



Data Science – Opportunities and Challenges
United Nations Headquarters

March 25, 2019
Felix Naumann

The Hasso Plattner Institute in Potsdam, Germany



Felix Naumann
Data Science 2019

Information Systems Team



Dr. Thorsten Papenbrock



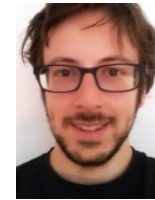
Diana Stephan



Prof. Felix Naumann



Dr. Ralf Krestel



Leon Bornemann



Hazar Harmouch



Konstantina Lazaridou



Tim Repke



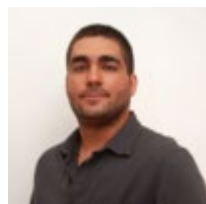
Felix Naumann
Data Science 2019



Julian Risch



Michael Loster



John Koumarelas



Tobias Bleifuß



Nitisha Jain



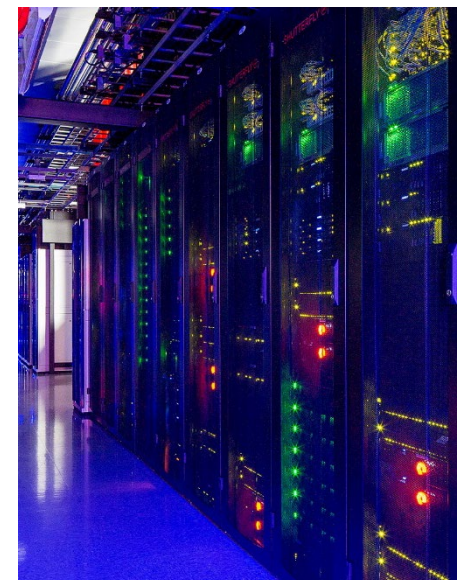
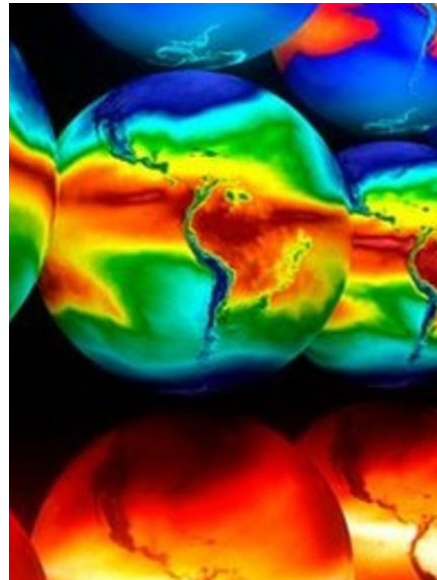
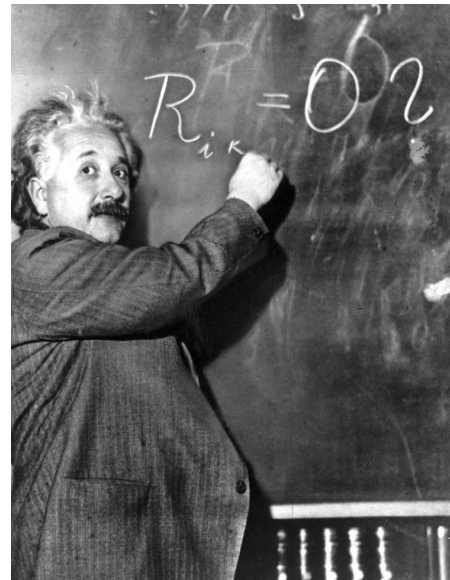
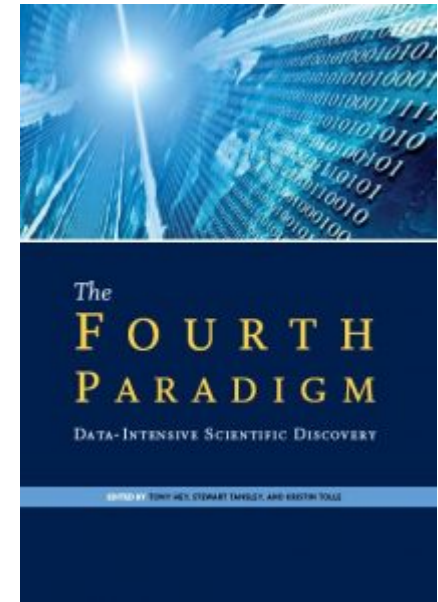
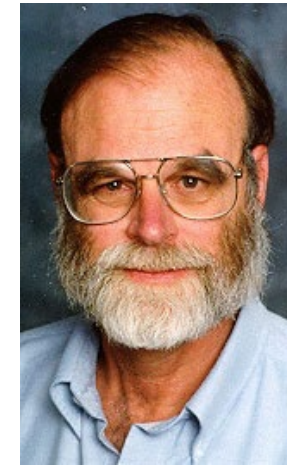
Lan Jiang

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

The Fourth Paradigm of Science

1. Empirical and experimental
2. Theoretical
3. Computational
4. Data-intensive

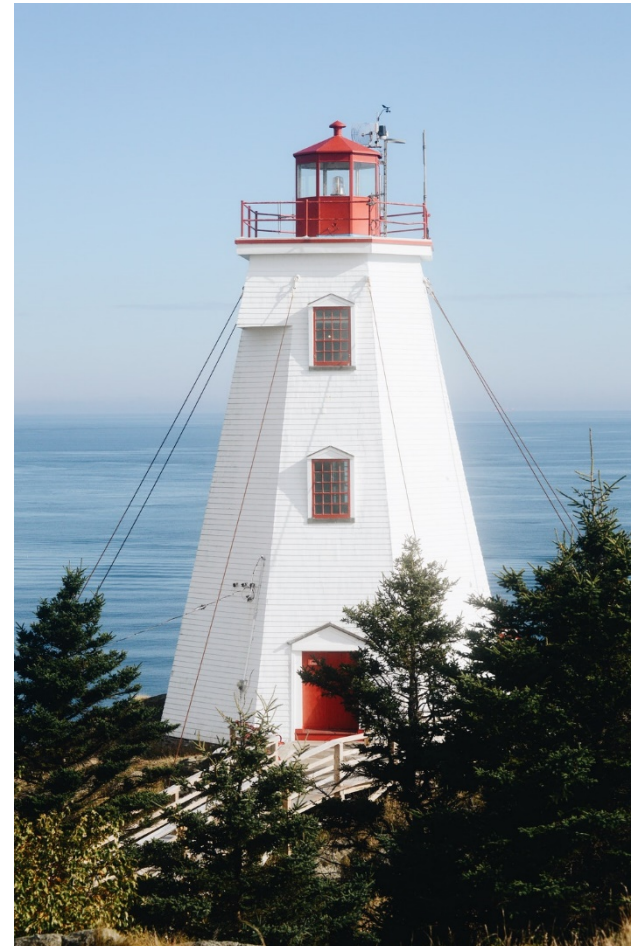
We have to do better producing tools to support the whole research cycle - from data capture and data curation to data analysis and data visualization. Jim Gray



Felix Naumann
Data Science 2019

Overview

- 1. Data Science**
2. Big Data
3. Data Profiling
4. Data Preparation
5. Data Cleaning



Felix Naumann
Data Science 2019

<https://unsplash.com/photos/vGefUiWm0xI>

More and more data are available to science and business

■ Drivers

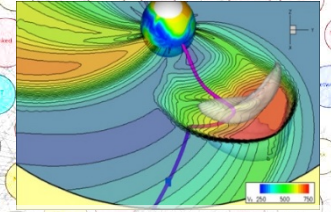
- Cloud computing
- Internet of services
- Internet of things
- Cyberphysical systems

■ Underlying trends

- Connectivity
- Collaboration
- Computer generated data



sensor data



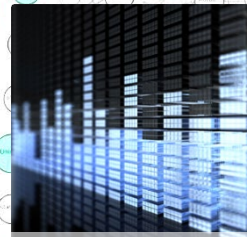
simulation data



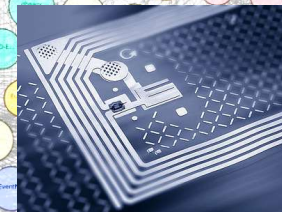
video streams



web archives



audio streams



RFID data



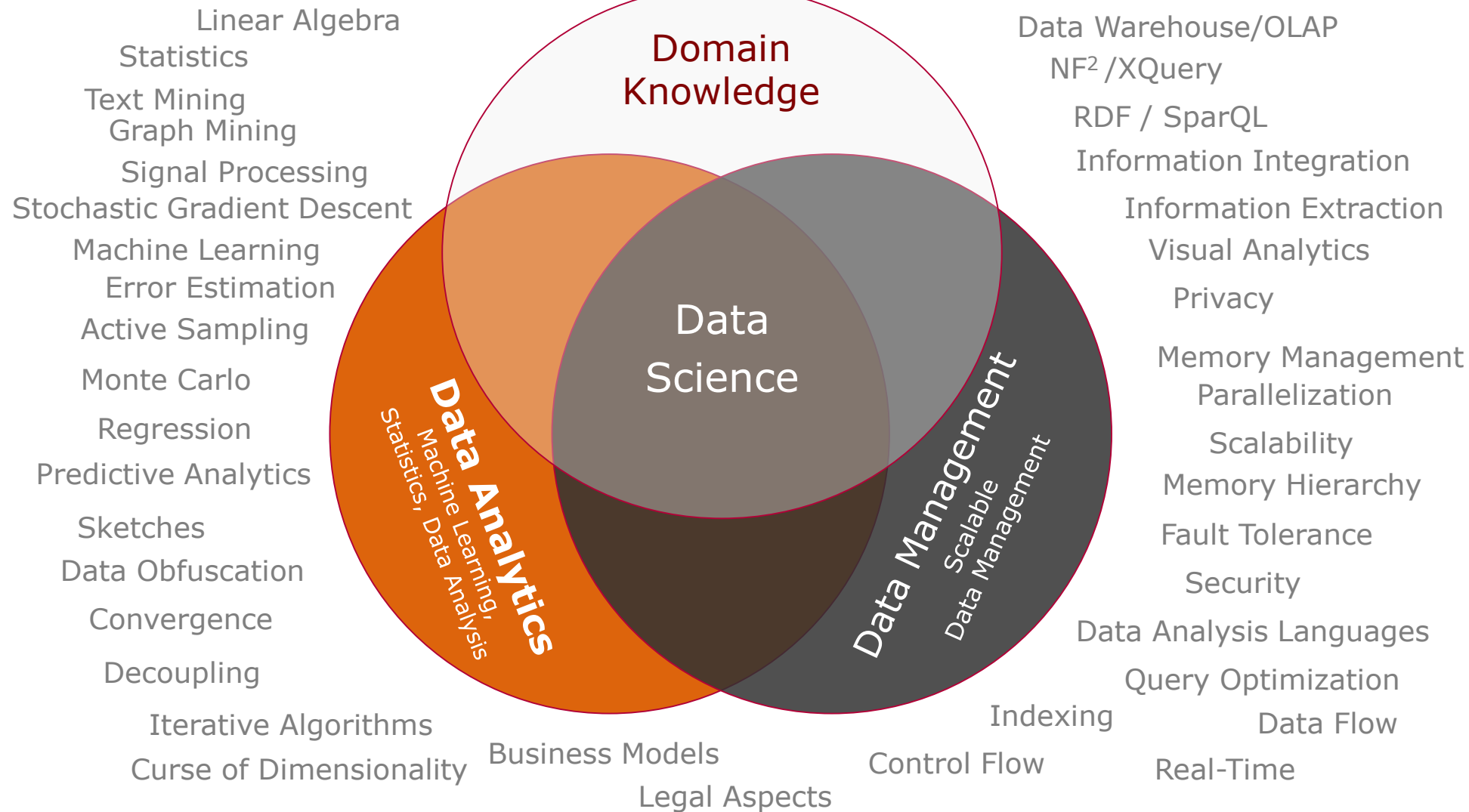
Government data

"Data Scientist" – "Jack of All Trades!"

Domain Expertise (e.g., Industry 4.0, Medicine, Physics, Engineering, Energy, Logistics)

Mathematical Programming

Relational Algebra / SQL



Felix Naumann
Data Science 2019

Positive Uses of Big Data

- **Prediction**
 - Weather, natural disaster, predictive maintenance, disease
- **Optimization**
 - Planning, traffic, logistics, machine efficiency, site selection
- **Individualization**
 - Digital health and personalized medicine, personalized learning, recommendations
- **Comfort**
 - Sharing, smart home, authentication (face, gait)
 - Happiness: HappyDB – a database of happy moments
 - Autonomous vehicles
- **Intelligence**
 - Fraud detection, translation, gaming
 - Robotics

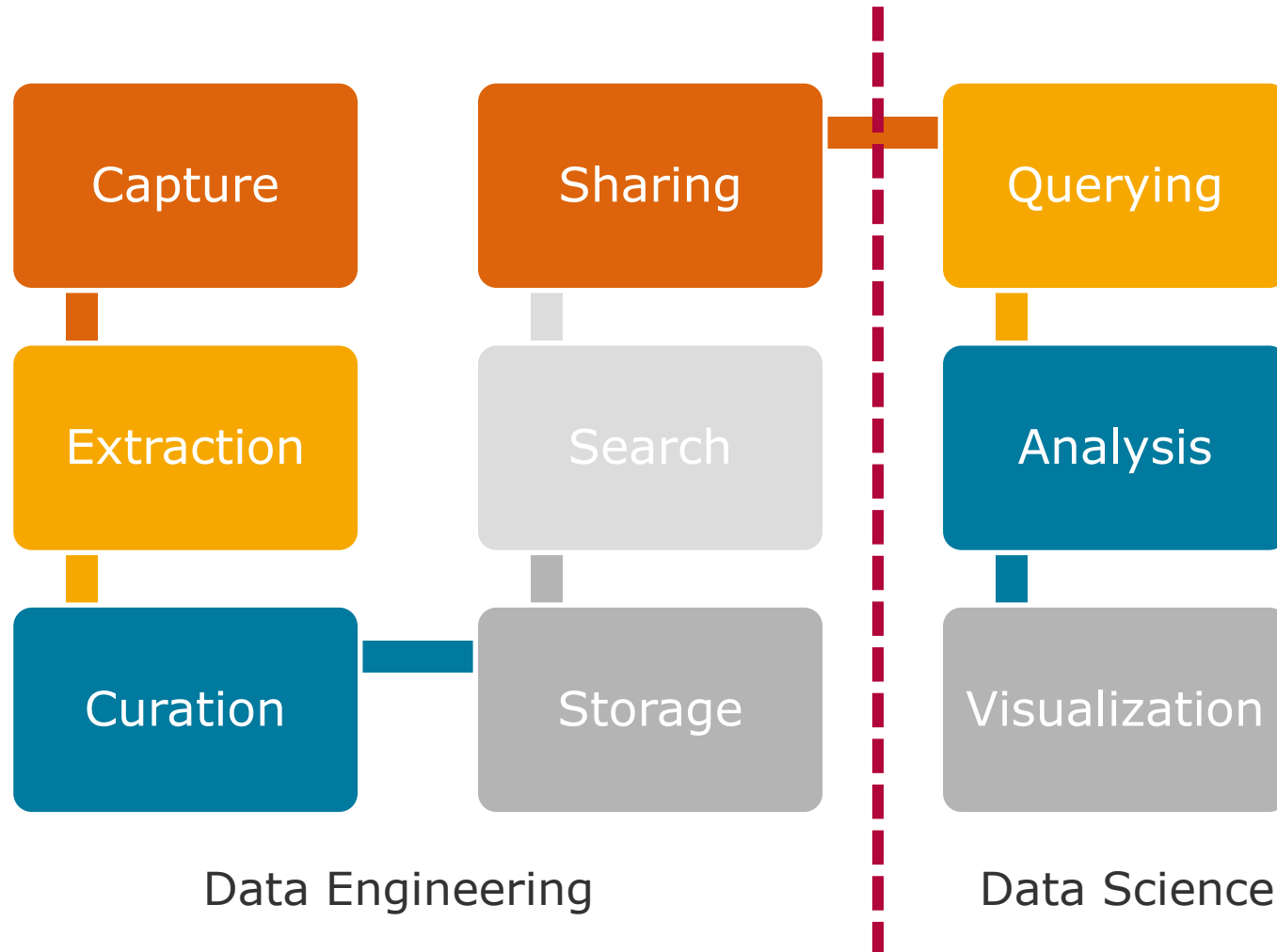


Questionable Uses of Big Data

- **Invading lives**
 - Tracking persons:
 - Direct: GPS location tracking
 - Indirect: face recognition / surveillance
 - Tracking behavior: social networks, sensors, smart homes
- **Classifying individuals**
 - Behavior prediction
 - Crime prediction
 - Social Scoring
- **Misinformation**
 - Filter bubble
 - Manipulating/inflaming opinion
- **Intervention**
 - Restricting free movement
 - Censorship
 - Autonomous drones



Data Science Pipeline



Felix Naumann
Data Science 2019

Overview

1. Data Science
- 2. Big Data**
3. Data Profiling
4. Data Preparation
5. Data Cleaning



Felix Naumann
Data Science 2019

<https://unsplash.com/photos/vGefUiWm0xI>

Gartner's 3 (+ 1) V's – Properties of Big Data

Volume

- Size of dataset

Velocity

- Speed at which data arrives and must be processed

Variety

- Different data modalities, models, schemata, semantics

Veracity

- Data quality: Correctness, completeness, consistency, up-to-dateness, etc.

Viscosity

- Integration and dataflow friction

Venue

- Different locations that require different access & extraction methods

Vocabulary

- Different language and vocabulary

Value

- Added-value of data to organization and use-case

Virality

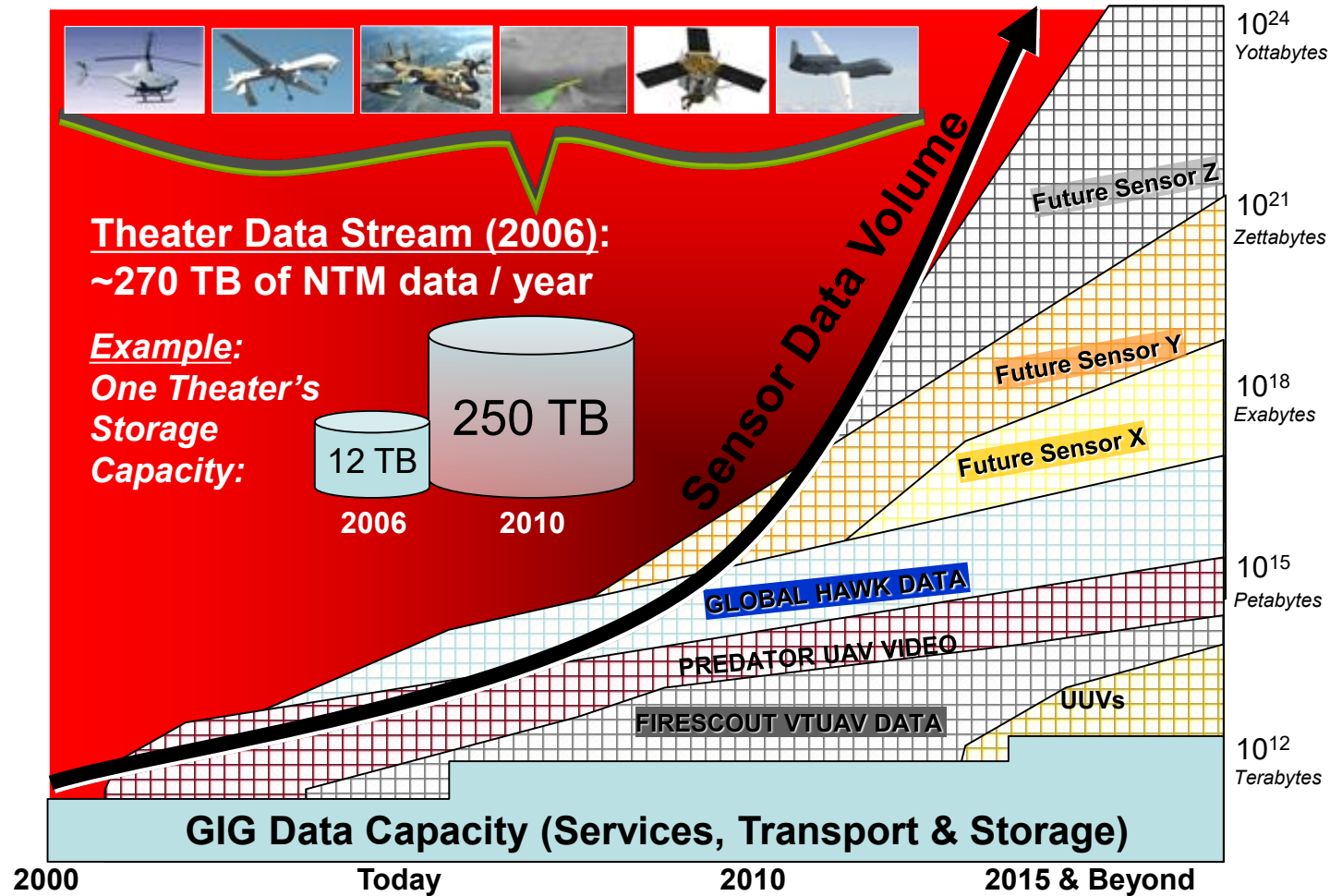
- Speed of dispersal among community

Variability

- Data, formats, schema, semantics change

Volatility, vagueness, validity, visualization, ...

Military Projection of Sensor Data Volume (later refuted)



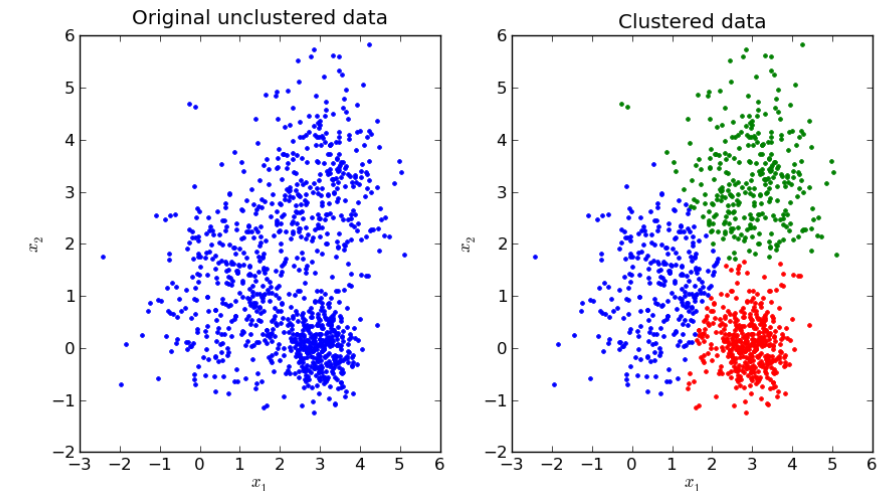
153 hard disks per person on the planet

Felix Naumann
Data Science 2019

Using 1TB drives, this would require 1 trillion (10¹²) drives!

Abridged History of Big Data Analytics

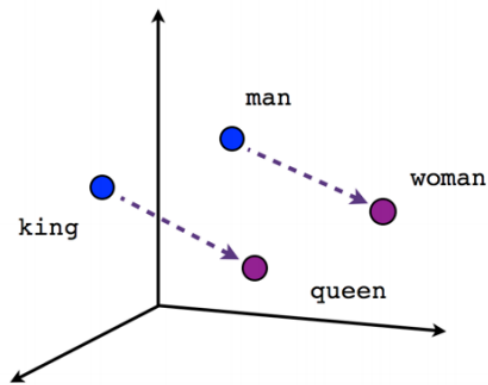
- **Aggregation**: Calculate statistics
 - Sum of sales, average cheese consumption (per state)
- **Data mining**: Identify useful rules
 - 35% of all customers who bought X, also bought Y (X=beer and Y=diapers)
- **Clustering**: Group similar items
 - Cluster patients into 10 groups based on a similarity measure (age, weight, income ...)
- **Classification**: Organize items into a set of known groups based on similarity
 - Assort products into categories
 - Collaborative filtering (for movies)
- **Machine learning**: Generalization of all of the above
 - Build a model that explains the data (for a given target dimension)
 - Apply the model to new data items to find out target dimension value



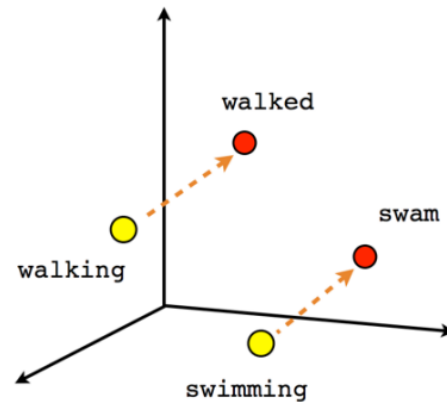
Felix Naumann
Data Science 2019

Appetite for Training Data

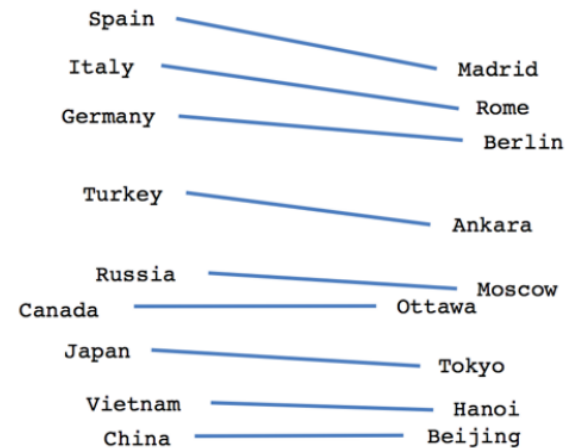
- Sophisticated models need many input dimensions
 - Few dimensions for spam filtering →
 - Tens of dimensions for intrusion detection →
 - Hundreds of dimensions for user classification →
 - Thousands of dimensions to understand text →
- ... and have many model parameters.
- Need at least as many input data items as parameters
 - Labeled spam emails
 - Annotated log entries
 - Detailed user profiles
 - Sample texts



Male-Female



Verb tense



Country-Capital

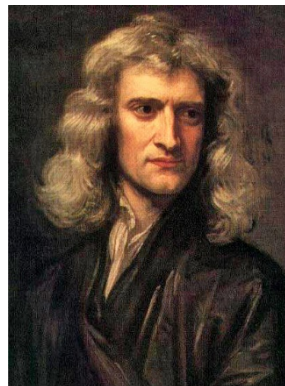
Felix Naumann
Data Science 2019

Big Data = Science?

- The End of Theory: The Data Deluge Makes the Scientific Method Obsolete (Chris Anderson, Wired, 2008)
 - All models are wrong, but some are useful. (George Box)
 - All models are wrong, and increasingly you can succeed without them. (Peter Norvig)
- Before Big Data: Correlation is not causation!
- With Big Data: Who cares?
 - Traditional approach to science — hypothesize, model, test — is becoming obsolete.
 - Petabytes allow us to say: **“Correlation is enough.”**

http://www.wired.com/science/discoveries/magazine/16-07/pb_theory

$$F = G \frac{m_1 m_2}{r^2}$$

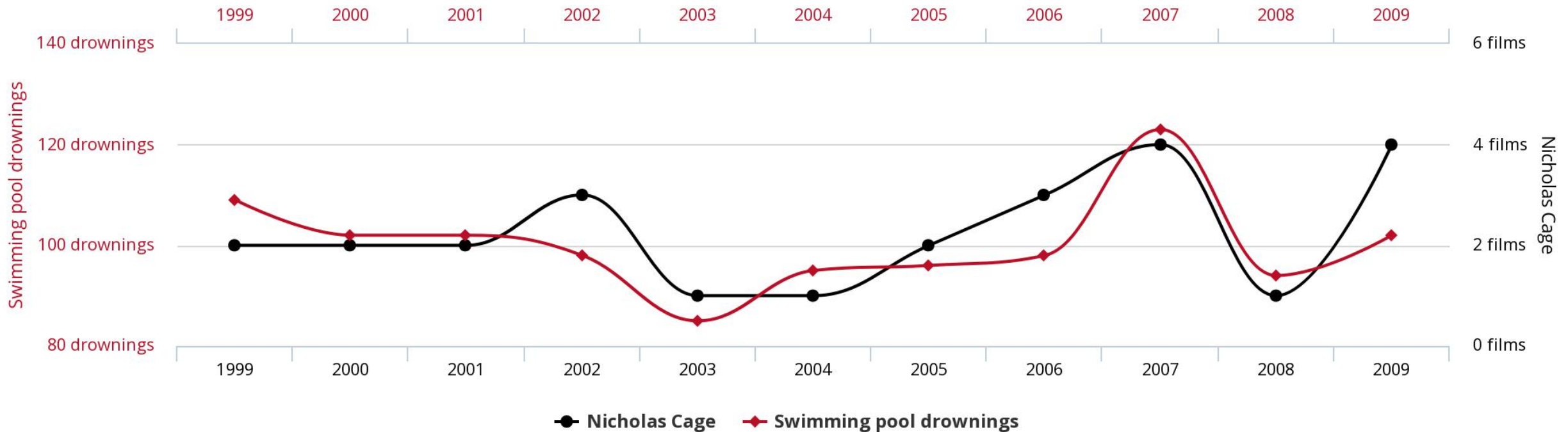


vs.



Correlation vs. Causation

Number of people who drowned by falling into a pool
 correlates with
Films Nicolas Cage appeared in



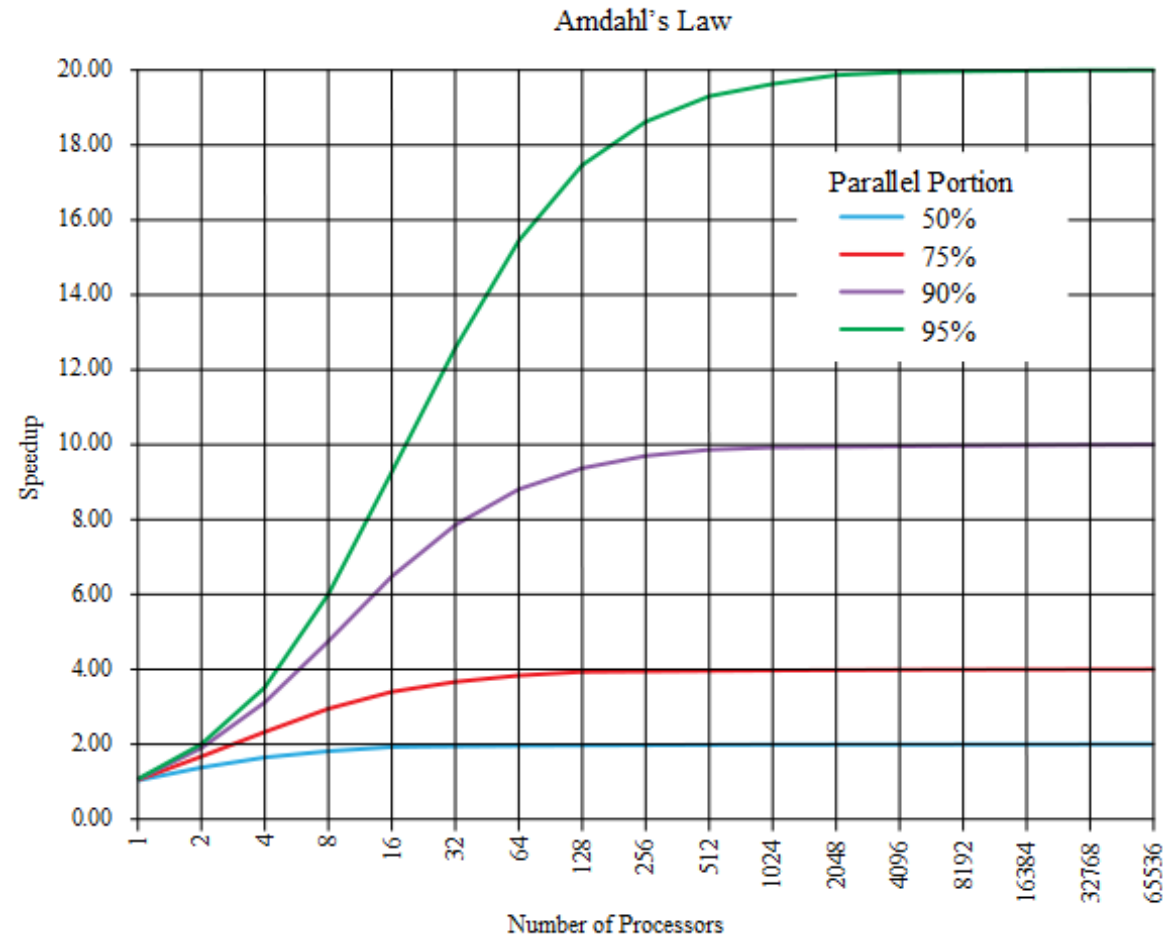
Parallelization obstacle: Amdahl's Law

■ Maximal speedup is determined by non-parallelizable part of program:

□ $S_{max} = \frac{1}{(1-f)+f/p}$

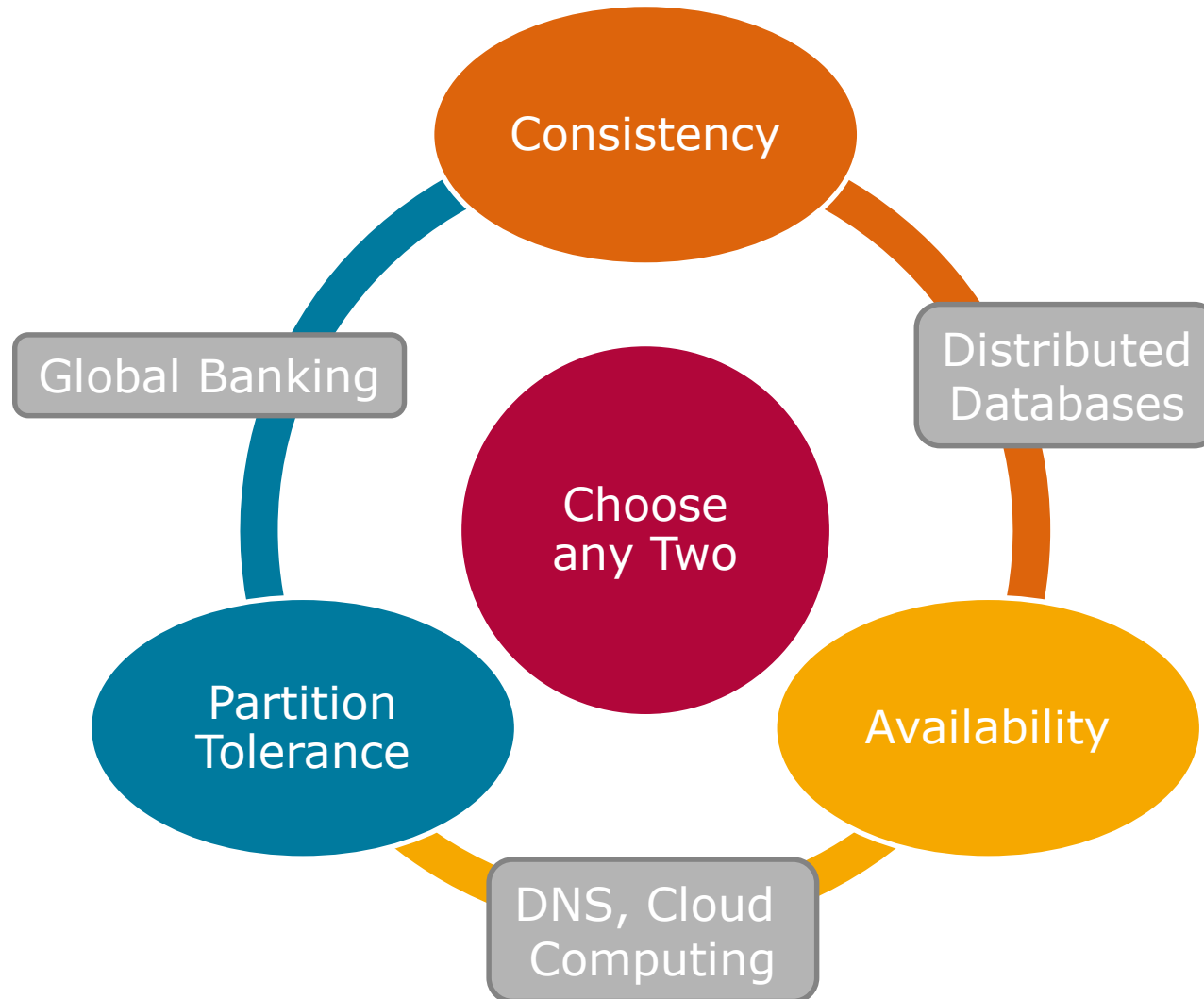
□ p processors

□ f parallelizable fraction



Felix Naumann
Data Science 2019

Distribution Obstacle: CAP Theorem



Overview

1. Data Science
2. Big Data
- 3. Data Profiling**
4. Data Preparation
5. Data Cleaning



Felix Naumann
Data Science 2019

<https://unsplash.com/photos/vGefUiWm0xI>

Definition Data Profiling

- Data profiling refers to the activity of creating small but informative summaries of a database.

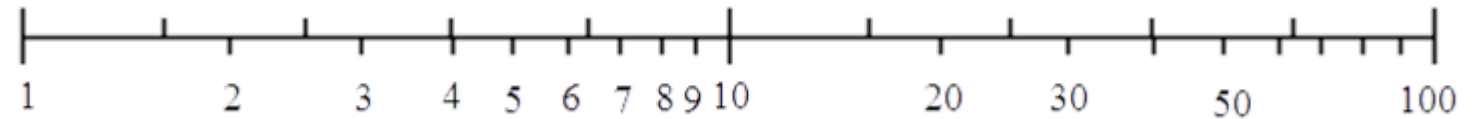
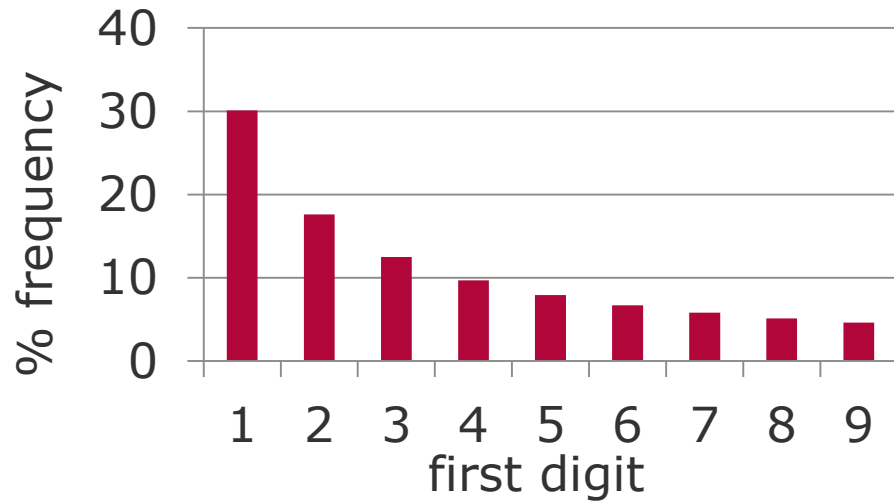
Ted Johnson, Encyclopedia of Database Systems

- Extracting metadata from given data
 - Basic statistics and histograms
 - Datatypes
 - Key and foreign keys
 - Dependencies and rules
- Data profiling is first step in any data management task
 - *“What shape does my data have?”*

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

- Statement about the distribution of first digits d in (many) *naturally occurring* numbers:

□ $P(d) = \log_{10}(d + 1) - \log_{10}(d) = \log_{10}(1 + 1/d)$

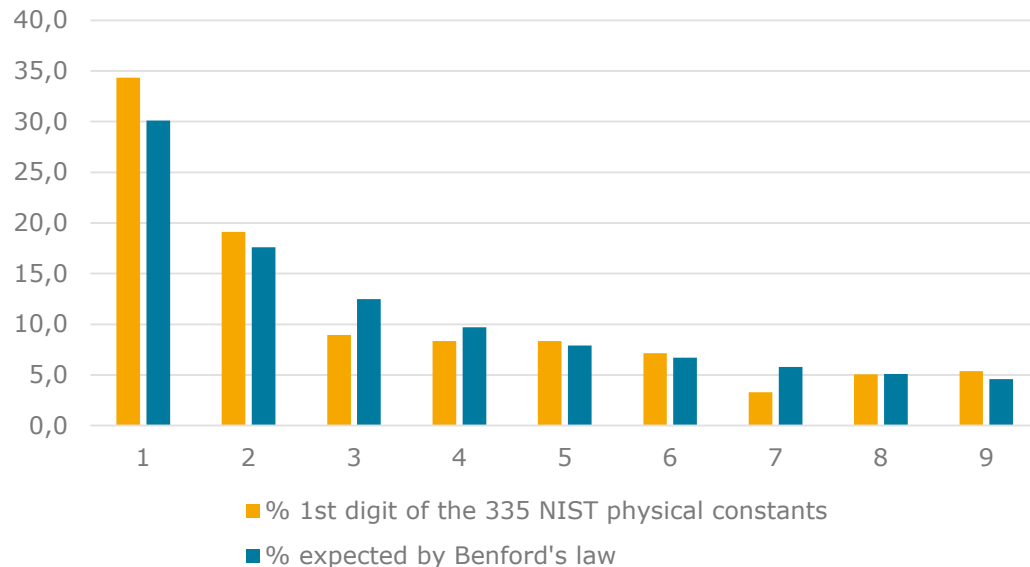


Is true if $\log(x)$ is uniformly distributed

Felix Naumann
Data Science 2019

Examples for Benford's Law

- Surface areas of 335 rivers
- Sizes of 3259 US populations
- 1800 molecular weights
- 5000 entries from a mathematical handbook
- 308 numbers in an issue of Reader's Digest
- Street addresses of the first 342 persons listed in American Men of Science
- Powers of 2: 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



Felix Naumann
Data Science 2019

CODATA RECOMMENDED VALUES OF THE FUNDAMENTAL PHYSICAL CONSTANTS: 2014

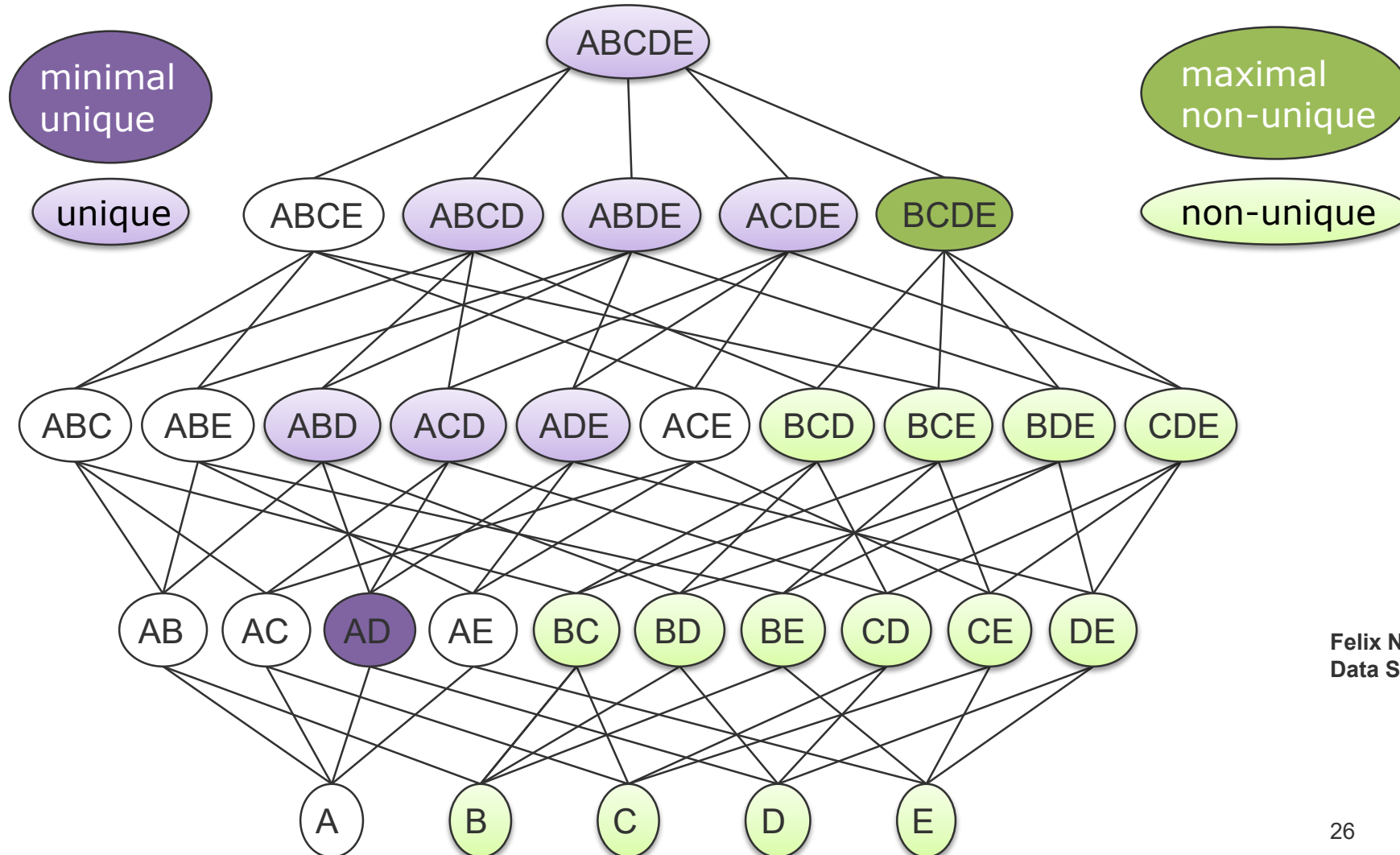
NIST SP 961 (Sept/2015) Values from: P. J. Mohr, D. B. Newell, and B. N. Taylor, arXiv:1507.07956

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

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)	m s^{-1}	muon g -factor $-2(1 + a_\mu)$	g_μ	$-2.002\,331\,8418(13)$	
magnetic constant	μ_0	$4\pi \times 10^{-7}$ (exact)	N A^{-2}	muon-proton magnetic moment ratio	μ_μ/μ_p	$-3.183\,345\,142(71)$	
electric constant $1/\mu_0 c^2$	ϵ_0	$8.854\,187\,817\dots \times 10^{-12}$	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	μ_p	$1.410\,606\,7873(97) \times 10^{-26}$	J T ⁻¹
in eV s		$6.582\,119\,514(40) \times 10^{-16}$	eV s	to nuclear magneton ratio	μ_p/μ_N	$2.792\,847\,3508(85)$	
elementary charge	e	$1.602\,176\,6208(98) \times 10^{-19}$	C	proton magnetic shielding correction $1 - \mu'_p/\mu_p$	σ'_p	$25.691(11) \times 10^{-6}$	
magnetic flux quantum $h/2e$	Φ_0	$2.067\,833\,831(13) \times 10^{-15}$	Wb	(H ₂ O, sphere, 25 °C)			
Josephson constant $2e/h$	K_J	$483\,597.8525(30) \times 10^9$	Hz V ⁻¹	proton gyromagnetic ratio $2\mu_p/\hbar$	γ_p	$2.675\,221\,900(18) \times 10^8$	s ⁻¹ T ⁻¹
von Klitzing constant $h/e^2 = \mu_0 c/2\alpha$	R_K	$25\,812.807\,4555(59)$	Ω		$\gamma_p/2\pi$	$42.577\,478\,92(29)$	MHz T ⁻¹
Bohr magneton $e\hbar/2m_e$	μ_B	$927.400\,9994(57) \times 10^{-26}$	J T ⁻¹	shielded proton gyromagnetic ratio $2\mu'_p/\hbar$	γ'_p	$2.675\,153\,171(33) \times 10^8$	s ⁻¹ T ⁻¹
in eV T ⁻¹		$5.788\,381\,8012(26) \times 10^{-5}$	eV T ⁻¹	(H ₂ O, sphere, 25 °C)			
nuclear magneton $e\hbar/2m_p$	μ_N	$5.050\,783\,699(31) \times 10^{-27}$	J T ⁻¹		$\gamma'_p/2\pi$	$42.576\,385\,07(53)$	MHz T ⁻¹
in eV T ⁻¹		$3.152\,451\,2550(15) \times 10^{-8}$	eV T ⁻¹	neutron mass in u	m_n	$1.008\,664\,915\,88(49)$	u
fine-structure constant $e^2/4\pi\epsilon_0\hbar c$	α	$7.297\,352\,5664(17) \times 10^{-3}$		energy equivalent in MeV	$m_n c^2$	$939.565\,4133(58)$	MeV
inverse fine-structure constant	α^{-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_∞	$10\,973\,731.568\,508(65)$	m ⁻¹	neutron magnetic moment	μ_n	$-0.966\,236\,50(23) \times 10^{-26}$	J T ⁻¹
energy equivalent in eV	$R_\infty c$	$3.289\,841\,960\,355(19) \times 10^{15}$	Hz	to nuclear magneton ratio	μ_n/μ_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	μ_d	$0.433\,073\,5040(36) \times 10^{-26}$	J T ⁻¹
in u		$5.485\,799\,090\,70(16) \times 10^{-4}$	u	to nuclear magneton ratio	μ_d/μ_N	$0.857\,438\,2311(48)$	
energy equivalent in MeV	$m_e c^2$	$0.510\,998\,9461(31)$	MeV	helion (³ He nucleus) mass in u	m_h	$3.014\,932\,246\,73(12)$	u
electron-muon mass ratio	m_e/m_μ	$4.836\,331\,70(11) \times 10^{-3}$		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	μ'_h	$-1.074\,553\,080(14) \times 10^{-26}$	J T ⁻¹
electron charge to mass quotient	$-e/m_e$	$-1.758\,820\,024(11) \times 10^{11}$	C kg ⁻¹	(gas, sphere, 25 °C)			
Compton wavelength $h/m_e c$	λ_C	$2.426\,310\,2367(11) \times 10^{-12}$	m	to Bohr magneton ratio	μ'_h/μ_B	$-1.158\,671\,471(14) \times 10^{-3}$	
$\lambda_C/2\pi = \alpha a_0 = \alpha^2/4\pi R_\infty$	λ_C	$386.159\,267\,64(18) \times 10^{-15}$	m	to nuclear magneton ratio	μ'_h/μ_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_α	$4.001\,506\,179\,127(63)$	u
Thomson cross section $(8\pi/3)r_e^2$	σ_e	$0.665\,245\,871\,58(91) \times 10^{-28}$	m ²	energy equivalent in MeV	$m_\alpha c^2$	$3727.379\,378(23)$	MeV
electron magnetic moment	μ_e	$-928.476\,4620(57) \times 10^{-26}$	J T ⁻¹	Avogadro constant	N_A, L	$6.022\,140\,857(74) \times 10^{23}$	mol ⁻¹
to Bohr magneton ratio	μ_e/μ_B	$-1.001\,159\,652\,180\,91(26)$		atomic mass constant $\frac{1}{12} m(^{12}\text{C}) = 1$ u	m_u	$1.660\,539\,040(20) \times 10^{-27}$	kg
to nuclear magneton ratio	μ_e/μ_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 ⁻¹
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 ⁻¹ K ⁻¹
electron-proton magnetic moment ratio	μ_e/μ_p	$-658.210\,6866(20)$		Boltzmann constant R/N_A	k	$1.380\,648\,52(79) \times 10^{-23}$	J K ⁻¹
muon mass in u	m_μ	$0.113\,428\,9257(25)$	u	in eV K ⁻¹		$8.617\,3303(50) \times 10^{-5}$	eV K ⁻¹
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 ³ mol ⁻¹
muon-electron mass ratio	m_μ/m_e	$206.768\,2826(46)$		($T = 273.15$ K, $p = 101.325$ kPa)			
muon magnetic moment	μ_μ	$-4.490\,448\,26(10) \times 10^{-26}$	J T ⁻¹	Stefan-Boltzmann constant $\pi^2 k^4/60\hbar^3 c^2$	σ	$5.670\,367(13) \times 10^{-8}$	W m ⁻² K ⁻⁴
to Bohr magneton ratio	μ_μ/μ_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 ²
to nuclear magneton ratio	μ_μ/μ_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_μ	$1.165\,920\,89(63) \times 10^{-3}$		$b = \lambda_{\text{max}} T = c_2/4.965\,114\,231\dots$	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
Energy equivalents							
$(1 \text{ m}^{-1})c = 299\,792\,458$ Hz	$(1 \text{ Hz})h/k = 4.799\,2447(28) \times 10^{-11}$ K	$(1 \text{ J}) = 6.241\,509\,126(38) \times 10^{18}$ eV	$(1 \text{ eV})/c^2 = 1.073\,544\,1105(66) \times 10^{-9}$ u				
$(1 \text{ m}^{-1})hc/k = 1.438\,777\,36(83) \times 10^{-2}$ K	$(1 \text{ Hz})h = 4.135\,667\,662(25) \times 10^{-15}$ eV	$(1 \text{ eV}) = 1.602\,176\,6208(98) \times 10^{-19}$ J	$(1 \text{ kg}) = 6.022\,140\,857(74) \times 10^{26}$ u				
$(1 \text{ m}^{-1})hc = 1.239\,841\,9739(76) \times 10^{-6}$ eV	$(1 \text{ K})k/hc = 69.503\,457(40) \text{ m}^{-1}$	$(1 \text{ eV})/hc = 8.065\,544\,005(50) \times 10^5 \text{ m}^{-1}$	$(1 \text{ u}) = 1.660\,539\,040(20) \times 10^{-27}$ kg				
$(1 \text{ m}^{-1})h/c = 1.331\,025\,049\,00(61) \times 10^{-15}$ u	$(1 \text{ K})k/h = 2.083\,6612(12) \times 10^{10}$ Hz	$(1 \text{ eV})/h = 2.417\,989\,262(15) \times 10^{14}$ Hz	$(1 \text{ u})c/h = 7.513\,006\,6166(34) \times 10^{14} \text{ m}^{-1}$				
$(1 \text{ Hz})/c = 3.335\,640\,951\dots \times 10^{-9} \text{ m}^{-1}$	$(1 \text{ K})k = 8.617\,3303(50) \times 10^{-5}$ eV	$(1 \text{ eV})/k = 1.160\,452\,21(67) \times 10^4$ K	$(1 \text{ u})c^2 = 931.494\,0954(57) \times 10^6$ eV				

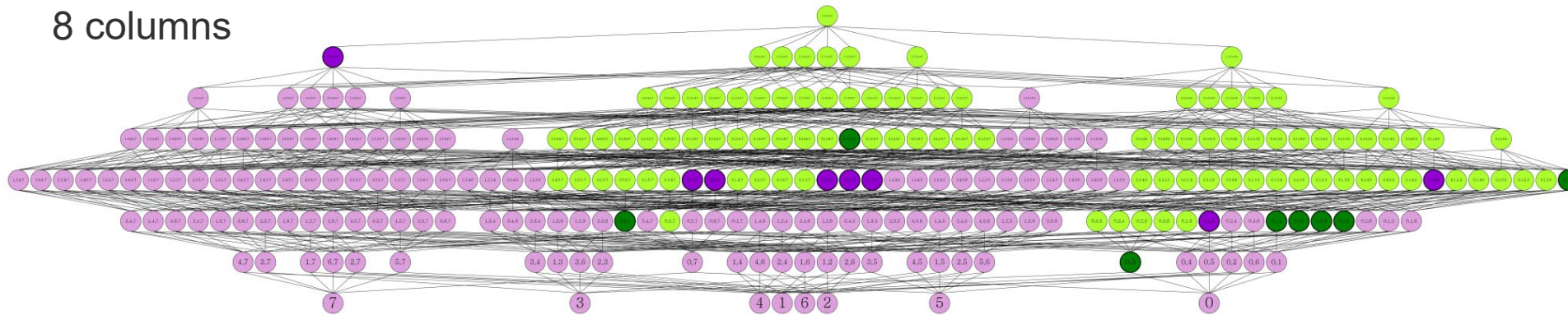
Detecting Unique Column Combinations (aka. keys)



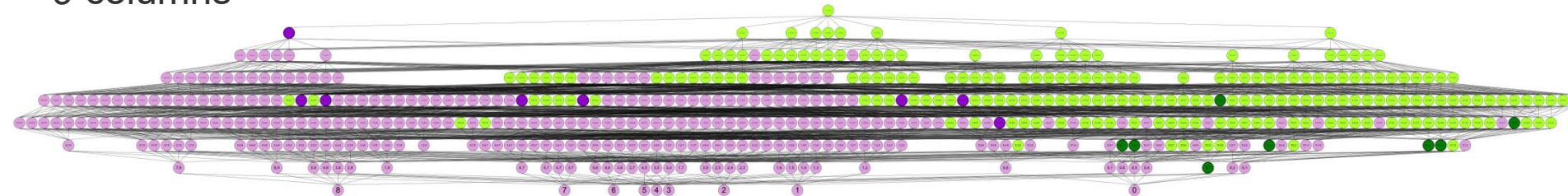
Felix Naumann
Data Science 2019

Large search space: $2^n - 1$
Large solution space: $\binom{n}{n/2}$

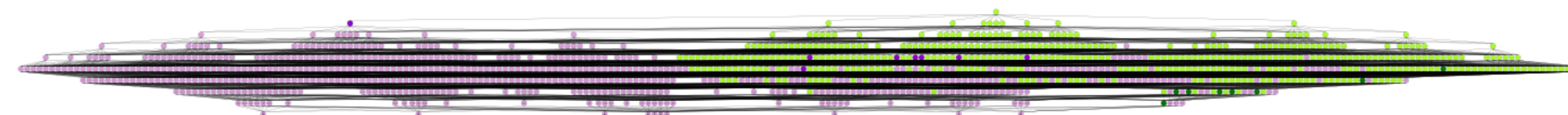
8 columns



9 columns

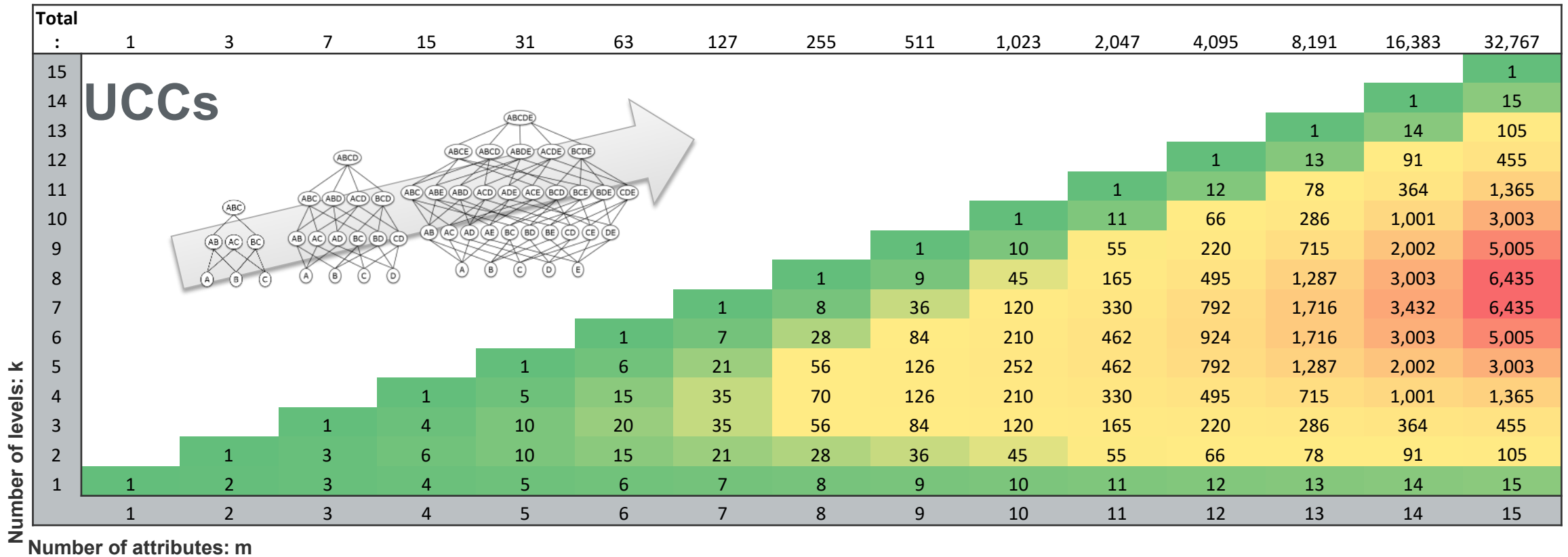


10 columns



Felix Naumann
Data Science 2019

Candidate Set Growth for Unique Column Combinations



Functional Dependencies



Game of Dependencies

Spoiler alert for Season 1

Felix Naumann
Data Science 2019

Functional Dependencies

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

Some Functional Dependencies:

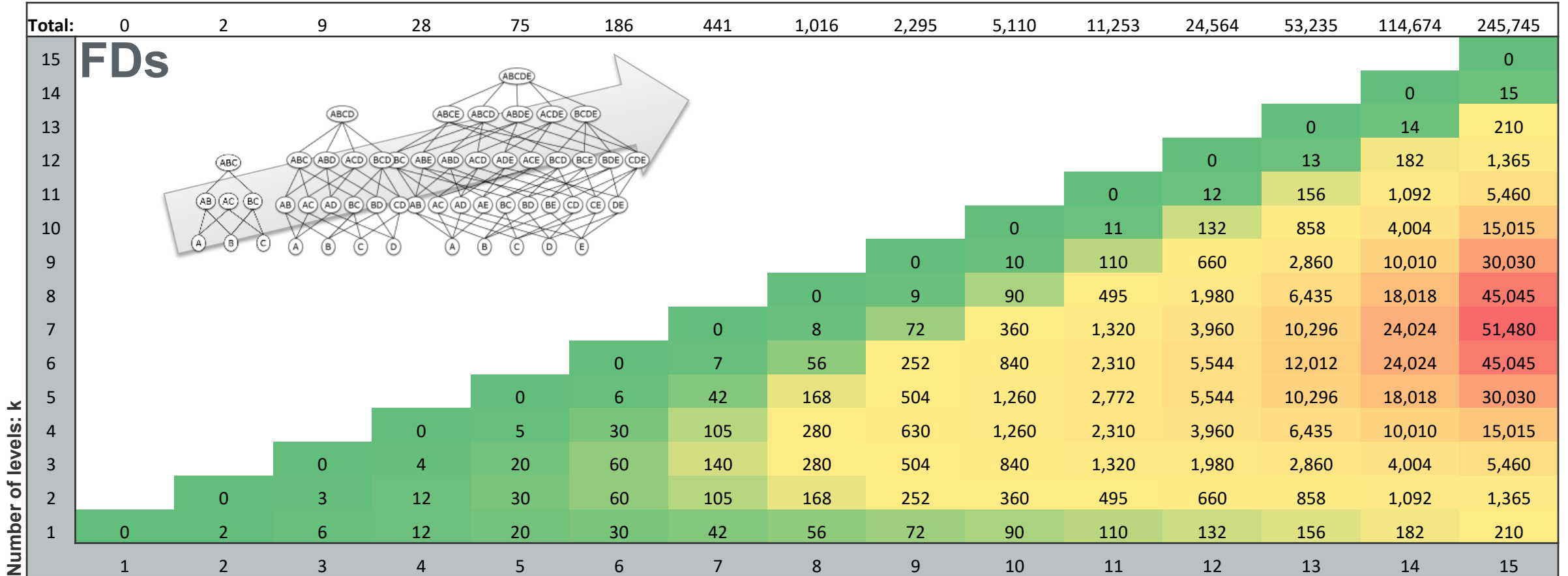
1. Person → Lineage
2. Person → Hair
3. Person → Religion
4. Lineage → Hair
5. Religion, Hair → Lineage
6. ...

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

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

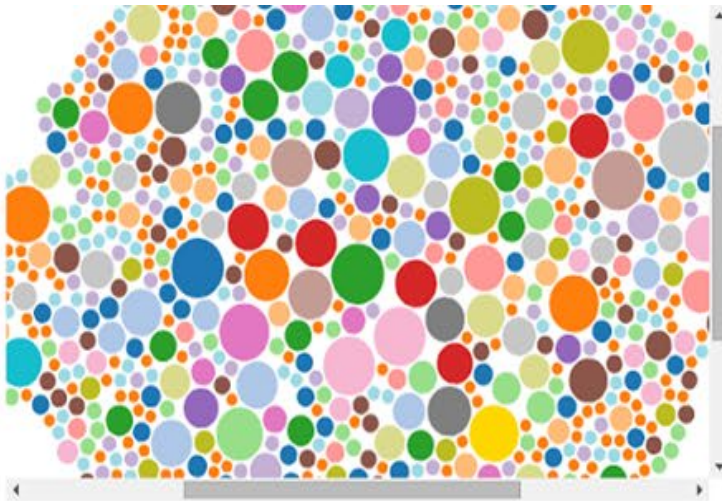
Felix Naumann
Data Science 2019

Candidate Set Growth for Functional Dependencies



Number of attributes: m

Linking up millions of web tables with inclusion dependencies



96242-1	96242-1.'Rotational_Engly / Rotational_Engly 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.'Astrolgy_and_the_classical_elements'. 'Triplicities 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

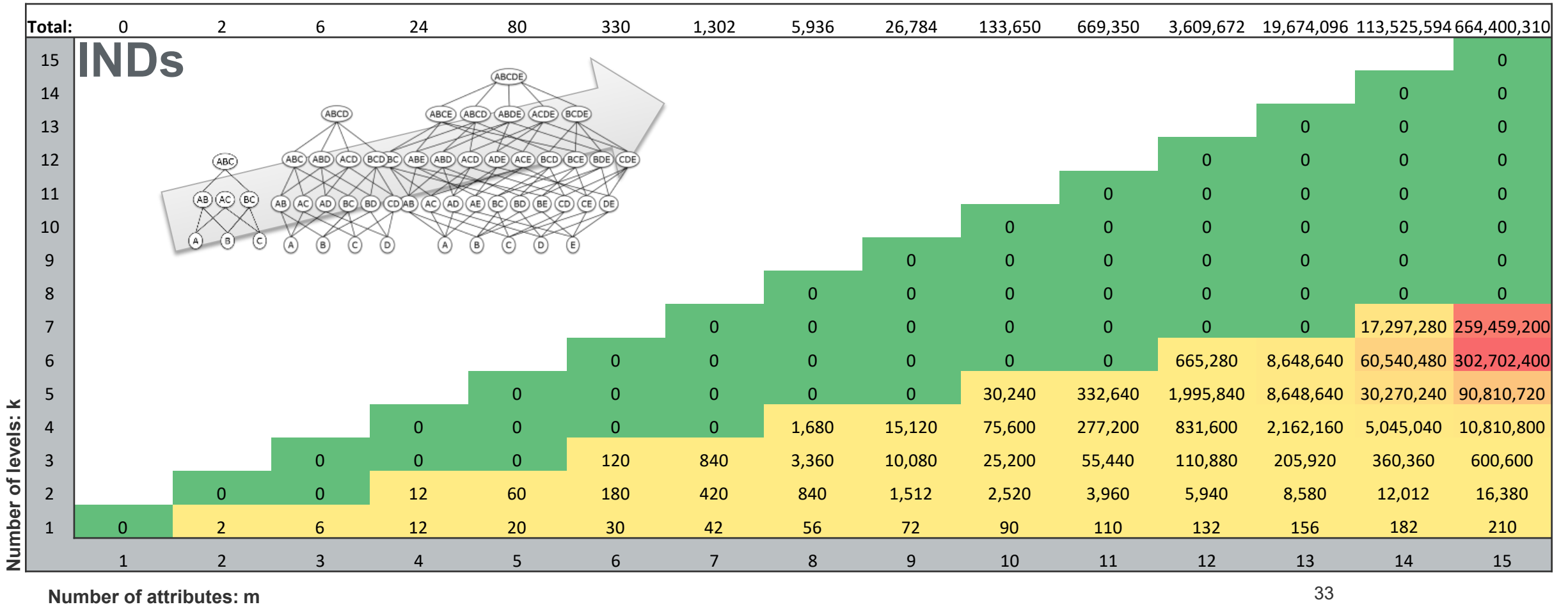
Zoom (1-5)

Range (logarithmic)

Dataset

allFilters

Candidate Set Growth for Inclusion Dependencies



Profiling Challenges

- Efficient profiling
 - Scalable profiling
 - Holistic profiling
 - Incremental profiling
 - Temporal profiling
 - Profiling query results
 - Profiling new types of data
-
- Hundreds of UCCs – which ones are keys?
 - Thousands of FDs – which ones are true?
 - Millions of INDs – which ones are foreign keys?

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**

Overview

1. Data Science
2. Big Data
3. Data Profiling
- 4. Data Preparation**
5. Data Cleaning



<https://unsplash.com/photos/vGefUiWm0xI>

Felix Naumann
Data Science 2019

36

Its late...

Data P-r_e\p+a|r_q|a.t/i~o-n

Title	Authors	Venue	Year
Immunogold labelling is a quantitative method as demonstrated by studies on aminopeptidase N in	GH Hansen, LL Wetterberg, H SjÃfÃ¶strÃfÃ¶m, O NorÃfÃ©n	The Histochemical Journal,	1992
Infectious Inmates and Releasees From Correctional Facilities	TM Hammett, P Harmon, W Rhodes	see	
World Population Prospects: The 1996 Revision	U Nations	New York,	
Consequences of Migration and Remittances for Mexican Transnational Communities.	D Conway, JH Cohen	Economic Geography,	1998

Wrong encoding

Incorrect title

Incorrect venue

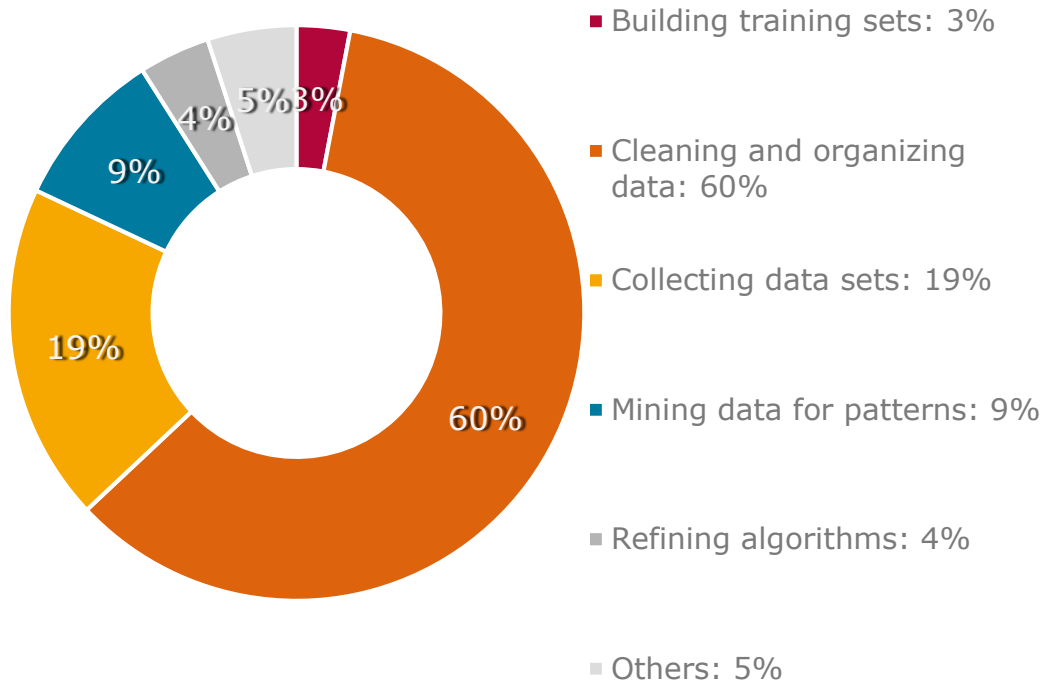
Missing values

redundant characters

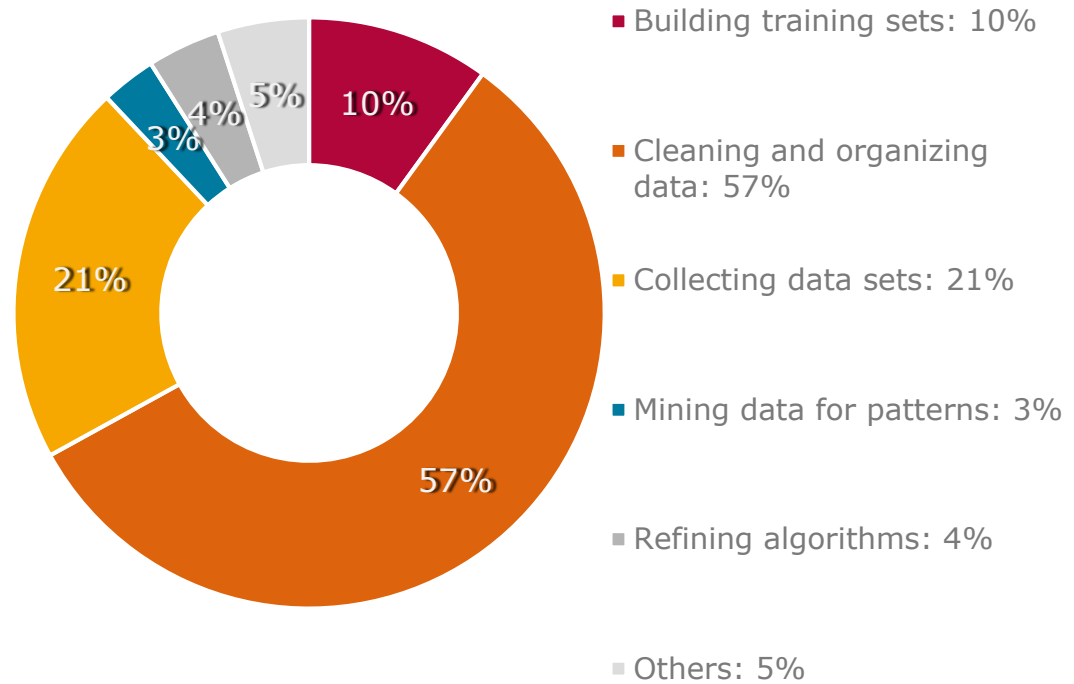
Felix Naumann
Data Science 2019

Data preparation in reality

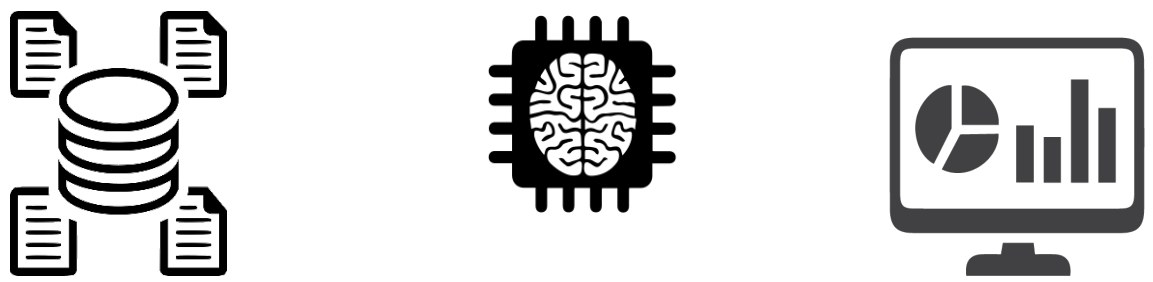
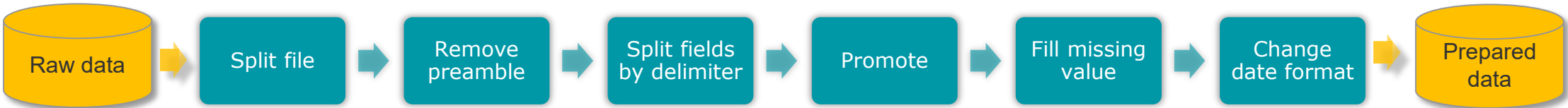
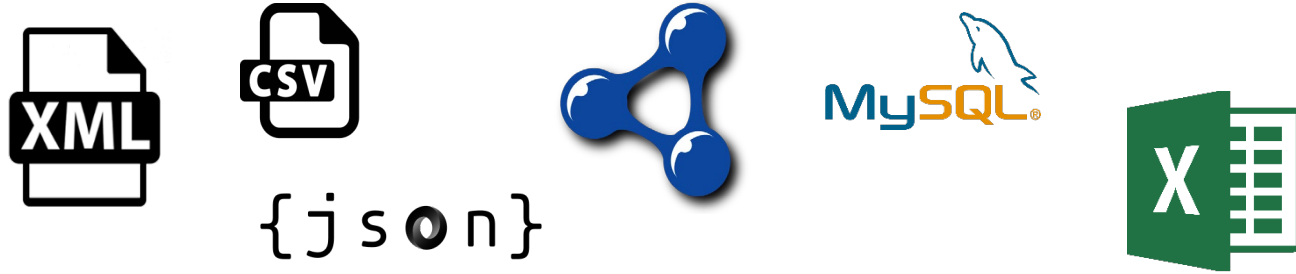
What data scientists spend the **most time** doing?



What is the **least enjoyable** part of data science?



Data preparation pipelines



Felix Naumann
Data Science 2019

Overview

1. Data Science
2. Big Data
3. Data Profiling
4. Data Preparation
- 5. Data Cleaning**



<https://unsplash.com/photos/vGefUiWm0xI>

Felix Naumann
Data Science 2019

40

Its late...

Difficult names

488941	britney spears	29	britent spears	9	brinttany spears	5	brney spears	3	britiy spears	2	brirreny spears
40134	brittany spears	29	brittnany spears	9	britanay spears	5	broitney spears	3	britmeny spears	2	brirtany spears
36315	brittney spears	29	britttany spears	9	britinany spears	5	brotny spears	3	britneeey spears	2	brirttany spears
24342	britany spears	29	btiney spears	9	britn spears	5	bruteny spears	3	britnehy spears	2	brirttney spears
7331	britny spears	26	birttney spears	9	britnew spears	5	btiyney spears	3	britnelly spears	2	britain spears
6633	briteny spears	26	breitney spears	9	britneyn spears	5	britttney spears	3	britnesy spears	2	britane spears
2696	brittney spears	26	brinity spears	9	britrney spears	5	gritney spears	3	britnetty spears	2	britaneny spears
1807	briney spears	26	britenay spears	9	brtiny spears	5	spritney spears	3	britnex spears	2	britannia spears
1635	brittny spears	26	britneyt spears	9	brtitttney spears	4	bittny spears	3	britneyxxx spears	2	britann spears
1479	brintey spears	26	brittan spears	9	brtny spears	4	bnritney spears	3	britnity spears	2	britannna spears
1479	britanny spears	26	brittne spears	9	brytny spears	4	brandy spears	3	brintey spears	2	brittannie spears
1338	britiny spears	26	btittany spears	9	rbitney spears	4	brbritley spears	3	brintyey spears	2	britannt spears
1211	britnet spears	24	beitney spears	8	birtiny spears	4	breatiny spears	3	britterny spears	2	britannu spears
1096	britiney spears	24	birteny spears	8	bithney spears	4	breetney spears	3	brittneey spears	2	britanyl spears
991	britaney spears	24	brightney spears	8	brattany spears	4	bretiney spears	3	britttney spears	2	britanyt spears
991	britnay spears	24	brintiny spears	8	breitny spears	4	brfitney spears	3	brittnyey spears	2	briteeny spears
811	brithney spears	24	brintany spears	8	breteny spears	4	briattany spears	3	brityen spears	2	britenany spears
811	brtiney spears	24	briteny spears	8	brightny spears	4	brieteny spears	3	briytney spears	2	britenet spears
664	birtney spears	24	britini spears	8	brintay spears	4	briety spears	3	brltney spears	2	briteniy spears
664	brintney spears	24	britnwy spears	8	brinttey spears	4	briitny spears	3	broteny spears	2	britenys spears
664	briteney spears	24	brittni spears	8	briotney spears	4	briittany spears	3	brtaney spears	2	britianey spears
601	bitney spears	24	brittnie spears	8	britanys spears	4	brinie spears	3	brtiiany spears	2	britin spears
601	brinty spears	21	britley spears	8	britley spears	4	brinteny spears	3	brtinay spears	2	brinary spears
544	brittaney spears	21	birtany spears	8	britneyb spears	4	brintne spears	3	brtinney spears	2	brity spears
544	brittnay spears	21	biteny spears	8	britrney spears	4	britaby spears	3	brtitany spears	2	britaney spears
364	britey spears	21	bratney spears	8	brinty spears	4	britaey spears	3	brtiteny spears	2	brtinat spears
364	brittiny spears	21	britani spears	8	brittner spears	4	britainey spears	3	brtnet spears	2	brtlnbey spears
329	brtney spears	21	britanie spears	8	brottany spears	4	britinie spears	3	brytiny spears	2	brtndy spears
269	bretney spears	21	briteany spears	7	baritney spears	4	britinney spears	3	btney spears	2	brtneh spears
269	britneys spears	21	brittay spears	7	birntey spears	4	britmney spears	3	drittney spears	2	brtneeny spears
244	britne spears	21	brittinay spears	7	biteney spears	4	britnear spears	3	pretney spears	2	brtney6 spears
244	brytney spears	21	brtany spears	7	bitiny spears	4	britel spears	3	rbritney spears	2	brtneye spears
220	breatney spears	21	brtiany spears	7	breateny spears	4	brtneuy spears	2	barittany spears	2	brtneyh spears
220	britiany spears	19	birney spears	7	brianty spears	4	brtnewy spears	2	bbbritney spears	2	brtneym spears
199	brttney spears	19	brirtney spears	7	brintye spears	4	brttnmey spears	2	kbitney spears	2	brtneyyy spears
163	brttny spears	19	brttnay spears	7	britianny spears	4	brittaby spears	2	bbritny spears	2	brtney spears

FIFA registration form (2010)

Nationality Select
Country of Residence Palestine
Mother Tongue Palestine
Preferred FIFA Language Palestine, British Mandate
Secondary FIFA Language Panama
Organisation Name Papua New Guinea
Organisation Role (Prof) Paraguay
 Notes (Max 2000 chars) Peru
 Philippines
 Poland
 Portugal
 Puerto Rico
 Qatar
 Representations of Czechs and Slovaks (RCS)
 Republic of Ireland
 Réunion
 Rhodesia
 Romania
 Russia
 Rwanda
 Saar
 Samoa
 San Marino
 São Tomé e Príncipe
 Saudi Arabia
 Scotland
 Senegal
 Serbia
 Serbia and Montenegro
 Seychelles
 Sierra Leone



Select
 German Democratic Republic
 German Democratic Republic
 Germany
 Germany Federal Republic
 Ghana
 Gibraltar
 Great Britain

with a public account such as Hotmail or

Select
 All Ireland (all-Ireland pre 1921)
 All Ireland (all-Ireland pre 1921)
 American Samoa
 Andorra
 Angola

Wales
 Yemen
 Yemen PDR
 Yugoslavia
 Zaire
 Zambia
 Zimbabwe

Felix Naumann
Data Science 2019

Directmarketing by The Economist

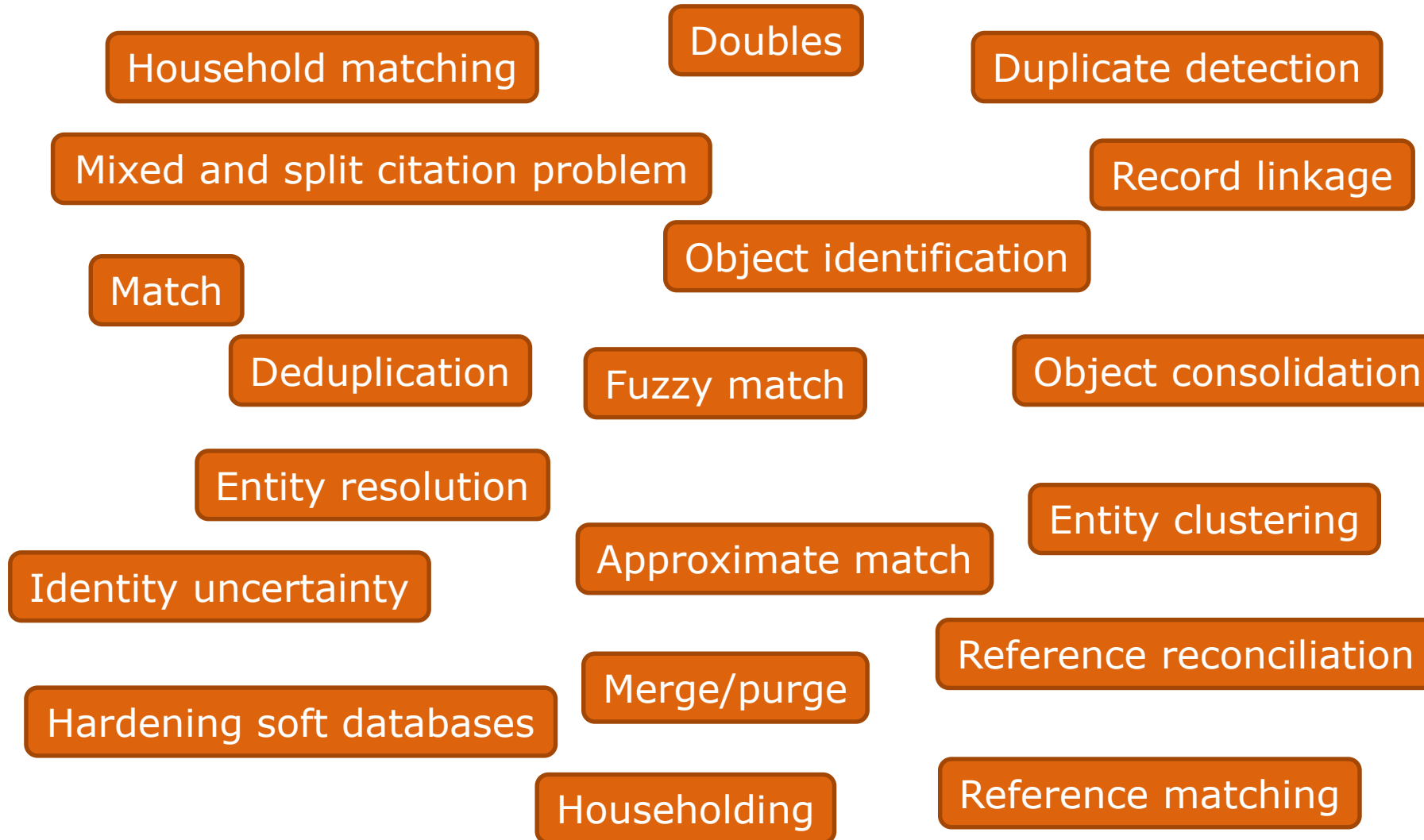


Felix Naumann
Data Science 2019

Data Cleaning: Duplicate Detection

- Duplicate detection is the discovery of multiple representations of the same real-world object.
- Problem 1: **Representations are not identical.**
 - *Fuzzy duplicates*
- Solution: Similarity measures / models
 - Value- and record-comparisons
 - Domain-dependent or domain-independent
- Problem 2: **Datasets are large.**
 - Quadratic complexity: Comparison of every pair of records.
- Solution: Algorithms
 - E.g., avoid comparisons by partitioning.

Ironically, "Duplicate Detection" has many Duplicates

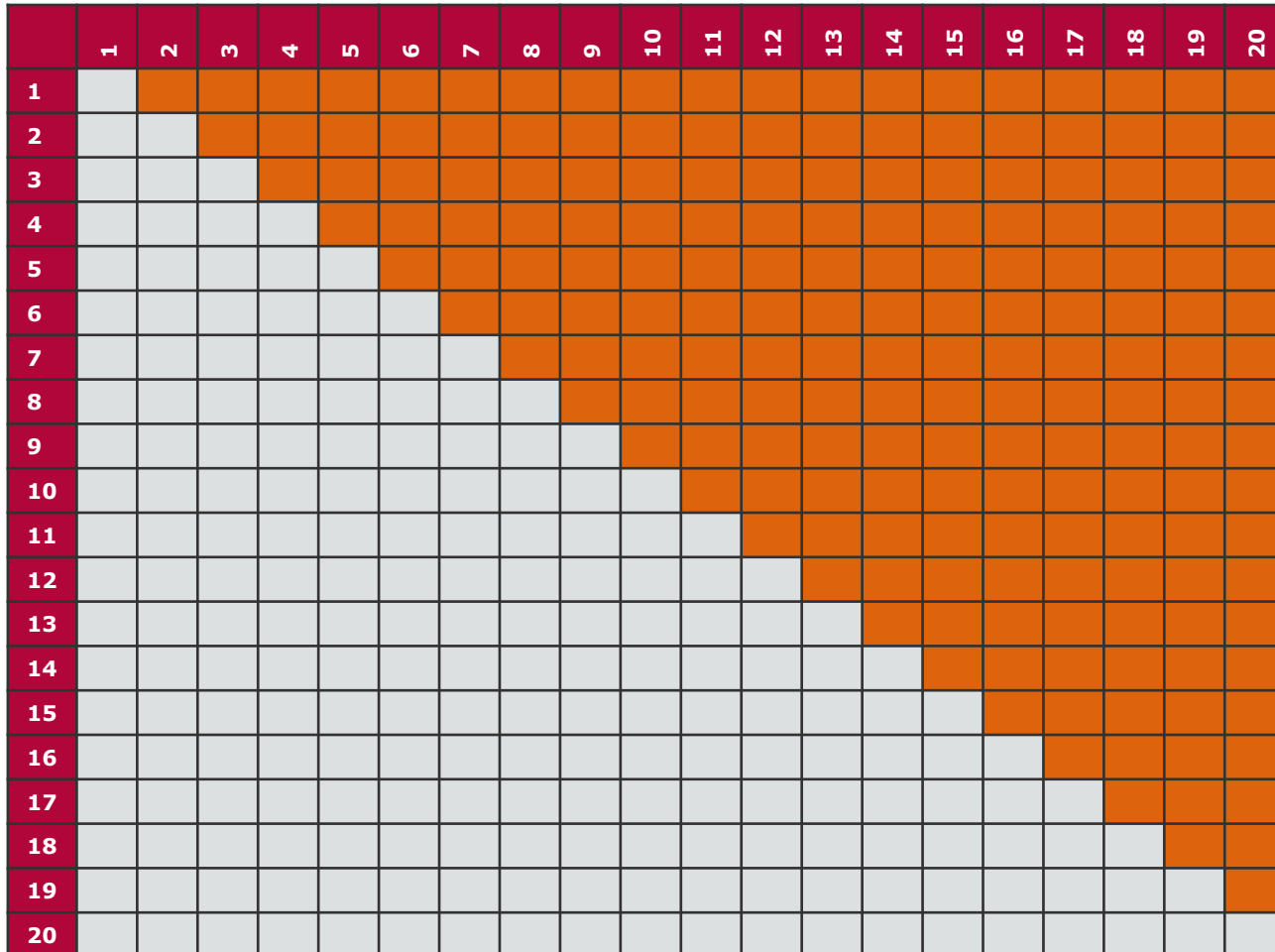


Number of comparisons: All pairs

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1																				
2																				
3																				
4																				
5																				
6																				
7																				
8																				
9																				
10																				
11																				
12																				
13																				
14																				
15																				
16																				
17																				
18																				
19																				
20																				

400
comparisons

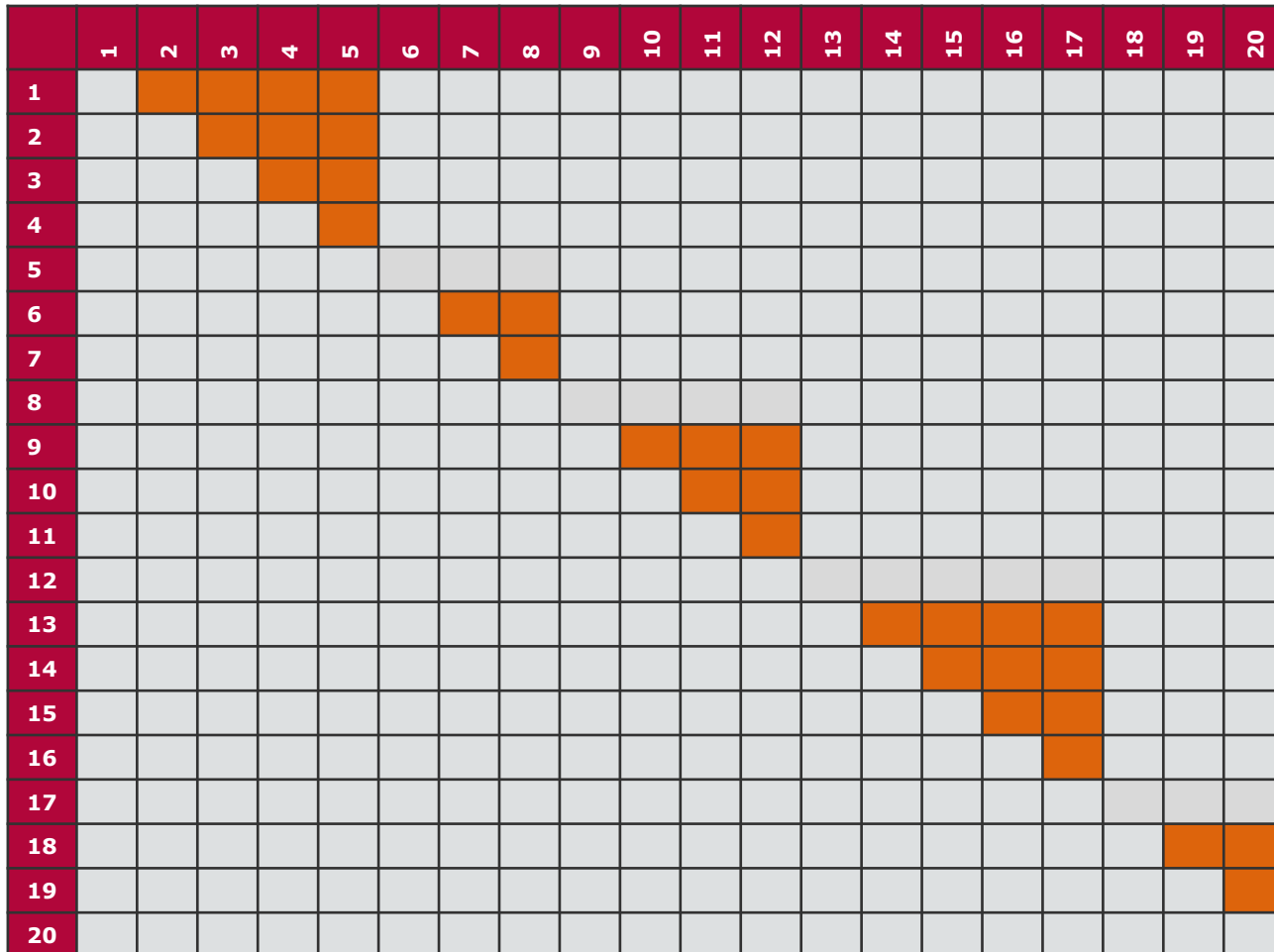
Symmetry of Similarity



190
comparisons

Felix Naumann
Data Science 2019

Blocking by zip-code



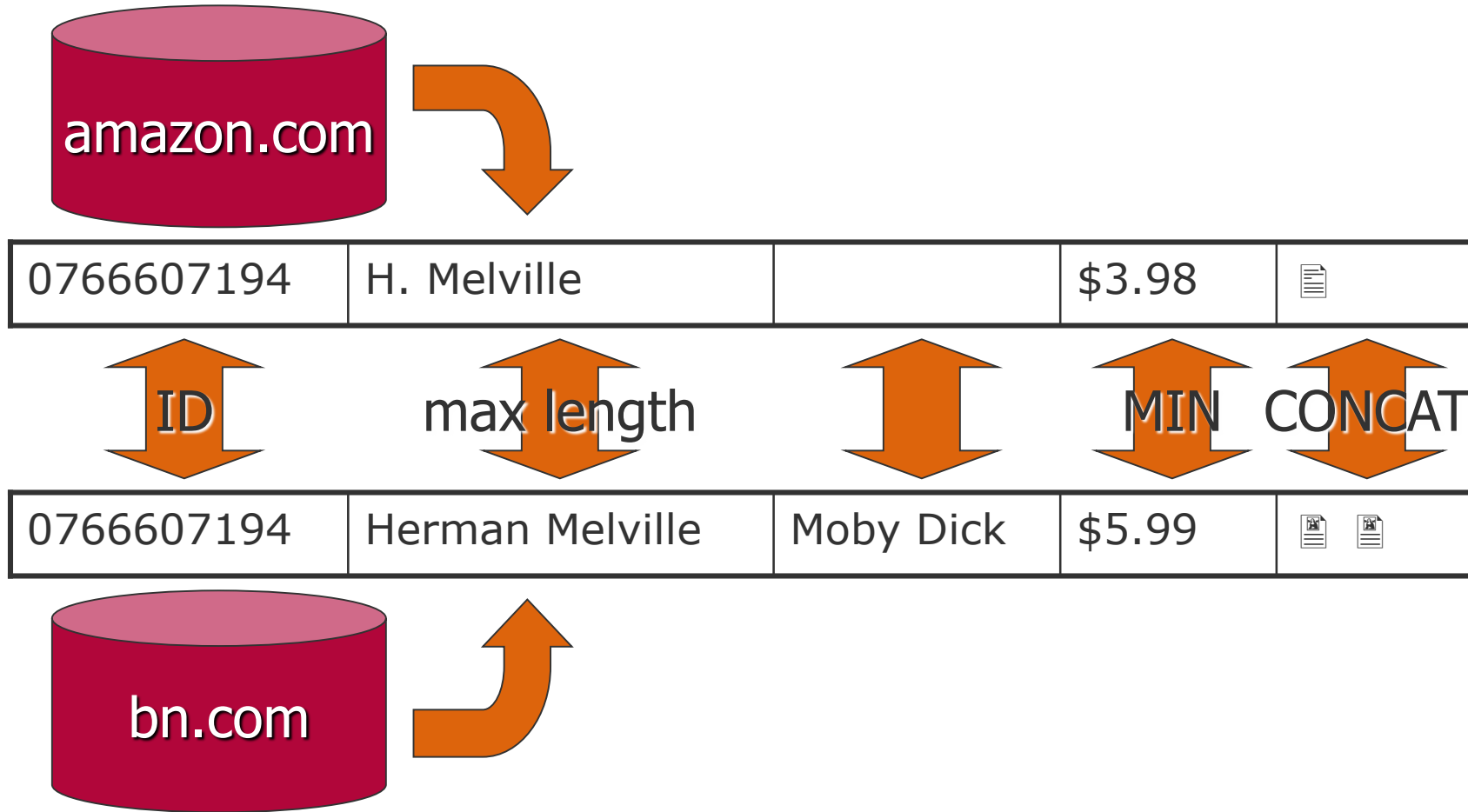
32
comparisons

Sorting by zip-code

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1		■	■	■																
2			■	■	■															
3				■	■	■														
4					■	■	■													
5						■	■	■												
6							■	■	■											
7								■	■	■										
8									■	■	■									
9										■	■	■								
10											■	■	■							
11												■	■	■						
12													■	■	■					
13														■	■	■				
14															■	■	■			
15																■	■	■		
16																	■	■	■	
17																		■	■	■
18																			■	■
19																				■
20																				

54
comparisons

Data Fusion



Summary

1. Data Science
2. Big Data
3. Data Profiling
4. Data Preparation
5. Data Cleaning



Felix Naumann
Data Science 2019

<https://unsplash.com/photos/vGefUiWm0xI>

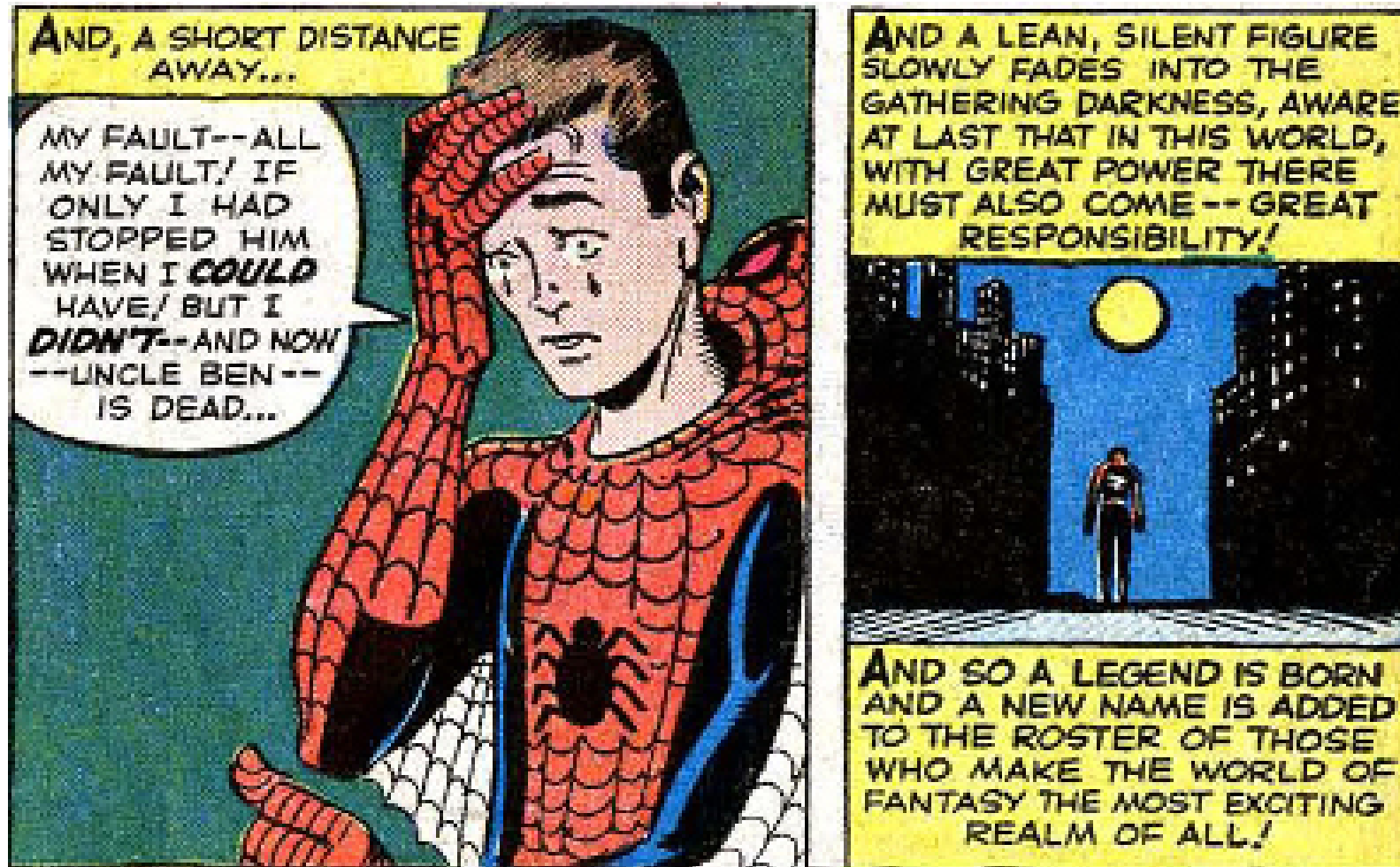
Big Data and Ethics



- Industry keynote speakers on credit ratings using big data
 - “If the data is out there, we will find it.”
 - “... and that is why I closed my Twitter account.”
 - “... and that is why I had my son close his Twitter account.”

Felix Naumann
Data Science 2019

With Great Power there must also come – Great Responsibility



Spider-Man from *Amazing Fantasy* #15, August 1962