

ASSOCIATION MINING

Definition

- Search for patterns recurring in the given data set
- Given a set of item sets or transactions, find rules predicting the occurrence of items based on the occurrences of other items in the transactions

Applications

- Netflix movie recommendations
{Breaking Bad, House of Cards}



-> Mad Men



- Google search
 - Suggested autofill
 - Google.com, search “how to”, top 3 suggestions
- Order of webpages / Ads
 - Search history
 - Location

Search history: “jobs”

Microsoft Edge

The screenshot shows the Microsoft Edge browser with a Google search for "jobs". The search results page displays "About 6,180,000,000 results (0.62 seconds)". The top result is "Jobs | Indeed.com" with a link to "www.indeed.com/jobs". Below this, there are two columns: "Find Jobs" and "Find Resumes". A blue box highlights a result for "Monster.com® - Official Site | Use Monster find jobs near you" with a link to "www.monster.com/v". Below this, there is a blue header for "Jobs" with the location "Near East Lansing, MI". The search filters include "Computer & IT", "Past 3 days", "Full-time", "Sales & Retail", "Education", and "Science & Engineering". The first job listing is for "Computer Science Faculty" at "Michigan State University" in "East Lansing, MI" via "IEEE Job Site", with a "Full-time" tag.

Mozilla Firefox

The screenshot shows the Mozilla Firefox browser with a Google search for "jobs". The search results page displays "About 6,180,000,000 results (0.68 seconds)". The top result is "Jobs | Indeed.com" with a link to "www.indeed.com/jobs". Below this, there are two columns: "Find Jobs" and "Find Resumes". Below this, there is a blue header for "Jobs" with the location "Near East Lansing, MI". The search filters include "Sales & Retail", "Past 3 days", "Full-time", "Computer & IT", "Education", and "Management". The first job listing is for "Computer Science Faculty" at "Michigan State University" in "East Lansing, MI" via "IEEE Job Site".

Location: "bounty"

United States

google bounty

All Shopping Images News Videos More Settings Tools

About 89,300,000 results (0.53 seconds)

Dictionary

Enter a word, e.g. "pie"

boun·ty
/'boun(t)jē/ ⓘ

noun

1. a sum paid for killing or capturing a person or animal.
"there was an increased bounty on his head"
2. **HISTORICAL**
a sum paid to encourage trade.
"bounties were paid to colonial producers of indigo dye"

Translations, word origin, and more definitions

Feedback

Bounty Paper Towels and Napkins | Home
<https://bountytowels.com/en-us>

Bounty paper towels and napkins products clean up the smallest spills and the biggest messes. Start cleaning with the absorbent quicker picker upper today!
[About Bounty](#) [Bounty Essentials](#) [Bounty FAQs](#) [Bounty History](#)

Bounty | Definition of Bounty by Merriam-Webster
<https://www.merriam-webster.com/dictionary/bounty>

Bounty definition is - something that is given generously. How to use bounty in a sentence.

See bounty Sponsored

Bounty Select-A-Size White Paper Towels 12 Mega Rolls
\$15.89
Google Express
Free shipping

Bounty Paper Towels, White, 12 Super Rolls = 22 Regular Rolls
\$17.47
Walmart
Store pickup
★★★★★ (1k+)

→ More on Google

Bounty
Brand

Bounty is an American paper towel product manufactured by Procter & Gamble in the United States. It was introduced in 1965. Wikipedia

Introduced: 1965
Product type: Paper Towels

United Kingdom

google bounty

All Images Shopping News Videos More Settings Tools

About 86,400,000 results (0.34 seconds)

A privacy reminder from Google

REMIND ME LATER REVIEW

Bounty | Pregnancy & Parenthood Advice
www.bounty.com/

Bounty offer support on getting pregnant, pregnancy and parenthood. Free Bounty packs are available to every member with exclusive offers and samples.

Bounty Portrait
Bounty Portrait - The first photos of your newest family member.

Join Bounty.com | Your details
Register. Join for free stuff, offers and info tailored to you, your ...

Bounty packs
Find out how to get hold of one of the world famous Bounty Packs ...

Pregnancy and birth
... pregnancy & birth including week by week developments ...

Show all products
Browse Bounty's fantastic product offers for pregnancy, baby ...

Bounty Mum-to-be Pack
Why you don't want to miss your Bounty mum-to-be pack ...

More results from bounty.com »

Bounty Portrait - capture the moment
<https://www.bountyweb.co.uk/>

Wall Decor View Range. Bounty Boutique Become a photographer! Do you have a soft spot for babies, and a keen eye for cute portraits? Join our amazing team, ...

Bounty (UK) Limited
Company

bounty.com

Founded: 1959
Parent organization: Romeo Midco Ltd.

Disclaimer
Claim this knowledge panel

Feedback

See results about

Bounty (Chocolate bar)
Bounty is a chocolate bar manufactured by Mars. Incorporated and sold internationally. It was ...

Bounty (Company)
Bounty is a promotions company, pregnancy and parenting club. The pregnancy club gives advice in ...

Applications (cont'd)

- Market basket analysis: what items do customers buy together

{Bread, Milk} => {Paper Towel}

- Recommender System: A sales manager at an electronic store talking to a customer who recently purchased a computer and a camera, what should he recommend next?

{Video camera} => {warranty, memory card}

- Customer relationship management: identify preferences of different customer groups {Home, 2 cars} => {Policy A}, {Home, Ann Arbor} => {Policy B}

- Medical Diagnosis: find associations in symptoms and observations to predict diagnosis {Fever, lethargic, vomiting} => {Food Poisoning}

Market Basket Analysis

Market-Basket 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

Each row in this table corresponds to a transaction, which contains a unique identifier and a set of items bought by a given customer. Retailers are interested in analyzing the data to learn about the purchasing behavior of their customers.

Example of Association Rules

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

**Implication means co-occurrence,
not causality!**

Issues

- Discovering patterns from large transaction data:
computationally extensive
- Discovery of fake patterns

Definitions

- Itemset: $I = \{i_1, i_2, \dots, i_d\}$
 - The set of all items in the market basket data

$I = \{Bread, Milk, Diaper, Beer, Eggs, Coke\}$

- K-itemset: an itemset containing k items
- Null/empty itemset: an itemset that does not contain any items
- Transaction set: $T = \{t_1, t_2, \dots, t_N\}$
- Each transaction t_j contains a subset of I

<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

Definitions

- **Support** of an itemset X: number of transactions containing X

- $\sigma(\{\text{Bread, Diapers}\})$

3

- $\sigma(\{\text{Diapers, Milk, Coke}\})$

2

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diapers, Beer, Eggs
3	Milk, Diapers, Beer, Coke
4	Bread, Milk, Diapers, Beer
5	Bread, Milk, Diapers, Coke

Definitions

- **Association Rule:**

- Implication of the form $X \rightarrow Y$, where X and Y are *disjoint* itemsets
- $\{\text{Bread, Diapers}\} \rightarrow \{\text{Milk}\}$
- $\{\text{Bread}\} \rightarrow \{\text{Milk}\}$

- **Support of a rule:**

- The fraction of transactions containing both X and Y

$$s(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N}$$

- **Confidence of a rule:**

- The fraction transactions containing X that also contains Y

$$c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$$

Example

Association Rule: $\{Bread, Diaper\} \rightarrow \{Milk\}$

$$s(\{Bread, Diaper\} \rightarrow \{Milk\}) = \frac{\sigma(\{Bread, Diaper, Milk\})}{N}$$
$$= \frac{2}{5}$$

$$c(\{Bread, Diaper\} \rightarrow \{Milk\}) = \frac{\sigma(\{Bread, Diaper, Milk\})}{\sigma(\{Bread, Diaper\})}$$
$$= \frac{2}{3}$$

TID	Items
1	Bread, Milk
2	Bread, Diapers, Beer, Eggs
3	Milk, Diapers, Beer, Coke
4	Bread, Milk, Diapers, Beer
5	Bread, Milk, Diapers, Coke

Interpretation

- Support:
 - Low: items may occur together by chance
 - Used to eliminate uninteresting rules
 - Transaction set contains 1000 transaction.
 - A single transaction contains the items {Band aids, TV}
 - No other transactions contain either item
 - What is the support for {TV} => {Band aids}?
 - What is the confidence for {TV} => {Band aids}?
- $\sigma(\{\text{TV}, \text{Band aids}\}) = \# \text{ trans. containing both} = 1$
- $s(\{\text{TV}\} \Rightarrow \{\text{Band aids}\}) = \sigma(\{\text{TV}, \text{Band aids}\}) / N = 1/1000 = 0.001$
- $c(\{\text{TV}\} \Rightarrow \{\text{Band aids}\}) = \sigma(\{\text{TV}, \text{Band aids}\}) / \sigma(\{\text{TV}\}) = 1/1 = 1$

Interpretation

- Confidence:
 - Measures reliability of implication
 - The higher the confidence, the more likely Y is present in transactions containing X
 - Transaction set contains 1000 transaction.
 - 200 transactions contain the items {Milk, Paper}
 - 250 transactions contain {Milk}
 - 800 transactions contain {Paper}
 - What is the support for {Milk} => {Paper}?
 - What is the confidence for {Milk} => {Paper}?

**Co-occurrence /
Not causality
relationship**

$\sigma(\{\text{Milk, Paper}\}) = \# \text{ trans. containing both} = 200$

$s(\{\text{Milk}\} \Rightarrow \{\text{Paper}\}) = \sigma(\{\text{Milk, Paper}\}) / N = 200/1000 = 0.2$

$c(\{\text{Milk}\} \Rightarrow \{\text{Paper}\}) = \sigma(\{\text{Milk, Paper}\}) / \sigma(\{\text{Milk}\}) = 200/250 = 0.8$

Association Rule Discovery Problem

- Given:
 - a set of transactions T
 - a minimum support $minsup$
 - A minimum confidence $minconf$
- Find all association rules having:
 - support $\geq minsup$
 - confidence $\geq minconf$

Association Rule Discovery

- Brute Force Approach: find all possible rules then filter.
- 2-Step Approach: find frequent items then generate rules.

Brute Force Approach

- Compute support and confidence of every possible rule
- Select only rules satisfying *minsup* and *minconf* threshold

- Possible number of rules:

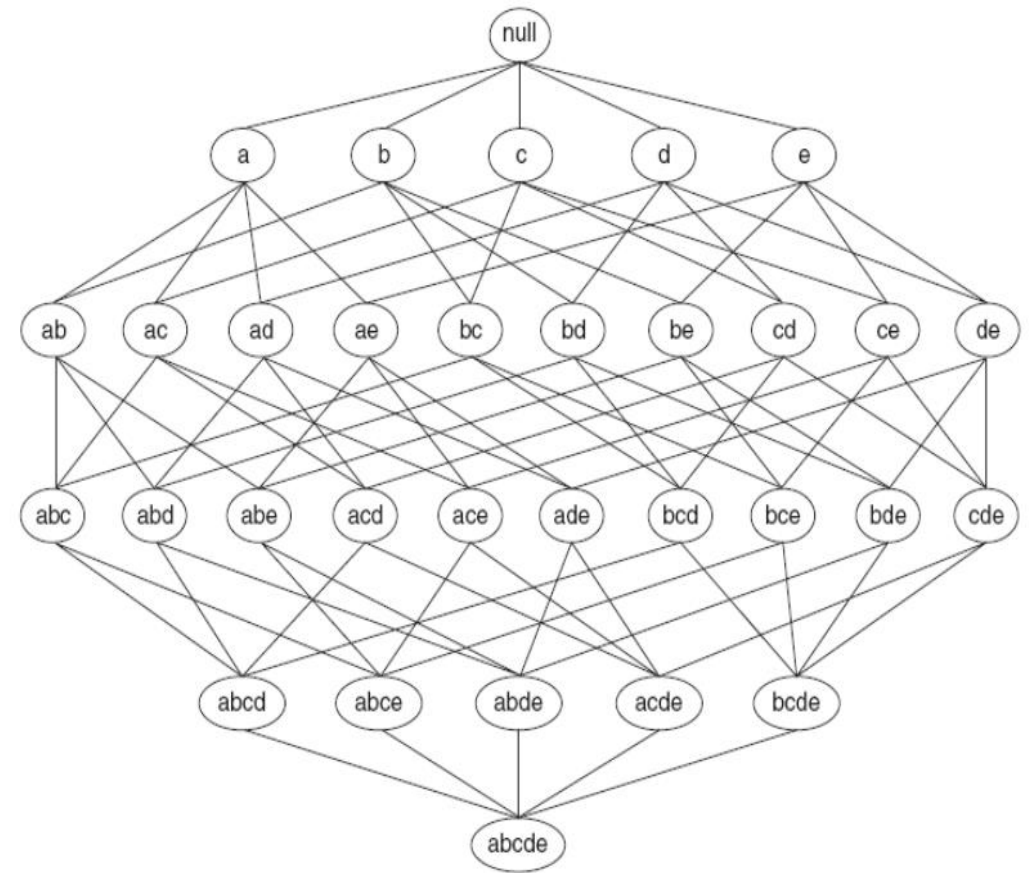
$$R = 3^d - 2^{d+1} + 1 \quad (d: \text{items set size})$$

Exponential $O(3^d)$

Prohibitively expensive!

- Example:

- $d = 6$ $R = 3^6 - 2^7 + 1 = 602$
- $d = 10$ $R = 3^{10} - 2^{11} + 1 = 57,002$
- $d = 15$ $R = 14,283,372$



If a store has 1000 different items:

```
> R:= 3^d - 2^(d+1) + 1;  
R := 132207081948080663689045525975214436596542203275214816766492036822682\  
859734670489954077831385060806196390977769687258235595095458210061891186\  
534272525795367402762022519832080385658460208523949485542189913638858182\  
990085158151242228918655840564706565199178372826151276809162989184514543\  
746640243175117323050698398769180942714083126679043875234727009443324274\  
972748463095841593799563260256679785505805331152530775858354670997698520\  
7992693374081496689754702826924329094491519081250
```

If 1000 rules can be generated per second

```
> R/1000 ./3600/24/365;  
 .4192259068 10467
```

Observation

- Support of a rule $X \rightarrow Y$ depends on support of itemset $\{X \cup Y\}$

$$s(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N}$$

- Six possible rules from itemset $\{Bread, Diaper, Milk\}$:

$\{Bread, Diaper\} \rightarrow \{Milk\}, \{Milk\} \rightarrow \{Bread, Diaper\}$

$\{Bread, Milk\} \rightarrow \{Diaper\}, \{Diaper\} \rightarrow \{Bread, Milk\}$

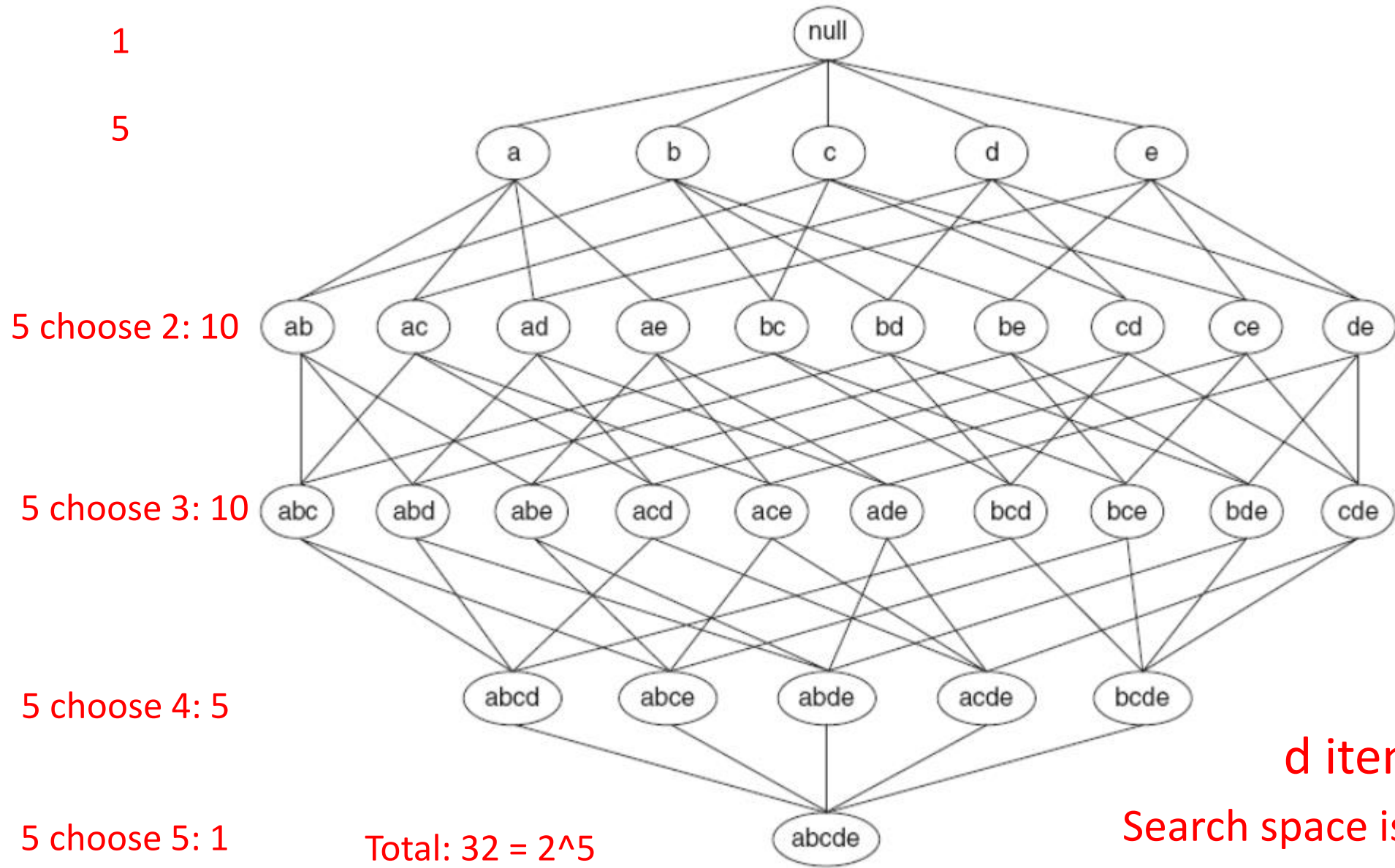
$\{Diaper, Milk\} \rightarrow \{Bread\}, \{Bread\} \rightarrow \{Diaper, Milk\}$

- If the itemset has low support:
 - All six rules have low support
 - Can be pruned

Better Approach

2-Step Approach: find frequent items then generate rules.

1. **Frequent itemset** generation
 - Generate all itemsets satisfying *minsup*
2. Rule Generation:
 - Extract rules satisfying *minconf* from **frequent itemsets**



1

5

5 choose 2: 10

5 choose 3: 10

5 choose 4: 5

5 choose 5: 1

Total: 32 = 2⁵

d items: 2^d itemsets

Search space is exponentially large

How to reduce computational complexity

- Reduce the number of candidate itemsets using the *Apriori* principles.
- Reduce the number of comparison: instead of matching each candidate itemset against every transaction, we can reduce the number by using more advanced data structures.
- Reduce the number of transactions

Apriori Principle

Anti-monotone property:

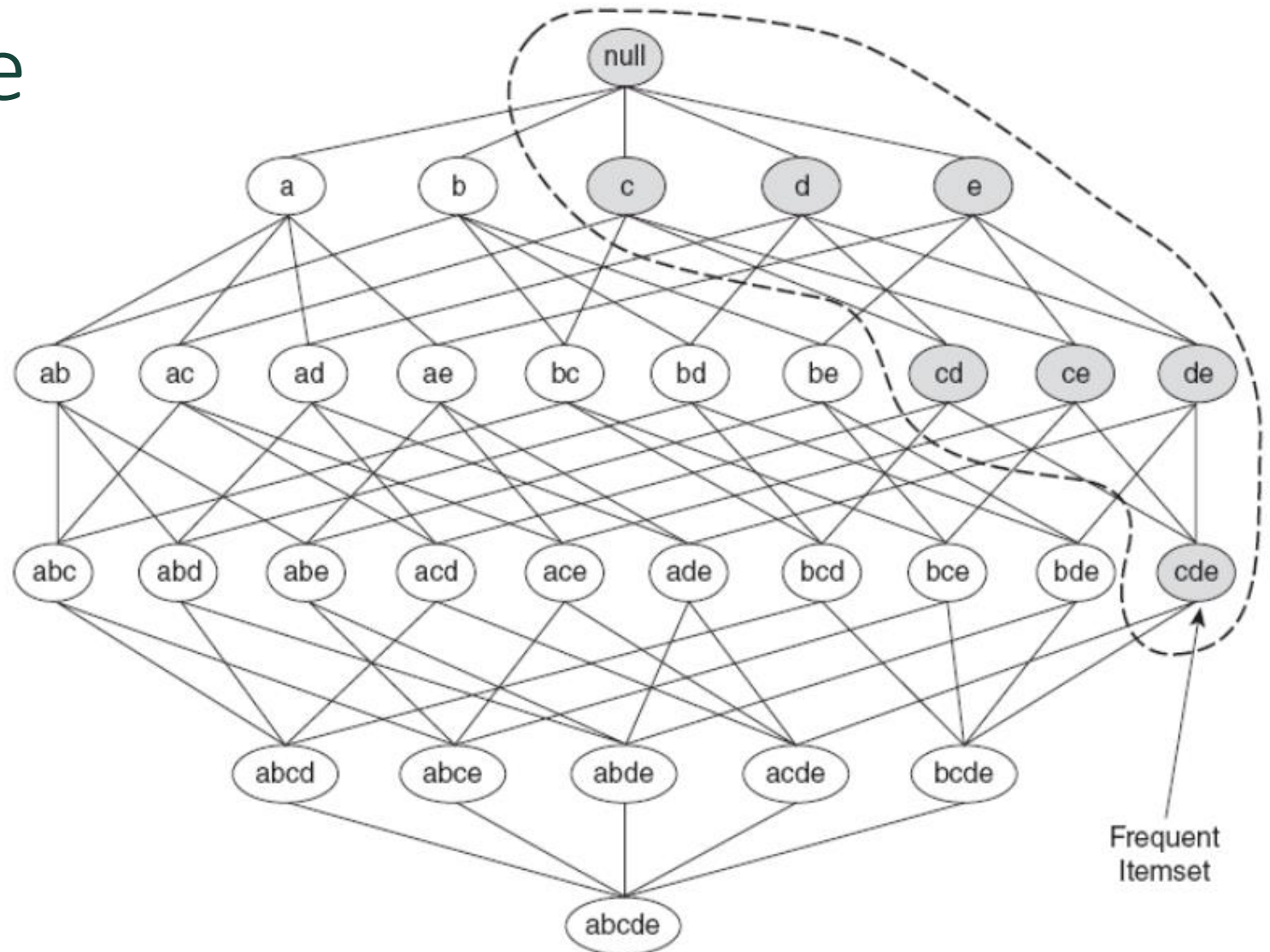
The support of an itemset never exceeds the support for the subsets

i.e., support of an itemset \leq support for its subsets

Apriori Principle

Anti-monotone property:
support of an itemset \leq
support for its subsets

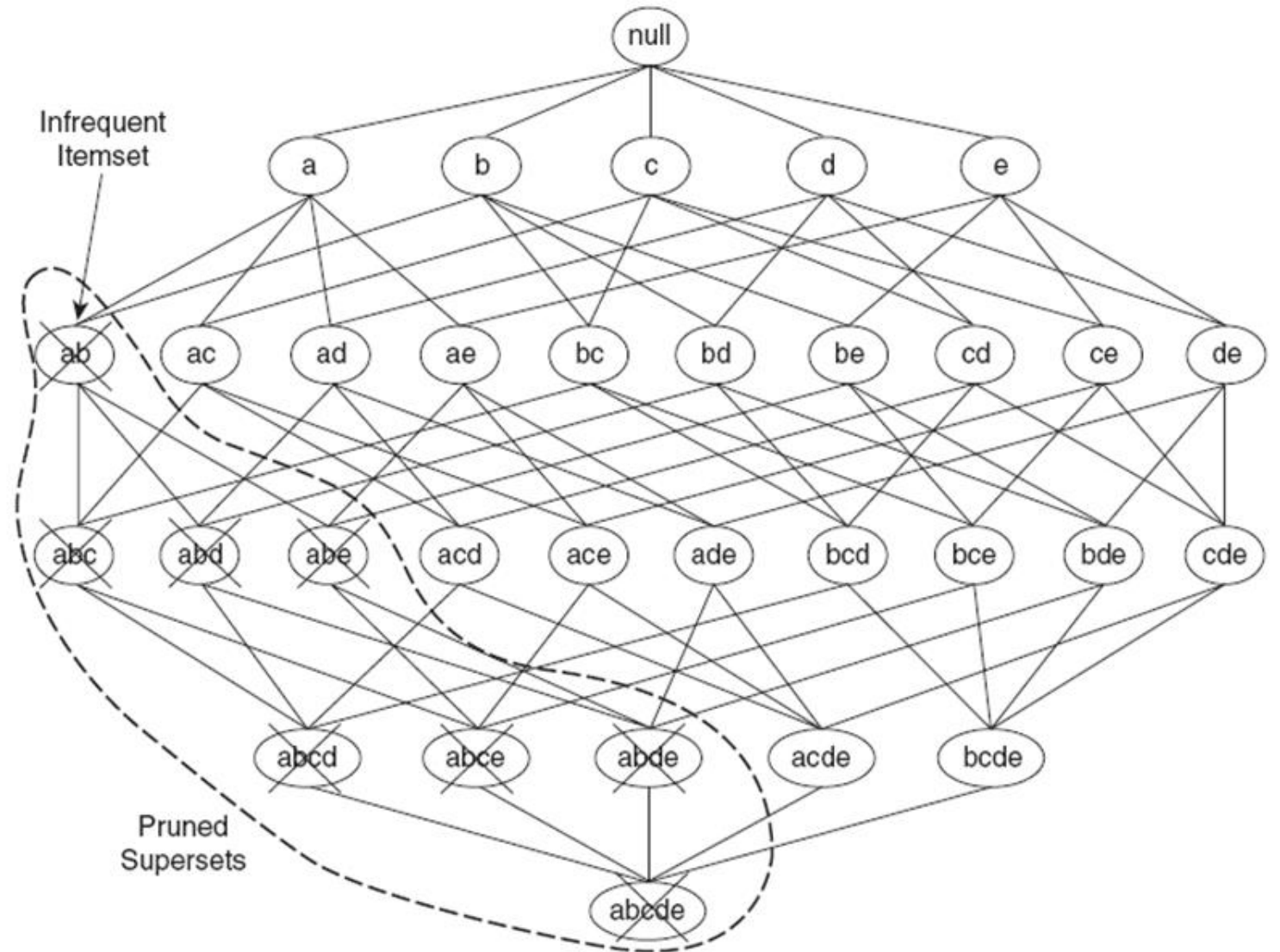
- If an itemset is frequent, then all its subsets are frequent



Apriori Principle

Anti-monotone property:
support of an itemset \leq
support for its subsets

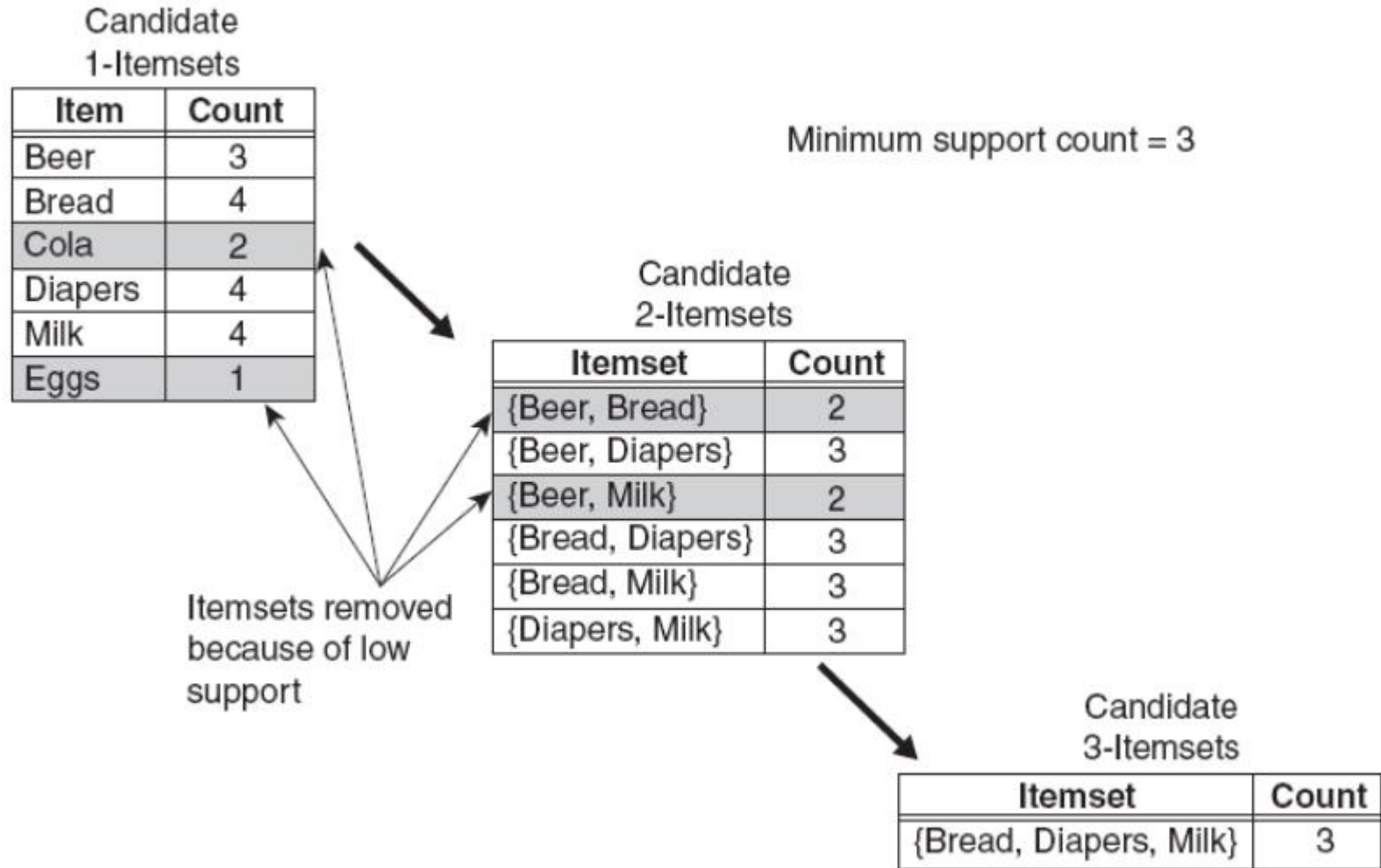
- If an itemset is infrequent, then all its supersets are infrequent



Apriori Algorithm

- C_k : Candidate itemsets of size k (itemsets possibly frequent)
- F_k : Frequent itemsets of size k
- Compute F_1 :
 - A single pass over the transactions table to count support of individual items
- Iteratively, use F_{k-1} to compute C_k and then F_k
- Stop when F_k is empty
- A pass over the transactions is needed to count the support of every C_k

Example



Candidate Set Generation

- Avoid generating too many unnecessary candidates
- Ensure the set is complete: no frequent itemset is left out
- Do not generate duplicate itemsets

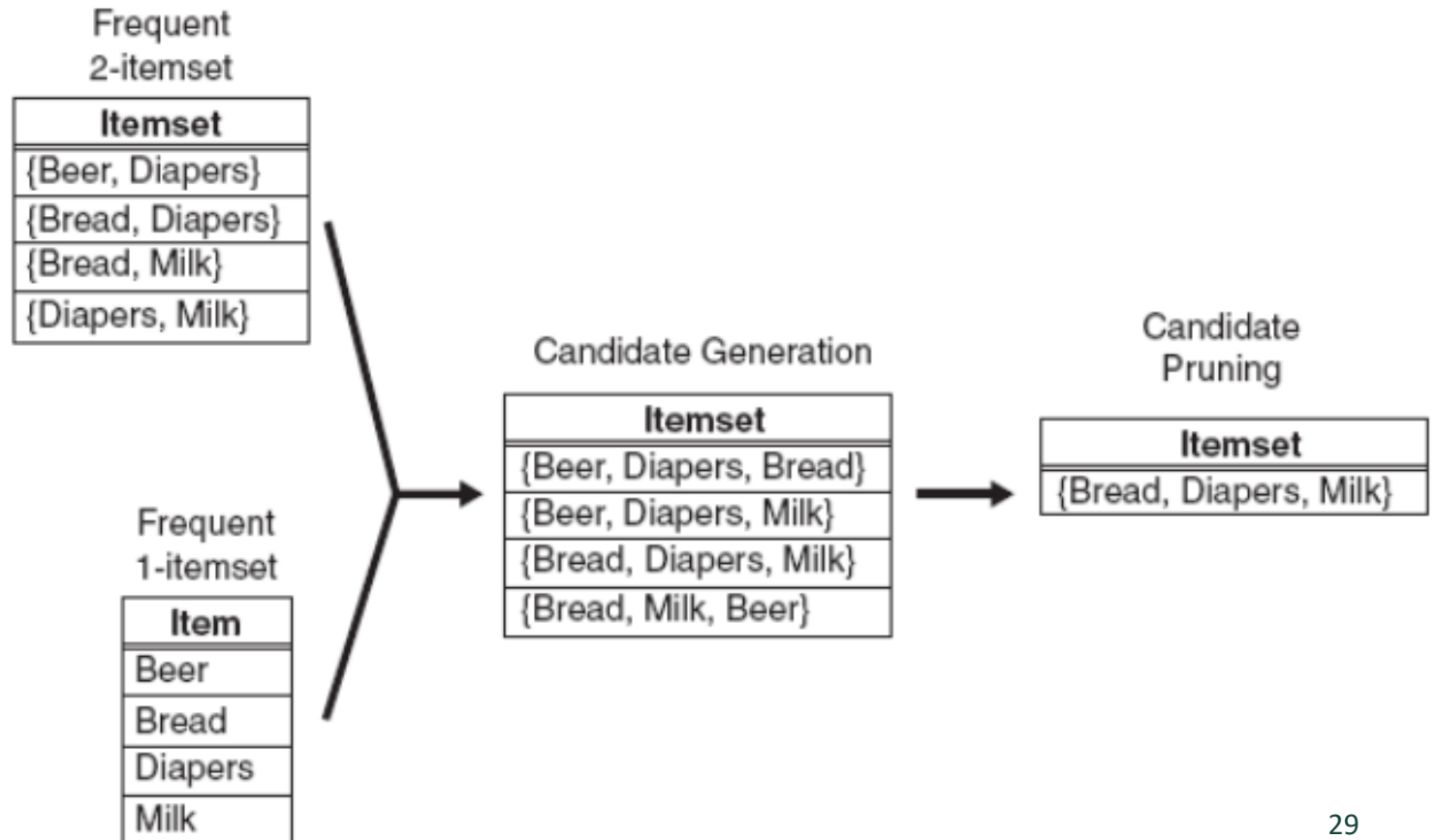
$\{a, b, c, d\}$ can be generated by merging:

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

$F_{k-1} \times F_1$ Method

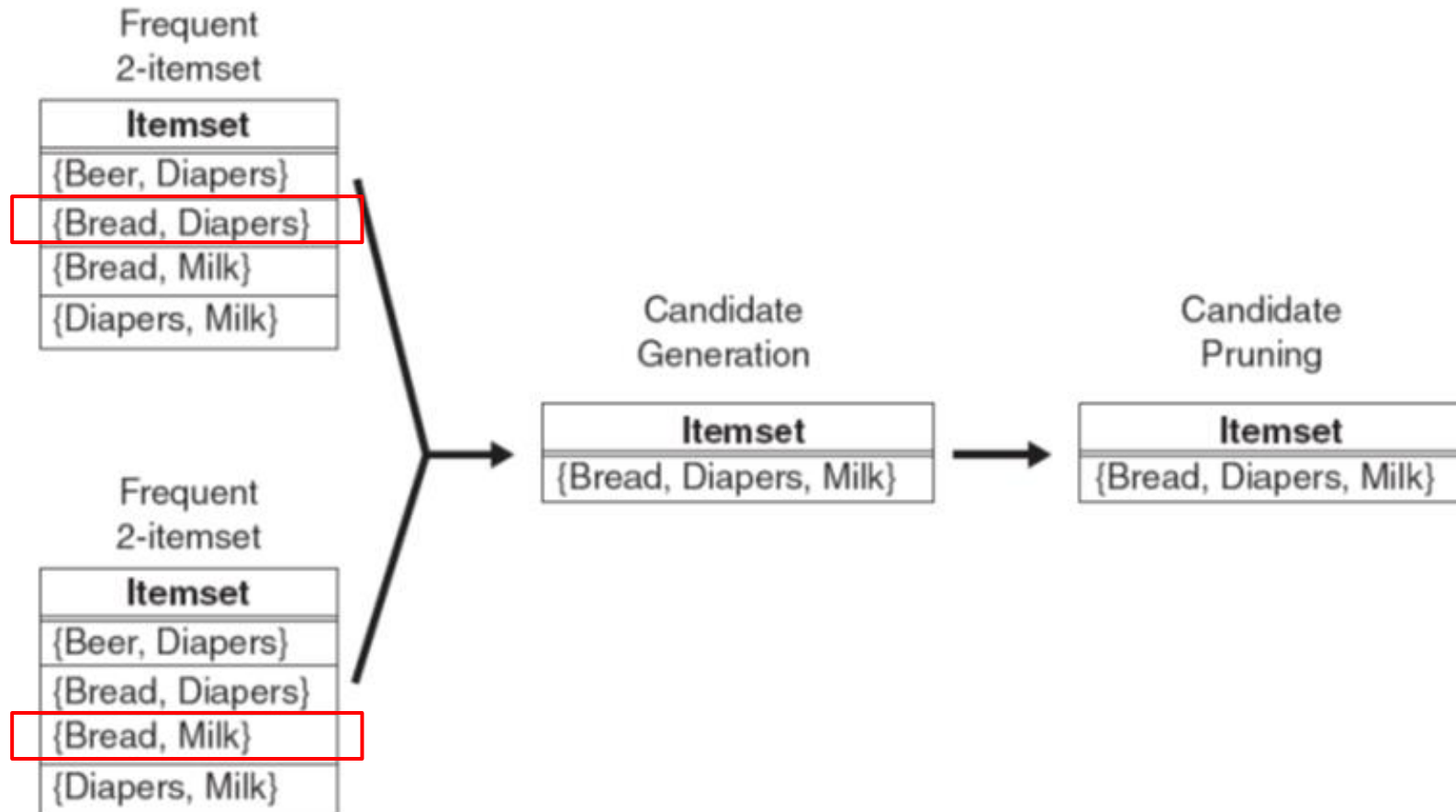
- Extend each itemset in F_{k-1} by a frequent item in F_1

Use F_{k-1} to compute C_k



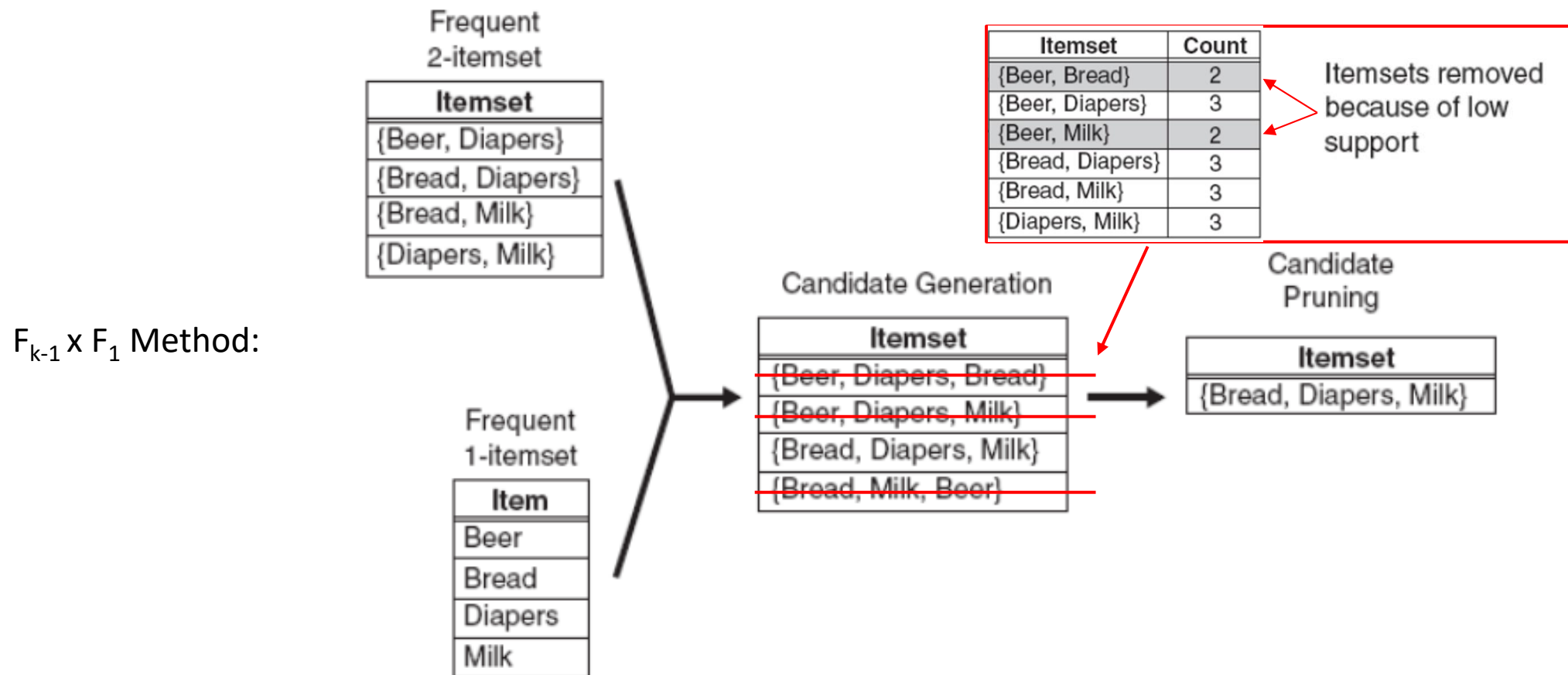
$F_{k-1} \times F_{k-1}$ Method (Apriori Gen)

- Merge two itemsets in F_{k-1} if their first $k-2$ items are identical



Candidate Pruning

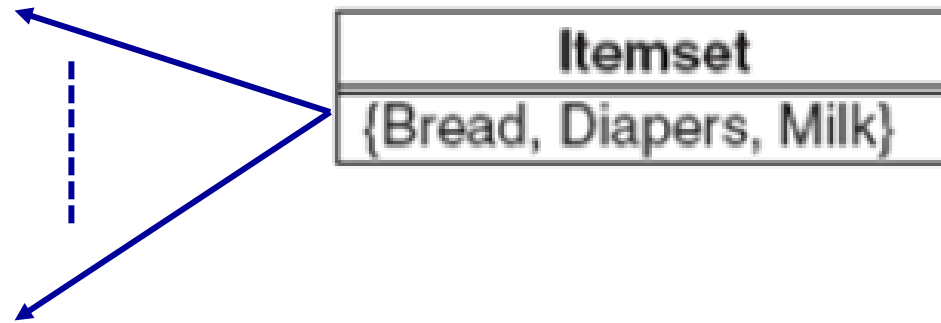
- Remove itemsets containing infrequent subsets:



Support Counting

- Compare each transaction against every itemset: **computationally expensive!**

<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



Improving Efficiency

- **Transaction Reduction:**
 - a transaction that does not contain any frequent k -itemset cannot contain any frequent $(k+1)$ -itemset
- **Sampling:** pick a random sample and find frequent itemsets on sample. Trading accuracy for efficiency

Computational Complexity

- **Support threshold:** lower support implies:
 - More frequent itemsets, more candidate itemsets
 - Larger frequent itemsets (larger k)
- **Number of items:**
 - More space needed to store support counts
 - Increases the number of candidate itemsets
- **Number of transactions:**
 - Increases the time needed for a pass of the data
- **Transaction Width:**
 - Increases the maximum size of frequent itemsets

Rule Generation

2-Step Approach: find frequent items then generate rules.

- ✓ **Frequent itemset** generation
 - Generate all itemsets satisfying *minsup*
- Rule Generation:
 - Extract rules satisfying *minconf* from **frequent itemsets**

Rule Generation

- Given the minimum confidence $minconf$, generating association rules by going through all possible combinations of **frequent item** sets and pruning the rules according to confidence criterion.
- Given a frequent k-itemset Z:
 - There are $2^k - 2$ possible association rules
 - Ignoring $\emptyset \rightarrow Z$ and $Z \rightarrow \emptyset$
- When considering rule , both $X \cup Y$ and X are frequent
 - Support is already computed
 - All rules satisfy $minsup$
 - Do not need to traverse transaction table
- Select only those satisfying $minconf$

$$c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$$

Rule Generation

- Generate all nonempty subsets for each frequent itemset
- For every nonempty subset S of Itemset I , output of the rule:
 - $S \rightarrow (I - S)$
 - **If** support_count (I) / support_count (S) \geq minimum confidence threshold **then** rule is a **strong Association Rule**.

TID	Items Bought
1	{ A, C }
2	{ A, B, C, E }
3	{ A, D }
4	{ A, B, C, E }
5	{ A, B, C, D, E }

minimum threshold support = 60% ,

minimum threshold confidence = 80% ,

Item set	Support Count	Support
{A, B, C }	3	60%

Rule Generation

minimum threshold confidence = 80%

Item set	Support Count	Support
{A, B, C }	3	60%

Step 1:

- Generate all nonempty subsets for each frequent itemset
 - For Itemset - { A, B, C } , all non empty subsets are {A,B}, {B,C}, {A,C}, {A}, {B}, {C}

Step 2.1:

- For every nonempty subset S of Itemset I , output of the rule:
 - $S \rightarrow (I - S)$
 - {A,B} \rightarrow {C}
 - {B,C} \rightarrow {A}
 - {A,C} \rightarrow {C}
 - {A} \rightarrow {B,C}
 - {B} \rightarrow {A,C}
 - {C} \rightarrow {A,B}

Rule Generation

minimum threshold confidence = 80%

Item set	Support Count	Support
{A, B, C }	3	60%

Step 2.2: ○ If $\text{support_count}(I) / \text{support_count}(S) \geq \text{minimum confidence threshold}$ then rule is a **strong Association Rule**.

- {A,B} -> {C}, Confidence = $3/3 * 100 = 100\%$ - **Yes**, it is a strong association rules
- {B,C} -> {A}, Confidence = $3/3 * 100 = 100\%$ - **Yes**, it is a strong association rules
- {A,C} -> {C}, Confidence = $3/4 * 100 = 80\%$ - **Yes**, it is a strong association rules
- {A} -> {B,C}, Confidence = $3/5 * 100 = 60\%$ - **No**, it is not a strong association rules
- {B} -> {A,C}, Confidence = $3/3 * 100 = 100\%$ - **Yes**, it is a strong association rules
- {C} -> {A,B}, Confidence = $3/4 * 100 = 80\%$ - **Yes**, it is a strong association rules

Credit Card Promotion Database

- 10 samples
- Single itemset can be twice as large than previous example

Magazine Promo	Watch Promo	Life Ins Promo	Credit Card Ins.	Sex
Yes	No	No	No	Male
Yes	Yes	Yes	No	Female
No	No	No	No	Male
Yes	Yes	Yes	Yes	Male
Yes	No	Yes	No	Female
No	No	No	No	Female
Yes	No	Yes	Yes	Male
No	Yes	No	No	Male
Yes	No	No	No	Male
Yes	Yes	Yes	No	Female

Single item sets at a 40% coverage threshold:

single item sets	Number of items
A. Magazine Promo=Yes	7
B. Watch Promo=Yes	4
C. Watch Promo=No	6
D. Life Ins Promo=Yes	5
E. Life Ins Promo=No	5
F. Credit Card Ins=No	8
G. Sex=Male	6
H. Sex=Female	4

Credit Card Promotion Database

Single item sets at a 40% coverage threshold:

single item sets	Number of items
A. Magazine Promo=Yes	7
B. Watch Promo=Yes	4
C. Watch Promo=No	6
D. Life Ins Promo=Yes	5
E. Life Ins Promo=No	5
F. Credit Card Ins=No	8
G. Sex=Male	6
H. Sex=Female	4

Now begin pairing up combinations with the same minimal support threshold (40%)

	A	B	C	D	E	F	G	H
B	3	-						
C	4		-					
D	5			-				
E	2		4		-			
F	5		5		5	-		
G	4		4		4	4	-	
H						4		-

Credit Card Promotion Database

	A	B	C	D	E	F	G	H
B	3	-						
C	4		-					
D	5			-				
E	2		4		-			
F	5		5		5	-		
G	4		4		4	4	-	
H						4		-

Resulting rules from two item sets. Consider rules in both directions:

1. (A \rightarrow D)
(MagazinePromo=Yes) \rightarrow (LifeInsPromo=Yes) at 5/7 confidence
2. (D \rightarrow A)
(LifeInsPromo=Yes) \rightarrow (MagazinePromo=Yes) at 5/5 confidence
3. twenty more rules from the 10 two-item-sets (A then C, C then A, A then F, F then A, etc.)

Now apply minimum confidence threshold

If confidence threshold would be 80%, then the first rule (A \rightarrow D) is eliminated.