



Detecting Inclusion Dependencies

15.6.2017
Felix Naumann

Overview

1. Inclusion Dependencies
2. SQL
3. De Marchi et al.
4. MIND
5. SPIDER
6. MANY



Felix Naumann
Data Profiling
Summer 2017

Inclusion Dependencies: Definition

- INDs (typically) involve more than one relation.
- Let D be a relational schema and let I be an instance of D .
- $R[A_1, \dots, A_n]$ denotes projection of I on attributes A_1, \dots, A_n of relation R :
 $R[A_1, \dots, A_n] = \pi_{A_1, \dots, A_n}(R)$

- IND $R[A_1, \dots, A_n] \subseteq S[B_1, \dots, B_n]$, where R, S are (possibly identical) relations of D .
 - Projection on R and S must have same number of attributes.
 - An instance I of D satisfies an IND if $I(R)[A_1, \dots, A_n] \subseteq I(S)[B_1, \dots, B_n]$
 - IND is maximal if $R[XA] \subseteq S[YB]$ is invalid for any $A \in R, B \in S$
 - Values of R : “dependent values”
 - Values of S : “referenced values”

- Task: Find all maximal, non-trivial INDs
 - Typical assumptions: No repeating attributes, disjoint LHS and RHS

Example

- Each Title in Showings should appear as a Title in Movies
 - Showings[Title] ⊆ Movie[Title]

<i>Movies</i>	<i>Title</i>	<i>Director</i>	<i>Actor</i>
	The Birds	Hitchcock	Hedren
	The Birds	Hitchcock	Taylor
	Bladerunner	Scott	Hannah
	Apocalypse Now	Coppola	Brando

<i>Showings</i>	<i>Theater</i>	<i>Screen</i>	<i>Title</i>	<i>Snack</i>
	Rex	1	The Birds	coffee
	Rex	1	The Birds	popcorn
	Rex	2	Bladerunner	coffee
	Rex	2	Bladerunner	popcorn
	Le Champo	1	The Birds	tea
	Le Champo	1	The Birds	popcorn
	Cinoche	1	The Birds	Coke
	Cinoche	1	The Birds	wine
	Cinoche	2	Bladerunner	Coke
	Cinoche	2	Bladerunner	wine
	Action Christine	1	The Birds	tea
	Action Christine	1	The Birds	popcorn

Felix Naumann
Data Profiling
Summer 2017

- Aka. “referential integrity”
 - Referenced attributes need not be a key (or unique)
 - Foreign key: helps prune candidates

Inference rules for INDs

- Reflexivity: $R[X] \subseteq R[X]$

- Projection:
 - $R[A_1, \dots, A_n] \subseteq S[B_1, \dots, B_n] \Rightarrow R[A_{i_1}, \dots, A_{i_m}] \subseteq S[B_{i_1}, \dots, B_{i_m}]$
for each sequence i_1, \dots, i_m of Integers in $\{1, \dots, n\}$

- Transitivity:
 - $R[X] \subseteq S[Y]$ and $S[Y] \subseteq T[Z] \Rightarrow R[X] \subseteq T[Z]$
 - Example: “transitive foreign keys” for 1:1 relationships

IND types

- Unary INDs
 - INDs on single attributes: $R[A] \subseteq S[B]$
- n-ary INDs
 - INDs on multiple attributes: $R[X] \subseteq S[Y]$
 - $|X| = |Y|$
- Partial INDs
 - IND $R[A] \subseteq S[B]$ is satisfied for $x\%$ of all tuples in R
 - IND $R[A] \subseteq S[B]$ is satisfied for all but x tuples in R
- Approximate INDs
 - IND $R[A] \subseteq S[B]$ is satisfied with probability p .
 - Based on sampling or other heuristics

Examples

- Unary: $R[C] \subseteq S[F]$
- N-ary: $R[B,C] \subseteq S[G,F]$
- Partial: $R[A] \subseteq_{75\%} S[F]$
- Approximate: $R[BA] \subseteq S[GH]$

R	A	B	C
	1	x	1
	2	x	1
	3	y	2
	5	z	4

S	F	G	H
	1	x	1
	2	y	3
	3	z	4
	4	z	4

Felix Naumann
Data Profiling
Summer 2017

Further IND types

■ Prefix/Suffix INDs

- IND $R[A] \subseteq S[B]$ is satisfied after removing a fixed (or variable) prefix/suffix from each value of A.
- Twist: A dependent value can now match multiple referenced values

■ Example

A	B
bbc	b
	bb

$A \subseteq_s B$
(suffix with variable length)

Further IND types

- Conditional INDs
 - Only useful for partial INDs
 - More another time

Catalog

Unit cost	DBName	ProdID
200 USD	ToyDB	17
50 EUR	ToyDB	18
1000 QAR	FashionDB	19

ToyDB

EntityID	further data
17	abcd...
18	efgh...

FashionDB

EntityID	further data
18	abcd...
19	efgh...

Unary IND detection complexity

Name	Type	Equatorial diameter	Mass	Orbital radius	Orbital period	Rotation period	Confirmed moons	Rings	Atmosphere
Mercury	Terrestrial	0.382	0.06	0.47	0.24	58.64	0	no	minimal
Venus	Terrestrial	0.949	0.82	0.72	0.62	-243.02	0	no	CO ₂ , N ₂
Earth	Terrestrial	1.000	1.00	1.00	1.00	1.00	1	no	N ₂ , O ₂ , Ar
Mars	Terrestrial	0.532	0.11	1.52	1.88	1.03	2	no	CO ₂ , N ₂ , Ar
Jupiter	Giant	11.209	317.8	5.20	11.86	0.41	67	yes	H ₂ , He
Saturn	Giant	9.449	95.2	9.54	29.46	0.43	62	yes	H ₂ , He
Uranus	Giant	4.007	14.6	19.22	84.01	-0.72	27	yes	H ₂ , He
Neptune	Giant	3.883	17.2	30.06	164.8	0.67	14	yes	H ₂ , He

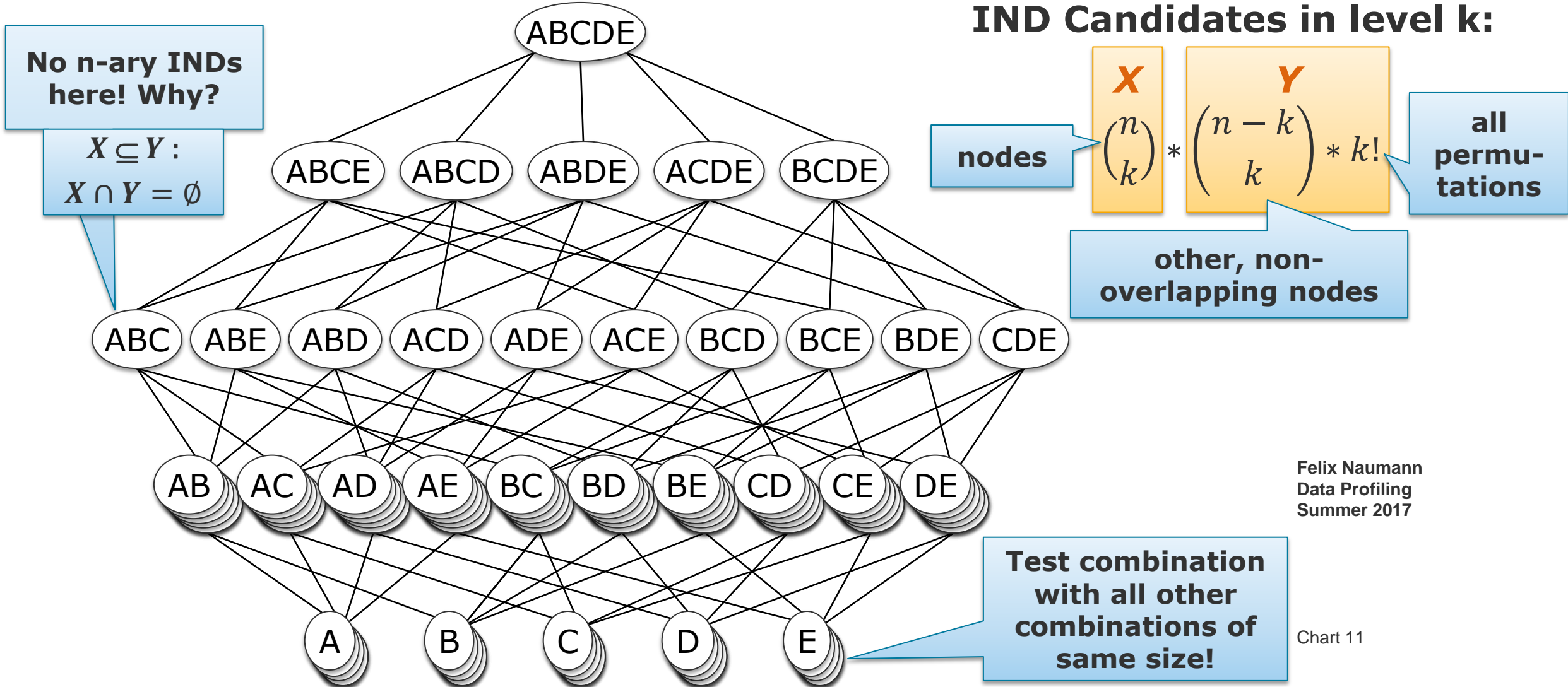
Complexity: $O(n^2-n)$
for n attributes

Example:
10 attr ~ 90 checks
1,000 attr ~ 999,000 checks

$n_{IND} \gg n_{UCC}$

- Name \subseteq Type ?
- Type \subseteq Name ?
- Mass \subseteq Name ?
- Name \subseteq Equatorial_diameter ?
- Type \subseteq Equatorial_diameter ?
- Mass \subseteq Type ?
- Name \subseteq Mass ?
- Type \subseteq Mass ?
- Mass \subseteq Equatorial_diameter ?
- Name \subseteq Orbital_radius ?
- Type \subseteq Orbital_radius ?
- ...
- Name \subseteq Orbital_period ?
- Type \subseteq Orbital_period ?
- Name \subseteq Rotation_period ?
- Type \subseteq Rotation_period ?
- Name \subseteq Confirmed_moons ?
- Type \subseteq Confirmed_moons ?
- Name \subseteq Rings ?
- Type \subseteq Rings ?
- Name \subseteq Atmosphere ?
- Type \subseteq Atmosphere ?

N-ary IND detection complexity

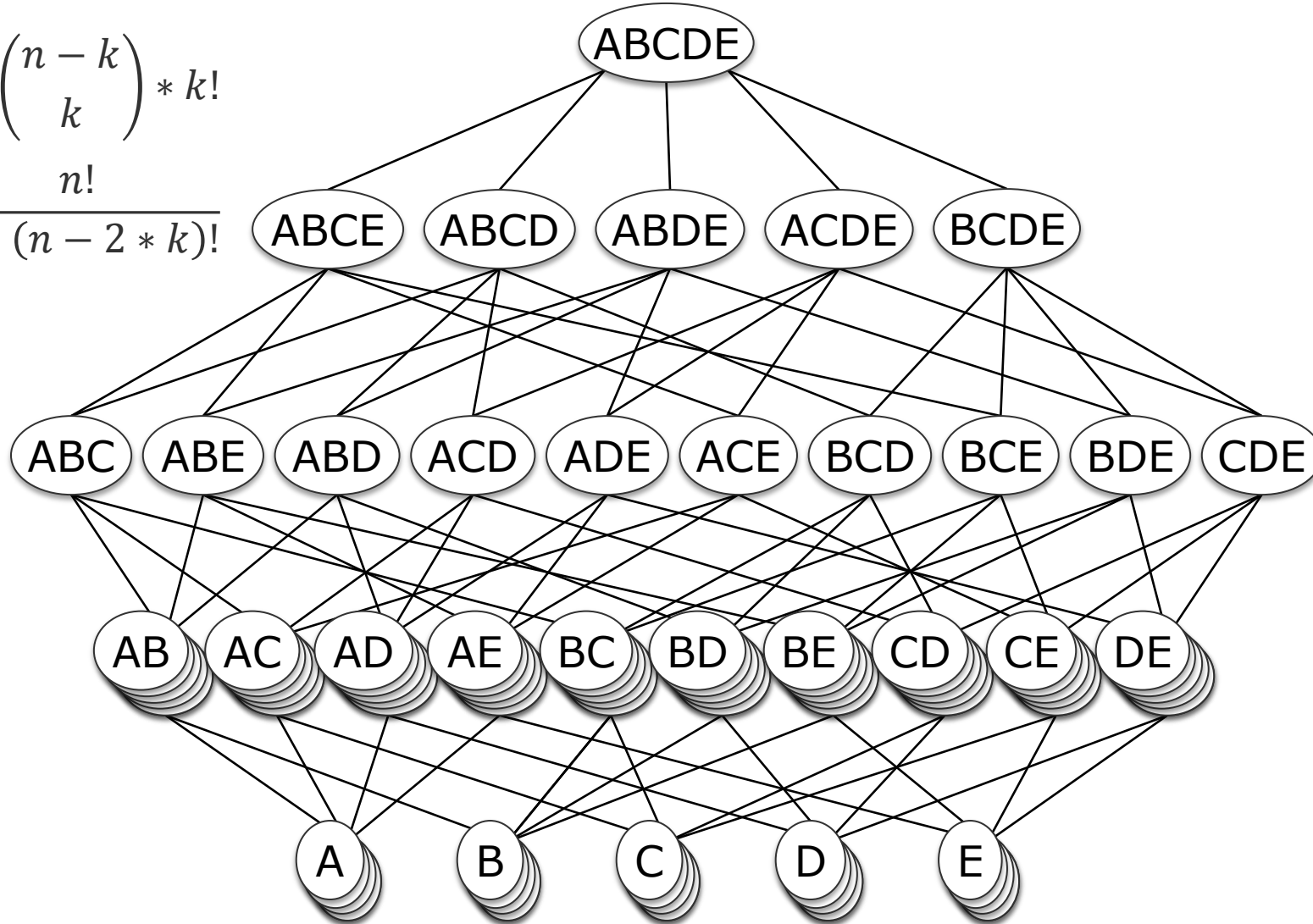


Felix Naumann
Data Profiling
Summer 2017

N-ary IND detection complexity

$$\binom{n}{k} * \binom{n-k}{k} * k!$$

$$= \frac{n!}{k! * (n - 2 * k)!}$$



$$\binom{5}{5} * \binom{5-5}{5} * 5! \sim 0$$

$$\binom{5}{4} * \binom{5-4}{4} * 4! \sim 0$$

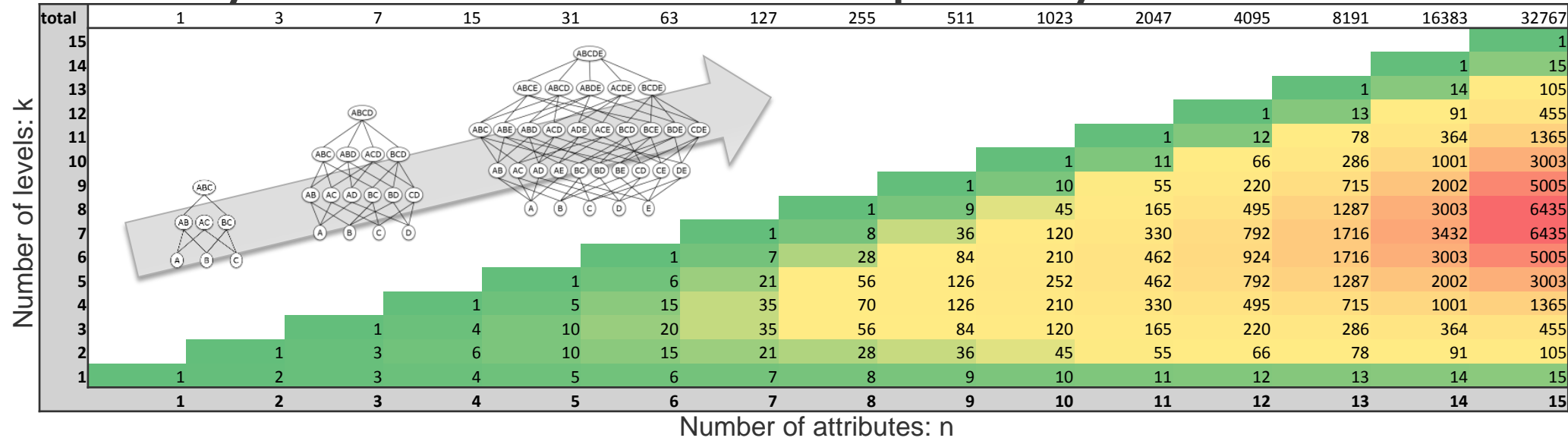
$$\binom{5}{3} * \binom{5-3}{3} * 3! \sim 0$$

$$\binom{5}{2} * \binom{5-2}{2} * 2! = 60$$

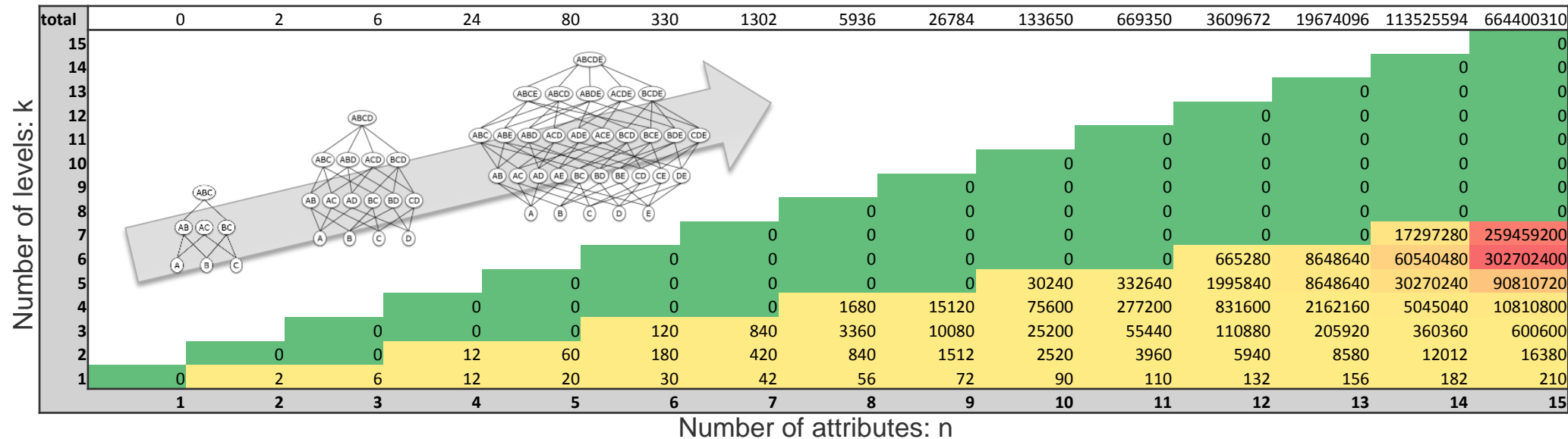
$$\binom{5}{1} * \binom{5-1}{1} * 1! = 20 = n^2 - n$$

BINDER – divide & conquer based IND detection

N-ary IND detection complexity



Unique Column Combinations



Inclusion Dependencies

Felix Naumann
Data Profiling
Summer 2017

Chart 13

Motivation for IND discovery

- General insight into data
- Detect unknown foreign keys
- Example
 - PDB: Protein Data Bank
 - OpenMMS provides relational schema
 - Parses protein and nucleic acid macromolecular structure data from the standard mmCIF format.
 - 175 tables with primary key constraints
 - 2705 attributes
 - But: Not a single foreign key constraint!

```

_pdbx_poly_seq_scheme.pdb_strand_id
_pdbx_poly_seq_scheme.pdb_ins_code
_pdbx_poly_seq_scheme.hetero
A 1 1 DC 1 1 1 DC C A . n
A 1 2 DC 2 2 2 DC C A . n
A 1 3 DG 3 3 3 DG G A . n
A 1 4 DT 4 4 4 DT T A . n
A 1 5 DA 5 5 5 DA A A . n
A 1 6 DC 6 6 6 DC C A . n
A 1 7 DG 7 7 7 DG G A . n
A 1 8 DT 8 8 8 DT T A . n
A 1 9 DA 9 9 9 DA A A . n
A 1 10 DC 10 10 10 DC C A . n
A 1 11 DG 11 11 11 DG G A . n
A 1 12 DG 12 12 12 DG G A . n
#
loop_
_refine_B_iso.class
_refine_B_iso.details
_refine_B_iso.treatment
_refine_B_iso.pdbx_refine_id
'ALL ATOMS' TR isotropic 'X-RAY DIFFRACTION'
'ALL WATERS' TR isotropic 'X-RAY DIFFRACTION'
#
loop_
_refine_occupancy.class
_refine_occupancy.treatment
_refine_occupancy.pdbx_refine_id
'ALL ATOMS' fix 'X-RAY DIFFRACTION'
'ALL WATERS' fix 'X-RAY DIFFRACTION'
#
loop_
_pdbx_version.entry_id
_pdbx_version.revision_date
_pdbx_version.major_version
_pdbx_version.minor_version
_pdbx_version.revision_type
_pdbx_version.details
116D 2008-05-22 3 2 'Version format complianc
116D 2011-07-13 4 0000 'Version format complianc
#
software_name NUCLSO

```

Felix Naumann
Data Profiling
Summer 2017

Motivation for IND discovery

- Ensembl – genome database
 - Shipped as MySQL dump files
 - More than 200 tables
 - Not a single foreign key constraint!

- Web tables: No schema, no constraints, but many connections

- Why are FKs missing?
 - Lack of support for checking foreign key constraints in the host system
 - Example: Oracle did not support FKs up to v6
 - Fear that checking such constraints would impede database performance
 - Lack of database knowledge within the development team
 - Dirty data prevents setting the constraint

Overview

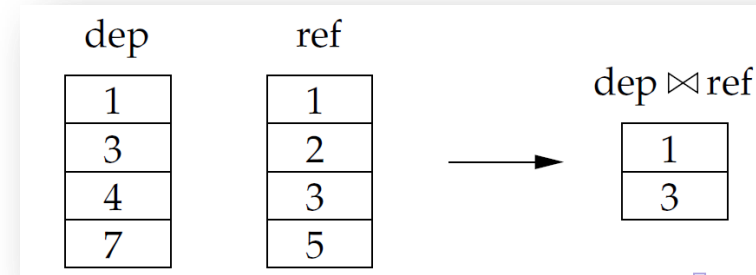
1. Inclusion Dependencies
2. **SQL**
3. De Marchi et al.
4. MIND
5. SPIDER
6. MANY



Felix Naumann
Data Profiling
Summer 2017

The JOIN

- `SELECT COUNT(depColumn) AS numDeps`
`FROM depTable`
- `SELECT COUNT(*) AS matchedDeps`
`FROM depTable JOIN refTable`
`ON depTable.depColumn = refTable.refColumn`



- Assumption: ref is a unique column
- $(\text{numDeps} = \text{matchedDeps}) \Leftrightarrow \text{depColumn} \subseteq \text{refColumn}$

- Missed opportunity
 - DBMS could stop early: As soon as we observe a dependent value without a join partner

The EXCEPT

- ```
SELECT count(*) AS unmatchedDeps FROM
 ((SELECT to char (depColumn)
 FROM depTable
 WHERE depColumn IS NOT NULL
 EXCEPT
 SELECT to char (refColumn)
 FROM refTable
)
 FETCH FIRST 1 ROWS ONLY
)
```
- $\text{unmatchedDeps} = 0 \Leftrightarrow \text{depColumn} \subseteq \text{refColumn}$

## The “Antijoin”

---

- `SELECT COUNT(*) AS unmatched  
FROM R  
WHERE A IS NOT NULL  
AND A NOT IN  
    (SELECT B FROM S)  
FETCH FIRST 1 ROWS ONLY`

- `depColumn  $\subseteq$  refColumn  
 $\Leftrightarrow$  unmatched = 0`

- `SELECT COUNT(*) AS unmatched  
FROM R  
WHERE A IS NOT NULL  
AND NOT EXISTS  
    (SELECT * FROM S WHERE  
      R.A=S.B)  
FETCH FIRST 1 ROWS ONLY`

- `depColumn  $\subseteq$  refColumn  
 $\Leftrightarrow$  unmatched = 0`

## Measurements (2006)

- High efficiency of joins in DBMS
- Inability of DBMS optimizer to move STOP operator into inner queries
- Overall problems
  - Still too slow
  - One SQL statement per attribute pair
  - Each attribute joined  $n$  times (many sorts/hashes)

|                  | CATH      | SCOP      | UniProt   | TPC-H     | PDB       |
|------------------|-----------|-----------|-----------|-----------|-----------|
| DB size          | 20 MB     | 17,5 MB   | 900 MB    | 1.3 GB    | 2.8 GB    |
| # attributes     | 25        | 22        | 68        | 61        | 1,215     |
| # IND candidates | 68        | 43        | 910       | 477       | 139,807   |
| # INDS           | 0         | 11        | 36        | 33        | 4,972     |
| join             | 6 s       | 7 s       | 9 m 04 s  | 25 m 02 s | 16 h 14 m |
| except           | 15 m 27 s | 16 m 05 s | 27 m 35 s | 1 h 09 m  | –         |
| not in           | 5 s       | 52 m 11 s | 6 h 33 m  | 7 m 45 s  | –         |
| not exists       | 5 s       | 6 s       | 3 m 57 s  | 7 m 51 s  | 10 h 20 m |

Table 4.1: Runtime performance of the SQL approaches. IND candidates are restricted to cover unique referenced attributes. We used only a fraction of PDB.

## Discussion on data profiling experiments

---

- What can we assume? What is the scenario?
  - Index every column
  - Statistics for each table and column
    - Example: min and max values for pruning
  - Where is the data originally
    - In a database
    - In files
  - Do I count importing the data?
    - Could then do statistics on the fly

## Overview

---

1. Inclusion Dependencies
2. SQL
- 3. De Marchi et al.**
4. MIND
5. SPIDER
6. MANY



Felix Naumann  
Data Profiling  
Summer 2017

## Pruning for Unary INDs

### ■ Datatype

- Incompatible datatypes cannot be included in one-another

### ■ Statistics

- $\text{card}(R[A]) > \text{card}(S[B]) \Rightarrow R.A \not\subseteq S.B$ 
  - $\text{card}(R[X]) > \text{card}(S[Y]) \Rightarrow R.X \not\subseteq S.Y$
- $\text{max}(R[A]) > \text{max}(S[B]) \Rightarrow R.A \not\subseteq S.B$
- $\text{min}(R[A]) < \text{min}(S[B]) \Rightarrow R.A \not\subseteq S.B$

### ■ Bloom filters

|                                | UniProt<br>900 MB | TPC-H<br>1.3 GB | PDB     |         |
|--------------------------------|-------------------|-----------------|---------|---------|
|                                |                   |                 | 2.8 GB  | 32 GB   |
| # attributes                   | 68                | 61              | 1,215   | 1,297   |
| # IND candidates               | 1,393             | 877             | 219,106 | 245,562 |
| # satisfied INDs               | 36                | 33              | 4,972   | 5,431   |
| # attributes in INDs           | 31                | 20              | 448     | 478     |
| distinct                       |                   |                 |         |         |
| # IND candidates               | 910               | 477             | 139,807 | 157,818 |
| # attributes in IND candidates | 68                | 58              | 1,208   | 1,297   |
| distinct & max                 |                   |                 |         |         |
| # IND candidates               | 541               | 295             | 72,016  | 83,321  |
| # attributes in IND candidates | 59                | 54              | 997     | 1,080   |
| distinct & min                 |                   |                 |         |         |
| # IND candidates               | 345               | 275             | 61,920  | 68,664  |
| # attributes in IND candidates | 54                | 57              | 990     | 1,042   |
| distinct & max & min           |                   |                 |         |         |
| # IND candidates               | 174               | 137             | 22,655  | 25,821  |
| # attributes in IND candidates | 49                | 52              | 670     | 709     |

## Data preprocessing

---

- Key idea: For a given domain, associate each value with every attribute having this value.
  - Create binary relation  $B \subseteq \text{Values} \times \text{Attributes}$  with  $(v, A) \in B$  iff  $v \in \pi_A(R)$
  - Analogy: Inverted index

Efficient Algorithms for Mining Inclusion Dependencies,  
Fabien De Marchi, Stéphane Lopes, and Jean-Marc Petit,  
In: EDBT 2002

Felix Naumann  
Data Profiling  
Summer 2017



# Example

- Three domains: int, real, and string
- Example for domain "int"
  - Values = {1,2,3,4,6,7,9}
  - Attributes = {A,C,E,G,K}
  - Examples for relation B: (1,A), (1,E), (1,K)
- Can build relation with single full scan of each base relation

r

| A | B | C | D    |
|---|---|---|------|
| 1 | X | 3 | 11.0 |
| 1 | X | 3 | 12.0 |
| 2 | Y | 4 | 11.0 |
| 1 | X | 3 | 13.0 |

s

| E | F | G | H    |
|---|---|---|------|
| 1 | X | 3 | 11.0 |
| 2 | Y | 4 | 12.0 |
| 4 | Z | 6 | 14.0 |
| 7 | W | 9 | 14.0 |

t

| I    | J    | K | L |
|------|------|---|---|
| 11.0 | 11.0 | 1 | X |
| 12.0 | 12.0 | 2 | Y |
| 11.0 | 14.0 | 4 | Z |
| 11.0 | 9.0  | 7 | W |
| 13.0 | 13.0 | 9 | R |

# Example „Extraction contexts“

r

| A | B | C | D    |
|---|---|---|------|
| 1 | X | 3 | 11.0 |
| 1 | X | 3 | 12.0 |
| 2 | Y | 4 | 11.0 |
| 1 | X | 3 | 13.0 |

s

| E | F | G | H    |
|---|---|---|------|
| 1 | X | 3 | 11.0 |
| 2 | Y | 4 | 12.0 |
| 4 | Z | 6 | 14.0 |
| 7 | W | 9 | 14.0 |

t

| I    | J    | K | L |
|------|------|---|---|
| 11.0 | 11.0 | 1 | X |
| 12.0 | 12.0 | 2 | Y |
| 11.0 | 14.0 | 4 | Z |
| 11.0 | 9.0  | 7 | W |
| 13.0 | 13.0 | 9 | R |

int

| V | U       |
|---|---------|
| 1 | A E K   |
| 2 | A E K   |
| 3 | C G     |
| 4 | C E G K |
| 6 | G       |
| 7 | E K     |
| 9 | G K     |

real

| V    | U       |
|------|---------|
| 9.0  | J       |
| 11.0 | D H I J |
| 12.0 | D H I J |
| 13.0 | D I J   |
| 14.0 | H J     |

string

| V | U     |
|---|-------|
| R | L     |
| X | B F L |
| Y | B F L |
| Z | F L   |
| W | F L   |

# IND discovery

- Insight: If all values of attribute  $A$  can be found in values of  $B$  (i.e.,  $A \subseteq B$ ), then by construction  $B$  will be present in all lines of the binary relation containing  $A$ .

$$A \subseteq B \iff B \in \bigcap_{v \in \mathbb{V} | (v, A) \in \mathbb{B}} \{C \in \mathbb{U} \mid (v, C) \in \mathbb{B}\}$$

int

| V | U       |
|---|---------|
| 1 | A E K   |
| 2 | A E K   |
| 3 | C G     |
| 4 | C E G K |
| 6 | G       |
| 7 | E K     |
| 9 | G K     |

real

| V    | U       |
|------|---------|
| 9.0  | J       |
| 11.0 | D H I J |
| 12.0 | D H I J |
| 13.0 | D I J   |
| 14.0 | H J     |

string

| V | U     |
|---|-------|
| R | L     |
| X | B F L |
| Y | B F L |
| Z | F L   |
| W | F L   |

# IND discovery algorithm

**Input:** the triplet  $\mathbb{V}, \mathbb{U}, \mathbb{B}$ , associated with  $\mathbf{d}$  and  $t$ .

**Output:**  $\mathcal{I}_1$  the set of unary INDs verified by  $\mathbf{d}$  between attributes of type  $t$ .

1: for all  $A \in \mathbb{U}$  do  $rhs(A) = \mathbb{U}$ ;

All attributes are ref candidates

2: for all  $v \in \mathbb{V}$  do

3:     for all  $A$  s.t.  $(v, A) \in \mathbb{B}$  do

4:          $rhs(A) = rhs(A) \cap \{B \mid (v, B) \in \mathbb{B}\}$ ;

Remove candidates

5: for all  $A \in \mathbb{U}$  do

6:     for all  $B \in rhs(A)$  do

7:          $\mathcal{I}_1 = \mathcal{I}_1 \cup \{A \subseteq B\}$ ;

Generate output

8: return  $\mathcal{I}_1$ .

| int          |              |
|--------------|--------------|
| $\mathbb{V}$ | $\mathbb{U}$ |
| 1            | A E K        |
| 2            | A E K        |
| 3            | C G          |
| 4            | C E G K      |
| 6            | G            |
| 7            | E K          |
| 9            | G K          |

| real         |              |
|--------------|--------------|
| $\mathbb{V}$ | $\mathbb{U}$ |
| 9.0          | J            |
| 11.0         | D H I J      |
| 12.0         | D H I J      |
| 13.0         | D I J        |
| 14.0         | H J          |

| string       |              |
|--------------|--------------|
| $\mathbb{V}$ | $\mathbb{U}$ |
| R            | L            |
| X            | B F L        |
| Y            | B F L        |
| Z            | F L          |
| W            | F L          |

# IND discovery algorithm: Example

- Step 0:  $\text{rhs}(A) = \dots = \text{rhs}(K) = \{A, C, E, G, K\}$
- Step 1 ( $v=1$ ):  $\text{rhs}(A)=\{A, E, K\}$ ,  $\text{rhs}(E)=\{A, E, K\}$ ,  $\text{rhs}(K)=\{A, E, K\}$ ,  $\text{rhs}(C) = \text{rhs}(G) = \{A, C, E, G, K\}$
- Step 2 ( $v=2$ ): unchanged
- Step 3 ( $v=3$ ):  $\text{rhs}(C)=\{C, G\}$ ,  $\text{rhs}(G)=\{C, G\}$
- ...
- Step 9:  $\text{rhs}(A)=\{A, E, K\}$ ,  $\text{rhs}(C)=\{C, G\}$ ,  $\text{rhs}(E)=\{E, K\}$ ,  $\text{rhs}(G)=\{G\}$ ,  $\text{rhs}(K)=\{K\}$
- $A \subseteq E$ ,  $A \subseteq K$ ,  $C \subseteq G$ , and  $E \subseteq K$

int

| V | U       |
|---|---------|
| 1 | A E K   |
| 2 | A E K   |
| 3 | C G     |
| 4 | C E G K |
| 6 | G       |
| 7 | E K     |
| 9 | G K     |

real

| V    | U       |
|------|---------|
| 9.0  | J       |
| 11.0 | D H I J |
| 12.0 | D H I J |
| 13.0 | D I J   |
| 14.0 | H J     |

string

| V | U     |
|---|-------|
| R | L     |
| X | B F L |
| Y | B F L |
| Z | F L   |
| W | F L   |

Question: Why distinguish domains?

## Overview

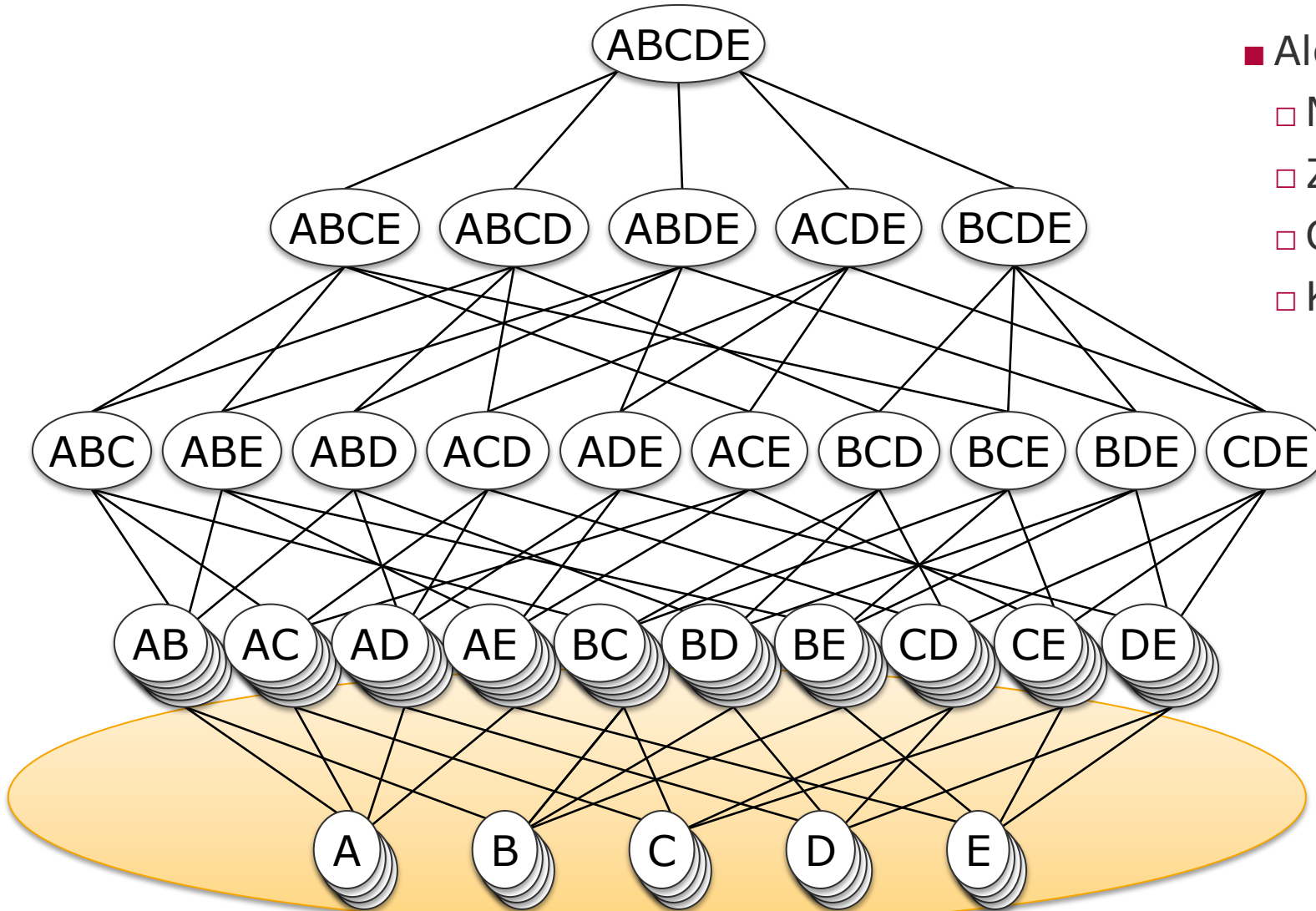
---

1. Inclusion Dependencies
2. SQL
3. De Marchi et al.
4. **MIND**
5. SPIDER
6. MANY



Felix Naumann  
Data Profiling  
Summer 2017

# N-ary IND detection



- Algorithms for n-ary detection:
  - MIND (apriori bottom-up)
  - ZigZag (bottom-up and top-down)
  - Clim (closed item set mining)
  - Koeller (aggressive bottom up)

Felix Naumann  
Data Profiling  
Summer 2017

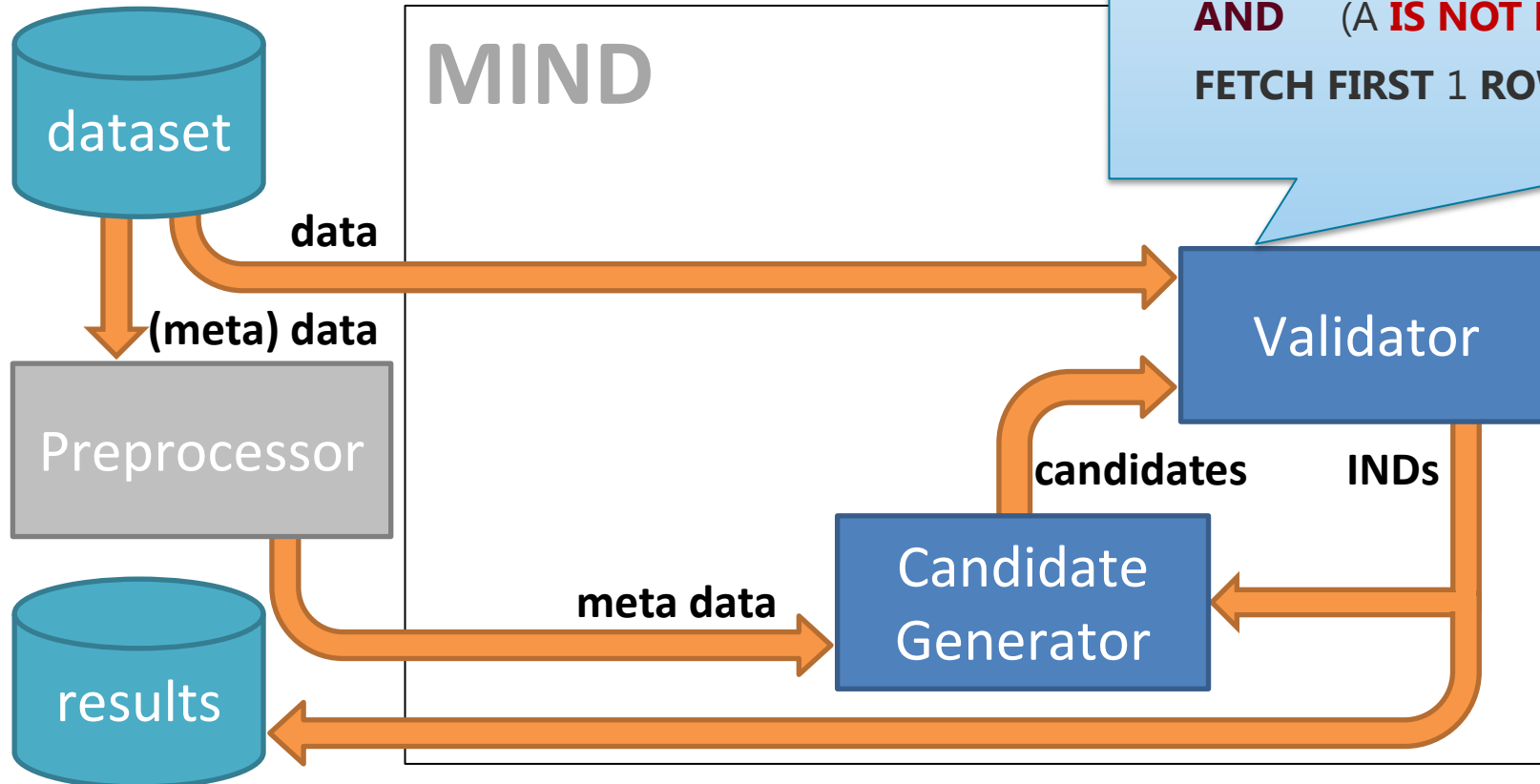
## Pruning for INDs

---

- $R[A_1, \dots, A_n] \subseteq S[B_1, \dots, B_n] \Rightarrow R[A_{i1}, \dots, A_{im}] \subseteq S[B_{i1}, \dots, B_{im}]$
  
- Pruning down
  - Example:  $R[AB] \subseteq S[DE] \Rightarrow R[A] \subseteq S[D]$  and  $R[B] \subseteq S[E]$
  
- Pruning up
  - Example:  $R[AB] \not\subseteq S[DE] \Rightarrow R[ABC] \not\subseteq S[DEF]$
  - Apriori: Use only satisfied INDs to generate next level candidates
  
- Pruning laterally
  - $R[AB] \subseteq S[DE] \Rightarrow R[BA] \subseteq S[ED]$
  - Define permutation strategy
    - Lexicographic ordering of attribute labels for LHS



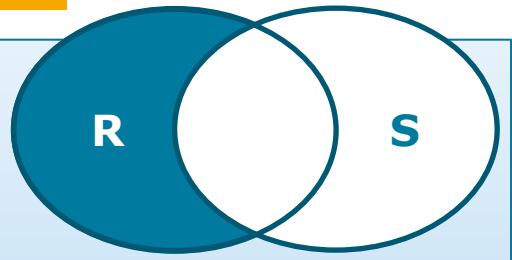
# MIND – workflow



**MIND**

**Validate**  $A, B \subseteq C, D$ :

```
SELECT R.A, R.B, S.C, S.D
FROM R LEFT OUTER JOIN S ON A = C AND B = D
WHERE (C IS NULL AND D IS NULL)
AND (A IS NOT NULL OR B IS NOT NULL)
FETCH FIRST 1 ROWS ONLY;
```



```
SELECT depTable.AAGE, depTable.ACLSWKR,
 depTable.ADTIND, depTable.ADTOCC,
 depTable.AGI, depTable.AHGA
FROM CENSUS6 depTable LEFT OUTER JOIN
CENSUS refTable ON
 depTable.AAGE = refTable.AAGE AND
 depTable.ACLSWKR = refTable.ACLSWKR AND
 depTable.ADTIND = refTable.ADTIND AND
 depTable.ADTOCC = refTable.ADTOCC AND
 depTable.AGI = refTable.AGI AND
 depTable.AHGA = refTable.AHGA
WHERE refTable.AAGE IS NULL AND
 refTable.ACLSWKR IS NULL AND
 refTable.ADTIND IS NULL AND
 refTable.ADTOCC IS NULL AND
 refTable.AGI IS NULL AND
 refTable.AHGA IS NULL
AND (depTable.AAGE IS NOT NULL OR
 depTable.ACLSWKR IS NOT NULL OR
 depTable.ADTIND IS NOT NULL OR
 depTable.ADTOCC IS NOT NULL OR
 depTable.AGI IS NOT NULL OR
 depTable.AHGA IS NOT NULL)
FETCH FIRST 1 ROWS ONLY;
```

# The MIND Algorithm

---

## Algorithm 2 MIND

---

**Input:**  $\mathbf{d}$  a database, and  $\mathcal{I}_1$  the set of unary INDs satisfied by  $\mathbf{d}$ .

**Output:** Inclusion dependencies satisfied by  $\mathbf{d}$

```

1: $\mathcal{C}_2 := \text{GenNext}(\mathcal{I}_1)$;
2: $i := 2$;
3: while $\mathcal{C}_i \neq \emptyset$ do
4: forall $I \in \mathcal{C}_i$ do
5: if $\mathbf{d} \models I$ then
6: $\mathcal{I}_i := \mathcal{I}_i \cup \{I\}$;
7: $\mathcal{C}_{i+1} := \text{GenNext}(\mathcal{I}_i)$;
8: $i := i + 1$;
9: end while
10: return $\cup_{j < i} \mathcal{I}_j$

```

Start with unary INDs from previous algorithm

Use SQL queries

Felix Naumann  
Data Profiling  
Summer 2017

# The MIND Algorithm – Apriori-based generation

**Algorithm 3** *GenNext* : Generation of candidate INDs of size  $i + 1$

**Input:**  $\mathcal{I}_i$ , inclusion dependencies of size  $i$ .

**Output:**  $\mathcal{C}_{i+1}$ , sequence of candidate inclusion dependencies of size  $i+1$

```

1: insert into \mathcal{C}_{i+1}
2: select $p.lhs.rel$ “[$p.lhs[1], p.lhs[2], \dots, p.lhs[i], q.lhs[i]$]” “ \subseteq ”
 $p.rhs.rel$ “[$p.rhs[1], p.rhs[2], \dots, p.rhs[i], q.rhs[i]$]”
3: from \mathcal{I}_i p, \mathcal{I}_i q
4: where $p.lhs.rel = q.lhs.rel$ and $p.rhs.rel = q.rhs.rel$
5: and $p.lhs[1] = q.lhs[1]$ and $p.rhs[1] = q.rhs[1]$
6: and ...
7: and $p.lhs[i - 1] = q.lhs[i - 1]$ and $p.rhs[i - 1] = q.rhs[i - 1]$
8: and $p.lhs[i] < q.lhs[i]$
9: and $p.rhs[i] <> q.rhs[i]$
10: for all $I \in \mathcal{C}_{i+1}$ do
11: for all $J < I$ and J of size i do
12: if $J \notin \mathcal{I}_i$ then
13: $\mathcal{C}_{i+1} = \mathcal{C}_{i+1} \setminus \{I\}$
14: end if
15: end for
16: end for

```

All LHS (and RHS) attributes from same relation

Same prefix of size  $i-1$

Avoid superfluous permutations

Disallow repeated attributes

Generation

Pruning

Remove candidate if not all (ordered) subsets are in lower level

Felix Naumann  
Data Profiling  
Summer 2017

## Overview

---

1. Inclusion Dependencies
2. SQL
3. De Marchi et al.
4. MIND
- 5. SPIDER**
6. MANY



Felix Naumann  
Data Profiling  
Summer 2017

## Making use of order

- Idea: Order each column only once
  - As index in DBMS
  - Or as sorted columns in individual files
    - `SELECT DISTINCT A FROM R ORDER BY A`
    - Store result in file

| A | B |
|---|---|
| 1 | 1 |
| 3 | 2 |
| 4 | 3 |
| 7 | 5 |

- Simulate merging procedure (merge join, merge sort)
  - Move cursor along both columns
  - Stop after first dependent value that is not in referenced attribute

## Testing a single IND candidate

- 2 ordered lists of distinct values: depValues and refValues
- **while** (depValues has next)
  - currentDep = depValues.next();
  - **while** (true)
    - currentRef = refValues.next();
    - **if** (currentDep = currentRef) **then** break;
    - **else if** (currentDep < currentRef) **then return** false;
    - **else if** not(refValues has next) **then return** false;
- **return** true;

| A | B | C |
|---|---|---|
| 1 | 1 | 1 |
| 3 | 2 | 2 |
| 4 | 3 | 3 |
|   | 5 | 4 |

## Brute force approach

---

- Sequentially check each column pair
- At least: Re-use order for each attribute
  
- Problem: Run through data multiple times
  - $A \subseteq C$
  - $A \subseteq D$
  - $B \subseteq C$
  - $B \subseteq D$

| A | B | C | D |
|---|---|---|---|
| 1 | 2 | 1 | 1 |
| 3 | 3 | 3 | 2 |
| 4 | 5 | 7 | 3 |
|   |   | 8 | 5 |

# SPIDER

## Single Pass Inclusion DEpendency Recognition

---

### ■ Main ideas

- Test all IND-candidate pairs simultaneously.
  - In both directions
- Read attribute values only once.
- Stop test of an IND-candidate after first counter-example.
- Reduce number of value comparisons by specialized data structure.
- No need to build inverted index.

### ■ Two steps:

- Sort and distinct all attribute's values and write them to disk
  - For each attribute: **SELECT DISTINCT A FROM R ORDER BY A**
- Test all IND candidate pairs simultaneously

- Jana Bauckmann and Ulf Leser and Felix Naumann. *Efficiently Computing Inclusion Dependencies for Schema Discovery*. In Proceedings of the International Conference on Data Engineering Workshops (ICDE workshops), 2006.
- Jana Bauckmann, Ulf Leser, Felix Naumann, Véronique Tietz: Efficiently Detecting Inclusion Dependencies. In: ICDE , 2007.



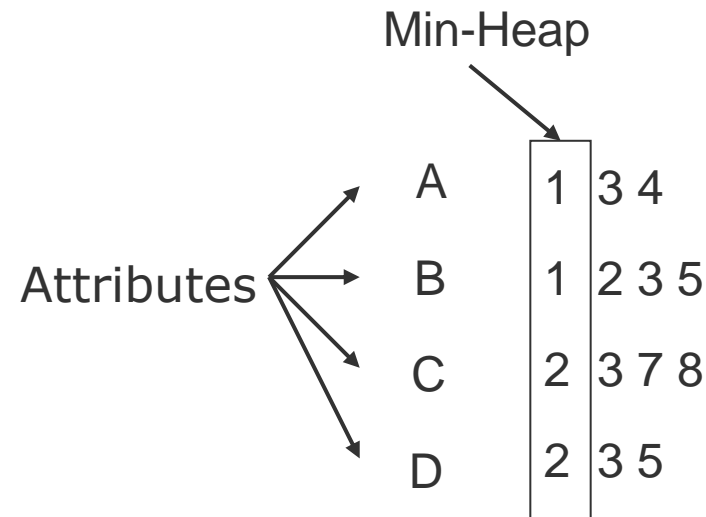
# SPIDER

- Parallel generation and test of all IND candidates
  - Reads each value at most once
- Challenge: Synchronize reading of values of all attributes
  - Each dependent attribute value influences when a referenced attribute value can be read.
  - Each referenced attribute value influences when a dependent attribute value can be read.
- Move cursor  $r$  on a referenced file R when all cursors to dependent files point to values that are greater than the current value pointed to by  $r$ .
- Move a cursor  $d$  on a dependent file D one step further, when  $d$ 's value is smaller than all values currently pointed to in referenced files.

|   | dep1 |   | dep2 |   | ref1 |   | ref2 |
|---|------|---|------|---|------|---|------|
| → | 1    | → | 2    | → | 1    | → | 1    |
|   | 3    |   | 3    |   | 2    |   | 3    |
|   | 4    |   | 5    |   | 3    |   | 7    |
|   |      |   |      |   | 5    |   | 8    |

# SPIDER: Idea

- All values within each attribute are sorted.
- Attributes themselves are sorted by current minimum value (in a min-heap).
  - Order changes all the time
- IND candidates represented as a list for each dependent attribute, containing all referenced attributes.



# SPIDER by example

- In each step: Intersect „attributes to process“ with each refs list of previous step
  - Intersection only with refs-list of “attributes to process”
- Attributes with empty refs-list can be removed from consideration

attributes A, B, C

|          |          |          |
|----------|----------|----------|
| <b>A</b> | <b>B</b> | <b>C</b> |
| s        |          | s        |
| t        | t        | t        |
| x        |          |          |
| y        | y        | y        |
|          |          | z        |

|        | <b>attributes to process</b> | <b>dep A refs</b> | <b>dep B refs</b> | <b>dep C refs</b> |
|--------|------------------------------|-------------------|-------------------|-------------------|
| Init   |                              | B,C               | A,C               | A,B               |
| Step 1 | A,C                          | C                 | A,C               | A                 |
| Step 2 | A,B,C                        | C                 | A,C               | A                 |
| Step 3 | A                            | ∅                 | A,C               | A                 |
| Step 4 | A,B,C                        | ∅                 | A,C               | A                 |
| Step 5 | C                            | ∅                 | A,C               | ∅                 |

Felix Naumann  
Data Profiling  
Summer 2017

# SPIDER results (data in DBMS)

|                     | <b>UniProt</b> | <b>TPC-H</b> | <b>PDB</b>    |              |
|---------------------|----------------|--------------|---------------|--------------|
| <b>DB size</b>      | 900 MB         | 1.3 GB       | 2.8 GB        | 32GB         |
| <b># attributes</b> | 68             | 61           | 1215          | 1297         |
| <b># IND cand.</b>  | 910            | 477          | 139,807       | 157,818      |
| <b># INDs</b>       | 36             | 33           | 4,972         | 5,431        |
| join                | 9m04s          | 25m02s       | 16h14m        | > 7 days     |
| Marchi et al.       | 9h 58m         | -            | -             | -            |
| Brute force         | 2m11s          | 6m30s        | 3h29m         | 19h51m       |
| <b>SPIDER</b>       | <b>1m51s</b>   | <b>6m25s</b> | <b>23m36s</b> | <b>6h07m</b> |

Felix Naumann  
Data Profiling  
Summer 2017

## Analysis and Extension

---

- Complexity:  $O(nt \log t)$  comparisons for  $n$  attributes and  $t$  tuples
  - Sorting all columns:  $O(nt \log t)$
  - Insertion into minHeap (of size  $n$ ):  $O(\log n)$  for each value
    - $O(nt \log n)$  for all values
  - Popping from heap again  $O(nt \log n)$
  - Intersections in constant time (bit vectors), so  $O(nt)$  for all
  - Assuming  $t \gg n$ :  $O(nt \log t)$
  - I/O complexity is also dominated by sorting
- Extension for partial INDs
  - During intersection:
    - Count how many times intersection removed and attributes.
    - Remove only after  $k$  unsuccessful intersections

## Overview

---

1. Inclusion Dependencies
2. SQL
3. De Marchi et al.
4. MIND
5. SPIDER
6. **MANY**



Felix Naumann  
Data Profiling  
Summer 2017

**Detecting Inclusion Dependencies on Very Many Tables**  
F. Tschirschnitz, T. Papenbrock, F. Naumann (TODS 2017)






# Planets

## Planeten Tabelle

Drucken

Geschrieben von Denise, Veröffentlicht in Planeten

### Die Planeten des Sonnensystems im Vergleich

|          | Abstand zur Sonne in Mio km                                                         | Durchmesser in km                                                                   | Masse in kg                                                                         | Umlaufzeit um die Sonne in Tagen                                                    | Geschwindigkeit                                                                     |
|----------|-------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------|
| Merkur   | 58                                                                                  | 4.879                                                                               | $3,30 \times 10^{23}$                                                               | 88                                                                                  | 172.332 km/h                                                                        |
| Venus    | 108                                                                                 | 12.103                                                                              | $4,87 \times 10^{24}$                                                               | 225                                                                                 | 126.072 km/h                                                                        |
| Erde     | 150                                                                                 | 12.734                                                                              | $5,97 \times 10^{24}$                                                               | 365                                                                                 | 107.208 km/h                                                                        |
| Mars     | 228                                                                                 | 6.772                                                                               | $6,42 \times 10^{23}$                                                               | 687                                                                                 | 86.868 km/h                                                                         |
| Jupiter  | 778                                                                                 | 138.346                                                                             | $1,90 \times 10^{27}$                                                               | 4329<br>(11 Jahre, 314 Tage)                                                        | 47.052 km/h                                                                         |
| Saturn   | 1.433                                                                               | 114.632                                                                             | $5,69 \times 10^{26}$                                                               | 10751<br>(29 Jahre, 166 Tage)                                                       | 34.884 km/h                                                                         |
| Uranus   | 2.872                                                                               | 50.532                                                                              | $8,68 \times 10^{25}$                                                               | 30664<br>(84 Jahre, 4 Tage)                                                         | 24.516 km/h                                                                         |
| Neptun   | 4.495                                                                               | 49.105                                                                              | $1,02 \times 10^{26}$                                                               | 60148<br>(164 Jahre, 288 Tage)                                                      | 19.548 km/h                                                                         |
| Diagramm |  |  |  |  |  |

### Die gerundeten (und genauen) Verhältnisse zwischen den Umlaufzeiten der Planeten

|         |   |      |            |   |         |
|---------|---|------|------------|---|---------|
| Merkur  | ♿ | 2:5  | (2:5,11)   | ♀ | Venus   |
| Venus   | ♀ | 8:13 | (8:13,004) | ♂ | Erde    |
| Erde    | ♂ | 1:2  | (1:1,88)   | ♂ | Mars    |
| Mars    | ♂ | 1:6  | (1:6,31)   | ♃ | Jupiter |
| Jupiter | ♃ | 2:5  | (2:4,97)   | ♄ | Saturn  |
| Saturn  | ♄ | 1:3  | (1:2,85)   | ♅ | Uranus  |
| Uranus  | ♅ | 1:2  | (1:1,96)   | ♆ | Neptun  |

### Beispiele: Einige Fluchtgeschwindigkeiten

| Himmelskörper                      | $v_2$ am Äquator in km/s |
|------------------------------------|--------------------------|
| Merkur                             | 4,3                      |
| Venus                              | 10,2                     |
| Erde                               | 11,2                     |
| Mond                               | 2,3                      |
| Mars                               | 5,0                      |
| Jupiter                            | 59,6                     |
| Saturn                             | 35,5                     |
| Uranus                             | 21,3                     |
| Neptun                             | 23,3                     |
| Pluto                              | 1,1                      |
| Sonne                              | 617,3                    |
| Sonne im <sup>48</sup> Erdbabstand | 42,0                     |

# Moons of Planets

| Zentralkörper | Nr.   | Name              | a<br>Halbachse<br>/km | e<br>Exz. | i<br>Neig.<br>/° | T<br>Umlauf<br>/d | d<br>/km | m<br>Masse<br>/kg     |
|---------------|-------|-------------------|-----------------------|-----------|------------------|-------------------|----------|-----------------------|
| Erde          | I     | (Erd-)Mond (Luna) | 384.400               | 0,0055    | 5,145            | 27,3217           | 3476     | $7,35 \times 10^{22}$ |
| Mars          | I     | Phobos            | 9.378                 | 0,0151    | 1,08             | 0,3189            | 22       | $1,07 \times 10^{16}$ |
| Mars          | II    | Deimos            | 23.459                | 0,0003    | 0,93             | 1,262             | 12       | $1,8 \times 10^{15}$  |
| Jupiter       | I     | Io                | 421.800               | 0,004     | 0,036            | 1,77              | 3643     | $8,9 \times 10^{22}$  |
| Jupiter       | II    | Europa            | 671.100               | 0,009     | 0,467            | 3,55              | 3122     | $4,8 \times 10^{22}$  |
| Jupiter       | III   | Ganymed           | 1.070.400             | 0,001     | 0,172            | 7,16              | 5262     | $1,5 \times 10^{23}$  |
| Jupiter       | IV    | Kallisto          | 1.882.700             | 0,007     | 0,307            | 16,69             | 4821     | $1,1 \times 10^{23}$  |
| Jupiter       | V     | Amalthea          | 111.000               | 0,001     | 0,004            | 0,42              | 130      | $1,3 \times 10^{18}$  |
| Jupiter       | VI    | Himalia           | 11.100.000            | 0,001     | 0,004            | 16,7              | 11100000 | $1,1 \times 10^{18}$  |
| Jupiter       | VII   | Elara             | 11.100.000            | 0,001     | 0,004            | 16,7              | 11100000 | $1,1 \times 10^{18}$  |
| Jupiter       | VIII  | Pasiphae          | 23.000.000            | 0,001     | 0,004            | 16,7              | 23000000 | $1,1 \times 10^{18}$  |
| Jupiter       | IX    | Sinope            | 23.000.000            | 0,001     | 0,004            | 16,7              | 23000000 | $1,1 \times 10^{18}$  |
| Jupiter       | X     | Lysithea          | 11.100.000            | 0,001     | 0,004            | 16,7              | 11100000 | $1,1 \times 10^{18}$  |
| Jupiter       | XI    | Carme             | 23.000.000            | 0,001     | 0,004            | 16,7              | 23000000 | $1,1 \times 10^{18}$  |
| Jupiter       | XII   | Ananke            | 21.276.000            | 0,244     | 148,9            | 610,5             | 28       | $3,0 \times 10^{16}$  |
| Jupiter       | XIII  | Leda              | 11.165.000            | 0,164     | 27,46            | 240,9             | 18       | $1,1 \times 10^{16}$  |
| Jupiter       | XIV   | Thebe             | 221.900               | 0,018     | 1,070            | 0,68              | 98       | $7,6 \times 10^{17}$  |
| Jupiter       | XV    | Adrastea          | 128.900               | 0,002     | 0,027            | 0,30              | 16       | $1,9 \times 10^{16}$  |
| Jupiter       | XVI   | Metis             | 128.100               | 0,001     | 0,021            | 0,30              | 44       | $9,6 \times 10^{16}$  |
| Jupiter       | XVII  | Callirrhoe        | 24.102.000            | 0,283     | 147,1            | 758,8             | 9        | $8,7 \times 10^{14}$  |
| Jupiter       | XVIII | Themisto          | 7.507.000             | 0,242     | 43,08            | 130,0             | 9        | $6,9 \times 10^{14}$  |
| Jupiter       | XIX   | Menoclite         | 23.806.000            | 0,271     | 152,8            | 752,8             | 6        | $2,1 \times 10^{14}$  |

Project goal: Combine (=join) sets of related tables to gather (all) information about a specific entity type.\*

\*using unary INs



# A Typical Web Table

10 rows

5 columns

| Year                        | Date                                                           | “                                   | Name           | Team           |
|-----------------------------|----------------------------------------------------------------|-------------------------------------|----------------|----------------|
| 2006                        | 1                                                              | Ōizumi                              | I am a colspan |                |
| 2005                        | 2                                                              | F                                   |                | Canada         |
| 2010, (I am destroying 1NF) | 3                                                              | M                                   | Democratic     | Germany        |
| 2008                        | 4                                                              | heterogeneous “null” representation |                | France         |
| 2003                        | I am a real cool row- and colspan, you won't find me in a DBMS |                                     | Re-Elected     | -              |
| 2000                        |                                                                |                                     |                | “ “            |
| 2011                        | 6                                                              |                                     |                | United Kingdom |
| 2002                        | 7                                                              |                                     |                | Japan          |
|                             | 8                                                              | “uncommon” in relational world      |                |                |
| 2009                        |                                                                |                                     |                | Australia      |

no defined foreign key constraints

# A first approach

\*.html (WikiTables project)

```

<table>
<tr>
<th>Planet</th>
<th>Durchmesser</th>
</tr>
<tr>
<td>Merkur</td>
<td>4879 km</td>
</tr>
<tr>
<td>Venus</td>
<td>12103 km</td>
</tr>
...

```



\*.csv

```

Planet,Durchmesser
Merkur,4.879
Venus,12.103
Erde,12.734
Mars,6.772
Jupiter,138.346
Saturn,114.632
Uranus,50.532
Neptun,49.105

```

(N = **1,398,105**)



INDs

```

[table2.column1]⊆
[table13.column3],
[table2.column7]⊆
[table13.column9],
[table1.column2]⊆
[table13.column8],
[table42.column1]⊆
[table13.column1],
[table6.column1]⊆
...

```

Felix Naumann  
Data Profiling  
Summer 2017

worst case:  $\approx (1,398,105 * 5 \text{ columns})^2$  INDS

# The Result – too many open files (SPIDER) or out of memory (BINDER)

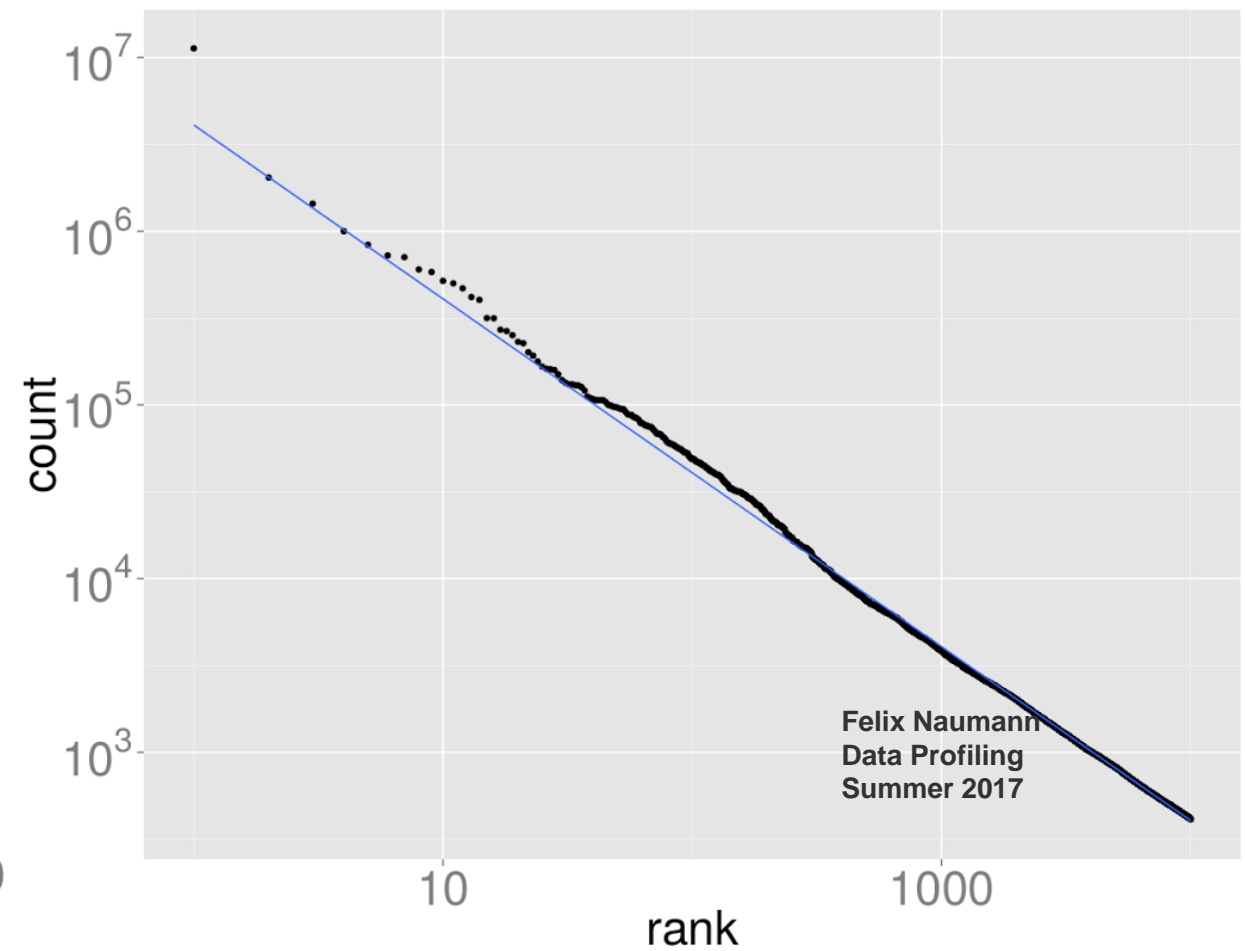
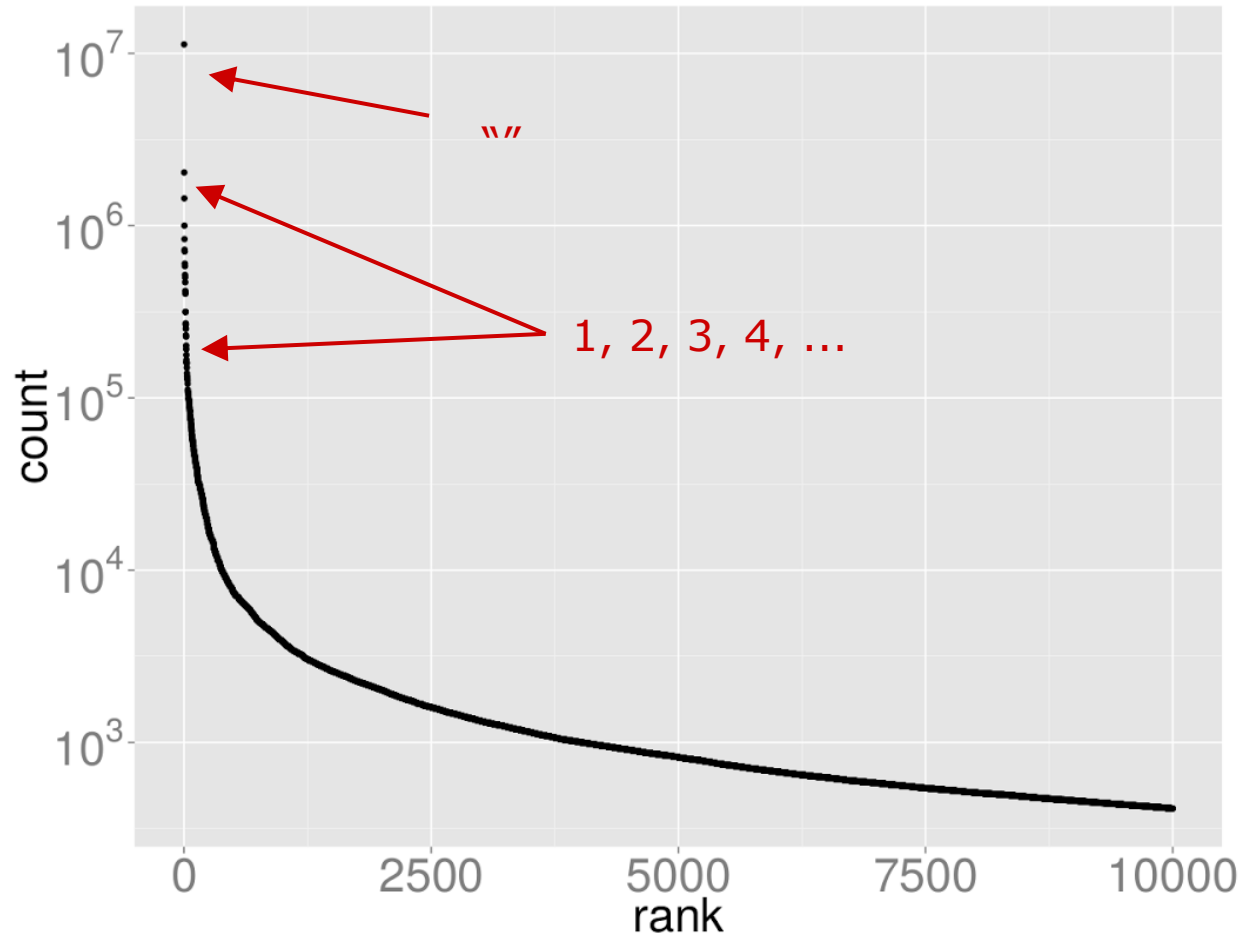
```

Problems @ Javadoc Declaration Console Call Hierarchy
<terminated> Main [Java Application] /usr/lib/jvm/java-7-openjdk-amd64/bin/java (Sep 16, 2014, 2:35:34 PM)
Config [algorithm=SPIDER, databaseURL=null, userName=null, password=null, databaseName=WIKITABLES, databaseType=FILE, inputRowLimit=-1, inputFolderPath=
java.io.IOException: Too many open files
 at java.io.UnixFileSystem.createFileExclusively(Native Method)
 at java.io.File.createNewFile(File.java:1006)
 at de.uni_potsdam.hpi.utils.FileUtils.createFile(FileUtils.java:73)
 at de.uni_potsdam.hpi.utils.FileUtils.buildFileWriter(FileUtils.java:93)
 at de.uni_potsdam.hpi.metanome.algorithms.spider.TPMS.write(TPMS.java:92)
 at de.uni_potsdam.hpi.metanome.algorithms.spider.TPMS.sortToDisk(TPMS.java:60)
 at de.uni_potsdam.hpi.metanome.algorithms.spider.Attribute.<init>(Attribute.java:103)
 at de.uni_potsdam.hpi.metanome.algorithms.spider.SPIDER.initializeAttributes(SPIDER.java:313)
 at de.uni_potsdam.hpi.metanome.algorithms.spider.SPIDER.execute(SPIDER.java:204)
 at de.uni_potsdam.hpi.metanome_test_runner.mocks.MetanomeMock.executeSPIDER(MetanomeMock.java:262)
 at de.uni_potsdam.hpi.metanome_test_runner.MetanomeTestRunner.run(MetanomeTestRunner.java:84)
 at de.uni_potsdam.hpi.metanome_test_runner.MetanomeTestRunner.runWikitable(MetanomeTestRunner.java:186)
 at de.uni_potsdam.hpi.metanome_test_runner.Main.main(Main.java:7)
de.uni_potsdam.hpi.metanome.algorithm_integration.AlgorithmExecutionException: Too many open files
 at de.uni_potsdam.hpi.metanome.algorithm_integration.AlgorithmExecutionException.<init>(AlgorithmExecutionException.java:101)
 at de.uni_potsdam.hpi.metanome.algorithms.spider.SPIDER.execute(SPIDER.java:204)
 at de.uni_potsdam.hpi.metanome_test_runner.mocks.MetanomeMock.executeSPIDER(MetanomeMock.java:262)
 at de.uni_potsdam.hpi.metanome_test_runner.MetanomeTestRunner.run(MetanomeTestRunner.java:84)
 at de.uni_potsdam.hpi.metanome_test_runner.MetanomeTestRunner.runWikitable(MetanomeTestRunner.java:186)
 at de.uni_potsdam.hpi.metanome_test_runner.Main.main(Main.java:7)
<terminated> Main [Java Application] /usr/lib/jvm/java-7-openjdk-amd64/bin/java (Sep 16, 2014, 2:40:43 PM)
Config [algorithm=BINDER, databaseURL=null, userName=null, password=null, databaseName=WIKITABLES, databaseType=FILE, inputRowLimit=-1, inputFolderPath=
Exception in thread "main" java.lang.OutOfMemoryError: GC overhead limit exceeded
 at java.lang.Integer.valueOf(Integer.java:642)
 at it.unimi.dsi.fastutil.ints.AbstractIntIterator.next(AbstractIntIterator.java:56)
 at it.unimi.dsi.fastutil.ints.AbstractIntIterator.next(AbstractIntIterator.java:51)
 at de.uni_potsdam.hpi.metanome.algorithms.binder.IntSingleLinkedList.<init>(IntSingleLinkedList.java:71)
 at de.uni_potsdam.hpi.metanome.algorithms.binder.BINDER.fetchCandidates(BINDER.java:888)
 at de.uni_potsdam.hpi.metanome.algorithms.binder.BINDER.checkViaTwoStageIndexAndLists(BINDER.java:805)
 at de.uni_potsdam.hpi.metanome.algorithms.binder.BINDER.execute(BINDER.java:279)
 at de.uni_potsdam.hpi.metanome_test_runner.mocks.MetanomeMock.executeBINDER(MetanomeMock.java:59)
 at de.uni_potsdam.hpi.metanome_test_runner.MetanomeTestRunner.run(MetanomeTestRunner.java:87)
 at de.uni_potsdam.hpi.metanome_test_runner.MetanomeTestRunner.runWikitable(MetanomeTestRunner.java:186)
 at de.uni_potsdam.hpi.metanome_test_runner.Main.main(Main.java:7)
(Spider) Runtime

```

Felix Naumann  
Data Profiling  
Summer 2017

# Value distribution

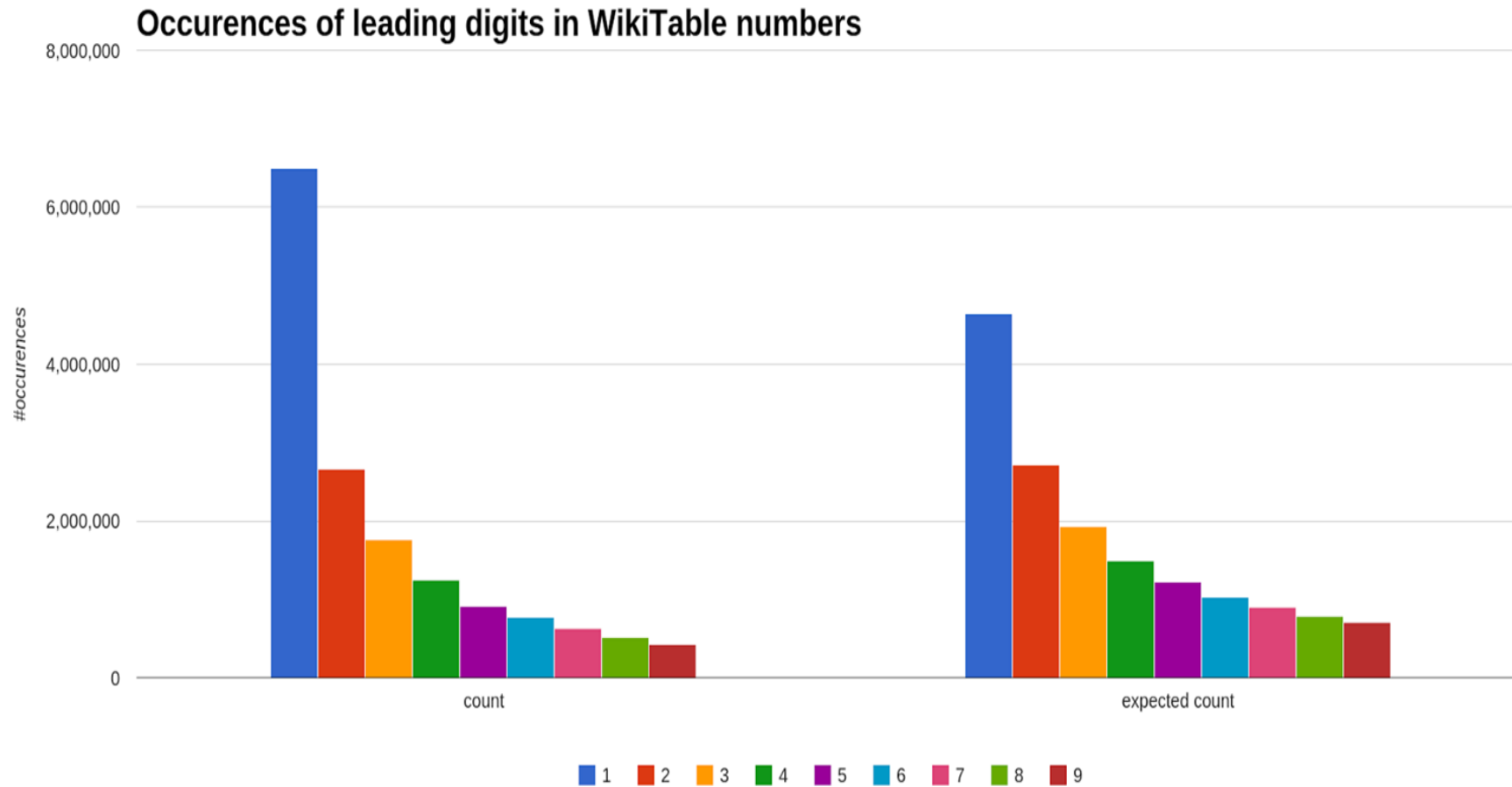


# Null values

| null repr. | <i>tf</i> rank | occurences | null repr. | <i>tf</i> rank | occurences |
|------------|----------------|------------|------------|----------------|------------|
| “__”       | 6th            | 724,917    | “n/a”      | 182nd          | 26,596     |
| “_”        | 11th           | 501,114    | “•”        | 198st          | 23,574     |
| “_”        | 19th           | 252,008    | “- - -”    | 578th          | 6,619      |
| “N/A”      | 33th           | 131,339    | “.”        | 604th          | 6,367      |
| “?”        | 49th           | 97,236     | “??”       | 645th          | 6,053      |
| “Unknown”  | 146th          | 32,484     | “(n/a)”    | 915th          | 4,228      |
| “- -”      | 176th          | 27,898     |            |                |            |

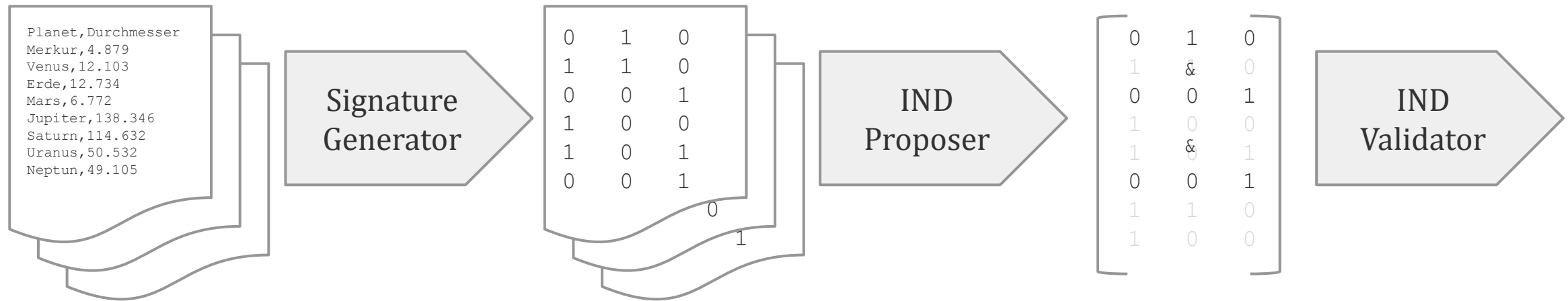
Felix Naumann  
Data Profiling  
Summer 2017

# Newcomb-Benford's law



Felix Naumann  
Data Profiling  
Summer 2017

# The MANY algorithm



## Algorithm 1 MANY

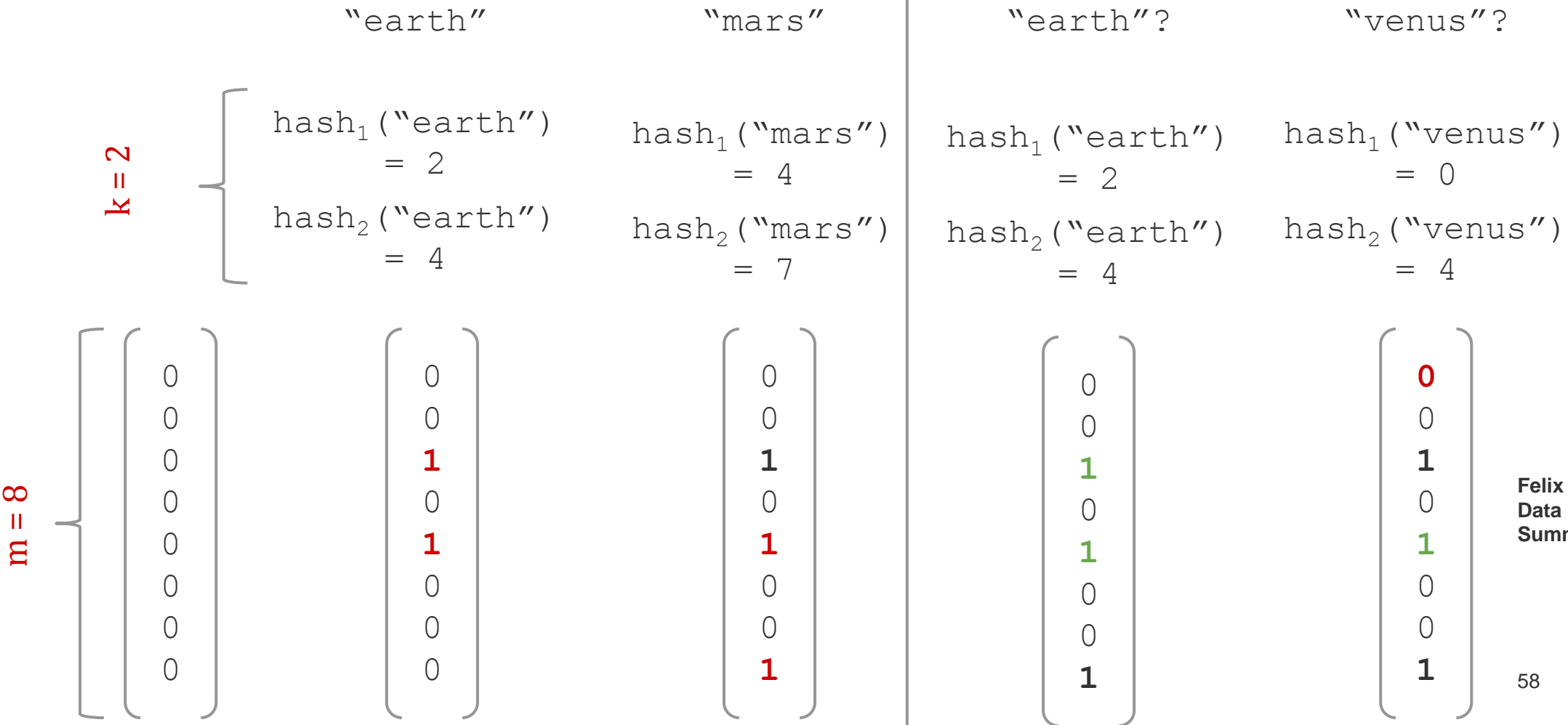
**Require:** Database schema  $D = \{R_1, R_2, \dots, R_s\}$ , instance of database schema  $d = \mathcal{I}(D) = \{r_1, r_2, \dots, r_s\}$ , Bloom filter size  $m \in \mathbb{N}$ , hash-function count  $k \in \{1, \dots, m\}$ , *strategy*  $\in \{\text{dep2refs}, \text{ref2deps}\}$

**Ensure:** All valid INDs *uinds*

- 1: *colIndicesMap*  $\leftarrow$  BUILD COLINDICESMAP( $D$ )
- 2: **function** MANY( $D, d$ )
- 3:     *uinds*  $\leftarrow \emptyset$
- 4:     *sigMatrix*  $\leftarrow$  GENERATESIGNATURES( $D, d$ )
- 5:     *candidates*  $\leftarrow$  GENERATECANDIDATES(*sigMatrix*)
- 6:     **for all** *candidate*  $\in$  *candidates* **do**
- 7:         **if** ISVALID(*candidate*,  $d$ ) **then**
- 8:             *uinds*  $\leftarrow$  *uinds*  $\cup$  {*candidate*}
- 9:     **return** *uinds*
- 10: **return** MANY( $D, d$ )



# Reminder: Bloom filters to check presence of value



# Signature generator

table1

| Planet | Symbol | Earth-like |
|--------|--------|------------|
| Earth  | ♁      | Y          |
| Mars   | ♂      | Y          |
| Venus  | ♀      | Y          |

Bloom filter  
m = 8, k = 1

|   |
|---|
| 1 |
| 1 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |

|   |
|---|
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |

|   |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |

table2

| Orb     | Atmosphere |
|---------|------------|
| Mercury | Y          |
| Venus   | Y          |
| Earth   | Y          |
| Mars    | Y          |

|   |
|---|
| 1 |
| 1 |
| 0 |
| 1 |
| 1 |
| 0 |
| 0 |
| 0 |

|   |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |

table1.Planet  $\subseteq$  table2.Orb

$\Rightarrow$

|   |
|---|
| 1 |
| 1 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |

$\cap$

|   |
|---|
| 1 |
| 1 |
| 0 |
| 1 |
| 1 |
| 0 |
| 0 |
| 0 |

table1.Symbol  $\not\subseteq$  table2.Orb

$\Leftarrow$

|   |
|---|
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 1 |
| 1 |

$\neq$

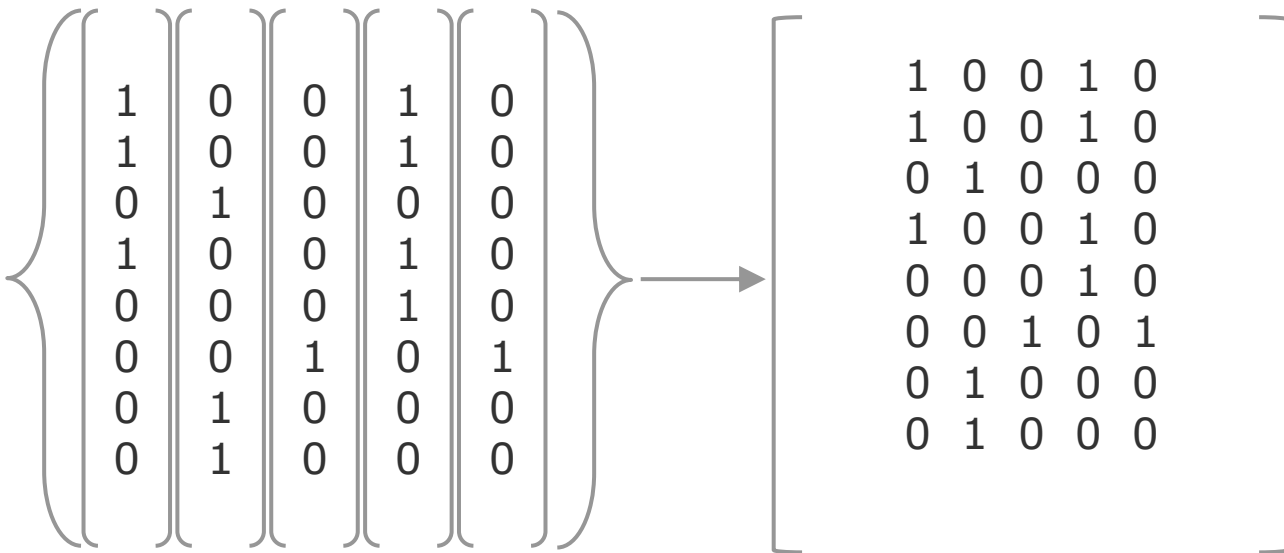
|   |
|---|
| 1 |
| 1 |
| 0 |
| 1 |
| 1 |
| 0 |
| 0 |
| 0 |

# IND Proposer

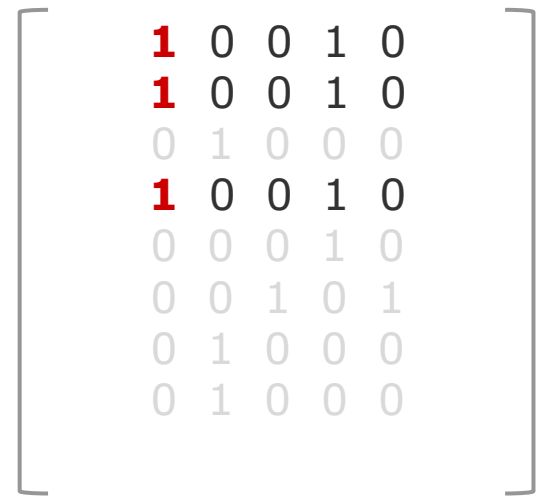
table1

table2

|       | Planet | Symbol | Earth-like | Orb | Atmosphere |
|-------|--------|--------|------------|-----|------------|
| index | 0      | 1      | 2          | 3   | 4          |



$C_0 \subseteq ?$



**1** 0 1 1 0 &&  
**1** 0 0 1 0 &&  
**1** 0 0 1 0 =  
**1** **0** **0** **1** **0**

$C_0 \sqcup C_0'$   
 $C_0 \sqcup C_3$

# IND Validator

Bloom filter are probabilistic DS:

$$\text{bf}(\{\text{"Foo"}, \text{"ABC"}, \text{"123"}\}) = (1, 0, 1, 1)$$

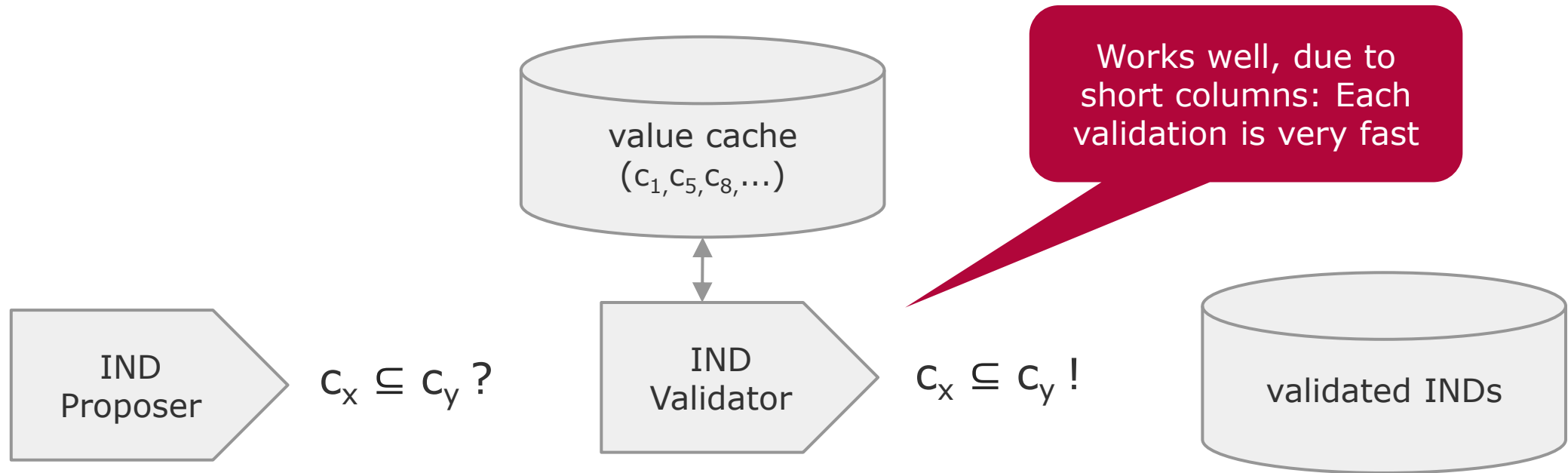
$$\text{bf}(\{\text{"Bar"}, \text{"XYZ"}\}) = (1, 0, 0, 1)$$

Problem:

$$\{\text{"Bar"}, \text{"XYZ"}\} \not\subseteq \{\text{"Foo"}, \text{"ABC"}, \text{"123"}\}$$

$$\text{recall} = 1$$

$$\text{precision} \leq 1$$



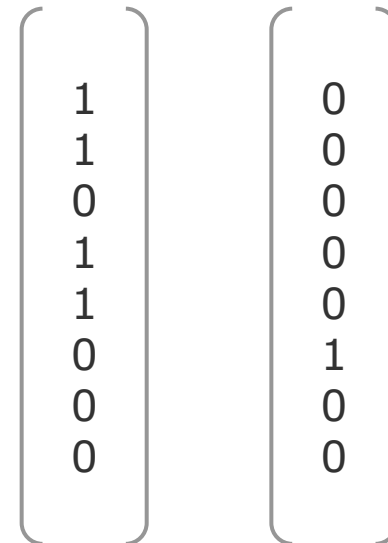
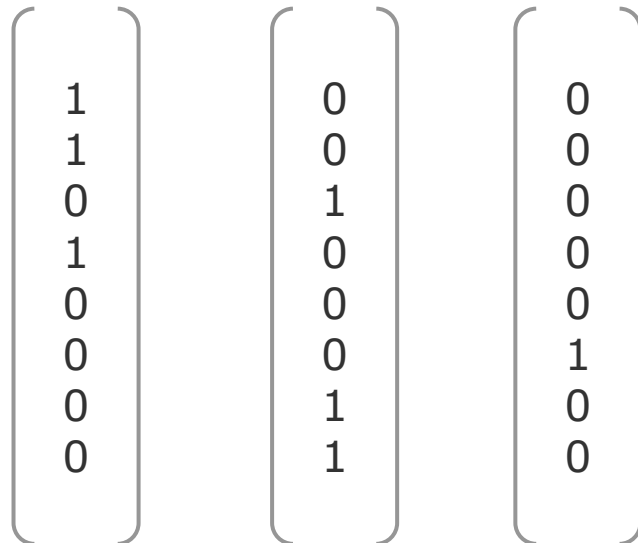
# Insight: Check superset relationship, too

table1

| Planet | Symbol | Earth-like |
|--------|--------|------------|
| Earth  | ♁      | Y          |
| Mars   | ♂      | Y          |
| Venus  | ♀      | Y          |

table2

| Orb     | Atmosphere |
|---------|------------|
| Mercury | Y          |
| Venus   | Y          |
| Earth   | Y          |
| Mars    | Y          |

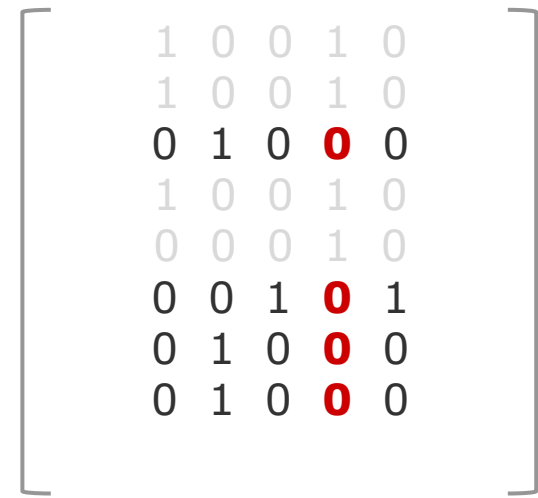
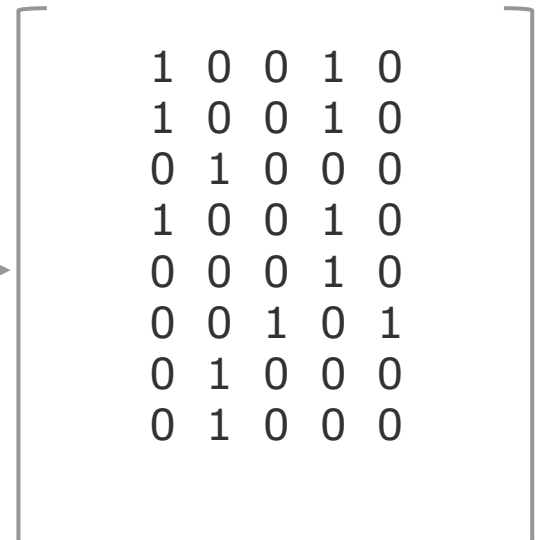
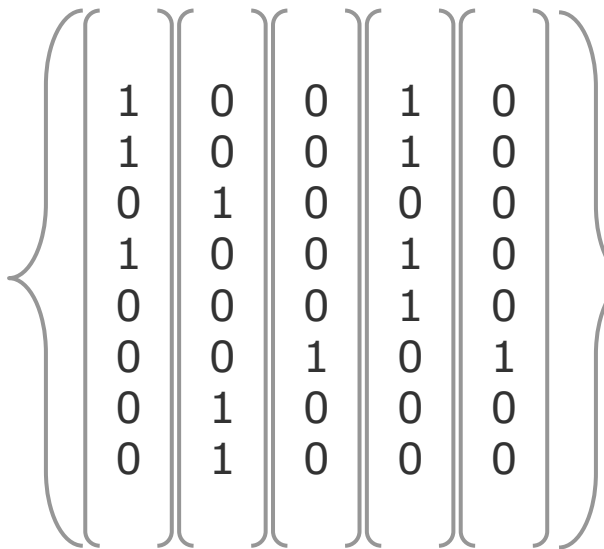


# Checking superset relationships to generate candidates

table1

table2

|       | Planet | Symbol | Earth-like | Orb | Atmosphere |
|-------|--------|--------|------------|-----|------------|
| index | 0      | 1      | 2          | 3   | 4          |



1 0 1 **1** 1 &&  
 1 1 0 **1** 0 &&  
 1 0 1 **1** 1 &&  
 1 0 1 **1** 1 =  
**1** 0 0 **1** 0

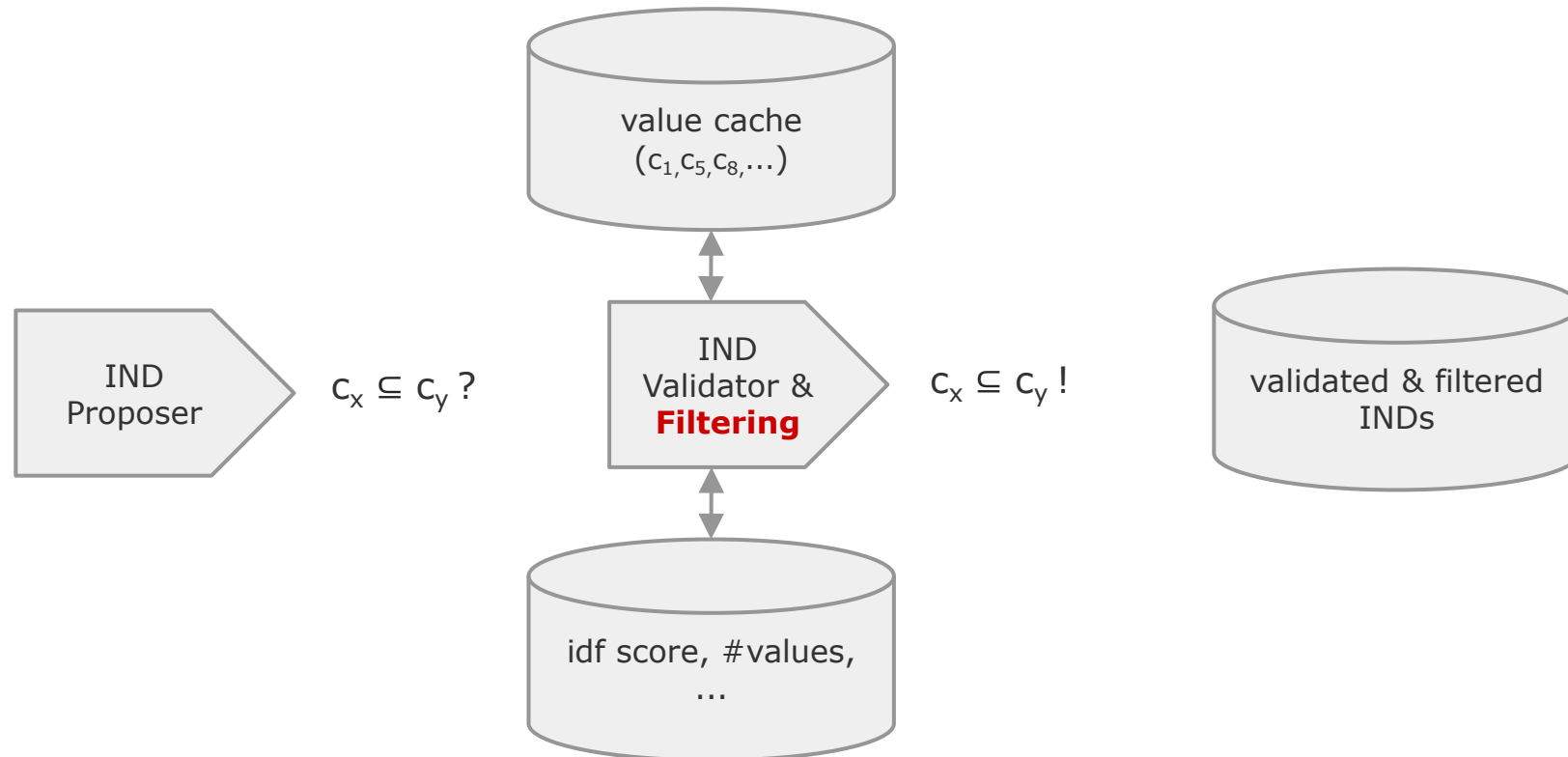
$C_3 \supseteq C_0,$   
 $C_3 \supseteq C_3$

Felix Naumann  
Data Profiling  
Summer 2017

Incremental detection: For each new column, easily generate all candidates

# Problem: Too many INDs

| #tables | 1,000   | 2,500   | 5,000     | 10,000     | 25,000     | 50,000      |
|---------|---------|---------|-----------|------------|------------|-------------|
| #INDs   | 122,755 | 774,375 | 2,586,716 | 10,214,320 | 63,917,448 | 259,429,353 |



Felix Naumann  
Data Profiling  
Summer 2017

## Filters

---

- Filter INDs with non-unique referenced column
  
- Filter all transitive INDs  $\{A \subseteq B, B \subseteq C, A \subseteq C\}$ 
  - Goal: Reference largest possible column
  - Lose some information
  
- Filter “uninteresting” (e.g., numeric) columns
  
- Filter empty columns
  
- Filter INDs with low coverage on the referenced values

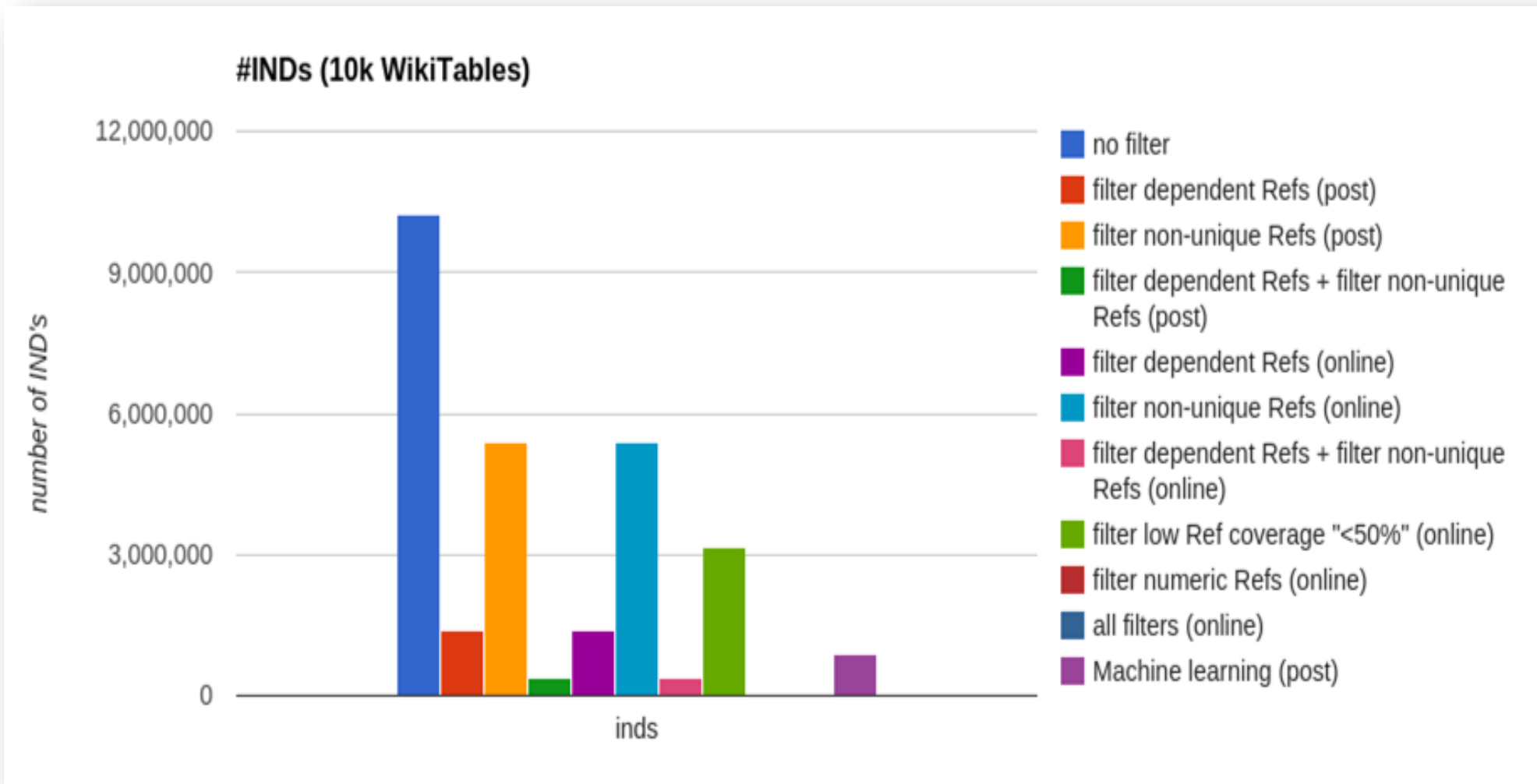


## “Online” filter implementation: Use both directions of candidate generation

---

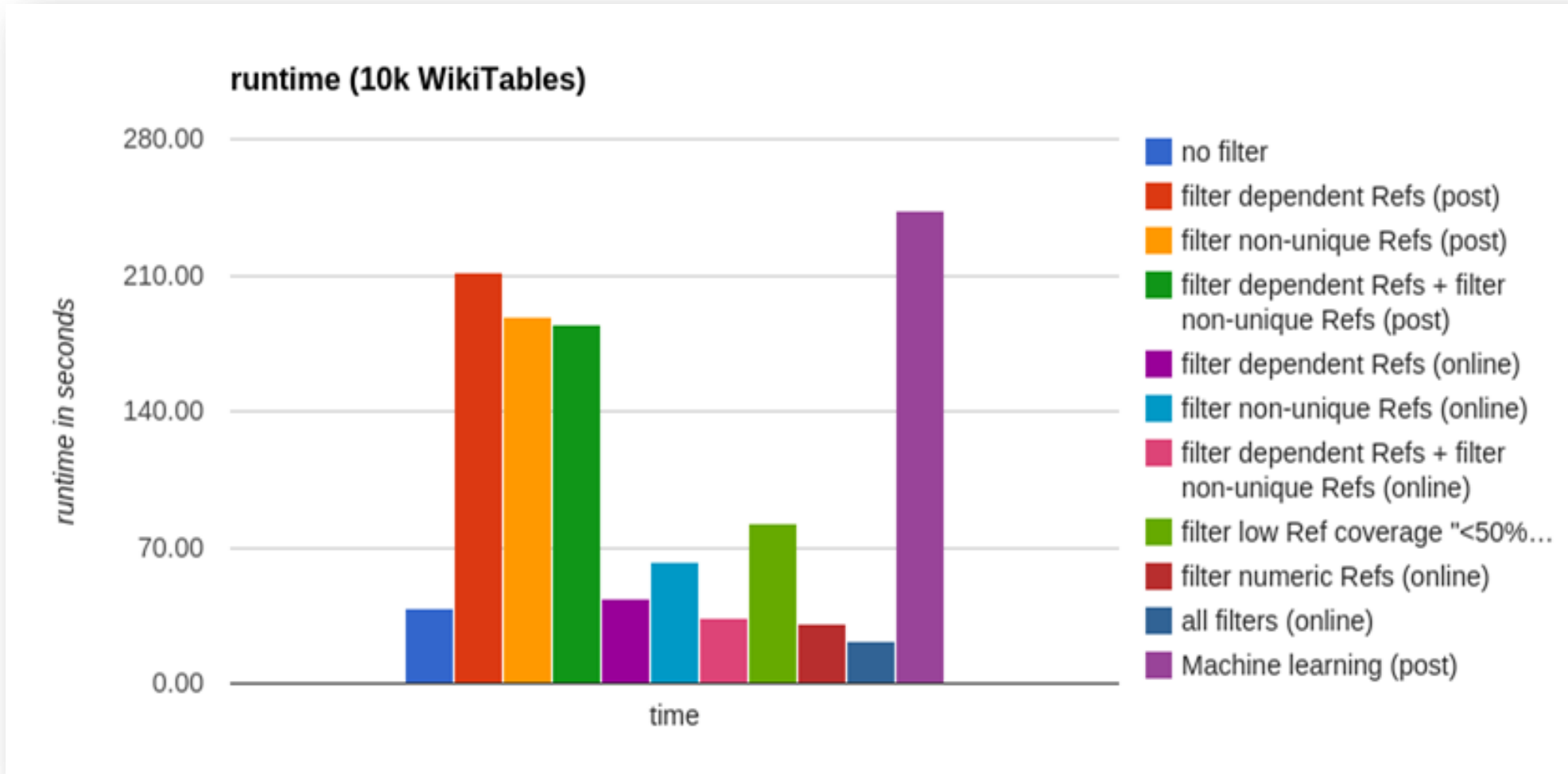
- Non-unique columns
  - Remove as referenced candidate
  
- Numeric and null-columns
  - Remove as referenced and as dependent candidate
  
- Columns referencing other columns
  - Remove as referenced candidate
  
- Low coverage filter
  - Expensive -> Postprocessing to reduce result-size

# Filter Evaluation



Felix Naumann  
Data Profiling  
Summer 2017

# Filter Evaluation



Felix Naumann  
Data Profiling  
Summer 2017

# Visualisation

```
[1011066.Name] =] [1011057.Name]
[129284.Reference] =] [1223862.null] [586920.Ref.] [1030730.RCDB page] [108435.No.] [1248790.Source]
 [983315.References] [207338.Home railway (external link)] [975850.Ref] [1375996.Source] [1129539.References]
 [1168707.References] [744488.Ref] [1169311.Ref] [1068498.Ref] [163214.Reference] [604676.References] [1002900.Ref]
 [749972.Reference] [951640.References] [939700.Page] [900853.Ref] [788203.Ref] [788409.References] [978758.Ref]
 [652885.Link] [652377.Ref] [1320358.Reference] [1287392.Ref] [1012269.Report] [1180077.References] [1274408.Ref]
 [856227.NFL Recap] [1286480.Ref] [1354142.null] [525501.References] [630016.Notes] [762537.Refs] [902406.Report]
 [1005369.Link] [1255682.Source] [1157534.Source] [1065320.Ref] [956840.Ref] [775466.References] [988811.Ref]
 [1005838.Link] [1005593.Link] [576411.References] [1134428.Ref] [1170953.Reference(s)] [699144.Note]
 [268733.References] [931606.Notes] [1284557.Ref.] [1357973.Source] [1238931.Report] [867400.Reference]
 [794774.Ref] [716064.Refs] [377521.References] [995370.Ref] [1282132.References] [1358158.Ref.] [1120007.Ref]
 [1342522.Ref] [1319381.null] [889114.Ref] [1004839.Link] [697527.Website] [980509.Ref(s)] [1078901.Ref]
[1390416.Rank] =] [1169921.Rank] [1183098.Rank] [1011765.Rank] [1225076.Rank] [454782.Rank] [1186535.Rank]
 [1209635.Rank] [1161665.Rank] [708465.Rank] [708648.Rank]
[637307.Date] =] [1311505.Date] [1337020.Date]
[1083420.Event] =] [976659.Event] [976901.Event] [975917.Event] [1060037.Event] [1068182.Event] [1067251.Event]
 [1067097.Event] [1000067.Event] [972968.Event] [1058267.Event] [988323.Event] [1003312.Event] [1063506.Event]
 [1027145.Event] [1078507.Event] [1062268.Event]
[302006.Role:] =] [391330.Role:] [703281.Role:] [387497.Role:] [735612.Role:] [151885.Role:] [150598.Role:]
[1083410.Event] =] [983546.Event] [975773.Event] [1071989.Event] [1068219.Event] [1002900.Event] [1074984.Event]
 [967160.Event] [1052352.Event] [1066949.Event] [1082562.Event] [1151162.Event] [1042660.Event] [1056643.Event]
 [950860.Event] [958921.Event] [1063309.Event] [973967.Event] [1027145.Event] [1062263.Event]
[73362.State] =] [1185141.State]
[1083402.Event] =] [1083339.Event] [1068498.Event] [1060027.Event] [1002823.Event] [1046135.Event] [1249836.Event]
 [1000145.Event] [994576.Event] [990543.Event]
[854590.Venue] =] [883202.Venue] [890993.Venue] [1104659.Venue]
[648260.TEAM] =] [1286540.Club] [1308745.Club]
[627822.Division Record] =] [466958.Sets W - L]
[1236345.Match] =] [1231569.Match]
```

...

# Visualisation

```

INDs = {
 R1.A ⊆ R2.B,
 R3.A ⊆ R1.D,
 R3.C ⊆ R2.A,
 R3.B ⊆ R4.A
}

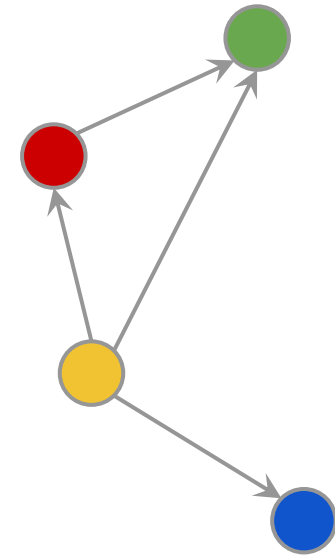
```



```

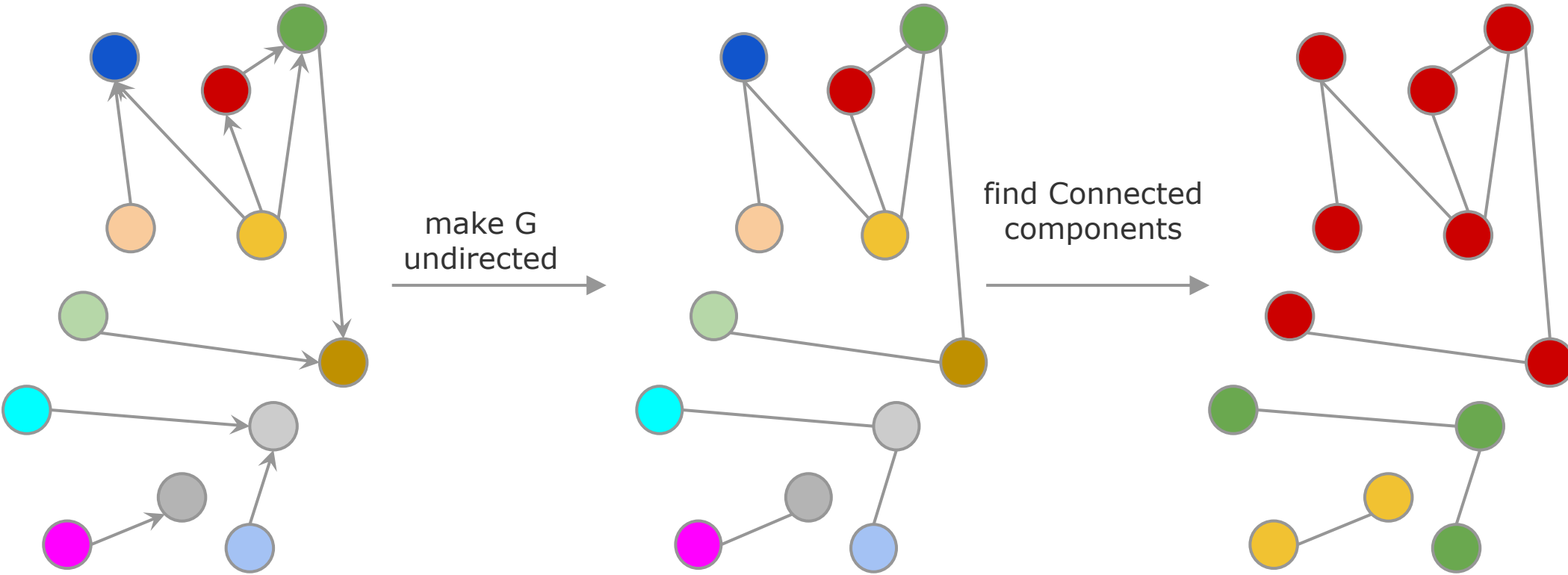
G = (
 V = {
 R1, R2, R3, R4
 },
 E = { (R1, R2), (R3, R1),
 (R3, R2), (R3,
 R4)
 }
)

```



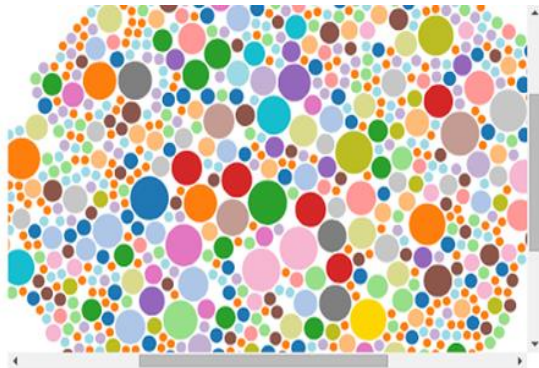
Felix Naumann  
Data Profiling  
Summer 2017

# Visualisation



Felix Naumann  
Data Profiling  
Summer 2017

# Application



|          |                                                                                           |
|----------|-------------------------------------------------------------------------------------------|
| 96242-1  | ...                                                                                       |
| 43666-3  | 43666-3.'BBC_Radio_Stoke'. 'Programming'.csv                                              |
| 53064-1  | 53064-1.'Rotation_period';'Rotation period of selected objects'.csv                       |
| 562884-4 | 562884-4.'Planets_in_astrology';'Ruling planets of the astrological signs and houses'.csv |
| 175797-1 | 175797-1.'Sun_sign_astrology';'Sun signs'.csv                                             |
| 177750-2 | 177750-2.'BBC_Radio_Manchester'. 'Programming'.csv                                        |
| 89462-4  | 89462-4.'Astrology_and_the_classical_elements';'Triplcities by season'.csv                |
| 213213-1 | 213213-1.'Dalton_Park';'Opening times'.csv                                                |
| 470402-  | 470402-                                                                                   |

| Celestial Objects | Rotation period                                                                          | Rotation period                                             |
|-------------------|------------------------------------------------------------------------------------------|-------------------------------------------------------------|
| Sun               | 25.379995 days (equatorial) 35 days (high latitude)                                      | 25 d 9 h 7 m 11.6 s 35 d                                    |
| Mercury           | 58.6462 days                                                                             | 58 d 15 h 30 m 30 s                                         |
| Venus             | ?243.0187 days                                                                           | ?243 d 0 h 26 m                                             |
| Earth             | 0.99726968 days                                                                          | 0 d 23 h 56 m 4.100 s                                       |
| Moon              | 27.321661 days ( synchronous toward Earth)                                               | 27 d 7 h 43 m 11.5 s                                        |
| Mars              | 1.02595675 days                                                                          | 1 d 0 h 37 m 22.663 s                                       |
| Ceres             | 0.37809 days                                                                             | 0 d 9 h 4 m 27.0 s                                          |
| Jupiter           | 0.4135344 days (deep interior) 0.41007 days (equatorial) 0.41369942 days (high latitude) | 0 d 9 h 55 m 29.37 s 0 d 9 h 50 m 30 s 0 d 9 h 55 m 43.63 s |
| Saturn            | 0.44403 days (deep interior) 0.426 days (equatorial) 0.443 days (high latitude)          | 0 d 10 h 39 m 24 s 0 d 10 h 14 m 0 d 10 h 38 m              |

Zoom (1-5)

Range (logarithmic)

Dataset

allFilters

Felix Naumann  
Data Profiling  
Summer 2017

## Summary

---

1. Inclusion Dependencies
2. SQL
3. De Marchi et al.
4. MIND
5. SPIDER
6. MANY



Felix Naumann  
Data Profiling  
Summer 2017