

FREQUENT PATTERN MINING

BASED ON: INTRODUCTION TO DATA MINING
BY TAN, STEINBACH, KUMAR

ITEM SETS

A NEW TYPE OF DATA

- Some notation:
 - All possible items: $I = \{i_1, \dots, i_d\}$
 - Database: T is a bag of transactions t_1, \dots, t_N
 - Transaction $t = (TID, X)$
 - transaction identifier: TID
 - item set $X \subseteq I$

ITEMSET DATABASE	
<i>TID</i>	<i>Items</i>
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

DEFINITIONS

ITEMSET DATABASE	
<i>TID</i>	<i>Items</i>
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

DEFINITIONS

- **Itemset**
 - A collection of one or more items
 - Example: {Milk, Bread, Diaper}
 - k-itemset
 - An itemset that contains k items

ITEMSET DATABASE	
<i>TID</i>	<i>Items</i>
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

DEFINITIONS

- **Itemset**
 - A collection of one or more items
 - Example: {Milk, Bread, Diaper}
 - k-itemset
 - An itemset that contains k items
- **Support count $\sigma()$**
 - Frequency of occurrence of an itemset
 - E.g. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

DEFINITIONS

- **Itemset**

- A collection of one or more items
 - Example: {Milk, Bread, Diaper}
- k-itemset
 - An itemset that contains k items

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

- Frequency of occurrence of an itemset
- E.g. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

- **Support $s()$**

- Fraction of transactions that contain an itemset
- E.g. $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

DEFINITIONS

- **Itemset**

- A collection of one or more items
 - Example: {Milk, Bread, Diaper}
- k-itemset
 - An itemset that contains k items

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

- Frequency of occurrence of an itemset
- E.g. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

- **Support $s()$**

- Fraction of transactions that contain an itemset
- E.g. $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$

- **Frequent Itemset**

- An itemset whose support is greater than or equal to a *minsup* threshold

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

ASSOCIATION RULE MINING

- Given a set of transactions, find rules that predict the occurrence of an item based on occurrences of other items in the transaction

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

Example of Association Rules

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

\rightarrow means co-occurrence

DEFINITIONS

ITEMSET DATABASE	
<i>TID</i>	<i>Items</i>
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

DEFINITIONS

- Association Rule
 - Expression of the form $X \rightarrow Y$
(X and Y are itemsets)

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

DEFINITIONS

- Association Rule

- Expression of the form $X \rightarrow Y$
(X and Y are itemsets)
- Example:
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

- Rule Evaluation Measures

- Support (s)
 - Fraction of transactions that contain both X and Y
- Confidence (c)

DEFINITIONS

- Association Rule

- Expression of the form $X \rightarrow Y$
(X and Y are itemsets)
- Example:
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

- Rule Evaluation Measures

- Support (s)
 - Fraction of transactions that contain both X and Y
- Confidence (c)
 - Measures how often items in Y appear in transactions that contain X

DEFINITIONS

- Association Rule

- Expression of the form $X \rightarrow Y$
(X and Y are itemsets)
- Example:
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

- Rule Evaluation Measures

- Support (s)
 - Fraction of transactions that contain both X and Y
- Confidence (c)
 - Measures how often items in Y appear in transactions that contain X

Example:

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

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

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

ASSOCIATION RULE MINING TASK

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

Example of Rules:

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

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

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

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

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

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

ASSOCIATION RULE MINING TASK

- Given set of transactions T, find all rules having
 - support \geq *minsup* threshold
 - confidence \geq *minconf* threshold

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

Example of Rules:

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

MINING ASSOCIATION RULES

- 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
 - \Rightarrow **Computationally prohibitive (combinatorial explosion)**
- Can we do better?

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

{Milk,Diaper} \rightarrow {Beer} (s=0.4, c=0.67)

{Milk,Beer} \rightarrow {Diaper} (s=0.4, c=1.0)

{Diaper,Beer} \rightarrow {Milk} (s=0.4, c=0.67)

{Beer} \rightarrow {Milk,Diaper} (s=0.4, c=0.67)

{Diaper} \rightarrow {Milk,Beer} (s=0.4, c=0.5)

{Milk} \rightarrow {Diaper,Beer} (s=0.4, c=0.5)

MINING ASSOCIATION RULES

Observations:

- All rules are binary partitions of the same itemset:
{Milk, Diaper, Beer}
- Rules originating from the same itemset have **identical support** but can have **different confidence**
- Thus, we may **decouple** the support and confidence requirements

ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

{Milk,Diaper} → {Beer} (s=0.4, c=0.67)

{Milk,Beer} → {Diaper} (s=0.4, c=1.0)

{Diaper,Beer} → {Milk} (s=0.4, c=0.67)

{Beer} → {Milk,Diaper} (s=0.4, c=0.67)

{Diaper} → {Milk,Beer} (s=0.4, c=0.5)

{Milk} → {Diaper,Beer} (s=0.4, c=0.5)

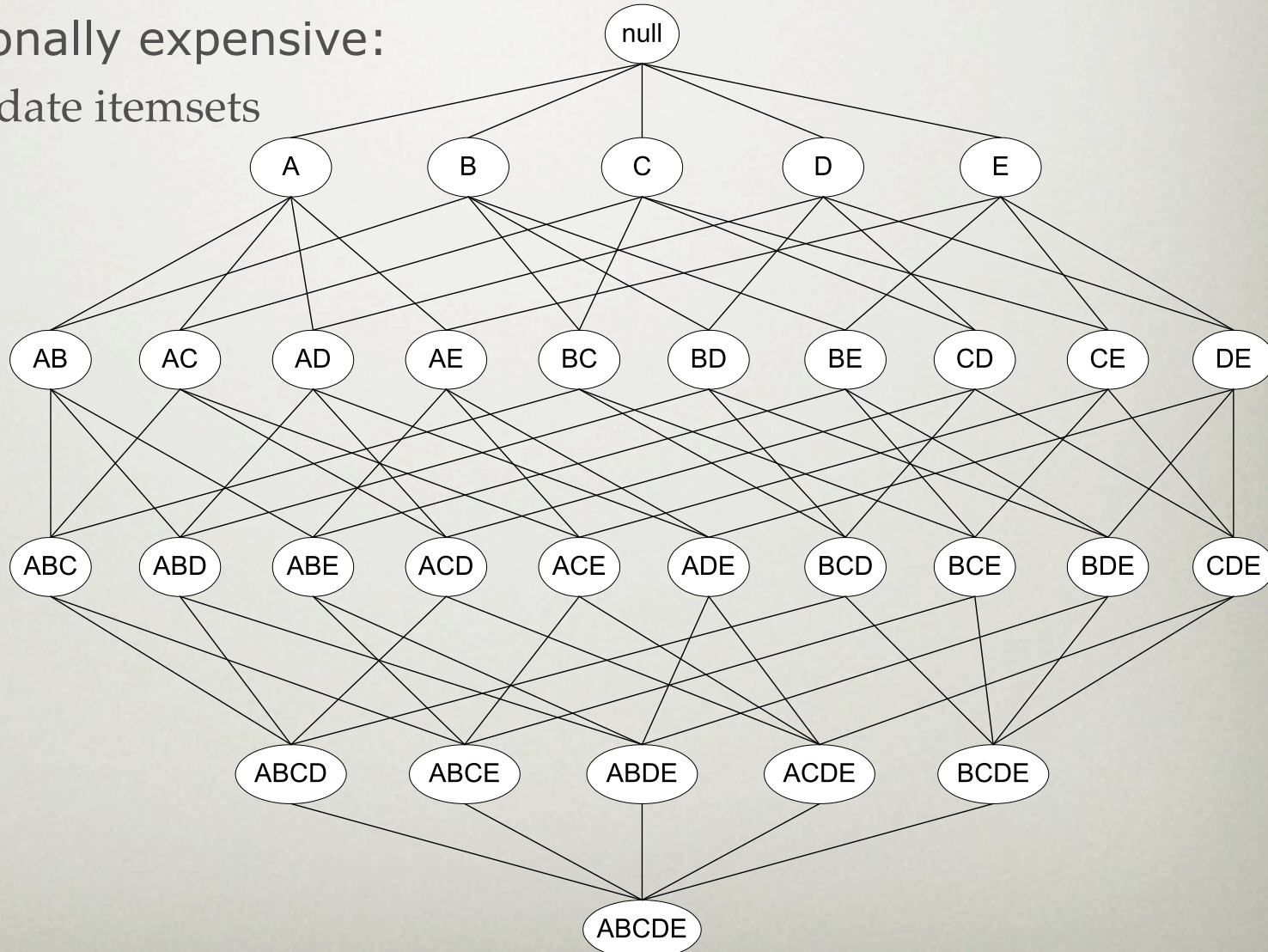
MINING ASSOCIATION RULES

- Two-step approach:
 - **Frequent Itemset Generation (count support)**
 - Generate *all* itemsets whose support $\geq \text{minsup}$
 - **Rule Generation (count confidence)**
 - Generate high confidence rules from each frequent itemset, where each rule is a *binary partitioning* of a frequent itemset

FREQUENT ITEMSET GENERATION

Still computationally expensive:

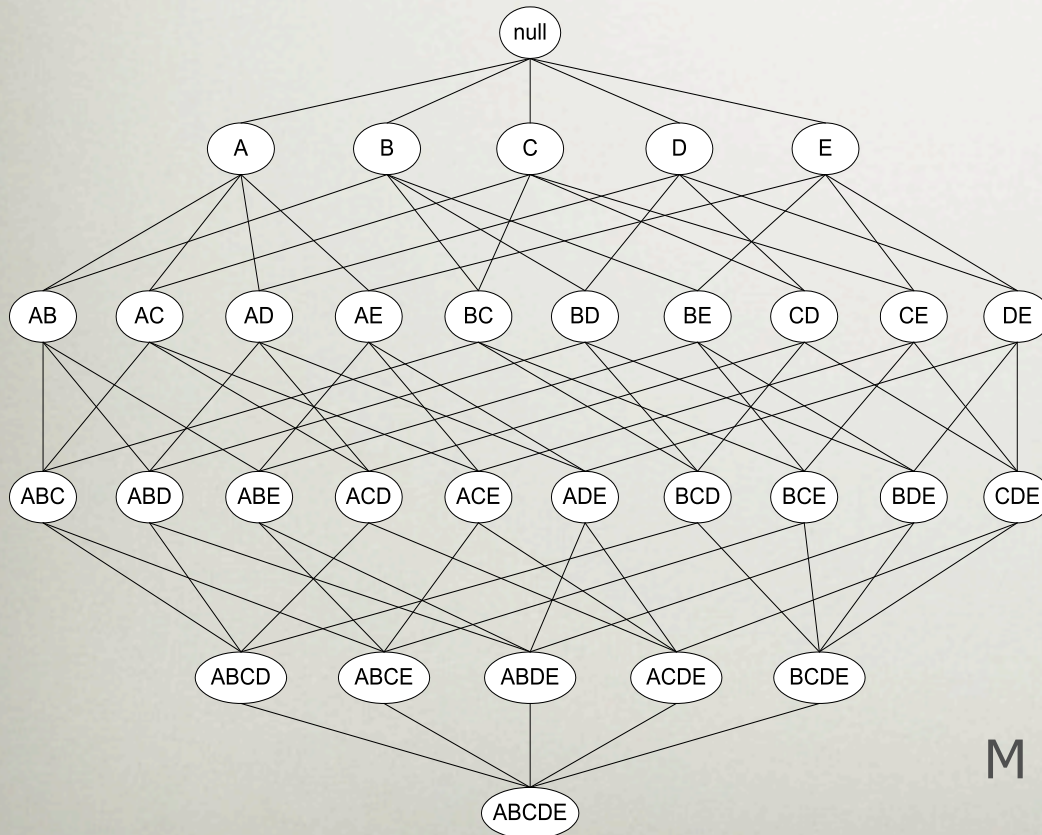
d items: 2^d candidate itemsets



FREQUENT ITEMSET GENERATION

Brute force:

Match every itemset against transaction database

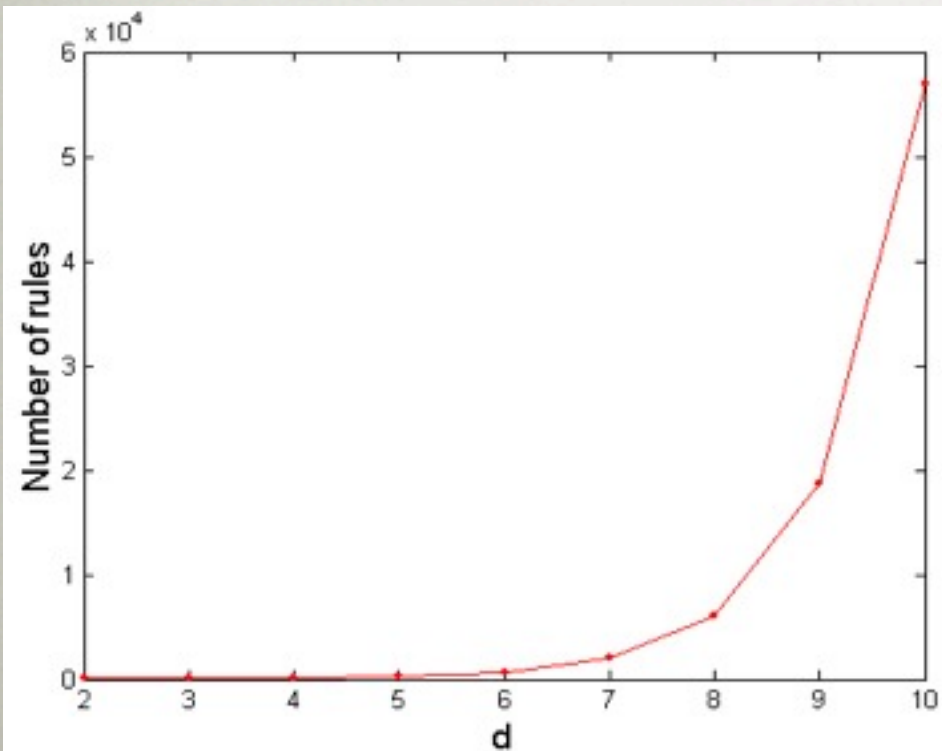


ITEMSET DATABASE	
TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

M (2^d) itemsets, N transactions:
Complexity $\sim O(NMw)$!

COMPUTATIONAL COMPLEXITY

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



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

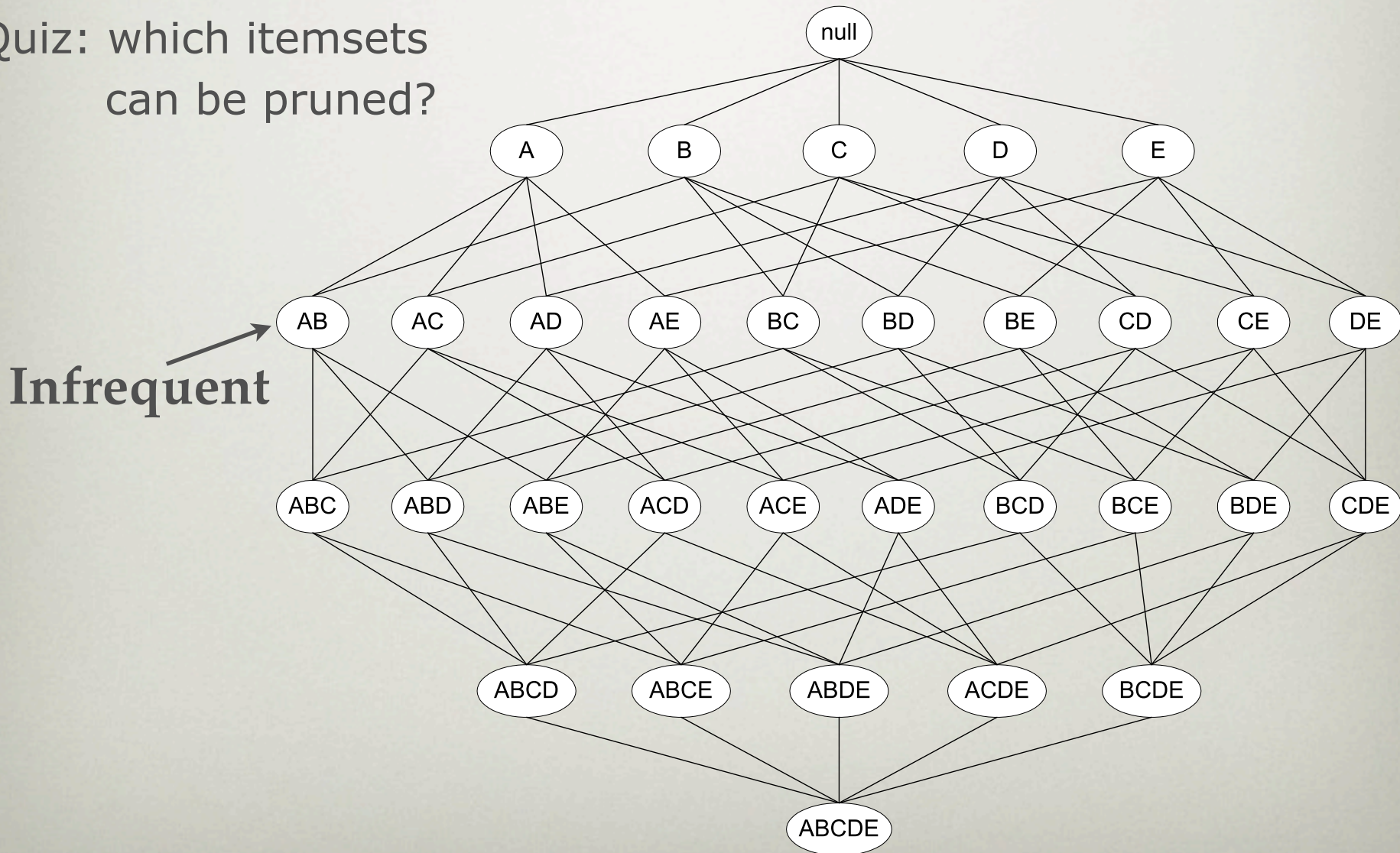
If $d=6$, $R = 602$ rules

FREQUENT ITEMSET GENERATION STRATEGIES

- Need to:
 - Reduce the **number of candidates** (M)
 - *Pruning*
 - Reduce the **number of comparisons** (NM)
 - *Efficient support counting*

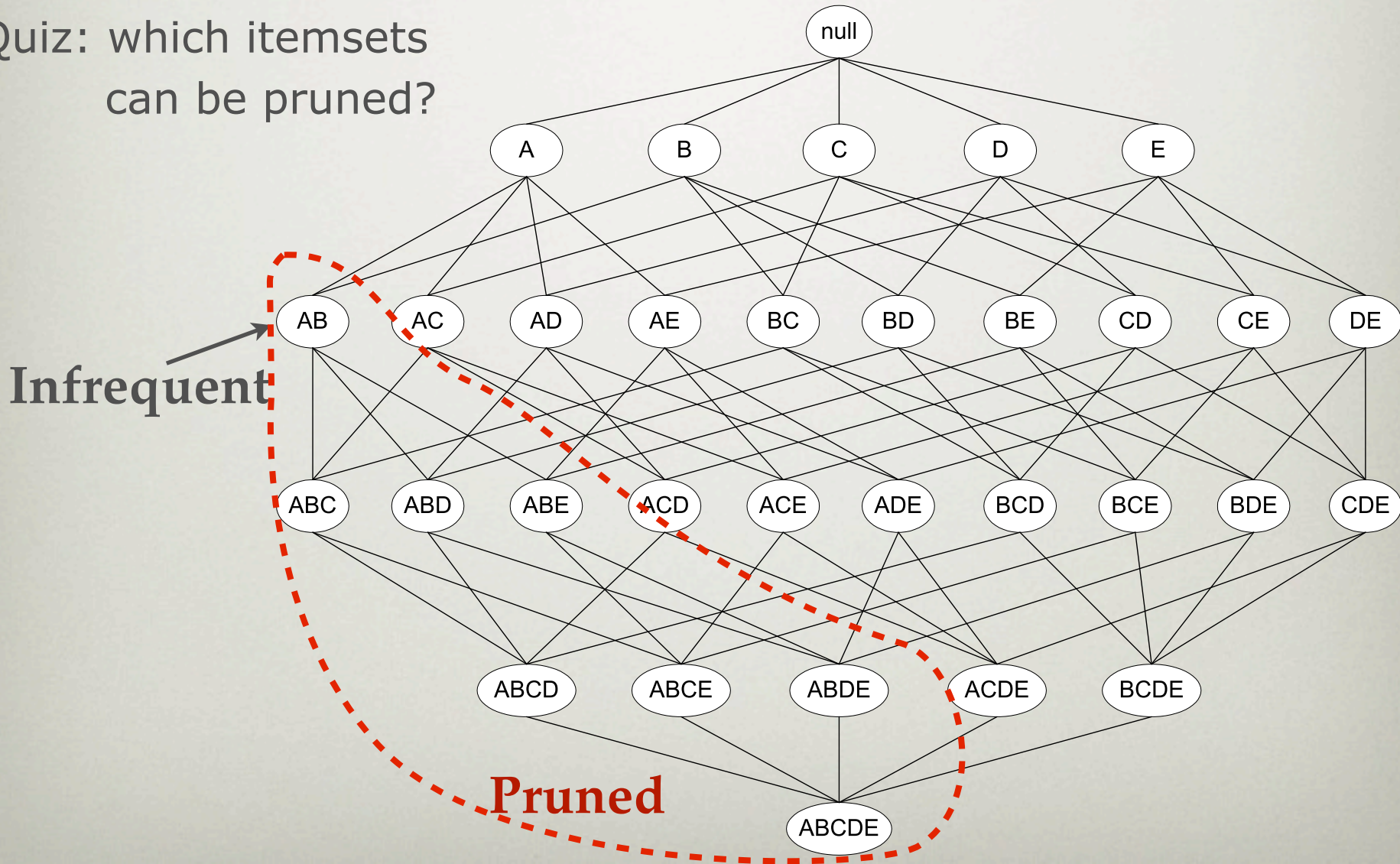
1) REDUCE NUMBER OF CANDIDATES

Quiz: which itemsets can be pruned?



1) REDUCE NUMBER OF CANDIDATES

Quiz: which itemsets can be pruned?



APRIORI PRINCIPLE

- If an itemset is frequent, then all of its subsets must also be frequent

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

- Support of itemset never exceeds the support of its subsets
- *Anti-monotonicity*

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



Triplets (3-itemsets)

Itemset	Count
{Bread,Milk,Diaper}	2



If every subset is considered,

$${}^6C_1 + {}^6C_2 + {}^6C_3 = 41$$

With support-based pruning,

$$6 + 6 + 1 = 13$$

APRIORI ALGORITHM

- Method:
 - Let $k=1$
 - Generate frequent itemsets of length 1
 - Repeat until no new frequent itemsets are identified
 - **Generate length $(k+1)$ candidate itemsets** from length k frequent itemsets
 - **Prune** candidate itemsets containing subsets of length k that are infrequent
 - **Count the support** of each candidate by scanning the DB
 - **Eliminate infrequent candidates**, leaving only those that are frequent

APRIORI ALGORITHM: LEVEL 1

candidate	support	Frequent?
A	6	Y
B	7	Y
C	6	Y
D	2	Y
E	2	Y

Minimum Support = 2

Level 2 candidates?

TID	Items
1	ABE
2	BD
3	BC
4	ABD
5	AC
6	BC
7	AC
8	ABCE
9	ABC

APRIORI ALGORITHM: LEVEL 2

candidate	support	Frequent?
AB	4	Y
AC	4	Y
AD	1	N
AE	2	Y
BC	4	Y
BD	2	Y
BE	2	Y
CD	0	N
CE	1	N
DE	0	N

TID	Items
1	ABE
2	BD
3	BC
4	ABD
5	AC
6	BC
7	AC
8	ABCE
9	ABC

Level 3?

APRIORI ALGORITHM: LEVEL 3

candidate	support	Frequent?
ABC	2	Y
ABE	2	Y

TID	Items
1	ABE
2	BD
3	BC
4	ABD
5	AC
6	BC
7	AC
8	ABCE
9	ABC

Level 4?

2) REDUCE NUMBER OF COMPARISONS

candidate	support	Frequent
AB	4	Y
AC	4	Y
AD	1	N
AE	2	Y
BC	4	Y
BD	2	Y
BE	2	Y
CD	0	N
CE	1	N
DE	0	N

Look-up

TID	Items
1	ABE
2	BD
3	BC
4	ABD
5	AC
6	BC
7	AC
8	ABCE
9	ABC

Quiz: efficient data structure to look up items in a list?

2) REDUCE NUMBER OF COMPARISONS

candidate	support	Frequent?
AB	4	Y
AC	4	Y
AD	1	N
AE	2	Y
BC	4	Y
BD	2	Y
BE	2	Y
CD	0	N
CE	1	N
DE	0	N

Look-up

TID	Items
1	ABE
2	BD
3	BC
4	ABD
5	AC
6	BC
7	AC
8	ABC
9	ABC

Hash maps!

Key	Value
A,C,E	AB, AC, AD, AE, CD, CE
B,D	BC, BD, BE, DE

FACTORS AFFECTING COMPLEXITY

- Minimum support threshold
 - lower support threshold: more frequent itemsets
 - increases number and length of candidate itemsets
- Dimensionality (number of items)
 - more space needed to store support counts
 - if many frequent items: computation and I/O cost increase
- Size of database (number of transactions)
 - increases run time (Apriori makes multiple passes)
- Transaction width (dense datasets)
 - increases number of subsets
 - increases length of frequent itemsets and hash tree traversals

BACK TO ASSOCIATION RULES

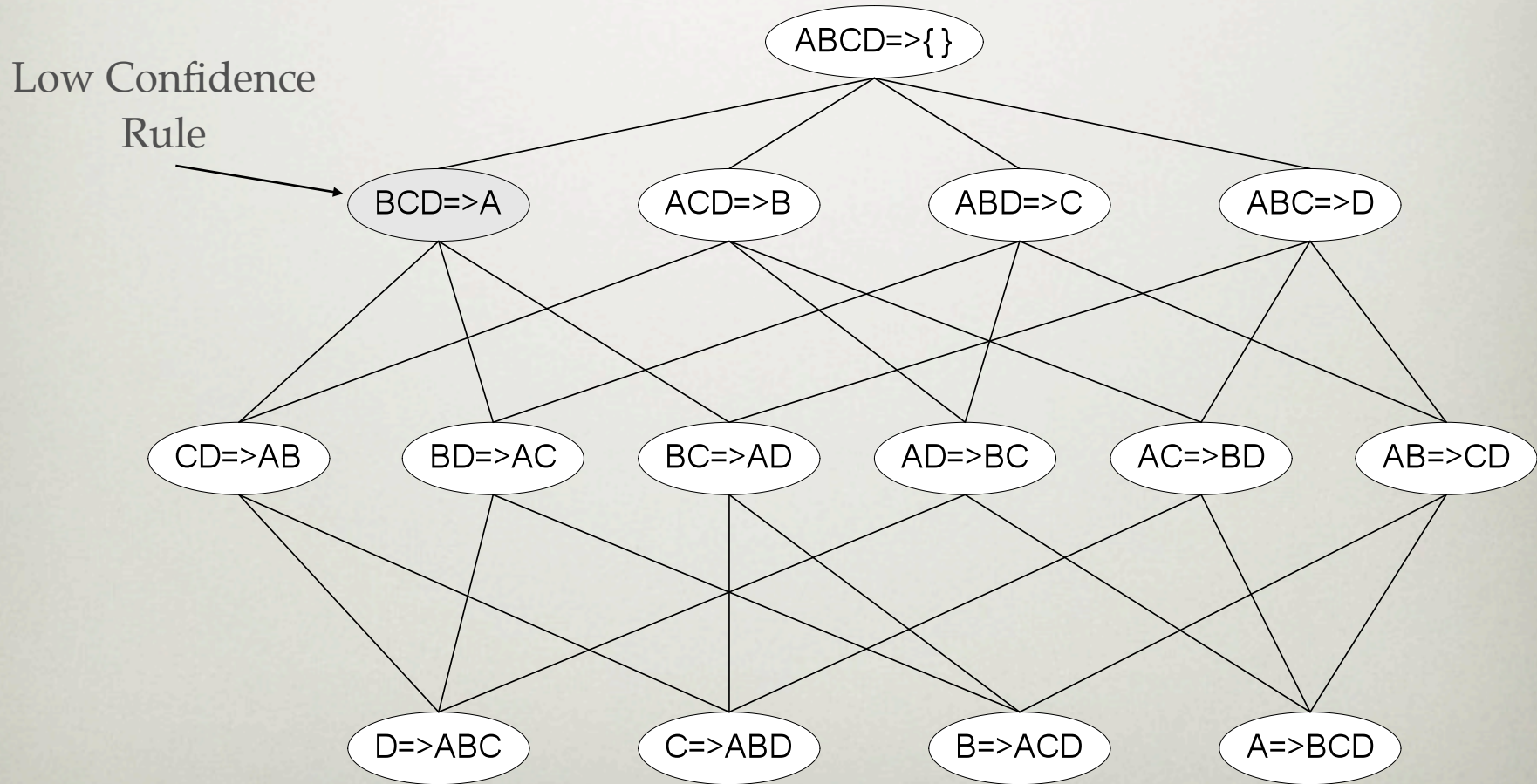
- Given a frequent itemset L , find all non-empty subsets $f \subset L$ such that $f \rightarrow L - f$ satisfies the minimum confidence requirement
- If $\{A,B,C,D\}$ is a frequent itemset, candidate rules
 - $ABC \rightarrow D, ABD \rightarrow C, ACD \rightarrow B, BCD \rightarrow A,$
 $A \rightarrow BCD, B \rightarrow ACD, C \rightarrow ABD, D \rightarrow ABC$
 $AB \rightarrow CD, AC \rightarrow BD, AD \rightarrow BC, BC \rightarrow AD,$
 $BD \rightarrow AC, CD \rightarrow AB$
- If $|L| = k$, then there are $2^k - 2$ candidate association rules (ignoring $L \rightarrow \emptyset$ and $\emptyset \rightarrow L$)

RULE GENERATION

- How to efficiently generate rules?
 - **Confidence** is *not* anti-monotone
 - $c(ABC \rightarrow D)$ can be larger or smaller than $c(AB \rightarrow D)$
 - **Confidence of rules** generated from the **same** itemset *is!*
 - e.g., $L = \{A, B, C, D\}$:
$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$
 - Anti-monotone w.r.t. number of items on the RHS of the rule

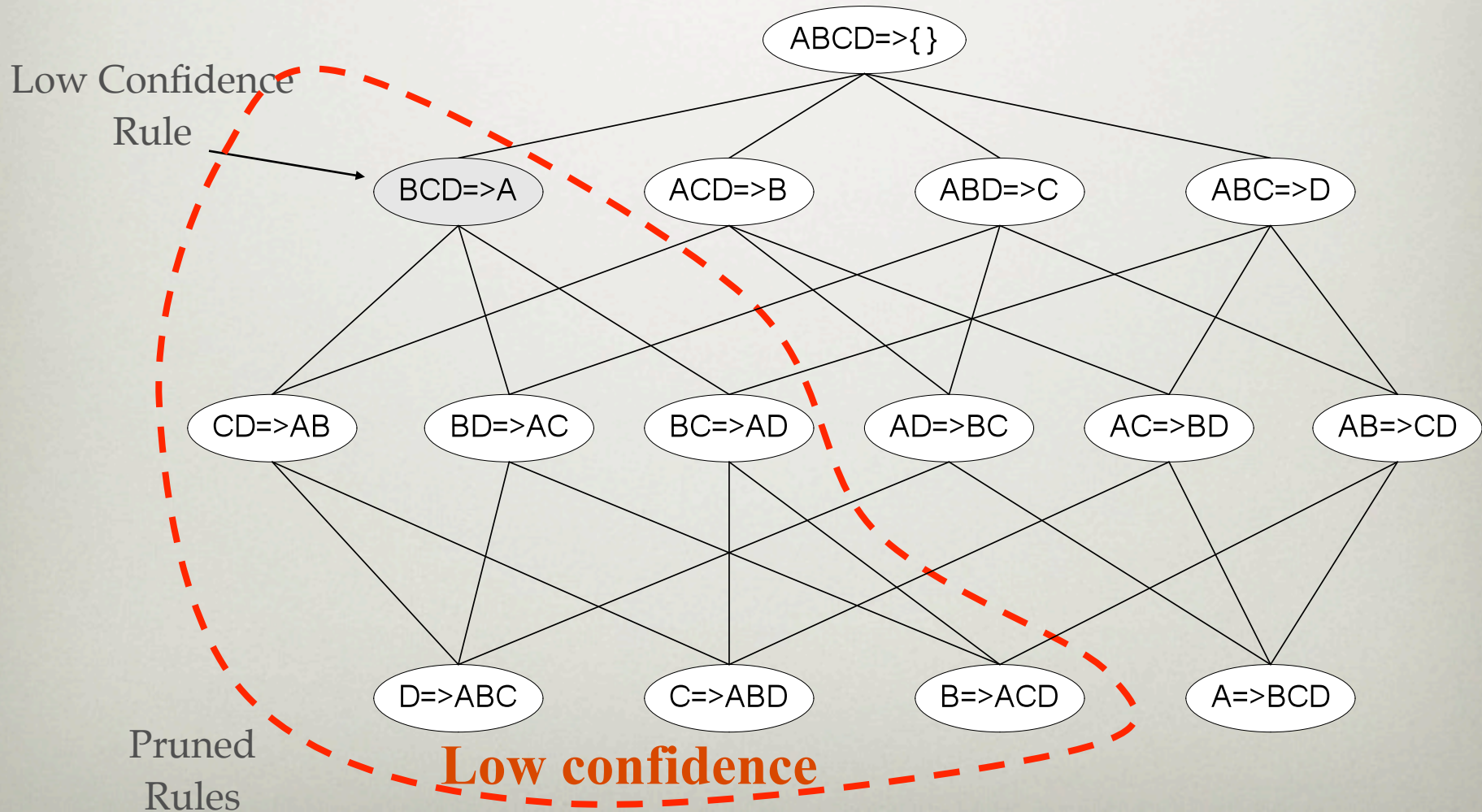
RULE GENERATION

Quiz: which rules have low confidence?



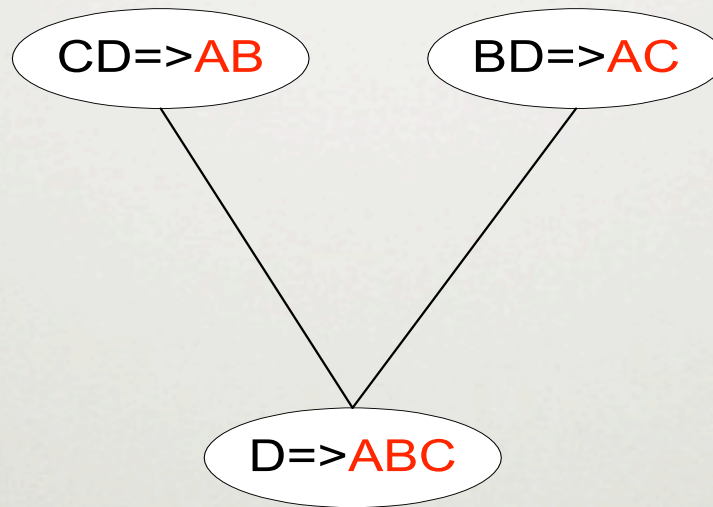
RULE GENERATION

Quiz: which rules have low confidence?



RULE GENERATION FOR APRIORI ALGORITHM

- Candidate rule generated by merging two high-confidence rules with same prefix on right hand side

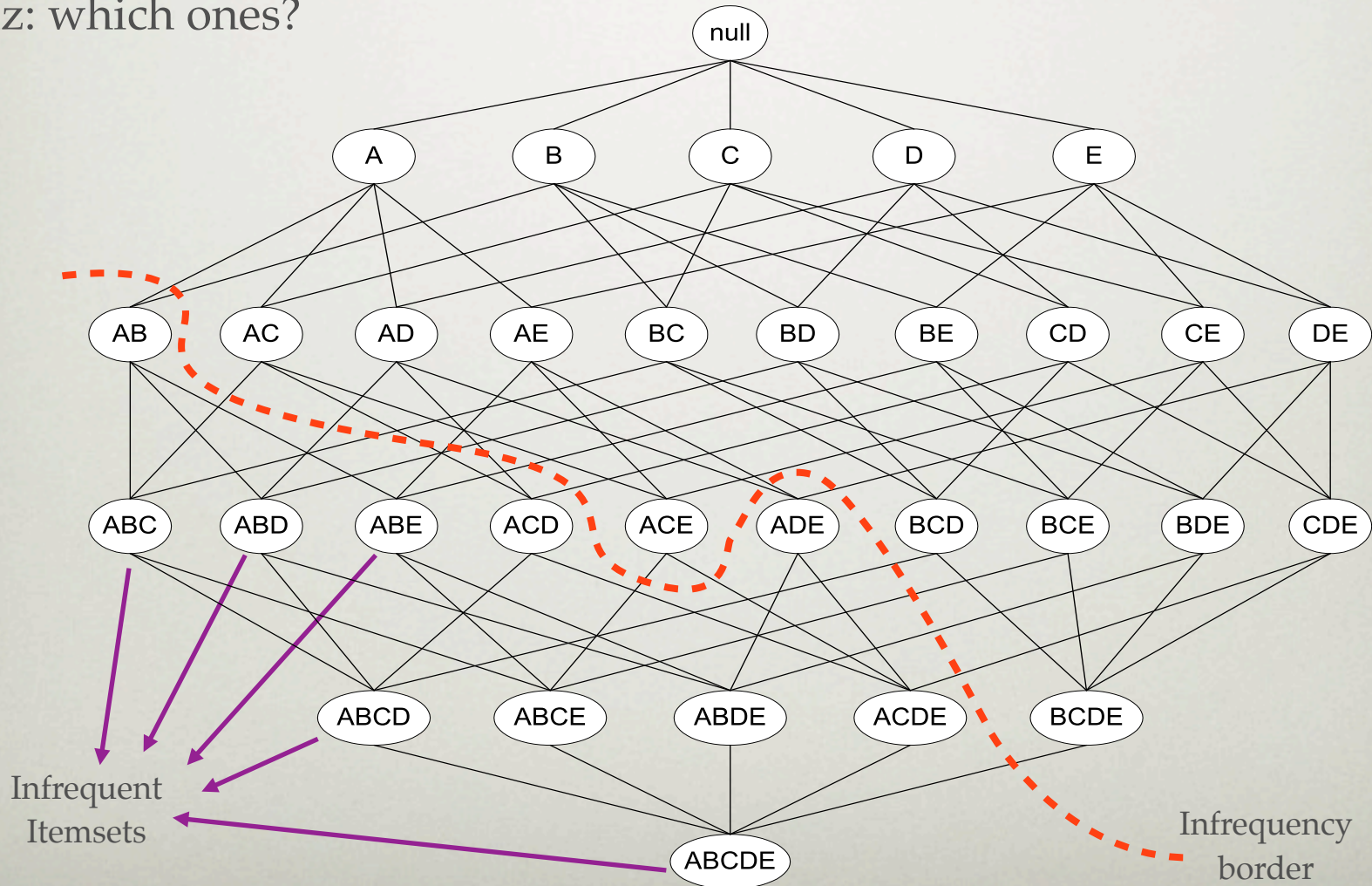


- Prune $D \Rightarrow ABC$ if subset $AD \Rightarrow BC$ has low confidence

MAXIMAL FREQUENT ITEMSET

Itemset is *maximal frequent* if **none** of its immediate supersets is frequent

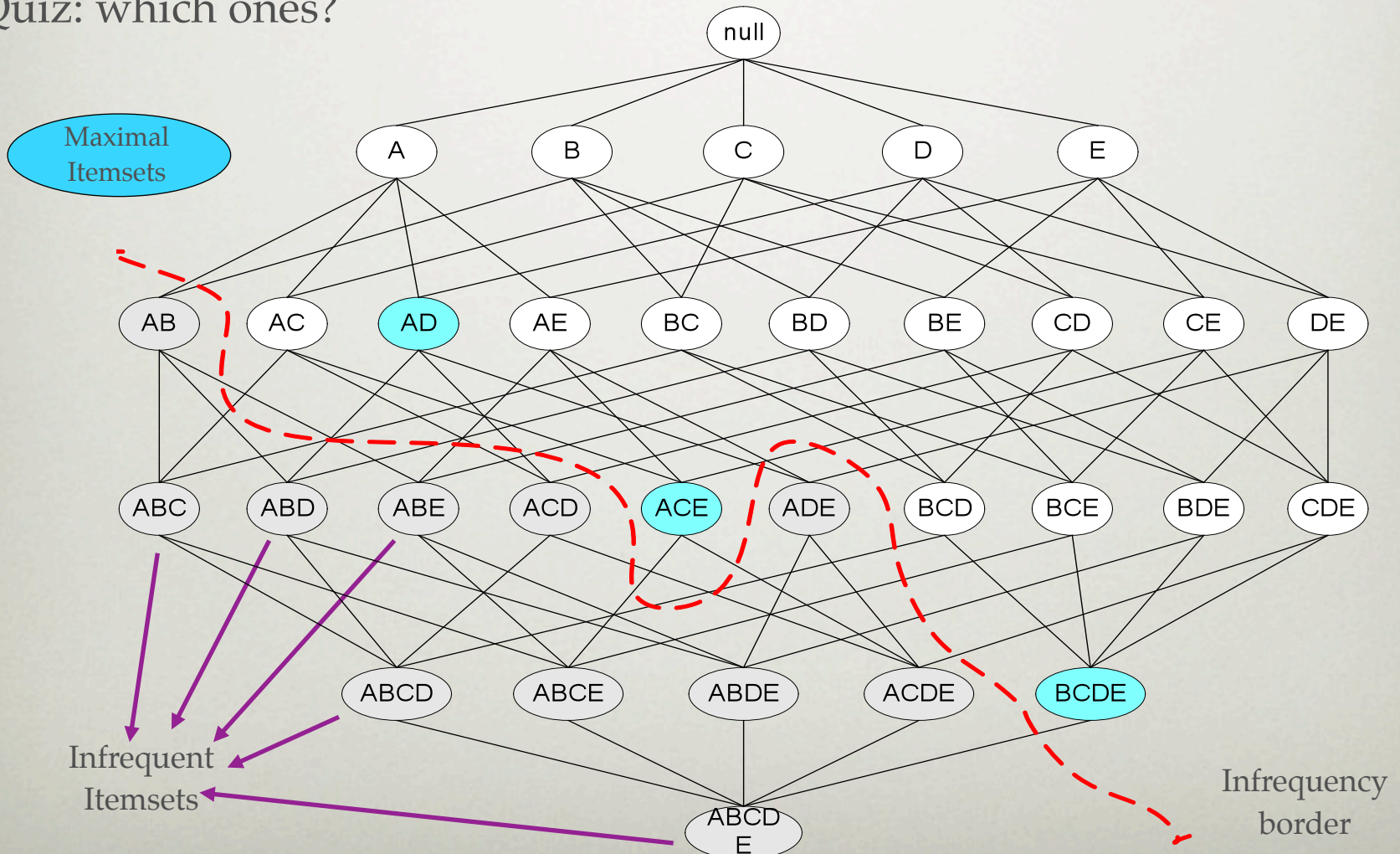
Quiz: which ones?



MAXIMAL FREQUENT ITEMSET

Itemset is *maximal frequent* if **none** of its immediate supersets is frequent

Quiz: which ones?



CLOSED ITEMSET

- An itemset is *closed* if **none** of its (immediate) supersets has the same support
- Quiz: which ones are closed?

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

CLOSED ITEMSET

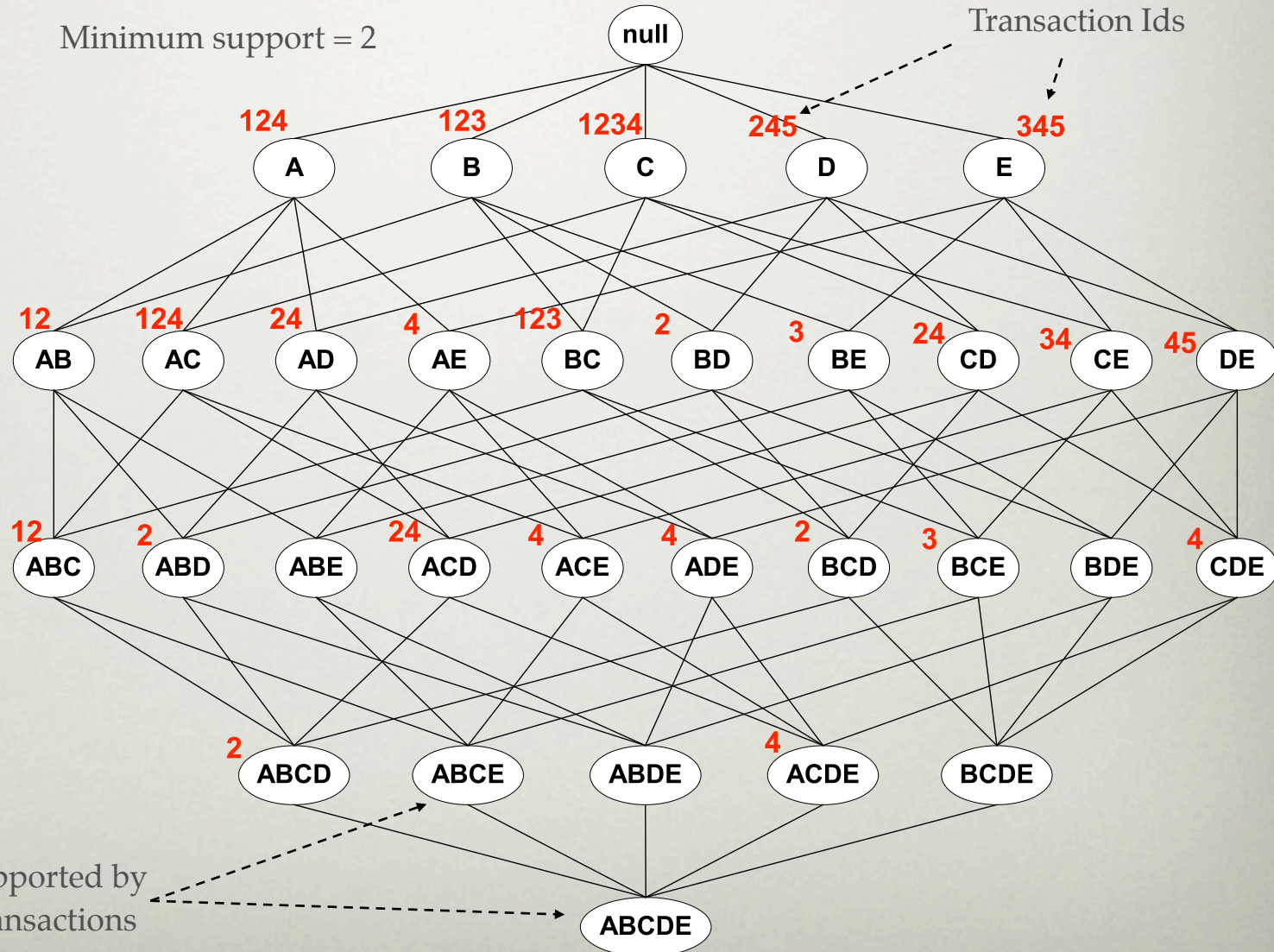
- An itemset is *closed* if **none** of its (immediate) supersets has the same support
- Quiz: which ones are closed?

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

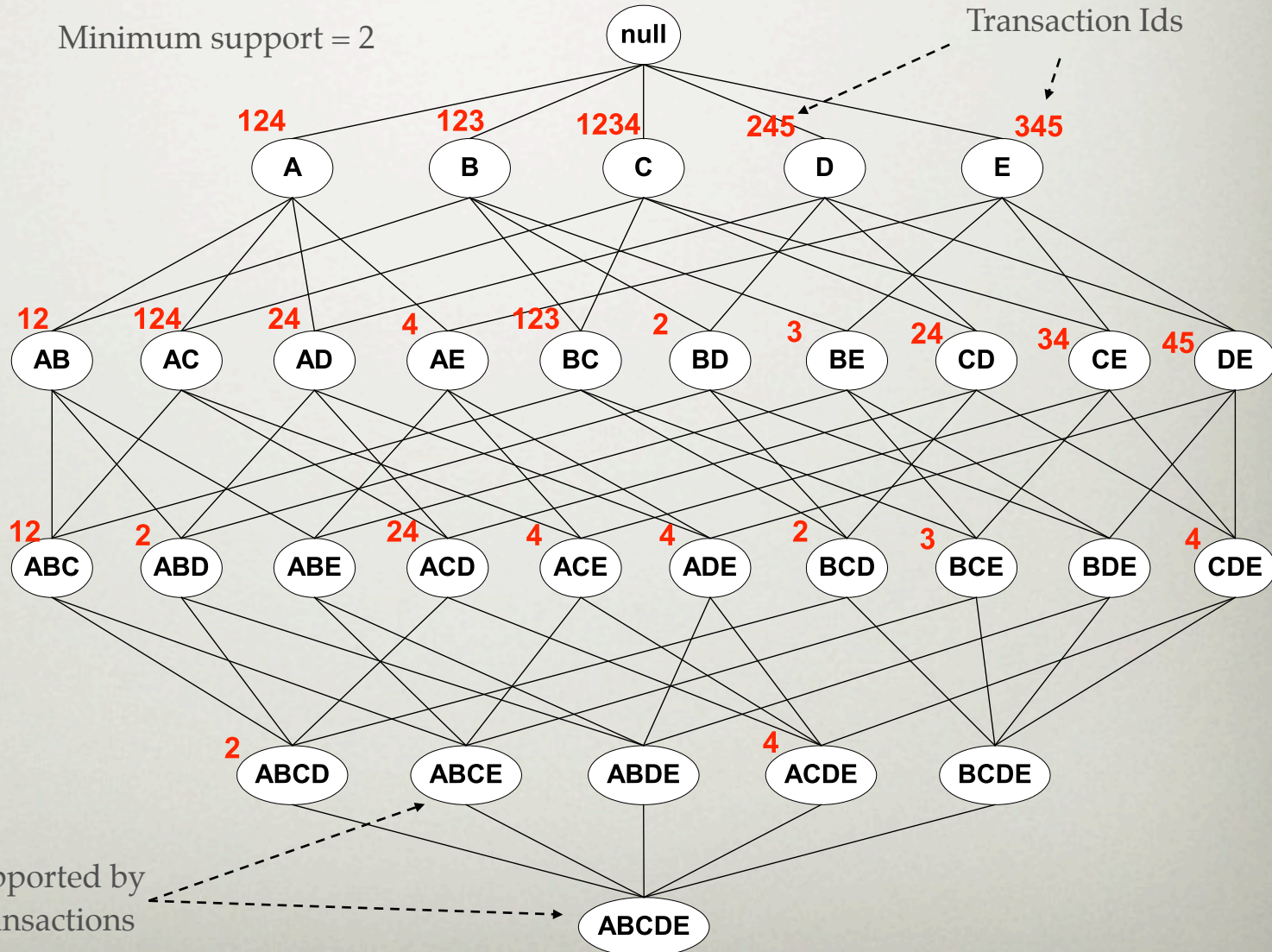
QUIZ: WHICH ARE FREQUENT, CLOSED, MAXIMAL?

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



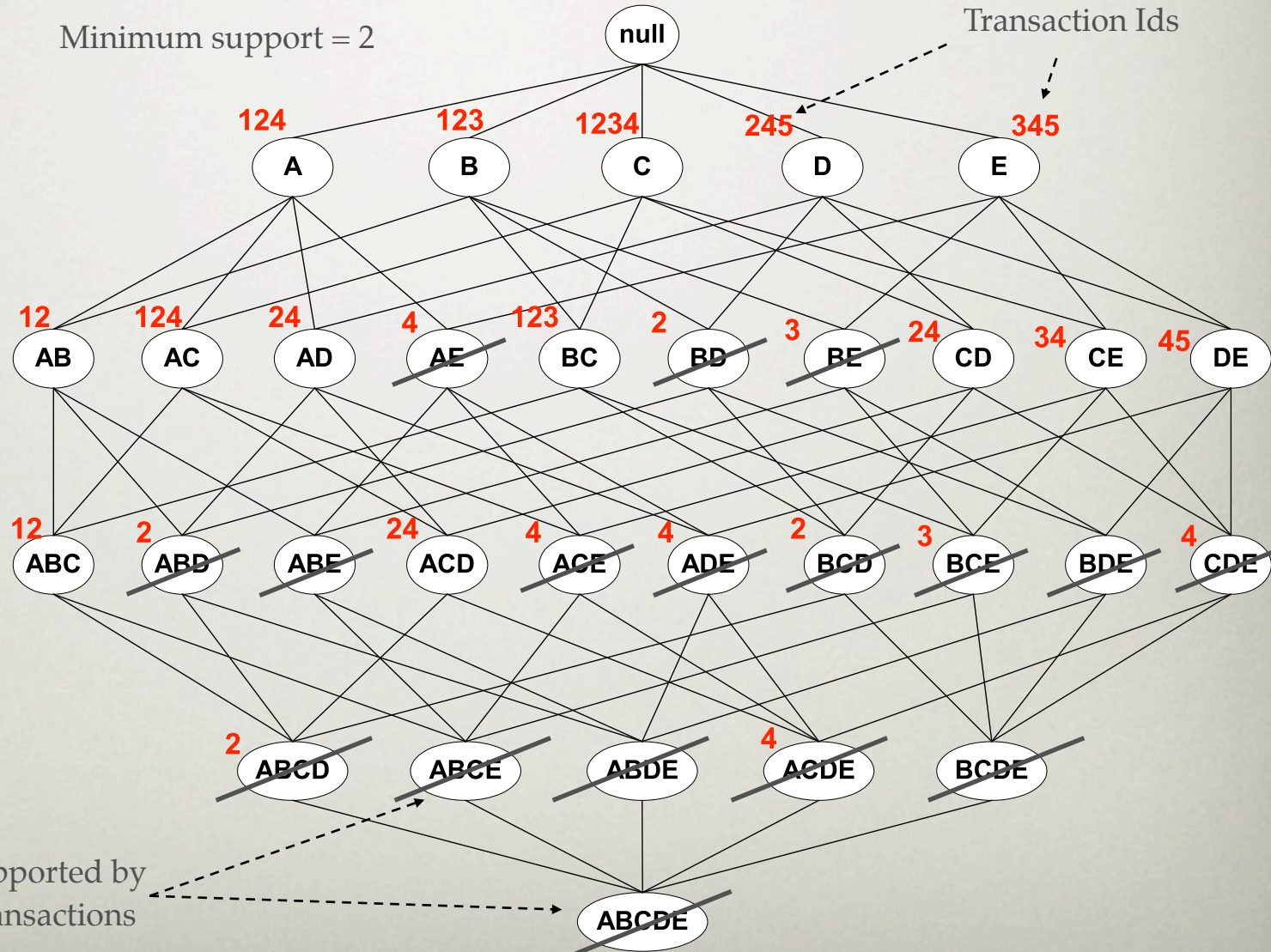
QUIZ: WHICH ARE FREQUENT, CLOSED, MAXIMAL?

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



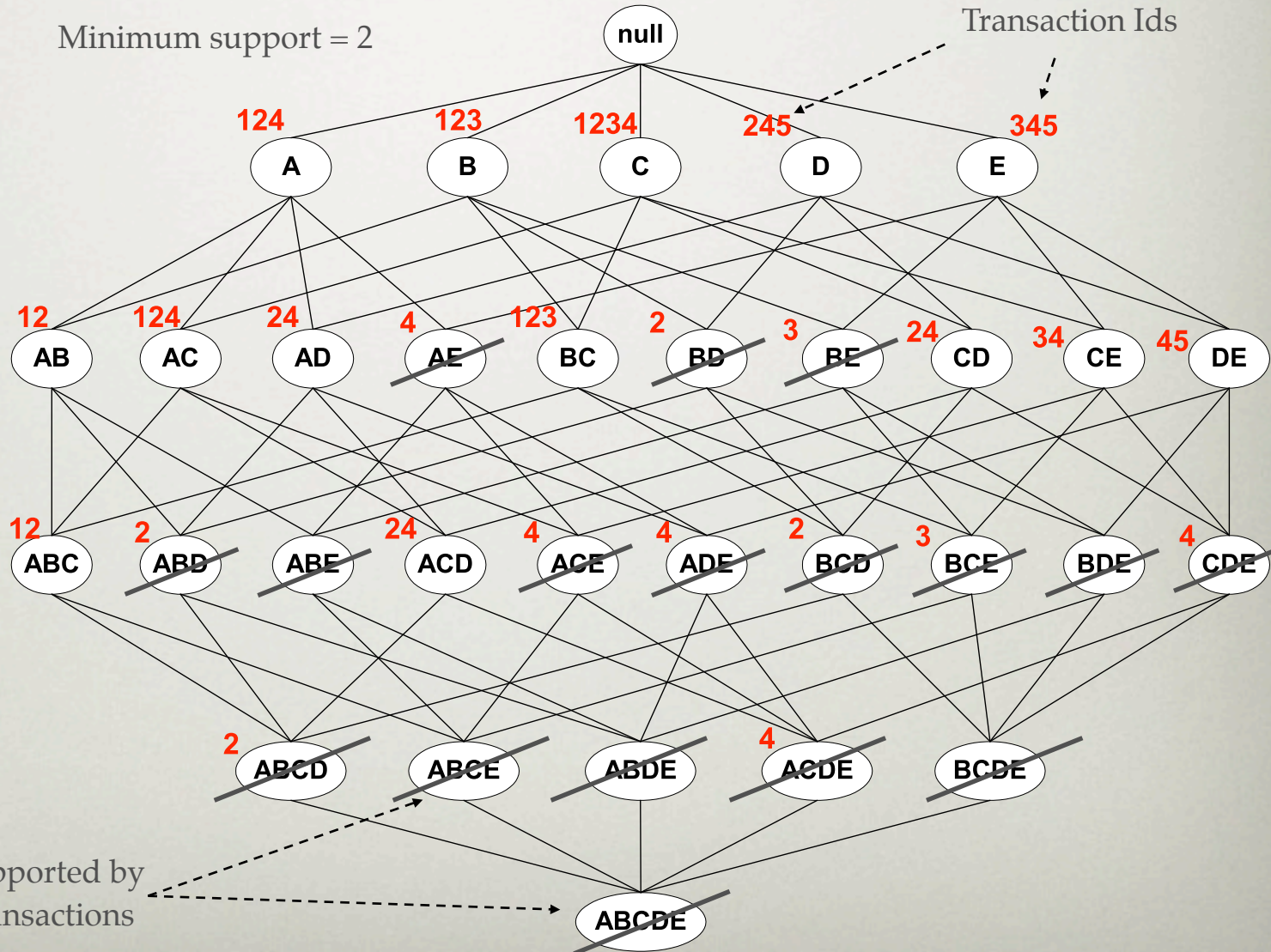
QUIZ: WHICH ARE FREQUENT, CLOSED, MAXIMAL?

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



QUIZ: WHICH ARE FREQUENT, CLOSED, MAXIMAL?

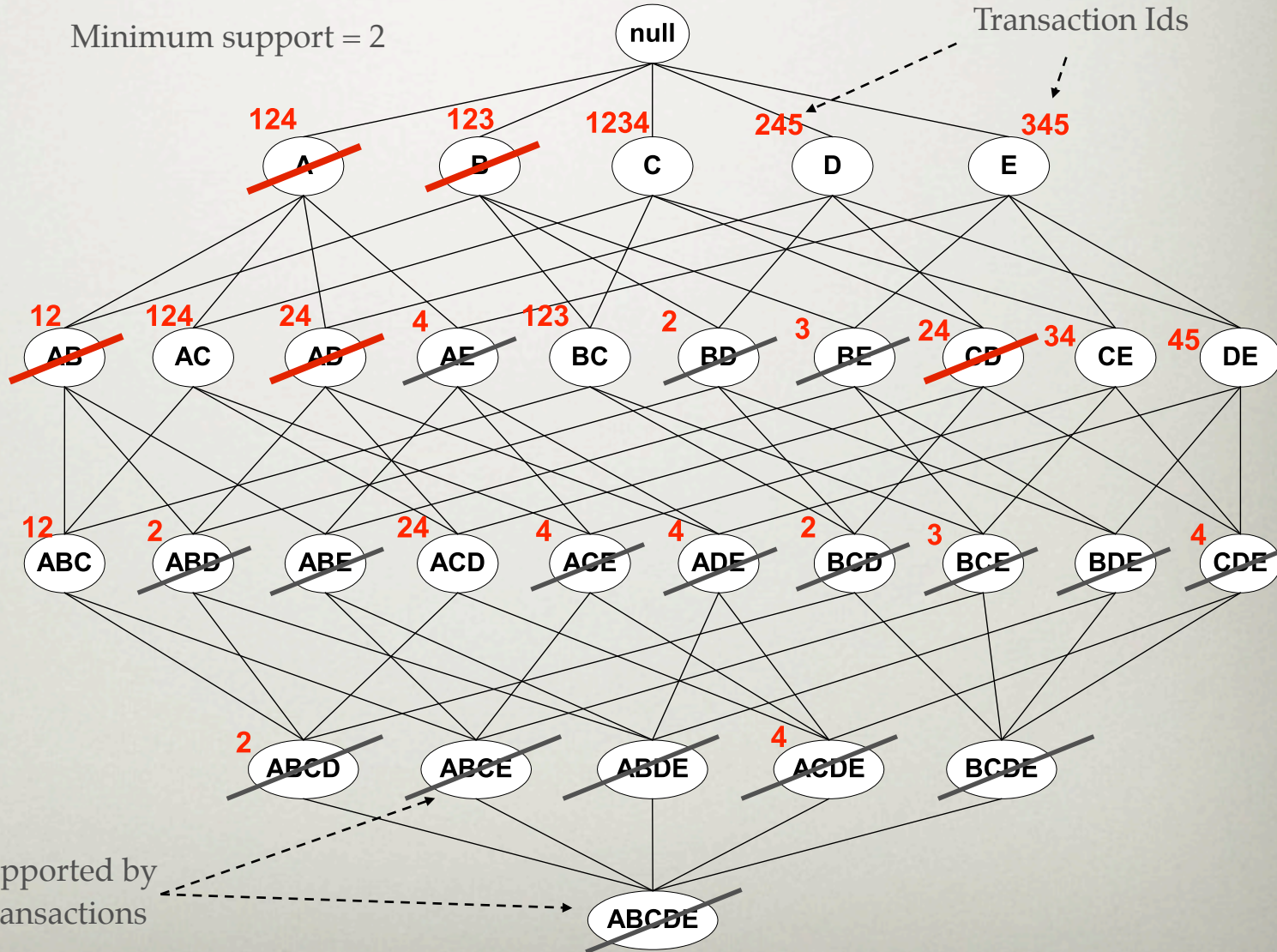
TID	Items
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE



QUIZ: WHICH ARE FREQUENT, CLOSED, MAXIMAL?

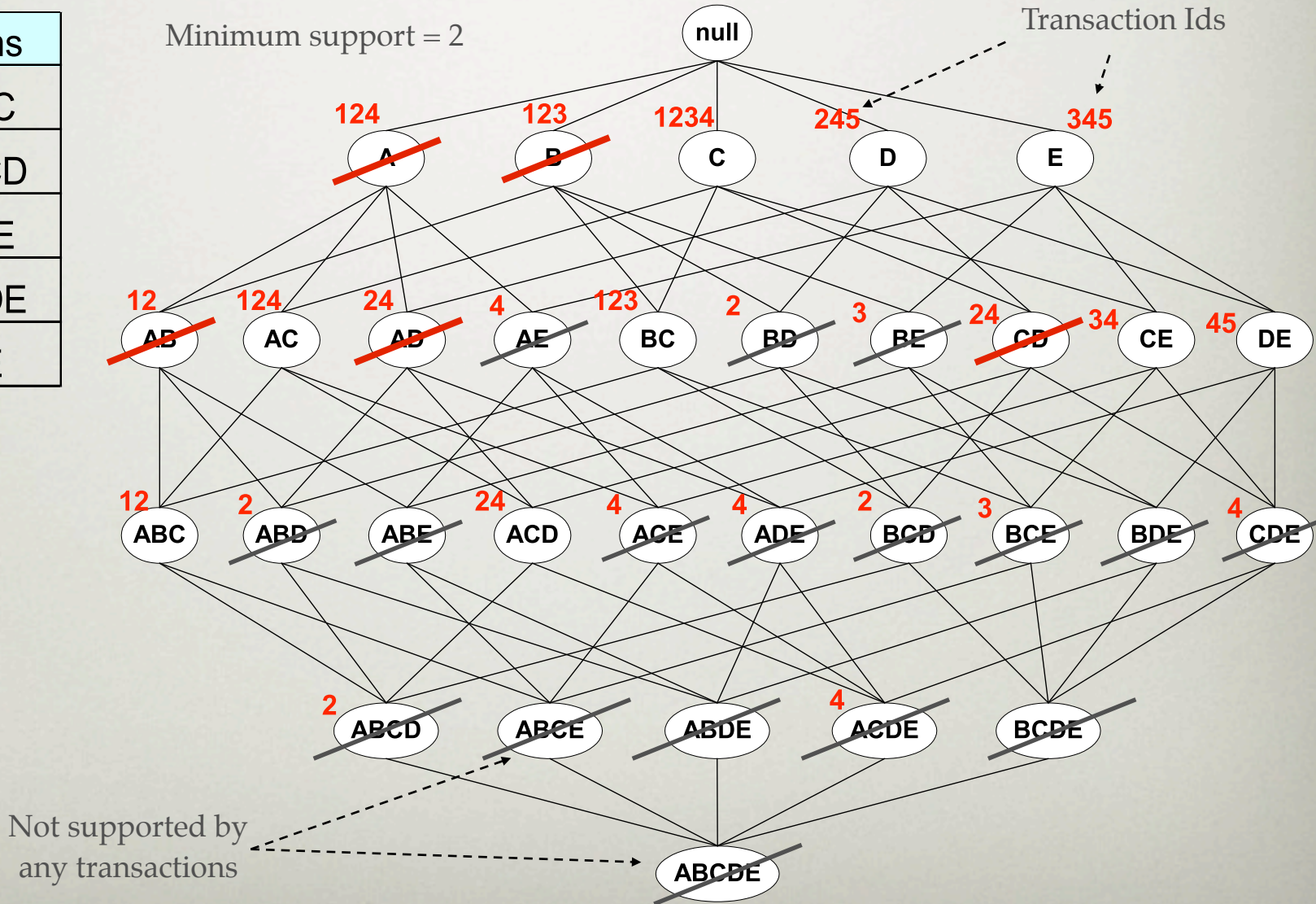
TID	Items
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE

Minimum support = 2



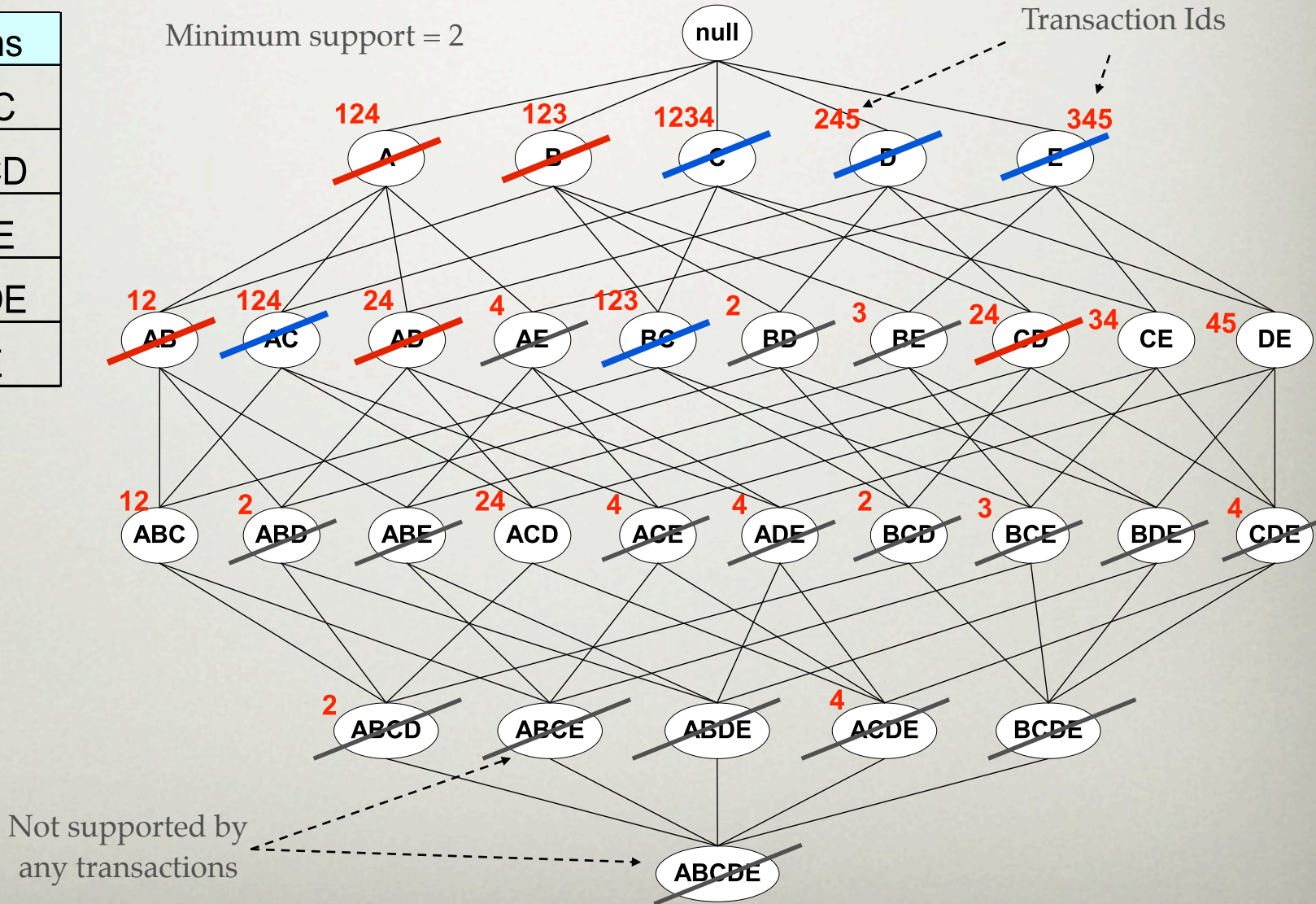
QUIZ: WHICH ARE FREQUENT, CLOSED, MAXIMAL?

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



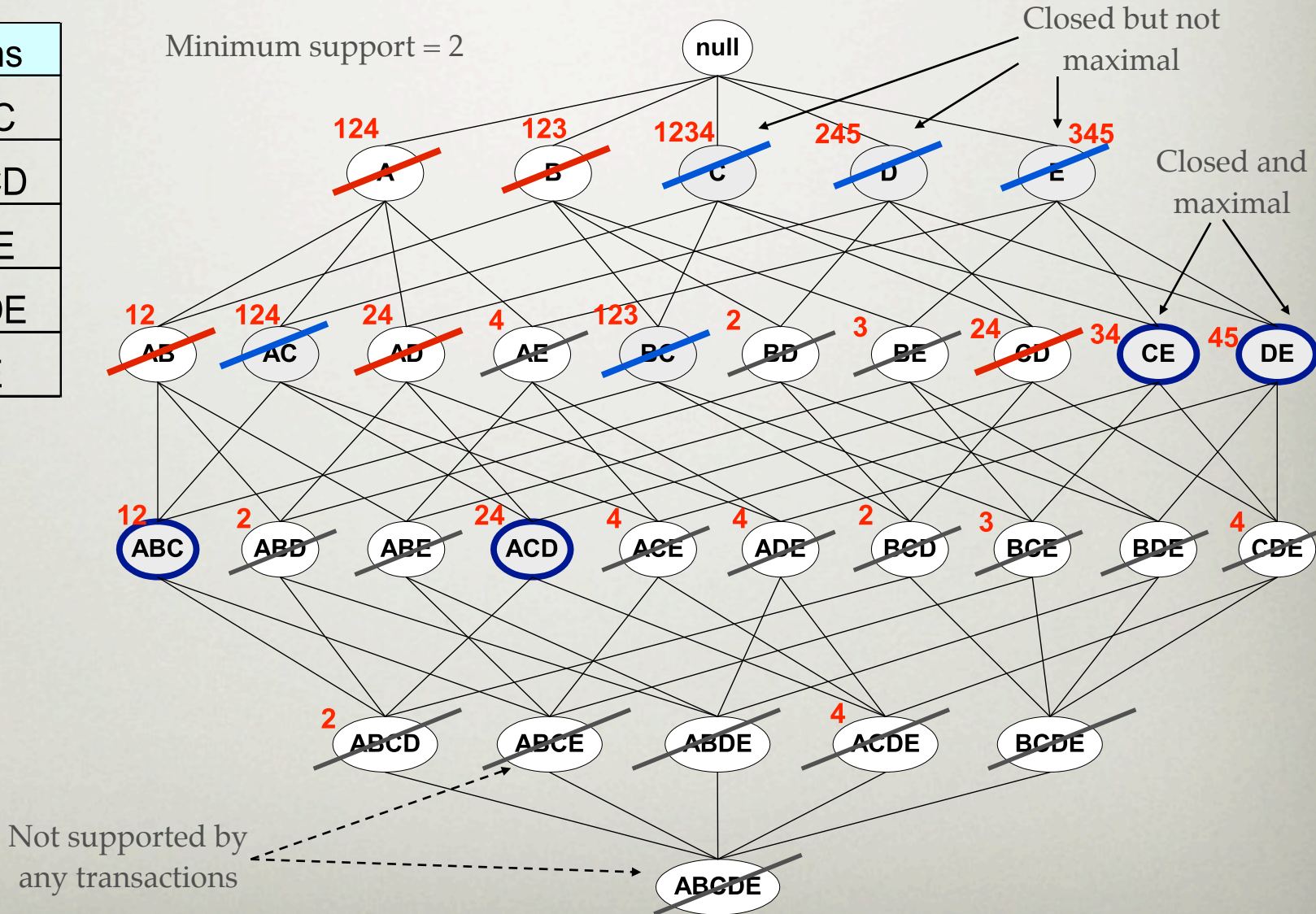
QUIZ: WHICH ARE FREQUENT, CLOSED, MAXIMAL?

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

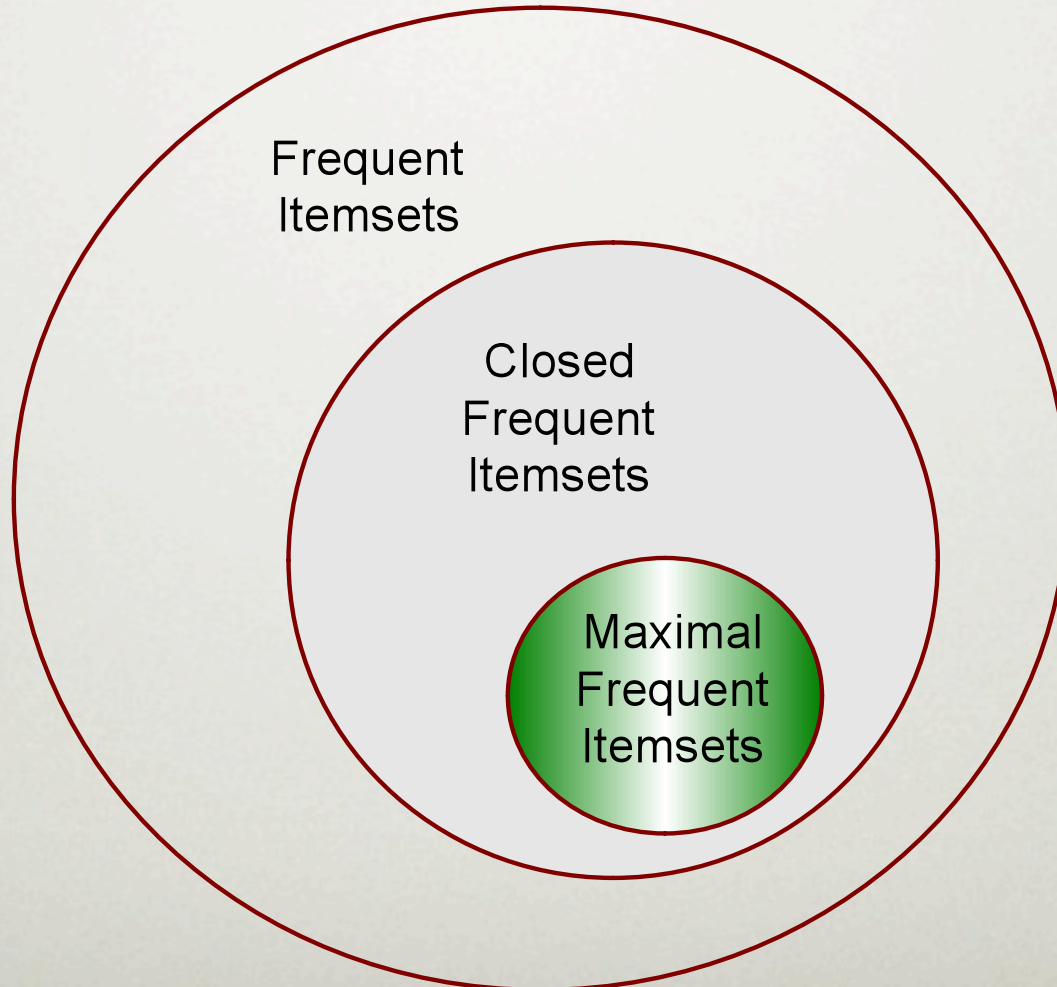


QUIZ: WHICH ARE FREQUENT, CLOSED, MAXIMAL?

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



MAXIMAL VS CLOSED ITEMSETS



EXAM QUESTION

For $minsup=0.3$, draw the itemset lattice and label each node with:

I = infrequent itemset

C = closed itemset

M = maximal itemset

tid	Items
1	{a, b, c}
2	{a, d, e}
3	{a, c}
4	{d, e}
5	{b, c}
6	{a, c, d, e}
7	{c, d, e}
8	{b, c}
9	{a, c, d, e}
10	{b}

Compute support and confidence for:

$$\{b\} \rightarrow \{c\}, \{a, d\} \rightarrow \{e\}, \text{ and } \{c\} \rightarrow \{d, e\}$$