# Elements of Machine Learning & Data Science

# Frequent Itemsets

Lecture 9

Prof. Wil van der Aalst
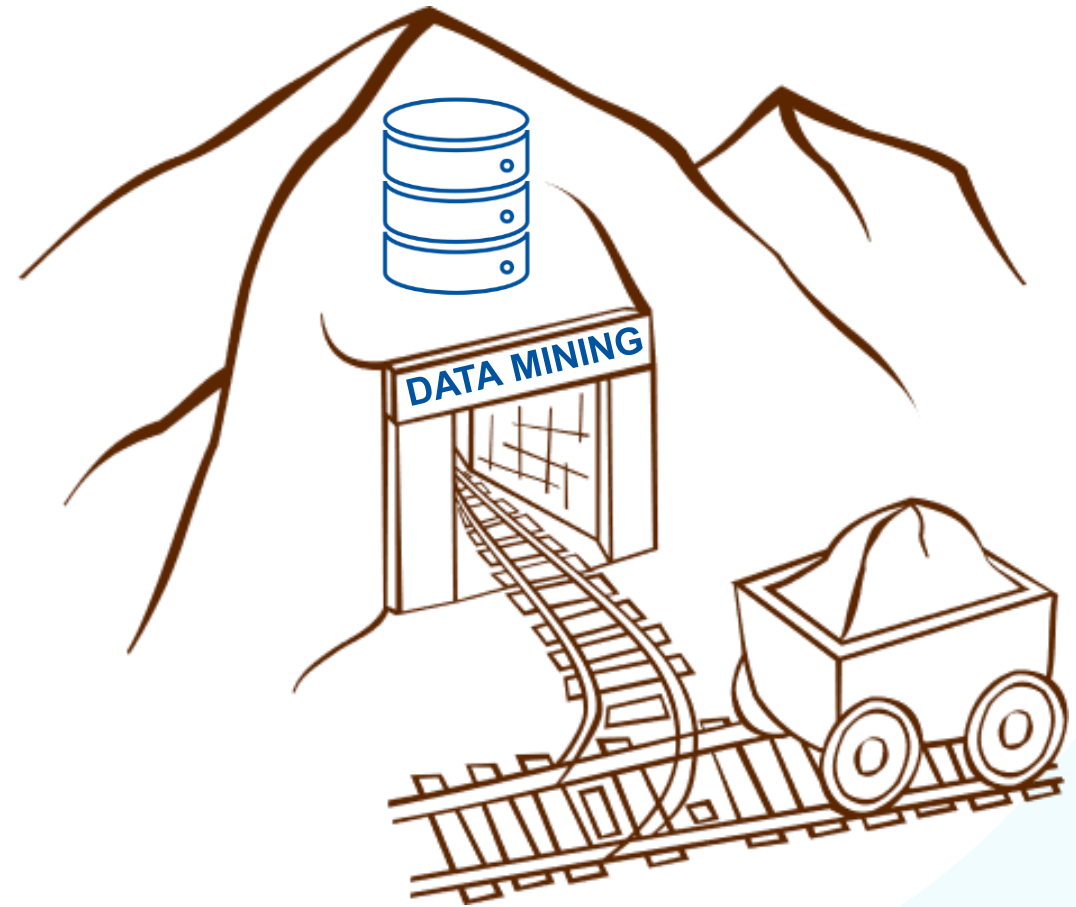
Marco Pegoraro, M.Sc.
Christopher Schwanen, M.Sc.
Tsunghao Huang, M.Sc.

# Frequent Itemsets

# Pattern Mining

- Finding surprising patterns in the input data


- Types of patterns:
  - Frequent itemsets
  - Association rules
  - Sequential patterns
  - Partial orders
  - Subgraphs

# Itemset Data

| ID | $f_1$ | $f_2$ | $f_3$ | $f_4$ | ... | $f_D$ |
|----|-------|-------|-------|-------|-----|-------|
| 1 | | | | | | |
| 2 | | | | | | |
| 3 | | | | | | |
| 4 | | | | | | |
| 5 | | | | | | |
| ... | | | | | | |

Each feature refers to an item (e.g., a product, disease, song, course, or error code)
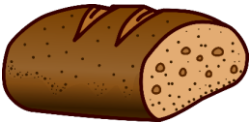
Each instance is a **transaction**

Each cell describes the presence of the item (Boolean 0/1 or Positive Integer Count)

# Itemset Data – Example

# Itemset Data – Example

| ID | 🧀 | 🍞 | 🍟 | 🥛 | ... | 🍝 |
|----|----|----|----|----|-----|----|
| 1  | 2  | 2  | 0  | 3  |     | 2  |
| 2  | 0  | 0  | 1  | 1  |     | 0  |
| 3  | 2  | 1  | 0  | 0  |     | 0  |
| 4  | 0  | 1  | 0  | 0  |     | 0  |
| 5  | 0  | 0  | 0  | 0  |     | 2  |
| ...| ...| ...| ...| ...|     | ...|

# Itemset Data – Example

| ID | 🧀 | 🍞 | 🍟 | 🥛 | ... | 🍝 |
|----|------|------|------|------|-----|------|
| 1 | True | True | False | True | | True |
| 2 | False | False | True | True | | False |
| 3 | True | True | False | False | | False |
| 4 | False | True | False | False | | False |
| 5 | False | False | False | False | | True |
| ... | ... | ... | ... | ... | | ... |

# Other Itemset Data Examples

| Rows | Columns |
|------|---------|
| EdX users | Courses taken |
| Spotify users | Songs Played |
| Netflix users | Movies Watched |
| Patients in hospital | Diseases |
| Repair bills | Components replaced |
| ... | ... |

# Application of Frequent Itemsets

| ID |  |  |  |  | ... |  |
|---|---|---|---|---|---|---|
| 1 | True | True | False | True | | True |
| 2 | False | False | True | True | | False |
| 3 | True | True | False | True | | False |
| 4 | False | True | False | True | | False |
| 5 | False | False | False | False | | True |
| ... | ... | ... | ... | ... | | ... |

Frequent Itemsets (movies)

NETFLIX

# Application of Frequent Itemsets

| ID |  |  |  |  | ... |  |
|----|------|------|------|------|-----|------|
| 1 | True | True | False | True | | True |
| 2 | True | True | True | False | | False |
| 3 | True | True | False | True | | False |
| 4 | True | False | False | False | | False |
| 5 | False | False | False | False | | True |
| ... | ... | ... | ... | ... | | ... |

Frequent Itemsets

amazon

- A notorious success story: the Tesco Clubcard

- Introduced in 1995, it was the first loyalty card with automatic data collection

- Widely regarded as responsible for Tesco's supremacy in the UK

- 1bn£ of increase in sales (4%) in one year

- Today, the Clubcard program is still incredibly profitable, even though Tesco gives away about 1bn£ in rewards and discounts each year!



"You know more about my customers after three months than I know after 30 years."

- Lord MacLaurin, chairman for Tesco, talking to the data scientists of the Clubcard program

# Frequent Itemsets – Notation

- $\mathcal{I} = \{I_1, I_2, \ldots, I_D\}$ is the set of all possible items

- $\mathcal{A} \subseteq \mathcal{I}$ is an itemset

- A transaction $\mathcal{T}$ is a non-empty itemset

- A dataset $\mathcal{X}$ is a collection of transactions

- Technically $\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$ such that $\emptyset \notin \mathcal{X}$
  ($\mathbb{M}$ is the multiset and $\mathbb{P}$ is the powerset operator)

# Frequent Itemsets – Notation Example

| ID | Che | Bre | Chi | Mil | ... | Pas |
|---|---|---|---|---|---|---|
| 1 | 2 (true) | 0 (false) | 0 (false) | 3 (true) | 0 (false) | 2 (true) |
| 2 | 0 (false) | 0 (false) | 1 (true) | 1 (true) | 0 (false) | 0 (false) |
| 3 | 2 (true) | 1 (true) | 0 (false) | 0 (false) | 0 (false) | 0 (false) |
| 4 | 0 (false) | 1 (true) | 0 (false) | 0 (false) | 0 (false) | 0 (false) |

$$\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$$

- Set of all items $\mathcal{I} = \{Che, Bre, Chi, Mil, \ldots, Pas\}$

- Transaction $\mathcal{T}_1 = \{Che, Mil, Pas\} \subseteq \mathcal{I}$

- Dataset with four transactions $\mathcal{X} = [\{Che, Mil, Pas\}, \{Chi, Mil\}, \{Che, Bre\}, \{Bre\}]$

- Dataset with ten transactions $\mathcal{X} = [\{Che, Mil, Pas\}^4, \{Chi, Mil\}^3, \{Che, Bre\}^2, \{Bre\}^1]$

# Frequent Itemsets – Notation Generalization

| ID | Che | Bre | Chi | Mil | ... | Pas |
|----|-----|-----|-----|-----|-----|-----|
| 1 | 2 (true) | 0 (false) | 0 (false) | 3 (true) | 0 (false) | 2 (true) |
| 2 | 0 (false) | 0 (false) | 1 (true) | 1 (true) | 0 (false) | 0 (false) |
| 3 | 2 (true) | 1 (true) | 0 (false) | 0 (false) | 0 (false) | 0 (false) |
| 4 | 0 (false) | 1 (true) | 0 (false) | 0 (false) | 0 (false) | 0 (false) |

$$\mathcal{X} \in \mathbb{M}(\mathbb{M}(\mathcal{I}))$$

- Set of all items $\mathcal{I} = \{Che, Bre, Chi, Mil, \ldots, Pas\}$

- Transaction $\mathcal{T}_1 = [Che^2, Mil^3, Pas^2] \in \mathbb{M}(\mathcal{I})$

- Dataset with four transactions $\mathcal{X} = [[Che^2, Mil^3, Pas^2], [Chi, Mil], [Che^2, Bre], [Bre]]$

We will consider only itemsets that are proper sets (not multisets). However, generalization is trivial.

## Frequent Itemsets – Support

$$\text{support}(\mathcal{A}) = \frac{|[\mathcal{T} \in \mathcal{X} | \mathcal{A} \subseteq \mathcal{T}]|}{|\mathcal{X}|}$$

(relative)

> Fraction of transactions $\mathcal{T}$ in dataset $\mathcal{X}$ that cover the itemset $\mathcal{A}$

$$\text{support\_count}(\mathcal{A}) = |[\mathcal{T} \in \mathcal{X} | \mathcal{A} \subseteq \mathcal{T}]|$$

(absolute, also called frequency or count)

## Frequent Itemsets – Support

$$\text{support}(\mathcal{A}) = \frac{|[\mathcal{T} \in \mathcal{X} | \mathcal{A} \subseteq \mathcal{T}]|}{|\mathcal{X}|}$$

(relative)

Fraction of transactions $\mathcal{T}$ in dataset $\mathcal{X}$ that cover the itemset $\mathcal{A}$

$$\text{support\_count}(\mathcal{A}) = |[\mathcal{T} \in \mathcal{X} | \mathcal{A} \subseteq \mathcal{T}]|$$

(absolute, also called frequency or count)

- Minimum **support threshold** $min\_sup$: lower bound for $support(A)$

- An itemset is **frequent** if its support is higher than $min\_sup$

- Frequent itemsets are used to find **association rules**

# Support – Example

| ID | Che | Bre | Chi | Mil | ... | Pas |
|----|-----|-----|-----|-----|-----|-----|
| 1 | 2 (true) | 0 (false) | 0 (false) | 3 (true) | 0 (false) | 2 (true) |
| 2 | 0 (false) | 0 (false) | 1 (true) | 1 (true) | 0 (false) | 0 (false) |
| 3 | 2 (true) | 1 (true) | 0 (false) | 0 (false) | 0 (false) | 0 (false) |
| 4 | 1 (true) | 1 (true) | 0 (false) | 1 (true) | 0 (false) | 0 (false) |

$$\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$$

Dataset $\mathcal{X} = [\{Che, Mil, Pas\}, \{Chi, Mil\}, \{Che, Bre\}, \{Che, Bre, Mil\}]$

# Support – Example

| ID | Che | Bre | Chi | Mil | ... | Pas |
|----|-----|-----|-----|-----|-----|-----|
| 1 | 2 (true) | 0 (false) | 0 (false) | 3 (true) | 0 (false) | 2 (true) |
| 2 | 0 (false) | 0 (false) | 1 (true) | 1 (true) | 0 (false) | 0 (false) |
| 3 | 2 (true) | 1 (true) | 0 (false) | 0 (false) | 0 (false) | 0 (false) |
| 4 | 1 (true) | 1 (true) | 0 (false) | 1 (true) | 0 (false) | 0 (false) |

$$\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$$

Dataset $\mathcal{X} = [\{Che, Mil, Pas\}, \{Chi, Mil\}, \{Che, Bre\}, \{Che, Bre, Mil\}]$

Itemset $\mathcal{A} = \{Che, Mil\} \subseteq \mathcal{I}$

$\text{support\_count}(\mathcal{A}) = |[\mathcal{T} \in \mathcal{X} \mid \mathcal{A} \subseteq \mathcal{T}]| = |[\mathcal{T}_1, \mathcal{T}_4]| = 2$

$\text{support}(\mathcal{A}) = \frac{|[\mathcal{T} \in \mathcal{X} \mid \mathcal{A} \subseteq \mathcal{T}]|}{|\mathcal{X}|} = \frac{|[\mathcal{T}_1, \mathcal{T}_4]|}{4} = \frac{2}{4}$

$\mathcal{A}$ is **frequent** if $\text{min\_sup} \leq 0.5$

# Support – Example

| ID | Che | Bre | Chi | Mil | ... | Pas |
|----|-----|-----|-----|-----|-----|-----|
| 1 | 2 (true) | 0 (false) | 0 (false) | 3 (true) | 0 (false) | 2 (true) |
| 2 | 0 (false) | 0 (false) | 1 (true) | 1 (true) | 0 (false) | 0 (false) |
| 3 | 2 (true) | 1 (true) | 0 (false) | 0 (false) | 0 (false) | 0 (false) |
| 4 | 1 (true) | 1 (true) | 0 (false) | 1 (true) | 0 (false) | 0 (false) |

$$\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$$

Dataset $\mathcal{X} = [\{Che, Mil, Pas\}, \{Chi, Mil\}, \{Che, Bre\}, \{Che, Bre, Mil\}]$

Itemset $\mathcal{A} = \{Che, Mil\} \subseteq \mathcal{I}$

$\text{support\_count}(\mathcal{A}) = |[\mathcal{T} \in \mathcal{X} \mid \mathcal{A} \subseteq \mathcal{T}]| = |[\mathcal{T}_1, \mathcal{T}_4]| = 2$

$\text{support}(\mathcal{A}) = \frac{|[\mathcal{T} \in \mathcal{X} \mid \mathcal{A} \subseteq \mathcal{T}]|}{|\mathcal{X}|} = \frac{|[\mathcal{T}_1, \mathcal{T}_4]|}{4} = \frac{2}{4}$

$\mathcal{A}$ is **frequent** if min_sup $\leq 0.5$

Itemset $\mathcal{B} = \{Mil\} \subseteq \mathcal{I}$

$\text{support\_count}(\mathcal{B}) = |[\mathcal{T} \in \mathcal{X} \mid \mathcal{A} \subseteq \mathcal{T}]| = |[\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_4]| = 3$

$\text{support}(\mathcal{B}) = \frac{|[\mathcal{T} \in \mathcal{X} \mid \mathcal{B} \subseteq \mathcal{T}]|}{|\mathcal{X}|} = \frac{|[\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_4]|}{4} = \frac{3}{4}$

$\mathcal{B}$ is **frequent** if min_sup $\leq 0.75$

# Support – Example

| ID | Che | Bre | Chi | Mil | ... | Pas |
|----|-----|-----|-----|-----|-----|-----|
| 1 | 2 (true) | 0 (false) | 0 (false) | 3 (true) | 0 (false) | 2 (true) |
| 2 | 0 (false) | 0 (false) | 1 (true) | 1 (true) | 0 (false) | 0 (false) |
| 3 | 2 (true) | 1 (true) | 0 (false) | 0 (false) | 0 (false) | 0 (false) |
| 4 | 1 (true) | 1 (true) | 0 (false) | 1 (true) | 0 (false) | 0 (false) |

$$\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$$

Itemset $\mathcal{A} = \{Che, Mil\} \subseteq \mathcal{I}$

$$\text{support}(\mathcal{A}) = \frac{|[\mathcal{T} \in \mathcal{X} | \mathcal{A} \subseteq \mathcal{T}]|}{|\mathcal{X}|} = \frac{|[\mathcal{T}_1, \mathcal{T}_4]|}{4} = \frac{2}{4}$$

Itemset $\mathcal{B} = \{Mil\} \subseteq \mathcal{I}$

$$\text{support}(\mathcal{B}) = \frac{|[\mathcal{T} \in \mathcal{X} | \mathcal{B} \subseteq \mathcal{T}]|}{|\mathcal{X}|} = \frac{|[\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_4]|}{4} = \frac{3}{4}$$

$$\mathcal{B} \subseteq \mathcal{A} \implies \text{support}(\mathcal{B}) \geq \text{support}(\mathcal{A})$$

General rule:
Subset of an itemset has higher or equal support than this itemset

# Support – Summary

Support

- A measure of the popularity (frequency) of an itemset.

- Calculated as the fraction of transactions in a dataset that contain the itemset.

$$\text{support}(\mathcal{A}) = \frac{|[\mathcal{T} \in \mathcal{X} | \mathcal{A} \subseteq \mathcal{T}]|}{|\mathcal{X}|}$$

- Any itemset with a support below the threshold is considered to be infrequent.

- Support is also used to find association rules

# Frequent Itemsets

## Problem Statement

Given dataset $\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$ and minimum support threshold $\boldsymbol{min\_sup}$, find all frequent non-empty itemsets:

$$\{\mathcal{A} \subseteq \mathcal{I} \mid \text{support}(\mathcal{A}) \geq \text{min\_sup}\}$$

# Naïve Approach

- Given $\mathcal{A} \subseteq \mathcal{I}$, it is possible to check whether $\boldsymbol{support}(\mathcal{A}) \geq \boldsymbol{min\_sup}$ by testing all transactions

- If there are $\boldsymbol{D}$ unique items, then there are $\boldsymbol{2^D - 1}$ candidate itemsets that can all be tested individually

- However, this can be very time consuming…

# Assume $D = 50\ 000$ products

$2^D - 1 =$

Good luck!

Of course most of them are not frequent

# Subsets of Frequent Itemsets Are Also Frequent

- Assume $\mathcal{A} = \{I_1, I_2, \ldots, I_{100}\}$ and $\textbf{\textit{support}}(\mathcal{A}) \geq \textbf{\textit{min\_sup}}$

- All subsets of $\mathcal{A}$ are also frequent

- There are $\binom{\textbf{100}}{\textbf{1}} = \textbf{100}$ frequent itemsets having 1 item

- There are $\binom{\textbf{100}}{\textbf{k}}$ frequent itemsets having $\textbf{\textit{k}}$ items ("100 choose $k$")

- There are $\binom{\textbf{100}}{\textbf{1}} + \binom{\textbf{100}}{\textbf{2}} + \cdots \binom{\textbf{100}}{\textbf{99}} = \textbf{2}^{\textbf{100}} - \textbf{2} = \textbf{1.27} \times \textbf{10}^{\textbf{30}}$ smaller frequent itemsets contained in $\mathcal{A}$

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} = \frac{n(n-1)\ldots(n-k+1)}{k(k-1)\ldots 1}$$

# Summary

- We should avoid exhaustively testing all candidate itemsets

- We need to focus on the "interesting" ones

→ Closed itemsets

# Closed Itemsets

- An itemset $\mathcal{A}$ is closed if there is no proper superset $\mathcal{B} \supset \mathcal{A}$ that has the same support

- If $\mathcal{A}$ is closed, then $support(\mathcal{A}) > support(\mathcal{B})$ for any $\mathcal{B} \supset \mathcal{A}$

adding any item to $\mathcal{A}$ will always **reduce support**

# Closed Frequent Itemsets

- An itemset $\mathcal{A}$ is closed if there is no proper superset $\mathcal{B} \supset \mathcal{A}$ that has the same support

- If $\mathcal{A}$ is closed, then $support(\mathcal{A}) > support(\mathcal{B})$ for any $\mathcal{B} \supset \mathcal{A}$

- $\mathcal{A}$ is frequent if its support is higher than threshold $min\_sup$

closed frequent itemsets are **closed and frequent**

# Maximal Frequent Itemsets

An itemset $\mathcal{A}$ is a maximal frequent itemset if:

- $\mathcal{A}$ is frequent

- there is no proper superset $\mathcal{B} \supset \mathcal{A}$ that is also frequent

**smaller** itemsets

**more** frequent

Maximal

$min\_sup$

**bigger** itemsets

**less** frequent

# Relationships

An itemset $\mathcal{A}$ is a closed frequent itemset if:

- $\mathcal{A}$ is frequent

- there is no proper superset $\mathcal{B} \supset \mathcal{A}$ that has the same support

An itemset $\mathcal{A}$ is a maximal frequent itemset if:

- $\mathcal{A}$ is frequent

- there is no proper superset $\mathcal{B} \supset \mathcal{A}$ that is also frequent

Hence, maximal frequent itemsets are closed by definition.

all frequent itemsets

closed
frequent
itemsets

maximal
frequent
itemsets

# Example

Assume:

$$\mathcal{I} = \{I_1, I_2, \ldots, I_{100}\}, \boldsymbol{min\_sup} = \frac{5}{20} = 0.25$$
$$\mathcal{X} = [\{I_1, I_2, \ldots, I_{50}\}^{10}, \{I_1, I_2, \ldots, I_{100}\}^{10}]$$

- There are $\boldsymbol{2^{100} - 1 = 1.27 \times 10^{30}}$ itemsets; all are frequent.

- There are two closed frequent itemsets:

    - $\mathcal{A} = \{I_1, I_2, \ldots, I_{50}\}$ *with* $\boldsymbol{support}(\mathcal{A}) = \frac{\mathbf{20}}{\mathbf{20}}$

    - $\mathcal{B} = \{I_1, I_2, \ldots, I_{100}\}$ *with* $\boldsymbol{support}(\mathcal{B}) = \frac{\mathbf{10}}{\mathbf{20}}$

- There is only one maximal frequent itemset: $\mathcal{B} = \{I_1, I_2, \ldots, I_{100}\}$

# Example

Assume:

$$\mathcal{I} = \{I_1, I_2, \ldots, I_{100}\}, \boldsymbol{min\_sup} = \frac{15}{20} = 0.75$$

$$\mathcal{X} = [\{I_1, I_2, \ldots, I_{50}\}^{10}, \{I_1, I_2, \ldots, I_{100}\}^{10}]$$

- There are $\mathbf{2^{50} - 1 = 3.17 \times 10^{15}}$ itemsets that are frequent.

- There is one closed frequent itemsets:

  - $\mathcal{A} = \{I_1, I_2, \ldots, I_{50}\} \, with \, \boldsymbol{support}(\mathcal{A}) = \frac{\mathbf{20}}{\mathbf{20}}$

- There is only one maximal frequent itemset: $\mathcal{A} = \{I_1, I_2, \ldots, I_{50}\}$

# Example

Assume:

$$\mathcal{I} = \{I_1, I_2, \ldots, I_{100}\}, \boldsymbol{min\_sup} = \frac{15}{20}$$

$$\mathcal{X} = [\{I_1, I_2, \ldots, I_{99}\}^{10}, \{I_2, I_3, \ldots, I_{100}\}^{10}]$$

- There are $\mathbf{2^{98} - 1 = 1.17 \times 10^{29}}$ itemsets that are frequent.

- There is one closed frequent itemset:

    - $\mathcal{A} = \{I_2, I_3, \ldots, I_{99}\} \, with \, \boldsymbol{support}(\mathcal{A}) = \frac{\mathbf{20}}{\mathbf{20}}$

- There is only one maximal frequent itemset: $\mathcal{A} = \{I_2, I_3, \ldots, I_{99}\}$

# Observations

- The supports of the closed frequent itemsets provide complete information about the supports of all frequent item sets

- Formally, assume:
    - $\mathcal{A} \subset \mathcal{B},$
    - $\mathcal{B}$ is a closed frequent itemset, and
    - there is no closed frequent itemset $\mathcal{B}'$ such that $\mathcal{A} \subseteq \mathcal{B}' \subset \mathcal{B}.$

    Then $\boldsymbol{support}(\mathcal{A}) = \boldsymbol{support}(\mathcal{B}).$

→ It suffices to store closed frequent itemsets
  (maximal frequent itemsets provide less information)

all frequent itemsets

closed frequent itemsets

maximal frequent itemsets

# **Summary**

- Both maximal frequent itemsets and closed frequent itemsets are subsets of frequent itemsets.

- Maximal frequent itemsets are closed by definition.

- Closed frequent itemsets provide a more comprehensive list of frequent patterns

all frequent itemsets

closed frequent itemsets

maximal frequent itemsets

no proper superset has the same support

no proper superset that is also frequent

# Frequent Itemsets

# Apriori Algorithm

- Introduced by Rakesh Agrawal and Ramakrishnan Srikant in "Fast Algorithms for Mining Association Rules in Large Databases. VLDB 1994: 487-499"

- Computes frequent itemsets / association rules in a dataset

- Uses a "bottom up" approach (starts with candidate itemsets of size one)

- Extends frequent subsets one item at a time (candidate generation)

- Avoids unnecessary checks by re-using information from smaller subsets and exploiting frequent itemsets' properties

# Apriori Algorithm – Basic Idea

$$\mathcal{L}_k = \{\mathcal{A} \subseteq \mathcal{I} \mid support(\mathcal{A}) \geq \mathbf{min\_sup} \wedge |\mathcal{A}| = k\}$$ frequent itemsets of length $k$

1. Candidate generation: use the set $\mathcal{L}_k$ of frequent itemsets of length k to generate the candidate set $\mathcal{C}_{k+1}$ of candidate itemsets with length $k+1$

# Apriori Algorithm – Basic Idea

$$\mathcal{L}_k = \{\mathcal{A} \subseteq \mathcal{I} \mid support(\mathcal{A}) \geq \mathbf{min\_sup} \wedge |\mathcal{A}| = k\}$$ frequent itemsets of length *k*

1. Candidate generation: use the set $\mathcal{L}_k$ of frequent itemsets of length k to generate the candidate set $\mathcal{C}_{k+1}$ of candidate itemsets with length *k+1*

does not need the input data (efficient)

# Apriori Algorithm – Basic Idea

$$\mathcal{L}_k = \{\mathcal{A} \subseteq \mathcal{I} \mid support(\mathcal{A}) \geq \mathbf{min\_sup} \wedge |\mathcal{A}| = k\}$$ ← frequent itemsets of length *k*

1. Candidate generation: use the set $\mathcal{L}_k$ of frequent itemsets of length k to generate the candidate set $\mathcal{C}_{k+1}$ of candidate itemsets with length *k+1*

2. Pruning (antimonotonicity): all nonempty subsets of a frequent itemset must also be frequent → superset of an infrequent itemset cannot be frequent

does not need the input data (efficient)

# Apriori Algorithm – Basic Idea

$$\mathcal{L}_k = \{\mathcal{A} \subseteq \mathcal{I} \mid \boldsymbol{support}(\mathcal{A}) \geq \mathbf{min\_sup} \wedge |\mathcal{A}| = k\}$$

frequent itemsets of length $k$

1. Candidate generation: use the set $\mathcal{L}_k$ of frequent itemsets of length k to generate the candidate set $\mathcal{C}_{k+1}$ of candidate itemsets with length $k+1$

2. Pruning (antimonotonicity): all nonempty subsets of a frequent itemset must also be frequent → superset of an infrequent itemset cannot be frequent

does not need the input data (efficient)

3. Testing candidates: use the dataset to filter the infrequent itemsets from $\mathcal{C}_{k+1}$ and obtain $\mathcal{L}_{k+1}$

needs the input data (inefficient)

# Apriori Algorithm – Basic Idea

$$\mathcal{L}_k = \{\mathcal{A} \subseteq \mathcal{I} \mid support(\mathcal{A}) \geq \mathbf{min\_sup} \wedge |\mathcal{A}| = k\}$$ ← frequent itemsets of length *k*

1. Candidate generation: use the set $\mathcal{L}_k$ of frequent itemsets of length k to generate the candidate set $\mathcal{C}_{k+1}$ of candidate itemsets with length *k+1*

   does not need the input data (efficient)

2. Pruning (antimonotonicity): all nonempty subsets of a frequent itemset must also be frequent → superset of an infrequent itemset cannot be frequent

3. Testing candidates: use the dataset to filter the infrequent itemsets from $\mathcal{C}_{k+1}$ and obtain $\mathcal{L}_{k+1}$

   needs the input data (inefficient)

# Apriori Algorithm – Basic Idea

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$$\mathcal{L}_1$$

# Apriori Algorithm – Basic Idea

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$$\mathcal{L}_1$$

generate
using $\mathcal{L}_1$

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$$\mathcal{C}_2$$

Prune using $\mathcal{L}_1$

# Apriori Algorithm – Basic Idea

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$\mathcal{L}_1$

generate
using $\mathcal{L}_1$

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$\mathcal{C}_2$

Prune using $\mathcal{L}_1$

scan dataset
keep frequent

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$\mathcal{L}_2$

# Apriori Algorithm – Basic Idea

# Apriori Algorithm – Basic Idea



| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$\mathcal{L}_1$

generate using $\mathcal{L}_1$

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$\mathcal{C}_2$

Prune using $\mathcal{L}_1$

scan dataset keep frequent

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$\mathcal{L}_2$

generate using $\mathcal{L}_2$

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$\mathcal{C}_3$

Prune using $\mathcal{L}_2$

$\ldots$

Prune using $\mathcal{L}_{k-1}$

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$\mathcal{L}_k$

scan dataset keep frequent

| Itemset | Support |
|---------|---------|
|         |         |
|         |         |
|         |         |

$\mathcal{C}_k$

# Apriori Algorithm – Basic Idea

# Candidate Generation – Leveling

Leveling is used to generate candidate itemset $\mathcal{C}_{k+1}$ from $\mathcal{L}_k$:

For any $\mathcal{A} \in \mathcal{L}_{k+1}$ there exist $\mathcal{A}', \mathcal{A}'' \in \mathcal{L}_k$ such that $\mathcal{A} = \mathcal{A}' \cup \mathcal{A}''$

Hence, we can obtain the candidate itemsets $\mathcal{C}_{k+1} \supseteq \mathcal{L}_{k+1}$ by joining suitable $\mathcal{A}', \mathcal{A}'' \in \mathcal{L}_k$!

all itemsets of length $k+1$

$\mathcal{C}_{k+1}$ before pruning

$\mathcal{C}_{k+1}$ after pruning

$\mathcal{L}_{k+1}$

# Candidate Generation – Leveling

Leveling is used to generate candidate itemset $\mathcal{C}_{k+1}$ from $\mathcal{L}_k$:

For any $\mathcal{A} \in \mathcal{L}_{k+1}$ there exist $\mathcal{A}', \mathcal{A}'' \in \mathcal{L}_k$ such that $\mathcal{A} = \mathcal{A}' \cup \mathcal{A}''$

Assume that the items are ordered $(I_1, I_2, \ldots)$ and that
$\mathcal{A} = \{I_1, I_2, \ldots, I_{k-1}, I_k, I_{k+1}\} \in \mathcal{L}_{k+1}$

If $\mathcal{A}$ is frequent, its subsets must be frequent, in particular:
$\mathcal{A}' = \{I_1, I_2, \ldots, I_{k-1}, I_k\} \in \mathcal{L}_k$
$\mathcal{A}'' = \{I_1, I_2, \ldots, I_{k-1}, I_{k+1}\} \in \mathcal{L}_k$
$\mathcal{A}' \cup \mathcal{A}'' = \{I_1, I_2, \ldots, I_{k-1}, I_k, I_{k+1}\} \in \mathcal{L}_{k+1} = \mathcal{A}$

all itemsets of length $k + 1$

$\mathcal{C}_{k+1}$ before pruning

$\mathcal{C}_{k+1}$ after pruning

$\mathcal{L}_{k+1}$

# Candidate Generation – Leveling

Leveling is used to generate candidate itemset $\mathcal{C}_{k+1}$ from $\mathcal{L}_k$:

For any $\mathcal{A} \in \mathcal{L}_{k+1}$ there exist $\mathcal{A}', \mathcal{A}'' \in \mathcal{L}_k$ such that $\mathcal{A} = \mathcal{A}' \cup \mathcal{A}''$

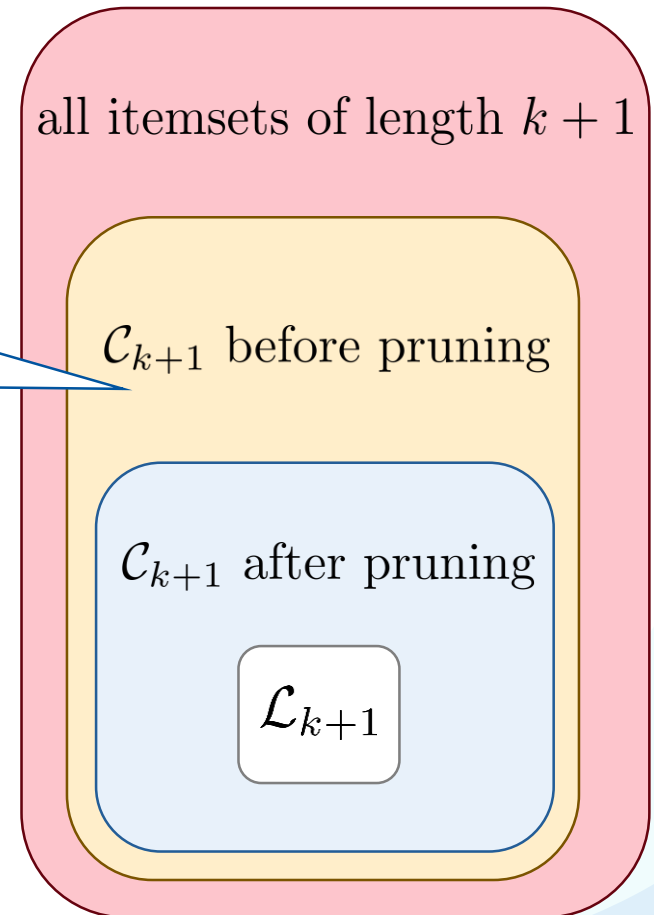Assume that the items are ordered $(I_1, I_2, \ldots)$ and that
$\mathcal{A} = \{I_1, I_2, \ldots, I_{k-1}, I_k, I_{k+1}\} \in \mathcal{L}_{k+1}$

If $\mathcal{A}$ is frequent, its subsets must be frequent, in particular:
$\mathcal{A}' = \{I_1, I_2, \ldots, I_{k-1}, I_k\} \in \mathcal{L}_k$
$\mathcal{A}'' = \{I_1, I_2, \ldots, I_{k-1}, I_{k+1}\} \in \mathcal{L}_k$
$\mathcal{A}' \cup \mathcal{A}'' = \{I_1, I_2, \ldots, I_{k-1}, I_k, I_{k+1}\} \in \mathcal{L}_{k+1} = \mathcal{A}$

$\Rightarrow$ We can generate $\mathcal{C}_{k+1}$ by joining itemsets $\mathcal{A}', \mathcal{A}'' \in \mathcal{L}_k$ which differ in one item

all itemsets of length $k + 1$

$\mathcal{C}_{k+1}$ before pruning

$\mathcal{C}_{k+1}$ after pruning

$\mathcal{L}_{k+1}$

# Candidate Generation

Thanks to leveling:

- Apriori creates the set of candidate itemsets of length k+1, $\mathcal{C}_{k+1}$, by joining two frequent itemsets of length *k*

- This can be done efficiently without creating duplicates

- Next, we prune the set $\mathcal{C}_{k+1}$ based on infrequent subsets

If $\mathcal{A}$ is frequent, its subsets must be frequent, in particular:
$$\mathcal{A}' = \{I_1, I_2, \ldots, I_{k-1}, I_k\} \in \mathcal{L}_k$$
$$\mathcal{A}'' = \{I_1, I_2, \ldots, I_{k-1}, I_{k+1}\} \in \mathcal{L}_k$$
$$\mathcal{A}' \cup \mathcal{A}'' = \{I_1, I_2, \ldots, I_{k-1}, I_k, I_{k+1}\} \in \mathcal{L}_{k+1} = \mathcal{A}$$

$\Rightarrow$ We can generate $\mathcal{C}_{k+1}$ by joining itemsets $\mathcal{A}', \mathcal{A}'' \in \mathcal{L}_k$ which differ in one item

all itemsets of length $k+1$

$\mathcal{C}_{k+1}$ before pruning

$\mathcal{C}_{k+1}$ after pruning

$\mathcal{L}_{k+1}$

# Pruning – Antimonotonicity

Antimonotonicity is used to prune the candidate set:

If $\mathcal{B}$ is a frequent itemset, any subset $\mathcal{A} \subseteq \mathcal{B}$ must be frequent
$\Rightarrow$ If a subset $\mathcal{A} \subseteq \mathcal{B}$ is infrequent, then $\mathcal{B}$ is infrequent

For any $\mathcal{A} \subseteq \mathcal{I}$ and $\mathcal{B} \subseteq \mathcal{I}$:

1. If $\mathcal{A} \subseteq \mathcal{B}$, then $\mathbf{support}(\mathcal{A}) \geq \mathbf{support}(\mathcal{B})$

2. If $\mathcal{A} \subseteq \mathcal{B}$ and $\mathbf{support}(\mathcal{B}) \geq \mathbf{min\_sup}$,
   then $\mathbf{support}(\mathcal{A}) \geq \mathbf{min\_sup}$

3. If $\mathcal{A} \subseteq \mathcal{B}$ and $\mathbf{support}(\mathcal{A}) < \mathbf{min\_sup}$,
   then $\mathbf{support}(\mathcal{B}) < \mathbf{min\_sup}$

all itemsets of length $k+1$

$\mathcal{C}_{k+1}$ before pruning

$\mathcal{C}_{k+1}$ after pruning

$\mathcal{L}_{k+1}$

# Pruning – Antimonotonicity

Antimonotonicity is used to prune the candidate set:

If $\mathcal{B}$ is a frequent itemset, any subset $\mathcal{A} \subseteq \mathcal{B}$ must be frequent
$\Rightarrow$ If a subset $\mathcal{A} \subseteq \mathcal{B}$ is infrequent, then $\mathcal{B}$ is infrequent

1. If {Apple} or {Banana} is infrequent,
   then {Apple, Banana} is infrequent

2. If {Grapes} or {Banana} is infrequent,
   then {Grapes, Banana} is infrequent

3. If {Apples} or {Grapes} is infrequent,
   then {Apples, Grapes} is infrequent

all itemsets of length $k + 1$

$\mathcal{C}_{k+1}$ before pruning
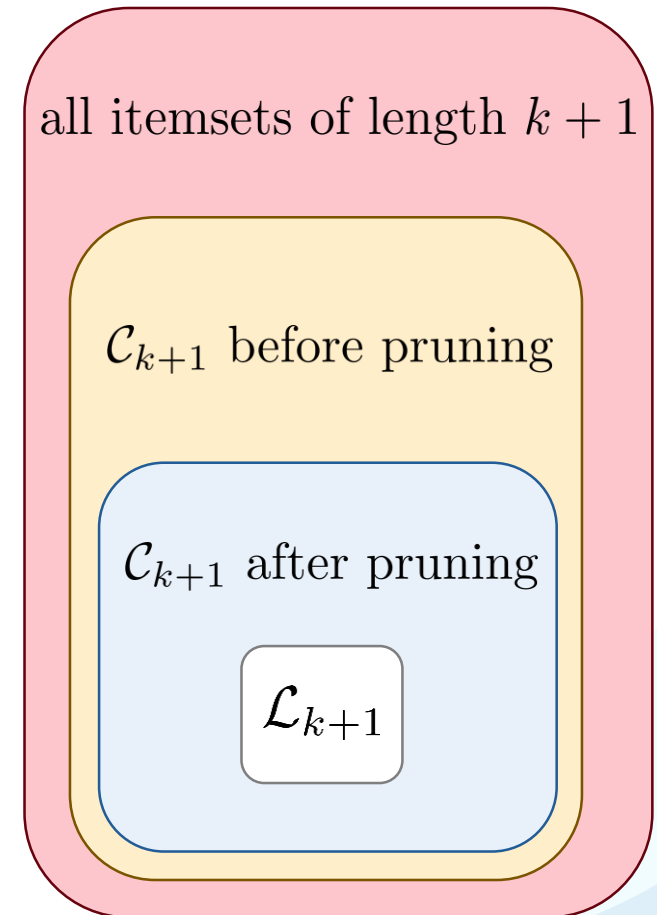
$\mathcal{C}_{k+1}$ after pruning

$\mathcal{L}_{k+1}$

# Pruning – Antimonotonicity

Antimonotonicity is used to prune the candidate set:

If $\mathcal{B}$ is a frequent itemset, any subset $\mathcal{A} \subseteq \mathcal{B}$ must be frequent
$\Rightarrow$ If a subset $\mathcal{A} \subseteq \mathcal{B}$ is infrequent, then $\mathcal{B}$ is infrequent

1.  If {Apple} or {Banana} is infrequent,
    then {Apple, Banana} is infrequent

2.  If {Grapes} or {Banana} is infrequent,
    then {Grapes, Banana} is infrequent

3.  If {Apples} or {Grapes} is infrequent,
    then {Apples, Grapes} is infrequent

| Bought Fruits | Support _count |
|---|---|
| {Banana} | 4 |
| {Grapes} | 1 |
| {Apple} | 2 |

$\text{min\_sup\_count} = 2$

all itemsets of length $k+1$

$\mathcal{C}_{k+1}$ before pruning
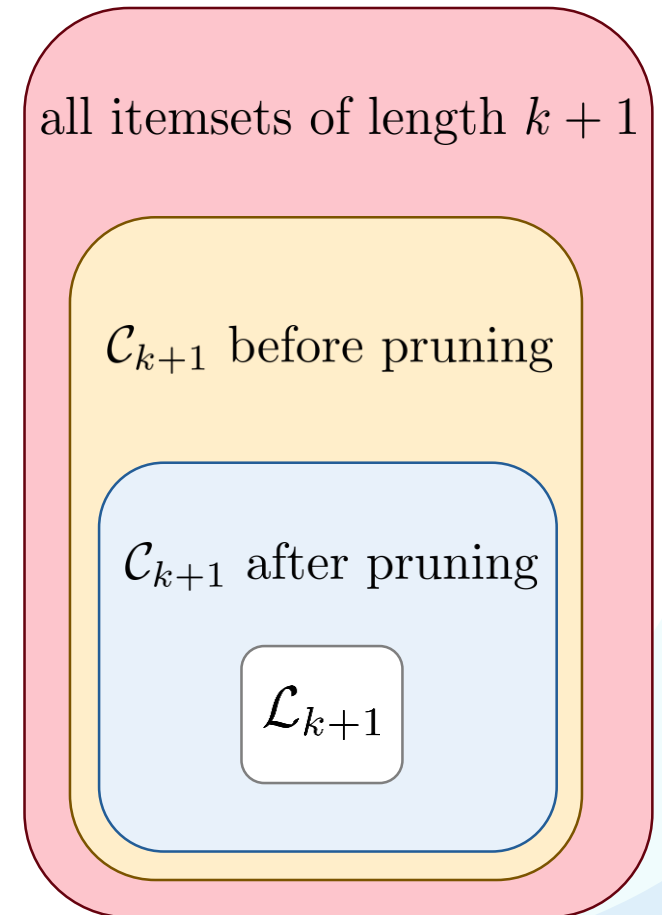
$\mathcal{C}_{k+1}$ after pruning
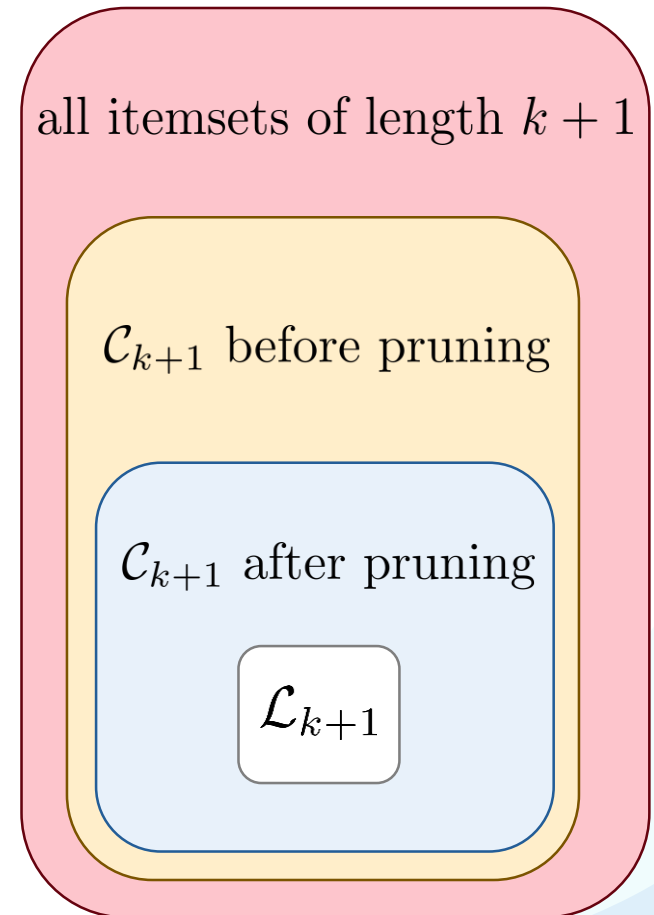
$\mathcal{L}_{k+1}$

# Pruning – Antimonotonicity

Antimonotonicity is used to prune the candidate set:

If $\mathcal{B}$ is a frequent itemset, any subset $\mathcal{A} \subseteq \mathcal{B}$ must be frequent
$\Rightarrow$ If a subset $\mathcal{A} \subseteq \mathcal{B}$ is infrequent, then $\mathcal{B}$ is infrequent

| ID | Bought Fruits |
|----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Grapes, Banana} |
| 3 | {Apple, Grapes} |

$$\text{min\_sup\_count} = 2$$

| ID | Bought Fruits |
|----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Grapes, Banana} |
| 3 | {Apple, Grapes} |

all itemsets of length $k + 1$

$\mathcal{C}_{k+1}$ before pruning

$\mathcal{C}_{k+1}$ after pruning

$\mathcal{L}_{k+1}$

# Testing Candidates

- After candidate generation and pruning test the remaining candidate itemsets

- We scan the dataset $\mathcal{X}$ and remove all infrequent candidate itemsets from $\mathcal{C}_{k+1}$ to obtain $\mathcal{L}_{k+1}$

all itemsets of length $k+1$

$\mathcal{C}_{k+1}$ before pruning

$\mathcal{C}_{k+1}$ after pruning

$\mathcal{L}_{k+1}$

# Testing Candidates

- After candidate generation and pruning test the remaining candidate itemsets

- We scan the dataset $\mathcal{X}$ and remove all infrequent candidate itemsets from $\mathcal{C}_{k+1}$ to obtain $\mathcal{L}_{k+1}$

- Consider all transactions $\mathcal{T} \in \mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$

- For each candidate itemset $\mathcal{A} \in \mathcal{C}_k$ increment the corresponding counter if $\mathcal{A} \subseteq \mathcal{T}_k$

- This returns the frequencies (**support_count**) of the candidate itemsets and we can compute $\mathcal{L}_{k+1}$ from $\mathcal{C}_{k+1}$

$$\mathcal{L}_{k+1} = \{\mathcal{A} \in \mathcal{C}_{k+1} | \mathbf{support}(\mathcal{A}) \geq \mathbf{min\_sup}\}$$

all itemsets of length $k+1$

$\mathcal{C}_{k+1}$ before pruning

$\mathcal{C}_{k+1}$ after pruning

$\mathcal{L}_{k+1}$

# Algorithm

**Apriori algorithm:**

1. Find the frequent itemsets of length 1 ($\mathcal{L}_1$)

2. Let $k \leftarrow 1$

3. **Repeat until** $\mathcal{L}_k \neq \varnothing$:

   (a) Generate set of candidate itemsets of length k+1 ($\mathcal{C}_{k+1}$) based on $\mathcal{L}_k$

   (b) Prune $\mathcal{C}_{k+1}$ from itemsets that have infrequent subsets

   (c) Test all remaining candidates from $\mathcal{C}_{k+1}$ and remove infrequent itemsets to obtain $\mathcal{L}_{k+1}$

   (d) Assign $k \leftarrow k + 1$

(Or until we find frequent itemsets of pre-defined length)

all itemsets of length $k + 1$

$\mathcal{C}_{k+1}$ before pruning

$\mathcal{C}_{k+1}$ after pruning

$\mathcal{L}_{k+1}$

# Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Grapes, Apple, Pineapple} |
| 2 | {Orange, Apple, Banana} |
| 3 | {Grapes, Orange, Apple, Banana} |
| 4 | {Orange, Banana} |
| 5 | {Grapes, Apple, Banana} |

$$\mathcal{X}$$

| Itemset | Count |
|---------|-------|
| {Grapes} | 3 |
| {Apple} | 4 |
| {Pineapple} | 1 |
| {Orange} | 3 |
| {Banana} | 4 |

$$\mathcal{C}_1$$

# Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Grapes, Apple, Pineapple} |
| 2 | {Orange, Apple, Banana} |
| 3 | {Grapes, Orange, Apple, Banana} |
| 4 | {Orange, Banana} |
| 5 | {Grapes, Apple, Banana} |

$\mathcal{X}$

| Itemset | Count |
|---------|-------|
| {Grapes} | 3 |
| {Apple} | 4 |
| {Pineapple} | 1 |
| {Orange} | 3 |
| {Banana} | 4 |

$\mathcal{C}_1$

$\mathbf{min\_sup\_count} = 2$

| Itemset | Count |
|---------|-------|
| {Grapes} | 3 |
| {Apple} | 4 |
| {Pineapple} | 1 |
| {Orange} | 3 |
| {Banana} | 4 |

$\mathcal{L}_1$

# Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Grapes, Apple, Pineapple} |
| 2 | {Orange, Apple, Banana} |
| 3 | {Grapes, Orange, Apple, Banana} |
| 4 | {Orange, Banana} |
| 5 | {Grapes, Apple, Banana} |

$$\mathcal{X}$$

| Itemset | Count |
|---------|-------|
| {Grapes} | 3 |
| {Apple} | 4 |
| {Pineapple} | 1 |
| {Orange} | 3 |
| {Banana} | 4 |

$$\mathcal{C}_1$$

$$\text{min\_sup\_count} = 2$$

| Itemset | Count |
|---------|-------|
| {Grapes} | 3 |
| {Apple} | 4 |
| {Pineapple} | 1 |
| {Orange} | 3 |
| {Banana} | 4 |

$$\mathcal{L}_1$$

generate candidates from $\mathcal{L}_1$

| Itemset |
|---------|
| {Grapes, Apple} |
| {Grapes, Orange} |
| {Grapes, Banana} |
| {Apple, Orange} |
| {Apple, Banana} |
| {Orange, Banana} |

$$\mathcal{C}_2$$

pruning based on $\mathcal{L}_1$

| Itemset | Count |
|---------|-------|
| {Grapes, Apple} | 3 |
| {Grapes, Orange} | 1 |
| {Grapes, Banana} | 2 |
| {Apple, Orange} | 2 |
| {Apple, Banana} | 3 |
| {Orange, Banana} | 3 |

$$\mathcal{C}_2$$

# Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Grapes, Apple, Pineapple} |
| 2 | {Orange, Apple, Banana} |
| 3 | {Grapes, Orange, Apple, Banana} |
| 4 | {Orange, Banana} |
| 5 | {Grapes, Apple, Banana} |

$\mathcal{X}$

| Itemset | Count |
|---------|-------|
| {Grapes} | 3 |
| {Apple} | 4 |
| {Pineapple} | 1 |
| {Orange} | 3 |
| {Banana} | 4 |

$\mathcal{C}_1$

$\mathbf{min\_sup\_count} = 2$

| Itemset | Count |
|---------|-------|
| {Grapes} | 3 |
| {Apple} | 4 |
| ~~{Pineapple}~~ | ~~1~~ |
| {Orange} | 3 |
| {Banana} | 4 |

$\mathcal{L}_1$

generate candidates from $\mathcal{L}_1$

| Itemset |
|---------|
| {Grapes, Apple} |
| {Grapes, Orange} |
| {Grapes, Banana} |
| {Apple, Orange} |
| {Apple, Banana} |
| {Orange, Banana} |

$\mathcal{C}_2$

pruning based on $\mathcal{L}_1$

| Itemset | Count |
|---------|-------|
| {Grapes, Apple} | 3 |
| {Grapes, Orange} | 1 |
| {Grapes, Banana} | 2 |
| {Apple, Orange} | 2 |
| {Apple, Banana} | 3 |
| {Orange, Banana} | 3 |

$\mathcal{C}_2$

$\mathbf{min\_sup\_count} = 2$

scan dataset test candidates

| Itemset | Count |
|---------|-------|
| {Grapes, Apple} | 3 |
| ~~{Grapes, Orange}~~ | ~~1~~ |
| {Grapes, Banana} | 2 |
| {Apple, Orange} | 2 |
| {Apple, Banana} | 3 |
| {Orange, Banana} | 3 |

$\mathcal{L}_2$

# Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Grapes, Apple, Pineapple} |
| 2 | {Orange, Apple, Banana} |
| 3 | {Grapes, Orange, Apple, Banana} |
| 4 | {Orange, Banana} |
| 5 | {Grapes, Apple, Banana} |

$$\mathcal{X}$$

$$\ldots$$

$$\text{min\_sup\_count} = 2$$

| Itemset | Count |
|---------|-------|
| {Grapes, Apple} | 3 |
| {Grapes, Banana} | 2 |
| {Apple, Orange} | 2 |
| {Apple, Banana} | 3 |
| {Orange, Banana} | 3 |

$$\mathcal{L}_2$$

# Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Grapes, Apple, Pineapple} |
| 2 | {Orange, Apple, Banana} |
| 3 | {Grapes, Orange, Apple, Banana} |
| 4 | {Orange, Banana} |
| 5 | {Grapes, Apple, Banana} |

$$\mathcal{X}$$

$$\cdots \qquad \mathrm{min\_sup\_count} = 2$$

| Itemset | Count |
|---------|-------|
| {Grapes, Apple} | 3 |
| {Grapes, Banana} | 2 |
| {Apple, Orange} | 2 |
| {Apple, Banana} | 3 |
| {Orange, Banana} | 3 |

$$\mathcal{L}_2$$

generate candidates from $\mathcal{L}_2$

| Itemset |
|---------|
| {Grapes, Apple, Banana} |
| {Grapes, Apple, Orange} |
| {Grapes, Banana, Orange} |
| {Apple, Banana, Orange} |

$$\mathcal{C}_3$$

pruning based on $\mathcal{L}_2$

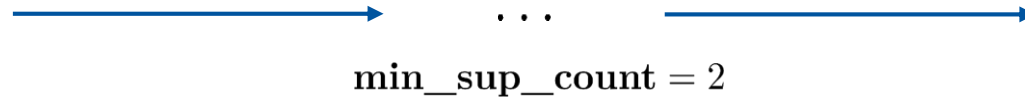| Itemset | Subsets of Length 2 |
|---------|---------------------|
| {Grapes, Apple, Banana} | {Grapes, Apple}, {Grapes, Banana}, {Apple, Banana} |
| {Grapes, Apple, Orange} | {Grapes, Apple}, {Grapes, Orange}, {Apple, Orange} |
| {Grapes, Banana, Orange} | {Grapes, Banana}, {Grapes, Orange}, {Banana, Orange} |
| {Apple, Banana, Orange} | {Apple, Banana}, {Apple, Orange}, {Banana, Orange} |

$$\mathcal{C}_3$$

$$\mathrm{min\_sup\_count} = 2$$

scan dataset test candidates

| Itemset | Count |
|---------|-------|
| {Grapes, Apple, Banana} | 2 |
| {Apple, Banana, Orange} | 2 |

$$\mathcal{L}_3$$

# Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Grapes, Apple, Pineapple} |
| 2 | {Orange, Apple, Banana} |
| 3 | {Grapes, Orange, Apple, Banana} |
| 4 | {Orange, Banana} |
| 5 | {Grapes, Apple, Banana} |

$$\mathcal{X}$$

$$\text{min\_sup\_count} = 2$$

$$\ldots$$

| Itemset | Count |
|---------|-------|
| {Grapes, Apple, Banana} | 2 |
| {Apple, Banana, Orange} | 2 |

$$\mathcal{L}_3$$

# Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Grapes, Apple, Pineapple} |
| 2 | {Orange, Apple, Banana} |
| 3 | {Grapes, Orange, Apple, Banana} |
| 4 | {Orange, Banana} |
| 5 | {Grapes, Apple, Banana} |

$\mathcal{X}$

. . .

$\text{min\_sup\_count} = 2$

| Itemset | Count |
|---------|-------|
| {Grapes, Apple, Banana} | 2 |
| {Apple, Banana, Orange} | 2 |

$\mathcal{L}_3$

generate candidates from $\mathcal{L}_3$

| Itemset |
|---------|
| {Grapes, Apple, Banana, Orange} |

$\mathcal{C}_4$

pruning based on $\mathcal{L}_3$

| Itemset | Subsets of length 3 |
|---------|---------------------|
| {Grapes, Apple, Banana, Orange} | {Grapes, Apple, Banana}, {Grapes, Apple, Orange}, {Grapes, Banana, Orange}, {Apple, Banana, Orange} |

No more candidates of length 4 – algorithm terminates

# Optimizations

**Further optimizations**

- Distributing the data (can be done in various ways)

- Gradually removing transactions not containing any frequent itemset of length $k$

- Sampling

**Limitations**

- It may remain challenging to generate the candidate sets (may be huge)

- Each candidate needs to be tested against the whole dataset

→ FP-Growth is an approach that aims to overcome these limitations

# Frequent Itemsets

1. Introduction

2. Properties of Frequent Itemsets

3. Apriori Algorithm

4. **FP-Growth Algorithm**

# Frequent Pattern Growth Algorithm

- Introduced by Jiawei Han, Jian Pei, Yiwen Yin in "Mining Frequent Patterns without Candidate Generation. SIGMOD Conference 2000: 1-12"

- Based on constructing the Frequent Pattern Tree (FP-Tree)

- Avoids generation of many candidates

- Depth-first rather than breadth-first

- Requires only two passes over the (potentially huge) dataset

# Motivation



**FP-Growth vs Apriori Runtime**

Graph is based on Jiawei Han, Jian Pei, Yiwen Yin in "Mining Frequent Patterns without Candidate Generation. SIGMOD Conference 2000: 1-12"

# FP-Growth Steps

1. Determine the frequency of each item (first pass through the dataset)

2. Sort $\mathcal{I} = \{I_1, \ldots, I_D\}$ based on their frequencies ($I_1$ is most frequent, $I_D$ is the least frequent)

3. Remove the non-frequent items

4. The remaining items in each transactions are ordered by frequency (same as above)

5. This can be used to build a so-called prefix tree (second pass trough the dataset)

# FP-Growth Steps

1. Determine the frequency of each item (first pass through the dataset)

2. Sort $\mathcal{I} = \{I_1, \ldots, I_D\}$ based on their frequencies ($I_1$ is most frequent, $I_D$ is the least frequent)

3. Remove the non-frequent items

4. The remaining items in each transactions are ordered by frequency (same as above)

5. This can be used to build a so-called prefix tree (second pass trough the dataset)

6. The resulting FP-tree contains all information needed to find the frequent itemsets of any length (no need to traverse the dataset again)

# Constructing FP-Tree – Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Grapes, Banana, Pineapple} |
| 3 | {Apple, Banana} |
| 4 | {Apple, Grapes, Pineapple} |
| 5 | {Grapes, Banana} |
| 6 | {Apple, Banana, Grapes} |

Dataset $\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{J}))$

🍌 5

🍎 4

🍇 4

🍍 2

1. Determine the frequencies of items
2. Order the items based on frequency

# Constructing FP-Tree – Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Grapes, Banana, Pineapple} |
| 3 | {Apple, Banana} |
| 4 | {Apple, Grapes, Pineapple} |
| 5 | {Grapes, Banana} |
| 6 | {Apple, Banana, Grapes} |

Dataset $\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{J}))$



$\text{min\_sup\_count} = 3$

1. Determine the frequencies of items
2. Order the items based on frequency

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Banana, Grapes} |
| 3 | {Banana, Apple} |
| 4 | {Apple, Grapes} |
| 5 | {Banana, Grapes} |
| 6 | {Banana, Apple, Grapes} |

3. Remove non-frequent items
4. Sort the items in the transactions

# Constructing FP-Tree – Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Banana, Grapes} |
| 3 | {Banana, Apple} |
| 4 | {Apple, Grapes} |
| 5 | {Banana, Grapes} |
| 6 | {Banana, Apple, Grapes} |

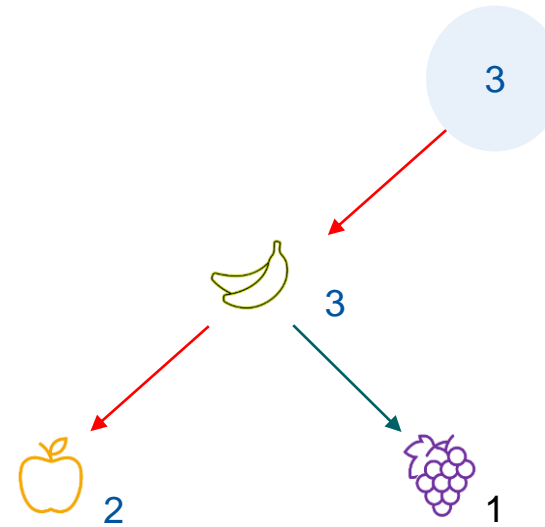5. Build the FP-tree going through each transaction

0

🍌 5

🍎 4

🍇 4

# Constructing FP-Tree – Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | **{Banana, Apple}** |
| 2 | {Banana, Grapes} |
| 3 | {Banana, Apple} |
| 4 | {Apple, Grapes} |
| 5 | {Banana, Grapes} |
| 6 | {Banana, Apple, Grapes} |

5. Build the FP-tree going through each transaction

# Constructing FP-Tree – Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Banana, Apple} |
| 2 | **{Banana, Grapes}** |
| 3 | {Banana, Apple} |
| 4 | {Apple, Grapes} |
| 5 | {Banana, Grapes} |
| 6 | {Banana, Apple, Grapes} |

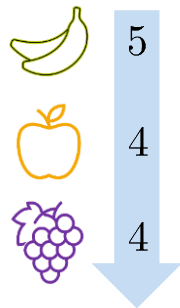5. Build the FP-tree going through each transaction

# Constructing FP-Tree – Example

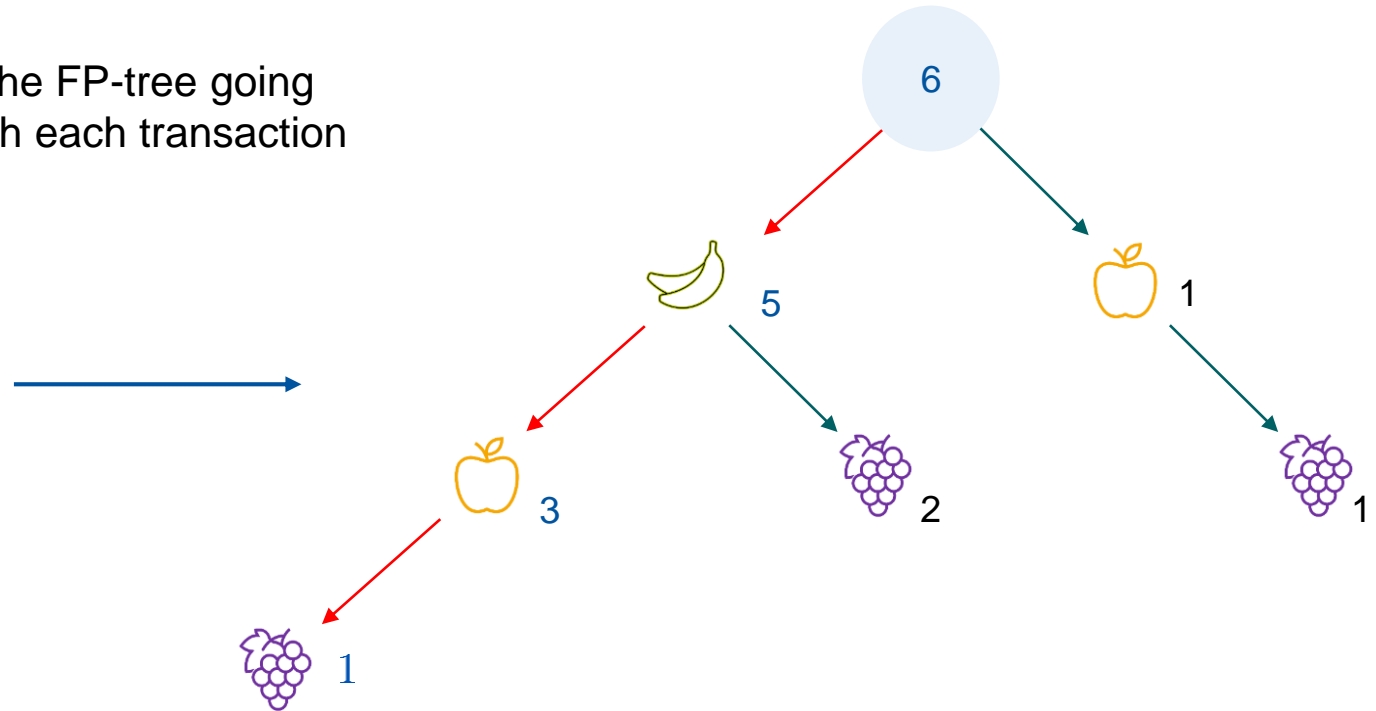| TID | Bought Fruits |
|-----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Banana, Grapes} |
| 3 | **{Banana, Apple}** |
| 4 | {Apple, Grapes} |
| 5 | {Banana, Grapes} |
| 6 | {Banana, Apple, Grapes} |

5. Build the FP-tree going through each transaction

# Constructing FP-Tree – Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Banana, Grapes} |
| 3 | {Banana, Apple} |
| 4 | **{Apple, Grapes}** |
| 5 | {Banana, Grapes} |
| 6 | {Banana, Apple, Grapes} |

5. Build the FP-tree going through each transaction

# Constructing FP-Tree – Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Banana, Grapes} |
| 3 | {Banana, Apple} |
| 4 | {Apple, Grapes} |
| 5 | **{Banana, Grapes}** |
| 6 | {Banana, Apple, Grapes} |

5. Build the FP-tree going through each transaction

# Constructing FP-Tree – Example

| TID | Bought Fruits |
|-----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Banana, Grapes} |
| 3 | {Banana, Apple} |
| 4 | {Apple, Grapes} |
| 5 | {Banana, Grapes} |
| 6 | **{Banana, Apple, Grapes}** |

5. Build the FP-tree going through each transaction

# Constructing FP-Tree – Example

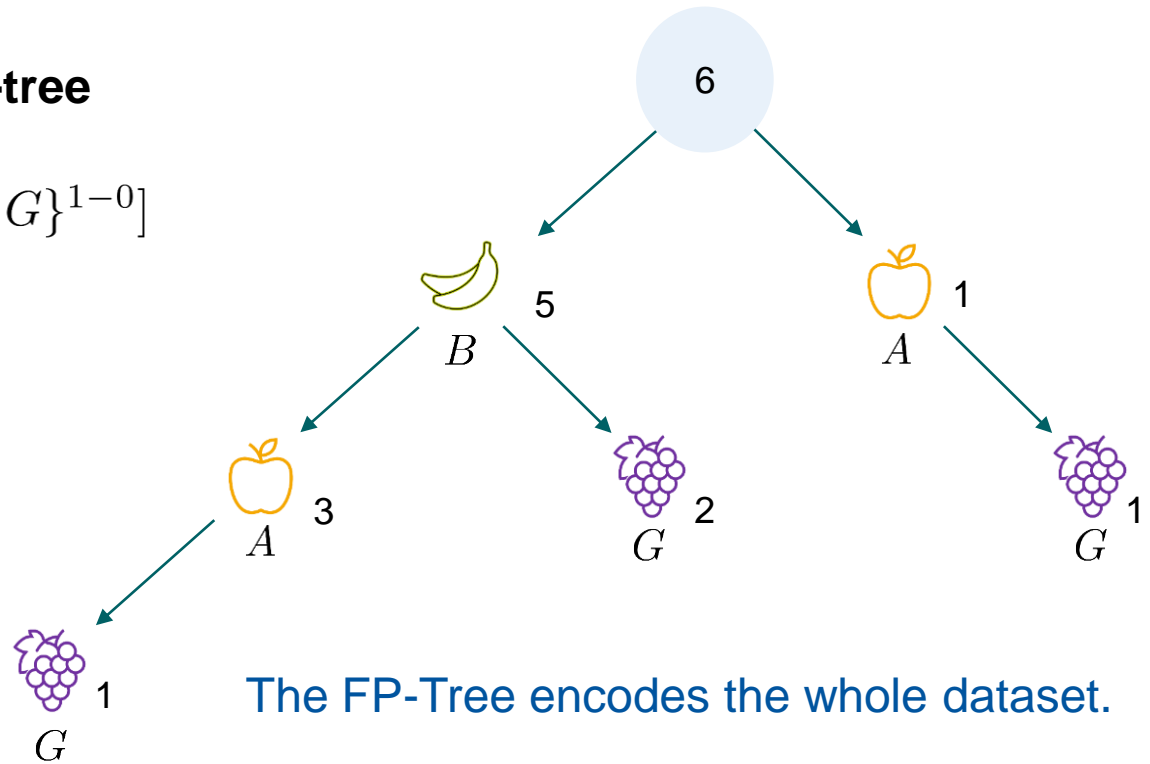| TID | Bought Fruits |
|-----|---------------|
| 1 | {Banana, Apple} |
| 2 | {Banana, Grapes} |
| 3 | {Banana, Apple} |
| 4 | {Apple, Grapes} |
| 5 | {Banana, Grapes} |
| 6 | {Banana, Apple, Grapes} |

5. Build the FP-tree going through each transaction



The FP-Tree encodes the whole dataset.

# FP-Tree – Encodes The Dataset
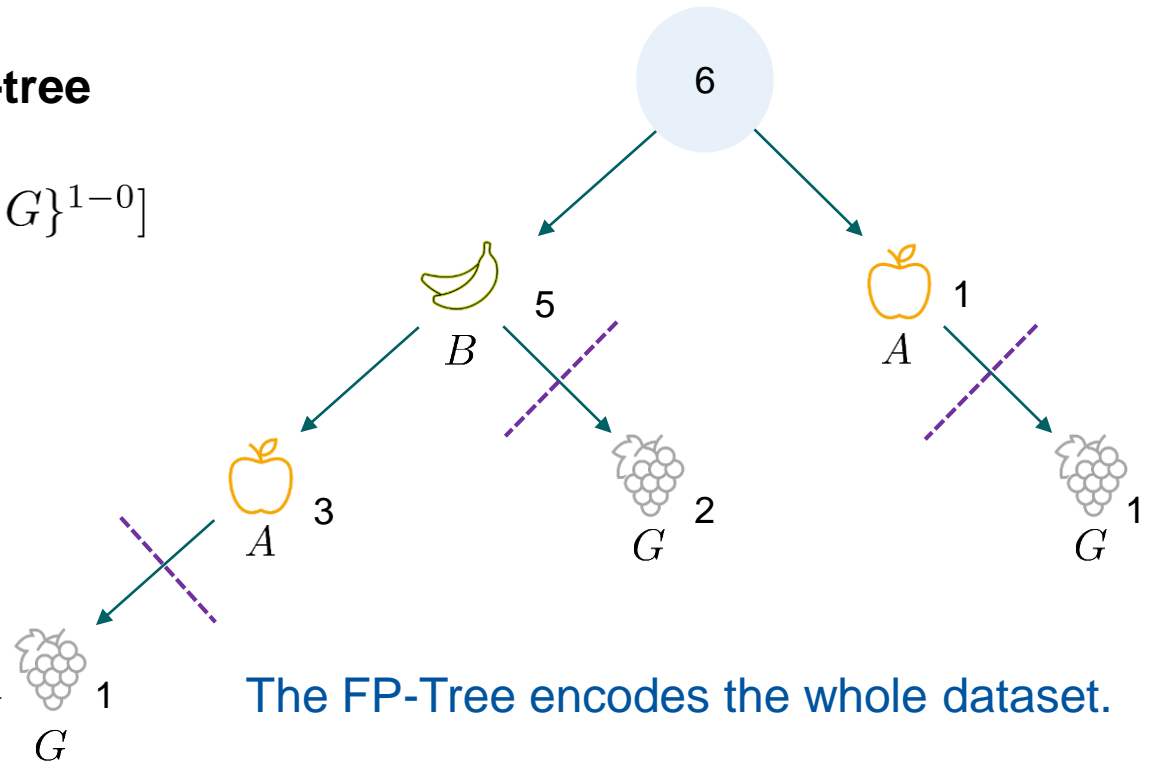
**We can read the transactions from the FP-tree**

$$\mathcal{X} = [\{B, A, G\}^{1-0}, \{B, A\}^{3-1}, \{B, G\}^{2-0}, \{A, G\}^{1-0}]$$
$$= [\{B, A, G\}, \{B, A\}^2, \{B, G\}^2, \{A, G\}]$$



The FP-Tree encodes the whole dataset.

# FP-Tree – Cannot Cut Naïvely

**We can read the transactions from the FP-tree**

$$\mathcal{X} = [\{B, A, G\}^{1-0}, \{B, A\}^{3-1}, \{B, G\}^{2-0}, \{A, G\}^{1-0}]$$
$$= [\{B, A, G\}, \{B, A\}^2, \{B, G\}^2, \{A, G\}]$$



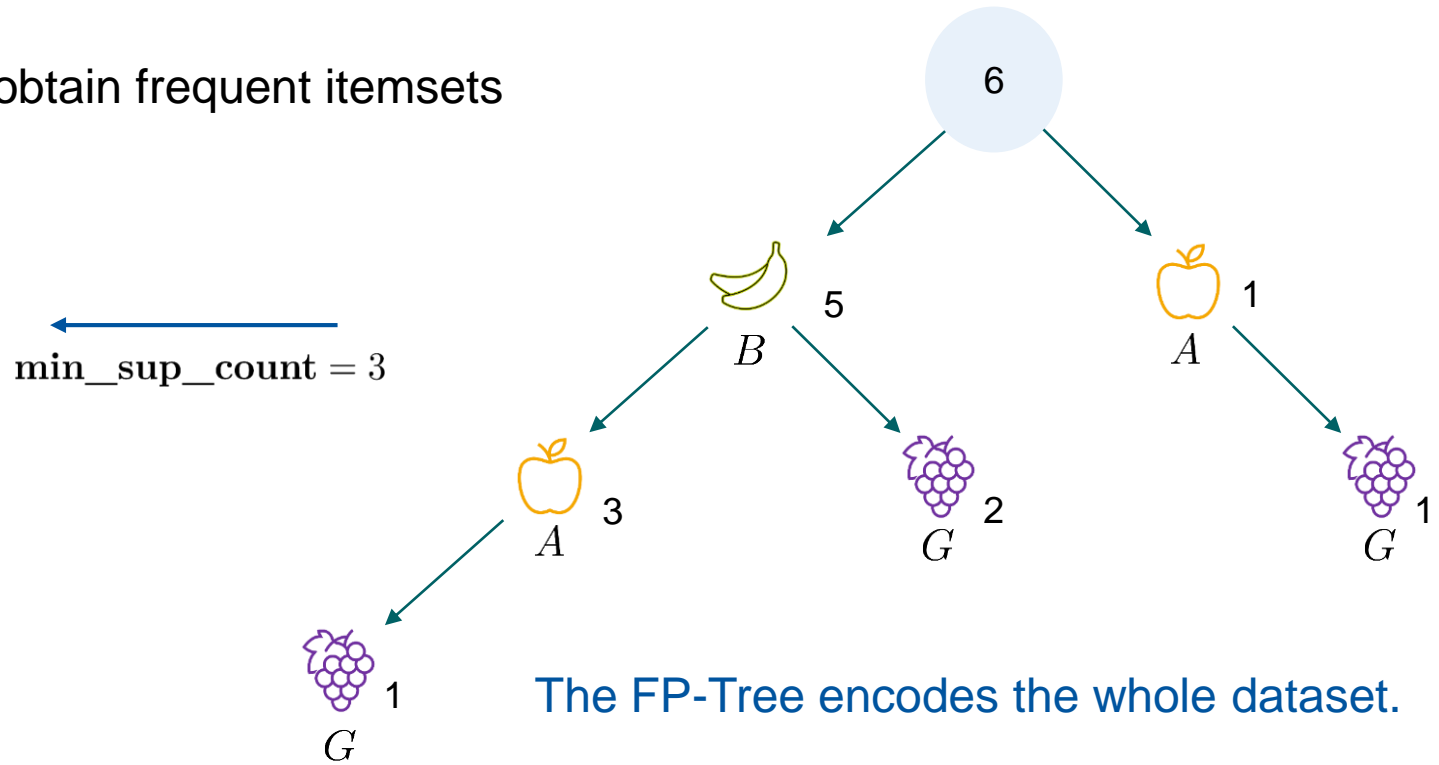The FP-Tree encodes the whole dataset.

Even though *G* exists only once in this subtree, we can't cut it naively, because its support may be higher than the threshold (3)
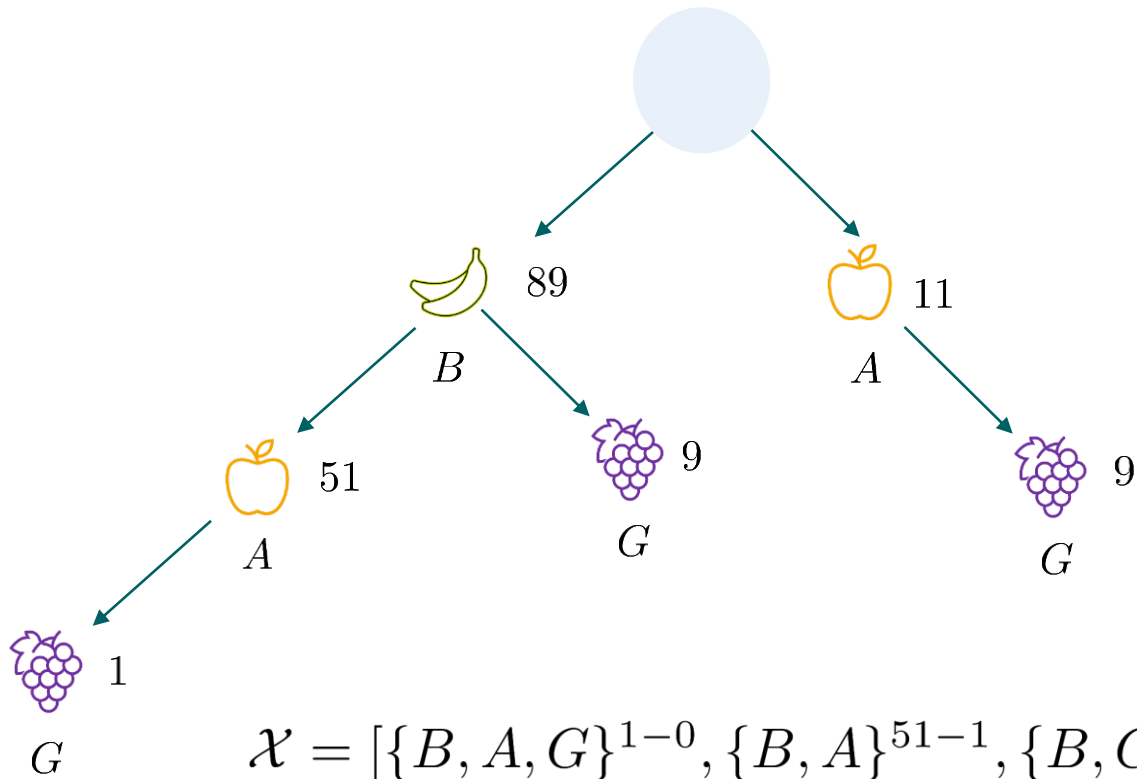
# FP-Tree – Frequent Itemsets

**Next:** 6. Mining the FP-tree to obtain frequent itemsets

| Frequent Itemsets | Support Count |
|---|---|
| {B} | 5 |
| {A} | 4 |
| {G} | 4 |
| {B, A} | 3 |
| {B, G} | 3 |

$\text{min\_sup\_count} = 3$



The FP-Tree encodes the whole dataset.

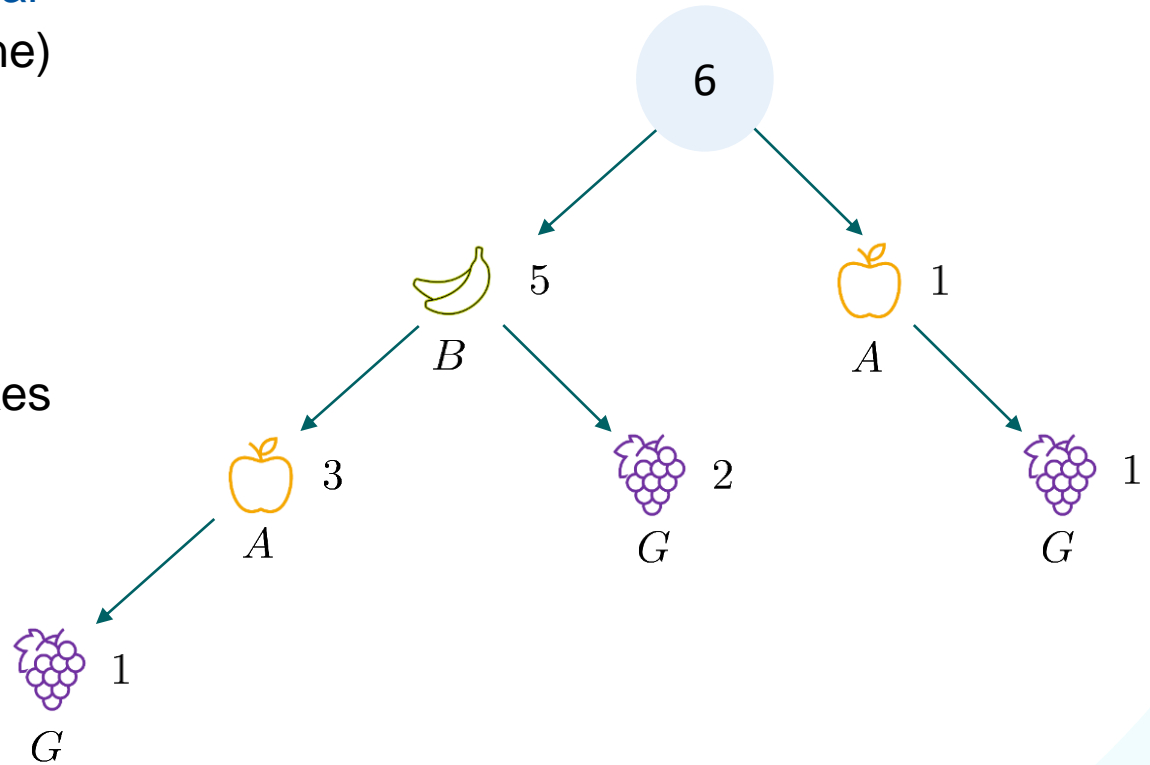## FP-Tree Encodes Dataset – Another Example



**We can read the transactions from the FP-tree**

$$\mathcal{X} = [\{B, A, G\}^{1-0}, \{B, A\}^{51-1}, \{B, G\}^{9-0}, \{B\}^{89-(51+9)}, \{A, G\}^{9-0}, \{A\}^{11-9}]$$
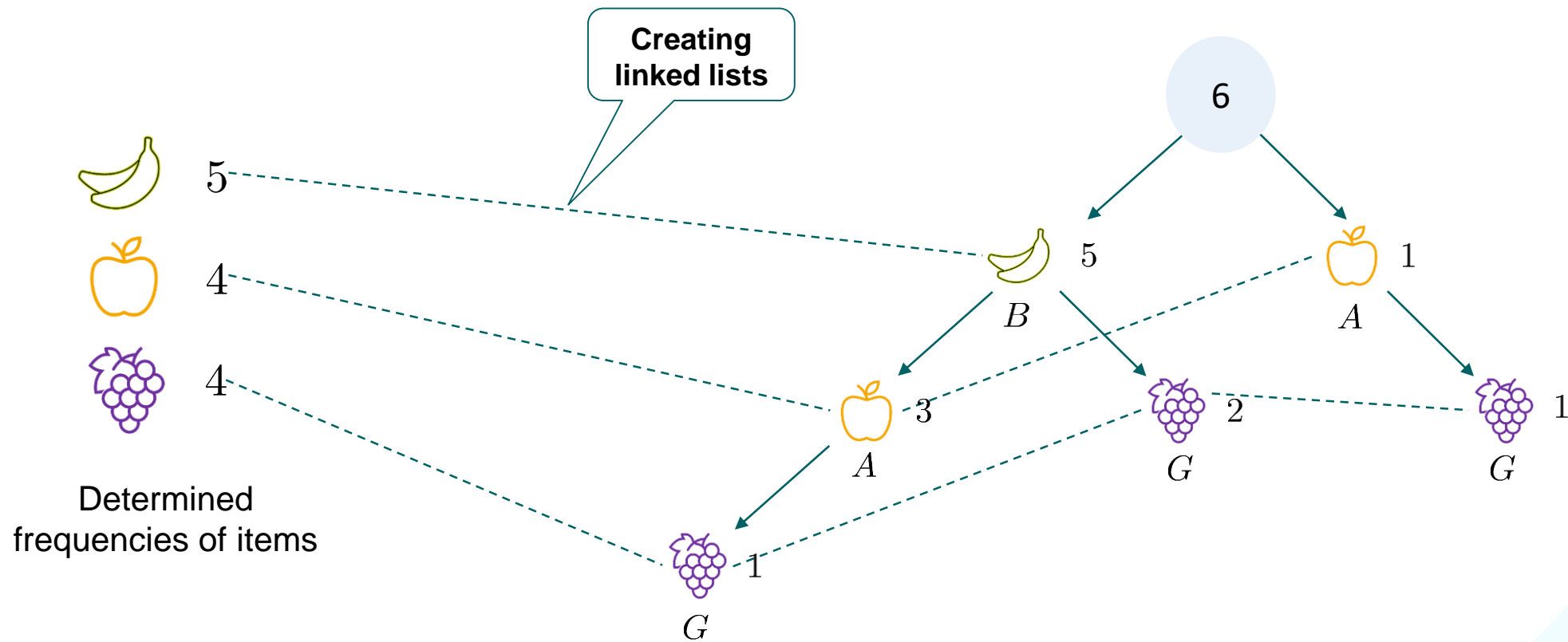
$$\mathcal{X} = [\{B, A, G\}^{1}, \{B, A\}^{50}, \{B, G\}^{9}, \{B\}^{29}, \{A, G\}^{9}, \{A\}^{2}]$$

# Mining the FP-Tree – Overview

- For each frequent item, create a conditional FP-tree (starting with the least frequent one)

- The conditional FP-tree considers all transactions ending with this item

- Apply this recursively

- Due to recursion, we also consider postfixes that contain multiple elements

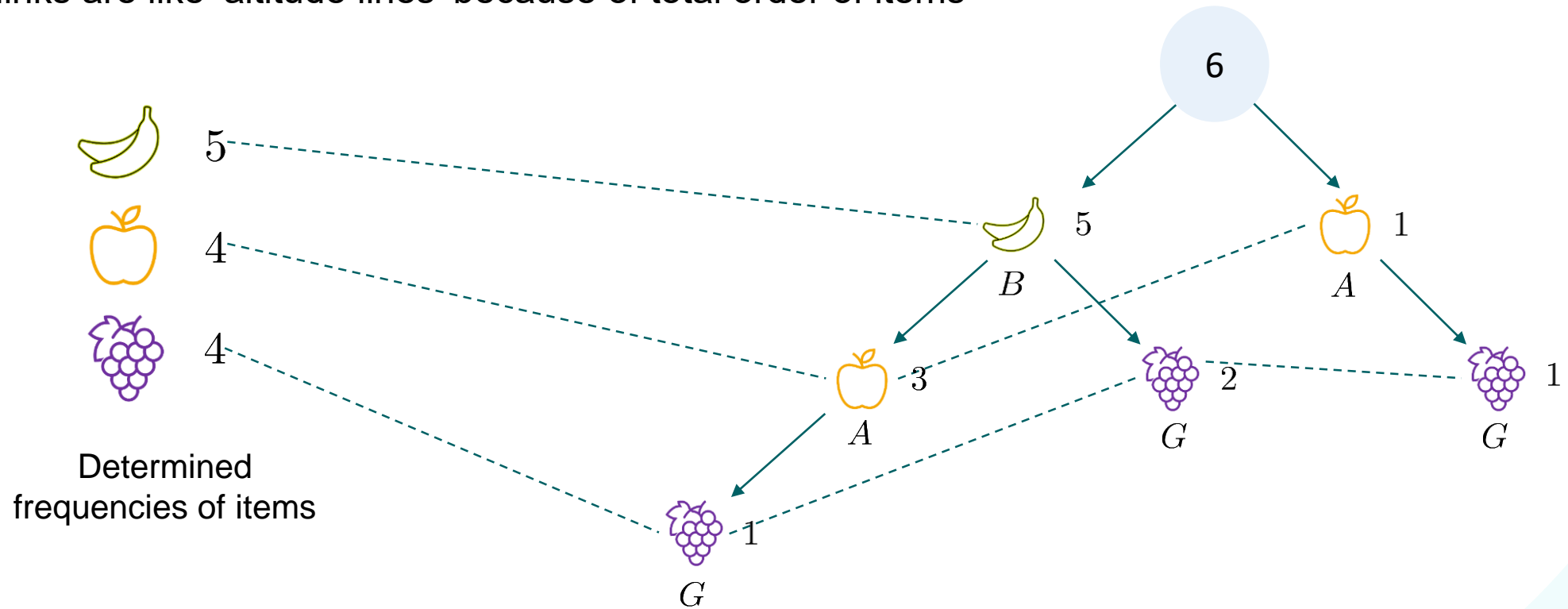- The ordering ensures that postfixes are considered only once
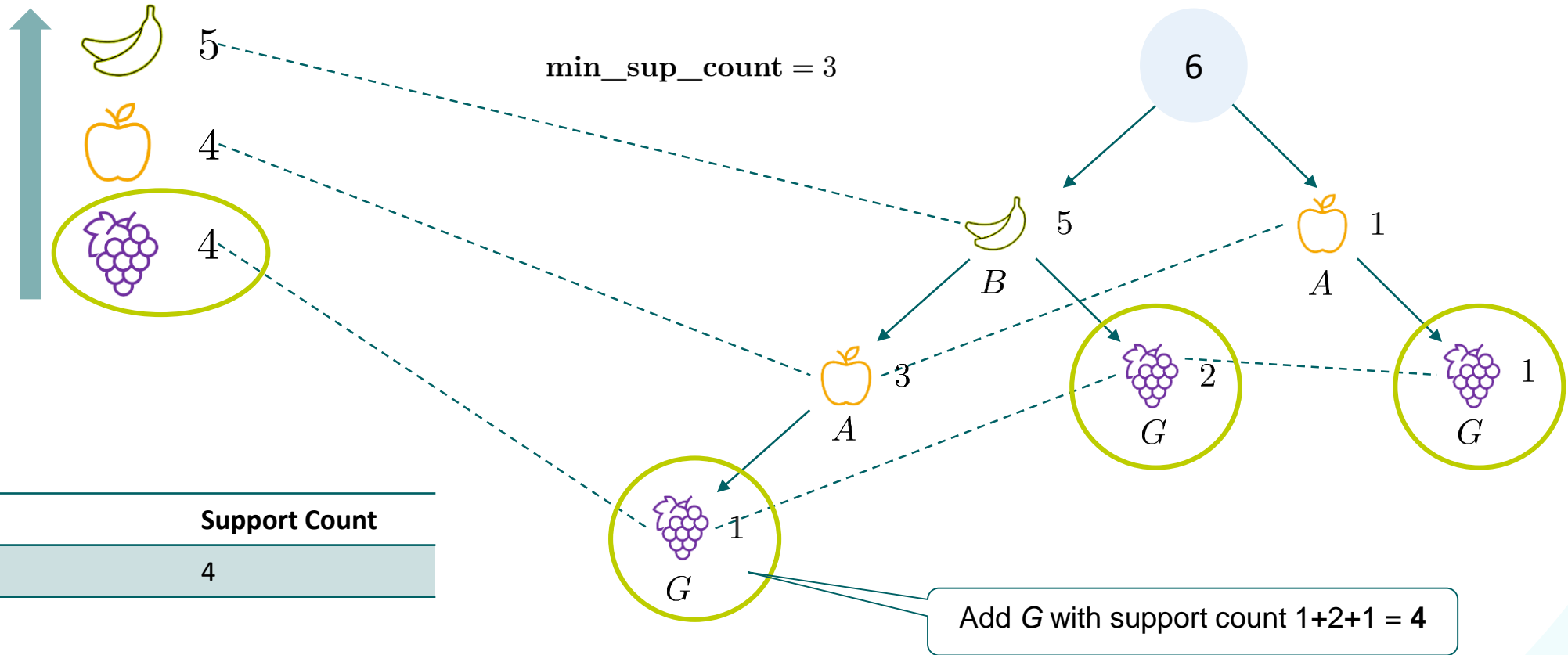
# Node Links

# Node Links

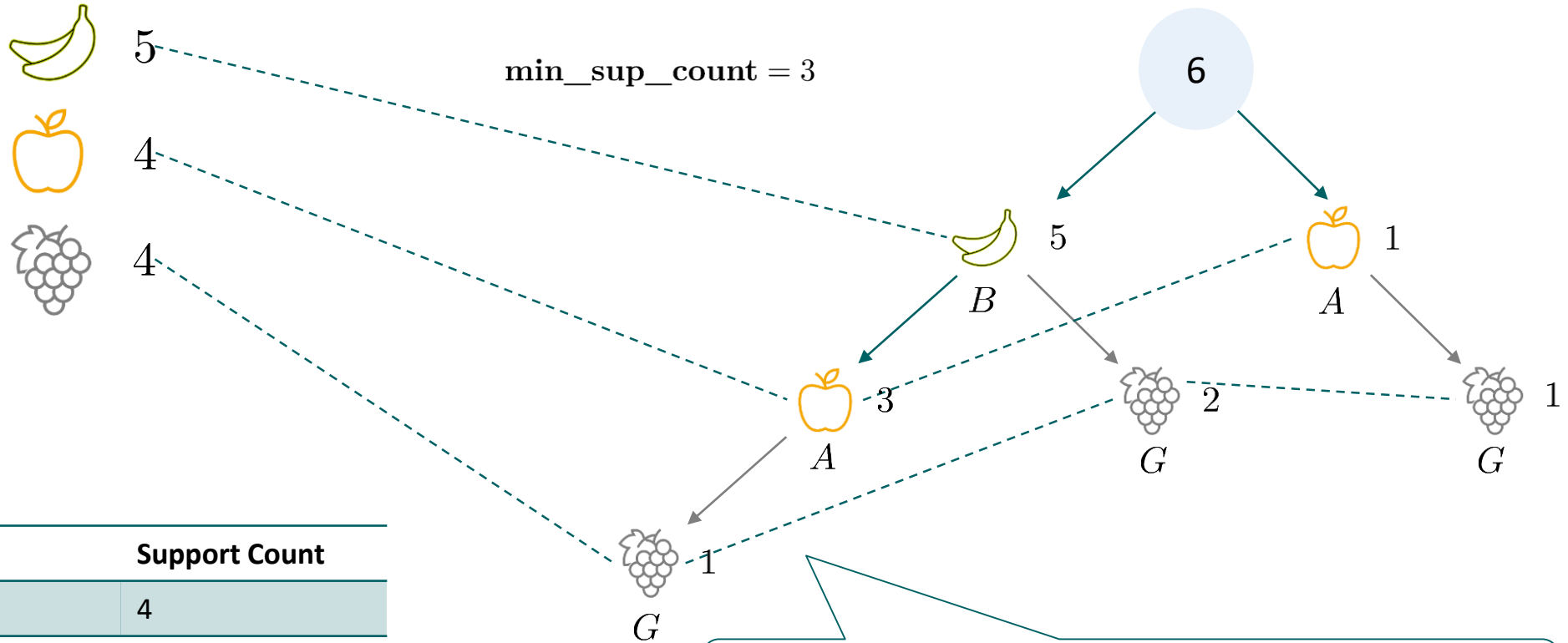Node links are like 'altitude lines' because of total order of items



Determined
frequencies of items

# Consider Postfix G



$\text{min\_sup\_count} = 3$

| Itemsets | Support Count |
|----------|---------------|
| {G} | 4 |

Add *G* with support count 1+2+1 = **4**

# Consider Postfix G

$min\_sup\_count = 3$

5

4

4

6

🍌 5
B

🍎 1
A

🍎 3
A

�grapes 2
G

�grapes 1
G

�grapes 1
G

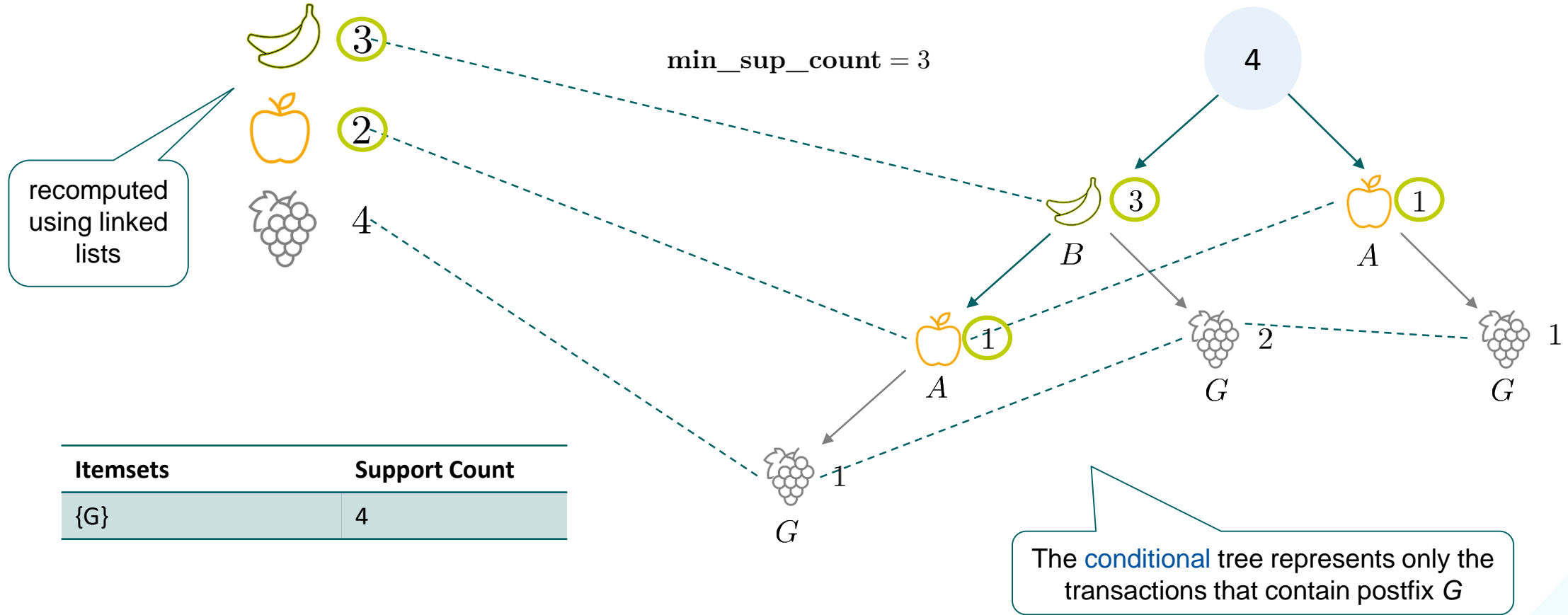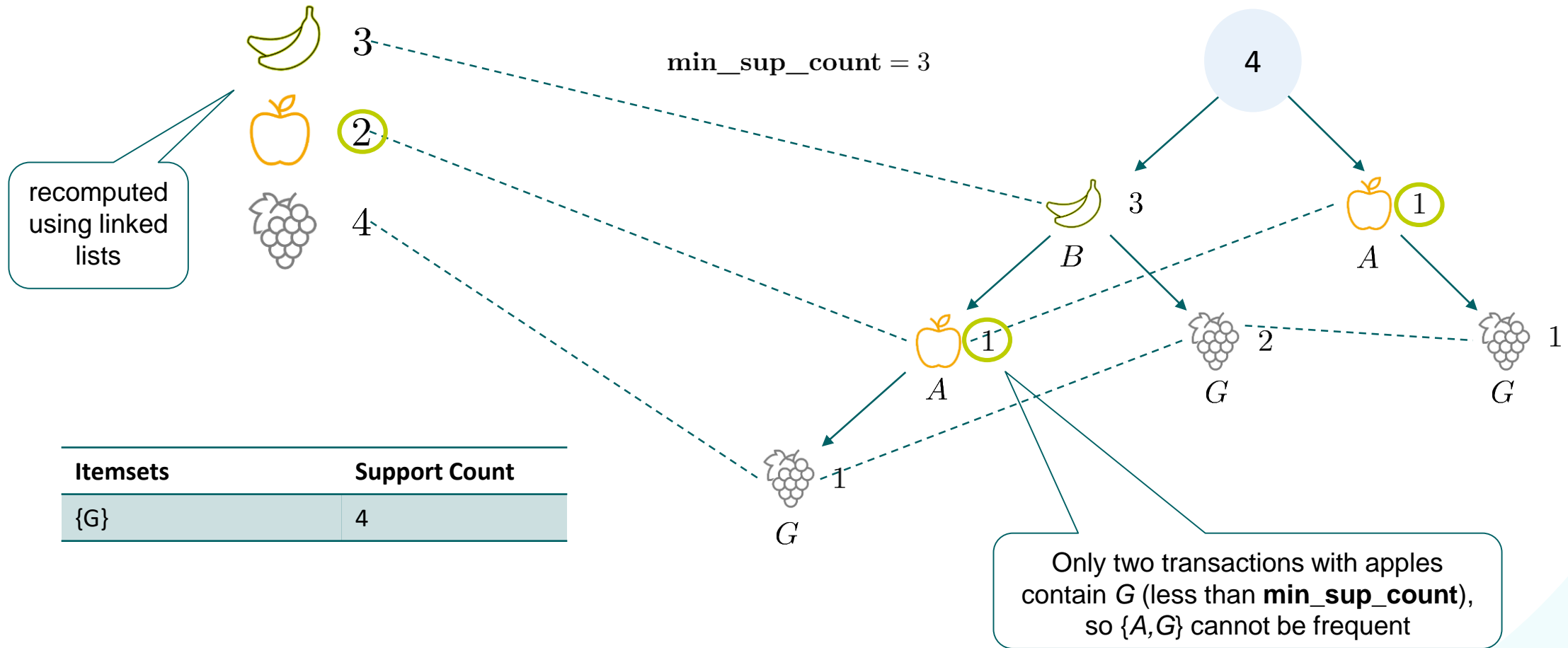| Itemsets | Support Count |
|----------|---------------|
| {G}      | 4             |

Now only consider **all the paths** leading from the root to *G*
(representing transactions that include *G*)

# Towards the Conditional FP-Tree for Postfix G
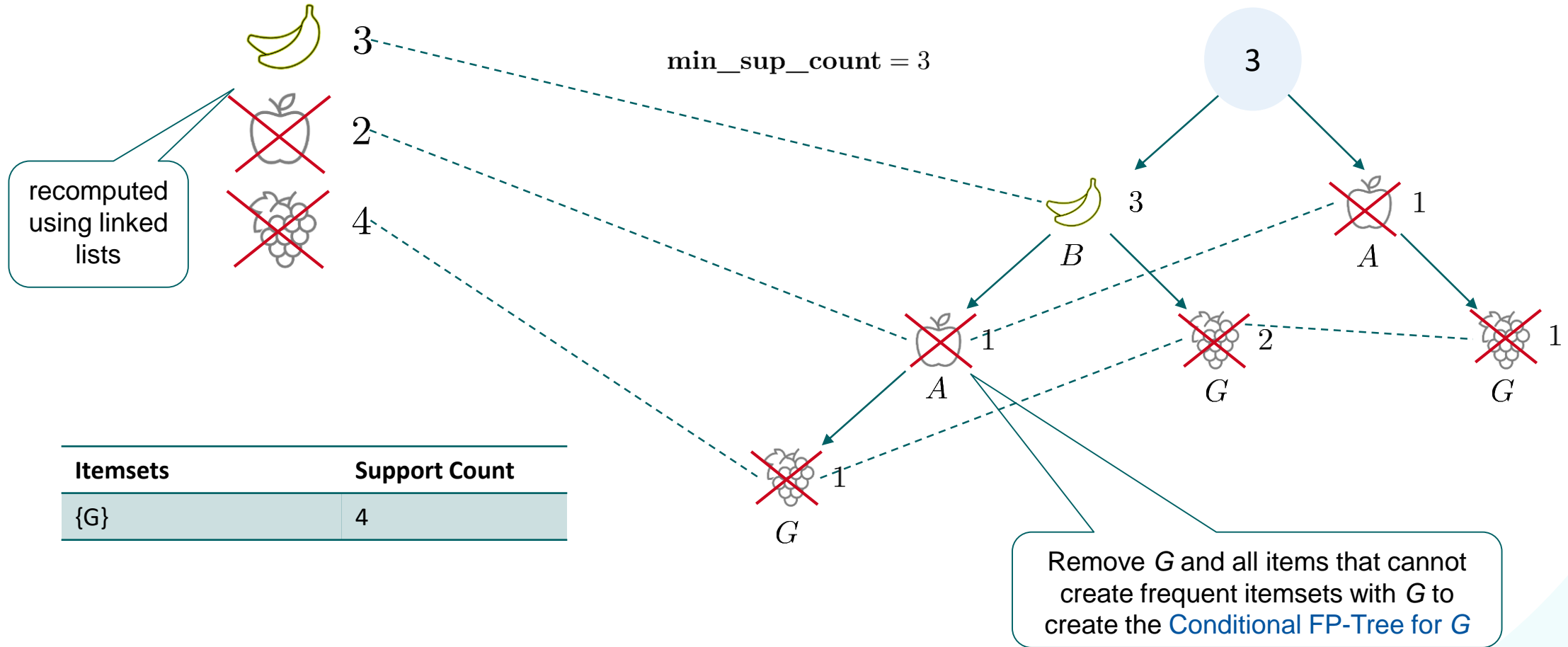


$\text{min\_sup\_count} = 3$

recomputed using linked lists

The conditional tree represents only the transactions that contain postfix *G*

| Itemsets | Support Count |
|----------|---------------|
| {G}      | 4             |

# Towards the Conditional FP-Tree for Postfix G



$min\_sup\_count = 3$

recomputed using linked lists

Only two transactions with apples contain $G$ (less than **min_sup_count**), so $\{A,G\}$ cannot be frequent

| Itemsets | Support Count |
|----------|---------------|
| {G}      | 4             |

# Towards the Conditional FP-Tree for Postfix G

# Conditional FP-Tree for Postfix G

$$\text{min\_sup\_count} = 3$$

3

3

3

$B$

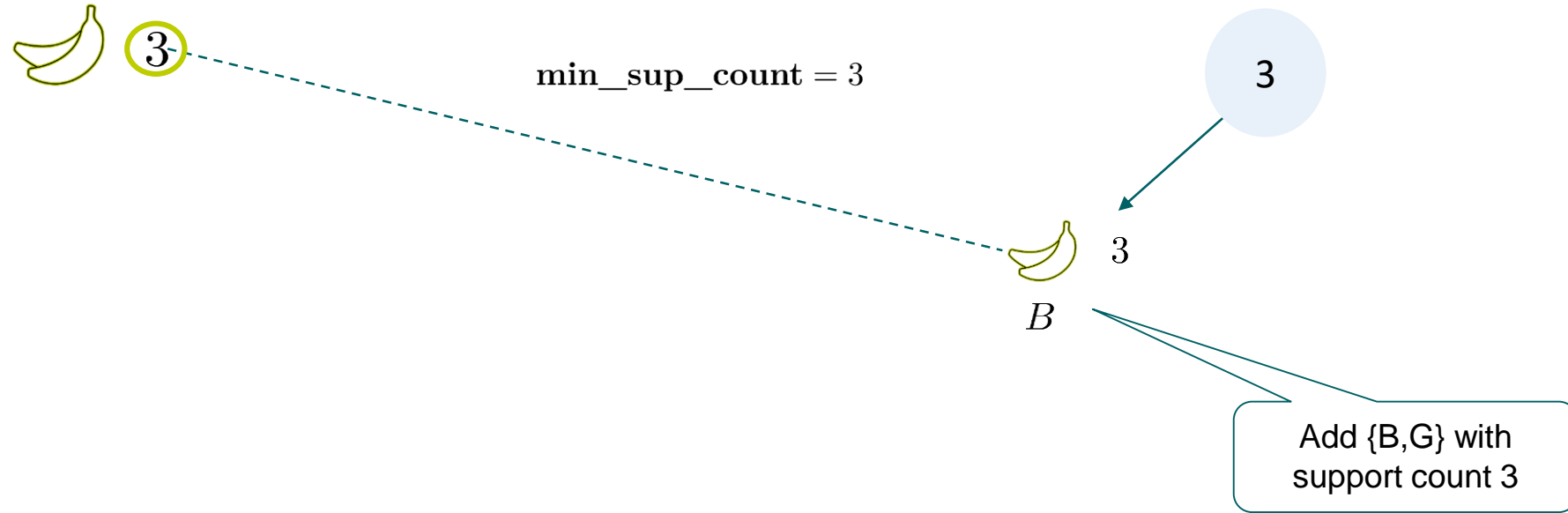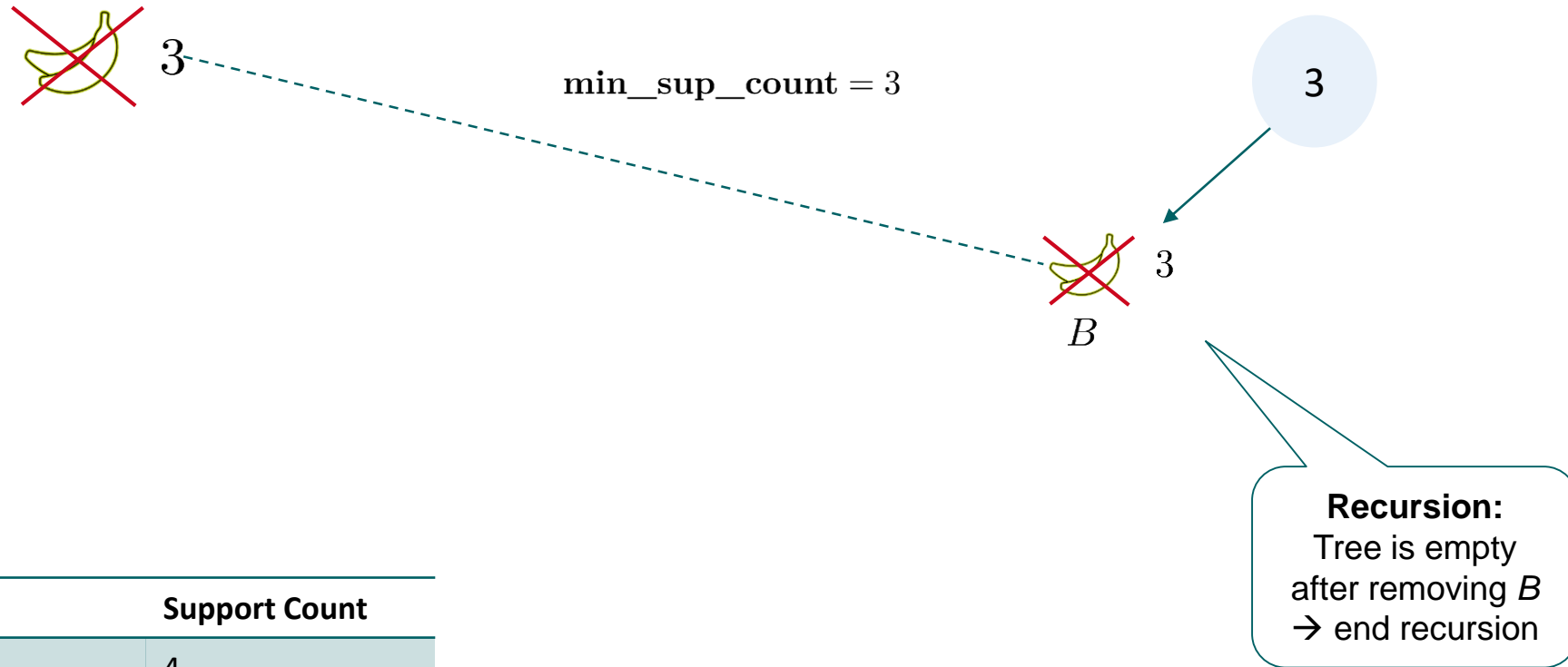| Itemsets | Support Count |
|----------|---------------|
| {G} | 4 |
| { ... , G } | |

recurse

Mine the Conditional FP-Tree for *G* to find frequent itemsets that contain *G* (Recursion)

# Conditional FP-Tree for Postfix G

③

$min\_sup\_count = 3$

3

3

$B$

Add {B,G} with
support count 3

| Itemsets | Support Count |
|----------|---------------|
| {G} | 4 |
| {B, G} | 3 |

# Conditional FP-Tree for Postfix G



$min\_sup\_count = 3$

**Recursion:**
Tree is empty after removing *B* → end recursion

| Itemsets | Support Count |
|----------|---------------|
| {G} | 4 |
| {B, G} | 3 |

# Consider Postfix A

$$\text{min\_sup\_count} = 3$$

Determined frequencies of items

| Itemsets | Support Count |
|----------|---------------|
| {G} | 4 |
| {B, G} | 3 |

We already added **all** the frequent itemsets with *G*

# Consider Postfix A

$$\text{min\_sup\_count} = 3$$

Determined
frequencies of items

| Itemsets | Support Count |
|----------|---------------|
| {G}      | 4             |
| {B, G}   | 3             |

A is the next frequent item

# Consider Postfix A

$\text{min\_sup\_count} = 3$

Determined
frequencies of items

5

4

4

6

5
B

1
A

3
A

2
G

1
G

1
G

Add *A* with support 3+1 = 4

| Itemsets | Support Count |
|----------|---------------|
| {G}, {A} | 4 |
| {B, G} | 3 |

# Towards the Conditional FP-Tree for Postfix A

# Towards Conditional FP-Tree for Postfix A



$$\text{min\_sup\_count} = 3$$

Determined
frequencies of items

| Itemsets | Support Count |
|----------|---------------|
| {G}, {A} | 4 |
| {B, G} | 3 |

# Conditional FP-Tree for Postfix A

$$\text{min\_sup\_count} = 3$$

3

$B$

**Recursion:**
Add {$B$, $A$} with support 3 (then end recursion)

Determined
frequencies of items

| Itemsets | Support Count |
|----------|---------------|
| {G}, {A} | 4 |
| {B, G}, {B, A} | 3 |

# Consider Postfix B



$\text{min\_sup\_count} = 3$

Determined frequencies of items

| Itemsets | Support Count |
|---|---|
| {G}, {A} | 4 |
| {B, G}, {B, A} | 3 |

we already added **all** the frequent itemsets with *G* and *A*

# Consider Postfix B



$min\_sup\_count = 3$

Determined frequencies of items

| Itemsets | Support Count |
|----------|---------------|
| {B} | 5 |
| {G}, {A} | 4 |
| {B, G}, {B, A} | 3 |

Add itemset {B} with support count 5

# Towards Conditional FP-Tree for Postfix B

$$\text{min\_sup\_count} = 3$$



5

B   5

5

Now only consider **all the paths** leading to B
(transactions including B)

Determined
frequencies of items

| Itemsets | Support Count |
|---|---|
| {B} | 5 |
| {G}, {A} | 4 |
| {B, G}, {B, A} | 3 |

# Conditional FP-Tree for Postfix B



5

4

4

0

**Conditional FP-tree is empty**, i.e., no additional frequent itemsets with *B*

| Itemsets | Support Count |
|---|---|
| {B} | 5 |
| {G}, {A} | 4 |
| {B, G}, {B, A} | 3 |

# Conditional FP-Tree for Postfix B

🍌 5

🍎 4

🍇 4

0

> **No more frequent items to consider**
> →algorithm terminates

| Itemsets | Support Count |
|---|---|
| {B} | 5 |
| {G}, {A} | 4 |
| {B, G}, {B, A} | 3 |

# All Frequent Itemsets Generated



FP-tree

| Itemsets | Support Count |
|---|---|
| {B} | 5 |
| {G}, {A} | 4 |
| {B, G}, {B, A} | 3 |

Frequent itemsets mined

# FP-Growth Algorithm – Summary

- Idea: frequent pattern growth based on FP-tree

- Method:
  - Construct the FP-tree from the dataset (previous video)
  - For each frequent item, construct its conditional pattern-base, and then its conditional FP-tree
  - Recursively repeat the process on each newly created conditional FP-tree until the tree is empty

# FP-Growth Algorithm – Summary

- Advantages:
  - ✓ Only two passes through the dataset are needed (when constructing the tree)
  - ✓ Avoiding testing many hopeless candidates
  - ✓ Very fast when FP-tree fits in main memory

- However: approach has problems when FP-tree is too large to fit into memory

# Frequent Itemsets – Summary

- Pattern mining is a form of unsupervised learning

- Frequent itemsets are the basis for finding patterns (ideas can be transferred to other patterns)

- Two well-known algorithms using generally applicable concepts:

  - Apriori algorithm
  - FP-growth algorithm

- Outlook

  - There may be many frequent "patterns"
  - How to determine which ones are surprising / interesting?

# Association Rules – Preview

(one of the topics of the next lecture)

- {Cheese, Bread} $\Rightarrow$ {Milk}

  People that buy Cheese and Bread also tend to buy Milk.

- {Track1, Track2} $\Rightarrow$ {Track3}

  Students that take the Track 1 and Track 2 modules of BridgingAI also tend to take the Track 3 courses. (We hope you do!)

- {Bitburger} $\Rightarrow$ {Heineken, Palm}

  People that buy Bitburger beer tend to buy both Heineken and Palm beer.

- {Carbonara, Margherita } $\Rightarrow$ {Espresso, Tiramisu}

  People that buy Carbonara and Margherita also tend to buy Espresso and Tiramisu.

- {part-245, part-345, part-456} $\Rightarrow$ {part-372}

  When Parts 245, 345, and 456 are replaced, then often also Part 372 is replaced.