



# Basic Data Structures and Algorithms for Data Profiling

8.5.2017  
Felix Naumann

## Overview

---

- 1. The lattice**
2. Apriori lattice traversal
3. Position List Indices
4. Bloom filters

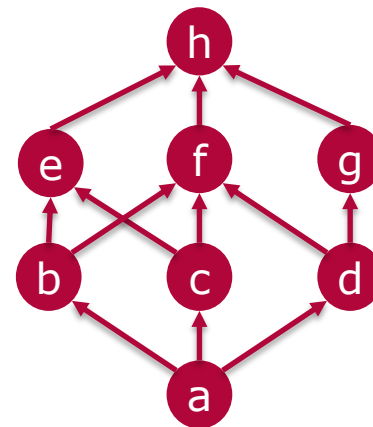
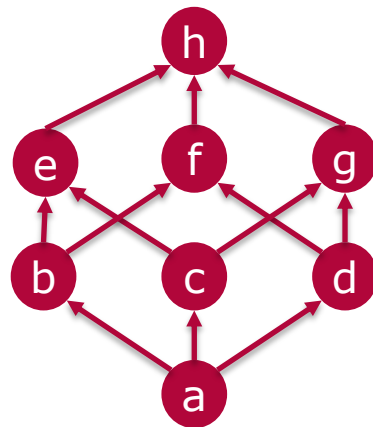
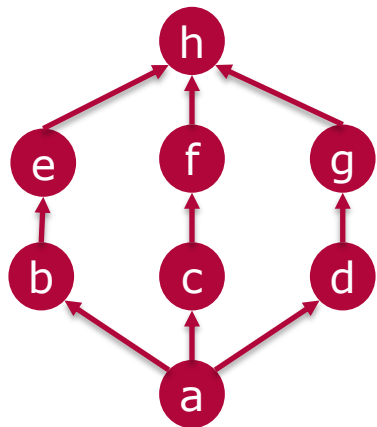


**Felix Naumann**  
Data Profiling  
Summer 2017

# Definitions

## ■ Lattice

- Partially ordered set (poset)
- Each pair of elements has unique supremum and infimum



Not a lattice

## ■ Hasse Diagram

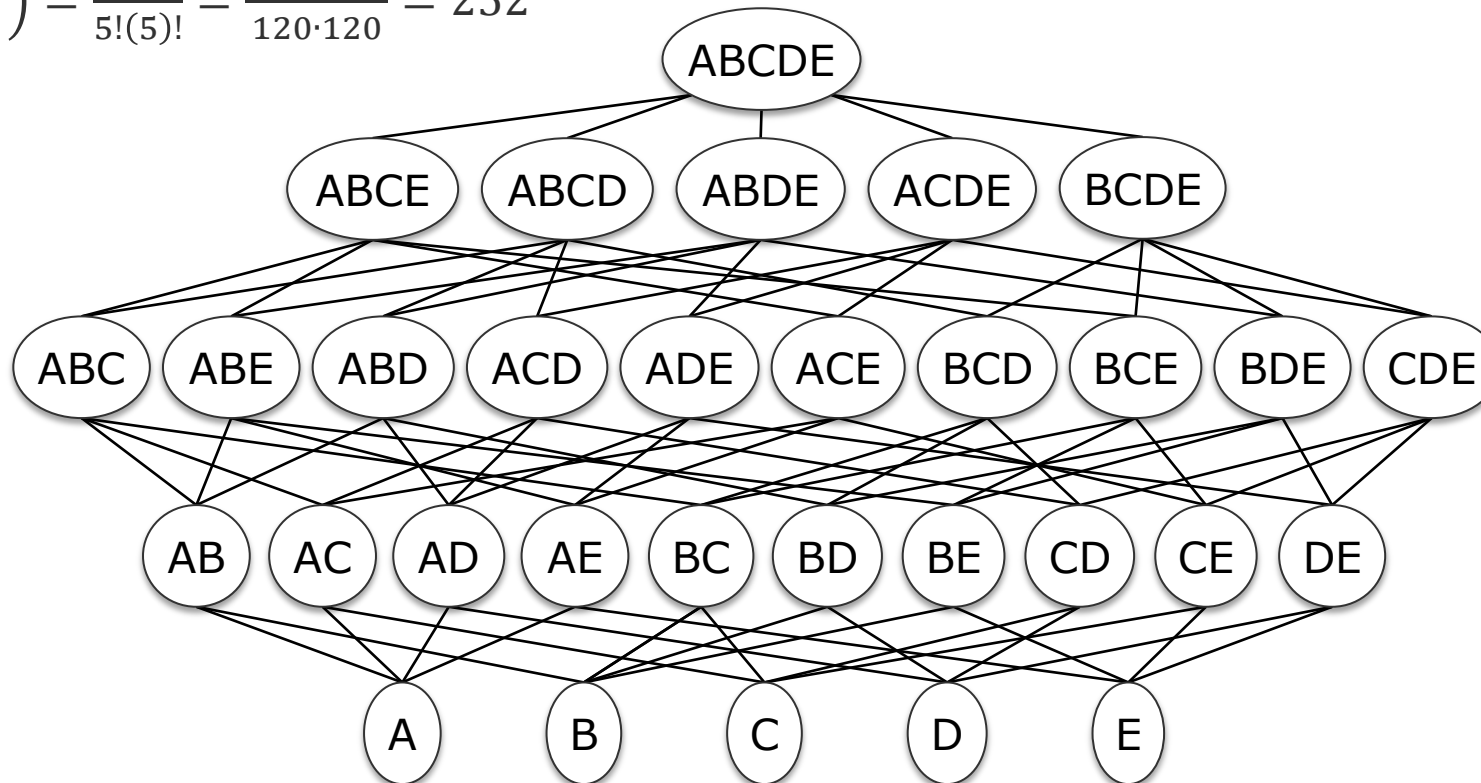
- Drawing of partially ordered set
- Each element is node
- Edges upward to smallest larger element
  - Upward: Arrows no longer necessary

# Basic lattice

- Represents each combination of elements

- $\binom{n}{k} = \frac{n!}{k!(n-k)!}$

- $\binom{10}{5} = \frac{10!}{5!(5)!} = \frac{3628800}{120 \cdot 120} = 252$



$$\binom{5}{5} = 1$$

$$\binom{5}{4} = 5$$

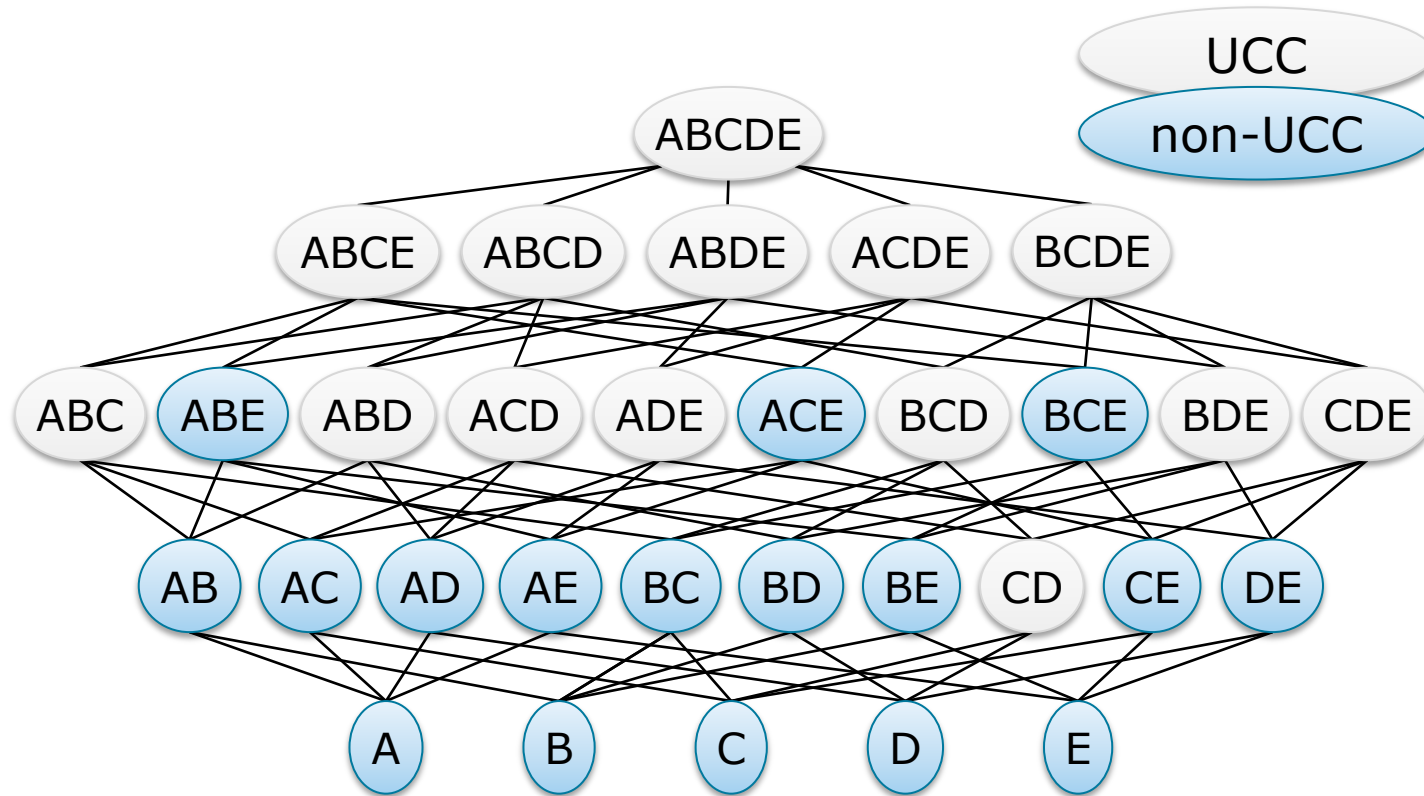
$$\binom{5}{3} = \frac{5 \cdot 4}{2}$$

$$\binom{5}{2} = \frac{5 \cdot 4 \cdot 3}{2 \cdot 3}$$

$$\binom{5}{1} = \frac{5 \cdot 4 \cdot 3 \cdot 2}{2 \cdot 3 \cdot 4}$$

Felix Naumann  
Data Profiling  
Summer 2017

# UCCs and non-UCCs in a Lattice

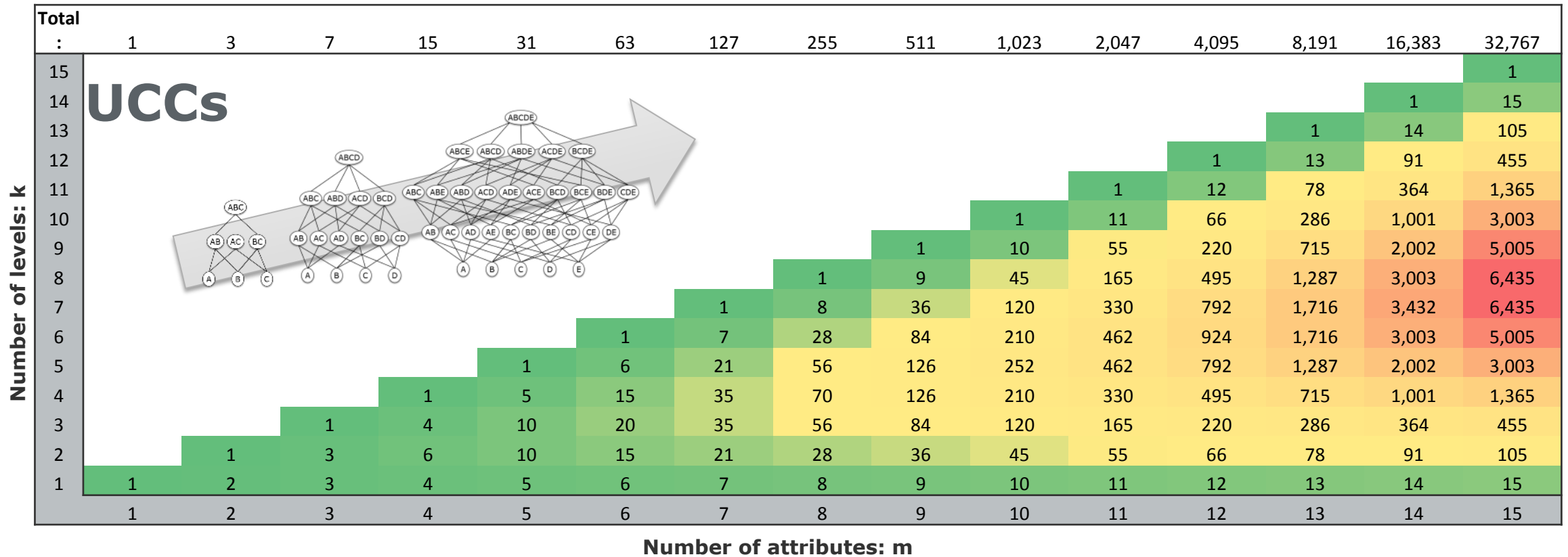


UCC  
non-UCC

Prune BCD, BDE and CDE, because CD is unique

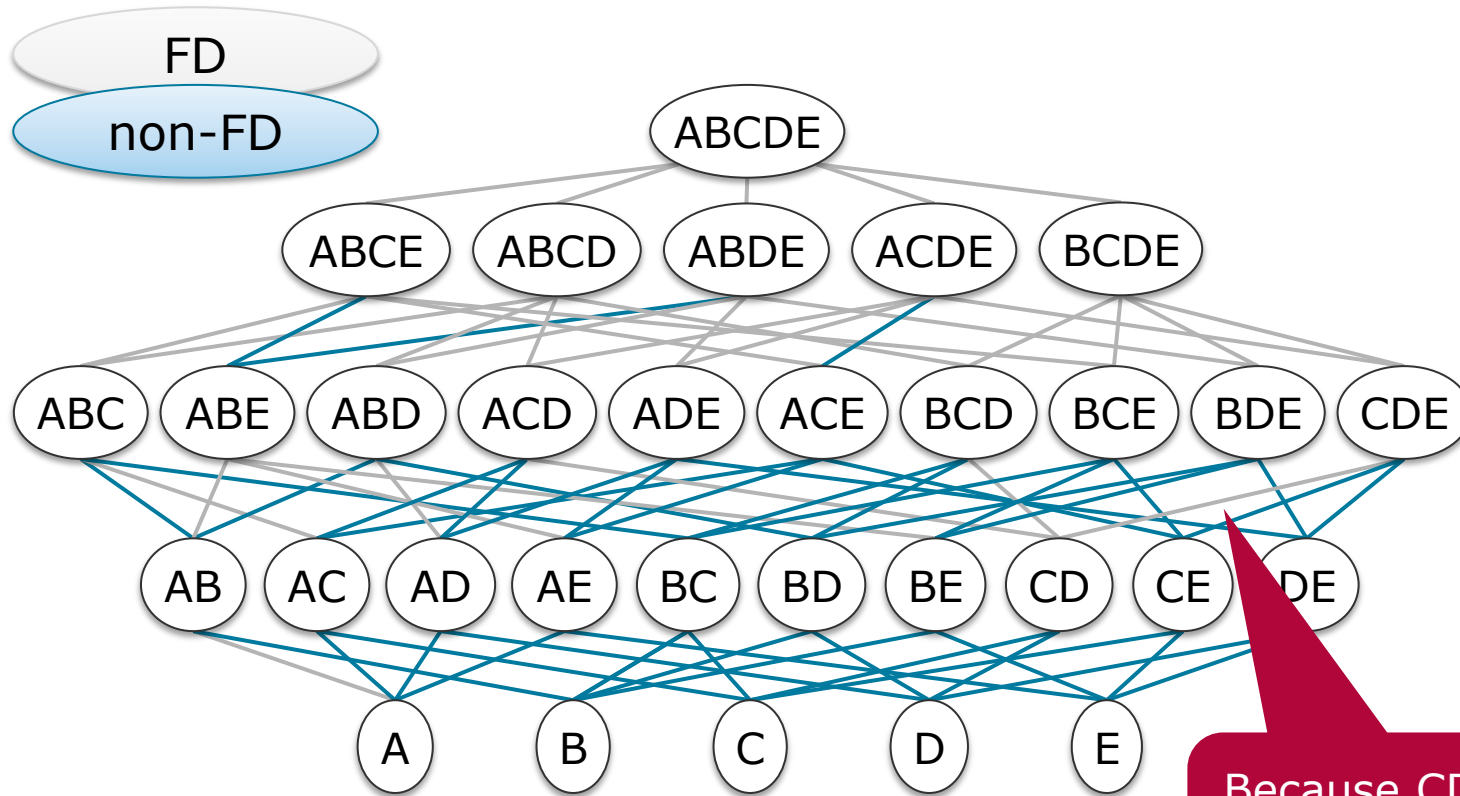
Felix Naumann  
Data Profiling  
Summer 2017

# Candidate Set Growth for UCCs



# FDs and non-FDs in a Lattice

- Candidates for level  $k$  (lhs) and  $m$  attributes:  $\binom{m}{k} \cdot (m - k)$



$$\binom{5}{4} \cdot (5 - 4) = 5$$

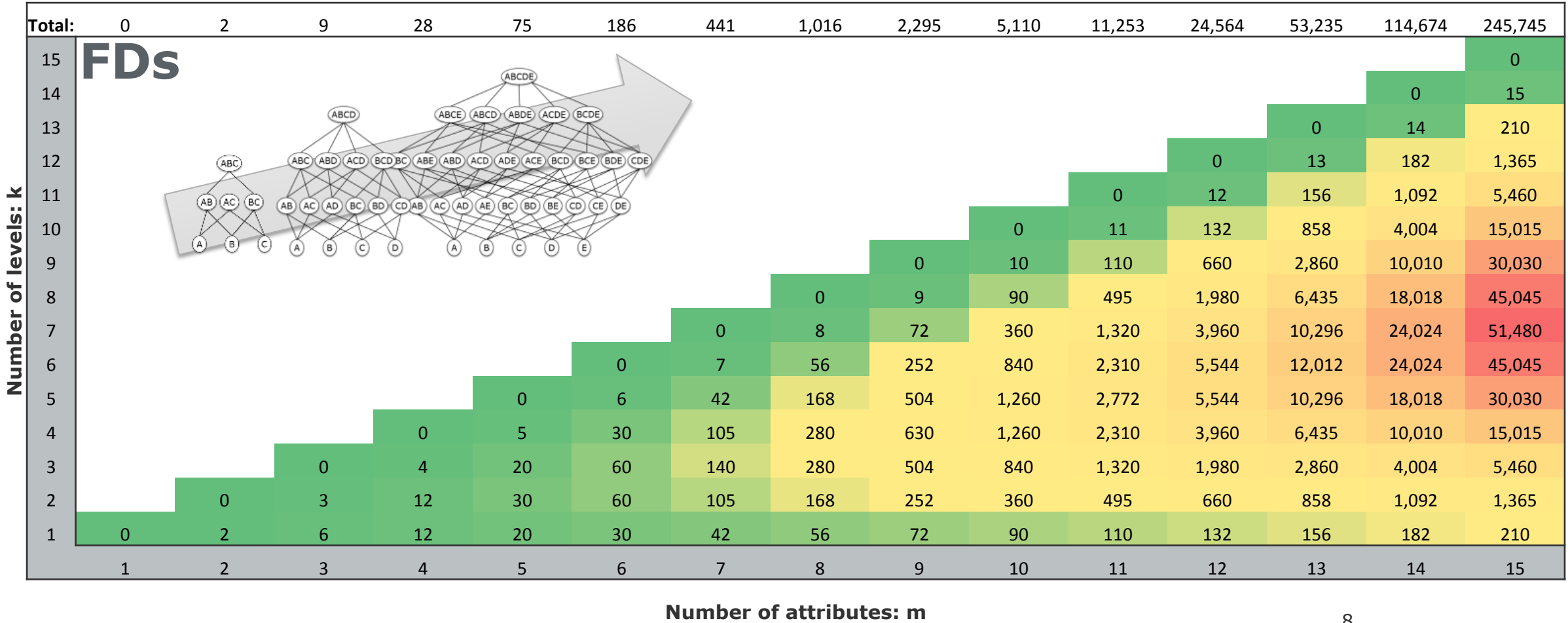
$$\binom{5}{3} \cdot (5 - 3) = 20$$

$$\binom{5}{2} \cdot (5 - 2) = 30$$

$$\binom{5}{1} \cdot (5 - 1) = 20$$

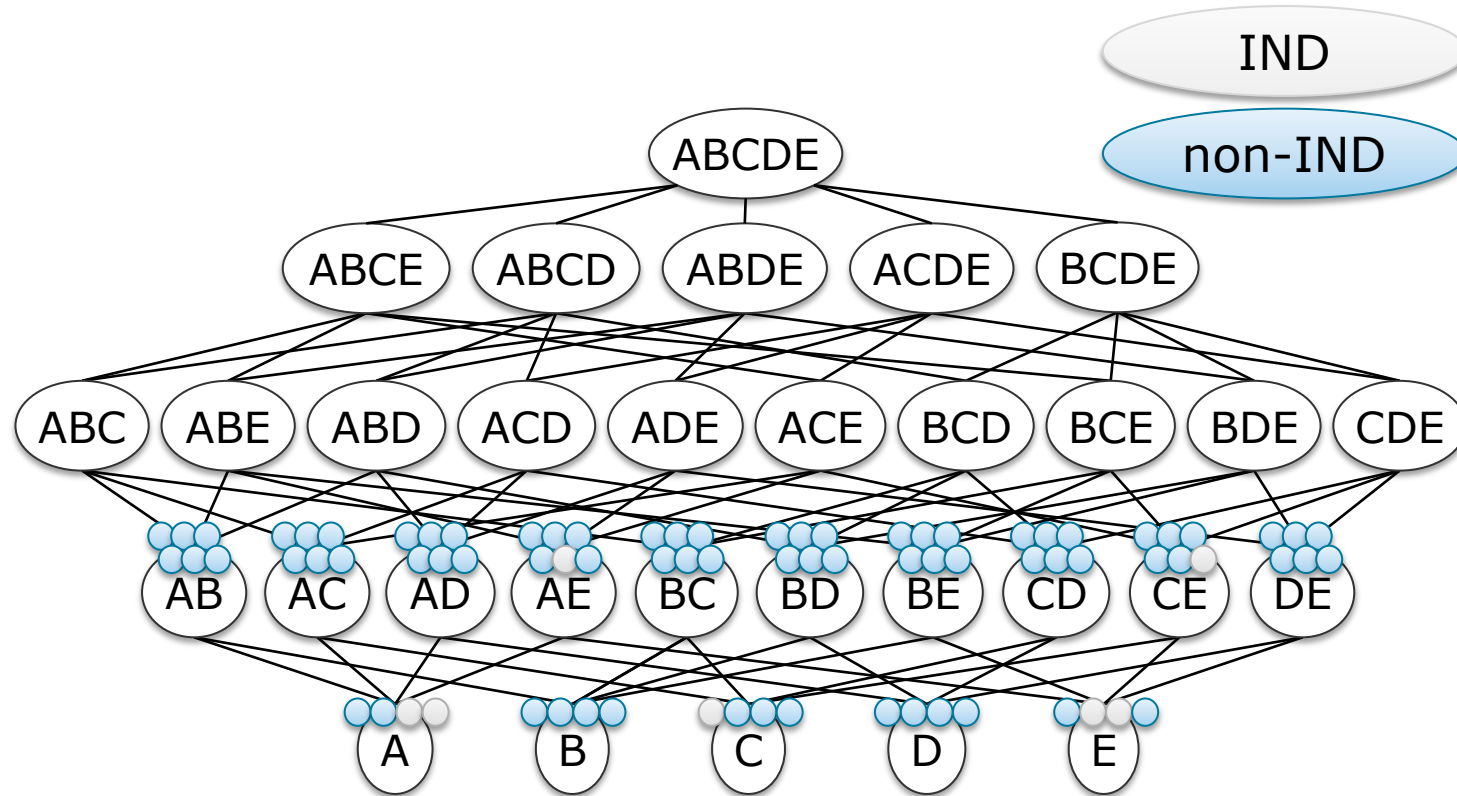
Because  $CD \rightarrow E$ , no need to check  $BCD \rightarrow E$  and others

# Candidate Set Growth for FDs



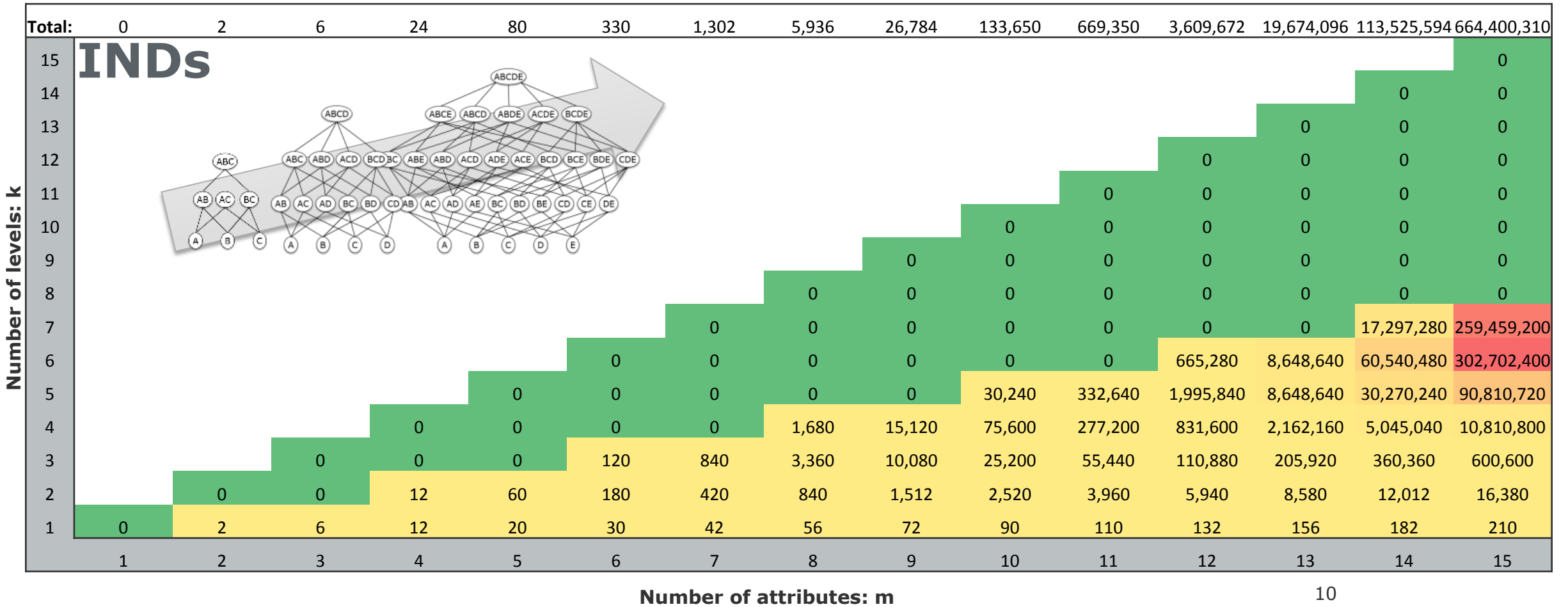


# INDs and non-INDs in a Lattice



- For  $X \subseteq Y$  we assume here that  $X \cap Y = \emptyset$ . Other IND definitions are possible.
- INDs live only in bottom half of lattice

# Candidate Set Growth for INDs



## Overview

1. The lattice
2. **Apriori lattice traversal**
3. Position List Indices
4. Bloom filters

R. Agrawal, R. Srikant  
"Fast Algorithms for Mining  
Association Rules"  
Proc. of the Int'l Conference  
on Very Large Databases  
(VLDB), 1994



Felix Naumann  
Data Profiling  
Summer 2017

## The Authors

---

### ■ Rakesh Agrawal

- 1983 Ph.D. at University of Wisconsin, Madison
- 1983-1989 Bell Laboratories
- 1989-2006 IBM Research
- 2006-2014 Microsoft Research
- Data Mining Pioneer



### ■ Ramakrishnan Srikant

- Ph.D. at University of Wisconsin, Madison
- IBM Almaden Research Center
- Since 2006 Distinguished Scientist at Google



**Felix Naumann**  
Data Profiling  
Summer 2017

# Terminology

---

- Itemset  $I$ , Transactionset  $D$
  
- Association rule
  - Statement of the form  $X \rightarrow Y$ , where  $X \subset I$ ,  $Y \subset I$  and  $X \cap Y = \{ \}$
  
- Support (frequency of a subset):
  - In what percentage of transactions is  $X \cup Y$  present?
  - $\rightarrow$  minSup
  
- Confidence (of rule):
  - In what percentage of transactions in which  $X$  appears does  $Y$  also appear?
  - $\rightarrow$  minConf
  
- Algorithm
  1. Find all itemsets with minSup (large itemsets).
  2. From these, derive association rules with minConf.

## Example

---

- Products:
  - $I = \{\text{Cola, Saft, Bier, Wein, Wasser, Schokolade, Brot, Schinken, Chips}\}$
- Transactions T1 - T9:
  - $T1 = \{\text{Saft, Cola, Bier}\}$
  - $T2 = \{\text{Saft, Cola, Wein}\}$
  - $T3 = \{\text{Saft, Wasser}\}$
  - $T4 = \{\text{Saft, Cola, Bier, Wein}\}$
  - $T5 = \{\text{Wasser}\}$
  - $T6 = \{\text{Schokolade, Cola, Chips}\}$
  - $T7 = \{\text{Cola, Bier}\}$
  - $T8 = \{\text{Schokolade, Schinken, Brot}\}$
  - $T9 = \{\text{Brot, Bier}\}$
- $\text{minSup} = 2$

## Key idea

---

- Problem
  - Many candidates (all subsets)
  - Checking transactions is expensive
  
- Idea: Support monotonically decreases with larger itemsets.
  - If a subset of some set  $M$  is small, then  $M$  itself is also small (= not large)
    1. Generate candidates using already discovered large itemsets.
    2. Delete all candidate that contain non-large subsets.
  - This suggests a bottom-up candidate generation approach

# Apriori

```
 $L_1 = \{ \text{large 1-itemsets} \}$   
For (  $k = 2; L_{k-1} \neq \emptyset ; k++$  ) do begin  
     $C_k = \text{apriori-gen}(L_{k-1} )$ ;  
    forall transactions  $t \in D$  do begin  
         $C_t = \text{subset}(C_k, t)$   
        forall candidates  $c \in C_t$  do  
             $c.\text{count}++$ ;  
        end  
    end  
     $L_k = \{ c \in C_k \mid c.\text{count} \geq \text{minsup} \}$   
end  
 $\text{Answer} = \bigcup_k L_k$ ;
```



# Apriori-Gen

- 2 Steps:

1. „Join“

```

insert into  $C_k$ 
select  $p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$ 
from  $L_{k-1} p, L_{k-1} q$ 
where  $p.item_1 = q.item_1, \dots, p.item_{k-2} = q.item_{k-2},$ 
       $p.item_{k-1} < q.item_{k-1};$ 

```

p and q are identical in the first  $k-2$  items

2. „Prune“

```

forall itemsets  $c \in C_k$  do
  forall  $(k-1)$ -subsets  $s$  of  $c$  do
    if  $(s \notin L_{k-1})$  then
      delete  $c$  from  $C_k;$ 

```

Delete all candidates with a subset not in  $L_{k-1}$

## Example

---

- $T1 = \{\text{Saft, Cola, Bier}\}$
- $T2 = \{\text{Saft, Cola, Wein}\}$
- $T3 = \{\text{Saft, Wasser}\}$
- $T4 = \{\text{Saft, Cola, Bier, Wein}\}$
- $T5 = \{\text{Wasser}\}$
- $T6 = \{\text{Schokolade, Cola, Chips}\}$
- $T7 = \{\text{Cola, Bier}\}$
- $T8 = \{\text{Schokolade, Schinken, Brot}\}$
- $T9 = \{\text{Brot, Bier}\}$
  
- $\text{minSup} = 2$
  
- Let  $L2 = \{\{\text{Cola, Saft}\}, \{\text{Cola, Bier}\}, \{\text{Cola, Wein}\}, \{\text{Saft, Bier}\}, \{\text{Saft, Wein}\}\}$
  
- Join:  $\{\text{Cola, Saft, Bier}\}, \{\text{Cola, Saft, Wein}\}, \{\text{Cola, Bier, Wein}\}, \{\text{Saft, Bier, Wein}\}$
  
- Prune deletes:  $\{\text{Cola, Bier, Wein}\}$  and  $\{\text{Saft, Bier, Wein}\}$

## Creation of rules (not relevant for lattice traversal)

---

- Let  $L$  be a large itemset and  $A \subseteq L$

$$A \rightarrow (L - A) \quad \leftrightarrow \quad \frac{\text{sup}(L)}{\text{sup}(A)} > \min \text{Conf}$$

- Idea: If  $XY \rightarrow Z$  not true, then  $X \rightarrow YZ$  is also not true
  - Because  $\text{sup}(XY) \leq \text{sup}(X)$
- Again inductive generation (over number of elements on rhs)
  - Check all rules with one element on rhs
  - From these, generate 2-element rhs's and check these
  - ...

## Extensions

---

- Consider quantity of items within transaction
- Consider order of transactions
- Consider taxonomy / categories
  - Outerwear → Hiking Boots
- Remove useless/non-actionable rules for pruning
- Determine interesting of discovered rules

## Overview

---

1. The lattice
2. Apriori lattice traversal
- 3. Position List Indices**
4. Bloom filters



Felix Naumann  
Data Profiling  
Summer 2017

## PLI Motivation

---

- UCC detection and FD detection need to know groups of rows with same value
  - UCCs: For a given attribute or attribute combination, are there ANY groups of such rows?
    - I.e., duplicates
  - FDs: For any such group (LHS), are there any dependent attributes (RHS) with also same values
  
- PLI for an attribute or set of attributes is a compact representation of such groups
  - Insight 1: Actual values are not needed, only row-ids
  - Insight 2: Singleton groups are not needed, only groups of size  $\geq 2$
  - Insight 3: PLIs for attribute sets can be efficiently built based on PLIs of their subsets.

# Position List Indices (PLIs)

	first	last
<b>0</b>	James	Smith
<b>1</b>	John	Smith
<b>2</b>	Robert	Johnson
<b>3</b>	John	Smith
<b>4</b>	John	Johnson
<b>5</b>	James	Williams

PLI: first
0, 5
1, 3, 4
2

PLI: last
0, 1, 3
2, 4
5

## Position List Indices (PLIs)

	first	last	PLI: first	PLI: last
<b>0</b>	James	Smith	0, 5	0, 1, 3
<b>1</b>	John	Smith	1, 3, 4	2, 4
<b>2</b>	Robert	Johnson		
<b>3</b>	John	Smith		
<b>4</b>	John	Johnson		
<b>5</b>	James	Williams		

- If PLI is empty: Column is unique
- Next step: Determine PLI for larger attribute sets
  - Idea: PLI intersection = intersection for each pair of groups
  - If intersection is empty (after removing singletons), attribute combination is unique



# Position List Indices (PLIs) – Intersection

PLI first		Probing table		PLI last		Value class		Row numbers		PLI first, last	
1	0, 5			1	0, 1, 3						
2	1, 3, 4	0	1	2	2, 4	1, 1	0			1, 3	
		1	2			2, 1	1, 3				
		2	0			2, 2	4				
		3	2								
		4	2								
		5	1								

0 for singletons

	Firstname	Lastname
0	James	Smith
1	<b>John</b>	<b>Smith</b>
2	Robert	Johnson
3	<b>John</b>	<b>Smith</b>
4	John	Johnson
5	Richard	Williams

Felix Naumann  
Data Profiling  
Summer 2017

# Key Error

- Key-error of column or column combination: Number of records to be removed to that column becomes unique.
- $\text{keyerror}(X) = \sum_{c \in \text{pli}(X)} (|c|) - |\text{pli}(X)|$
- If  $\text{keyerror}(X) = 0$ , X is unique
- Trick later: Use key-error to infer FDs
  - TANE

	first	last
0	James	Smith
1	John	Smith
2	Robert	Johnson
3	John	Smith
4	John	Johnson
5	James	Williams
<b>keyerror</b>	3	3
<b>keyerror</b>	1	

Felix Naumann  
Data Profiling  
Summer 2017

## Analysis of PLIs

---

- Intersection is associative
  - Potential to optimize intersection order
  - Build larger PLIs from smaller ones
  
- Intersection is commutative
  - Hash bigger PLI and probe smaller PLI
  - To reduce number of tests
  
- If there is enough main memory
  - Keep PLI of columns in main memory
  - Going up in the lattice requires only to probe the current PLI
    - Becomes increasingly fast when going up
    - <1ms for most combinations
  
- Going down the lattice
  - Unfortunately, PLIs do not help
  - Start from scratch

## Overview

---

1. The lattice
2. Apriori lattice traversal
3. Position List Indices
4. **Bloom filters**



Felix Naumann  
Data Profiling  
Summer 2017

## Searching for element (or confirming its existence)

- Problem: Is  $x$  an element of column  $A$ ?
  - Column  $A$  is large – must be stored on disk
- Idea: Store small representation of  $A$  in main memory
- Bloom filter: Developed 1970 by Burton Howard Bloom
  - > 5.500 citations
- Build Boolean hash-table  $T$  on  $A$  using all available main memory.
- Use  $T$  to mark whether an element is  $k$  in  $A$
- Test can fail, but only in one direction
  - If  $k \in T$ , we cannot be sure whether  $k \in A$ .
  - If  $k \notin T$ , we know that  $k \notin A$ .
- $T$  acts as filter: Bloom-filter
- Improvement: Use  $j$  independent hash-function
  - Improves false positive rate

> 6,000 citations

### Space/Time Trade-offs in Hash Coding with Allowable Errors

BURTON H. BLOOM  
*Computer Usage Company, Newton Upper Falls, Mass.*

In this paper trade-offs among certain computational factors in hash coding are analyzed. The paradigm problem considered is that of testing a series of messages one-by-one for membership in a given set of messages. Two new hash-coding methods are examined and compared with a particular conventional hash-coding method. The computational factors considered are the size of the hash area (space), the time required to identify a message as a nonmember of the given set (reject time), and an allowable error frequency.

The new methods are intended to reduce the amount of space required to contain the hash-coded information from that associated with conventional methods. The reduction in space is accomplished by exploiting the possibility that a small fraction of errors of commission may be tolerable in some applications, in particular, applications in which a large amount of data is involved and a core resident hash area is consequently not feasible using conventional methods.

In such applications, it is envisaged that overall performance could be improved by using a smaller core resident hash area in conjunction with the new methods and, when necessary, by using some secondary and perhaps time-consuming test to "catch" the small fraction of errors associated with the new methods. An example is discussed which illustrates possible areas of application for the new methods.

Analysis of the paradigm problem demonstrates that allowing a small number of test messages to be falsely identified as members of the given set will permit a much smaller hash area to be used without increasing reject time.

KEY WORDS AND PHRASES: hash coding, hash addressing, scatter storage, searching, storage layout, retrieval trade-offs, retrieval efficiency, storage efficiency

CR CATEGORIES: 3.73, 3.74, 3.79