonical: 0111101011001000**

68