# ASSOCIATION RULE MINING

BEIYU LIN

- Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items

$$X \longrightarrow Y$$

$$X = \{Beer\} \quad Y = \{eggs\}$$

Market transactions

| TID | Items |
|-----|-------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

Example of Association Rules

{Beer} → {Eggs},
{Milk, Bread} → {Diaper, Beer},

$$\delta(x) = 3 \geq 2$$

$$S(x) = \frac{\delta(x)}{5} = \frac{3}{5}$$

$$C(x \to Y) = \frac{\delta(x \cup Y)}{\delta(x)} = \frac{1}{3}$$

$$\delta(Y) = 1 \not\geq 2$$

$$S(Y) = \frac{1}{5}$$

$$\min \delta = 2 \quad \min S = \min \frac{\delta}{\boxed{\phantom{xx}}}$$

$$\delta(x \cup Y) = 1$$

sup count $\delta$ ; sup $S$ ; cont $C$ ; freq $\geq \min \delta$ or $\min S$

# ASSOCIATION RULE MINING

- **Itemset (set / subset)**
  - A collection of one or more items
    - Example: {Milk, Bread, Diaper}
  - k-itemset
    - An itemset that contains k items

- **Support count ($\sigma$)**
  - Frequency of occurrence of an itemset
  - E.g. $\sigma$({Milk, Bread,Diaper}) = 2

| TID | Items |
|-----|-------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

**Support**
Fraction of transactions that contain an itemset
E.g.   s({Milk, Bread, Diaper}) = 2/5

**Frequent Itemset**
- An itemset whose support is greater than or equal to a *minsup* threshold

# DEFINITION: ASSOCIATION RULE

- Association Rule
  - An implication expression of the form X → Y, where X and Y are itemsets
  - Example:
    {Milk, Diaper} → {Beer}

| TID | Items |
|-----|-------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

- Rule Evaluation Metrics
  - Support (s)
    - Fraction of transactions that contain both X and Y
  - Confidence (c)
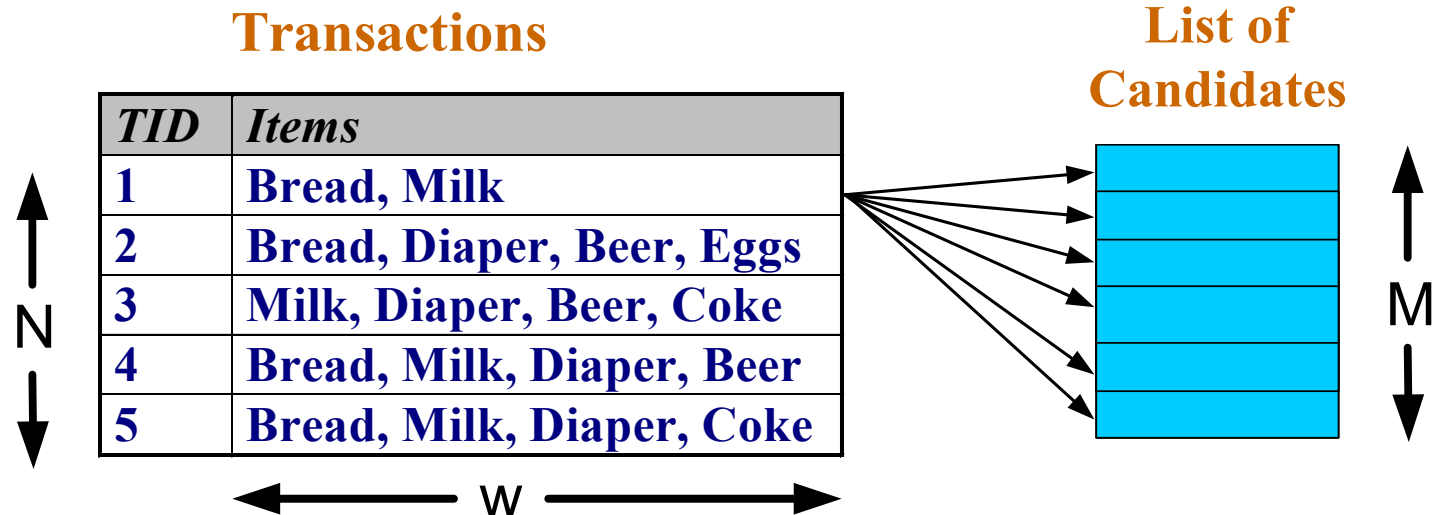    - Measures how often items in Y appear in transactions that contain X

Example:

$$\{Milk, Diaper\} \Rightarrow \{Beer\}$$

$$s = \frac{\sigma(Milk, Diaper, Beer)}{|T|} = \frac{2}{5} = 0.4 \quad =\# \text{ of itemset / total \# transac}$$

$$c = \frac{\sigma(Milk, Diaper, Beer)}{\sigma(Milk, Diaper)} = \frac{2}{3} = 0.67 \quad = \# \text{ of itemset of X and Y / \# of}$$
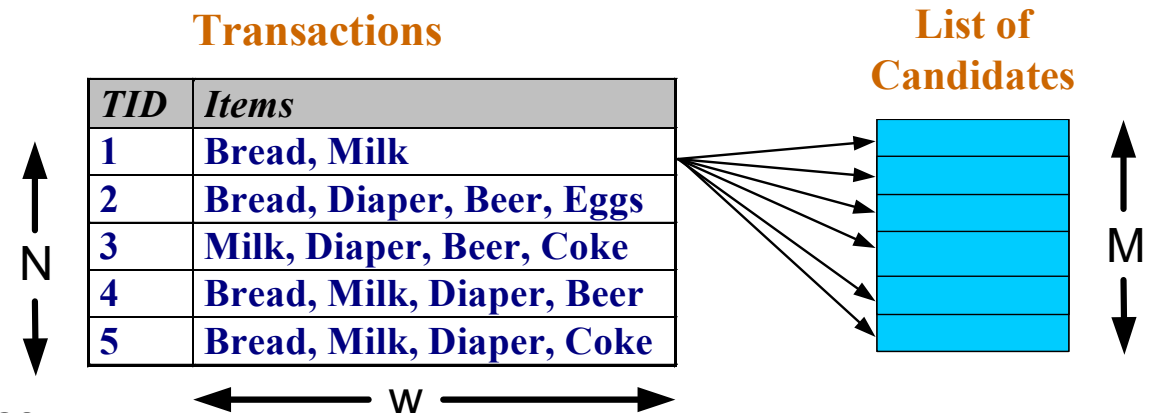
# FREQUENT ITEMSET GENERATION

- Brute-force approach:

    - Each itemset in the lattice is a candidate frequent itemset

    - Count the support of each candidate by scanning the database



**Transactions**

| TID | Items |
|-----|-------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

**List of Candidates**

    - Match each transaction against every candidate

    - Complexity ~ O(NMw) => Expensive since M = $2^d$ !!!

# FREQUENT ITEMSET GENERATION STRATEGIES

- Reduce the number of candidates (M)
  - Complete search: $M=2^d$    $d=6$    $2^6$
  - Use pruning techniques to reduce M

- Reduce the number of transactions (N)
  - Reduce size of N as the size of itemset increases
  - Used by DHP and vertical-based mining algorithms

$$\bar{I} = \{A, B, C\} \quad 2^3$$

$$E = \{B, M, D, B, \bar{G}, C\}$$

- Reduce the number of comparisons (NM)
  - Use efficient data structures to store the candidates or transactions
  - No need to match every candidate against every transaction

**Transactions**

| TID | Items |
|-----|-------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

N

W

**List of Candidates**

M

# REDUCING NUMBER OF CANDIDATES

| TID | Items |
|-----|-------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

- Apriori principle:
  - If an itemset is frequent, then all of its subsets must also be frequent

- Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

for any $x, Y$:     $\sigma(X) \geq \sigma(Y)$

  - Support of an itemset never exceeds the support of its subsets
  - This is known as the anti-monotone property of support

**Support**

s({Milk, Bread, Diaper}) = 2/5

$X_1$

$f:$   $\sigma(x_1) = 2 \geq min\sigma = 2$

$X_1 \rightarrow Y$

$X = \{Beers\}$   $f$   $\sigma(x) = 3$

$X_{11} = \{milk\} \subseteq X_1$   $\sigma(x_{11}) = 4 \geq min\sigma$

$X_{11} f \Leftarrow f$   $\sigma(x_{11}) \geq \sigma(x_1)$

# ILLUSTRATING APRIORI PRINCIPLE



$E = \{A, B, C, D, E\}$

Brute-force

$min\ \delta = 2$

$\delta(A) = 2$

$\delta(AB) = 1$ ← Found to be Infrequent

NF

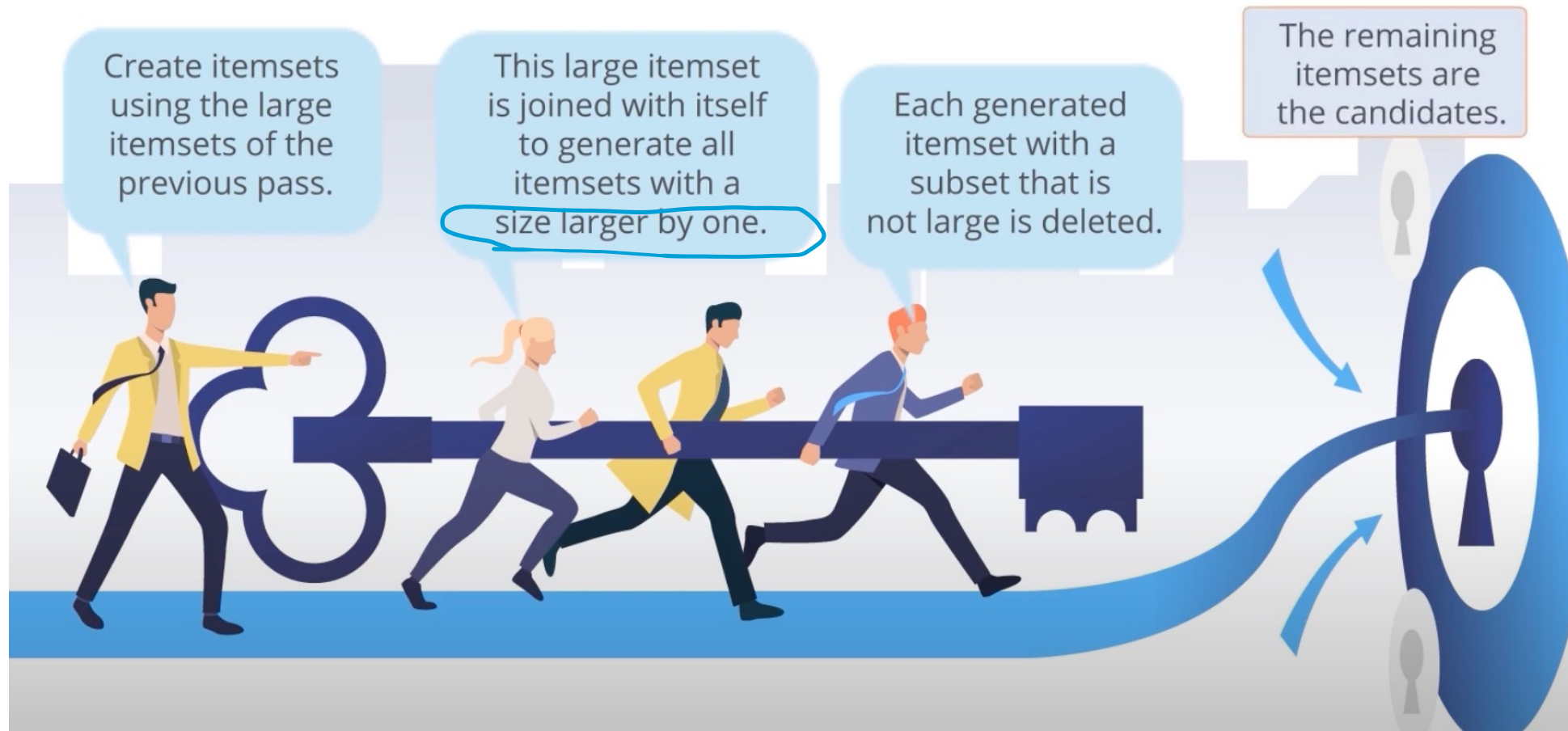$\delta(ABC) \leq \delta(AB) = 1$

Pruned supersets

Subset:
$\{\}, \{A\} \cdots \{E\}$
$--- \{A \cdots E\}$

Superset:
$\{AB\} \rightarrow ?$

$\{ABX\}$    $\{ABXXX\}$
$\{ABXX\}$    $CDE$
$CD, CE, DE$

# APRIORI PRINCIPLE

# APRIORI PRINCIPLE

Uses frequent itemsets to generate association rules

Support value of frequent itemsets is greater than the threshold value

The algorithm reduces the number of candidates being considered by only exploring the itemsets whose support count is greater than the minimum support count.

# APRIORI PRINCIPLE EXAMPLE

| TID | Items |
|-----|-------|
| 100 | 1 3 4 |
| 200 | 2 3 5 |
| 300 | 1 2 3 5 |
| 400 | 2 5 |
| 500 | 1 3 5 |

itemwe1

**CI1**

| Itemset | Support 6 |
|---------|-----------|
| {1} | 3 |
| {2} | 3 |
| {3} | 4 |
| {4} | ~~1~~ NF |
| {5} | 4 |

**F11**

| Itemset | Support |
|---------|---------|
| {1} | 3 |
| {2} | 3 |
| {3} | 4 |
| {5} | 4 |

$E = \{1, 2, 3, 4, 5\}$    $\min 6 = 2$

# APRIORI PRINCIPLE EXAMPLE

The length of the itemset is extended with 1 (k = k+1).

item = 2

**F11**

| Itemset | Support |
|---------|---------|
| {1} | 3 |
| {2} | 3 |
| {3} | 4 |
| {5} | 4 |

**CI2**

NF

| Itemset | Support |
|---------|---------|
| {1,2} | 1 |
| {1,3} | 3 |
| {1,5} | 2 |
| {2,3} | 2 |
| {2,5} | 3 |
| {3,5} | 3 |

**FI2**

| Itemset | Support |
|---------|---------|
| {1,3} | 3 |
| {1,5} | 2 |
| {2,3} | 2 |
| {2,5} | 3 |
| {3,5} | 3 |

# APRIORI PRINCIPLE EXAMPLE

The length of the itemset is extended with 1 (k = k+1).

**FI2**

| Itemset | Support |
|---------|---------|
| {1,3} | 3 |
| {1,5} | 2 |
| {2,3} | 2 |
| {2,5} | 3 |
| {3,5} | 3 |

**CI3**

| Itemset | Support |
|---------|---------|
| {1,2,3} | / |
| {1,2,5} | 1 |
| {1,3,5} | 2 |
| {2,3,5} | 2 |

| Itemset | Support |
|---------|---------|
| {1,3,5} | 2 |
| {2,3,5} | 2 |

min 5x

F { 1,3,5 }  {2,3,5}

List the F itemsets with k = 2

{1, 3}, {1, 5}, {3,5}, {2,3}, {2, 5} F

K=1

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

# APRIORI PRINCIPLE EXAMPLE

| TID | Items |
|-----|-------|
| 100 | 1 3 4 |
| 200 | 2 3 5 |
| 300 | 1 2 3 5 |
| 400 | 2 5 |
| 500 | 1 3 5 |

| Itemset | Support |
|---------|---------|
| {1,3,5} | 2 |
| {2,3,5} | 2 |

# APRIORI PRINCIPLE EXAMPLE

The length of the itemset is extended with 1 (k = k+1).

C4

| Itemset | Support |
|---------|---------|
| {1,3,5} | 2 |
| {2,3,5} | 2 |

| Itemset | Support |
|---------|---------|
| {1,2,3,5} | 1 |

MF

# ILLUSTRATING APRIORI PRINCIPLE

| TID | ITEMS |
| --- | --- |
| 1 | A, B |
| 2 | A, C, D, E |
| 3 | B, C, D |
| 4 | B, C, E |
| 5 | C, B, D |

# WEKA

- Dataset: https://storm.cis.fordham.edu/~gweiss/data-mining/weka-data/contact-lenses.arff