



Data Profiling

Télécom ParisTech, Paris  
Felix Naumann

# The Hasso Plattner Institute



Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

# Information Systems Team



Thorsten Papenbrock



Diana Stephan



Prof. Felix Naumann



Dr. Ralf Krestel



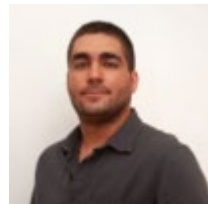
Tim Repke



Julian Risch



Tobias Bleifuß



John Koumarelas



Michael Loster



Leon Bornemann



Konstantina Lazaridou



Lan Jiang



Hazar Harmouch

Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

Data Change    Data Fusion    Duplicate Detection  
 project **DuDe**  
 Data Profiling    Information Integration    Web Science  
 project **DataChEx**    Data Scrubbing    Data as a Service  
 Entity Search  
 Information Quality    Data Cleansing    Text Mining  
 Web Data    Linked Open Data    RDF Data Mining  
 Dependency Detection    ETL Management  
 Service-Oriented Systems    Entity Recognition    Opinion Mining    Data Preparation  
 project **Metanome**    Change Exploration

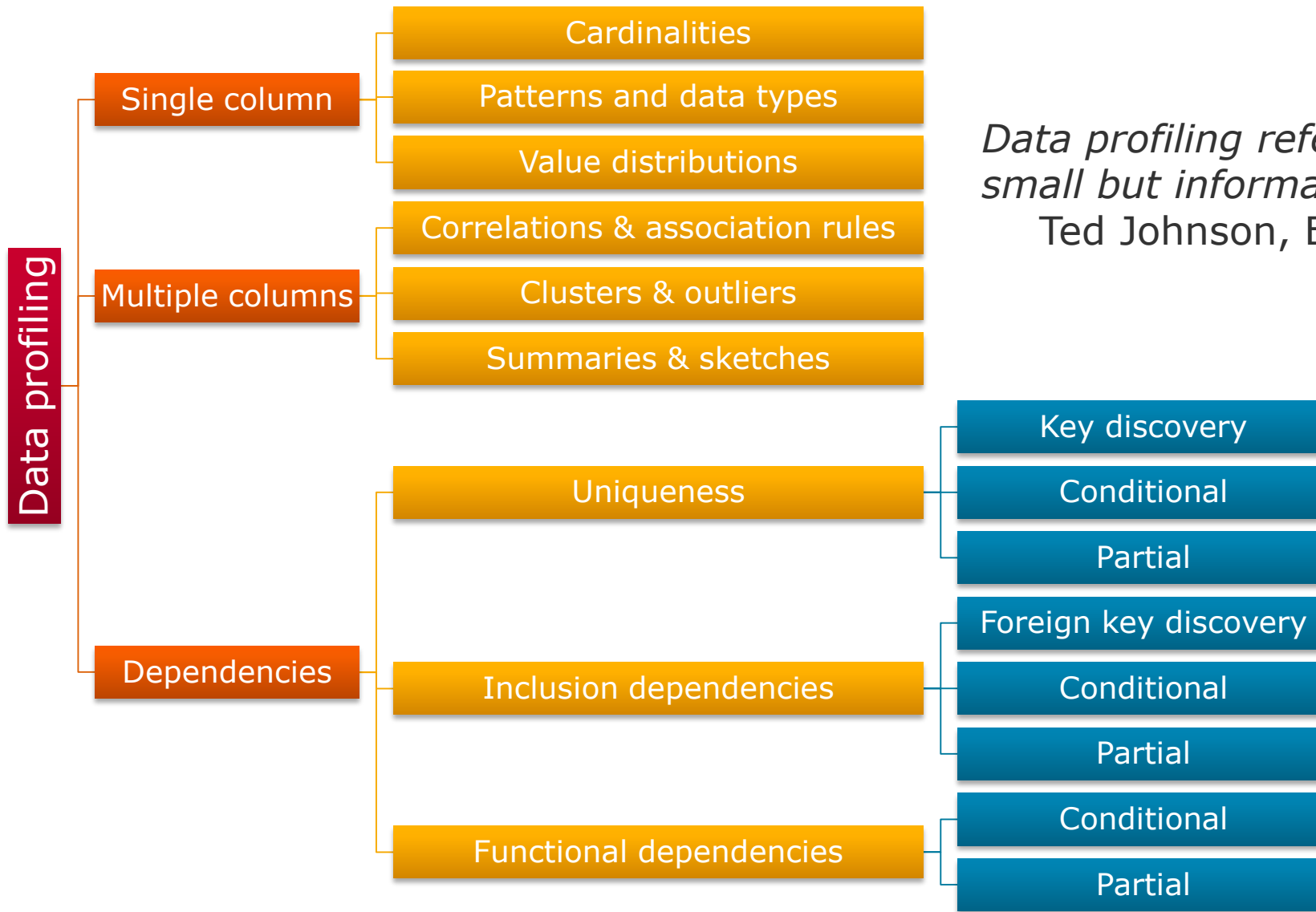
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X
	county_id	county_desc	voter_reg_no	status_cd	voter_status_desc	reason_cd	voter_status	last_name	first_name	midl_name	name	res_street	res_city_desc	state	zip_code	mail_addr1	mail_addr2	mail_city	mail_state	mail_zipcode	full_phone	race_code	ethnic_code	party_cd
2	1	ALAMANCE	9005990	A	ACTIVE	AV	VERIFIED	AABEL	EVELYN	LARSEN	4430 E GREENSBORO	GRAHAM	NC	27253	4430 E GREENSBORO-CHA		GRAHAM	NC	27253	000 0000	W	NL	UNA	
3	1	ALAMANCE	9048723	A	ACTIVE	AV	VERIFIED	AARON	CHRISTINA	CASTAGNA	421 WHITT AVE	BURLINGTON	NC	27215	PO BOX 4177		BURLINGTON	NC	27215	229 1110	W	UN	UNA	
4	1	ALAMANCE	9019674	A	ACTIVE	AV	VERIFIED	AARON	CLAUDIA	HAYDEN	1013 EDITH ST	BURLINGTON	NC	27215	1013 EDITH ST		BURLINGTON	NC	27215	222 8834	W	NL	UNA	
5	1	ALAMANCE	9129589	A	ACTIVE	AV	VERIFIED	AARON	JAMES	MICHAEL	1647 SAXAPAHAW	GRAHAM	NC	27253	PO BOX 98		SAXAPAHAW	NC	27340	336 525 2484	W	UN	DEM	
6	1	ALAMANCE	9041748	A	ACTIVE	AV	VERIFIED	AARON	NATHAN	EDWARD	421 WHITT AVE	BURLINGTON	NC	27215	PO BOX 4177		BURLINGTON	NC	27215	336 229 1110	W	UN	UNA	
7	1	ALAMANCE	9021947	A	ACTIVE	AV	VERIFIED	AARON	WILLIE	DALE	1013 EDITH ST	BURLINGTON	NC	27215	1013 EDITH ST		BURLINGTON	NC	27215	336 999 9999	W	NL	UNA	
8	1	ALAMANCE	9062002	A	ACTIVE	AV	VERIFIED	AARONSON	GENA	HOLT	107 TERRYWOOD	HAW RIVER	NC	27258	107 TERRYWOOD CT		HAW RIVER	NC	27258	336 578 9123	W	NL	REP	
9	1	ALAMANCE	9096423	A	ACTIVE	AV	VERIFIED	AARONSON	MICHAEL	CHARLES	107 TERRYWOOD	HAW RIVER	NC	27258	107 TERRYWOOD CT		HAW RIVER	NC	27258	336 266 7615	W	NL	UNA	
10	1	ALAMANCE	9117940	I	INACTIVE	IU	CONFIRMATI	ABAD	PRISCILLA	MARIE	100 COLONNADE	ELON	NC	27244	CAMPUS BOX 3008		ELON	NC	27244		O	HL	UNA	
11	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	COLLEEN	MIASHEL	1097 IVEY RD	#C GRAHAM	NC	27253	1097 IVEY RD	#C	GRAHAM	NC	27253		M	HL	REP	
12	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	COLLEEN	MIASHEL	1097 IVEY RD	#C GRAHAM	NC	27253	1097 IVEY RD	#C	GRAHAM	NC	27253	336 212 8140	W	NL	UNA	
13	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
14	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
15	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
16	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
17	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
18	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
19	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
20	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
21	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
22	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
23	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
24	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
25	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
26	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
27	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
28	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
29	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
30	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
31	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
32	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
33	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
34	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
35	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
36	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
37	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
38	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
39	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
40	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
41	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
42	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
43	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
44	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
45	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
46	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
47	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
48	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	
49	1	ALAMANCE	9034111	I	INACTIVE	IU	CONFIRMATI	ABADIE	MYRA	HOLLIFIELD	612 SIDEVIEW ST	GRAHAM	NC	27253	617 MITCHELL ST		BURLINGTON	NC	27217	336 212 8140	W	NL	UNA	

Column labels





# Data Profiling: Classification of Tasks



*Data profiling refers to the activity of creating small but informative summaries of a database.*  
Ted Johnson, Encyclopedia of Database Systems

## Use Cases for Data Profiling

---

- **Query optimization:** Counts and histograms, functional dependencies, ...
- **Data cleansing:** Patterns, rules, and violations
- **Data integration:** Cross-DB inclusion dependencies
- **Scientific data management:** Inspect new datasets
- **Data analytics and mining:** Profiling as preparation to decide on models and questions
- **Database reverse engineering**

In summary: **Data preparation**

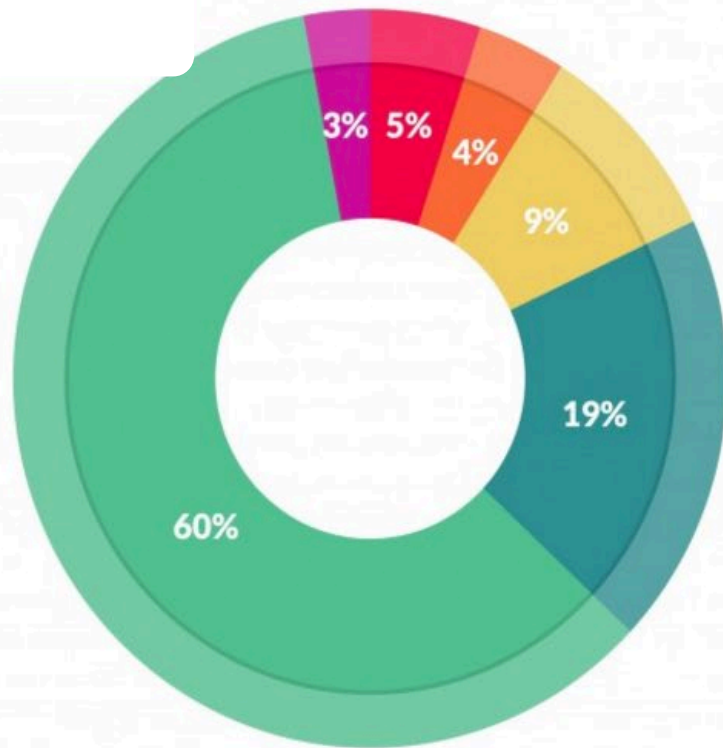
- “If we just have a bunch of data sets in a repository, it is unlikely anyone will ever be able to find, let alone reuse, any of this data. With adequate metadata, there is some hope, but even so, challenges will remain...”

Felix Naumann  
Data Profiling  
Télécom ParisTech 2018



## Data Profiling as Data Preparation

***Data preparation*** accounts for about 80% of the work of data scientists



What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

## Shortcomings of commercial and research tools

- Usability
  - Complex to configure
  - Results complex to view and interpret
- Scalability: SQL works poorly
- Efficiency
  - Coffee, Lunch, Overnight
- Functionality
  - Restricted to simplest tasks
  - Restricted to single columns or small column sets
  - „Checking“ vs. „discovery“
- Interpretation of profiling results

IBM Information Server File e.g., IBM Information Analyzer 9.43.86.77

IA\_OVERVIEW\_PROJECT

INVESTIGATE Foreign Key Analysis

Select Data Source to Work With

EMPLOYEE DEPARTMENT

Open Foreign Key Analysis View Details

You can use this pane to view analysis details about a primary key column and the foreign key column that is associated with the primary key column.

Frequency Values Analysis Details

Foreign Key Candidate Pair		
	Base Column	Paired Column
Column	EMPNO	MGRNO
Table	EMPLOYEE	DEPARTMENT
Source	IA	IA
Primary Key	Yes	No
Foreign Key	No	Yes
Data Class	Identifier	Quantity
Data Type	INT32	INT8
Length	0	0
Precision	0	0
Scale	0	0
Cardinality	48	9
Unique	No	No
Constant	No	No
Definition	No	No

Paired to Base:  
Common Data Values: 8 100.0000% Common Domain: Yes

Base to Paired:  
Common Data Values: 8 16.6667% Common Domain: No

Common Domain:

40 8 1

Base Column Paired Column

## Scalable profiling

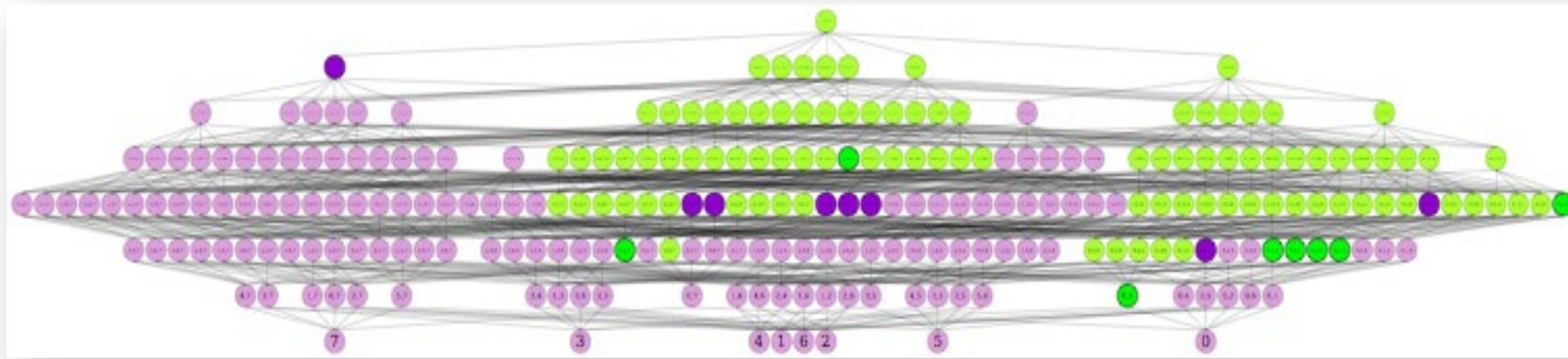
---

- Scalability in number of rows
- Scalability in number of columns
  - “Normal” table with 100 columns:  
 $2^{100} - 1 = 1,267,650,600,228,229,401,496,703,205,375$   
= 1.3 nonillion column combinations
  - Impossible to check or even enumerate
- Possible solutions
  - Scale up: More memory, faster CPUs
  - Scale in: More cores
  - Scale out: More machines
  - Scale smart: Intelligent enumeration and aggressive pruning



Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

# Large solution space



- Size of lattice:  $2^n - 1$  (empty set not considered)
- Nodes at level 1:  $n$
- Nodes at level  $n$ : 1
- Nodes at level  $k$ :  $\binom{n}{k} = \frac{n!}{(n-k)!k!}$
- Largest level at  $n/2$ :  $\binom{n}{n/2} = \frac{n!}{\left(\frac{n!}{2}\right)^2}$

Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

# Agenda

1. Basic statistics
2. Uniques and keys
3. Functional dependencies
4. Inclusion dependencies and foreign keys
5. Outlook: Other dependencies and more



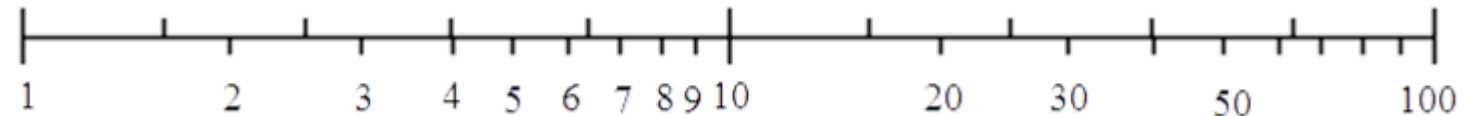
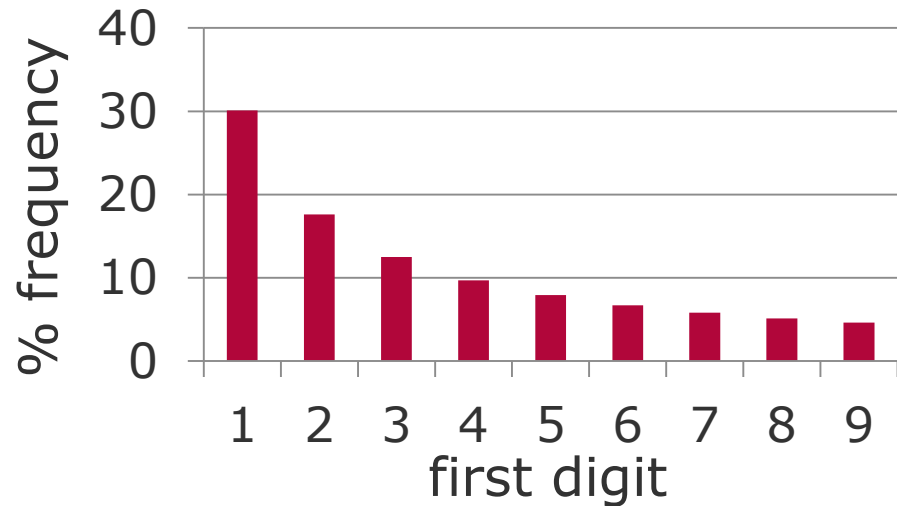
Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

## Cardinalities, distributions, and patterns

Category	Task	Description
<b>Cardinalities</b>	num-rows	Number of rows
	value length	Measurements of value lengths (min, max, median, and average)
	null values	Number or percentage of null values
	distinct	Number of distinct values; aka "cardinality"
	uniqueness	Number of distinct values divided by number of rows
<b>Value distributions</b>	histogram	Frequency histograms (equi-width, equi-depth, etc.)
	constancy	Frequency of most frequent value divided by number of rows
	quartiles	Three points that divide the (numeric) values into four equal groups
	soundex	Distribution of soundex codes
	first digit	Distribution of first digit in numeric values (Benford's law)
<b>Patterns, data types, and domains</b>	basic type	Generic data type: numeric, alphabetic, date, time
	data type	Concrete DBMS-specific data type: varchar, timestamp, etc.
	decimals	Maximum number of decimal places in numeric values
	precision	Maximum number of digits in numeric values
	patterns	Histogram of value patterns (Aa9...)
	data class	Semantic, generic data type: code, indicator, text, date/time, quantity, identifier, etc.
	domain	Classification of semantic domain: credit card, first name, city, phenotype, etc.

# Benford Law Frequency , a.k.a. “first digit law”

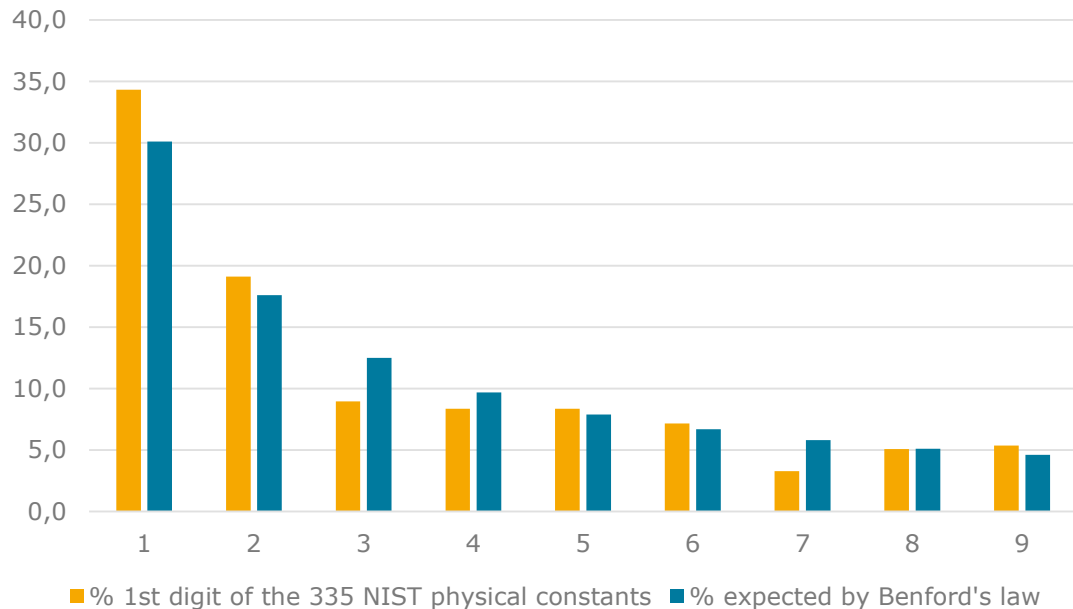
- Statement about the distribution of first digits  $d$  in (many) *naturally occurring* numbers:
  - $P(d) = \log_{10}(d + 1) - \log_{10}(d) = \log_{10}(1 + 1/d)$
  - Holds if  $\log(x)$  is uniformly distributed



Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

# Examples for Benford's Law

- Surface areas of 335 rivers
- Sizes of 3259 US populations
- 1800 molecular weights
- 5000 entries from a mathematical handbook
- Street addresses of the first 342 persons listed in American Men of Science
- $2^n$



## Heights of the 60 tallest structures

Leading digit	meters	
	Count	%
1	26	43.3%
2	7	11.7%
3	9	15.0%
4	6	10.0%
5	4	6.7%
6	1	1.7%
7	2	3.3%
8	5	8.3%
9	0	0.0%

In Benford's law
30.1%
17.6%
12.5%
9.7%
7.9%
6.7%
5.8%
5.1%
4.6%

[http://en.wikipedia.org/wiki/List\\_of\\_tallest\\_buildings\\_and\\_structures\\_in\\_the\\_world#Tallest\\_structure\\_by\\_category](http://en.wikipedia.org/wiki/List_of_tallest_buildings_and_structures_in_the_world#Tallest_structure_by_category)

**Felix Naumann**  
**Data Profiling**  
**Télécom ParisTech 2018**



# CODATA RECOMMENDED VALUES OF THE FUNDAMENTAL PHYSICAL CONSTANTS: 2014

NIST SP 961 (Sept/2015) Values from: P. J. Mohr, D. B. Newell, and B. N. Taylor, arXiv:1507.07956

A more extensive listing of constants is available in the above reference and on the NIST Physics Laboratory Web site [physics.nist.gov/constants](http://physics.nist.gov/constants).

The number in parentheses is the one-standard-deviation uncertainty in the last two digits of the given value.

Quantity	Symbol	Numerical value	Unit	Quantity	Symbol	Numerical value	Unit
speed of light in vacuum	$c, c_0$	299 792 458 (exact)	$\text{m s}^{-1}$	muon $g$ -factor $-2(1 + a_\mu)$	$g_\mu$	$-2.002\,331\,8418(13)$	
magnetic constant	$\mu_0$	$4\pi \times 10^{-7}$ (exact)	$\text{N A}^{-2}$	muon-proton magnetic moment ratio	$\mu_\mu/\mu_p$	$-3.183\,345\,142(71)$	
electric constant $1/\mu_0 c^2$	$\epsilon_0$	$8.854\,187\,817\dots \times 10^{-12}$	$\text{F m}^{-1}$	proton mass	$m_p$	$1.672\,621\,898(21) \times 10^{-27}$	kg
Newtonian constant of gravitation	$G$	$6.674\,08(31) \times 10^{-11}$	$\text{m}^3 \text{kg}^{-1} \text{s}^{-2}$	in u		$1.007\,276\,466\,879(91)$	u
Planck constant	$h$	$6.626\,070\,040(81) \times 10^{-34}$	J s	energy equivalent in MeV	$m_p c^2$	$938.272\,0813(58)$	MeV
in eV s		$4.135\,667\,662(25) \times 10^{-15}$	eV s	proton-electron mass ratio	$m_p/m_e$	$1836.152\,673\,89(17)$	
$h/2\pi$	$\hbar$	$1.054\,571\,800(13) \times 10^{-34}$	J s	proton magnetic moment	$\mu_p$	$1.410\,606\,7873(97) \times 10^{-26}$	J T <sup>-1</sup>
in eV s		$6.582\,119\,514(40) \times 10^{-16}$	eV s	to nuclear magneton ratio	$\mu_p/\mu_N$	$2.792\,847\,3508(85)$	
elementary charge	$e$	$1.602\,176\,6208(98) \times 10^{-19}$	C	proton magnetic shielding correction $1 - \mu'_p/\mu_p$	$\sigma'_p$	$25.691(11) \times 10^{-6}$	
magnetic flux quantum $h/2e$	$\Phi_0$	$2.067\,833\,831(13) \times 10^{-15}$	Wb	(H <sub>2</sub> O, sphere, 25 °C)			
Josephson constant $2e/h$	$K_J$	$483\,597.8525(30) \times 10^9$	Hz V <sup>-1</sup>	proton gyromagnetic ratio $2\mu_p/\hbar$	$\gamma_p$	$2.675\,221\,900(18) \times 10^8$	s <sup>-1</sup> T <sup>-1</sup>
von Klitzing constant $h/e^2 = \mu_0 c/2\alpha$	$R_K$	$25\,812.807\,4555(59)$	$\Omega$		$\gamma_p/2\pi$	$42.577\,478\,92(29)$	MHz T <sup>-1</sup>
Bohr magneton $e\hbar/2m_e$	$\mu_B$	$927.400\,9994(57) \times 10^{-26}$	J T <sup>-1</sup>	shielded proton gyromagnetic ratio $2\mu'_p/\hbar$	$\gamma'_p$	$2.675\,153\,171(33) \times 10^8$	s <sup>-1</sup> T <sup>-1</sup>
in eV T <sup>-1</sup>		$5.788\,381\,8012(26) \times 10^{-5}$	eV T <sup>-1</sup>	(H <sub>2</sub> O, sphere, 25 °C)			
nuclear magneton $e\hbar/2m_p$	$\mu_N$	$5.050\,783\,699(31) \times 10^{-27}$	J T <sup>-1</sup>		$\gamma'_p/2\pi$	$42.576\,385\,07(53)$	MHz T <sup>-1</sup>
in eV T <sup>-1</sup>		$3.152\,451\,2550(15) \times 10^{-8}$	eV T <sup>-1</sup>	neutron mass in u	$m_n$	$1.008\,664\,915\,88(49)$	u
fine-structure constant $e^2/4\pi\epsilon_0\hbar c$	$\alpha$	$7.297\,352\,5664(17) \times 10^{-3}$		energy equivalent in MeV	$m_n c^2$	$939.565\,4133(58)$	MeV
inverse fine-structure constant	$\alpha^{-1}$	$137.035\,999\,139(31)$		neutron-proton mass ratio	$m_n/m_p$	$1.001\,378\,418\,98(51)$	
Rydberg constant $\alpha^2 m_e c/2h$	$R_\infty$	$10\,973\,731.568\,508(65)$	m <sup>-1</sup>	neutron magnetic moment	$\mu_n$	$-0.966\,236\,50(23) \times 10^{-26}$	J T <sup>-1</sup>
energy equivalent in eV	$R_\infty c$	$3.289\,841\,960\,355(19) \times 10^{15}$	Hz	to nuclear magneton ratio	$\mu_n/\mu_N$	$-1.913\,042\,73(45)$	
Bohr radius $\alpha/4\pi R_\infty = 4\pi\epsilon_0\hbar^2/m_e e^2$	$a_0$	$0.529\,177\,210\,67(12) \times 10^{-10}$	m	deuteron mass in u	$m_d$	$2.013\,553\,212\,745(40)$	u
Hartree energy $e^2/4\pi\epsilon_0 a_0 = 2R_\infty hc = \alpha^2 m_e c^2$	$E_h$	$4.359\,744\,650(54) \times 10^{-18}$	J	energy equivalent in MeV	$m_d c^2$	$1875.612\,928(12)$	MeV
in eV		$27.211\,386\,02(17)$	eV	deuteron-proton mass ratio	$m_d/m_p$	$1.999\,007\,500\,87(19)$	
electron mass	$m_e$	$9.109\,383\,56(11) \times 10^{-31}$	kg	deuteron magnetic moment	$\mu_d$	$0.433\,073\,5040(36) \times 10^{-26}$	J T <sup>-1</sup>
in u		$5.485\,799\,090\,70(16) \times 10^{-4}$	u	to nuclear magneton ratio	$\mu_d/\mu_N$	$0.857\,438\,2311(48)$	
energy equivalent in MeV	$m_e c^2$	$0.510\,998\,9461(31)$	MeV	helion ( <sup>3</sup> He nucleus) mass in u	$m_h$	$3.014\,932\,246\,73(12)$	u
electron-muon mass ratio	$m_e/m_\mu$	$4.836\,331\,70(11) \times 10^{-3}$		energy equivalent in MeV	$m_h c^2$	$2808.391\,586(17)$	MeV
electron-proton mass ratio	$m_e/m_p$	$5.446\,170\,213\,52(52) \times 10^{-4}$		shielded helion magnetic moment	$\mu'_h$	$-1.074\,553\,080(14) \times 10^{-26}$	J T <sup>-1</sup>
electron charge to mass quotient	$-e/m_e$	$-1.758\,820\,024(11) \times 10^{11}$	C kg <sup>-1</sup>	(gas, sphere, 25 °C)			
Compton wavelength $h/m_e c$	$\lambda_C$	$2.426\,310\,2367(11) \times 10^{-12}$	m	to Bohr magneton ratio	$\mu'_h/\mu_B$	$-1.158\,671\,471(14) \times 10^{-3}$	
$\lambda_C/2\pi = \alpha a_0 = \alpha^2/4\pi R_\infty$	$\lambda_C$	$386.159\,267\,64(18) \times 10^{-15}$	m	to nuclear magneton ratio	$\mu'_h/\mu_N$	$-2.127\,497\,720(25)$	
classical electron radius $\alpha^2 a_0$	$r_e$	$2.817\,940\,3227(19) \times 10^{-15}$	m	alpha particle mass in u	$m_\alpha$	$4.001\,506\,179\,127(63)$	u
Thomson cross section $(8\pi/3)r_e^2$	$\sigma_e$	$0.665\,245\,871\,58(91) \times 10^{-28}$	m <sup>2</sup>	energy equivalent in MeV	$m_\alpha c^2$	$3727.379\,378(23)$	MeV
electron magnetic moment	$\mu_e$	$-928.476\,4620(57) \times 10^{-26}$	J T <sup>-1</sup>	Avogadro constant	$N_A, L$	$6.022\,140\,857(74) \times 10^{23}$	mol <sup>-1</sup>
to Bohr magneton ratio	$\mu_e/\mu_B$	$-1.001\,159\,652\,180\,91(26)$		atomic mass constant $\frac{1}{12} m(^{12}\text{C}) = 1 \text{ u}$	$m_u$	$1.660\,539\,040(20) \times 10^{-27}$	kg
to nuclear magneton ratio	$\mu_e/\mu_N$	$-1838.281\,972\,34(17)$		energy equivalent in MeV	$m_u c^2$	$931.494\,0954(57)$	MeV
electron magnetic moment anomaly $ \mu_e /\mu_B - 1$	$a_e$	$1.159\,652\,180\,91(26) \times 10^{-3}$		Faraday constant $N_A e$	$F$	$96\,485.332\,89(59)$	C mol <sup>-1</sup>
electron $g$ -factor $-2(1 + a_e)$	$g_e$	$-2.002\,319\,304\,361\,82(52)$		molar gas constant	$R$	$8.314\,4598(48)$	J mol <sup>-1</sup> K <sup>-1</sup>
electron-proton magnetic moment ratio	$\mu_e/\mu_p$	$-658.210\,6866(20)$		Boltzmann constant $R/N_A$	$k$	$1.380\,648\,52(79) \times 10^{-23}$	J K <sup>-1</sup>
muon mass in u	$m_\mu$	$0.113\,428\,9257(25)$	u	in eV K <sup>-1</sup>		$8.617\,3303(50) \times 10^{-5}$	eV K <sup>-1</sup>
energy equivalent in MeV	$m_\mu c^2$	$105.658\,3745(24)$	MeV	molar volume of ideal gas $RT/p$	$V_m$	$22.413\,962(13) \times 10^{-3}$	m <sup>3</sup> mol <sup>-1</sup>
muon-electron mass ratio	$m_\mu/m_e$	$206.768\,2826(46)$		( $T = 273.15 \text{ K}, p = 101.325 \text{ kPa}$ )			
muon magnetic moment	$\mu_\mu$	$-4.490\,448\,26(10) \times 10^{-26}$	J T <sup>-1</sup>	Stefan-Boltzmann constant $\pi^2 k^4/60\hbar^3 c^2$	$\sigma$	$5.670\,367(13) \times 10^{-8}$	W m <sup>-2</sup> K <sup>-4</sup>
to Bohr magneton ratio	$\mu_\mu/\mu_B$	$-4.841\,970\,48(11) \times 10^{-3}$		first radiation constant $2\pi\hbar c^2$	$c_1$	$3.741\,771\,790(46) \times 10^{-16}$	W m <sup>2</sup>
to nuclear magneton ratio	$\mu_\mu/\mu_N$	$-8.890\,597\,05(20)$		second radiation constant $hc/k$	$c_2$	$1.438\,777\,36(83) \times 10^{-2}$	m K
muon magnetic moment anomaly				Wien displacement law constant			
$ \mu_\mu /(e\hbar/2m_\mu) - 1$	$a_\mu$	$1.165\,920\,89(63) \times 10^{-3}$		$b = \lambda_{\text{max}} T = c_2/4.965\,114\,231\dots$	$b$	$2.897\,7729(17) \times 10^{-3}$	m K
				Cu x unit: $\lambda(\text{Cu K}\alpha_1)/1\,537.400$	$xu(\text{Cu K}\alpha_1)$	$1.002\,076\,97(28) \times 10^{-13}$	m
				Mo x unit: $\lambda(\text{Mo K}\alpha_1)/707.831$	$xu(\text{Mo K}\alpha_1)$	$1.002\,099\,52(53) \times 10^{-13}$	m
<b>Energy equivalents</b>							
$(1 \text{ m}^{-1})c = 299\,792\,458 \text{ Hz}$	$(1 \text{ Hz})h/k = 4.799\,2447(28) \times 10^{-11} \text{ K}$			$(1 \text{ J}) = 6.241\,509\,126(38) \times 10^{18} \text{ eV}$		$(1 \text{ eV})/c^2 = 1.073\,544\,1105(66) \times 10^{-9} \text{ u}$	
$(1 \text{ m}^{-1})hc/k = 1.438\,777\,36(83) \times 10^{-2} \text{ K}$	$(1 \text{ Hz})h = 4.135\,667\,662(25) \times 10^{-15} \text{ eV}$			$(1 \text{ eV}) = 1.602\,176\,6208(98) \times 10^{-19} \text{ J}$		$(1 \text{ kg}) = 6.022\,140\,857(74) \times 10^{26} \text{ u}$	
$(1 \text{ m}^{-1})hc = 1.239\,841\,9739(76) \times 10^{-6} \text{ eV}$	$(1 \text{ K})k/hc = 69.503\,457(40) \text{ m}^{-1}$			$(1 \text{ eV})/hc = 8.065\,544\,005(50) \times 10^5 \text{ m}^{-1}$		$(1 \text{ u}) = 1.660\,539\,040(20) \times 10^{-27} \text{ kg}$	
$(1 \text{ m}^{-1})h/c = 1.331\,025\,049\,00(61) \times 10^{-15} \text{ u}$	$(1 \text{ K})k/h = 2.083\,6612(12) \times 10^{10} \text{ Hz}$			$(1 \text{ eV})/h = 2.417\,989\,262(15) \times 10^{14} \text{ Hz}$		$(1 \text{ u})c/h = 7.513\,006\,6166(34) \times 10^{14} \text{ m}^{-1}$	
$(1 \text{ Hz})/c = 3.335\,640\,951\dots \times 10^{-9} \text{ m}^{-1}$	$(1 \text{ K})k = 8.617\,3303(50) \times 10^{-5} \text{ eV}$			$(1 \text{ eV})/k = 1.160\,452\,21(67) \times 10^4 \text{ K}$		$(1 \text{ u})c^2 = 931.494\,0954(57) \times 10^6 \text{ eV}$	



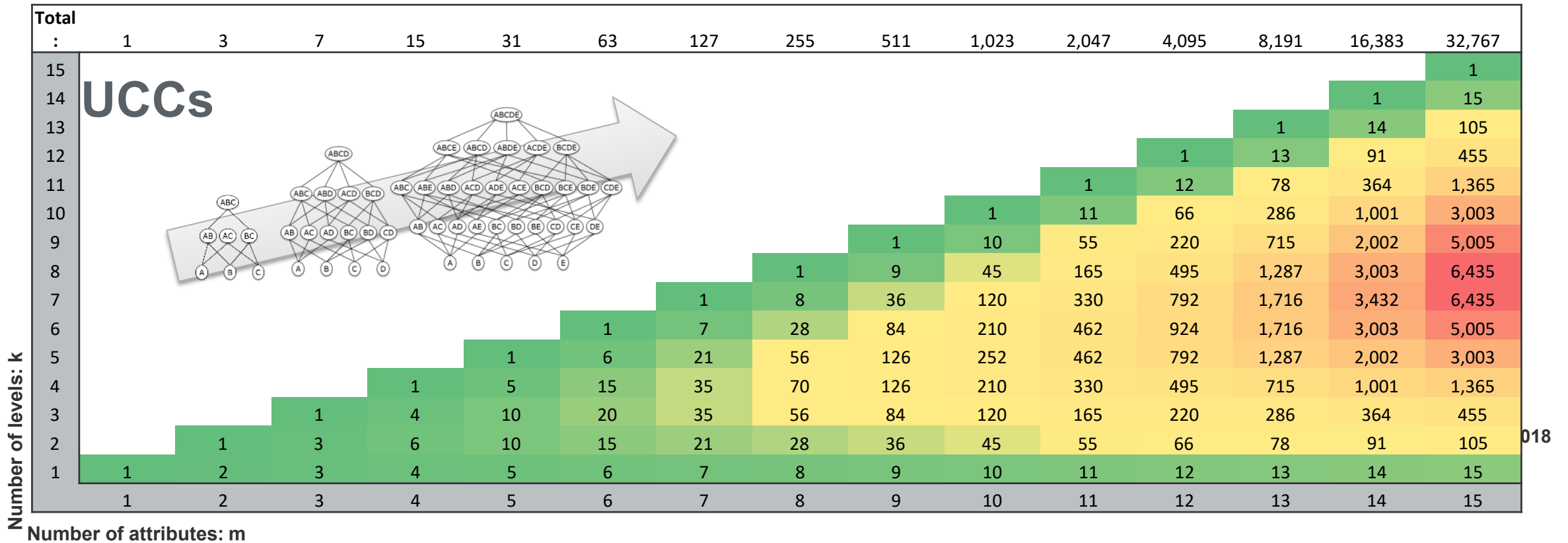
# Uniqueness, keys, and foreign keys

- Uniqueness and keys
  - Unique column: Only unique values
  - Unique column combination: Only unique value combinations
    - Minimality: No column subset is unique
  - Key candidate: No null values
  - Key: Only human expert can decide
    - UCC is prerequisite

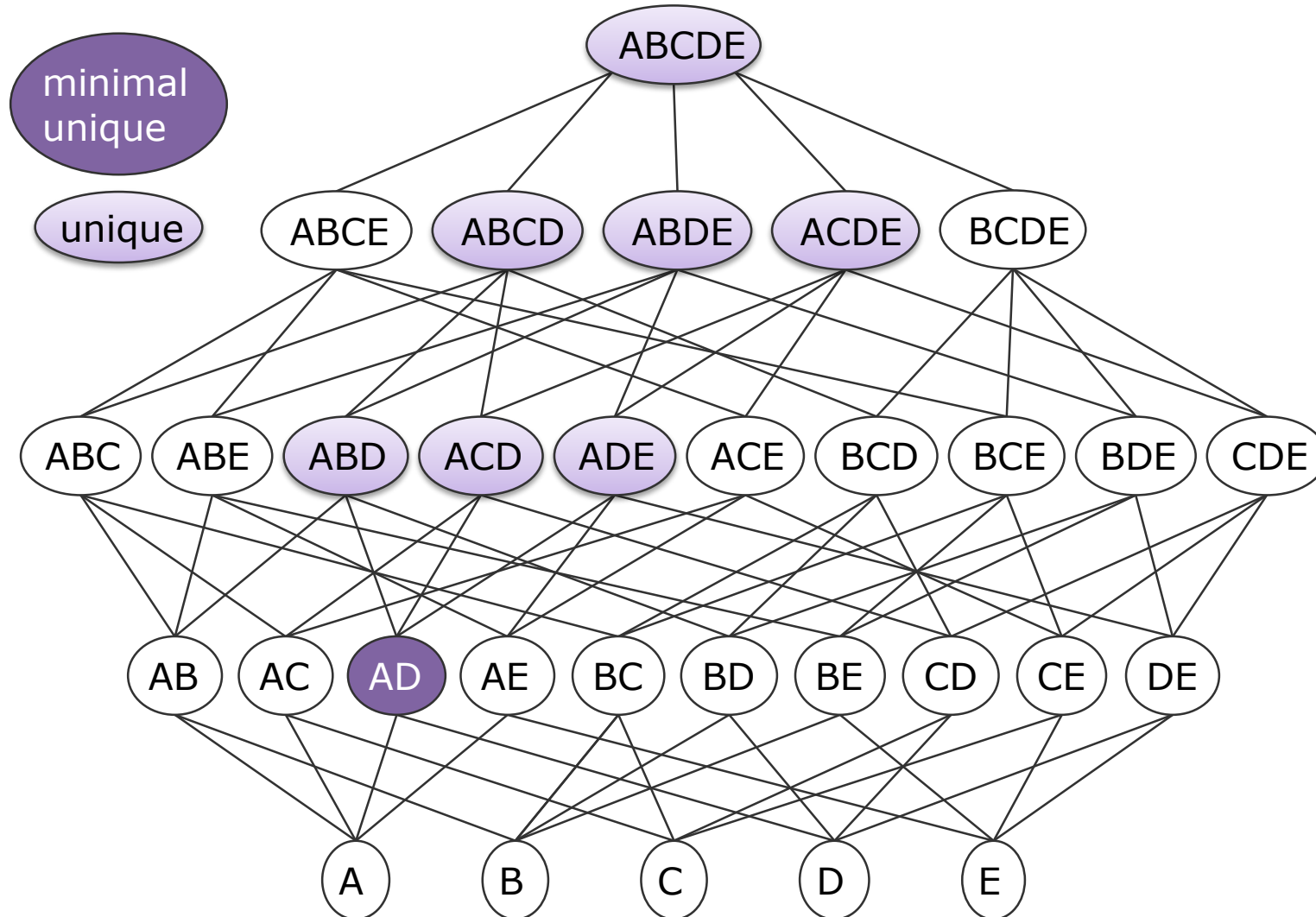
- Uniques: {A, AB, AC, BC, ABC}
- Minimal uniques: {A, BC}
- (Maximal) Non-uniques: {B, C}

A	B	C
a	1	x
b	2	x
c	2	y

# Candidate Set Growth for UCCs

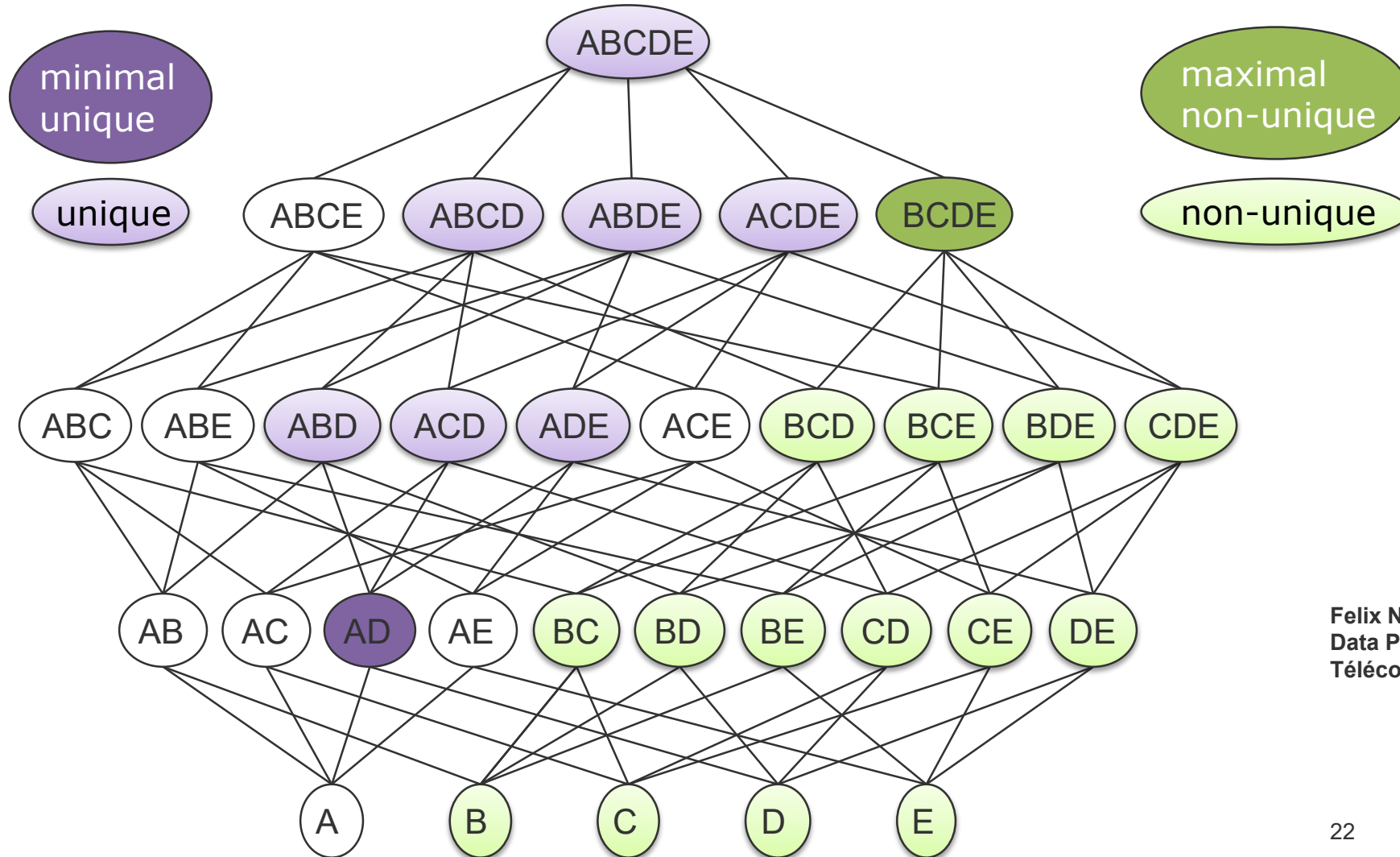


# Pruning effect of a pair



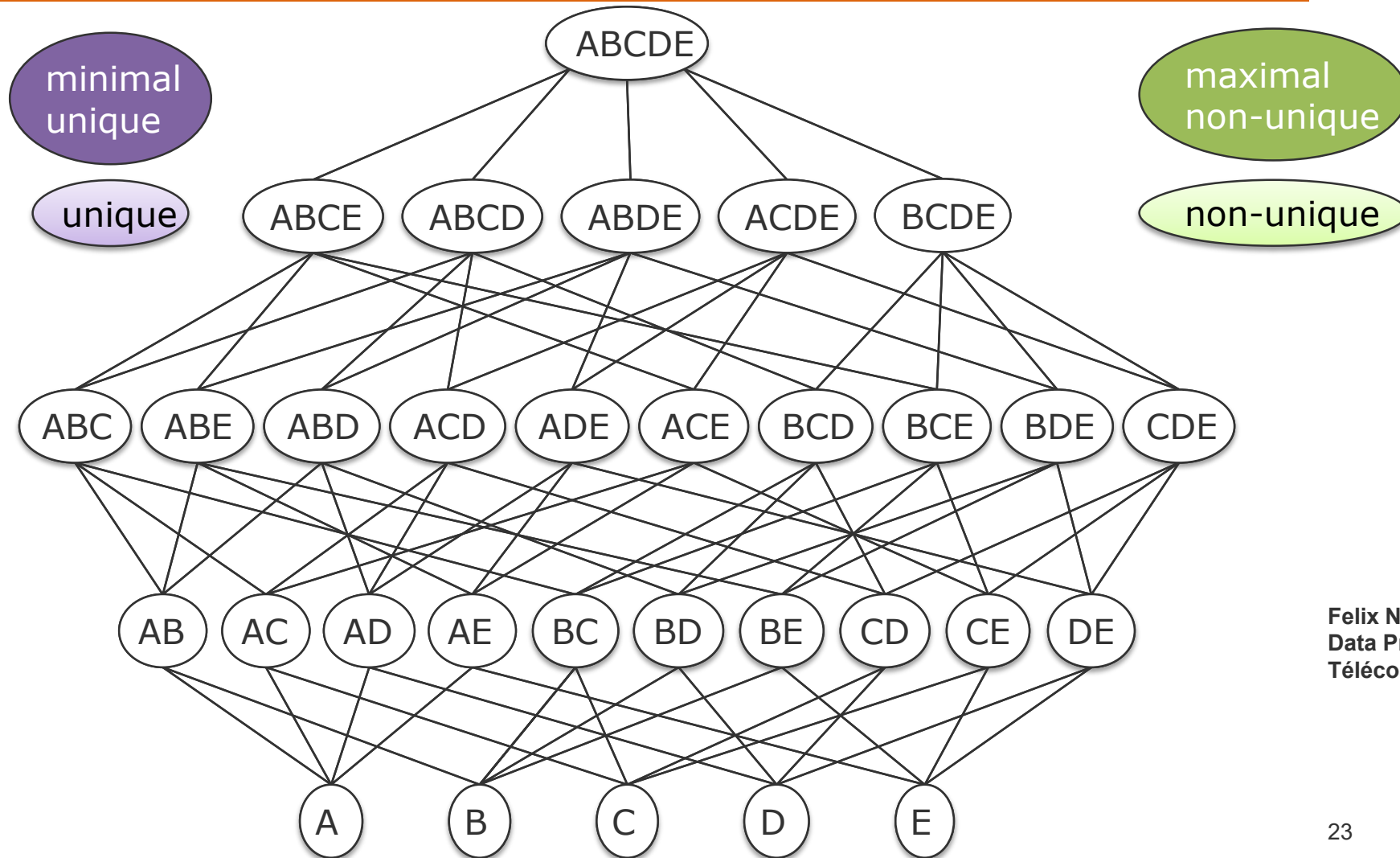
Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

# Pruning both ways



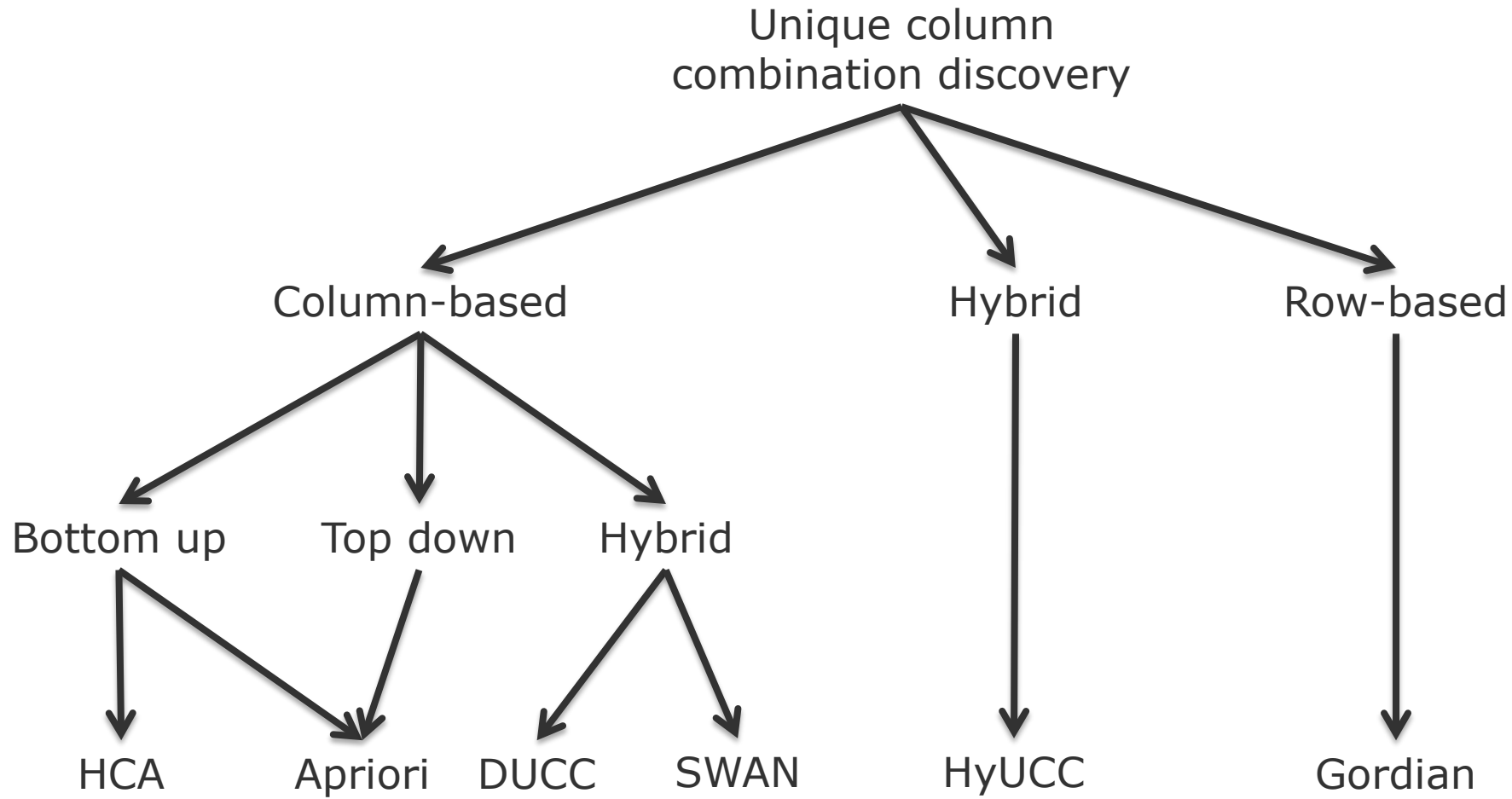
Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

# Apriori visualized



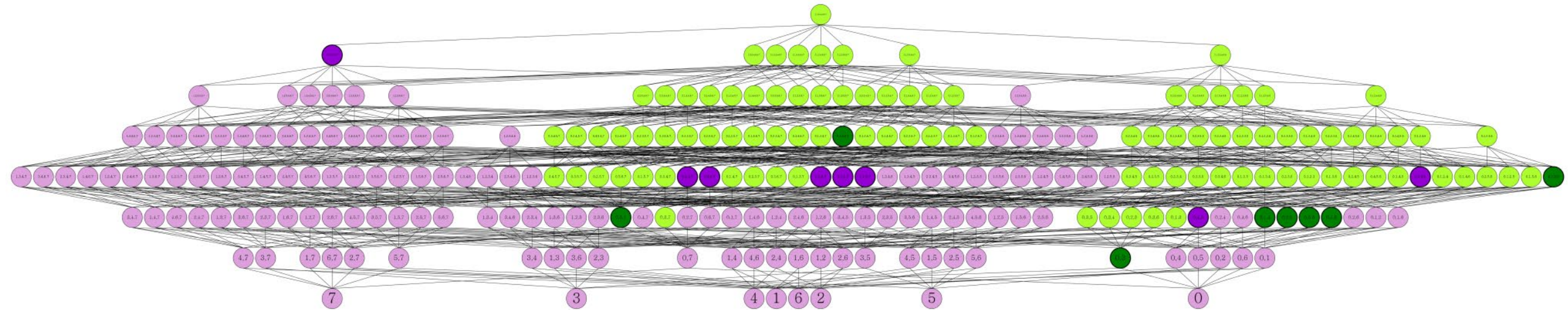
Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

# Discovery Algorithms





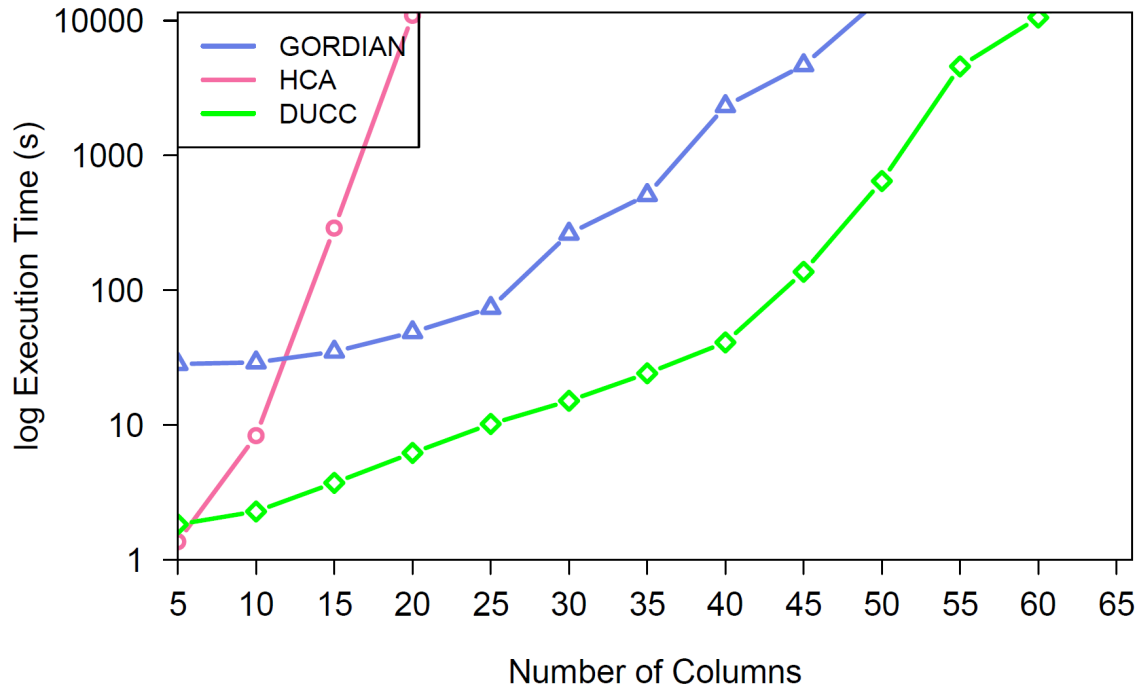
# DUCC – Detecting Unique Column Combinations



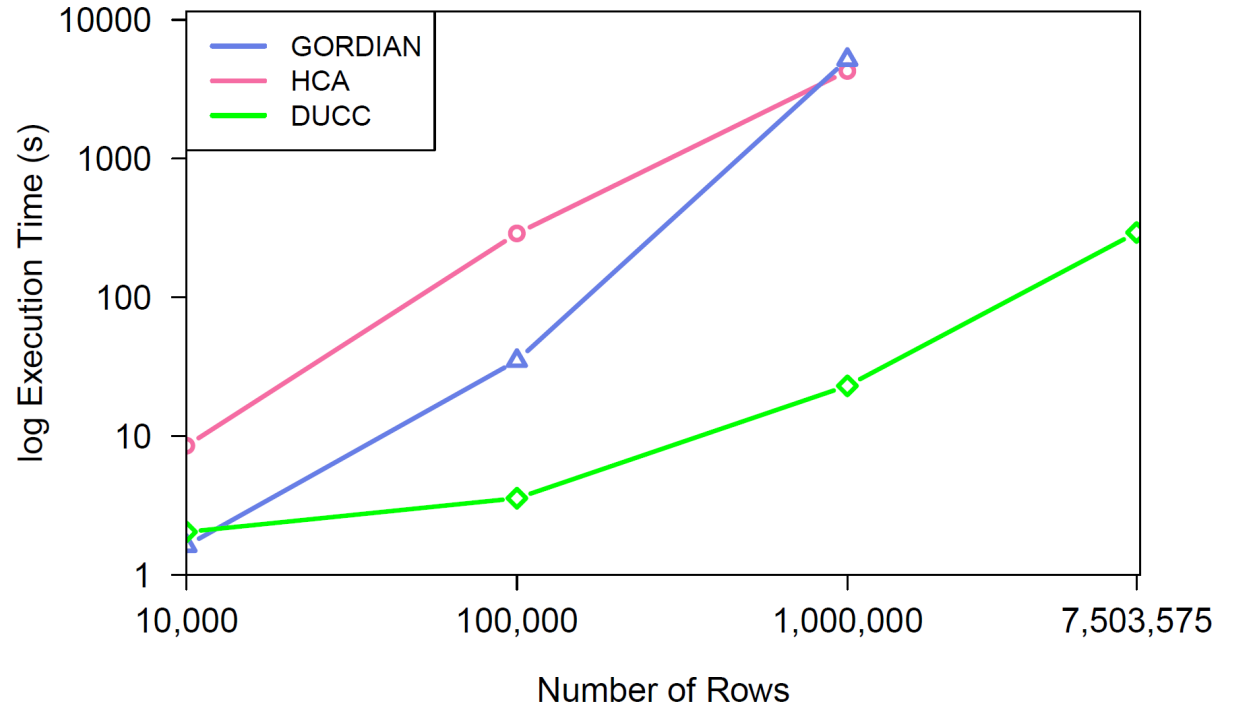
Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

# Scalability in the number of columns and rows

■ NCVoter data, 100k rows



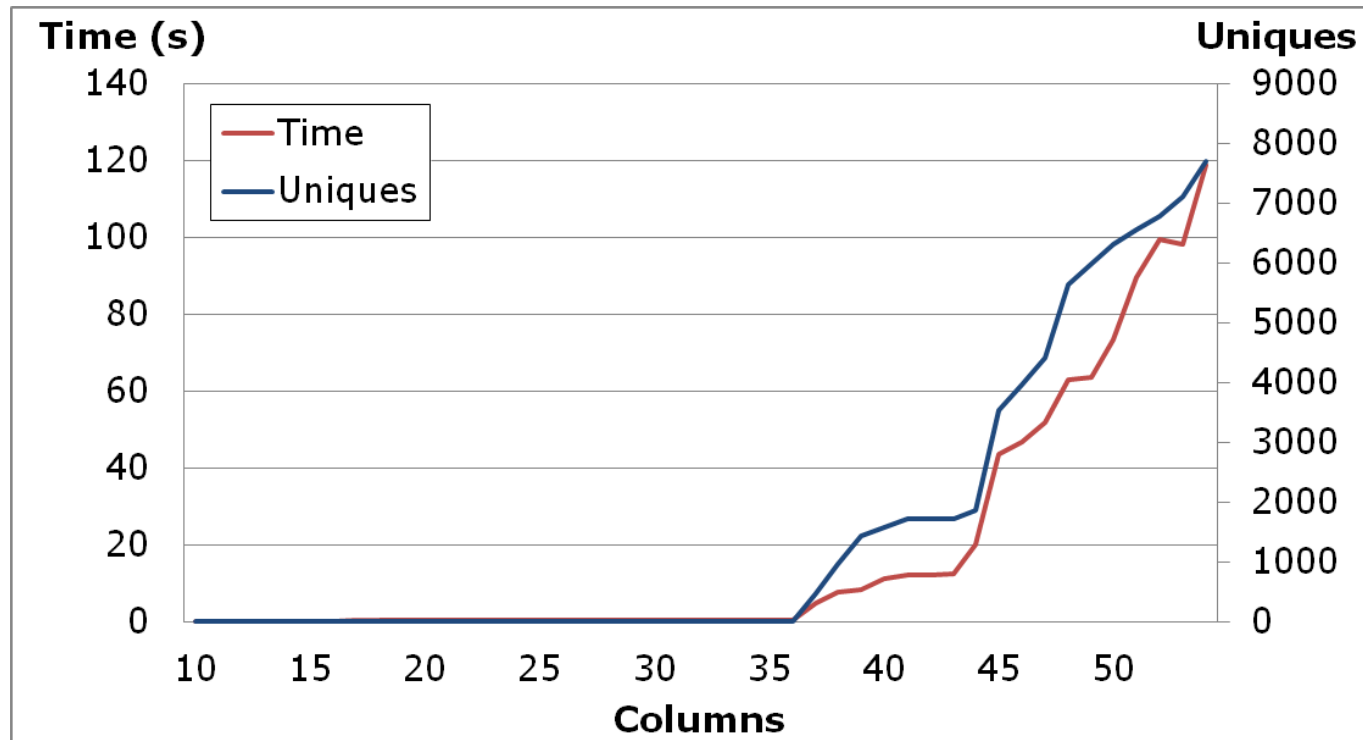
■ NCVoter, 15 columns



■ New hybrid version shaves off another order of magnitude

# Analysis of DUCC

- Runtime mainly depends on size of solution set



- Worst case: solution set in the middle of lattice:  $\binom{n}{n/2}$  uniques

## Uniques and non-uniques in NC-voter data

---

- **A minimal unique:** voter\_reg\_num, zip\_code, race\_code
- **A maximal non-unique:** voter\_reg\_num, status\_cd, voter\_status\_desc, reason\_cd, voter\_status\_reason\_desc, absent\_ind, name\_prefx\_cd, name\_sufx\_cd, half\_code, street\_dir, street\_type\_cd, street\_sufx\_cd, unit\_designator, unit\_num, state\_cd, mail\_addr2, mail\_addr3, mail\_addr4, mail\_state, area\_cd, phone\_num, full\_phone\_number, drivers\_lic, race\_code, race\_desc, ethnic\_code, ethnic\_desc, party\_cd, party\_desc, sex\_code, sex, birth\_place, precinct\_abbrev, precinct\_desc, municipality\_abbrev, municipality\_desc, ward\_abbrev, ward\_desc, cong\_dist\_abbrev, cong\_dist\_desc, super\_court\_abbrev, super\_court\_desc, judic\_dist\_abbrev, judic\_dist\_desc, nc\_senate\_abbrev, nc\_senate\_desc, nc\_house\_abbrev, nc\_house\_desc, county\_commiss\_abbrev, county\_commiss\_desc, township\_abbrev, township\_desc, school\_dist\_abbrev, school\_dist\_desc, fire\_dist\_abbrev, fire\_dist\_desc, water\_dist\_abbrev, water\_dist\_desc, sewer\_dist\_abbrev, sewer\_dist\_desc, sanit\_dist\_abbrev, sanit\_dist\_desc, rescue\_dist\_abbrev, rescue\_dist\_desc, munic\_dist\_abbrev, munic\_dist\_desc, dist\_1\_abbrev, dist\_1\_desc, dist\_2\_abbrev, dist\_2\_desc, confidential\_ind, age, vtd\_abbrev, vtd\_desc

# Agenda

1. Basic statistics
2. Uniques and keys
3. **Functional dependencies**
4. Inclusion dependencies and foreign keys
5. Outlook: Other dependencies and more



Felix Naumann  
Data Profiling  
Télécom ParisTech 2018



Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

# Functional Dependencies

Person	Lineage	Hair	Religion
			New gods
			New Gods
			Old gods
			New gods
			Old gods

Some Functional Dependencies:

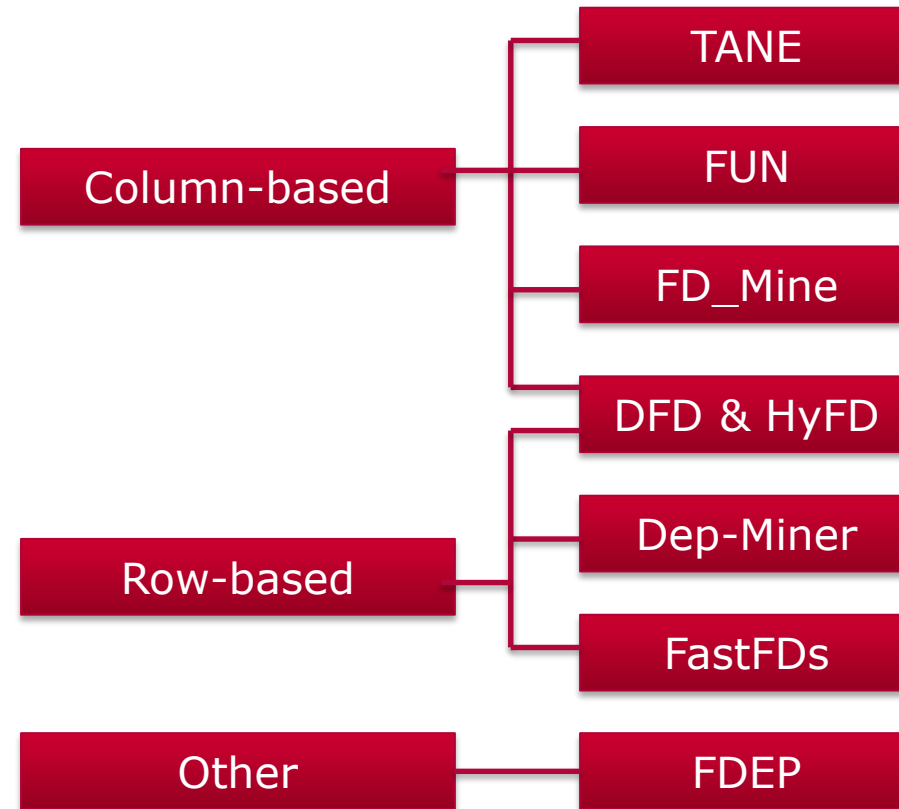
1. Person → Lineage
2. Person → Hair
3. Person → Religion
4. Lineage → Hair
5. Religion, Hair → Lineage
6. ...

Ned Stark: „#4 looks like a reasonable quality constraint“

Ned Stark: „I believe Joffrey violates my database constraint.“

# Uses and Algorithms for FDs

- Schema design
  - Normalization
  - Keys
- Data cleansing
- Query optimization
- Schema design and normalization
- Key discovery

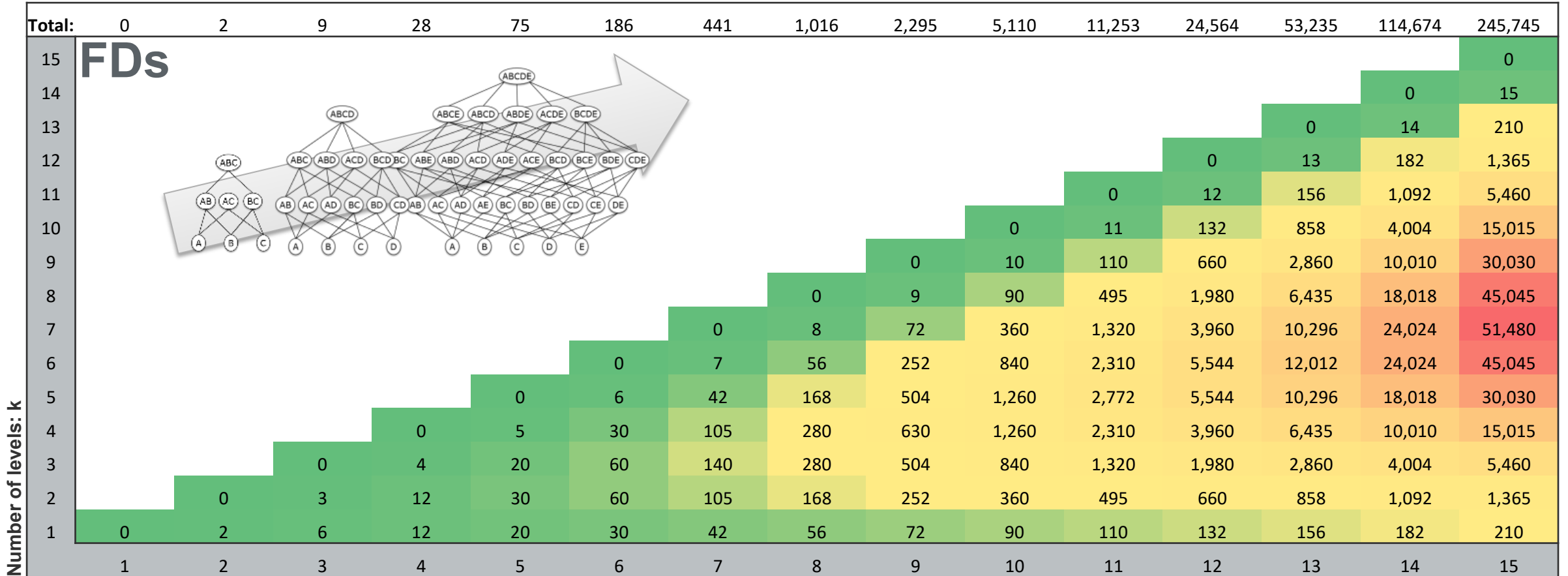


## Naïve discovery approach

- For each column combination X
  - For each pair of tuples (t1,t2)
    - If  $t1[X \setminus A] = t2[X \setminus A]$  and  $t1[A] \neq t2[A]$ : Break
- Exponential in number of attributes times number of rows squared

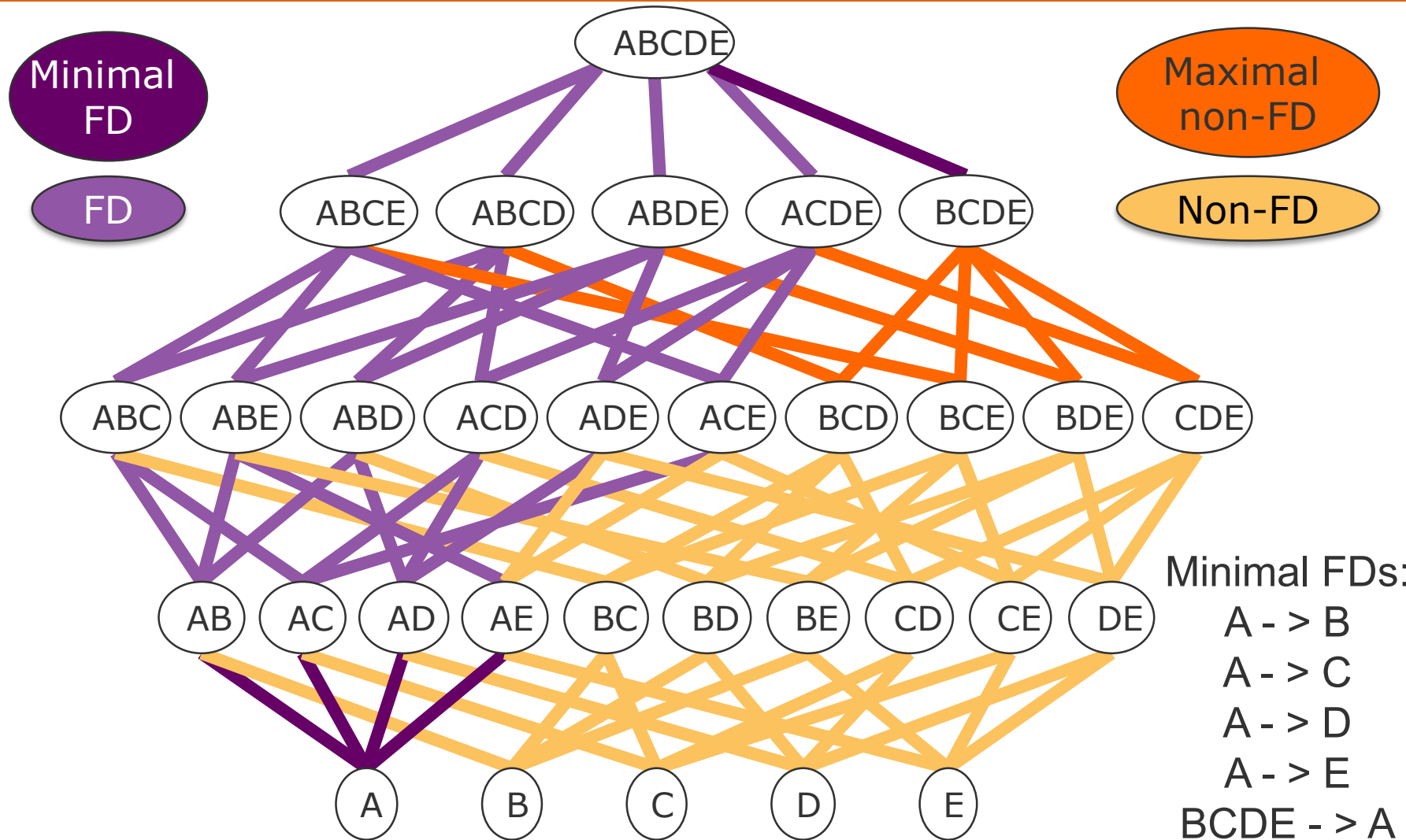


# Candidate Set Growth for FDs



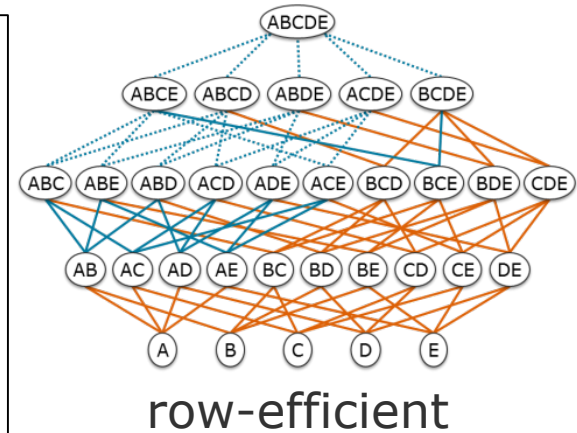
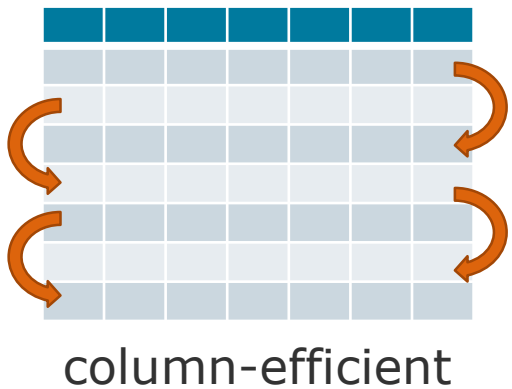
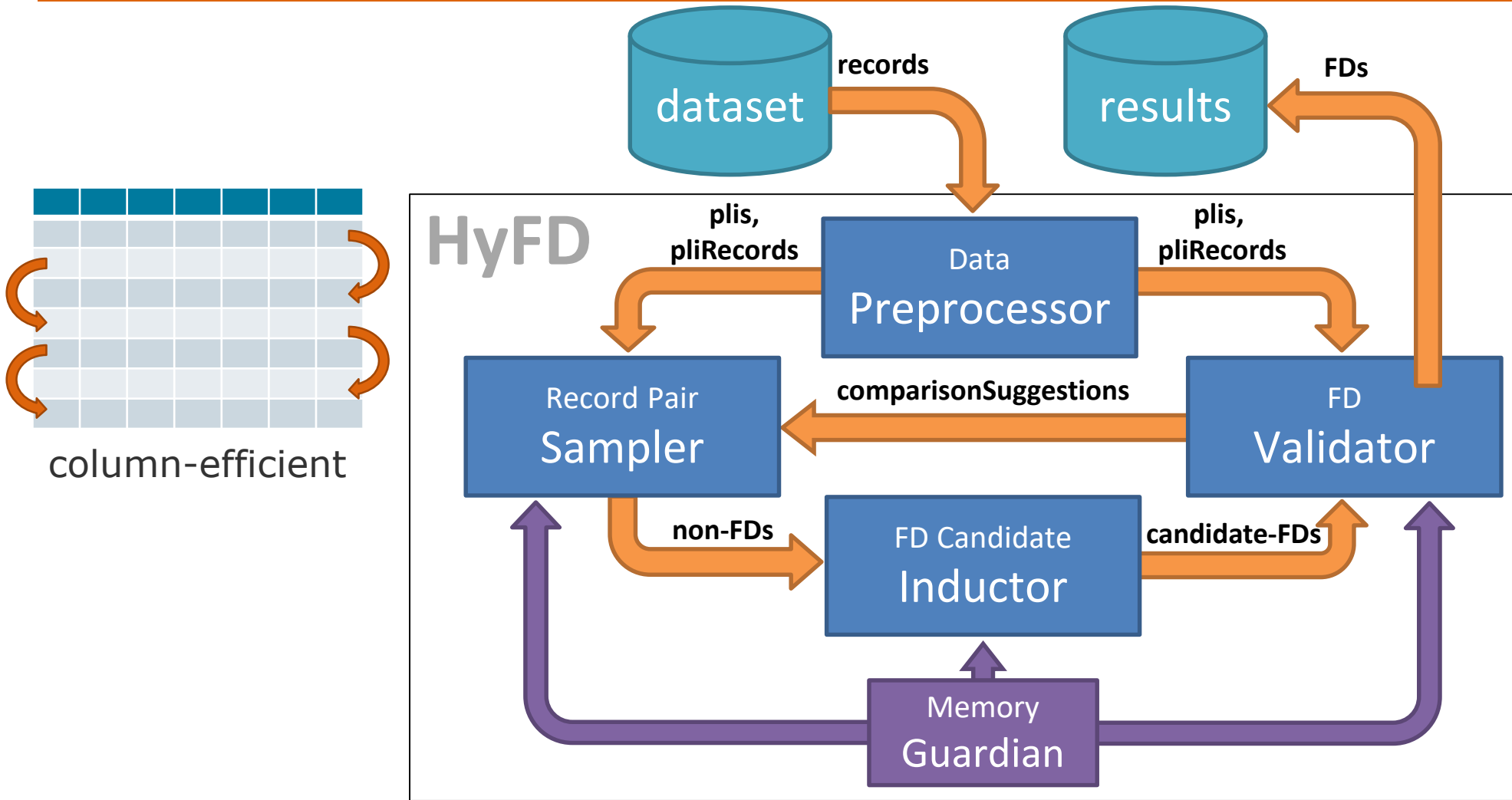
Number of attributes: m

Again: Model in lattice – edges represent FDs



Minimal FDs:  
 $A \rightarrow B$   
 $A \rightarrow C$   
 $A \rightarrow D$   
 $A \rightarrow E$   
 $BCDE \rightarrow A$

# HyFD: Hybrid FD Discovery



Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

# Functional Dependencies: State of the Art

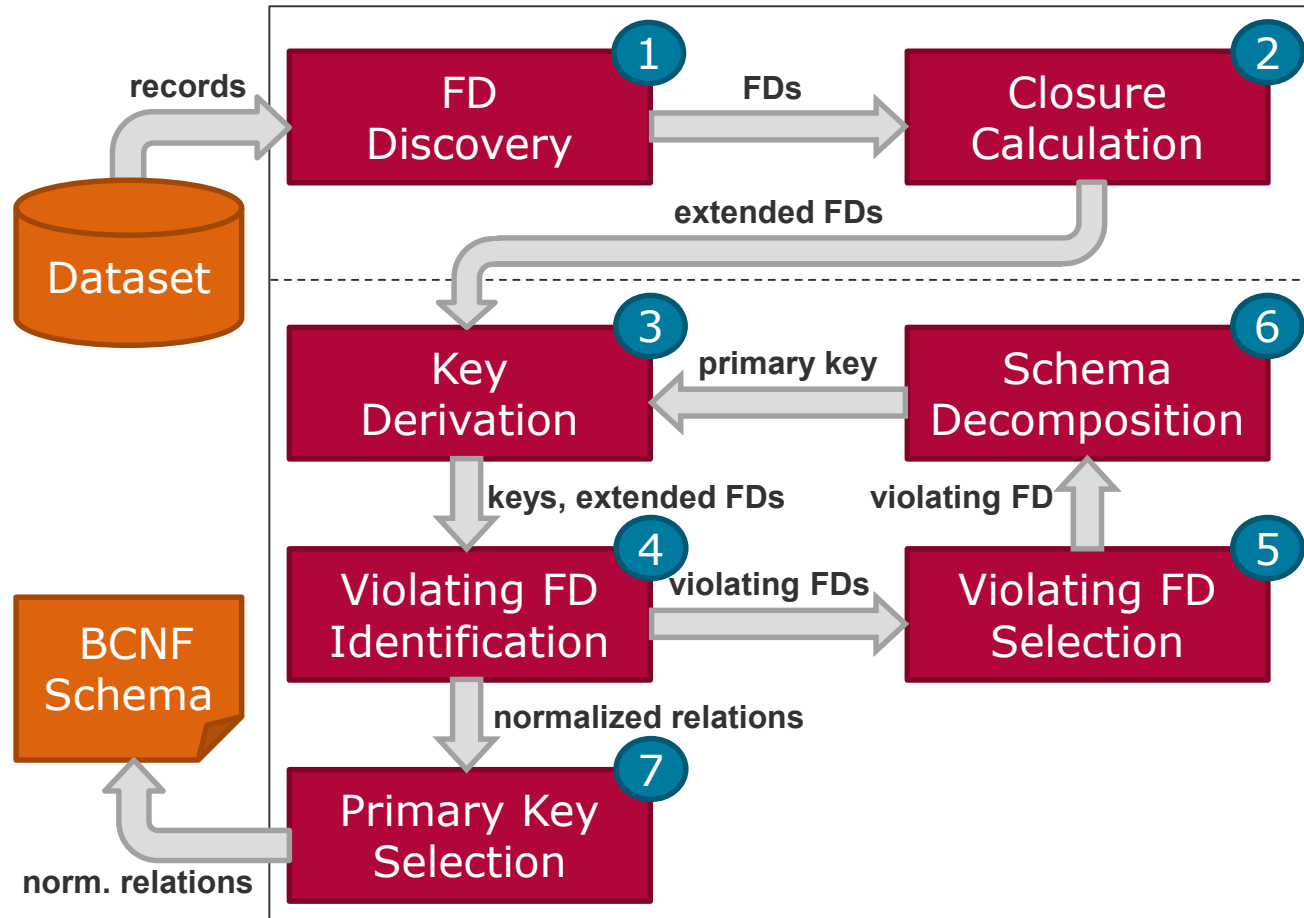
Dataset	Cols [#]	Rows [#]	Size [KB]	FDs [#]	TANE [12]	FUN [18]	FD_MINE [25]	DFD [1]	DEP-MINER [16]	FASTFDs [24]	FDEP [9]	HyFD
iris	5	150	5	4	1.1	<b>0.1</b>	0.2	0.2	0.2	0.2	<b>0.1</b>	<b>0.1</b>
balance-scale	5	625	7	1	1.2	<b>0.1</b>	0.2	0.3	0.3	0.3	0.2	<b>0.1</b>
chess	7	28,056	519	1	2.9	1.1	3.8	1.0	174.6	164.2	125.5	<b>0.2</b>
abalone	9	4,177	187	137	2.1	0.6	1.8	1.1	3.0	2.9	3.8	<b>0.2</b>
nursery	9	12,960	1,024	1	4.1	1.8	7.1	0.9	121.2	118.9	46.8	<b>0.5</b>
breast-cancer	11	699	20	46	2.3	0.6	2.2	0.8	1.1	1.1	0.5	<b>0.2</b>
bridges	13	108	6	142	2.2	0.6	4.2	0.9	0.5	0.6	0.2	<b>0.1</b>
echocardiogram	13	132	6	527	1.6	0.4	69.9	1.2	0.5	0.5	0.2	<b>0.1</b>
adult	14	48,842	3,528	78	67.4	111.6	531.5	5.9	6039.2	6033.8	860.2	1.1
letter	17	20,000	695	61	260.0	529.0	7204.8	6.0	1090.0	1015.5	291.3	<b>3.4</b>
ncvoter	19	1,000	151	758	4.3	4.0	ML	5.1	11.4	1.9	1.1	<b>0.4</b>
hepatitis	20	155	8	8,250	12.2	175.9	ML	326.7	5576.5	9.5	0.8	<b>0.6</b>
horse	27	368	25	128,727	457.0	TL	ML	TL	TL	385.8	7.2	<b>7.1</b>
fd-reduced-30	30	250,000	69,581	89,571	<b>41.1</b>	77.7	ML	TL	377.2	382.4	TL	513.0
plista	63	1,000	568	178,152	ML	ML	ML	TL	TL	TL	26.9	<b>21.8</b>
flight	109	1,000	575	982,631	ML	ML	ML	TL	TL	TL	216.5	<b>53.4</b>
uniprot	223	1,000	2,439	>2,437,556	ML	ML	ML	TL	TL	TL	ML	<b>&gt;5254.7</b>

Results larger than 1,000 FDs are only counted

TL: time limit of 4 hours exceeded

ML: memory limit of 100 GB exceeded

# Automatic BCNF Normalization



# Normalization results: TPC-H

( <u>linenumber</u> , extendedprice, discount, tax, returnflag, shipdate, commitdate, receiptdate, comment, <u>orderkey</u> , partkey)	<b>LINEITEM</b>
→( <u>linenumber</u> , <u>extendedprice</u> , <u>tax</u> , <u>commitdate</u> , <u>receiptdate</u> , shipinstruct)	
→( <u>extendedprice</u> , <u>discount</u> , shipmode, <u>orderkey</u> )	
→(quantity, <u>extendedprice</u> , <u>partkey</u> )	
→(linestatus, <u>shipdate</u> )	
→( <u>tax</u> , <u>returnflag</u> , <u>orderkey</u> , <u>partkey</u> , suppkey)	
↳( <u>availqty</u> , supplycost, comment, <u>partkey</u> , <u>suppkey</u> )	<b>PARTSUPP</b>
↳( <u>partkey</u> , name, brand, type, size, container, retailprice, comment)	<b>PART</b>
↳↳(mfgr, <u>brand</u> )	
↳( <u>suppkey</u> , name, address, phone, acctbal, comment, nationkey)	<b>SUPPLIER</b>
↳↳( <u>nationkey</u> , name, comment, regionkey)	<b>NATION</b>
↳↳(shippriority, <u>regionkey</u> , name, comment)	<b>REGION</b>
→( <u>orderkey</u> , totalprice, orderdate, orderpriority, clerk, comment, custkey)	<b>ORDERS</b>
↳(orderstatus, <u>totalprice</u> , <u>orderdate</u> )	
↳( <u>custkey</u> , name, address, phone, acctbal, mktsegment, comment)	<b>CUSTOMER</b>



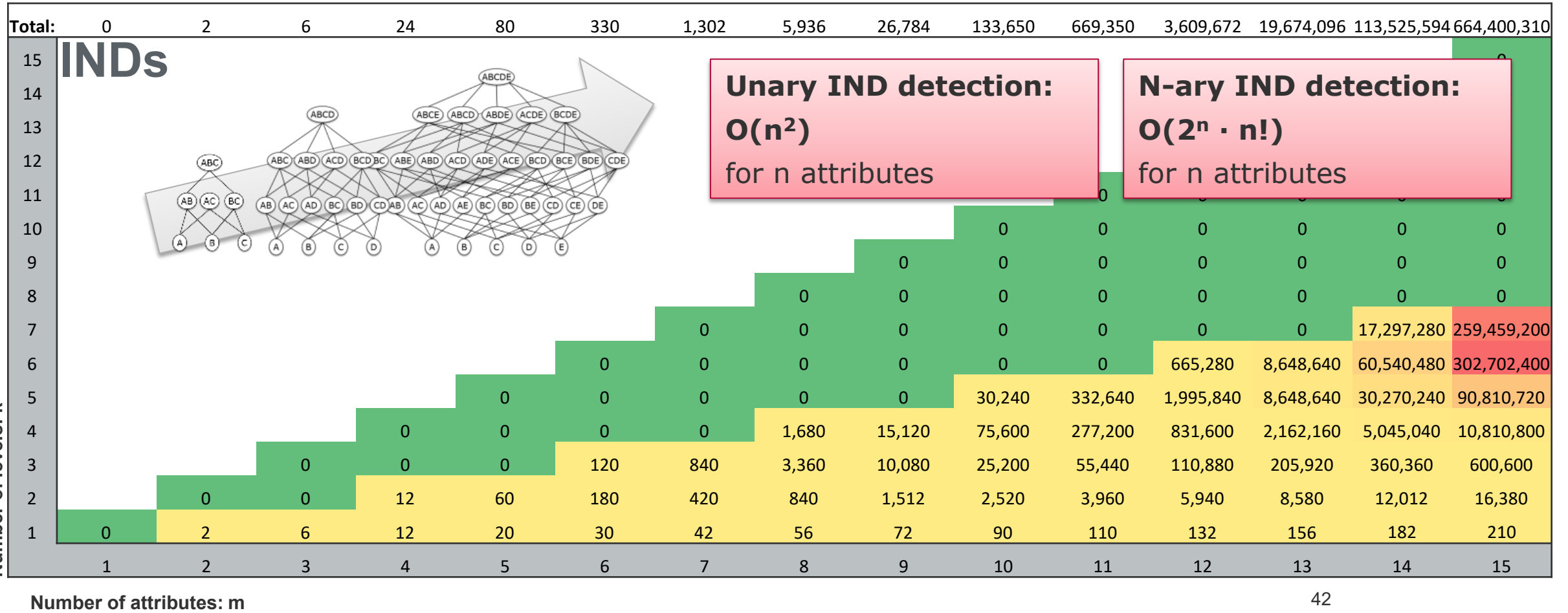
## IND discovery $R[X] \subseteq S[Y]$

---

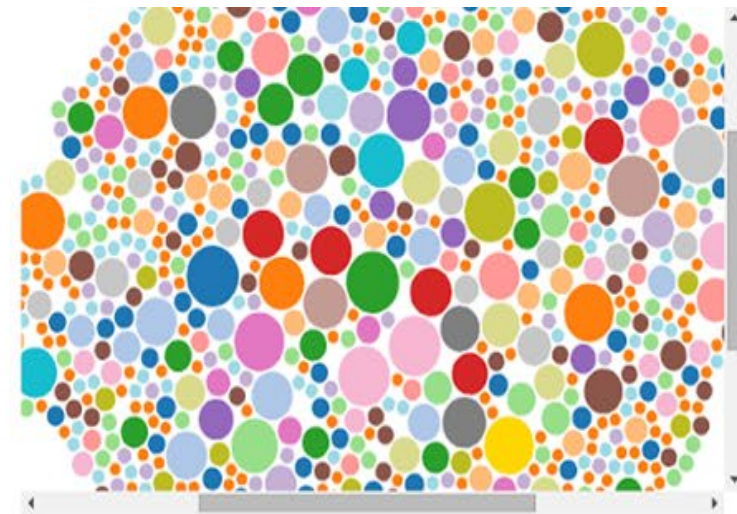
- Unary and n-ary INDs:  $R[A] \subseteq S[B]$  and  $R[ABC] \subseteq S[DEF]$
- Detect unknown foreign keys
  
- Example: PDB – Protein Data Bank
  - OpenMMS provides relational schema, 175 tables
  - Not a single foreign key constraint!
- Example: Ensembl – genome database
  - Shipped as MySQL dump files: >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 foreign key constraints in DBMS
  - Fear of performance drop for constraint checking
  - Lack of database knowledge



# Candidate Set Growth for INDs



# MANY: INDs among millions of web tables



96242-1	96242-1.'Rotational_Equinox / Rotational Equinox by House Association'.csv
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_astrolgy'. 'Ruling planets of the astrological signs and houses'.csv
175797-1	175797-1.'Sun_sign_astrolgy'. 'Sun signs'.csv
177750-2	177750-2.'BBC_Radio_Manchester'. 'Programming'.csv
89462-4	89462-4.'Astrology_and_the_classical_elements'. 'Triplicities by season'.csv
213213-1	213213-1.'Dalton_Park'. 'Opening times'.csv
470402-1	470402-1.'Dalton_Park'. 'Opening times'.csv

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 s 0 d 10 h 38 m

Zoom (1-5)

Range (logarithmic)

Dataset

allFilters



## Other dependencies

- Detecting multi-valued dependencies (MVDs) and join dependencies
- Detecting denial constraints (DCs)
- Detecting order dependencies (ODs)

□ `SELECT emp_name  
FROM employees  
ORDER BY rank, salary`

□ `SELECT emp_name  
FROM employees  
ORDER BY rank`

Remove rank

Replace with  
salary (if index  
only on salary)

emp_name	rank	salary
Smith	1	40k
Johnson	1	40k
Williams	1	45k
Brown	2	60k
Davis	2	60k
Miller	3	70k
Wilson	4	100k

salary „orders“ rank

Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

## Partial dependencies

- Aka. “approximate dependencies”
- Do not perfectly hold
  - For all but 10 of the tuples
  - Only for 80% of the tuples
  - Only for 1% of the tuples
- Also: Approximate dependencies
- Conditional dependencies
  - Concise description for which data the partial dependency is valid
- Matching dependencies
- Metric dependencies

RFD abbrev.	RFD name
ACOD	Approximate comparable dependency
ADD	Approximate differential dependency
AFD	Approximate functional dependency
COD	Comparable dependency
CFD	Conditional functional dependency
CFD <sup>P</sup>	CFD with built-in predicates
CFD <sup>C</sup>	CFD with cardinality constraints and synonym rules
CMD	Conditional matching dependency
CSD	Conditional sequential dependency
CD	Constrained functional dependency
DD	Differential dependency
ecFD	Extended conditional functional dependency
FFD	Fuzzy functional dependency
MD	Matching dependency
MFD	Metric functional dependency
ND	Neighborhood dependency
NUD	Numerical dependency
OD	Order dependency
OD <sub>K</sub>	OD satisfied within bound $k$
OD <sub>EA</sub>	OD satisfied almost everywhere
OFD	Ordered functional dependency
PD	Partial determination
POD	Polarized order dependencies
preFD	Preference functional dependency
PAC	Probabilistic approximate constraint
pFD	Probabilistic functional dependency
PuD	Purity dependency
RUD	Roll-up dependency
SD	Sequential dependency
SFD	Similarity functional dependency
soft FD	Soft functional dependency
TD	Trend dependency
TMFD	Type-M functional dependency
XCFD	XML conditional functional dependency
$\sigma\theta$ XFD	XML FD with $\sigma$ and $\theta$ approximation

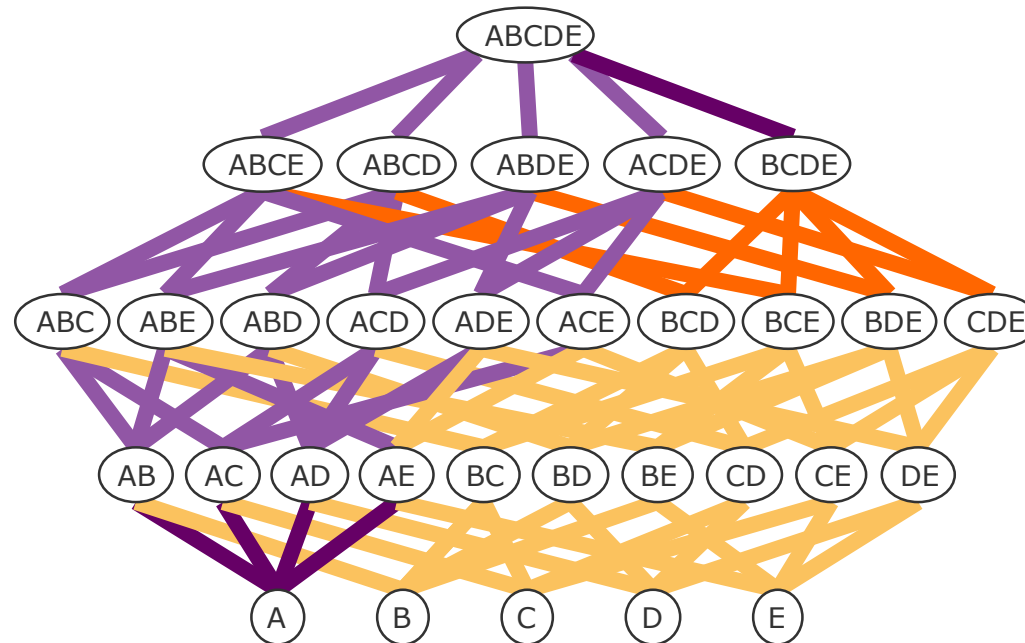
## Outlook: Profiling Challenges

---

- Efficient profiling
- Scalable profiling
- Holistic profiling
- Incremental profiling
- Online profiling
- Temporal profiling
- Profiling query results
- Profiling new types of data
- Data generation and testing
- Data profiling benchmark
- Hundreds of UCCs – which ones are keys?
- Thousands of FDs – which ones are true?
- Millions of INDs – which ones are foreign keys?
- User-driven interpretation:
  - Rank and visualize metadata
- Machine-driven interpretation
  - Machine learning

# Summary

1. Basic statistics
2. Uniques and keys
3. Functional dependencies
4. Inclusion dependencies and foreign keys
5. Outlook: Other dependencies and more



Felix Naumann  
Data Profiling  
Télécom ParisTech 2018

## References – work at HPI

- A Hybrid Approach for Efficient Unique Column Combination Discovery: Thorsten Papenbrock, Felix Naumann, BTW 2017
- Fast Approximate Discovery of Inclusion Dependencies, Sebastian Kruse, Thorsten Papenbrock, Christian Dullweber, Moritz Finke, Manuel Hegner, Martin Zabel, Christian Zöllner, Felix Naumann, BTW 2017
- Data-driven Schema Normalization, Thorsten Papenbrock, Felix Naumann, EDBT 2017
- Data Anamnesis: Admitting Raw Data into an Organization, Sebastian Kruse, Thorsten Papenbrock, Hazar Harmouch, Felix Naumann, IEEE Data Engineering Bulletin, 2016
- A Hybrid Approach to Functional Dependency Discovery, Thorsten Papenbrock, Felix Naumann, SIGMOD 2016
- Efficient Order Dependency Discovery, Philipp Langer and Felix Naumann, VLDB Journal 2016
- Holistic Data Profiling: Simultaneous Discovery of Various Metadata, Jens Ehrlich, Mandy Roick, Lukas Schulze, Jakob Zwiener, Thorsten Papenbrock, and Felix Naumann, EDBT 2016
- RDFind: Scalable Conditional Inclusion Dependency Discovery in RDF Datasets, Sebastian Kruse, Anja Jentzsch, Thorsten Papenbrock, Zoi Kaoudi, Jorge-Arnulfo Quiane-Ruiz, Felix Naumann, SIGMOD 2016
- Data Profiling (tutorial), Ziawasch Abedjan, Lukasz Golab and Felix Naumann, ICDE 2016
- Approximate Discovery of Functional Dependencies for Large Datasets, Tobias Bleifuß, Susanne Bülow, Johannes Frohnhofen, Julian Risch, Georg Wiese, Sebastian Kruse, Thorsten Papenbrock, Felix Naumann, CIKM 2016
- Divide & Conquer-based Inclusion Dependency Discovery, Thorsten Papenbrock, Sebastian Kruse, Jorge-Arnulfo Quiane-Ruiz, Felix Naumann, PVLDB 2015
- Functional Dependency Discovery: An Experimental Evaluation of Seven Algorithms, Thorsten Papenbrock, Jens Ehrlich, Jannik Marten, Tommy Neubert, Jan-Peer Rudolph, Martin Schönberg, Jakob Zwiener, Felix Naumann, PVLDB 2015
- Profiling relational data: a survey, Ziawasch Abedjan, Lukasz Golab, Felix Naumann, VLDB Journal 2015
- Scaling Out the Discovery of Inclusion Dependencies, Sebastian Kruse, Thorsten Papenbrock, Felix Naumann, BTW 2015
- Data Profiling with Metanome (demo), Thorsten Papenbrock, Tanja Bergmann, Moritz Finke, Jakob Zwiener, Felix Naumann, PVLDB 2015
- DFD: Efficient Discovery of Functional Dependencies, Ziawasch Abedjan, Patrick Schulze, Felix Naumann, CIKM 2014
- Detecting Unique Column Combinations on Dynamic Data, Ziawasch Abedjan, Jorge-Arnulfo Quiane-Ruiz, Felix Naumann, ICDE 2014
- Profiling and Mining RDF Data with ProLOD++, Ziawasch Abedjan, Toni Gruetze, Anja Jentzsch, Felix Naumann, ICDE Demo 2014
- LODOP - Multi-Query Optimization for Linked Data Profiling Queries., Benedikt Forchhammer, Anja Jentzsch, Felix Naumann, PROFILES 2014
- Scalable Discovery of Unique Column Combinations, Arvid Heise, Jorge-Arnulfo Quiane-Ruiz, Ziawasch Abedjan, Anja Jentzsch, Felix Naumann, PVLDB 2013
- Data Profiling Revisited, Felix Naumann, SIGMOD Record 2013
- Discovering Conditional Inclusion Dependencies. Jana Bauckmann, Ziawasch Abedjan, Heiko Müller, Ulf Leser, Felix Naumann, CIKM 2012
- Advancing the Discovery of Unique Column Combinations, Ziawasch Abedjan, Felix Naumann, CIKM 2011
- A Machine Learning Approach to Foreign Key Discovery, Alexandra Rostin, Oliver Albrecht, Jana Bauckmann, Felix Naumann, Ulf Leser, WebDB 2009
- Efficiently Detecting Inclusion Dependencies, Jana Bauckmann, Ulf Leser, Felix Naumann, Veronique Tietz, ICDE 2007
- Efficiently Computing Inclusion Dependencies for Schema Discovery, Jana Bauckmann, Ulf Leser, Felix Naumann, ICDE 2006