

CSCI 4502/5502

Data Mining

Fall 2020
Lecture 08 (Sep 17)

Reminders

◆ Homework 2

◆ due at 9:30am, Th, Sep 17

◆ Homework 3

◆ posted in Canvas, due at 9:30am, Th, Sep 24

Review

♦ Chapter 6: Mining Frequent Patterns

♦ basic concepts

- ♦ frequent patterns, association rules: support, confidence

♦ Apriori algorithm

- ♦ Apriori pruning, itemsets: $k \implies k+1$

♦ interestingness measure: correlation: lift

Frequent Pattern Mining

♦ Challenges

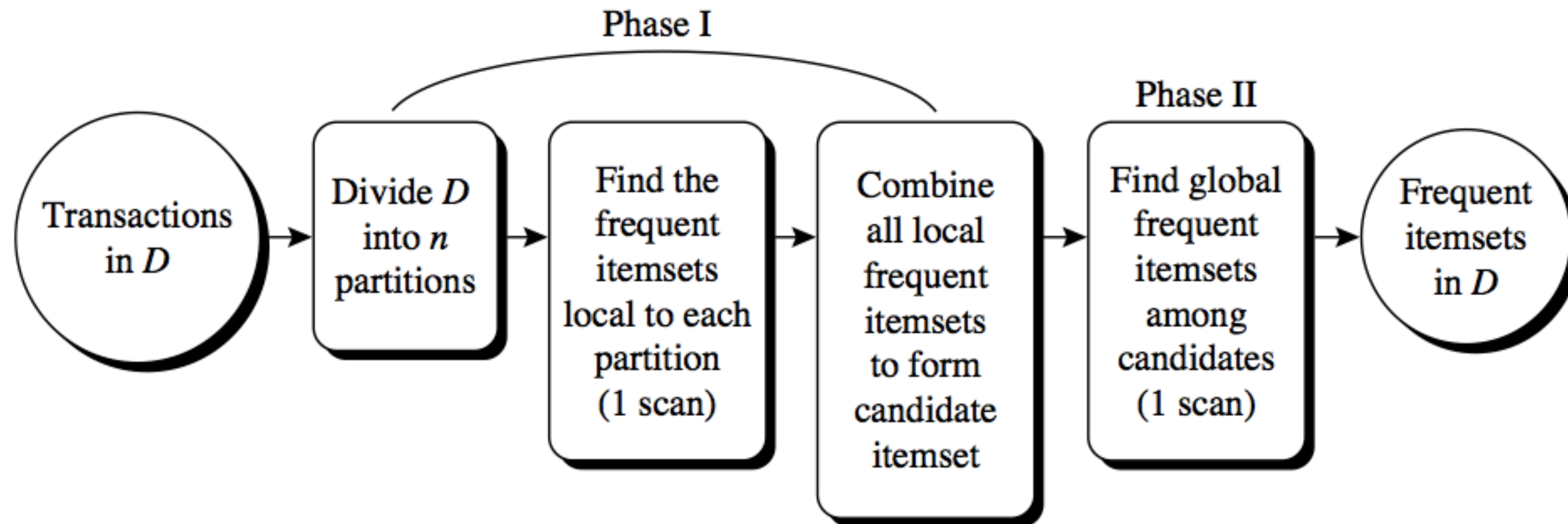
- ♦ multiple scans of the whole data set
- ♦ a huge number of candidates
- ♦ tedious support counting for candidates

♦ Improving Apriori: general ideas

- ♦ reduce data scans
- ♦ reduce number of candidates
- ♦ facilitate support counting of candidates

Partition: Two Data Scans

- ◆ A frequent itemset must be frequent in at least one partition
- ◆ Partition size? # of partitions?
 - ◆ each partition fits into main memory



Sampling for Freq. Patterns

- ◆ Select a sample data set
- ◆ Mine frequent patterns within sample
 - ◆ may use a lower min_sup
- ◆ Scan whole data set for actual support
 - ◆ only check closed patterns
 - ◆ e.g., check **abcd** instead of **ab**, **acd**, ..., etc.
- ◆ Scan again to find missed frequent patterns
- ◆ Sample size?

Transaction Reduction

- ◆ If a transaction T does not contain any frequent k -itemset
- ◆ then for any $h > k$, no need to check T when searching for frequent h -itemset
- ◆ Implementation
 - ◆ sequential scan
 - ◆ vs. random access

Reduce #Candidates

- ◆ Hash itemsets to buckets
- ◆ If a hash bucket count is below support threshold
- ◆ then itemsets in that hash bucket are not frequent itemsets

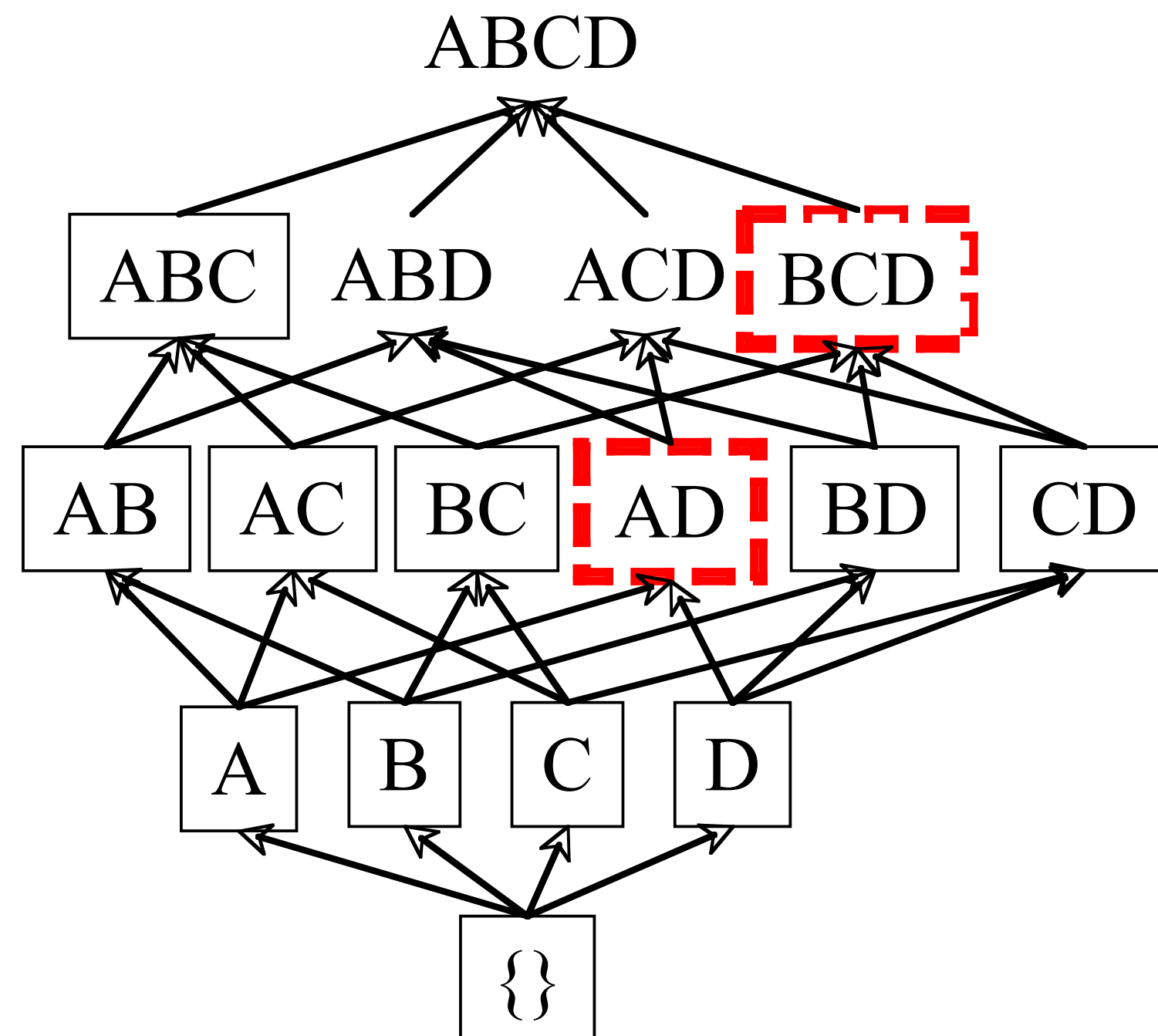
Create hash table H_2
using hash function
 $h(x, y) = ((\text{order of } x) \times 10$
 $+ (\text{order of } y)) \bmod 7$



H_2							
bucket address	0	1	2	3	4	5	6
bucket count	2	2	4	2	2	4	4
bucket contents	{I1, I4}	{I1, I5}	{I2, I3}	{I2, I4}	{I2, I5}	{I1, I2}	{I1, I3}
	{I3, I5}	{I1, I5}	{I2, I3}	{I2, I4}	{I2, I5}	{I1, I2}	{I1, I3}
			{I2, I3}			{I1, I2}	{I1, I3}
			{I2, I3}			{I1, I2}	{I1, I3}

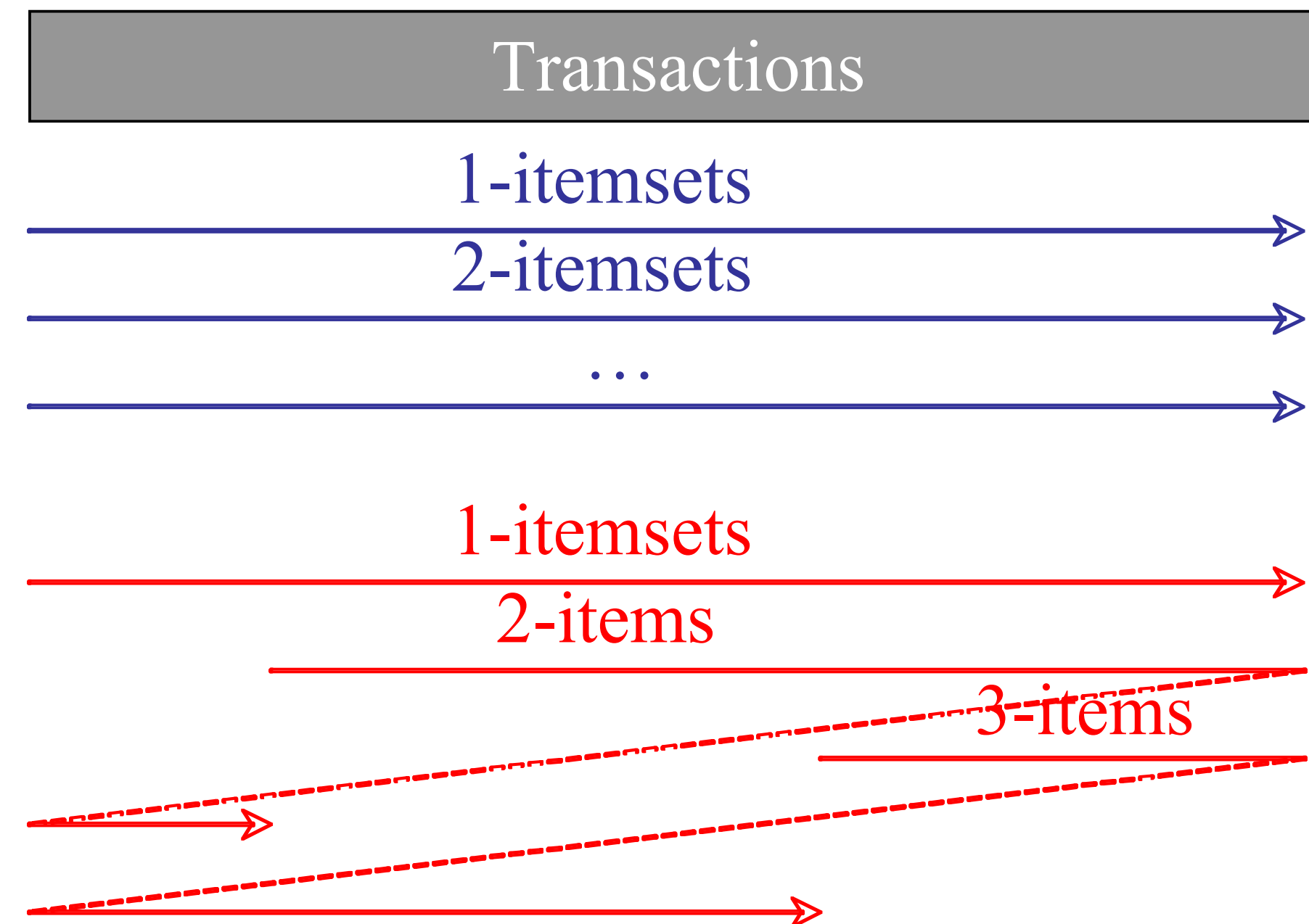
Dynamic Itemset Counting

- ♦ If A & D are freq., start count for AD
- ♦ If BC, BD, CD are freq., start count for BCD



Itemset lattice

Apriori

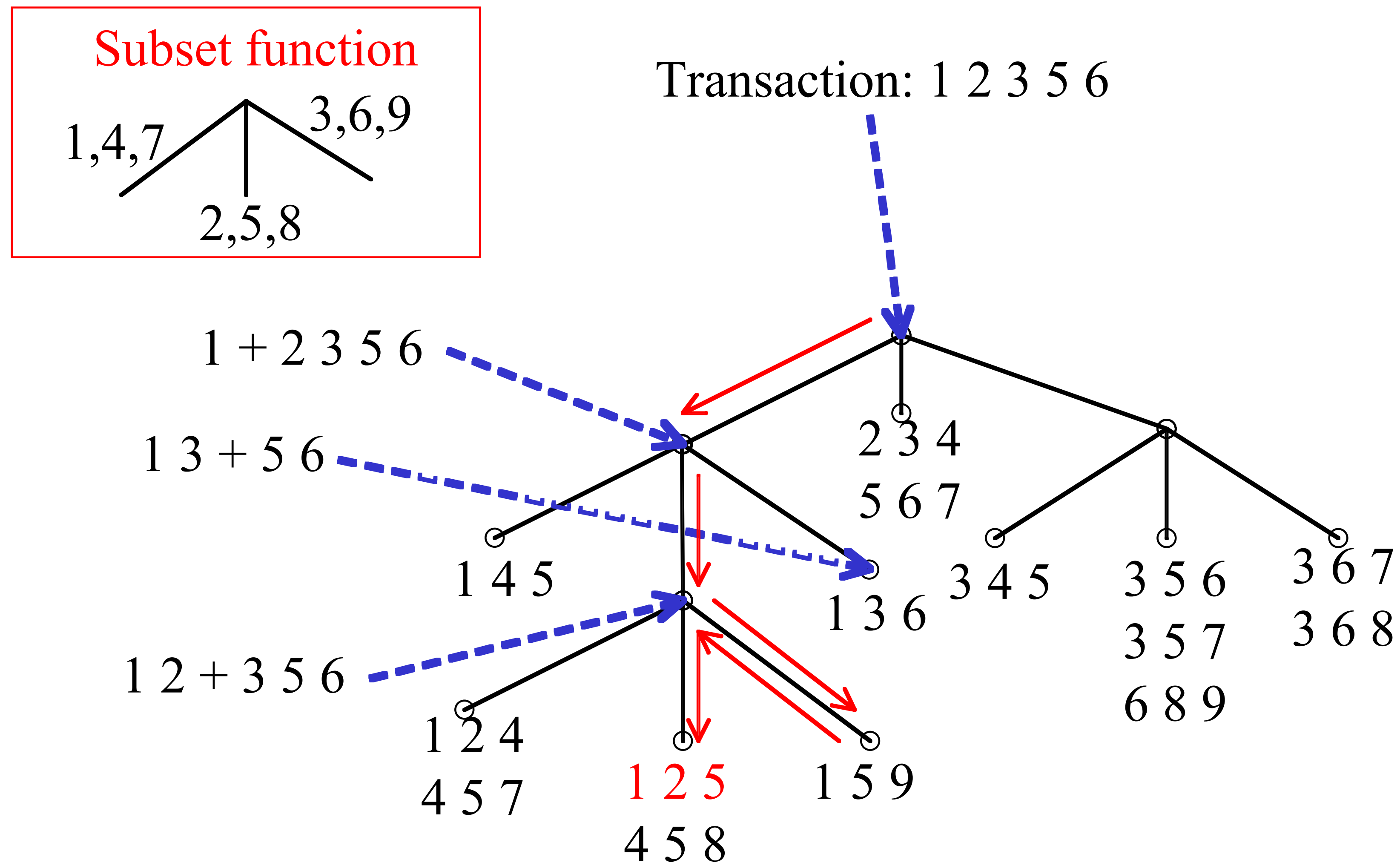


DIC

Count Support of Candidates

- ♦ Why counting candidate support a problem?
 - ♦ #candidates: total, per transaction
- ♦ Method
 - ♦ store candidate itemsets in a **hash-tree**
- ♦ **leaf-node** contains a list of itemsets and counts
- ♦ **interior node** contains a hash table
- ♦ **subset function**: finds all candidates contained in a transaction

Example



Vertical Data Format

- ◆ **Horizontal** data format

- ◆ $T1: \{A, D, E, F\}$

- ◆ **Vertical** data format

- ◆ $t(AD) = \{T1, T6, \dots\}$

- ◆ Derive closed pattern via **vertical intersection**

- ◆ $t(X) = \{T1, T2, T3\}$ and $t(Y) = \{T1, T3, T4\}$

- ◆ $t(XY) = \{T1, T3\}$

Frequent Itemset Mining

- ♦ Multiple **data scans** are costly
- ♦ Mining **long patterns** needs many scans and generates lots of candidates
 - ♦ e.g., 100 items: #scans, #candidates
- ♦ **Bottleneck**: candidate generation & test
- ♦ Can we avoid candidate generation?

FP-growth (I)

- ♦ Find frequent itemsets without candidate generation
- ♦ Grow long patterns from short ones using local frequent items
- ♦ Example
 - ♦ abc is a frequent itemset; get all transactions with abc : $DB \mid abc$
 - ♦ d is a local frequent item in $DB \mid abc$
 - ♦ then $abcd$ is a frequent itemset

FP-tree Construction

<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>
100	{ <i>f, a, c, d, g, i, m, p</i> }	{ <i>f, c, a, m, p</i> }
200	{ <i>a, b, c, f, l, m, o</i> }	{ <i>f, c, a, b, m</i> }
300	{ <i>b, f, h, j, o, w</i> }	{ <i>f, b</i> }
400	{ <i>b, c, k, s, p</i> }	{ <i>c, b, p</i> }
500	{ <i>a, f, c, e, l, p, m, n</i> }	{ <i>f, c, a, m, p</i> }

min_sup = 0.6

- ♦ Scan, find freq. 1-itemset
- ♦ Sort freq. items in descending frequency
- ♦ Scan, construct FP-tree

