# Responsible Data Science

## Association rule mining
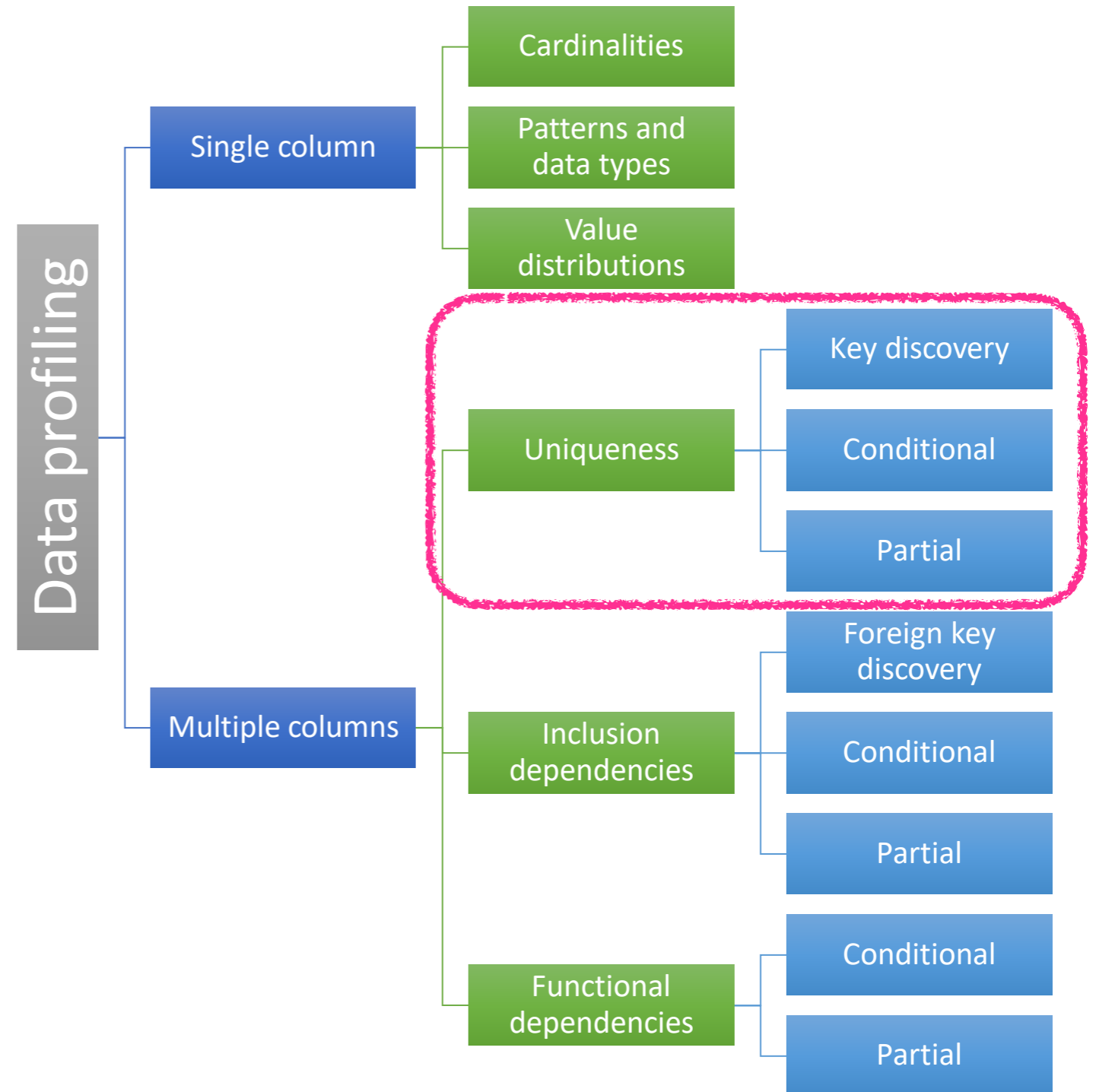## data profiling continued

**Prof. Julia Stoyanovich**

Center for Data Science &
Computer Science and Engineering
New York University

NYU

# Classification of data profiling tasks

[Abedjan, Golab, Naumann; *SIGMOD 2017*]

relational data (here: just one table)

Given a relation schema **R** *(A, B, C, D)* and a relation instance **r**, a **unique column combination** (or a **"unique"** for short) is a set of attributes **X** whose **projection** contains no duplicates in **r**

$Episodes(season, num, title, viewers)$

| season | num | title | viewers |
|--------|-----|-------|---------|
| 1 | 1 | Winter is Coming | 2.2 M |
| 1 | 2 | The Kingsroad | 2.2 M |
| 2 | 1 | The North Remembers | 3.9 M |

**Projection** is a relational algebra operation that takes as input relation **R** and returns a new relation **R'** with a subset of the columns of **R**.

$\pi_{season}(Episodes)$

| season |
|--------|
| 1 |
| 1 | non-unique |
| 2 |

$\pi_{season,num}(Episodes)$

| season | num |
|--------|-----|
| 1 | 1 |
| 1 | 2 | unique |
| 2 | 1 |

$\pi_{title}(Episodes)$

| title |
|-------|
| Winter is Coming |
| The Kingsroad | **unique** |
| The North Remembers |

# Discovering uniques

Given a relation schema **R** *(A, B, C, D)* and a relation instance **r**, a **unique column combination** (or a **"unique"** for short) is a set of attributes **X** whose **projection** contains no duplicates in **r**

*Episodes(season, num, title, viewers)*

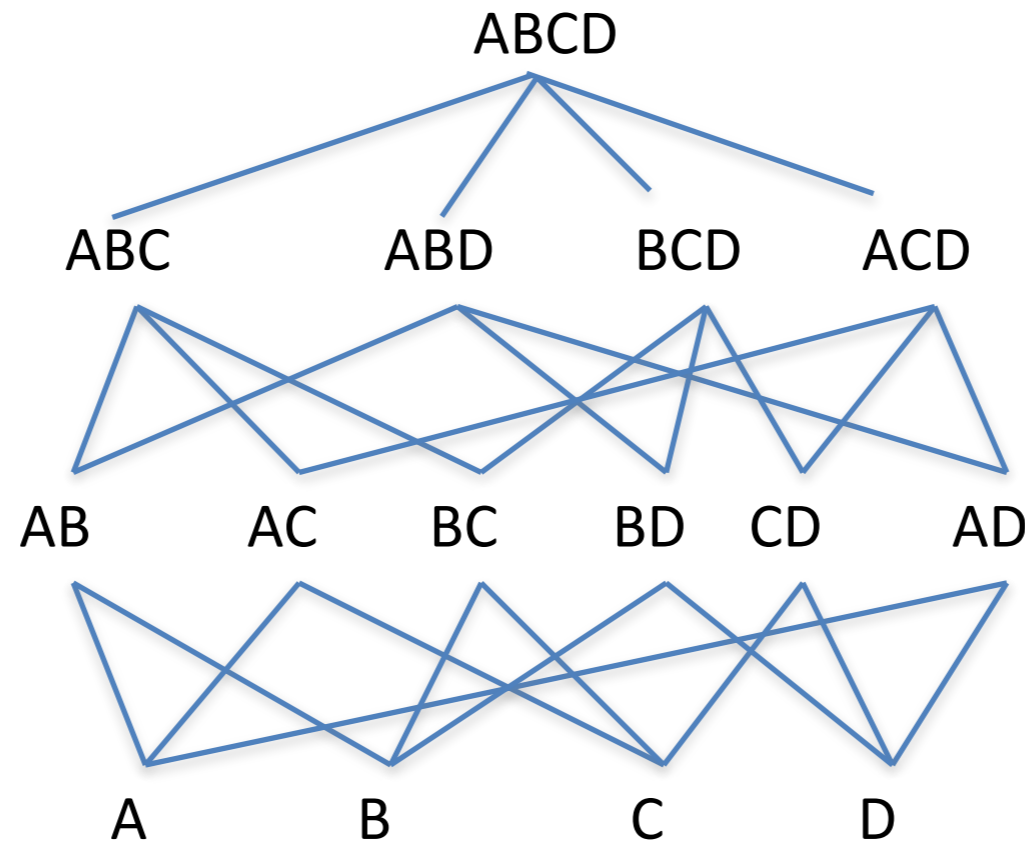| season | num | title | viewers |
|--------|-----|-------|---------|
| 1 | 1 | Winter is Coming | 2.2 M |
| 1 | 2 | The Kingsroad | 2.2 M |
| 2 | 1 | The North Remembers | 3.9 M |

**Projection** is a relational algebra operation that takes as input relation **R** and returns a new relation **R'** with a subset of the columns of **R**.

- Recall that more than one set of attributes **X** may be unique

- It may be the case that **X** and **Y** are both unique, and that they are not disjoint. When is this interesting?

R (A, B, C, D)  attribute lattice of **R**



$$\binom{4}{4} = 1$$

$$\binom{4}{3} = 4$$

$$\binom{4}{2} = 6$$
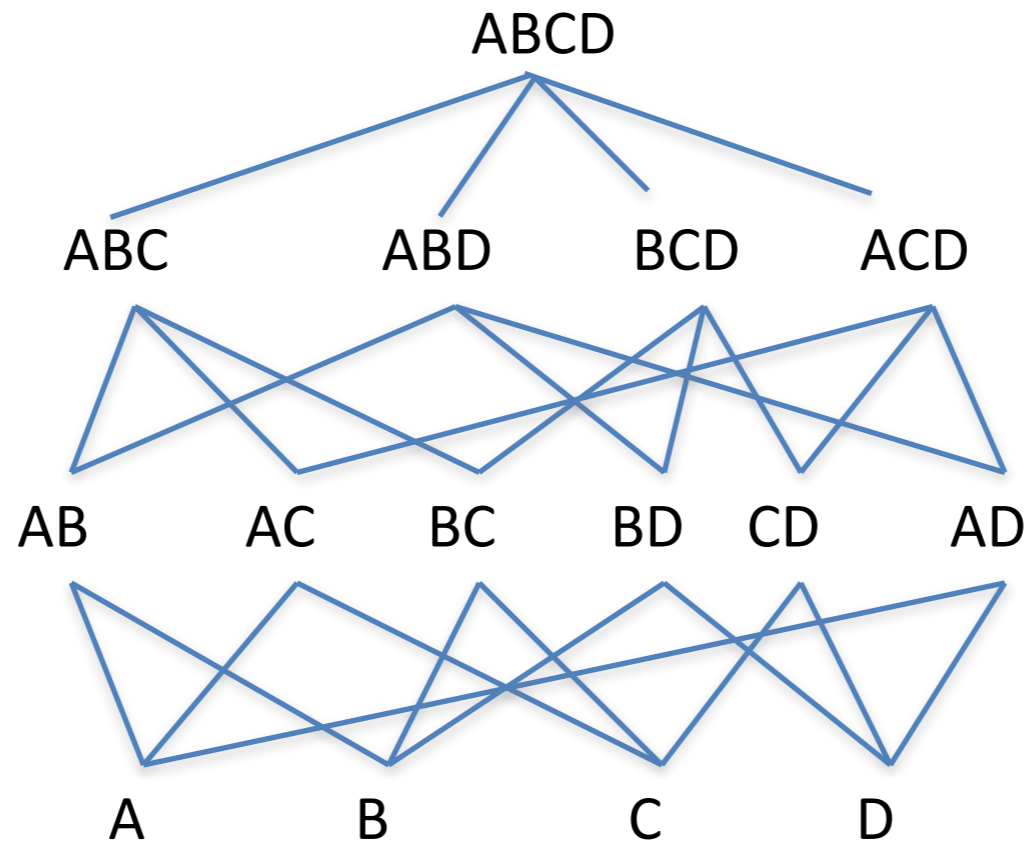
$$\binom{4}{1} = 4$$

What's the size of the attribute lattice of **R**?

**Look at all attribute combinations?**

# Discovering uniques

R (A, B, C, D)          attribute lattice of R



- If **X** is unique, then what can we say about its **superset Y**?

- If **X** is non-unique, then what can we say about its **subset Z**?

# Discovering uniques

Given a relation schema **R** *(A, B, C, D)* and a relation instance **r**, a **unique column combination** (or a **"unique"** for short) is a set of attributes **X** whose **projection** contains no duplicates in **r**

Given a relation schema **R** *(A, B, C, D)* and a relation instance **r**, a set of attributes **Y** is **non-unique** if its projection contains duplicates in **r**

**X** is **minimal unique** if every subset **Y** of **X** is non-unique

**Y** is maximal non-unique if every superset **X** of **Y** is unique

NYU

Given a relation schema **R** *(A, B, C, D)* and a relation instance **r**, a **unique column combination** is a set of attributes **X** whose **projection** contains no duplicates in **r**

*Episodes*(*season*, *num*, *title*, *viewers*)

| season | num | title | viewers |
|--------|-----|-------|---------|
| 1 | 1 | Winter is Coming | 2.2 M |
| 1 | 2 | The Kingsroad | 2.2 M |
| 2 | 1 | The North Remembers | 3.9 M |

A set of attributes is a **candidate key** for a relation if:
(1)   no two distinct tuples can have the same values for all key attributes (candidate key **uniquely identifies** a tuple), *and*
(2)  this is not true for any subset of the key attributes (candidate key **is minimal**)

**A minimal unique of a relation instance is a (possible) candidate key of the relation schema.** To find all possible candidate keys, find all minimal uniques in a relation instance.

# The early days of data mining

- Problem formulation due to Agrawal, Imielinski, Swami, SIGMOD 1993

- Solution: the **Apriori** algorithm by Agrawal & Srikant, VLDB 1994

- Initially for **market-basket data** analysis, has many other applications, we'll see one today

- We wish to answer two related questions:
  - **Frequent itemsets:** Which items are often purchased together, e.g., milk and cookies are often bought together
  - **Association rules:** Which items will likely be purchased, based on other purchased items, e.g., if diapers are bought in a transaction, beer is also likely bought in the same transaction

# Market-basket data

- $I = \{i_1, i_2, \ldots, i_m\}$ is the set of available items, e.g., a product catalog of a store

- $X \subseteq I$ is an **itemset**, e.g., {milk, bread, cereal}

- **Transaction** $t$ is a set of items purchased together, $t \subseteq I$, has a transaction id (TID)

    $t_1$: {bread, cheese, milk}

    $t_2$: {apple, eggs, salt, yogurt}

    $t_3$: {biscuit, cheese, eggs, milk}

- Database $T$ is a set of transactions $\{t_1, t_2, \ldots, t_n\}$

- A transaction $t$ **supports** an itemset $X$ if $X \subseteq t$

- Itemsets supported by at least *minSupp* transactions are called **frequent itemsets**

    **minSupp, which can be a number or a percentage, is specified by the user**

# Itemsets

| TID | Items |
|-----|-------|
| 1 | A |
| 2 | A C |
| 3 | A B D |
| 4 | A C |
| 5 | A B C |
| 6 | A B C |

**minSupp** = 2 transactions

How many possible itemsets are there (excluding the empty itemset)?

$$2^4 - 1 = 15$$

| itemset | support |
|---------|---------|
| ⭐ A | 6 |
| ⭐ B | 3 |
| ⭐ C | 4 |
| D | 1 |
| ⭐ A B | 3 |
| ⭐ A C | 4 |
| A D | 1 |
| ⭐ B C | 2 |
| B D | 1 |
| C D | 0 |
| ⭐ A B C | 2 |
| A B D | 1 |
| B C D | 0 |
| A C D | 0 |
| A B C D | 0 |

# Association rules

An **association rule** is an implication $X \rightarrow Y$, where $X, Y \subset I$, and $X \cap Y = \emptyset$

example: {milk, bread} $\rightarrow$ {cereal}

"A customer who purchased X is also likely to have purchased Y in the same transaction"

we are interested in rules with a **single item** in Y

can we represent {milk, bread} $\rightarrow$ {cereal, cheese}?

Rule $X \rightarrow Y$ holds with **support** *supp* in T if *supp* of transactions contain $X \cup Y$

Rule $X \rightarrow Y$ holds with confidence *conf* in T if *conf* % of transactions that contain X also contain Y

$conf \approx \Pr(Y \mid X)$

$conf (X \rightarrow Y) = supp (X \cup Y) / supp (X)$

# Association rules

**minSupp** = 2 transactions
**minConf** = 0.75

| itemset | support |
|---------|---------|
| A | 6 |
| B | 3 |
| C | 4 |
| D | 1 |
| A B | 3 |
| A C | 4 |
| A D | 1 |
| B C | 2 |
| B D | 1 |
| C D | 0 |
| A B C | 2 |
| A B D | 1 |
| B C D | 0 |
| A C D | 0 |
| A B C D | 0 |

supp = 3
A → B    conf = 3 / 6 = 0.5
B → A    conf = 3 / 3 = 1.0 ★

supp = 2
B → C    conf = 2 / 3 = 0.67
C → B    conf = 2 / 4 = 0.5

supp = 4
A → C    conf = 4 / 6 = 0.67
C → A    conf = 4 / 4 = 1.0 ★

supp = 2
AB → C    conf = 2 / 3 = 0.67
AC → B    conf = 2 / 4 = 0.5
BC → A    conf = 2 / 2 = 1.0 ★

**conf (X → Y) = supp (X U Y) / supp (X)**

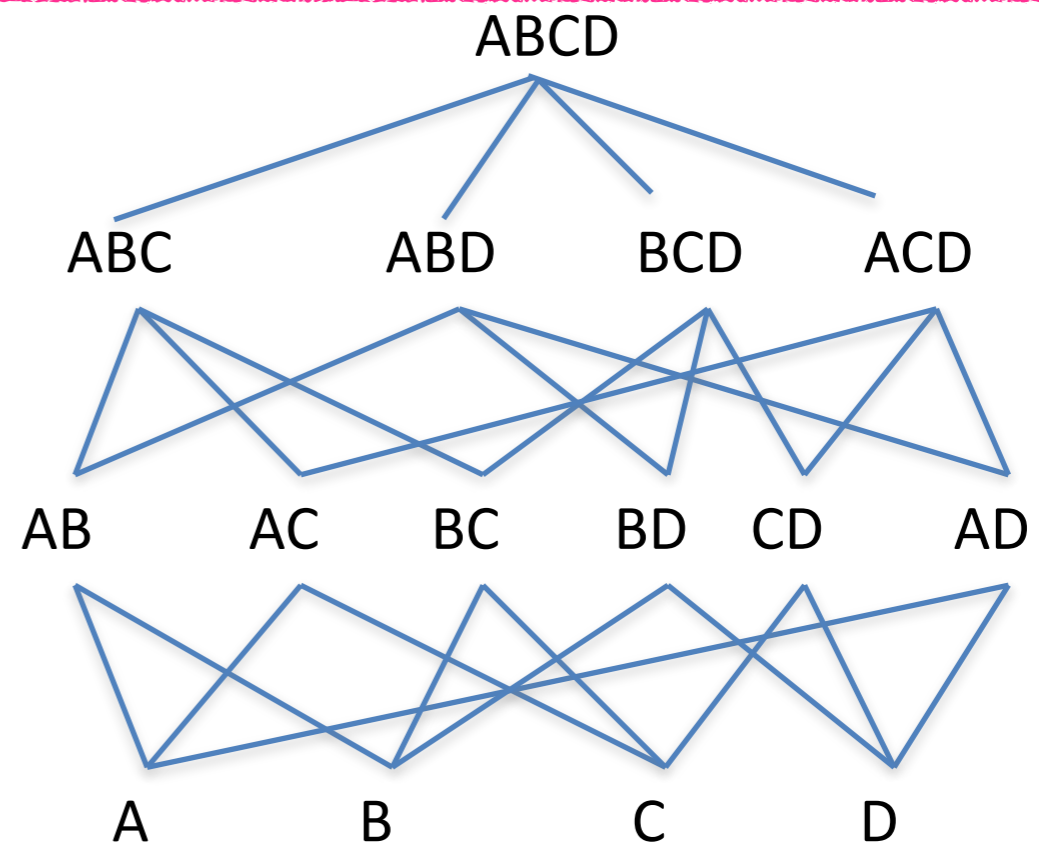Julia Stoyanovich

NYU

# Association rule mining

- Goal: find all association rules that satisfy the user-specified minimum support and minimum confidence

- Algorithm outline
  - Step 1: find all frequent itemsets
  - Step 2: find association rules

- Take 1: naïve algorithm for frequent itemset mining
  - Enumerate all subsets of *I*, check their support in *T*
  - **What is the complexity?**

# Key idea: downward closure

| itemset | support |
|---|---|
| ⭐ A | 6 |
| ⭐ B | 3 |
| ⭐ C | 4 |
| D | 1 |
| ⭐ A B | 3 |
| ⭐ A C | 4 |
| A D | 1 |
| ⭐ B C | 2 |
| B D | 1 |
| C D | 0 |
| ⭐ A B C | 2 |
| A B D | 1 |
| B C D | 0 |
| A C D | 0 |
| A B C D | 0 |

All subsets of a frequent itemset **X** are themselves frequent

So, if some subset of X is infrequent, then X cannot be frequent, we know this **apriori**



The converse is not true! If all subsets of **X** are frequent, **X** is not guaranteed to be frequent

**Algorithm Apriori(*T*, *minSupp*)**

    *$F_1$ = {frequent 1-itemsets};*

    **for** (*k* = 2; $F_{k-1}$ ≠ ∅; *k*++) **do**

        $C_k$ ← **candidate-gen**($F_{k-1}$);

        **for** each transaction *t* ∈ *T* **do**

            **for** each candidate *c* ∈ $C_k$ **do**

                **if** *c* is contained in *t* **then**

                    *c.count*++;

            **end**

        **end**

        $F_k$ ← {*c* ∈ $C_k$ | *c.count* ≥ *minSupp*}

    **end**

return *F* ← ⋃$_k$ $F_k$;

| itemset | support |
|---------|---------|
| ⭐ A | 6 |
| ⭐ B | 3 |
| ⭐ C | 4 |
| D | 1 |
| ⭐ A B | 3 |
| ⭐ A C | 4 |
| A D | 1 |
| ⭐ B C | 2 |
| B D | 1 |
| C D | 0 |
| ⭐ A B C | 2 |
| A B D | 1 |
| B C D | 0 |
| A C D | 0 |
| A B C D | 0 |

# Candidate generation

The **candidate-gen** function takes $F_{k-1}$ and returns a superset (called the candidates) of the set of all frequent k-itemsets. It has two steps:

Join: generate all possible candidate itemsets $C_k$ of length k

Prune: optionally remove those candidates in $C_k$ that have infrequent subsets

```
Insert into C_k (
  select    p.item_1, p.item_2, …, p.item_{k-1}, q.item_{k-1}
  from      F_{k-1} p, F_{k-1} q
  where     p.item_1 = q.item_1
    and        p.item_2 = q.item_2
  and      …
  and      p.item_{k-1} < q.item_{k-1} )
```

| itemset | support |
|---|---|
| A | 6 |
| B | 3 |
| C | 4 |
| D | 1 |
| A B | 3 |
| A C | 4 |
| A D | 1 |
| B C | 2 |
| B D | 1 |
| C D | 0 |
| A B C | 2 |
| A B D | 1 |
| B C D | 0 |
| A C D | 0 |
| A B C D | 0 |

**F₁ as p**

| |
|---|
| **A** |
| **B** |
| C |

**F₁ as q**

| |
|---|
| A |
| **B** |
| **C** |

**C₂**

| | |
|---|---|
| **A** | **B** |
| **A** | **C** |
| **B** | **C** |

```
Insert into Ck (
  select    p.item1, p.item2, …, p.itemk-1, q.itemk-1
  from      Fk-1 p, Fk-1 q
  where     p.item1 = q.item1
    and        p.item2 = q.item2
  and     …
  and     p.itemk-1 < q.itemk-1 )
```

Insert into $C_k$ (
 select $p.item_1, p.item_2, …, p.item_{k-1}, q.item_{k-1}$
 from $F_{k-1}\ p, F_{k-1}\ q$
 where $p.item_1 = q.item_1$
   and $p.item_2 = q.item_2$
 and …
 and $p.item_{k-1} < q.item_{k-1}$ )

| itemset | support |
|---------|---------|
| ⭐ A | 6 |
| ⭐ B | 3 |
| ⭐ C | 4 |
| D | 1 |
| ⭐ A B | 3 |
| ⭐ A C | 4 |
| A D | 1 |
| ⭐ B C | 2 |
| B D | 1 |
| C D | 0 |
| ⭐ A B C | 2 |
| A B D | 1 |
| B C D | 0 |
| A C D | 0 |
| A B C D | 0 |

$F_2$ as p

| | |
|---|---|
| **A** | **B** |
| A | C |
| B | C |

$F_2$ as q

| | |
|---|---|
| A | B |
| **A** | **C** |
| B | C |

$C_3$

| | | |
|---|---|---|
| **A** | **B** | **C** |

NYU

Assume a lexicographic ordering of the items

**Join**

```
Insert into Ck (
    select    p.item1, p.item2, …, p.itemk-1, q.itemk-1
    from      Fk-1 p, Fk-1 q
    where     p.item1 = q.item1
    and       p.item2 = q.item2
    and       …
    and       p.itemk-1 < q.itemk-1)
```
**why not p.item$_{k-1}$ ≠ q.item$_{k-1}$?**

**Prune**

> **for** each c in C$_k$ **do**
> > **for** each (k-1) subset s of c **do**
> > > **if** (s not in F$_{k-1}$) **then**
> > > > delete c from C$_k$

# Generating association rules

Rules = $\varnothing$

**for** each frequent *k-itemset* X **do**

    **for** each 1-itemset A $\subset$ X **do**

        compute conf (X / A $\rightarrow$ A) = supp(X) / sup (X / A)

        **if** conf (X / A $\rightarrow$ A) $\geq$ minConf **then**

          *Rules* $\leftarrow$ "X / A $\rightarrow$ A"

        **end**

    **end**

**end**

**return** Rules

# Performance of *Apriori*

- The possible number of frequent itemsets is exponential, O($2^m$), where **m** is the number of items

- Apriori exploits sparseness and locality of data

  - Still, it may produce a large number of rules: thousands, tens of thousands, ….

  - So, thresholds should be set carefully. What are some good heuristics?

NYU
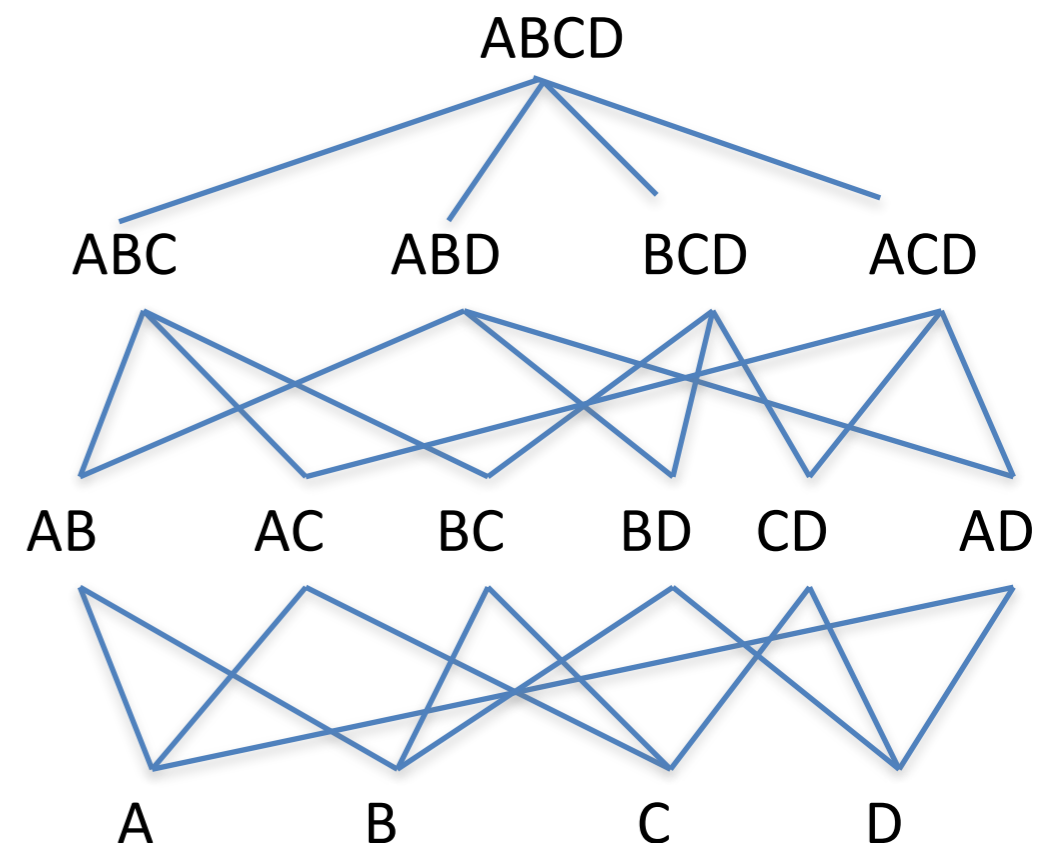
Given a relation schema **R** *(A, B, C, D)* and a relation instance **r**, a **unique column combination** (or a **"unique"** for short) is a set of attributes **X** whose **projection** contains no duplicates in **r**

Given a relation schema **R** *(A, B, C, D)* and a relation instance **r**, a set of attributes **Y** is **non-unique** if its projection contains duplicates in **r**
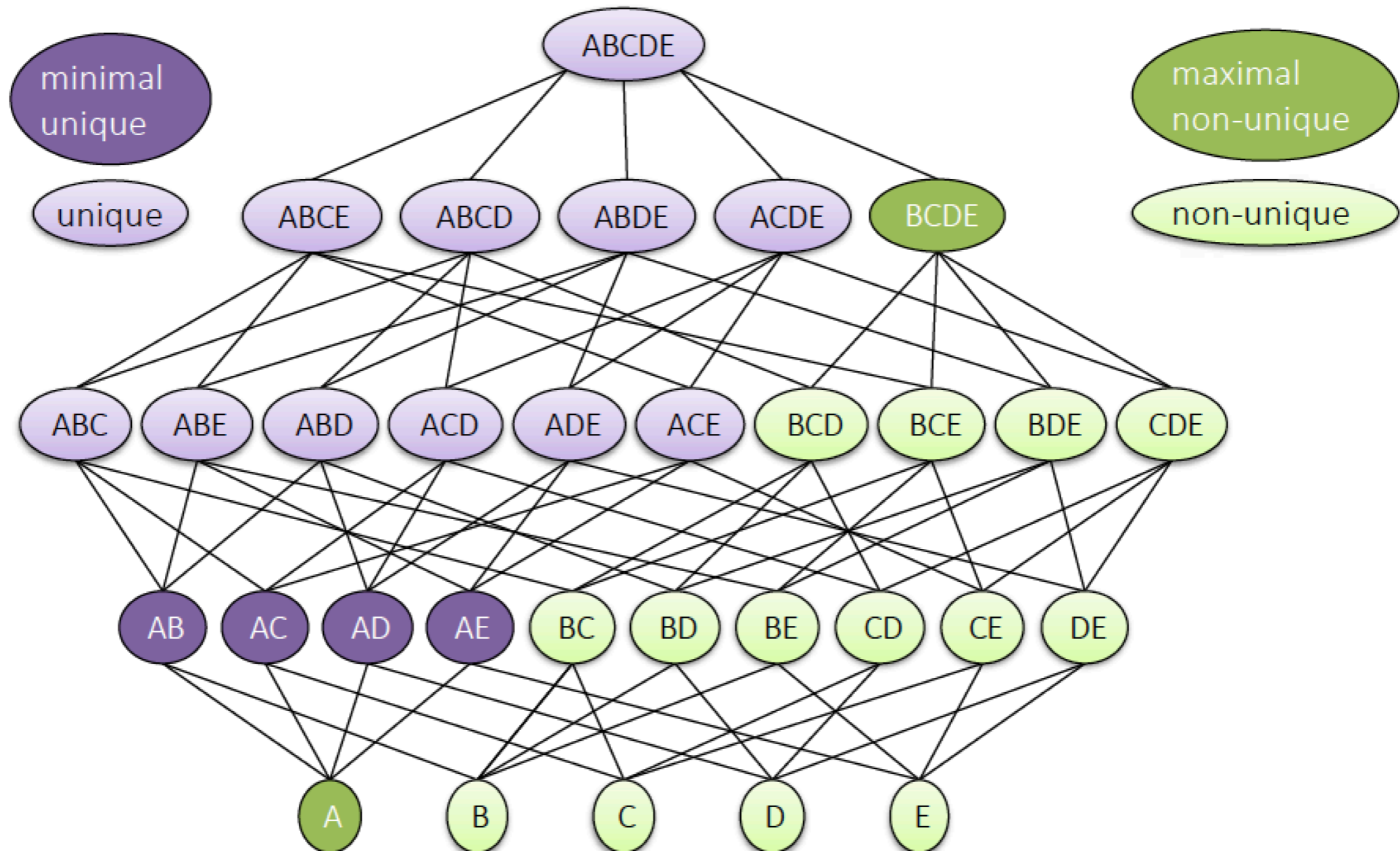
**X** is **minimal unique** if every subset **Y** of **X** is non-unique

**Y** is maximal non-unique if every superset **X** of **Y** is unique



ABCD

ABC    ABD    BCD    ACD

AB    AC    BC    BD    CD    AD

A    B    C    D

NYU

# Output

NYU

# From uniques to candidate keys

Given a relation schema **R** *(A, B, C, D)* and a relation instance **r**, a **unique column combination** is a set of attributes **X** whose **projection** contains no duplicates in **r**

$$Episodes(season, num, title, viewers)$$

| season | num | title | viewers |
|--------|-----|-------|---------|
| 1 | 1 | Winter is Coming | 2.2 M |
| 1 | 2 | The Kingsroad | 2.2 M |
| 2 | 1 | The North Remembers | 3.9 M |

A set of attributes is a **candidate key** for a relation if:
(1)       no two distinct tuples can have the value values for all key attributes (candidate key **uniquely identifies** a tuple), *and*
(2)  this is not true for any subset of the key attributes (candidate key **is minimal**)

**A minimal unique of a relation instance is a (possible) candidate key of the relation schema.** To find such possible candidate keys, find all minimal uniques in a given relation instance.

# Apriori-style uniques discovery

[Abedjan, Golab, Naumann; *SIGMOD 2017*]

**A minimal unique** of a relation instance is a **(possible) candidate key** of the relation schema.

**Algorithm Uniques** **// sketch, similar to HCA**

$U_1$ = *{1-uniques}*     $N_1$ = *{1-non-uniques}*

**for** ($k$ = 2; $N_{k-1} \neq \varnothing$; $k$++) **do**

  $C_k \leftarrow$ **candidate-gen**($N_{k-1}$)

  $U_k \leftarrow$ **prune-then-check** ($C_k$)

   // prune candidates with unique sub-sets, and with **value distributions that cannot be unique**

   // check each candidate in pruned set for uniqueness

  $N_k \quad \leftarrow C_k \setminus U_k$

 **end**

 **breadth-first bottom-up strategy for attribute lattice traversal**

return $U \leftarrow \bigcup_k U_k$;

NYU