

## Chapter 5

### Association Analysis: Basic Concepts

Introduction to Data Mining, 2<sup>nd</sup> Edition

by

Tan, Steinbach, Karpatne, Kumar

# Association Rule Mining

- Bir dizi işlem (*transactions*) verildiğinde, işlemlerdeki diğer öğelerin oluşumlarına bağlı olarak **bir öğenin oluşumunu tahmin** edecek kuralları bulma

## 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

## Example of Association Rules

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

Çıkarım, nedensellik değil, birlikte meydana gelme anlamına gelir!

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

# Definition: Frequent Itemset

- **Itemset**

- Bir veya daha fazla öğeden oluşan bir koleksiyon, öğe kümesi

- ◆ Example: {Milk, Bread, Diaper}

- k-itemset

- ◆ k tane öğe içeren bir öğe kümesi

- **Support count ( $\sigma$ )**

- Bir öğe kümesinin ortaya çıkma sıklığı (*frequency*)

- E.g.  $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

- **Support**

- Bir öğe kümesini içeren transaction'ların oranı

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

- **Frequent Itemset**

- Desteği belirli bir eşik değerinden (*minsup*) büyük veya ona eşit olan bir öğe kümesi

<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

# Definition: Association Rule

- **Association Rule**

- An implication expression of the form  $X \rightarrow Y$ , where  $X$  and  $Y$  are itemsets
- Example:  
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

<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

- **Rule Evaluation Metrics**

- Support (s)
  - ◆ Hem  $X$  hem de  $Y$  içeren işlemlerin oranı
- Confidence (c)
  - ◆  $X$  içeren işlemlerde  $Y$ 'deki öğelerin ne sıklıkla görüldüğünü ölçer

Example:

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

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

$$\text{Support, } s(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N};$$

$$\text{Confidence, } c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}.$$

# Why Use Support and Confidence?

- Destek önemli bir ölçüdür çünkü **desteği çok düşük** olan bir kural sadece **şans eseri (*by chance*)** ortaya çıkabilir.
- Düşük destek (**support**) kuralı, **müşterilerin nadiren birlikte satın aldıkları ürünleri tanıtmak karlı olmayabileceğinden**, işletme açısından da ilgi çekici olmayabilir (***uninteresting***) .
  - Bu nedenlerden dolayı, **ilgi çekici olmayan kuralları ortadan kaldırmak için genellikle destek kullanılır.**
- Öte yandan güven (**Confidence**), bir kural tarafından yapılan **çıkarımın güvenilirliğini** ölçer.
  - Belirli bir  $X \rightarrow Y$  kuralı için, güven ne kadar yüksekse,  $Y$ 'nin  $X$  içeren işlemlerde mevcut olma olasılığı o kadar yüksektir.

# Association Rule Mining Task

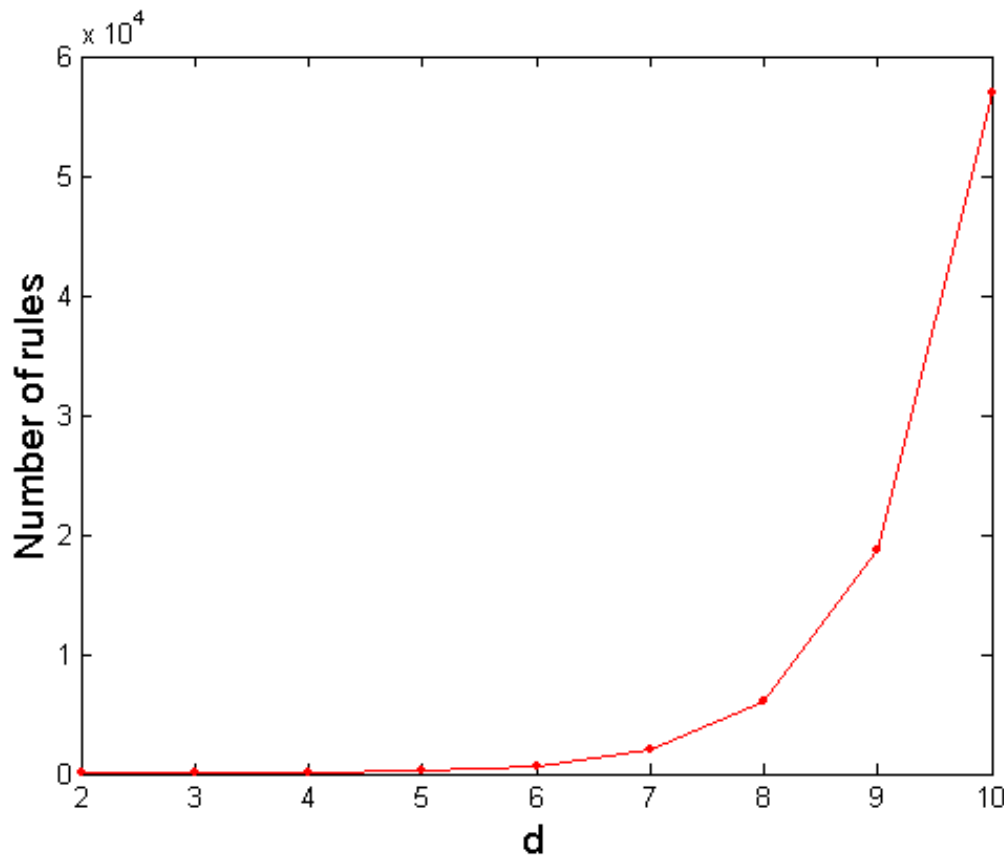
---

- Bir dizi transaction  $T$  verildiğinde, birliktelik kuralı madenciliğinin amacı, şunlara sahip olan tüm kuralları bulmaktır.
  - support  $\geq$  *minsup* threshold
  - confidence  $\geq$  *minconf* threshold
- Brute-force approach:
  - List all possible association rules
  - Compute the support and confidence for each rule
  - Prune rules that fail the *minsup* and *minconf* thresholds

⇒ **Computationally prohibitive!**

# Computational Complexity

- Given  $d$  unique items:
  - Total number of itemsets =  $2^d$
  - Total number of possible association rules:



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

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

# Mining Association Rules

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

## Example of Rules:

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

## Observations:

- Yukarıdaki kuralların tümü aynı öge kümesinin ikili bölümleridir (*binary partitions of the same itemset*):  
 $\{\text{Milk, Diaper, Beer}\}$
- Aynı öge setinden kaynaklanan kurallar aynı desteğe sahiptir ancak farklı güvenlere sahip olabilir
- Böylece, destek ve güven gereksinimlerini ayrıştırabiliriz

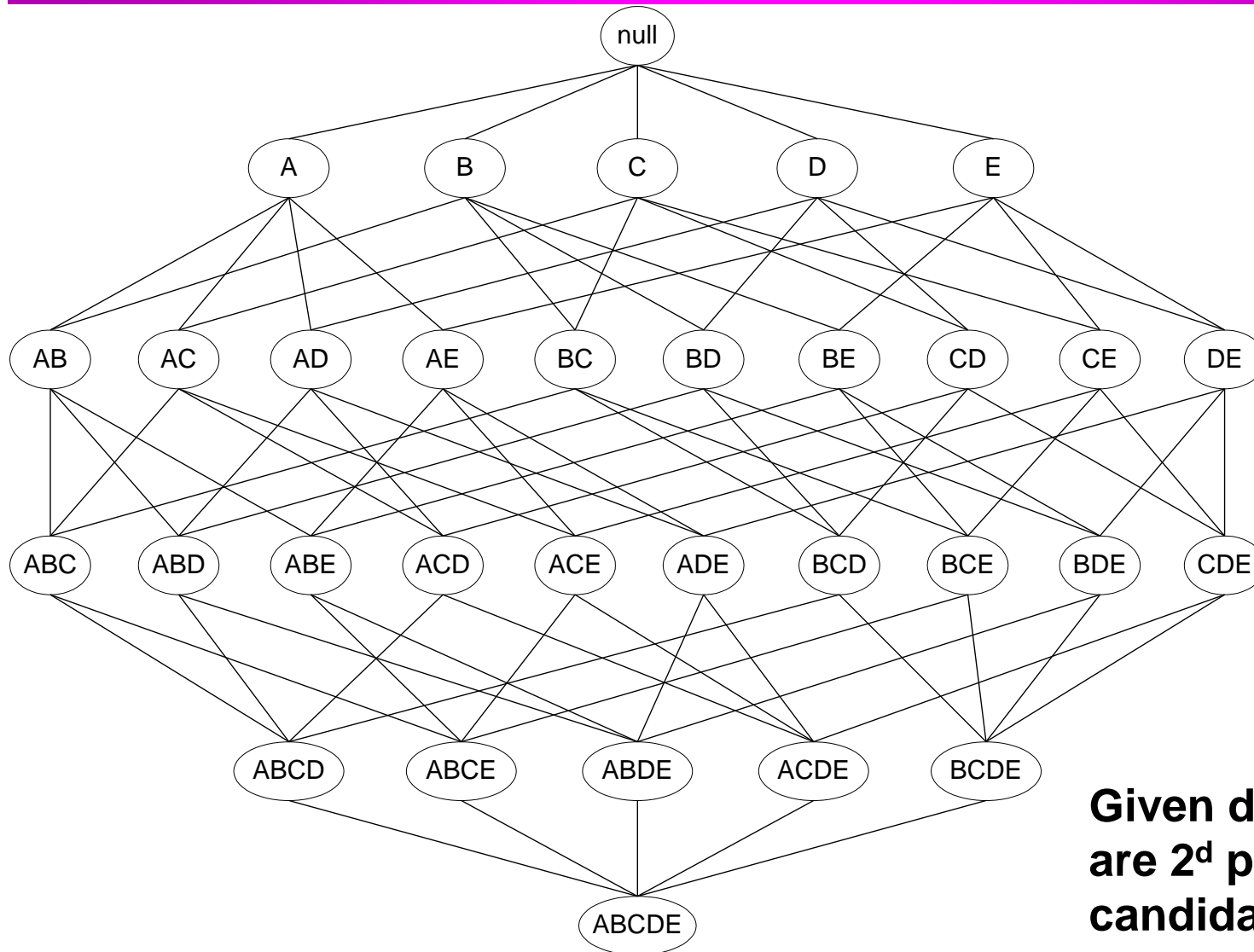


# Mining Association Rules

---

- Two-step approach:
  1. Frequent Itemset Generation
    - Generate all itemsets whose support  $\geq$  minsup
  2. Rule Generation
    - Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset
- Frequent itemset generation is still computationally expensive

# Frequent Itemset Generation



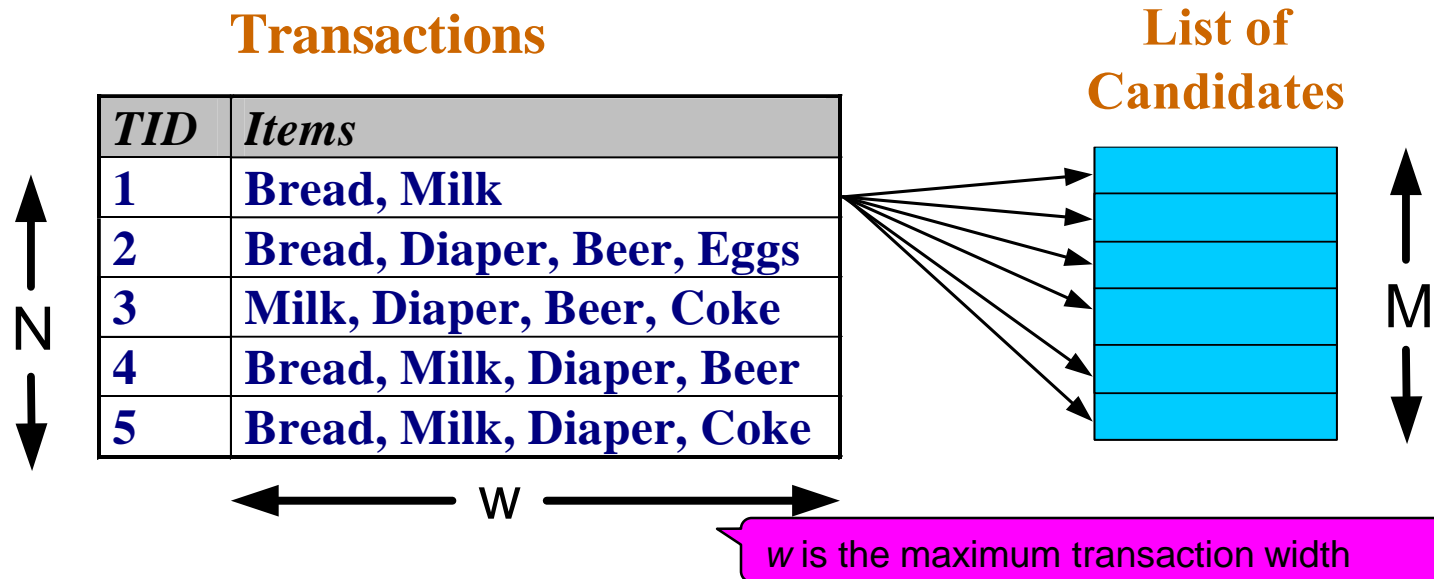
**Given  $d$  items, there are  $2^d$  possible candidate itemsets**

# Frequent Itemset Generation

- Brute-force approach:

- Each itemset in the lattice is a **candidate** frequent itemset
- Count the support of each candidate by scanning the database

If the candidate is contained in a transaction, its support count will be incremented.



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

# Frequent Itemset Generation Strategies

---

- Reduce the **number of candidates** (M)
  - Complete search:  $M=2^d$
  - Use pruning techniques to reduce M
    - ◆Örneğin **Apriori prensibi**, bazı aday öge setlerini destek değerlerini saymadan ortadan kaldırmanın etkili bir yoludur.
- Reduce the **number of comparisons** (NM)
  - Adayları veya transactionları depolamak için etkili veri yapılarını kullanın
  - Her adayı her transaction ile karşılaştırmaya gerek yok

# Reducing Number of Candidates

- **Apriori principle:**

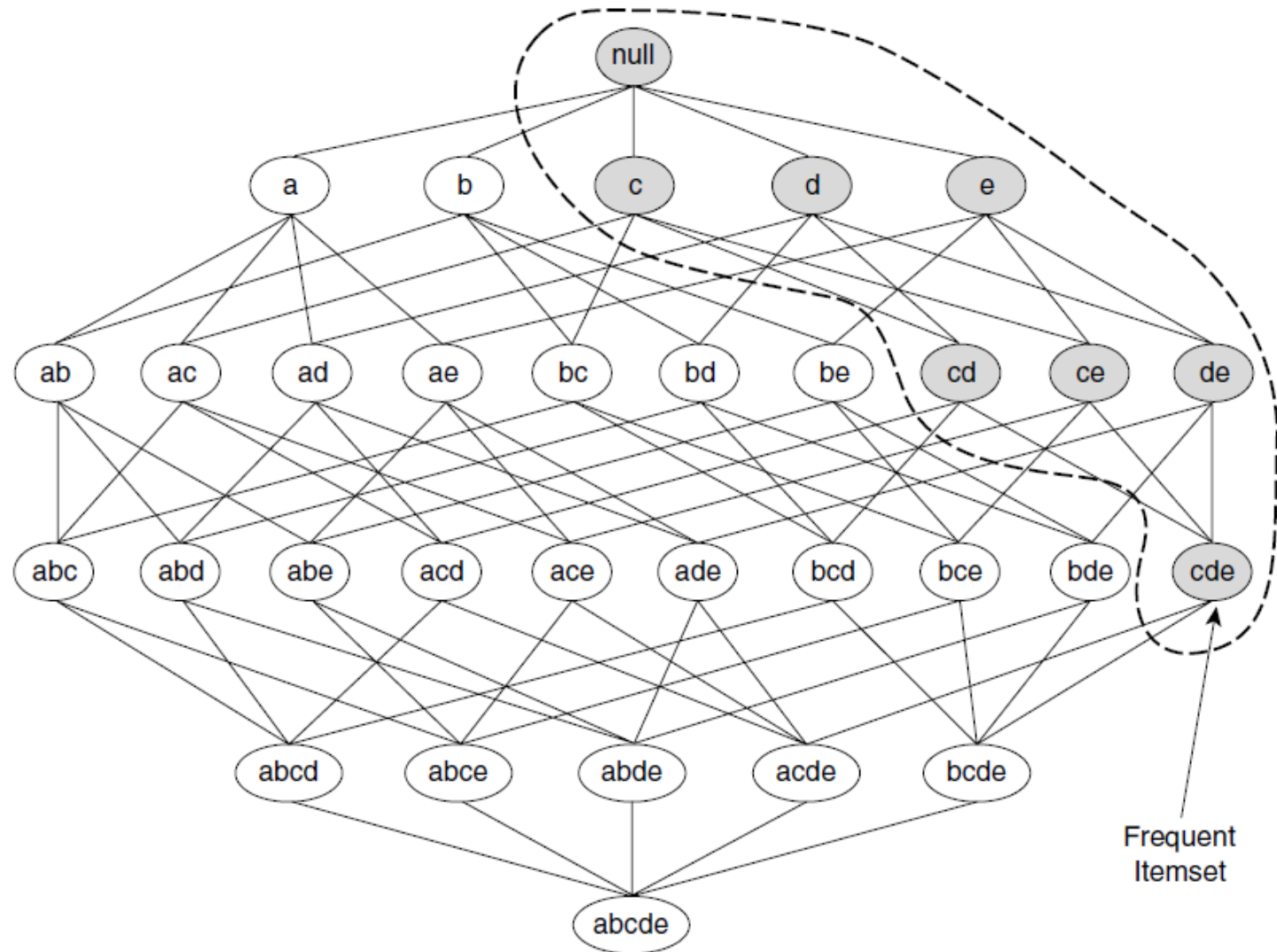
- If an itemset is frequent, then all of its subsets must also be frequent (*Bir öge kümesi «frequent» ise, onun alt-kümeleri de «frequent» olmalı*)

- Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

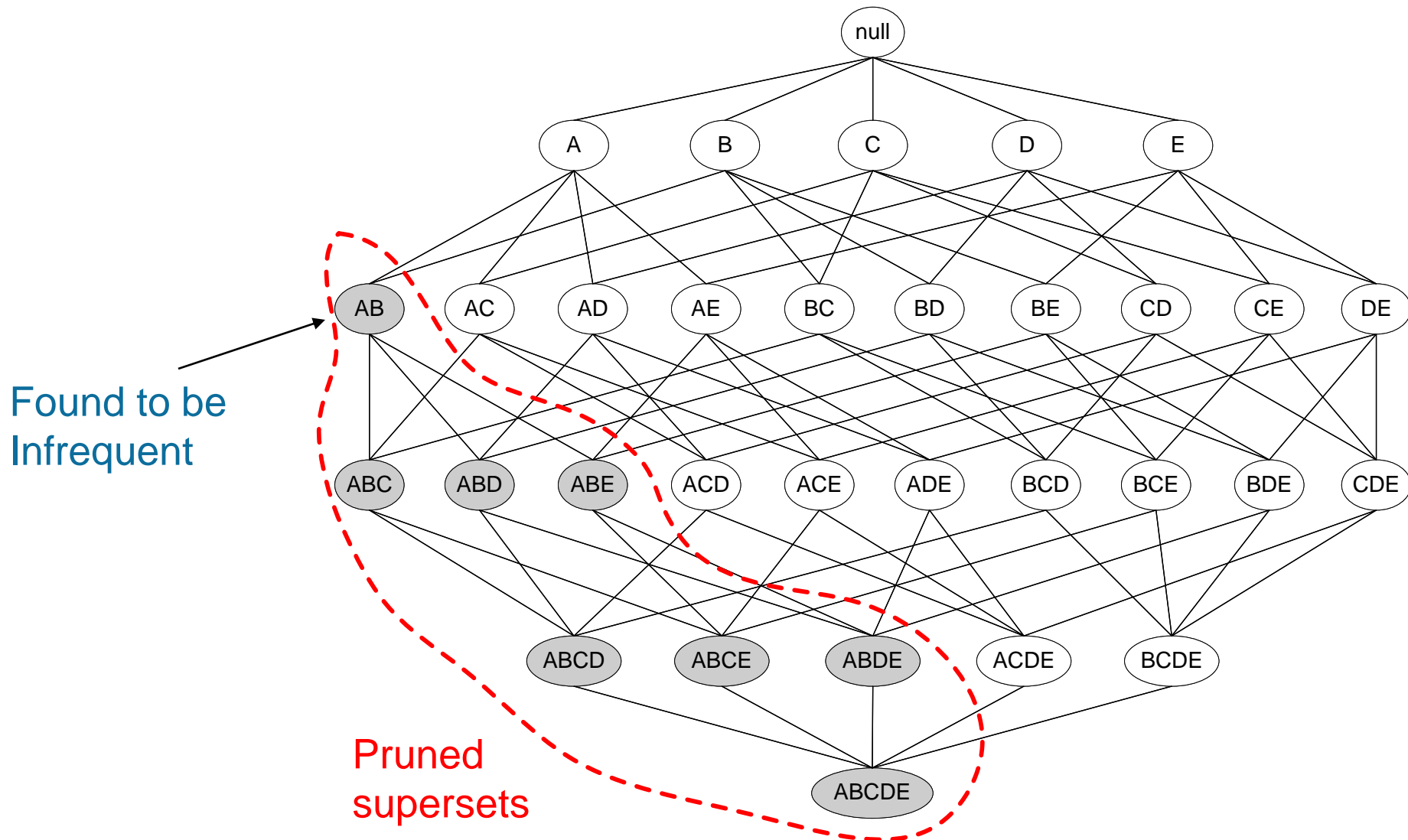
- Support of an itemset never exceeds the support of its subsets (*Bir öge ekümesinin desteği, asla alt kümelerinin desteğinden büyük olamaz*)
- This is known as the **anti-monotone** property of support

# Illustrating Apriori Principle



**Figure 6.3.** An illustration of the *Apriori* principle. If  $\{c, d, e\}$  is frequent, then all subsets of this itemset are frequent.

# Illustrating Apriori Principle



# Illustrating Apriori Principle

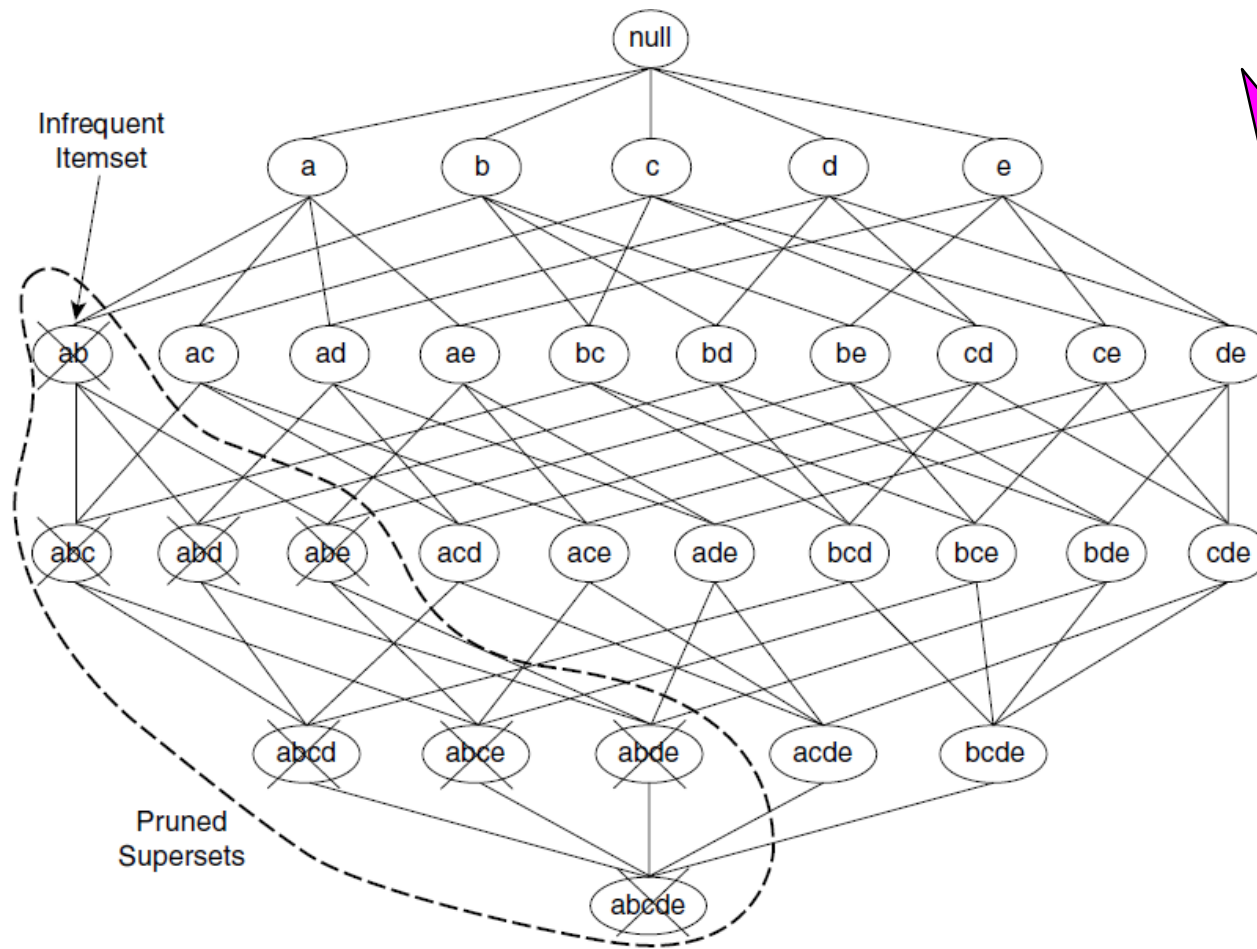


Figure 6.4. An illustration of support-based pruning. If  $\{a, b\}$  is infrequent, then all supersets of  $\{a, b\}$  are infrequent.

Destek ölçüsüne dayalı olarak üstel arama alanını küçültme stratejisi, desteğe dayalı budama (**support-based pruning**) olarak bilinir.

Böyle bir budama stratejisi, destek ölçüsünün temel bir özelliği ile, yani bir öge kümesine yönelik desteğin, alt gruplarının desteğini hiçbir zaman aşmaması ile mümkün kılınmaktadır. Bu özellik, destek ölçütünün **anti-monoton** özelliği olarak da bilinir.



# Illustrating Apriori Principle

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



Items (1-itemsets)

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Minimum Support = 3

If every subset is considered,  
 ${}^6C_1 + {}^6C_2 + {}^6C_3$   
 $6 + 15 + 20 = 41$   
 With support-based pruning,  
 $6 + 6 + 4 = 16$

We assume that the support threshold is 60%, which is equivalent to a **minimum support count** equal to 3.

Reminder: Combination formula

$$C(n, r) = \binom{n}{r} = \binom{n}{n-r} = \frac{P(n, r)}{r!} = \frac{n!}{r!(n-r)!}$$

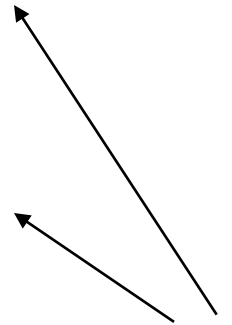
# Illustrating Apriori Principle

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



Items (1-itemsets)

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1



Minimum Support = 3

If every subset is considered,

$${}^6C_1 + {}^6C_2 + {}^6C_3$$

$$6 + 15 + 20 = 41$$

With support-based pruning,

$$6 + 6 + 4 = 16$$

$$\binom{6}{1} + \binom{6}{2} + \binom{6}{3}$$

$$\binom{6}{1} + \binom{4}{2} + \binom{4}{3}$$

Itemsets removed  
because of low  
support

# Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset
{Bread,Milk}
{Bread, Beer }
{Bread,Diaper}
{Beer, Milk}
{Diaper, Milk}
{Beer,Diaper}

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,

$${}^6C_1 + {}^6C_2 + {}^6C_3$$

$$6 + 15 + 20 = 41$$

With support-based pruning,

$$6 + 6 + 4 = 16$$

# Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Beer, Bread}	2
{Bread,Diaper}	3
{Beer,Milk}	2
{Diaper,Milk}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,

$${}^6C_1 + {}^6C_2 + {}^6C_3 \\ 6 + 15 + 20 = 41$$

With support-based pruning,

$$6 + 6 + 4 = 16$$

Itemsets removed  
because of low  
support

# Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)



Triplets (3-itemsets)

Itemset
{ Beer, Diaper, Milk}
{ Beer,Bread,Diaper}
{Bread, Diaper, Milk}
{ Beer, Bread, Milk}

Minimum Support = 3

If every subset is considered,

$${}^6C_1 + {}^6C_2 + {}^6C_3 \\ 6 + 15 + 20 = 41$$

With support-based pruning,

$$6 + 6 + 4 = 16$$

With the *Apriori* principle, we only need to keep candidate 3-itemsets **whose subsets are frequent**. The only candidate that has this property is {Bread, Diapers, Milk}.

# Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,

$${}^6C_1 + {}^6C_2 + {}^6C_3 \\ 6 + 15 + 20 = 41$$

With support-based pruning,

$$6 + 6 + 4 = 16 \\ 6 + 6 + 1 = 13$$



Triplets (3-itemsets)

Itemset	Count
{ Beer, Diaper, Milk}	2
{ Beer,Bread, Diaper}	2
{Bread, Diaper, Milk}	2
{Beer, Bread, Milk}	1

Apriori prensibi ile bu sayı 13 adaya düşüyor, bu da bu basit örnekte bile aday öge setlerinin sayısında %68'lik bir azalmayı temsil ediyor.

# Apriori Algorithm

---

- $F_k$ : frequent k-itemsets
- $L_k$ : candidate k-itemsets
- Algorithm
  - Let  $k=1$
  - Generate  $F_1 = \{\text{frequent 1-itemsets}\}$
  - Repeat until  $F_k$  is empty
    - ◆ **Candidate Generation:** Generate  $L_{k+1}$  from  $F_k$
    - ◆ **Candidate Pruning:** Prune candidate itemsets in  $L_{k+1}$  containing subsets of length  $k$  that are infrequent
    - ◆ **Support Counting:** Count the support of each candidate in  $L_{k+1}$  by scanning the DB
    - ◆ **Candidate Elimination:** Eliminate candidates in  $L_{k+1}$  that are infrequent, leaving only those that are frequent  $\Rightarrow F_{k+1}$

# Candidate Generation: Brute-force method

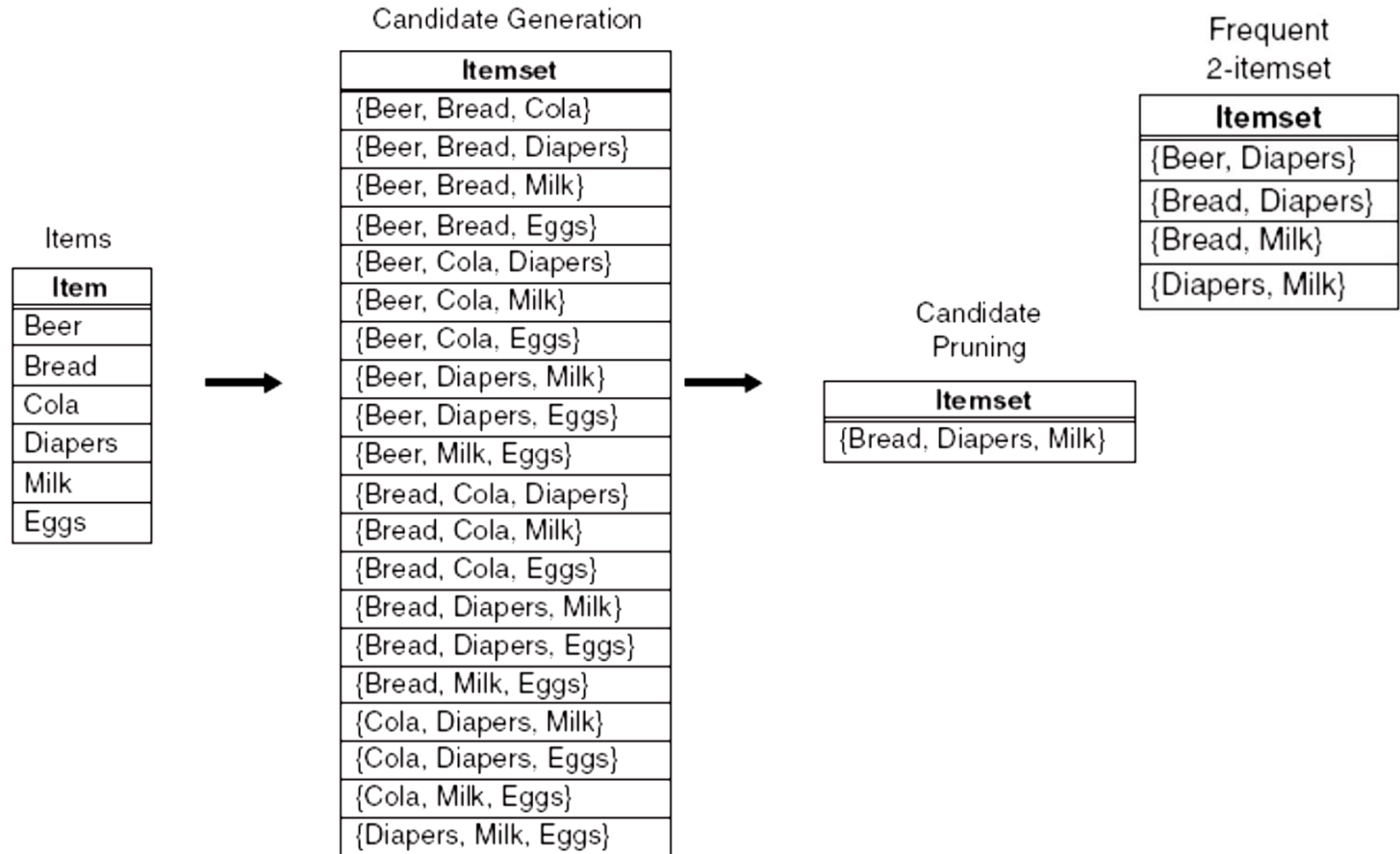


Figure 6.6. A brute-force method for generating candidate 3-itemsets.



# Candidate Generation: Merge Fk-1 and F1 itemsets

Items (2-itemsets)

Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Frequent  
2-itemset

Itemset
{Beer, Diapers}
{Bread, Diapers}
{Bread, Milk}
{Diapers, Milk}

Örneğin, budama adımını atlatan her aday  $k$ -itemset için adaydaki her öğenin, frequent  $(k - 1)$ -itemset'lerin en az  $k - 1$  'inde yer alması gerektiğini unutmayın. Aksi takdirde, adayın *infrequent* olması garanti edilir.

Items (1-itemsets)

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Frequent  
1-itemset

Item
Beer
Bread
Diapers
Milk

Candidate Generation

Itemset
{Beer, Diapers, Bread}
{Beer, Diapers, Milk}
{Bread, Diapers, Milk}
{Bread, Milk, Beer}

Candidate  
Pruning

Itemset
{Bread, Diapers, Milk}

Örneğin, { Beer, Diapers, Milk } ,yalnızca *Beer* da dahil olmak üzere adaydaki her öğe en az iki frequent 2-itemsets yer alıyorsa geçerli bir aday 3-itemset'tir. Beer içeren yalnızca bir tane frequent 2-itemset grubu olduğundan, Beer içeren tüm aday öğe kümeleri *infrequent* olmalıdır.

Figure 6.7. Generating and pruning candidate  $k$ -itemsets by merging a frequent  $(k - 1)$ -itemset with a frequent item. Note that some of the candidates are unnecessary because their subsets are infrequent.

# Candidate Generation: $F_{k-1} \times F_{k-1}$ Method

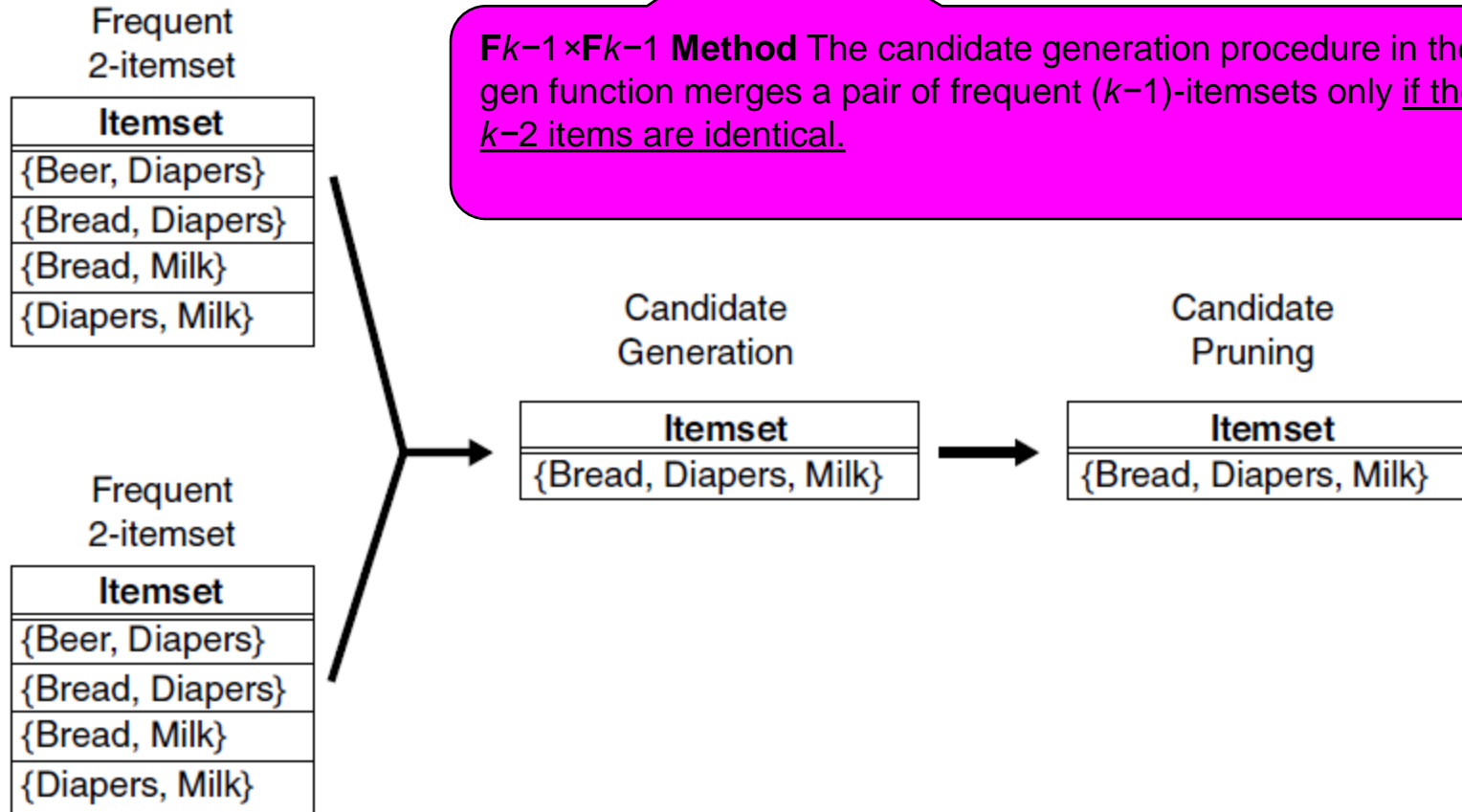


Figure 6.8. Generating and pruning candidate  $k$ -itemsets by merging pairs of frequent  $(k-1)$ -itemsets.

# Candidate Generation: $F_{k-1} \times F_{k-1}$ Method

---

- Merge two frequent  $(k-1)$ -itemsets if their first  $(k-2)$  items are identical
- $F_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$ 
  - Merge(ABC, ABD) = ABCD
  - Merge(ABC, ABE) = ABCE
  - Merge(ABD, ABE) = ABDE
  - Do not merge(ABD, ACD) because they share only prefix of length 1 instead of length 2

# Candidate Pruning

---

- Let  $F_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$  be the set of frequent 3-itemsets
- $L_4 = \{ABCD, ABCE, ABDE\}$  is the set of candidate 4-itemsets generated (from previous slide)
- Candidate pruning
  - Prune ABCE because ACE and BCE are infrequent
  - Prune ABDE because ADE is infrequent
- After candidate pruning:  $L_4 = \{ABCD\}$

# Alternate $F_{k-1} \times F_{k-1}$ Method

---

- Merge two frequent  $(k-1)$ -itemsets if the last  $(k-2)$  items of the first one is identical to the first  $(k-2)$  items of the second.
- $F_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$ 
  - Merge(ABC, BCD) = ABCD
  - Merge(ABD, BDE) = ABDE
  - Merge(ACD, CDE) = ACDE
  - Merge(BCD, CDE) = BCDE

## Candidate Pruning for Alternate $F_{k-1} \times F_{k-1}$ Method

---

- Let  $F_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$  be the set of frequent 3-itemsets
- $L_4 = \{ABCD, ABDE, ACDE, BCDE\}$  is the set of candidate 4-itemsets generated (from previous slide)
- Candidate pruning
  - Prune ABDE because ADE is infrequent
  - Prune ACDE because ACE and ADE are infrequent
  - Prune BCDE because BCE is infrequent
- After candidate pruning:  $L_4 = \{ABCD\}$

# Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)



Itemset	Count
{Bread, Diaper, Milk}	2

Triplets (3-itemsets)

Minimum Support = 3

If every subset is considered,

$${}^6C_1 + {}^6C_2 + {}^6C_3 \\ 6 + 15 + 20 = 41$$

With support-based pruning,

$$6 + 6 + 1 = 13$$

Use of  $F_{k-1} \times F_{k-1}$  method for candidate generation results in only one 3-itemset. This is eliminated after the support counting step.

# Support Counting of Candidate Itemsets

- Her adaya öğesinin desteğini belirlemek için transaction veritabanını tarayın
  - Her aday öğesini her transaction ile karşılaştırmalıdır, bu zaman alıcı bir işlemdir

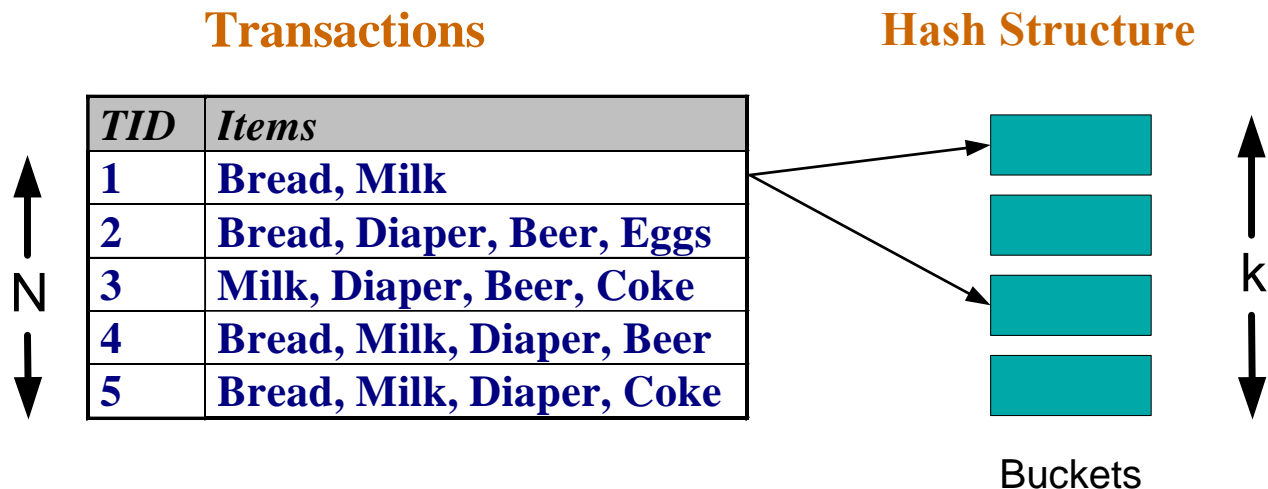
<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Coke, Diaper, Milk

Itemset
{ Beer, Diaper, Milk }
{ Beer, Bread, Diaper }
{ Bread, Diaper, Milk }
{ Beer, Bread, Milk }



# Support Counting of Candidate Itemsets

- Karşılaştırma sayısını azaltmak için, aday öge kümelerini bir hash yapıda saklayın
  - Her transaction'ı her adayla karşılaştırmak yerine, hashing uygulanmış kovalarda bulunan adaylarla karşılaştırın

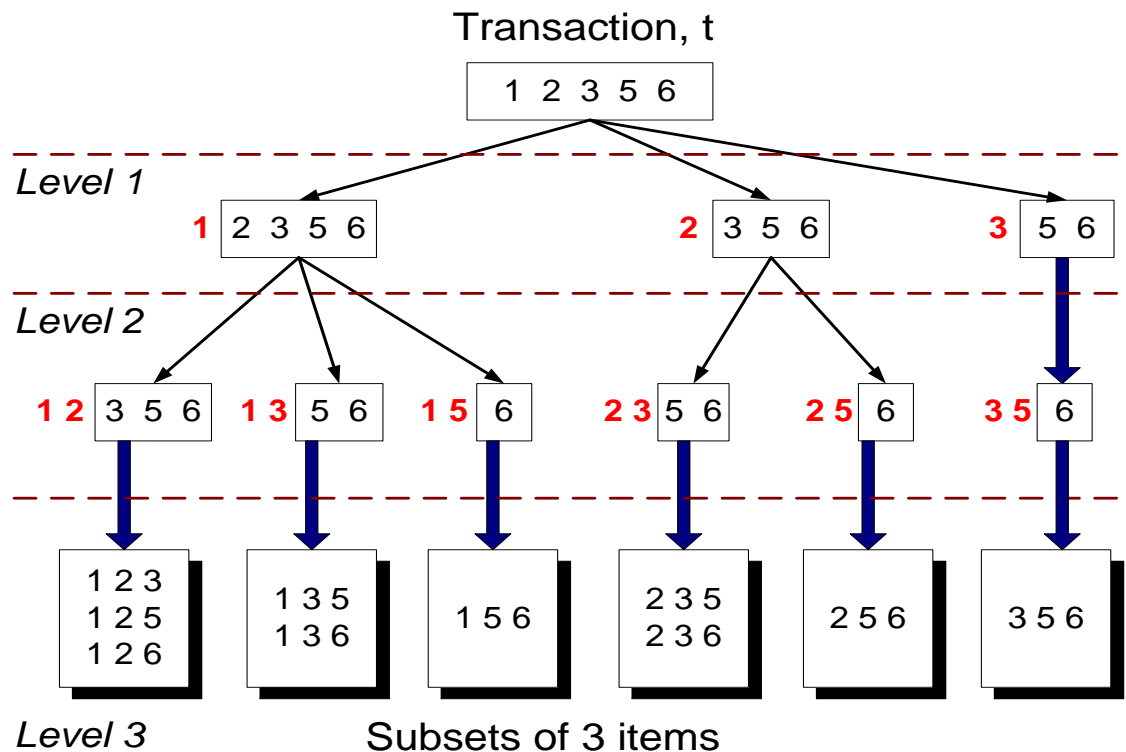


# Support Counting: An Example

Suppose you have 15 candidate itemsets of length 3:

{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4}, {5 6 7}, {3 4 5},  
{3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}

How many of these itemsets are supported by transaction (1,2,3,5,6)?



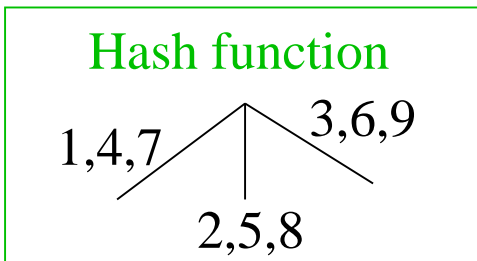
# Support Counting Using a Hash Tree

Suppose you have 15 candidate itemsets of length 3:

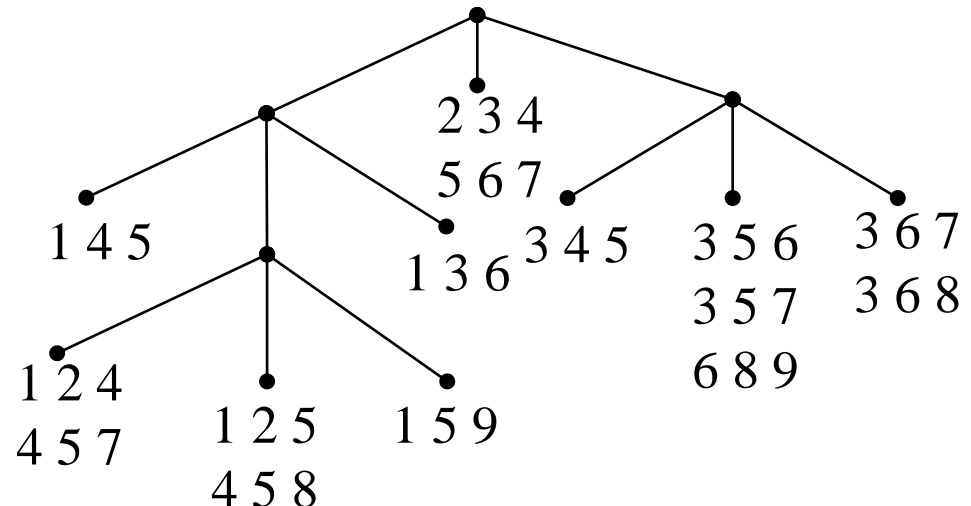
{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4}, {5 6 7}, {3 4 5},  
{3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}

You need:

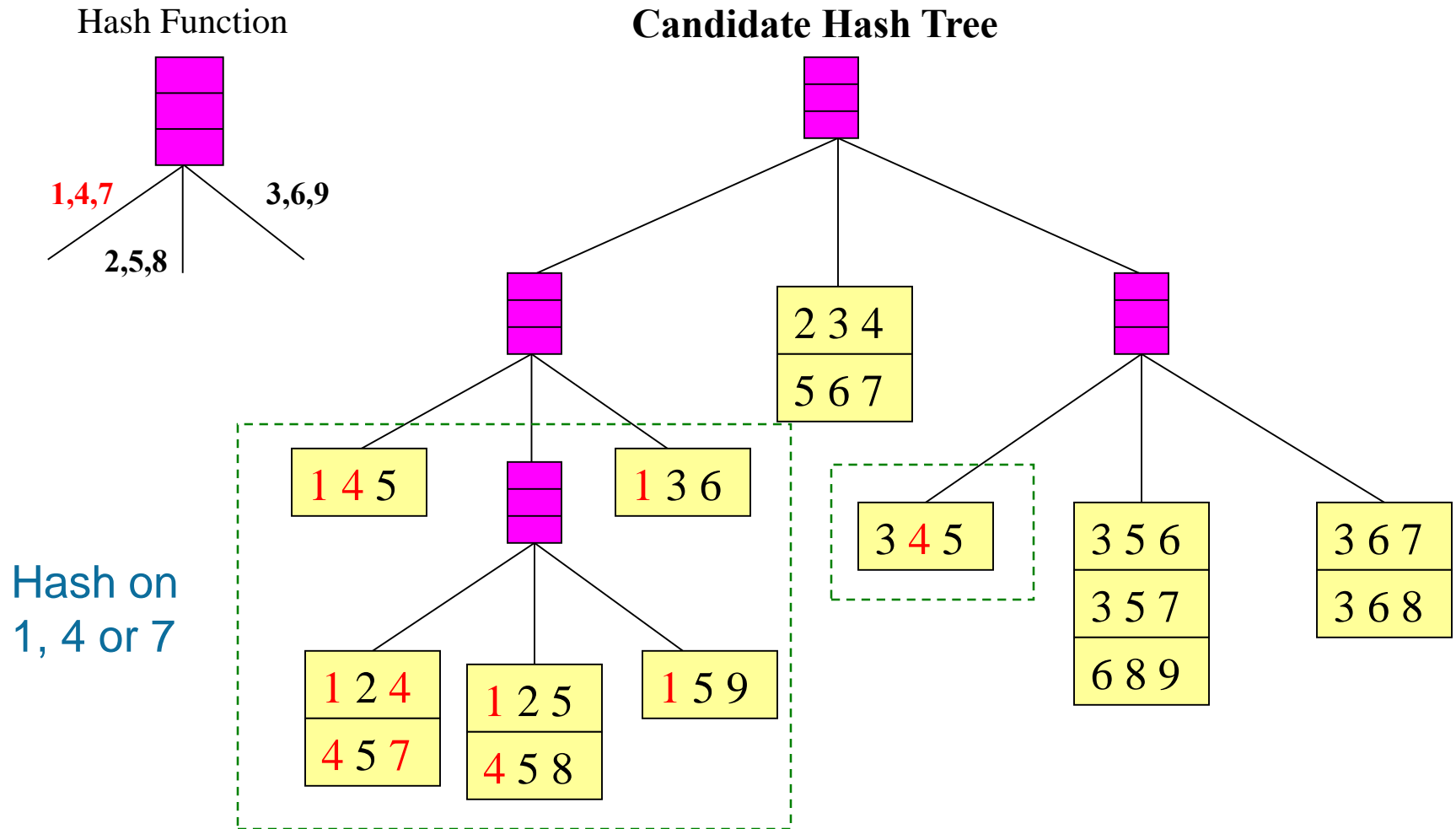
- Hash function
- Max leaf size: max number of itemsets stored in a leaf node (if number of candidate itemsets exceeds max leaf size, split the node)



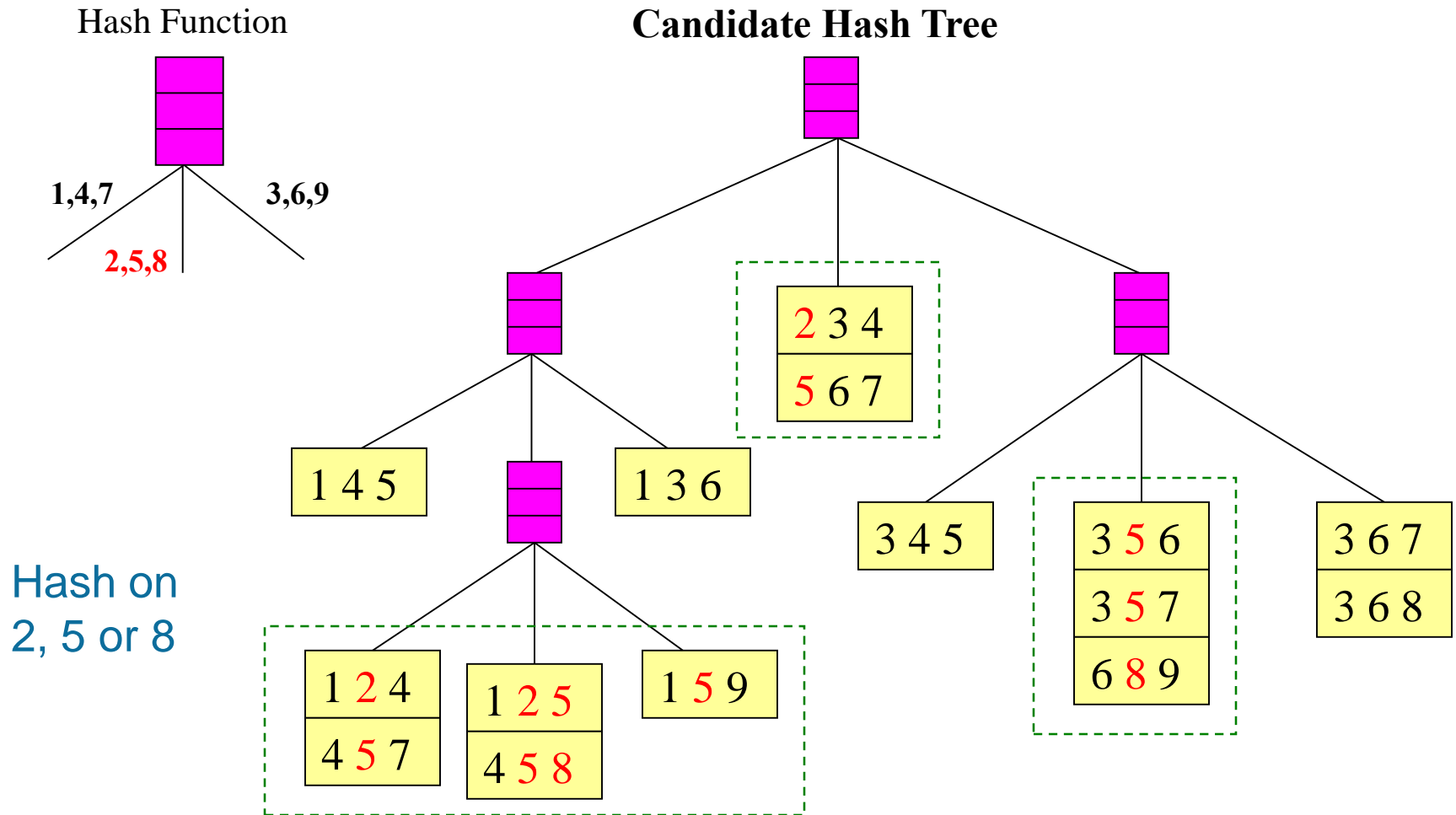
Ağacın her bir dahili düğümü, bir sonraki geçerli düğümün hangi dalının izleneceğini belirlemek için aşağıdaki **hash fonksiyonunu**,  $h(p) = p \bmod 3$  'ü kullanır.



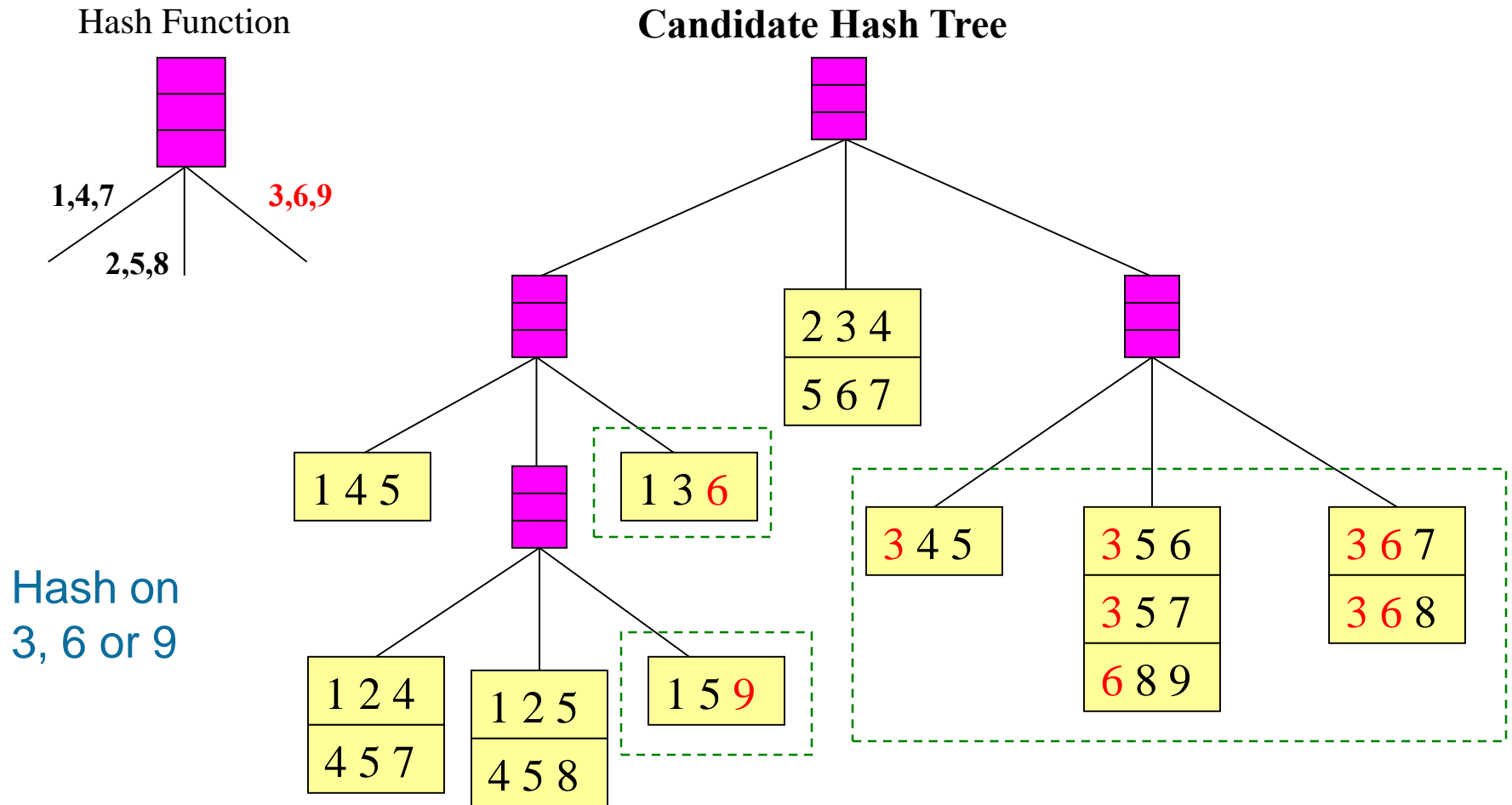
# Support Counting Using a Hash Tree



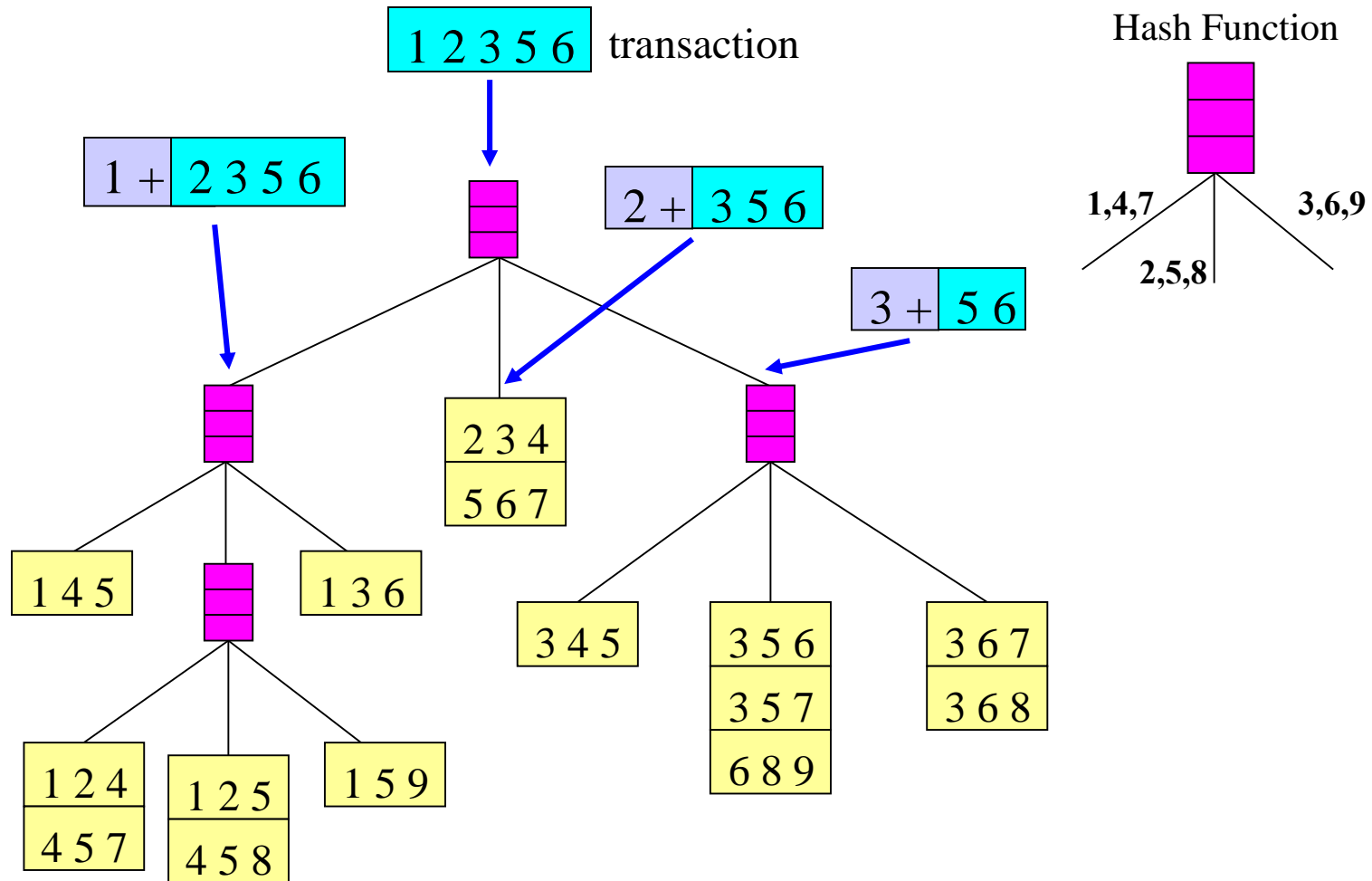
# Support Counting Using a Hash Tree



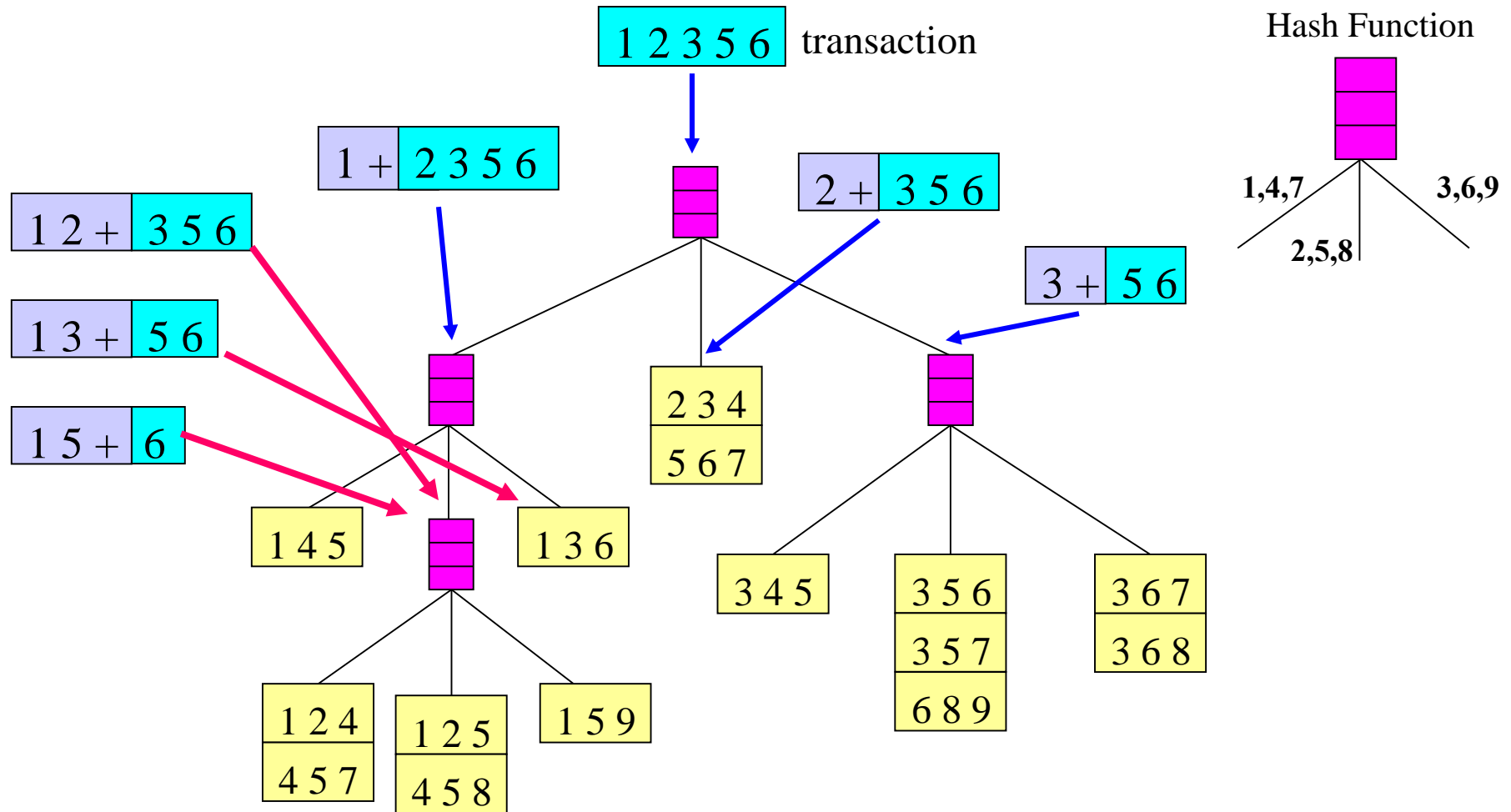
# Support Counting Using a Hash Tree



# Support Counting Using a Hash Tree

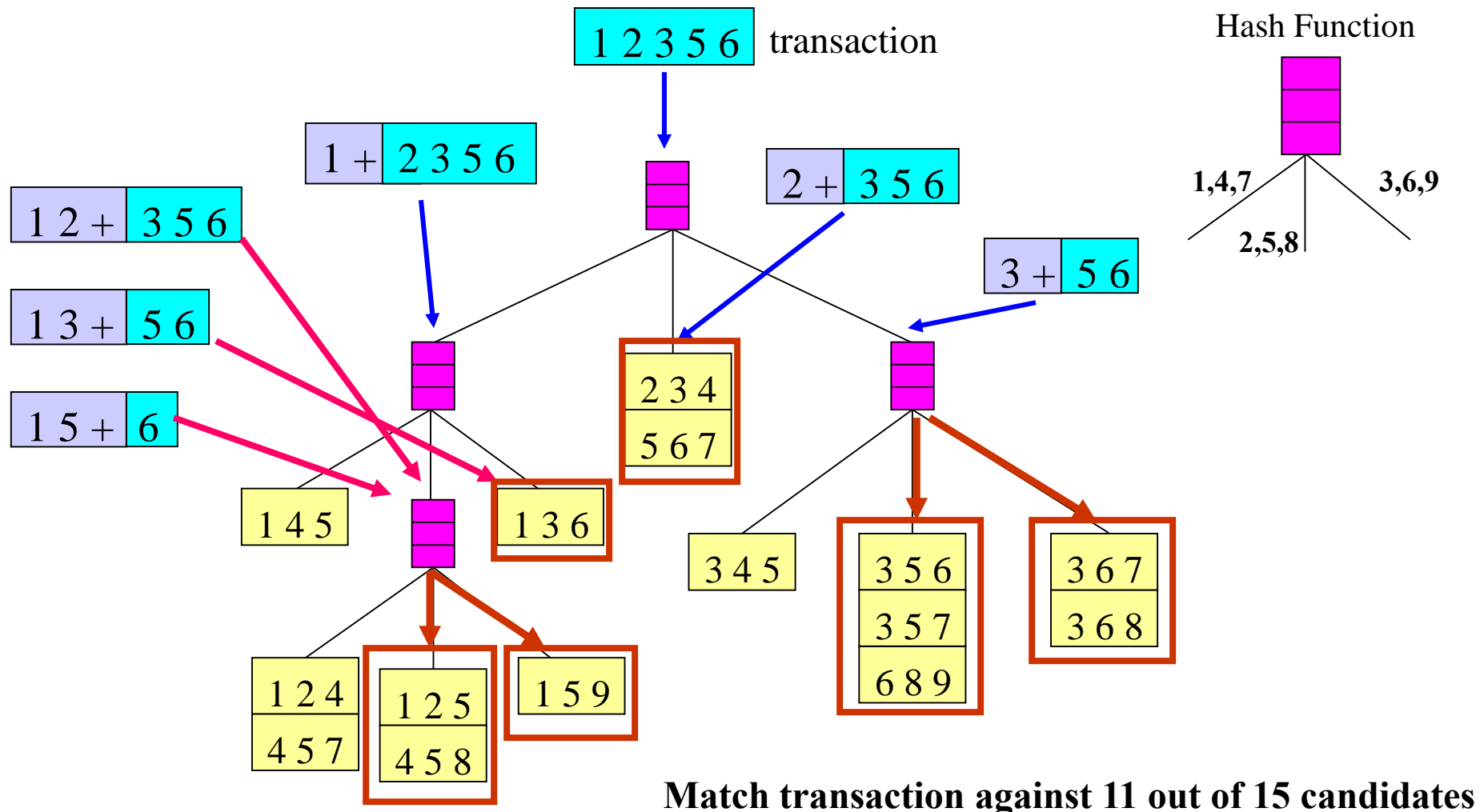


# Support Counting Using a Hash Tree





# Support Counting Using a Hash Tree



# Support Counting Using Hash Structure

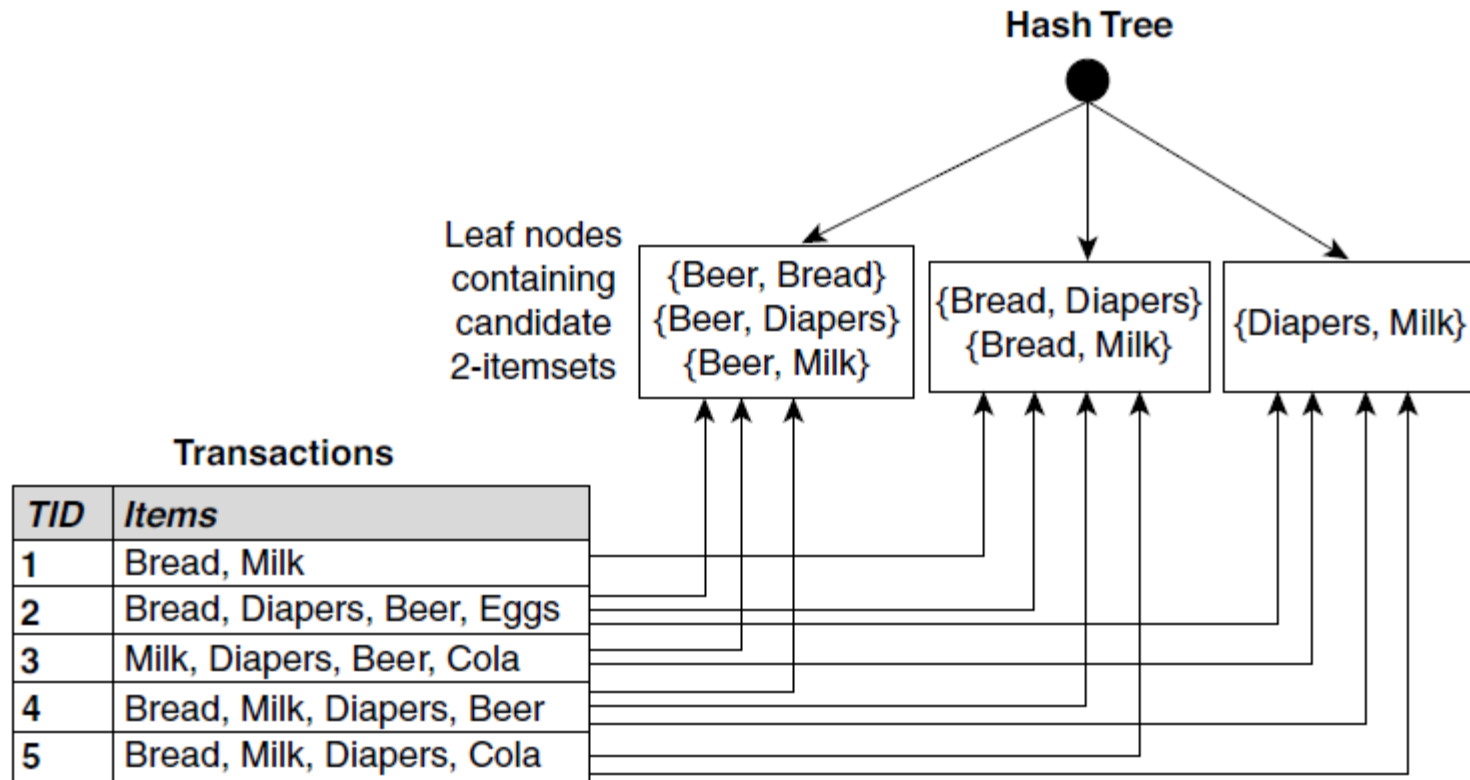


Figure 6.10. Counting the support of itemsets using hash structure.

# Rule Generation

*Bir Frequent L öge kümesi verildiğinde,  $f \rightarrow L - f$  kurallarının minimum güven gereksinimini karşılayacak şekilde tüm boş olmayan  $f \subset L$  alt kümelerini bulun*

- Given a frequent itemset  $L$ , find all non-empty subsets  $f \subset L$  such that  $f \rightarrow L - f$  satisfies the **minimum confidence** requirement

- If  $\{A,B,C,D\}$  is a frequent itemset, candidate rules:

$ABC \rightarrow D,$	$ABD \rightarrow C,$	$ACD \rightarrow B,$	$BCD \rightarrow A,$
$A \rightarrow BCD,$	$B \rightarrow ACD,$	$C \rightarrow ABD,$	$D \rightarrow ABC$
$AB \rightarrow CD,$	$AC \rightarrow BD,$	$AD \rightarrow BC,$	$BC \rightarrow AD,$
$BD \rightarrow AC,$	$CD \rightarrow AB,$		

- If  $|L| = k$ , then there are  $2^k - 2$  candidate association rules (ignoring  $L \rightarrow \emptyset$  and  $\emptyset \rightarrow L$ )

# Rule Generation

- In general, confidence does not have an anti-monotone property

$c(ABC \rightarrow D)$  can be larger or smaller than  $c(AB \rightarrow D)$

- But **confidence** of rules generated from the same itemset has an **anti-monotone** property

- E.g., Suppose  $\{A,B,C,D\}$  is a frequent 4-itemset:

$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$

$$\text{Confidence, } c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}.$$

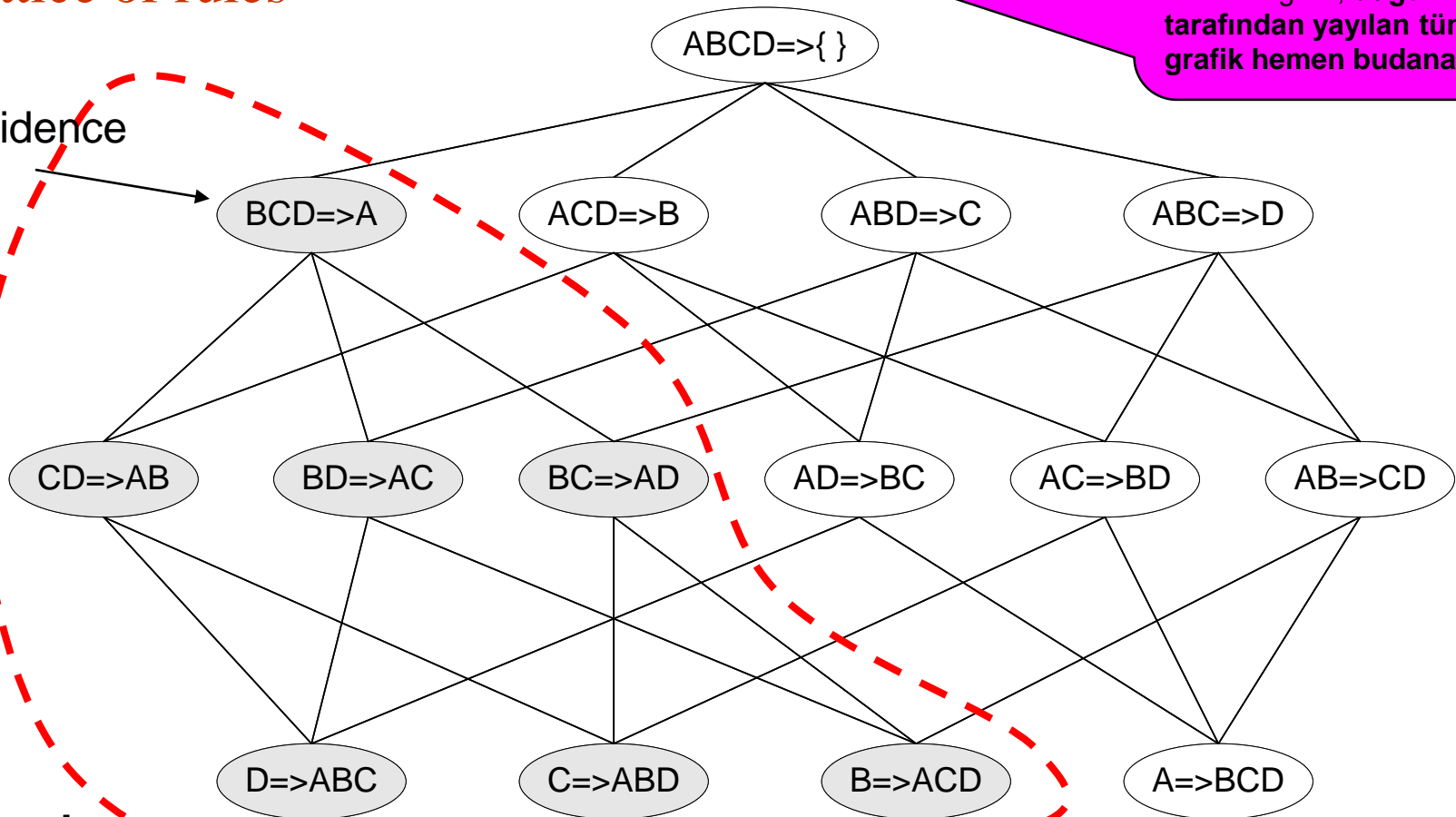
- Confidence is anti-monotone w.r.t. number of items on the RHS of the rule

# Rule Generation for Apriori Algorithm

## Lattice of rules

Low  
Confidence  
Rule

Ağdaki herhangi bir düğümün düşük güven değeri varsa, o zaman teoreme göre, **düğüm tarafından yayılan tüm alt grafik hemen budanabilir.**



Pruned  
Rules

$\{bcd\} \rightarrow \{a\}$  için güvenin düşük olduğunu varsayalım.  $\{cd\} \rightarrow \{ab\}$ ,  $\{bd\} \rightarrow \{ac\}$ ,  $\{bc\} \rightarrow \{ad\}$ ,  $\{d\} \rightarrow \{abc\}$ ,  $\{c\} \rightarrow \{abd\}$ , ve  $\{b\} \rightarrow \{acd\}$  dahil, sonuç kısmında a ögesini içeren tüm kurallar iptal edilebilir.

---

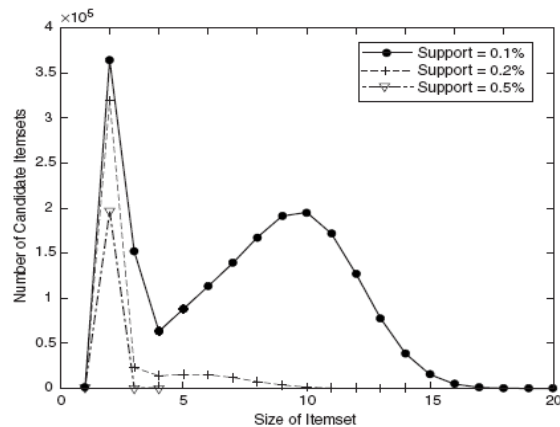
# **Association Analysis: Basic Concepts and Algorithms**

## Algorithms and Complexity

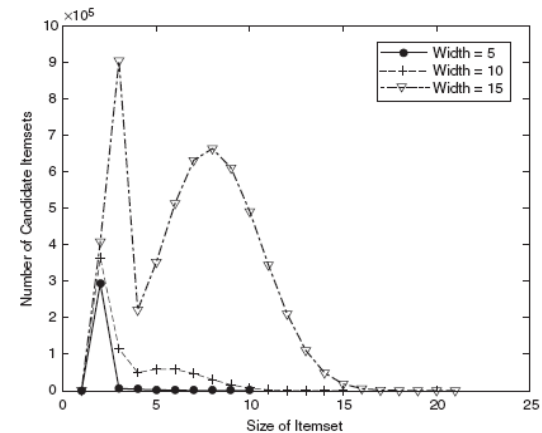
# Factors Affecting Complexity of Apriori

- Choice of minimum support threshold
  - destek eşiğini düşürmek daha fazla sayıda «frequent itemset»lere neden olur
  - Bu, aday sayısını ve frequent itemset'lerin maksimum uzunluğunu artırabilir
- Dimensionality (number of items) of the data set
  - her bir öğenin destek sayısını depolamak için daha fazla alana ihtiyaç vardır
  - frequent item'ların sayısı da artarsa, hem hesaplama hem de  $G / \mathcal{C}$  maliyetleri de artabilir
- Size of database
  - Apriori çoklu geçişler yaptığından, algoritmanın çalışma süresi transaction sayısı ile artabilir
- Average transaction width
  - daha yoğun veri kümeleri ile transaction genişliği artar
  - Bu, frequent itemset'lerin maksimum uzunluğunu ve hash ağacındaki gezinmeleri artırabilir (bir transaction'daki alt kümelerin sayısı, transaction genişliği ile artar)

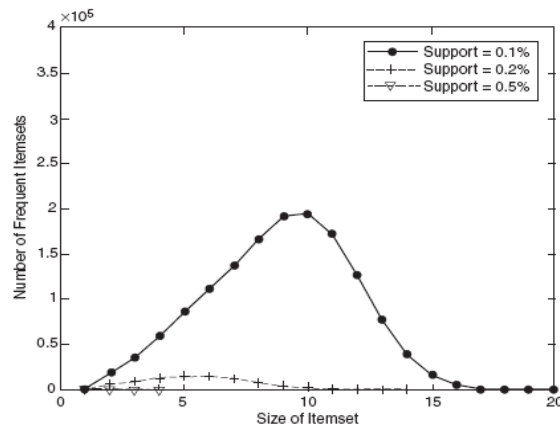
# Factors Affecting Complexity of Apriori



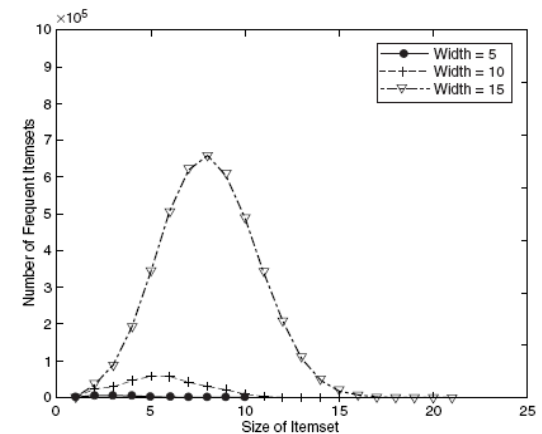
(a) Number of candidate itemsets.



(a) Number of candidate itemsets.



(b) Number of frequent itemsets.



(b) Number of Frequent Itemsets.

Figure 6.13. Effect of support threshold on the number of candidate and frequent itemsets.

Figure 6.14. Effect of average transaction width on the number of candidate and frequent itemsets.



# Compact Representation of Frequent Itemsets

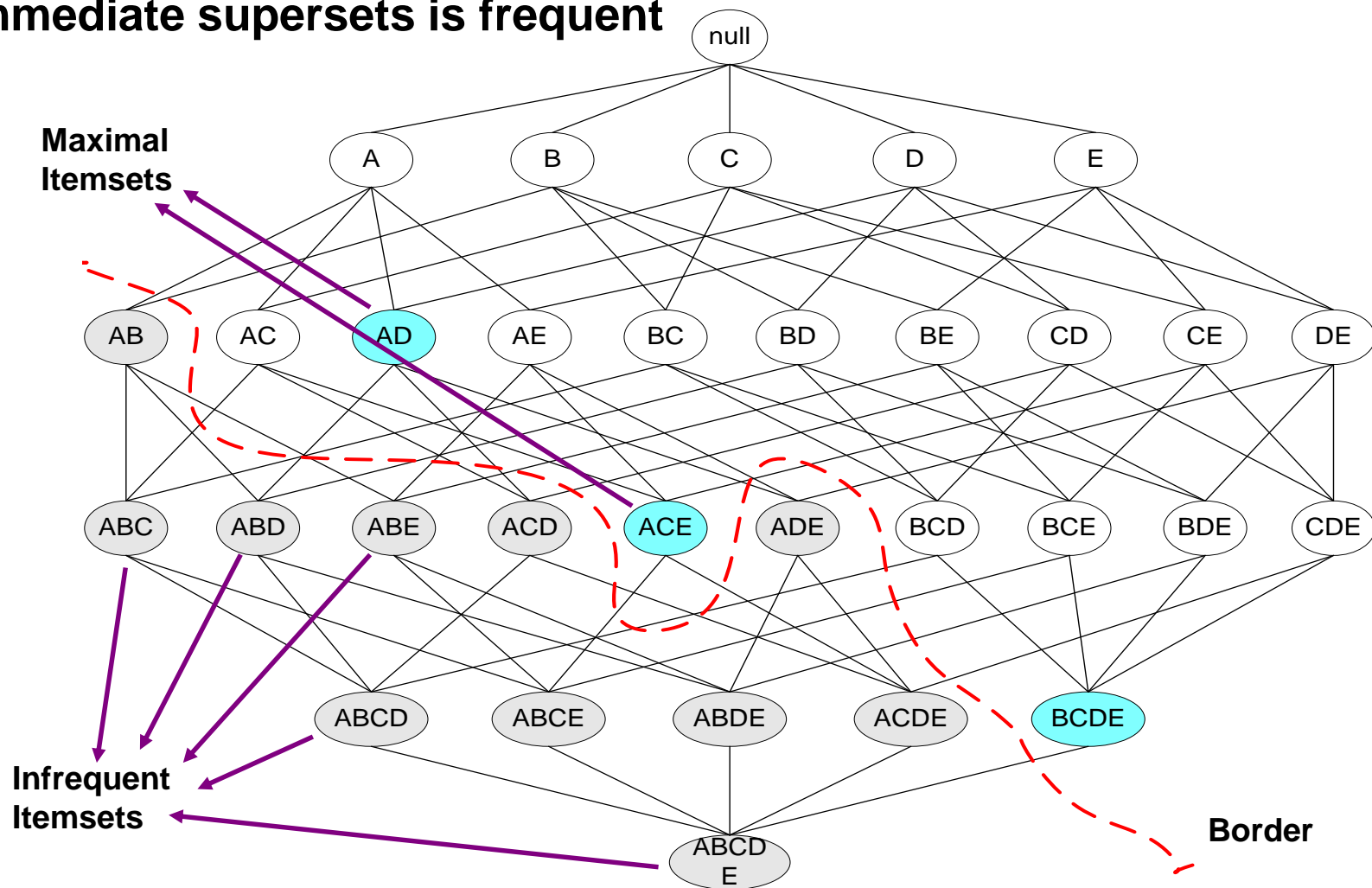
- Some itemsets are redundant because they have identical support as their supersets

TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1

- Number of frequent itemsets =  $3 \times \sum_{k=1}^{10} \binom{10}{k}$
- Need a compact representation

# Maximal Frequent Itemset

An itemset is maximal frequent if it is frequent and none of its immediate supersets is frequent



# What are the Maximal Frequent Itemsets in this Data?

TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1

Minimum support threshold = 5

# An illustrative example

Items

	A	B	C	D	E	F	G	H	I	J
1										
2										
3										
4										
5										
6										
7										
8										
9										
10										

Support threshold (by count) : 5  
Frequent itemsets: ?

# An illustrative example

		Items									
		A	B	C	D	E	F	G	H	I	J
Transactions	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Support threshold (by count) : 5  
Frequent itemsets: {F}

# An illustrative example

		Items									
		A	B	C	D	E	F	G	H	I	J
Transactions	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Support threshold (by count) : 5

Frequent itemsets: {F}

Support threshold (by count): 4

Frequent itemsets: ?

# An illustrative example

		Items									
		A	B	C	D	E	F	G	H	I	J
Transactions	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Support threshold (by count) : 5  
Frequent itemsets: {F}

Support threshold (by count): 4  
Frequent itemsets: {E}, {F}, {E,F}, {J}

# An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

**Support threshold (by count) : 5**  
Frequent itemsets: {F}

**Support threshold (by count): 4**  
Frequent itemsets: {E}, {F}, {E,F}, {J}

**Support threshold (by count): 3**  
Frequent itemsets: ?



# An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Support threshold (by count) : 5

Frequent itemsets: {F}

Support threshold (by count): 4

Frequent itemsets: {E}, {F}, {E,F}, {J}

Support threshold (by count): 3

Frequent itemsets:

All subsets of {C,D,E,F} + {J}

# An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

**Support threshold (by count) : 5**

Frequent itemsets: {F}

Maximal itemsets: ?

**Support threshold (by count): 4**

Frequent itemsets: {E}, {F}, {E,F}, {J}

Maximal itemsets: ?

**Support threshold (by count): 3**

Frequent itemsets:

All subsets of {C,D,E,F} + {J}

Maximal itemsets: ?

# An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

**Support threshold (by count) : 5**

Frequent itemsets: {F}

Maximal itemsets: {F}

**Support threshold (by count): 4**

Frequent itemsets: {E}, {F}, {E,F}, {J}

Maximal itemsets: ?

**Support threshold (by count): 3**

Frequent itemsets:

All subsets of {C,D,E,F} + {J}

Maximal itemsets: ?

# An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

**Support threshold (by count) : 5**

Frequent itemsets: {F}

Maximal itemsets: {F}

**Support threshold (by count): 4**

Frequent itemsets: {E}, {F}, {E,F}, {J}

Maximal itemsets: {E,F}, {J}

**Support threshold (by count): 3**

Frequent itemsets:

All subsets of {C,D,E,F} + {J}

Maximal itemsets: ?

# An illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

**Support threshold (by count) : 5**

Frequent itemsets: {F}

Maximal itemsets: {F}

**Support threshold (by count): 4**

Frequent itemsets: {E}, {F}, {E,F}, {J}

Maximal itemsets: {E,F}, {J}

**Support threshold (by count): 3**

Frequent itemsets:

All subsets of {C,D,E,F} + {J}

Maximal itemsets:

{C,D,E,F}, {J}

# Another illustrative example

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

**Support threshold (by count) : 5**

Maximal itemsets: {A}, {B}, {C}

**Support threshold (by count): 4**

Maximal itemsets: {A,B}, {A,C},{B,C}

**Support threshold (by count): 3**

Maximal itemsets: {A,B,C}

# Closed Itemset

- An itemset  $X$  is closed if none of its immediate supersets has exactly the same support count as  $X$ .
- $X$  is not closed if at least one of its immediate supersets has support count as  $X$ .

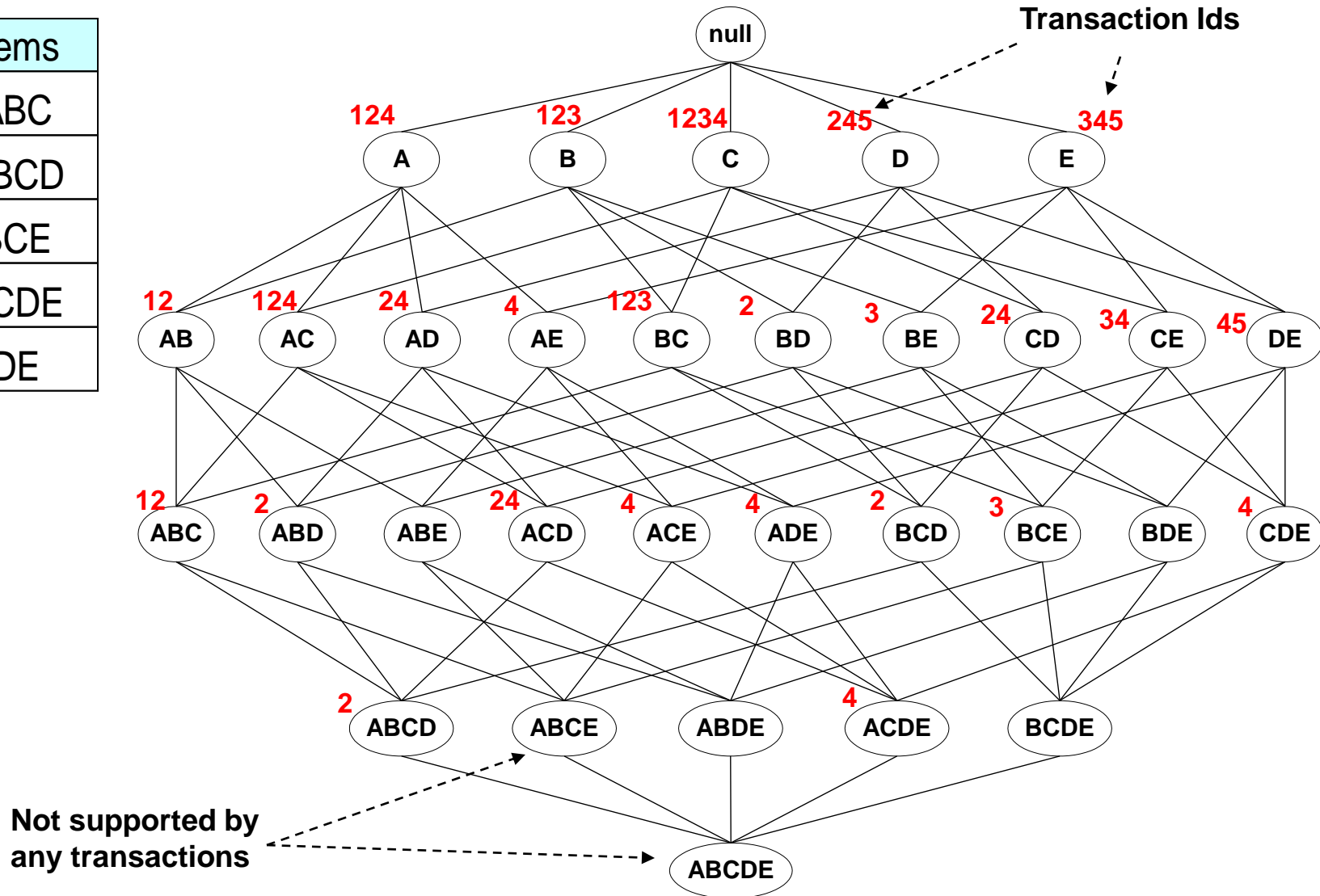
TID	Items
1	{A,B}
2	{B,C,D}
3	{A,B,C,D}
4	{A,B,D}
5	{A,B,C,D}

Itemset	Support
{A}	4
{B}	5
{C}	3
{D}	4
{A,B}	4
{A,C}	2
{A,D}	3
{B,C}	3
{B,D}	4
{C,D}	3

Itemset	Support
{A,B,C}	2
{A,B,D}	3
{A,C,D}	2
{B,C,D}	2
{A,B,C,D}	2

# Maximal vs Closed Itemsets

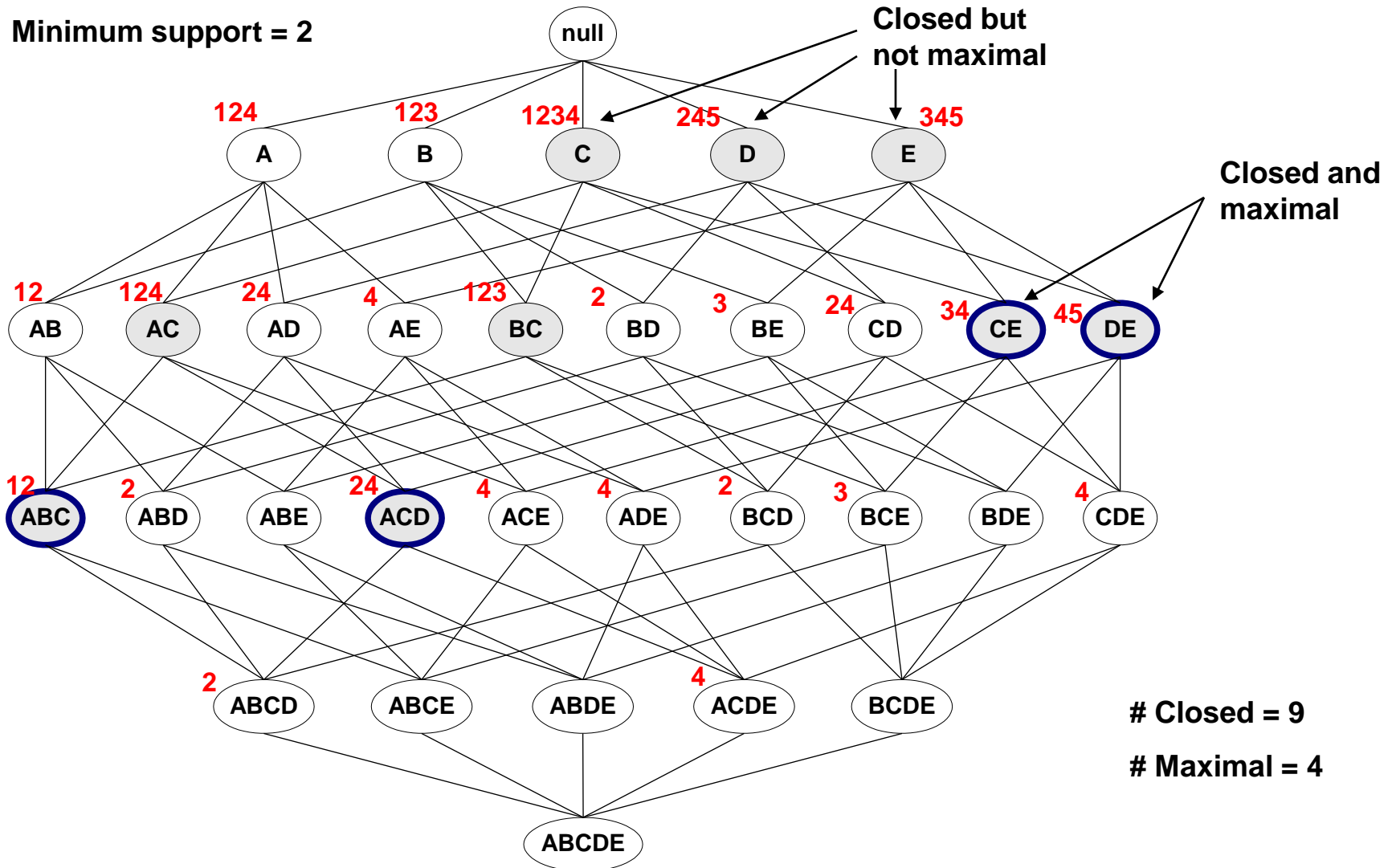
TID	Items
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE





# Maximal vs Closed Frequent Itemsets

Minimum support = 2



# What are the Closed Itemsets in this Data?

TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1

# Example 1

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Itemsets	Support (counts)	Closed itemsets
{C}	3	
{D}	2	
{C,D}	2	

# Example 1

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Itemsets	Support (counts)	Closed itemsets
<b>{C}</b>	<b>3</b>	✓
{D}	2	
<b>{C,D}</b>	<b>2</b>	✓

# Example 2

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Itemsets	Support (counts)	Closed itemsets
{C}	3	
{D}	2	
{E}	2	
{C,D}	2	
{C,E}	2	
{D,E}	2	
{C,D,E}	2	

# Example 2

		Items									
Transactions		A	B	C	D	E	F	G	H	I	J
	1										
	2										
	3										
	4										
	5										
	6										
	7										
	8										
	9										
	10										

Itemsets	Support (counts)	Closed itemsets
<b>{C}</b>	<b>3</b>	✓
{D}	2	
{E}	2	
{C,D}	2	
{C,E}	2	
{D,E}	2	
<b>{C,D,E}</b>	<b>2</b>	✓

# Example 3

Items

	A	B	C	D	E	F	G	H	I	J
1										
2										
3										
4										
5										
6										
7										
8										
9										
10										

Closed itemsets: {C,D,E,F}, {C,F}

# Example 4

Items

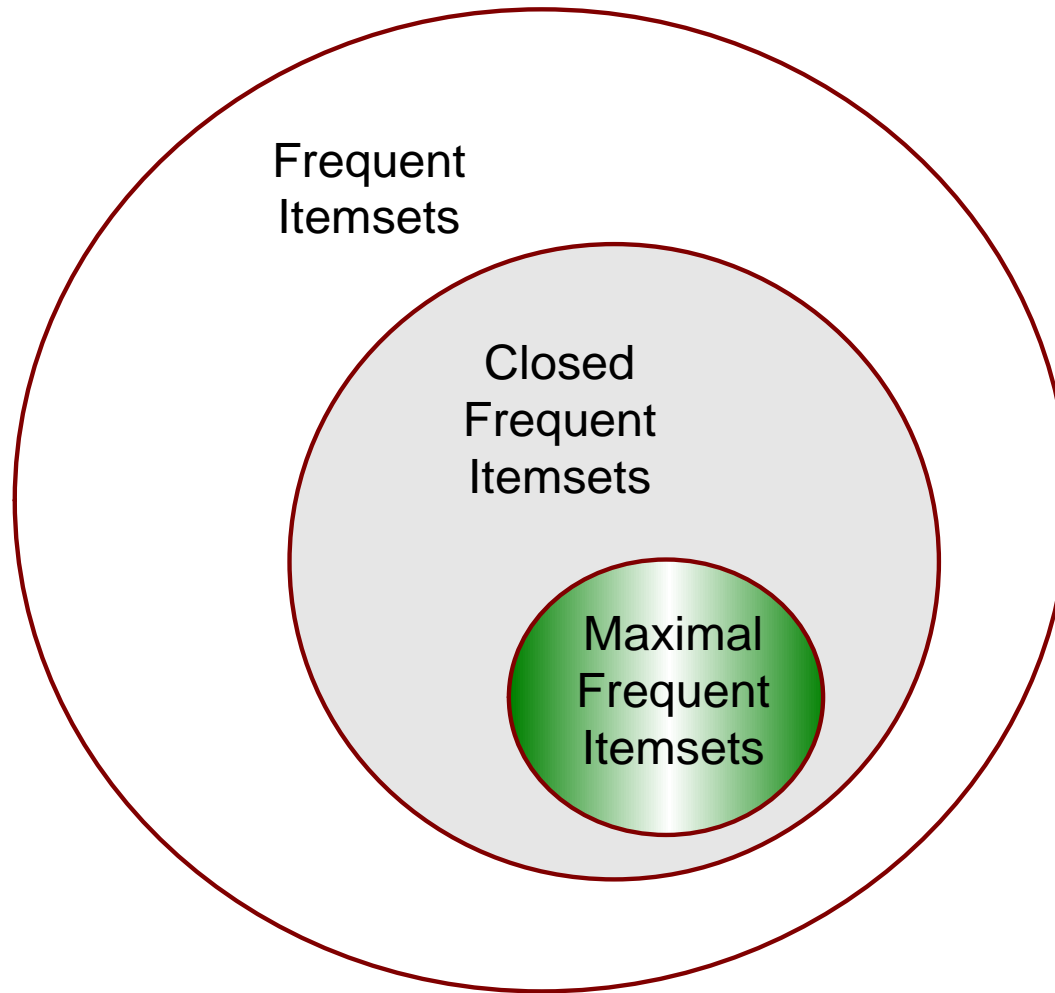
	A	B	C	D	E	F	G	H	I	J
1										
2										
3										
4										
5										
6										
7										
8										
9										
10										

Closed itemsets: {C,D,E,F}, {C}, {F}



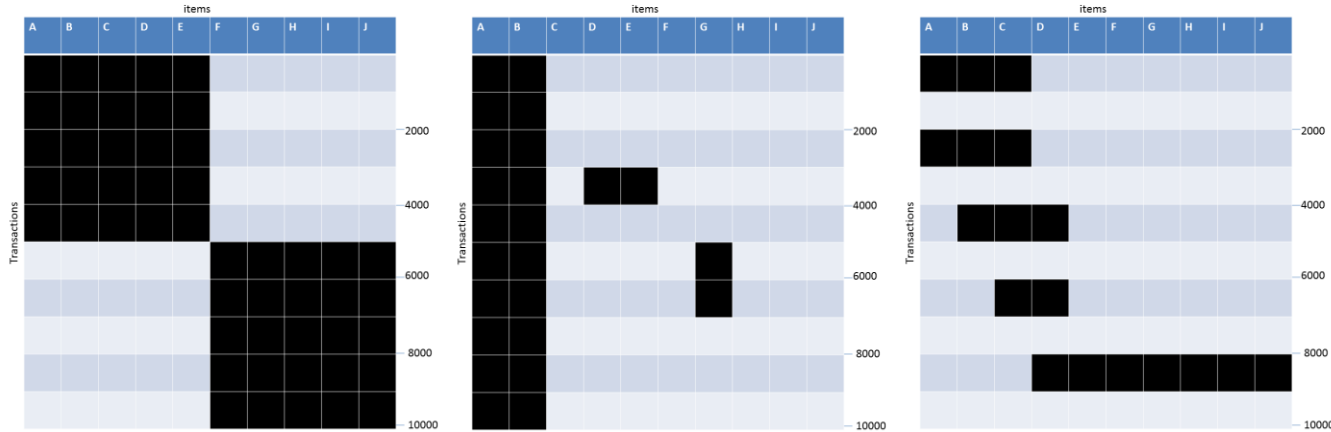
# Maximal vs Closed Itemsets

---



# Example question

- Aşağıdaki transaction veri kümeleri (koyu renkli hücreler, bir transaction'daki bir öğenin varlığını gösterir) ve %20'lik bir destek eşiği göz önüne alarak aşağıdaki soruları yanıtlayın



- Her bir veri kümesi için «*frequent itemset*» sayısı nedir? Hangi veri kümesi en çok sayıda «*frequent itemset*» üretir?
- En uzun (longest) frequent itemset'i hangi veri kümesi üretir?
- Hangi veri kümesi en yüksek maksimum desteğe sahip «*frequent itemset*»ler üretecektir?
- Hangi veri kümesi, çok çeşitli destek düzeylerine sahip öğeler içeren «*frequent itemset*»ler üretecektir (yani, % 20 ila % 70 arasında değişen, karışık destekli öğeler içeren öğe setleri)?
- Her veri kümesi için «*maximal frequent itemset*» sayısı nedir? Hangi veri kümesi en fazla sayıda *maximal frequent itemset* üretecektir?
- Her veri kümesi için «*closed frequent itemset*» sayısı nedir? Hangi veri kümesi en fazla sayıda *closed frequent itemset* üretir?

# Pattern Evaluation

---

- Birliktelik kuralı algoritmaları çok sayıda kural üretebilir
- Örüntüleri budamak / sıralamak için ***interestingness*** ölçütü kullanılabilir
  - Orijinal formülasyonda, destek ve güven (***support & confidence***) kullanılan tek ölçüttür

# Computing Interestingness Measure

- $X \rightarrow Y$  veya  $\{X, Y\}$  verildiğinde, «interestingness» hesaplamak için gereken bilgiler bir olasılık (*contingency*) tablosundan elde edilebilir

## Contingency table

	Y	$\overline{Y}$	
X	$f_{11}$	$f_{10}$	$f_{1+}$
$\overline{X}$	$f_{01}$	$f_{00}$	$f_{0+}$
	$f_{+1}$	$f_{+0}$	N

$f_{11}$ : support of X and Y

$f_{10}$ : support of  $\underline{X}$  and  $\overline{Y}$

$f_{01}$ : support of  $\overline{X}$  and  $\underline{Y}$

$f_{00}$ : support of  $\overline{X}$  and  $\overline{Y}$

Used to define various measures

- ◆ support, confidence, Gini, entropy, etc.

# Drawback of Confidence

Custo mers	Tea	Coffee	...
C1	0	1	...
C2	1	0	...
C3	1	1	...
C4	1	0	...
...			

	Coffee	<u>Coffee</u>	
Tea	15	5	20
<u>Tea</u>	75	5	80
	90	10	100

Association Rule: Tea  $\rightarrow$  Coffee

Confidence  $\cong P(\text{Coffee}|\text{Tea}) = 15/20 = 0.75$

Confidence  $> 50\%$ , meaning people who drink tea are more likely to drink coffee than not drink coffee

So rule seems reasonable

# Drawback of Confidence

	Coffee	<u>Coffee</u>	
Tea	15	5	20
<u>Tea</u>	75	5	80
	90	10	100

Association Rule: Tea  $\rightarrow$  Coffee

Confidence =  $P(\text{Coffee}|\text{Tea}) = 15/20 = 0.75$

but  $P(\text{Coffee}) = 0.9$ , which means knowing that a person drinks tea reduces the probability that the person drinks coffee!

$\Rightarrow$  Note that  $P(\text{Coffee}|\overline{\text{Tea}}) = 75/80 = 0.9375$

# Measure for Association Rules

---

- So, what kind of rules do we really want?
  - Confidence( $X \rightarrow Y$ ) should be sufficiently high
    - ◆ X satın alan kişilerin Y satın almamaktan çok Y satın almasını sağlamak için
  - Confidence( $X \rightarrow Y$ ) > support(Y)
    - ◆ Aksi takdirde, kural yanıltıcı olacaktır çünkü X öğesine sahip olmak, aynı transaction'da Y maddesine sahip olma şansını fiilen azaltır.
    - ◆ Bu kısıtı yakalayan herhangi bir ölçüt var mı?
      - Cevap: Evet. Çok sayıda var.

# Statistical Independence

---

- The criterion  
 $\text{confidence}(X \rightarrow Y) = \text{support}(Y)$

is equivalent to:

- $P(Y|X) = P(Y)$
- $P(X,Y) = P(X) \times P(Y)$

If  $P(X,Y) > P(X) \times P(Y)$  : X & Y are positively correlated

If  $P(X,Y) < P(X) \times P(Y)$  : X & Y are negatively correlated



# Measures that take into account statistical dependence

$$\text{Lift} = \frac{P(Y | X)}{P(Y)}$$

$$\text{Interest} = \frac{P(X, Y)}{P(X)P(Y)}$$

$$PS = P(X, Y) - P(X)P(Y)$$

$$\phi - \text{coefficient} = \frac{P(X, Y) - P(X)P(Y)}{\sqrt{P(X)[1 - P(X)]P(Y)[1 - P(Y)]}}$$

$$\text{Lift} = \frac{c(A \longrightarrow B)}{s(B)},$$

lift is used for rules while  
interest is used for itemsets

$$I(A, B) = \frac{s(A, B)}{s(A) \times s(B)}$$

For binary variables, lift is equivalent to another objective measure called **interest factor**,

# Example: Lift/Interest

	Coffee	<u>Coffee</u>	
Tea	15	5	20
<u>Tea</u>	75	5	80
	90	10	100

Association Rule: Tea  $\rightarrow$  Coffee

Confidence =  $P(\text{Coffee}|\text{Tea}) = 0.75$

but  $P(\text{Coffee}) = 0.9$

$\Rightarrow \text{Lift} = 0.75/0.9 = 0.8333 (< 1, \text{ therefore is negatively associated})$

So, is it enough to use confidence/lift for pruning?

$$\text{Lift} = \frac{c(A \rightarrow B)}{s(B)},$$

tea-coffee örneği, yüksek güvenilirlik kurallarının (high-confidence rules) bazen **yanıltıcı** olabileceğini gösterir çünkü **güven (confidence)** ölçüsü, kural sonuç kısmında ortaya çıkan öge setinin desteğini görmezden gelir. Bu sorunu çözmenin bir yolu, **lift** olarak bilinen bir metriği uygulamaktır:

# Lift or Interest

Contingency table

	Y	$\bar{Y}$	
X	$f_{11}$	$f_{10}$	$f_{1+}$
$\bar{X}$	$f_{01}$	$f_{00}$	$f_{0+}$
	$f_{+1}$	$f_{+0}$	N

	Y	$\bar{Y}$	
X	10	0	10
$\bar{X}$	0	90	90
	10	90	100

	Y	$\bar{Y}$	
X	90	0	90
$\bar{X}$	0	10	10
	90	10	100

$$Lift = \frac{0.1}{(0.1)(0.1)} = 10$$

$$I(A, B) = \frac{s(A, B)}{s(A) \times s(B)} = \frac{N f_{11}}{f_{1+} f_{+1}}.$$

$$I(A, B) \begin{cases} = 1, & \text{if } A \text{ and } B \text{ are independent;} \\ > 1, & \text{if } A \text{ and } B \text{ are positively correlated;} \\ < 1, & \text{if } A \text{ and } B \text{ are negatively correlated.} \end{cases}$$

$$Lift = \frac{0.9}{(0.9)(0.9)} = 1.11$$

**Statistical independence:**

**If  $P(X, Y) = P(X)P(Y) \Rightarrow Lift = 1$**

There are lots of measures proposed in the literature

#	Measure	Formula
1	$\phi$ -coefficient	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
2	Goodman-Kruskal's ( $\lambda$ )	$\frac{\sum_j \max_k P(A_j, B_k) + \sum_k \max_j P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
3	Odds ratio ( $\alpha$ )	$\frac{P(A,B)P(\bar{A},\bar{B})}{P(A,\bar{B})P(\bar{A},B)}$
4	Yule's $Q$	$\frac{P(A,B)P(\bar{A}\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A}\bar{B}) + P(A,\bar{B})P(\bar{A},B)} = \frac{\alpha-1}{\alpha+1}$
5	Yule's $Y$	$\frac{\sqrt{P(A,B)P(\bar{A}\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A}\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}} = \frac{\sqrt{\alpha}-1}{\sqrt{\alpha}+1}$
6	Kappa ( $\kappa$ )	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
7	Mutual Information ( $M$ )	$\frac{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\min(-\sum_i P(A_i) \log P(A_i), -\sum_j P(B_j) \log P(B_j))}$
8	J-Measure ( $J$ )	$\max \left( P(A, B) \log \left( \frac{P(B A)}{P(B)} \right) + P(\bar{A}\bar{B}) \log \left( \frac{P(\bar{B} \bar{A})}{P(\bar{B})} \right), \right. \\ \left. P(A, B) \log \left( \frac{P(A B)}{P(A)} \right) + P(\bar{A}\bar{B}) \log \left( \frac{P(\bar{A} \bar{B})}{P(\bar{A})} \right) \right)$
9	Gini index ( $G$ )	$\max \left( P(A)[P(B A)^2 + P(\bar{B} A)^2] + P(\bar{A})[P(B \bar{A})^2 + P(\bar{B} \bar{A})^2] \right. \\ \left. - P(B)^2 - P(\bar{B})^2, \right. \\ \left. P(B)[P(A B)^2 + P(\bar{A} B)^2] + P(\bar{B})[P(A \bar{B})^2 + P(\bar{A} \bar{B})^2] \right. \\ \left. - P(A)^2 - P(\bar{A})^2 \right)$
10	Support ( $s$ )	$P(A, B)$
11	Confidence ( $c$ )	$\max(P(B A), P(A B))$
12	Laplace ( $L$ )	$\max \left( \frac{NP(A,B)+1}{NP(A)+2}, \frac{NP(A,B)+1}{NP(B)+2} \right)$
13	Conviction ( $V$ )	$\max \left( \frac{P(A)P(\bar{B})}{P(\bar{A}B)}, \frac{P(B)P(\bar{A})}{P(\bar{B}A)} \right)$
14	Interest ( $I$ )	$\frac{P(A,B)}{P(A)P(B)}$
15	cosine ( $IS$ )	$\frac{P(A,B)}{\sqrt{P(A)P(B)}}$
16	Piatetsky-Shapiro's ( $PS$ )	$P(A, B) - P(A)P(B)$
17	Certainty factor ( $F$ )	$\max \left( \frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)} \right)$
18	Added Value ( $AV$ )	$\max(P(B A) - P(B), P(A B) - P(A))$
19	Collective strength ( $S$ )	$\frac{P(A,B) + P(\bar{A}\bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A,B) - P(\bar{A}\bar{B})}$
20	Jaccard ( $\zeta$ )	$\frac{P(A,B)}{P(A) + P(B) - P(A,B)}$
21	Klosgen ( $K$ )	$\sqrt{P(A, B)} \max(P(B A) - P(B), P(A B) - P(A))$

# Comparing Different Measures

10 examples of contingency tables:

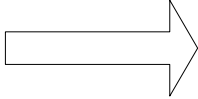
Example	$f_{11}$	$f_{10}$	$f_{01}$	$f_{00}$
E1	8123	83	424	1370
E2	8330	2	622	1046
E3	9481	94	127	298
E4	3954	3080	5	2961
E5	2886	1363	1320	4431
E6	1500	2000	500	6000
E7	4000	2000	1000	3000
E8	4000	2000	2000	2000
E9	1720	7121	5	1154
E10	61	2483	4	7452

Rankings of contingency tables using various measures:

#	$\phi$	$\lambda$	$\alpha$	$Q$	$Y$	$\kappa$	$M$	$J$	$G$	$s$	$c$	$L$	$V$	$I$	$IS$	$PS$	$F$	$AV$	$S$	$\zeta$	$K$
E1	1	1	3	3	3	1	2	2	1	3	5	5	4	6	2	2	4	6	1	2	5
E2	2	2	1	1	1	2	1	3	2	2	1	1	1	8	3	5	1	8	2	3	6
E3	3	3	4	4	4	3	3	8	7	1	4	4	6	10	1	8	6	10	3	1	10
E4	4	7	2	2	2	5	4	1	3	6	2	2	2	4	4	1	2	3	4	5	1
E5	5	4	8	8	8	4	7	5	4	7	9	9	9	3	6	3	9	4	5	6	3
E6	6	6	7	7	7	7	6	4	6	9	8	8	7	2	8	6	7	2	7	8	2
E7	7	5	9	9	9	6	8	6	5	4	7	7	8	5	5	4	8	5	6	4	4
E8	8	9	10	10	10	8	10	10	8	4	10	10	10	9	7	7	10	9	8	7	9
E9	9	9	5	5	5	9	9	7	9	8	3	3	3	7	9	9	3	7	9	9	8
E10	10	8	6	6	6	10	5	9	10	10	6	6	5	1	10	10	5	1	10	10	7

# Property under Variable Permutation

	B	$\bar{B}$
A	p	q
$\bar{A}$	r	s



	A	$\bar{A}$
B	p	r
$\bar{B}$	q	s

Does  $M(A,B) = M(B,A)$ ?

Symmetric measures:

- ◆ support, lift, collective strength, cosine, Jaccard, etc

Asymmetric measures:

- ◆ confidence, conviction, Laplace, J-measure, etc

# Property under Row/Column Scaling

Grade-Gender Example (Mosteller, 1968):

	Female	Male	
High	2	3	5
Low	1	4	5
	3	7	10

	Female	Male	
High	4	30	34
Low	2	40	42
	6	70	76

↓  
2x

↓  
10x

Mosteller:

Underlying association should be independent of the relative number of male and female students in the samples

# Property under Inversion Operation

	A	B	C	D	E	F
Transaction 1 →	1	0	0	1	0	0
■	0	0	1	1	1	0
■	0	0	1	1	1	0
■	0	1	1	0	1	1
■	0	0	1	1	1	0
■	0	0	1	1	1	0
■	0	0	1	1	1	0
Transaction N →	1	0	0	1	0	0

(a)
(b)
(c)

Effect of the inversion operation. The vectors *C* and *E* are inversions of vector *A*, while the vector *D* is an inversion of vectors *B* and *F*.



# Example: $\phi$ -Coefficient

$$\phi - \text{coefficient} = \frac{P(X,Y) - P(X)P(Y)}{\sqrt{P(X)[1 - P(X)]P(Y)[1 - P(Y)]}}$$

- $\phi$ -coefficient is analogous to correlation coefficient for continuous variables

	Y	$\overline{Y}$	
X	60	10	70
$\overline{X}$	10	20	30
	70	30	100

	Y	$\overline{Y}$	
X	20	10	30
$\overline{X}$	10	60	70
	30	70	100

$$\begin{aligned}\phi &= \frac{0.6 - 0.7 \times 0.7}{\sqrt{0.7 \times 0.3 \times 0.7 \times 0.3}} \\ &= 0.5238\end{aligned}$$

$$\begin{aligned}\phi &= \frac{0.2 - 0.3 \times 0.3}{\sqrt{0.7 \times 0.3 \times 0.7 \times 0.3}} \\ &= 0.5238\end{aligned}$$

**$\phi$  Coefficient is the same for both tables**

# Different Measures have Different Properties

Symbol	Measure	Inversion	Null Addition	Scaling
$\phi$	$\phi$ -coefficient	Yes	No	No
$\alpha$	odds ratio	Yes	No	Yes
$\kappa$	Cohen's	Yes	No	No
$I$	Interest	No	No	No
$IS$	Cosine	No	Yes	No
$PS$	Piatetsky-Shapiro's	Yes	No	No
$S$	Collective strength	Yes	No	No
$\zeta$	Jaccard	No	Yes	No
$h$	All-confidence	No	No	No
$s$	Support	No	No	No