# Frequent Pattern Analysis:
## Apriori Algorithm

CS5483 Data Warehousing and Data Mining

# Market basket analysis

| Transaction ID (TID) | Purchased Items |
|---|---|
| 1 | Tissue paper, bread, milk |
| 2 | Tissue paper, bread, milk |
| 3 | Tissue paper, bread, jam |
| 4 | Tissue paper, bread, jam |
| 5 | Tissue paper, milk, diaper |

- Supermarket may want to:
  - Increase the stocks of **frequently purchased items**.
  - Bundle together items that are *frequently purchased together*.
- How to learn from transactional data?

# Problem formulation

| Transaction ID (TID) | Purchased Items |
|---|---|
| 1 | Tissue paper, bread, milk |
| 2 | Tissue paper, bread, milk |
| 3 | Tissue paper, bread, jam |
| 4 | Tissue paper, bread, jam |
| 5 | Tissue paper, milk, diaper |

- $\mathcal{I} := \{I_1, \ldots, I_m\}$: Set of **items**, e.g.,

$$\mathcal{I} = \{\text{tissue paper,bread,milk,jam,diaper}\}$$

- $A \subseteq \mathcal{I}$: I_____, written as $A = \{I_{j_1}, I_{j_2}, \ldots\}$ with $1 \leq j_1 < j_2 < \cdots \leq m$.

- **Transactional data**: $D := \{T_i\}_{i=1}^{n}$ where $T_i \subseteq \mathcal{I}$, e.g.,

$$T_1 = \{\text{tissue paper,bread,milk}\}$$

# Problem formulation

| Transaction ID (TID) | Purchased Items |
|---|---|
| 1 | Tissue paper, bread, milk |
| 2 | Tissue paper, bread, milk |
| 3 | Tissue paper, bread, jam |
| 4 | Tissue paper, bread, jam |
| 5 | Tissue paper, milk, diaper |

- **(Support) count** of itemset $A \subseteq \mathcal{I}$ in $D$:
$$\text{count}(A) := |\{T \in D | A \subseteq T\}|$$

- Objective: Given min_sup > 0, obtain the set of **f_____** $k$-**itemsets**
$$L_k := \{A \subseteq \mathcal{I} | |A| = k, \text{count}(A) \geq \text{min\_sup}\}, \quad \text{for all } k \in \{1, \dots, m\}.$$

# Example

$$T_1 = \{\qquad I_3,\ I_4,\ I_5\}$$
$$T_2 = \{\qquad I_3,\ I_4,\ I_5\}$$
$$T_3 = \{\quad I_2, \qquad I_4,\ I_5\}$$
$$T_4 = \{\quad I_2, \qquad I_4,\ I_5\}$$
$$T_5 = \{I_1, \qquad I_3, \qquad I_5\}$$

- $\mathcal{I} = \{I_1, I_2, I_3, I_4, I_5\}$ and $D = \{T_1, T_2, T_3, T_4, T_5\}$.

- min_sup = 2.

  - $\{I_2, I_5\}$ is a frequent itemset. Any others? _____

  - $\{I_2, I_3, I_5\}$ is not a frequent itemset. Any others? _____

# Compute frequent 1-itemsets

$T_1 = \{\quad\ \ 3,\ 4,\ 5\}$

$T_2 = \{\quad\ \ 3,\ 4,\ 5\}$

$T_3 = \{\ \ 2,\quad\ 4,\ 5\}$

$T_4 = \{\ \ 2,\quad\ 4,\ 5\}$

$T_5 = \{1,\quad 3,\quad\ 5\}$

with $I_j := j$ for notational simplicity.

$C_1$

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

$L_1$

| | |
|---|---|
| {2} | |
| {3} | |
| {4} | |
| {5} | |

- Generate a **c_____ list** $C_1$ consisting of singleton sets.

- Look up the count of every $A \in C_1$

- Add $A$ from $C_1$ to $L_1$ if $\text{count}(A) \geq \text{min\_sup}$.

# Compute frequent 2-itemsets

$$T_1 = \{\quad\quad 3,\ 4,\ 5\}$$
$$T_2 = \{\quad\quad 3,\ 4,\ 5\}$$
$$T_3 = \{\ 2,\quad\quad 4,\ 5\}$$
$$T_4 = \{\ 2,\quad\quad 4,\ 5\}$$
$$T_5 = \{1,\quad 3,\quad\quad 5\}$$

$C_2$

| | |
|---|---|
| {2,3} | |
| {2,4} | |
| {2,5} | |
| {3,4} | |
| {3,5} | |
| {4,5} | |

$L_2$

| | |
|---|---|
| {2,4} | |
| {2,5} | |
| {3,4} | |
| {3,5} | |
| {4,5} | |

$L_1$

| | |
|---|---|
| {2} | |
| {3} | |
| {4} | |
| {5} | |

➤ Generate the candidate list $C_2$ of 2-itemsets with each item taken from $L_1$.

➤ Why not take from $\mathcal{I}$ instead of $L_1$? (i.e., why not consider item 1?)

# Apriori property

- count($A$) is **anti**_____:

$$\text{count}(A) \leq \text{count}(B), \qquad \text{for any non−empty } B \subseteq A \subseteq \mathcal{I}$$

  - All non-empty subsets of a frequent itemset must be <u>frequent/infrequent</u>.
  - All supersets of a infrequent itemset must be <u>frequent/infrequent</u>.
- Other such functions?
  - Minimum value of items in $A$. <u>Y/N</u>
  - Maximum value of items in $A$. <u>Y/N</u>
  - Total value of items in $A$. <u>Y/N</u>
  - $A \mapsto \sum_{T \in D: A \subseteq T} x_T$ when $x_T$'s are non-negative. <u>Y/N</u>

# Compute frequent 3-itemsets

| | $C_3$ | | $L_3$ | | $L_2$ | | $L_1$ | |
|---|---|---|---|---|---|---|---|---|
| $T_1 = \{\quad\ 3,\ 4,\ 5\}$ | {2,4,5} | | {2,4,5} | | {2,4} | | {2} | |
| $T_2 = \{\quad\ 3,\ 4,\ 5\}$ | {3,4,5} | | {3,4,5} | | {2,5} | | {3} | |
| $T_3 = \{\ 2,\quad 4,\ 5\}$ | | | | | {3,4} | | {4} | |
| $T_4 = \{\ 2,\quad 4,\ 5\}$ | | | | | {3,5} | | {5} | |
| $T_5 = \{1,\quad 3,\quad 5\}$ | | | | | {4,5} | | | |

- How to generate the candidate list $C_3$ of 3-itemsets?

    1. Join 3 items from $L_1$? E.g., join {2}, {3}, {4} $\in L_1$? Y/N

    2. Join 1 itemset from $L_2$ and 1 item from $L_1$? E.g., join {2,4} $\in L_2$ with {3} $\in L_1$? Y/N

    3. Join any 2 itemsets from $L_2$? E.g., join {2,4}, {3,5} $\in L_2$? Y/N

- **J_____ step:** Join two frequent $(k-1)$-itemsets with the sets of first $k-2$ items identical.

Why?

# Compute frequent 4-itemsets

$L_4 = \emptyset$      $L_3$      $L_2$      $L_1$

$T_1 = \{\quad\ 3,\ 4,\ 5\}$

$T_2 = \{\quad\ 3,\ 4,\ 5\}$

$T_3 = \{\ 2,\quad\ 4,\ 5\}$

$T_4 = \{\ 2,\quad\ 4,\ 5\}$

$T_5 = \{1,\quad 3,\quad\ 5\}$

| $L_3$ | |
|---|---|
| {2,4,5} | |
| {3,4,5} | |

| $L_2$ | |
|---|---|
| {2,4} | |
| {2,5} | |
| {3,4} | |
| {3,5} | |
| {4,5} | |

| $L_1$ | |
|---|---|
| {2} | |
| {3} | |
| {4} | |
| {5} | |

- **Join step:** Join two frequent $(k-1)$-itemsets with the sets of first $k-2$ items identical.
- $C_4$ and therefore $L_4$ are empty because _____.
- Any more tricks to shorten candidate lists?

# Compute frequent 1-itemsets

$T_1 = \{1, 2, 3\quad\}$

$T_2 = \{1, 2, 3\quad\}$

$T_3 = \{1, 2,\quad 4\ \}$

$T_4 = \{1, 2,\quad 4\ \}$

$T_5 = \{1,\quad 3,\quad 5\}$

$C_1$

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

$L_1$

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

➡ Rename item $i$ as $6 - i$ for the previous example.

# Compute frequent 2-itemsets

$T_1 = \{1, 2, 3 \quad \}$

$T_2 = \{1, 2, 3 \quad \}$

$T_3 = \{1, 2, \quad 4 \}$

$T_4 = \{1, 2, \quad 4 \}$

$T_5 = \{1, \quad 3, \quad 5\}$

$C_2$

| | |
|---|---|
| {1,2} | 4 |
| {1,3} | 3 |
| {1,4} | 2 |
| {2,3} | 2 |
| {2,4} | 2 |
| {3,4} | 0 |

$L_2$

| | |
|---|---|
| {1,2} | 4 |
| {1,3} | 3 |
| {1,4} | 2 |
| {2,3} | 2 |
| {2,4} | 2 |

$L_1$

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

➡ **Join step:** Join two frequent $(k-1)$-itemsets with the sets of first $k-2$ items identical.

# Compute frequent 3-itemsets

$T_1 = \{1, 2, 3\}$

$T_2 = \{1, 2, 3\}$

$T_3 = \{1, 2, \quad 4\}$

$T_4 = \{1, 2, \quad 4\}$

$T_5 = \{1, \quad 3, \quad 5\}$

$C_3$

| | |
|---|---|
| {1,2,3} | 2 |
| {1,2,4} | 2 |
| {1,3,4} | |
| {2,3,4} | |

$L_3$

| | |
|---|---|
| {1,2,3} | 2 |
| {1,2,4} | 2 |

$L_2$

| | |
|---|---|
| {1,2} | 4 |
| {1,3} | 3 |
| {1,4} | 2 |
| {2,3} | 2 |
| {2,4} | 2 |

$L_1$

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

- **Join step:** Join two frequent $(k-1)$-itemsets with the set of first $k-2$ items identical.
- **Prune step:** Remove an itemset from $C_k$ if any of its $(k-1)$-subsets is not in $L_{k-1}$.

# Compute frequent 4-itemsets

$$C_4 = \emptyset = L_4$$

$T_1 = \{1, 2, 3\}$

$T_2 = \{1, 2, 3\}$

$T_3 = \{1, 2, \quad 4\}$

$T_4 = \{1, 2, \quad 4\}$

$T_5 = \{1, \quad 3, \quad 5\}$

~~{1,2,3,4}~~

$L_3$

| | |
|---|---|
| {1,2,3} | 2 |
| {1,2,4} | 2 |

$L_2$

| | |
|---|---|
| {1,2} | 4 |
| {1,3} | 3 |
| {1,4} | 2 |
| {2,3} | 2 |
| {2,4} | 2 |

$L_1$

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

- **Join step:** Join two frequent $(k-1)$-itemsets with the sets of first $k-2$ items identical.

- **Prune step:** Remove an itemset from $C_k$ if any of its $(k-1)$-subsets is not in $L_{k-1}$.

  - (Check at most _____ subsets, which are obtained by removing _____.)

# Apriori algorithm

- Compute frequent $k$-itemsets for $k$ from 1 to $m$:
  - **Join step:** Join two frequent $(k-1)$-itemsets with the sets of first $k-2$ items identical.
  - **Prune step:** Remove an itemset from $C_k$ if any of its $(k-1)$-subsets is not in $L_{k-1}$.
- Complexity: _____.
  - $|L_k|$ can go up to _____.
  - The total number of (non-empty) frequent itemsets can go up to _____.
- Can we compute and store $L$ more efficiently?
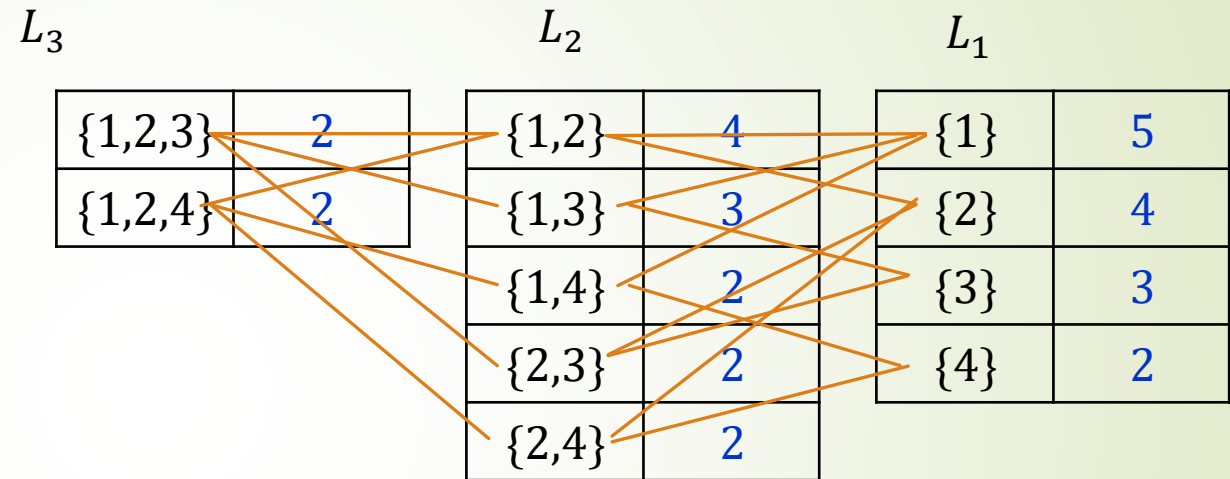
# How to store frequent itemsets

$T_1 = \{1, 2, 3\}$

$T_2 = \{1, 2, 3\}$

$T_3 = \{1, 2, \quad 4\}$

$T_4 = \{1, 2, \quad 4\}$

$T_5 = \{1, \quad 3, \quad 5\}$

$L_3$

| {1,2,3} | 2 |
|---------|---|
| {1,2,4} | 2 |

$L_2$

| {1,2} | 4 |
|-------|---|
| {1,3} | 3 |
| {1,4} | 2 |
| {2,3} | 2 |
| {2,4} | 2 |

$L_1$

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

- A frequent itemset is **m_____** iff all its proper supersets are not frequent.
- Can store only the maximal frequent itemsets because
  - All non-empty s_____ of a maximal frequent itemset are frequent.
  - Every frequent itemset is a s_____ of a maximal frequent itemset.
- How to store the counts of the frequent itemsets?
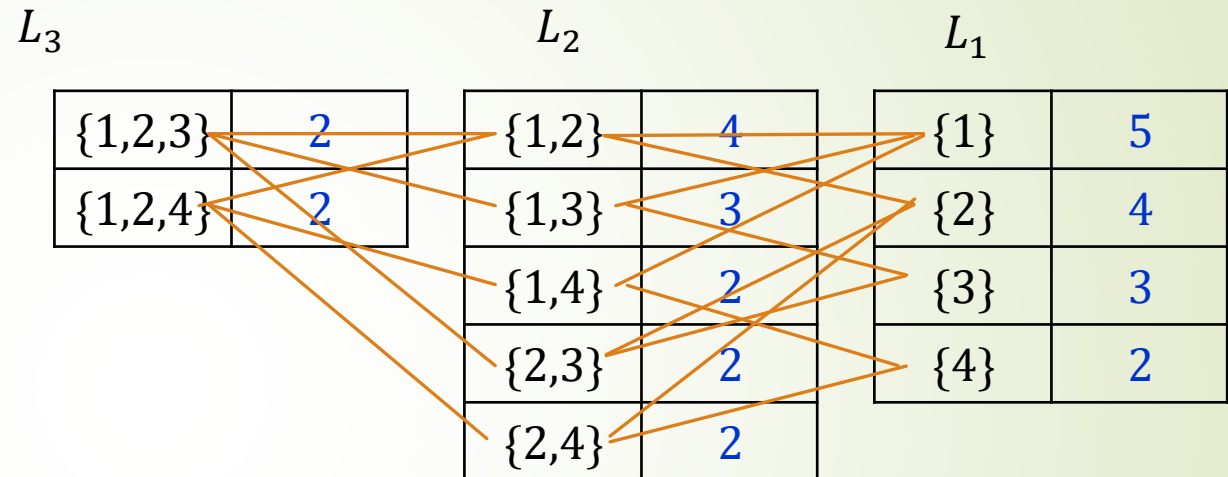
# How to store frequent itemsets with their counts

$T_1 = \{1, 2, 3\}$

$T_2 = \{1, 2, 3\}$

$T_3 = \{1, 2, \phantom{3,} 4\}$

$T_4 = \{1, 2, \phantom{3,} 4\}$

$T_5 = \{1, \phantom{2,} 3, \phantom{4,} 5\}$

$L_3$

| {1,2,3} | 2 |
|---------|---|
| {1,2,4} | 2 |

$L_2$

| {1,2} | 4 |
|-------|---|
| {1,3} | 3 |
| {1,4} | 2 |
| {2,3} | 2 |
| {2,4} | 2 |

$L_1$

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

- An itemset is **c**_____ iff no item can be added without decreasing the count.
- Can store only closed frequent itemsets and their counts because
  - maximal frequent itemsets are closed/not closed, so all frequent itemsets can be recovered.
  - a frequent itemset $A$ has the same count as a closed frequent itemset containing/contained by $A$ with largest/smallest count.
- [Challenge] How about storing itemsets where no item can be removed without increasing the count?

# Computational efficiency (Optional)

- Improving the efficiency of Apriori (6.2.3)
  - Hash-based technique
  - Transaction reduction
  - Partitioning
  - Sampling
  - Dynamic itemset counting
- Mining closed and Max Patterns (6.2.6)
  - Item merging
  - Sub-itemset pruning
  - Item skipping

# Association rules

- If a customer has bought all the items in $A$, is he/she likely to buy items in $B$? Why is it useful to know? **R**_____ **S**_____

- Association rules:

  - $A \subseteq \mathrm{T} \Rightarrow B \subseteq \mathrm{T}$ for a random transaction $\mathrm{T}$.

  - $A \Rightarrow B$ for short.

- Simplifying assumptions:

  - $A \cap B = \emptyset$: without loss of generality since $A \Rightarrow B$ is equivalent to $A \Rightarrow B \backslash A$.

  - $B \neq \emptyset$: avoid triviality since $A \Rightarrow \emptyset$ is always _____.

  - $A \neq \emptyset$ may be imposed because $\emptyset \Rightarrow B$ is not an a_____.

# Example
## Perfect rules

$T_1 = \{1, 2, 3\}$

$T_2 = \{1, 2, 3\}$

$T_3 = \{1, 2, \quad 4\}$

$T_4 = \{1, 2, \quad 4\}$

$T_5 = \{1, \quad 3, \quad 5\}$

- $\{2\} \Rightarrow \{1\}$? _____
- $\{3\} \Rightarrow \{1\}$? _____
- $\{4\} \Rightarrow \{1,2\}$? _____
- $\{3\} \Rightarrow \{1,2\}$? _____
- $\{5\} \Rightarrow \{1,3\}$? _____
- $\{3\} \Rightarrow \{5\}$? _____
- $\{1\} \Rightarrow \{5\}$? _____
- $\emptyset \Rightarrow \{1\}$? _____
- $3 \notin T \Rightarrow 4 \in T$? _____

# Measures of rule quality

- How many instances satisfy the antecedent?

$$\text{coverage}(A \Rightarrow B) := \frac{\text{count}(A)}{n}$$

- How many instances satisfy both the antecedent and consequence?

$$\text{support}(A \Rightarrow B) := \frac{\text{count}(A \cup B)}{n} \approx \Pr(A \cup B \subseteq \mathrm{T})$$

- Out of those instances that satisfy the antecedent, how many satisfy the consequence?

$$\text{confidence}(A \Rightarrow B) := \frac{\text{count}(A \cup B)}{\text{count}(A)} = \frac{\text{support}(A \Rightarrow B)}{\text{coverage}(A \Rightarrow B)} \approx \Pr(B \subseteq \mathrm{T} | A \subseteq \mathrm{T})$$

- **Support-confidence framework**: Prefer rules with high support and confidence.

# Example

$T_1 = \{1, 2, 3\}$
$T_2 = \{1, 2, 3\}$
$T_3 = \{1, 2, \quad 4\}$
$T_4 = \{1, 2, \quad 4\}$
$T_5 = \{1, \quad 3, \quad 5\}$

| Rule | Coverage | Support | Confidence |
|---|---|---|---|
| $\{2\} \Rightarrow \{1\}$ | 80% | 80% | 100% |
| $\{3\} \Rightarrow \{1\}$ | 60% | 60% | 100% |
| $\{4\} \Rightarrow \{1,2\}$ | 40% | 40% | 100% |
| $\{3\} \Rightarrow \{1,2\}$ | | | |
| $\{5\} \Rightarrow \{1,3\}$ | | | |
| $\{3\} \Rightarrow \{5\}$ | | | |
| $\{1\} \Rightarrow \{5\}$ | | | |

# Association rules from frequent itemsets

- Goal: Obtain all association rules with support $\geq s$ and confidence $\geq c$.

- How?

  1. Generate the list $L$ of frequent item sets with $\min\_sup = \lceil ns \rceil$.

  2. For $C \in L: |C| \geq 2$, find non-empty proper subset $A \subseteq C$ with $\text{count}(A) \leq \text{count}(C)/c$ to generate the rule

  $$A \Rightarrow B, \qquad \text{where } B := C \backslash A$$

- Correctness:

  - $\text{support}(A \Rightarrow B) \geq s$ iff $\text{count}(A \cup B) \geq \underline{\quad ns \quad}$

  - $\text{confidence}(A \Rightarrow B) \geq c$ iff $\text{count}(A) \leq \underline{\quad \text{count}(A \cup B)/c \quad}$

$$= \frac{\text{count}(A \cup B)}{\text{count}(A)} \geq c$$

$$C \in L$$

$$\frac{\text{count}(A \cup B)}{n} \geq \frac{\lceil ns \rceil}{n} \geq s$$

# Association rules from frequent itemsets

$T_1 = \{1, 2, 3\}$
$T_2 = \{1, 2, 3\}$
$T_3 = \{1, 2, \quad 4\}$
$T_4 = \{1, 2, \quad 4\}$
$T_5 = \{1, \quad 3, \quad 5\}$

50%

$L_3$

| | |
|---|---|
| {1,2,3} | 2 |
| {1,2,4} | 2 |

$L_2$

| | |
|---|---|
| {1,2} | 4 |
| {1,3} | 3 |
| {1,4} | 2 |
| {2,3} | 2 |
| {2,4} | 2 |

$L_1$

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

$s = 0.4$

- With min_sup=2, can generate all rules with support at least $\frac{2}{5} = 0.4$ _____.

- With $c = 0.6$ and $C = \{1,2,3\} \in L_3$, the desired association rules have count$(A) \leq \dfrac{count(\{1,2,3\})}{0.6}$

$= \dfrac{2}{0.6} = 3.\dot{3}$

{1,2} $\Rightarrow$ {3} Y/N          {1} $\Rightarrow$ {2,3} Y/N
{1,3} $\Rightarrow$ {2} Y/N          {2} $\Rightarrow$ {1,3} Y/N
{2,3} $\Rightarrow$ {1} Y/N          {3} $\Rightarrow$ {1,2} Y/N

- Exercise: Continue for other non-trivial choices of $C$: {1,2,4}, {1,2}, {1,3}, {1,4}, {2,3}, {2,4}.

# Limitation of support-confidence framework

$T_1 = \{1, 2, 3\}$

$T_2 = \{1, 2, 3\}$

$T_3 = \{1, 2, 4\}$

$T_4 = \{1, 2, 4\}$

$T_5 = \{1, 3, 5\}$

| # | Rule | Coverage | Support | Confidence |
|---|------|----------|---------|------------|
| 1 | $\{1,2\} \Rightarrow \{3\}$ | 80% | 40% | 50% |
| 2 | $\{1,3\} \Rightarrow \{2\}$ | 60% | 40% | 6̇6% |
| 3 | $\{2,3\} \Rightarrow \{1\}$ | 40% | 40% | 100% |
| 4 | $\{1\} \Rightarrow \{2,3\}$ | 100% | 40% | 40% |
| 5 | $\{2\} \Rightarrow \{1,3\}$ | 80% | 40% | 50% |
| 6 | $\{3\} \Rightarrow \{1,2\}$ | 60% | 40% | 6̇6% |
| ⋮ | | | | |

- Which rule above has the maximum confidence? Rule #___3___

- Is it the best rule that captures the strongest association? Y/N because
  _1 is purchased regardless of whether 2,3 are purchased._

# Limitation of support-confidence framework

$T_1 = \{1, 2, 3\}$

$T_2 = \{1, 2, 3\}$

$T_3 = \{1, 2, \quad 4\}$

$T_4 = \{1, 2, \quad 4\}$

$T_5 = \{1, \quad 3, \quad 5\}$

| # | Rule | Coverage | Support | Confidence | Prior | Lift |
|---|------|----------|---------|------------|-------|------|
| 1 | $\{1,2\} \Rightarrow \{3\}$ | 80% | 40% | 50% | 60% | 0.83 |
| 2 | $\{1,3\} \Rightarrow \{2\}$ | 60% | 40% | 66% | 80% | 0.83 |
| 3 | $\{2,3\} \Rightarrow \{1\}$ | 40% | 40% | 100% | 100% | 1 |
| 4 | $\{1\} \Rightarrow \{2,3\}$ | 100% | 40% | 40% | | |
| 5 | $\{2\} \Rightarrow \{1,3\}$ | 80% | 40% | 50% | | |
| 6 | $\{3\} \Rightarrow \{1,2\}$ | 60% | 40% | 66% | | |
| ⋮ | | | | | | |

$Pr(B \subseteq T \mid A \subseteq T)$

- How much more likely for the consequence to happen if the antecedent is satisfied?

$$\text{lift}(A \Rightarrow B) := \frac{\text{confidence}(A \Rightarrow B)}{\text{prior}(A \Rightarrow B)} \Bigg\{ \frac{\text{count}(B)}{n} \Bigg. = \frac{\text{count}(A \cup B) \cdot n}{\text{count}(A)\text{count}(B)}$$

$\frac{Pr(A \cup B \subseteq T)}{Pr(B \subseteq T)}$

- Lift is 1 if and only if $A \subseteq T$ and $B \subseteq T$ are independent.

- Give 2 rule(s) with positive association, i.e. lift>1: $\{5\} \Rightarrow \{1, 3\}$ with lift $= \frac{1}{0.6} = 1.6$

# References

- 6.2.1 Apriori Algorithm: Finding Frequent Itemsets by Confined Candidate Generation
- 6.2.2 Generating Association Rules from Frequent itemsets
- Optional:
  - 6.2.2 Generating Association Rules from Frequent itemsets
  - 6.2.3 Improving the Efficiency of Apriori
  - 6.2.6 Mining closed and Max Patterns
  - Azevedo, Paulo J., and Alípio M. Jorge. "Comparing rule measures for predictive association rules." *European Conference on Machine Learning*. Springer, Berlin, Heidelberg, 2007.