

LAST WEEK

- Euclidean distance
- Minkowski distance
- 1 ■ Simple Matching and Jaccard Coefficient
- Cosine similarity
- Correlation measures

	H	W	...	T
A				
S				

Handwritten notes on the table: A bracket under the first three columns is labeled 'f'. A bracket under the last column is labeled 'L'. To the right of the table, there is a brace grouping the last column and the label '1 PP'.

$$\frac{f_{10} + f_{11}}{\# f}$$

$$\frac{f_{11}}{f_{01} + f_{10} + f_{11}}$$

Diagram showing points $x_1 = (1, 2)$ and $x_2 = (4, 3)$ with a horizontal distance of 3 and a vertical distance of 1. The Euclidean distance is calculated as:

$$E = \sqrt{(4-1)^2 + (3-2)^2} = \sqrt{10}$$

$$= \left((4-1)^2 + (3-2)^2 \right)^{\frac{1}{2}}$$

$$M = \left(\right)^{\frac{1}{r}}$$

① $r=1$ $L_1 = M_{r=1} = 3 + 1 = 4$

② $r=2$ $L_2 = E$

③ $r=\infty$ $L_{\max} = L_{\infty} = \max(3, 1) = 3$



ASSOCIATION RULE MINING

BEIYU LIN



ASSOCIATION RULE MINING

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

Market transactions

<i>TID</i>	<i>Items</i>
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

$\{\text{Beer}\} \rightarrow \{\text{Eggs}\},$
 $\{\text{Milk, Bread}\} \rightarrow \{\text{Diaper, Beer}\},$

REVIEW: SET AND SUBSET

- $\{a, b, c, d\} \Leftrightarrow$ a set (there are one or more than one items)
- Subset \Leftrightarrow possible combinations of the items in a set

■ Possible sets:

6

- $\{a\}, \{b\}, \{c\}, \{d\}, \{a, b\}, \{a, c\}, \{a, d\}, \{b, c\}, \{b, d\}, \{c, d\}, \{a, b, c\}, \{a, b, d\}, \{a, c, d\}, \{b, c, d\}, \{a, b, c, d\}, \{\}$

■ What is the total number of the subset:

$$2^4 = 2^{\text{\#item}} = C_4^0 + C_4^1 + C_4^2 + C_4^3 + C_4^4$$

$$\text{No - Empty} : 2^4 - 1 = 2^{\text{\#of item}} - 1$$

subset:

- ① Empty $\{\} = C_4^0$
- ② 1 item $\{a\}, \{b\}, \{c\}, \{d\} = C_4^1$
- ③ 2 item $C_4^2 = \frac{4!}{2!2!} = \frac{4 \times 3 \times 2 \times 1}{2 \times 1 \times 2 \times 1} = 6$
- ③ 3 item $C_4^3 = 4$
- ④ 4 item $C_4^4 = 1$

ASSOCIATION RULE MINING

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

- Support count (σ)

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

of itemsets as a subset in datasets

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

$$\text{Support} = \frac{\sigma}{\text{size of Dataset}} = \frac{2}{5}$$

Fraction of transactions that contain an itemset

E.g. $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$

Frequent Itemset

An itemset whose support is greater than or equal to a minsup threshold

$$\text{minsup} = \min \sigma$$

$$\text{minsup} = 1 \leq f(\{\text{M, B, D}\}) = 2$$

DEFINITION: ASSOCIATION RULE

- Association Rule : $\text{conf of } X \rightarrow Y = \frac{\#X \cup Y}{\#X}$

– An implication expression of the form $X \rightarrow Y$, where X and Y are itemsets

– Example:

$\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

$\underbrace{\{\text{Milk, Diaper}\}}_{\text{2-itemset}} \rightarrow \underbrace{\{\text{Beer}\}}_{\text{1-itemset}}$

- Rule Evaluation Metrics

– Support (s) = $\{M, D, B\}$

◆ Fraction of transactions that contain both X and Y

– Confidence (c)

◆ Measures how often items in $Y = \{B\}$ appear in transactions that contain $X = \{M, D\}$

$$c = \frac{\# \{M, D, B\}}{\# \{M, D\}} = \frac{X \cup Y}{X} = \frac{2}{3}$$

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

$$S\{M, D, B\} = \frac{2}{5}$$

Example:

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

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

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

ASSOCIATION RULE MINING TASK

- Given a set of transactions T, the goal of association rule mining is to find all rules having

- support \geq minsup threshold *f*
- confidence \geq minconf threshold

Handwritten: $(B, M) \rightarrow \{D\}$

- Brute-force approach:

- List all possible association rules
- Compute the support and confidence for each rule
- Prune rules that fail the *minsup* and *minconf* thresholds

\Rightarrow **Computationally expensive / prohibitive!**

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

Handwritten: Given $X \rightarrow Y$
 $X = \{1, 2, 3, \dots, 6\}$

$\{a, b, c, d\} \rightarrow \# \text{ of subset:}$

Handwritten: $2^n - 1$
 $2^n = 2^6 = 64$
 $2^n - 1 = 64 - 1 = 63$

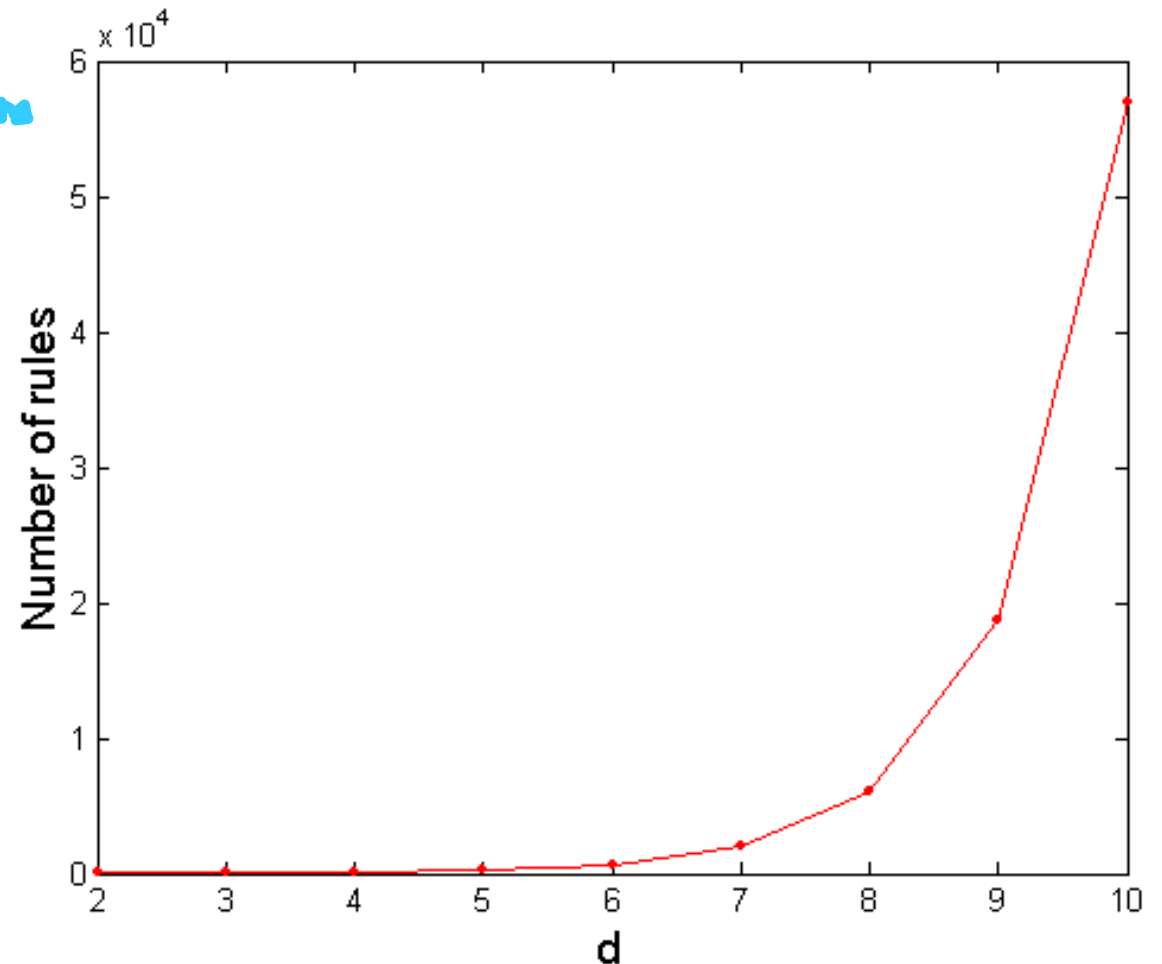
COMPUTATIONAL COMPLEXITY

- Given d unique items:
 - Total number of itemsets = 2^d *$d = \# \text{ of items}$*
 - Total number of possible association rules:

$$R = \sum_{k=1}^{d-1} \left[\binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right]$$
$$= 3^d - 2^{d+1} + 1$$

$X \rightarrow Y$

If $d=6$, $R = 602$ rules



MINING ASSOCIATION RULES

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 Rules:

$\{Milk, Diaper\} \rightarrow \{Beer\}$ (s=0.4, c=0.67)
 $\{Milk, Beer\} \rightarrow \{Diaper\}$ (s=0.4, c=1.0)
 $\{Diaper, Beer\} \rightarrow \{Milk\}$ (s=0.4, c=0.67)
 $\{Beer\} \rightarrow \{Milk, Diaper\}$ (s=0.4, c=0.67)
 $\{Diaper\} \rightarrow \{Milk, Beer\}$ (s=0.4, c=0.5)
 $\{Milk\} \rightarrow \{Diaper, Beer\}$ (s=0.4, c=0.5)

$\{a, b, c\}$
 $= \{b, c, a\}$
 $= \{c, a, b\}$
 $= \{c, b, a\}$

Observations: $s = \frac{\#(M, D, B)}{5} = \frac{2}{5} = 0.4$ $c = \frac{\#(M, D, B)}{\#(M)} = \frac{2}{4} = 0.5$

- All the above rules are binary partitions of the same itemset: {Milk, Diaper, Beer}
- Rules originating from the same itemset have identical support but can have different confidence
- Thus, we may decouple the support and confidence requirements

MINING ASSOCIATION RULES

- **Two-step approach:**

1. **Frequent Itemset Generation**

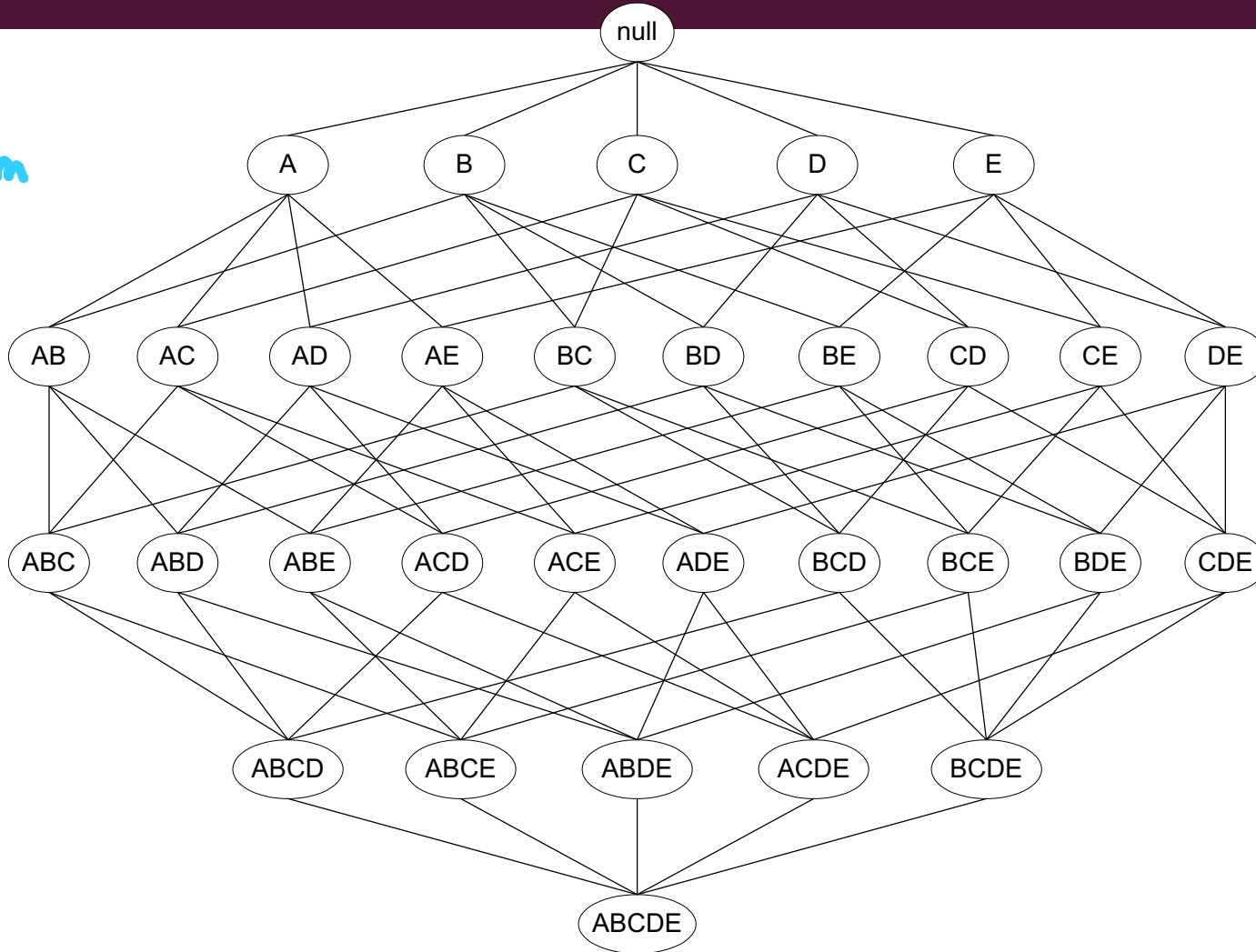
- Generate all itemsets whose support \geq minsup

2. **Rule Generation**

- Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

- **Frequent itemset generation is still computationally expensive**

FREQUENT ITEMSET GENERATION

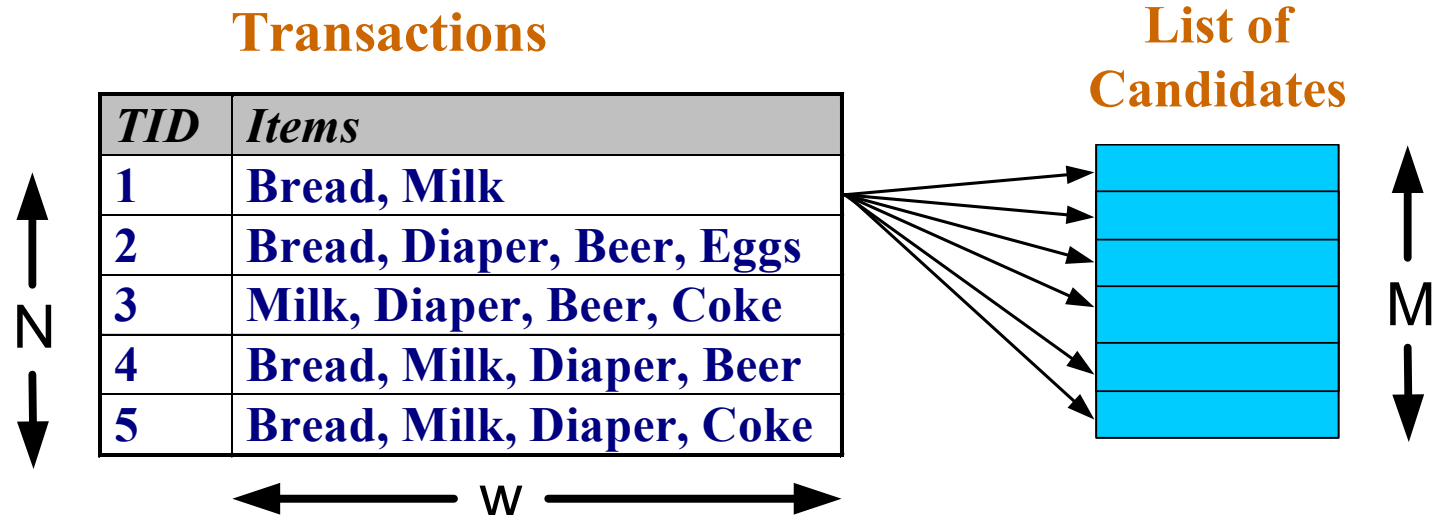


all pos. itemsets: $2^d = 2^5 = 32$

Given d items, there are 2^d possible candidate itemsets

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



- Match each transaction against every candidate
- Complexity $\sim O(NMw)$ => **Expensive since $M = 2^d$!!!**

FREQUENT ITEMSET GENERATION STRATEGIES

- Reduce the **number of candidates** (M) = # of itemsets = 2^d
 - Complete search: $M=2^d$
 - 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
- 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

