

Lecture

Data Science

Prof. Dr. Steffen Staab



Christoph Kling



Consulting



Dr. Christoph Tempich
Head of Consulting

 christoph.tempich@inovex.de

 Christoph Tempich bei Google+

 Christoph Tempich bei Twitter

 Christoph Tempich bei XING

- Module number: 04IN2043

Target audience

- Master Web Science, Inf, CV, eGov, WI, IM
- Bachelor (Wahlpflicht): Inf, CV, WI
- [http://www.uni-koblenz-landau.de/campus-koblenz/
fb4/west/teaching/ss14/proseminar](http://www.uni-koblenz-landau.de/campus-koblenz/fb4/west/teaching/ss14/proseminar)
- Klips-registration

- Programming with data
 - ◆ Statistics environment (e.g. R)
 - ◆ Cloud programming (e.g. Hadoop)
- Paper based excercises

Date: (tbd)

Will ask for:

- ◆ Knowledge, understanding and capability
 - Acquired from lecture
 - Acquired from pen-and-paper excercises
 - Acquired from programming

Subject to change

DATA SCIENCE: OVERVIEW

Most slides from

Hanspeter Pfister
pfister@seas.harvard.edu

Joe Blitzstein
blitzstein@stat.harvard.edu

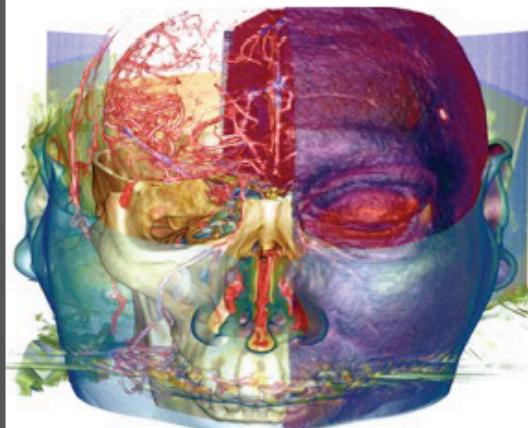
<https://drive.google.com/folderview?id=0BxYkKyLxfsNVd0xicUVDS1dIS0k>

<http://www.kdnuggets.com/2013/11/harvard-cs109-data-science-course-resources-free-online.html>

1. History and background, from statistics to programming; examples
2. Process
3. Background in statistics
4. Hypothesis driven research
5. Programming paradigms
6. Visualization
7. Simple machine learning on large scale data
8. Example application domain: text
9. Privacy

Data Science

To gain insights into data through computation, statistics, and visualization

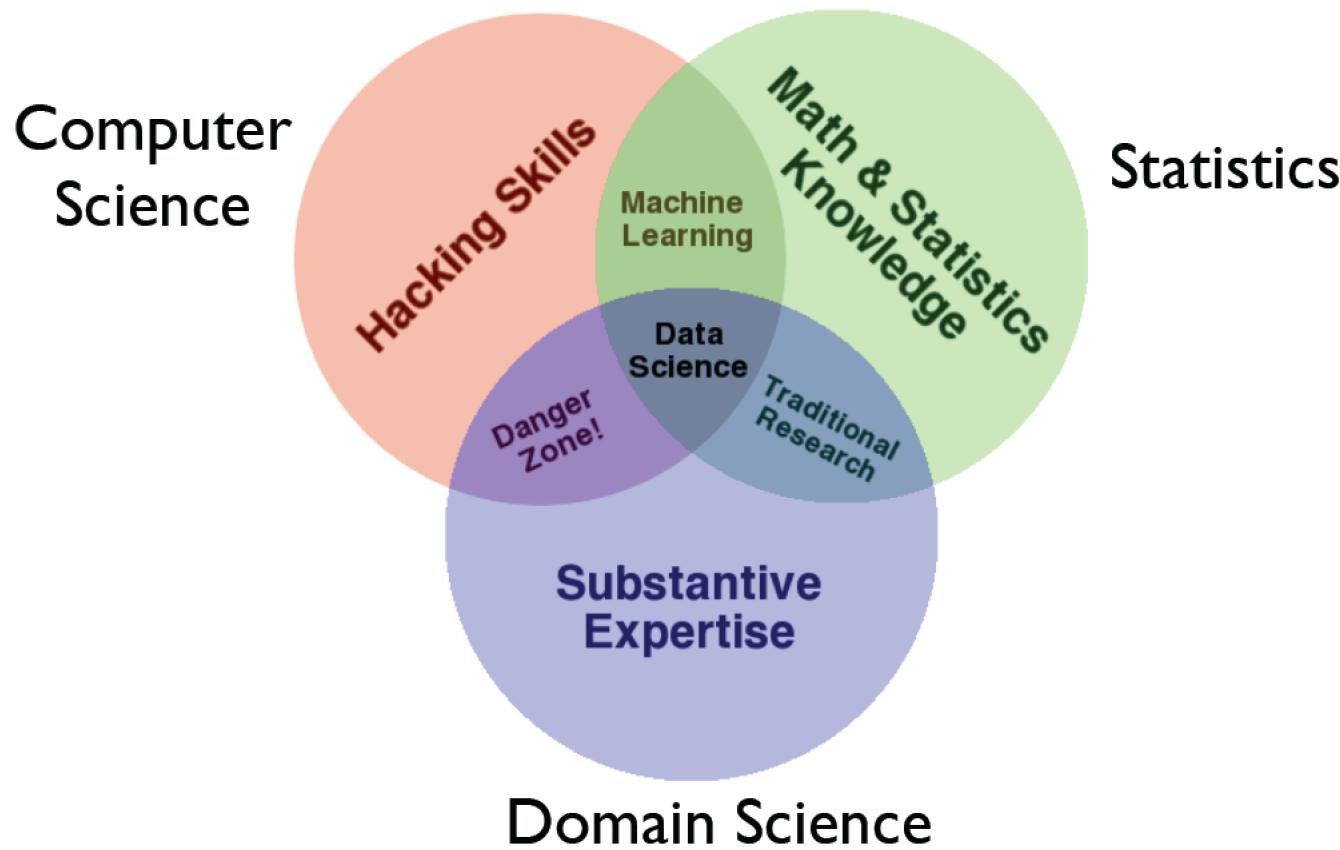


“A data scientist is someone who knows more statistics than a computer scientist and more computer science than a statistician.”

- Josh Blumenstock

“Data Scientist = statistician + programmer +
coach + storyteller + artist”

- Shlomo Aragmon



Drew Conway

Machine

Data Management

Data Mining

Machine Learning

Business Intelligence

Statistics

Data Science

Human

Human Cognition

Perception

Story Telling

Decision Making
Theory



Inspired by Daniel Keim, "Visual Analytics: Definition, Process, and Challenges"

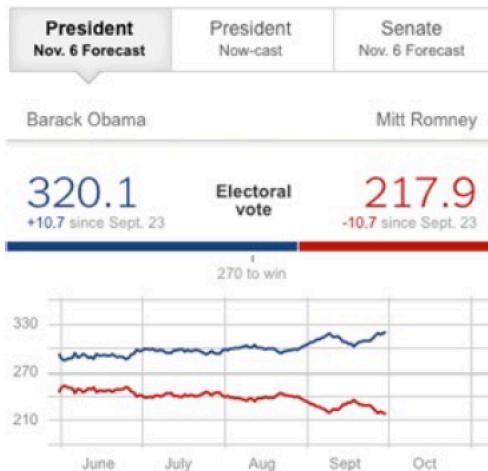
EXAMPLES

“Nate Silver won the election”

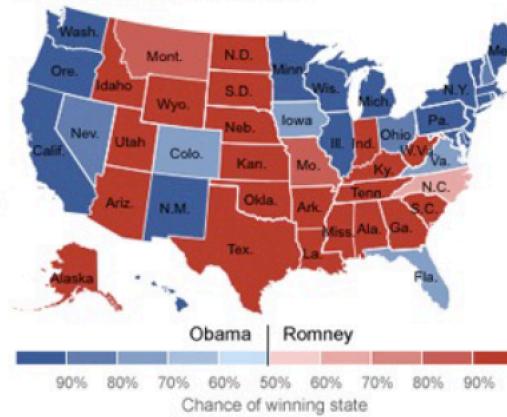
– Harvard Business Review

Five Thirty Eight Forecast

Updated 12:27 AM ET on Oct. 1



State-by-State Probabilities



Electoral Vote Distribution

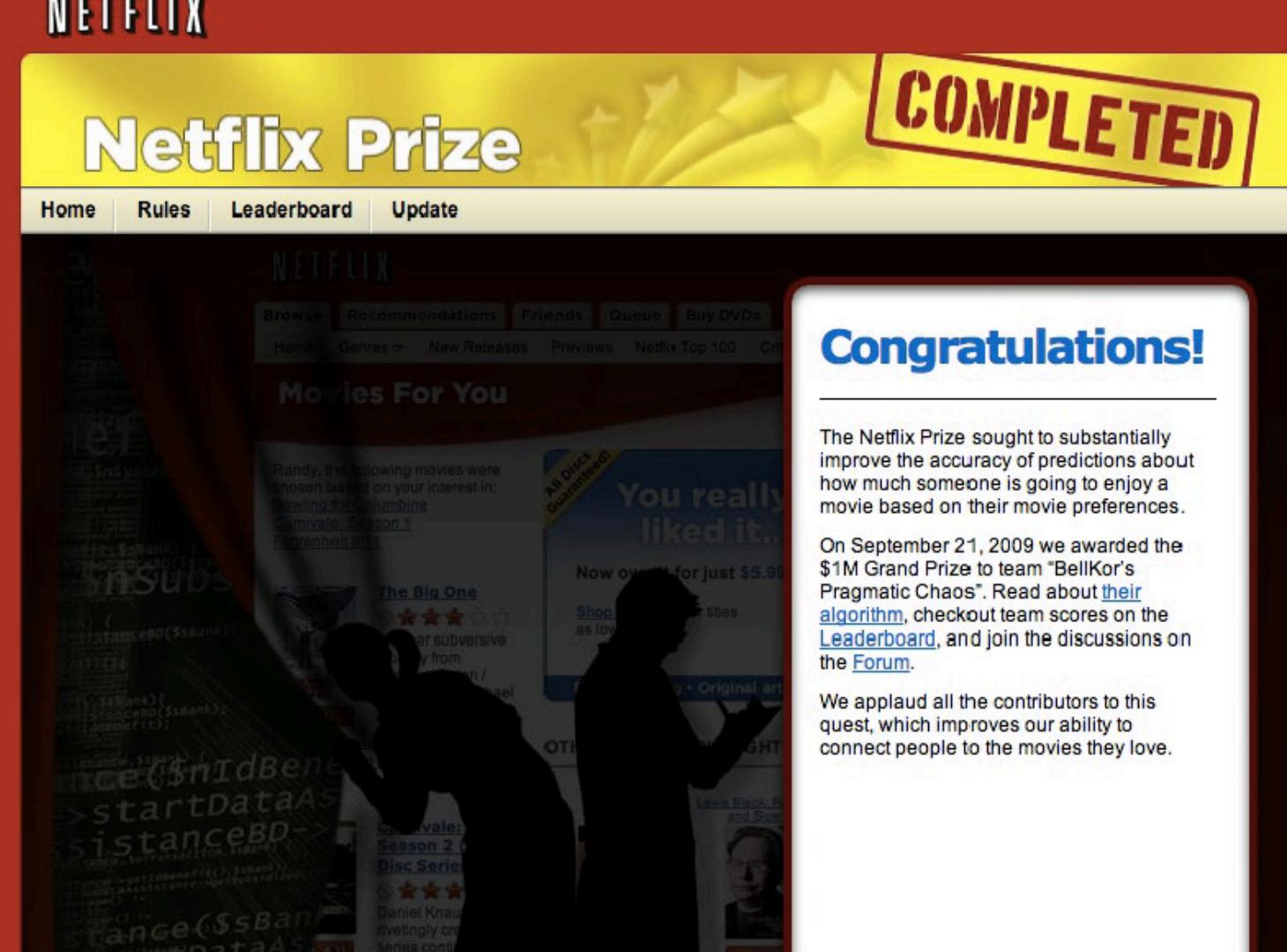
The probability that President Obama receives a given number of electoral votes.

Silver: „Pundits are no better than a coin toss.“



- use many data sources (the plural of anecdote is not data)
- understand how the data were collected (sampling is essential)
- weight the data thoughtfully (not all polls are equally good)
- use statistical models (not just hacking around in Excel)
- understand correlations (e.g., states that trend similarly)
- think like a Bayesian, check like a frequentist (reconciliation)
- have good communication skills
(What does a 60% probability even mean? How can we visualize, validate, and understand the conclusions?)

Netflix Prize



The image shows a screenshot of the Netflix Prize website. At the top, the Netflix logo is visible. Below it, a large yellow banner features the words "Netflix Prize" and a red "COMPLETED" stamp. A navigation bar with links for "Home", "Rules", "Leaderboard", and "Update" is present. The main content area displays a blurred screenshot of the Netflix interface, showing movie recommendations and user profiles. To the right, a white box contains the word "Congratulations!" in blue text. Below this, a paragraph explains the purpose of the prize and mentions the winning team, BellKor's Pragmatic Chaos. It also encourages users to explore the leaderboard and forum.

NETFLIX

Netflix Prize

COMPLETED

Home Rules Leaderboard Update

Movies For You

Randy, the following movies were chosen based on your interest in...
Bowling for Columbine
Carnivale: Season 1
Fahrenheit 9/11

You really liked it!

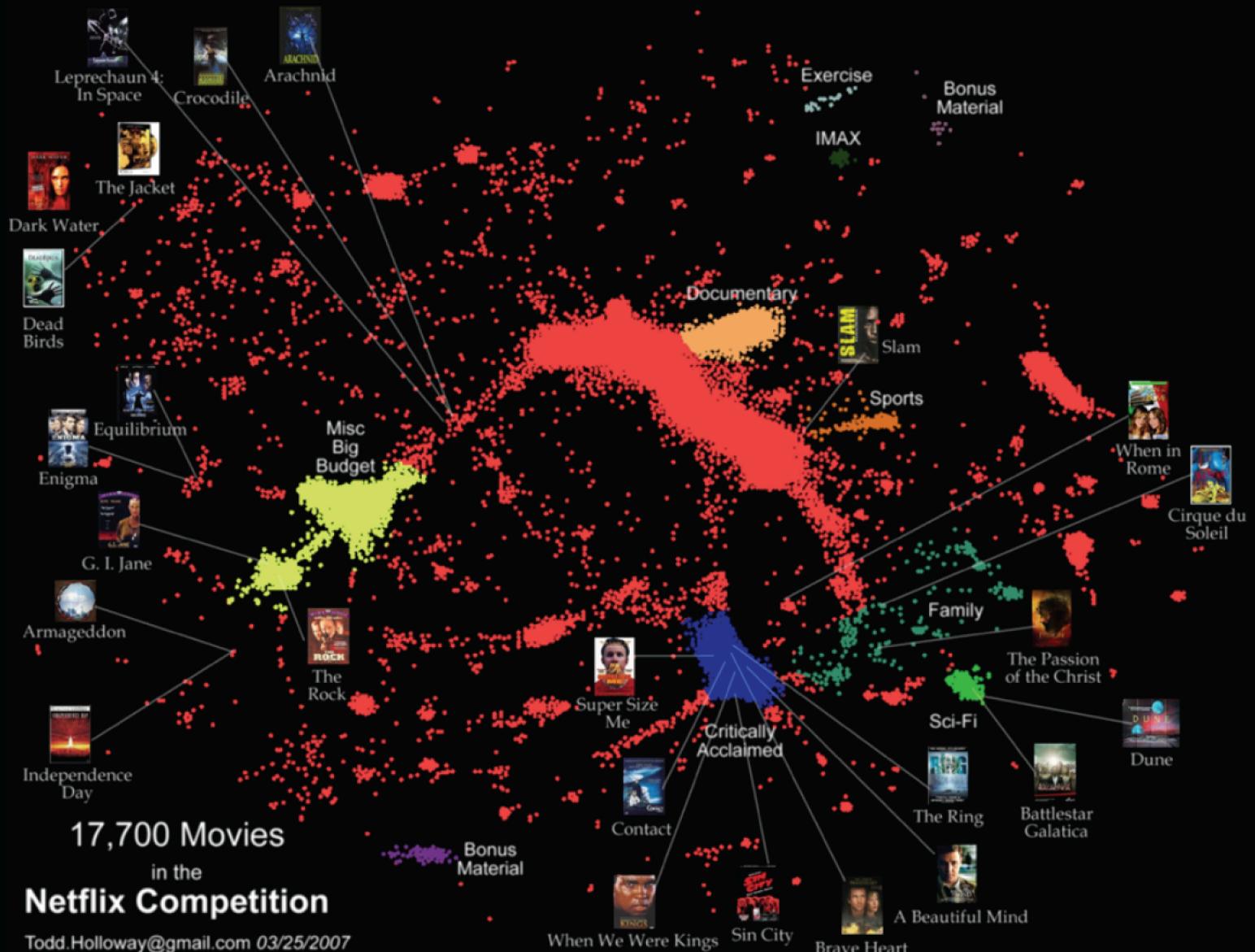
Now owned for just \$5.99

Congratulations!

The Netflix Prize sought to substantially improve the accuracy of predictions about how much someone is going to enjoy a movie based on their movie preferences.

On September 21, 2009 we awarded the \$1M Grand Prize to team "BellKor's Pragmatic Chaos". Read about [their algorithm](#), checkout team scores on the [Leaderboard](#), and join the discussions on the [Forum](#).

We applaud all the contributors to this quest, which improves our ability to connect people to the movies they love.





<http://blogs.hbr.org/2012/10/big-data-hype-and-reality/>

- *massive data* (500k users, 20k movies, 100m ratings)
- *curse of dimensionality* (very high-dimensional problem)
- *missing data* (99% of data missing; not missing at random)
- *extremely complicated set of factors that affect people's ratings of movies* (actors, directors, genre, ...)
- need to avoid *overfitting* (test data vs. training data)

No silver bullet: Netflix prize only semi-successful

<http://www.techdirt.com/blog/innovation/articles/20120409/03412518422/why-netflix-never-implemented-algorithm-that-won-netflix-1-million-challenge.shtml>

OPPORTUNITY: BIG DATA

What Happens in an Internet Minute?



And Future Growth is Staggering

Today, the number of networked devices



=

By 2015, the number of networked devices



= 2x

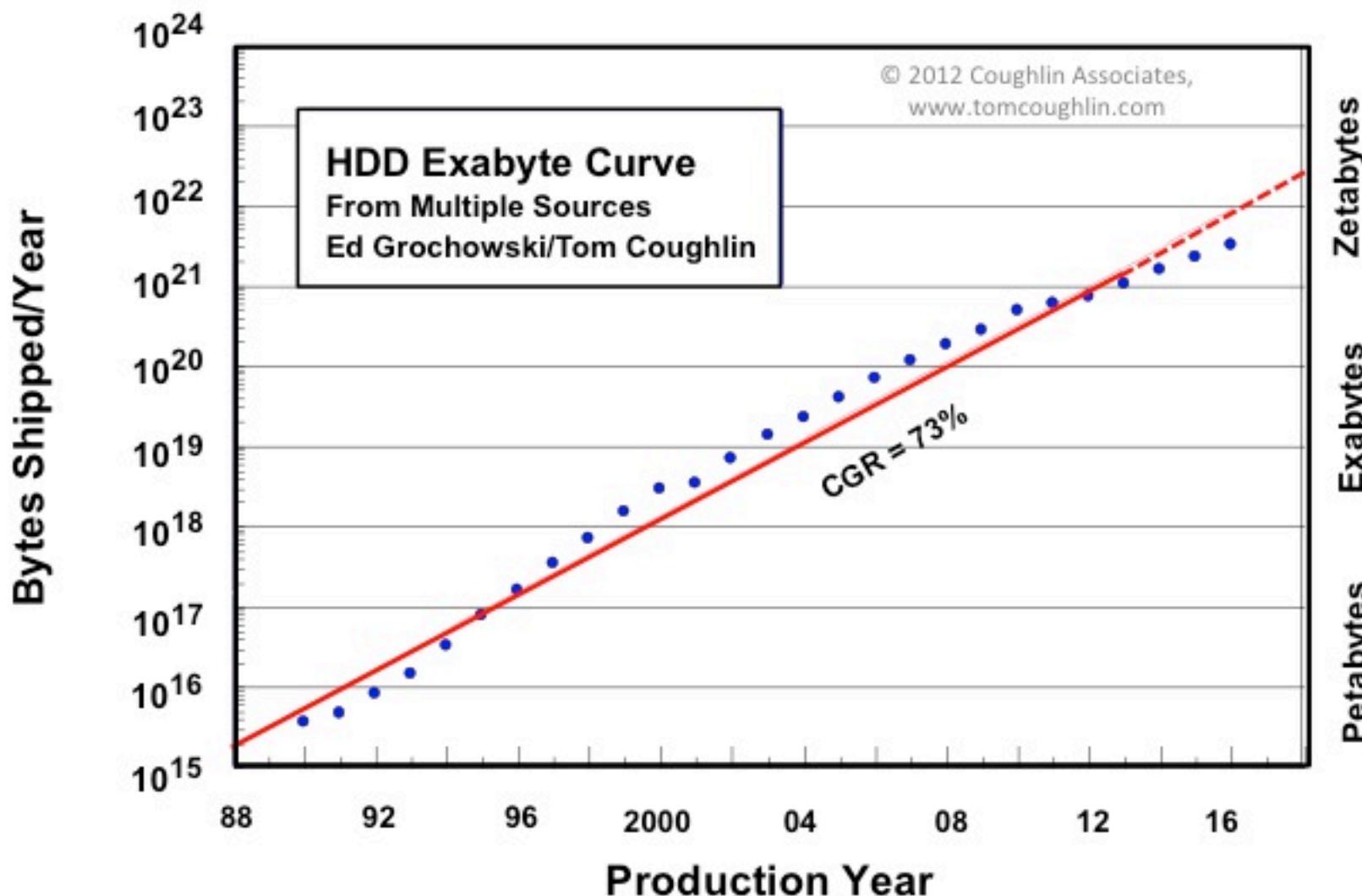
In 2015,
it would take
you 5 years

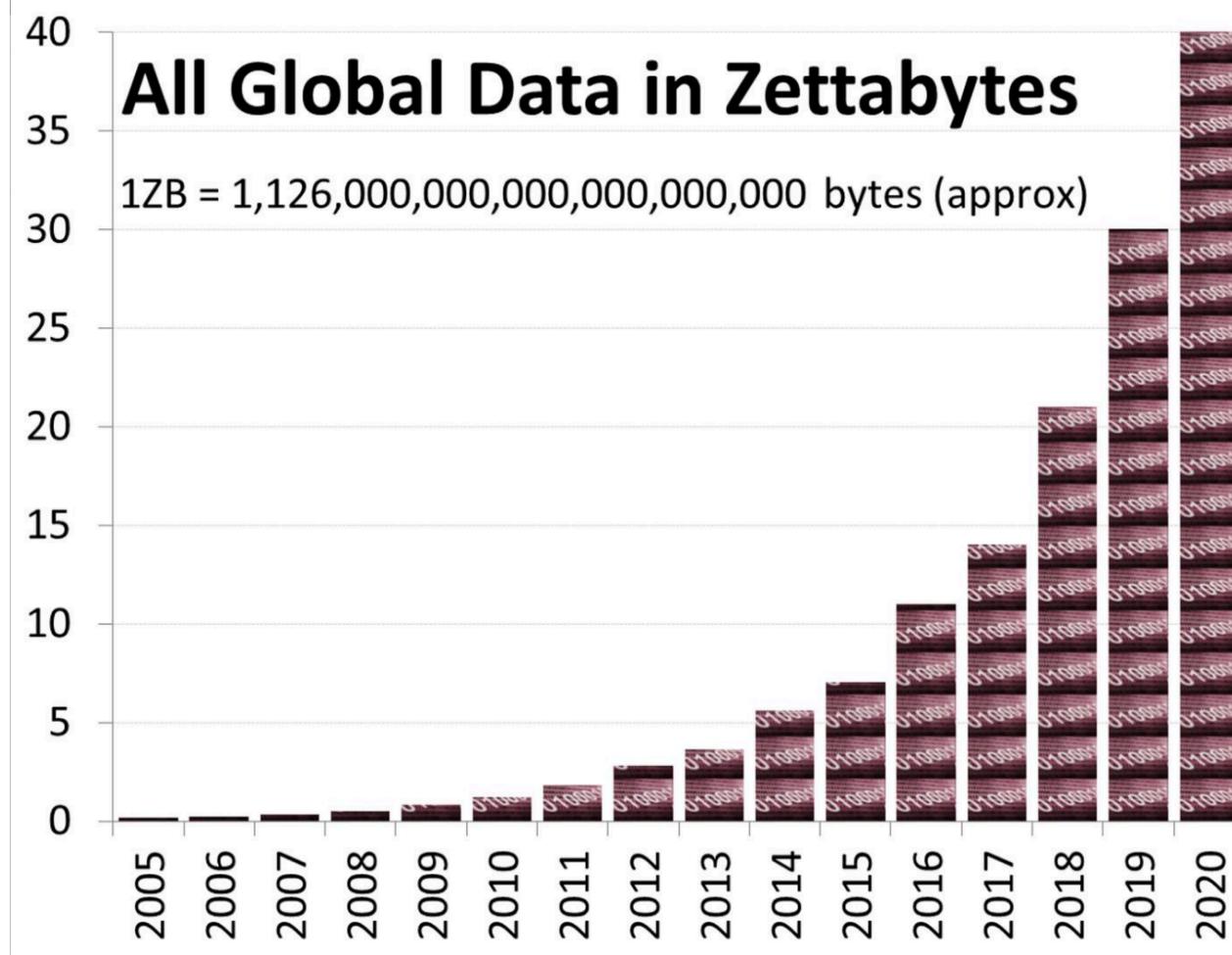


to view all video crossing IP networks each second

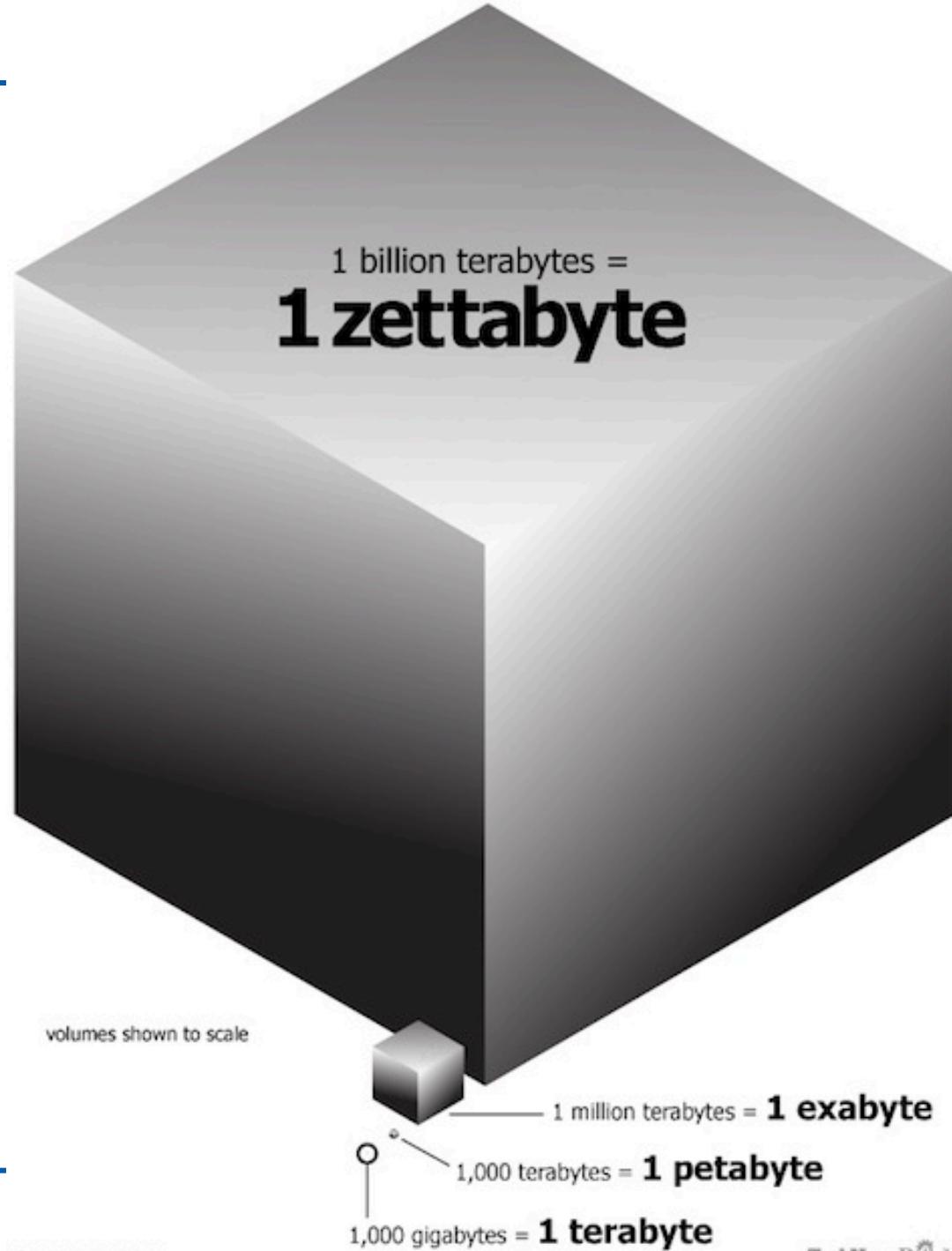
intel

WW HDD CAPACITY SHIPMENTS (1999 THROUGH 2016)





Big Data – what it means



“Between the dawn of civilization and 2003, we only created five exabytes of information; now we’re creating that amount every two days.”

Eric Schmidt, Google (and others)

A screenshot of a Google search results page. The search query "you tube cat videos" is entered in the search bar. Below the search bar, the "Web" tab is selected, and a red circle highlights the text "About 1,030,000,000 results (0.33 seconds)". The search results include an advertisement for "Vote For Funny Cat Videos - TheFriskies.com" and a video thumbnail for "Probably the Funniest Cat Video You'll Ever See - YouTube".

you tube cat videos

Web Images Maps Shopping News More Search tools

About 1,030,000,000 results (0.33 seconds)

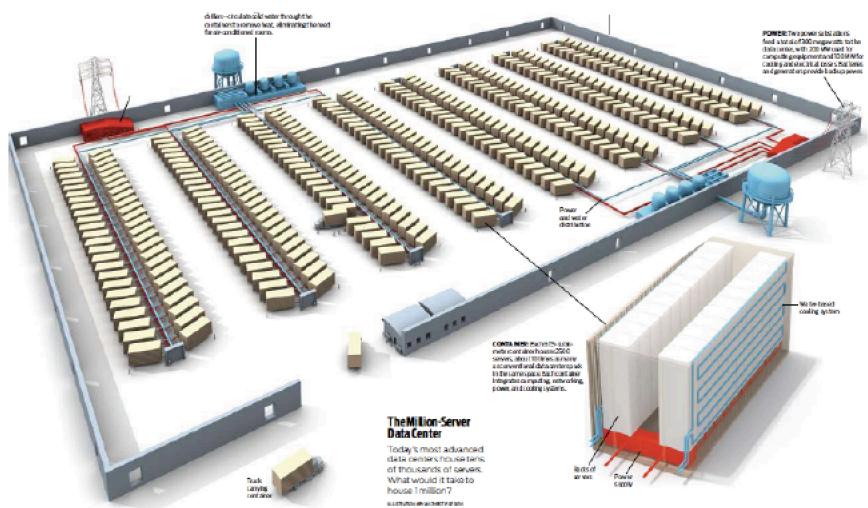
Ad related to you tube cat videos ⓘ

[Vote For Funny Cat Videos - TheFriskies.com](#)
www.thefriskies.com/ContestEntry ▾
The Friskies Will Honor The Best New Cat Videos. Cast Your Vote Now!

More About The Awards Visit Friskies.com
How To Make A Cat Video Help Us Give Back

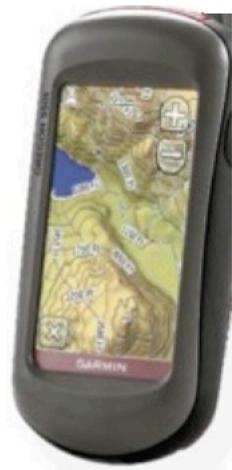
[Probably the Funniest Cat Video You'll Ever See - YouTube](#)
 www.youtube.com/watch?v=SUNmLuNdil8 ▾
Jan 12, 2007 - Uploaded by lanierloo
Now don't let the corny opening fool you, this is surely the most hilarious cat video you will ever see in

Commodity Computing



Michael Franklin, UC Berkeley

Smarter Devices



Michael Franklin, UC Berkeley

Ubiquitous Connectivity



- Volume
- Velocity
- Veracity
- Variety
- priVacy

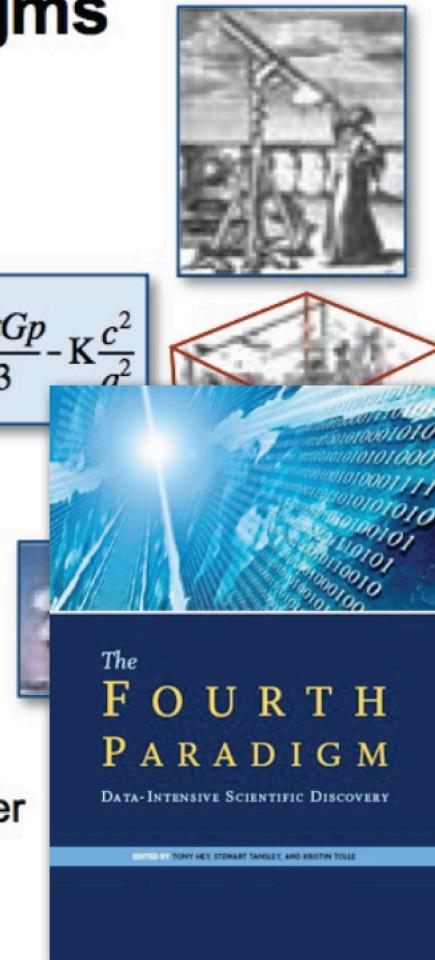


- To students of
 - ◆ Master in Informatik
 - ◆ Master in CV
- Please respond!

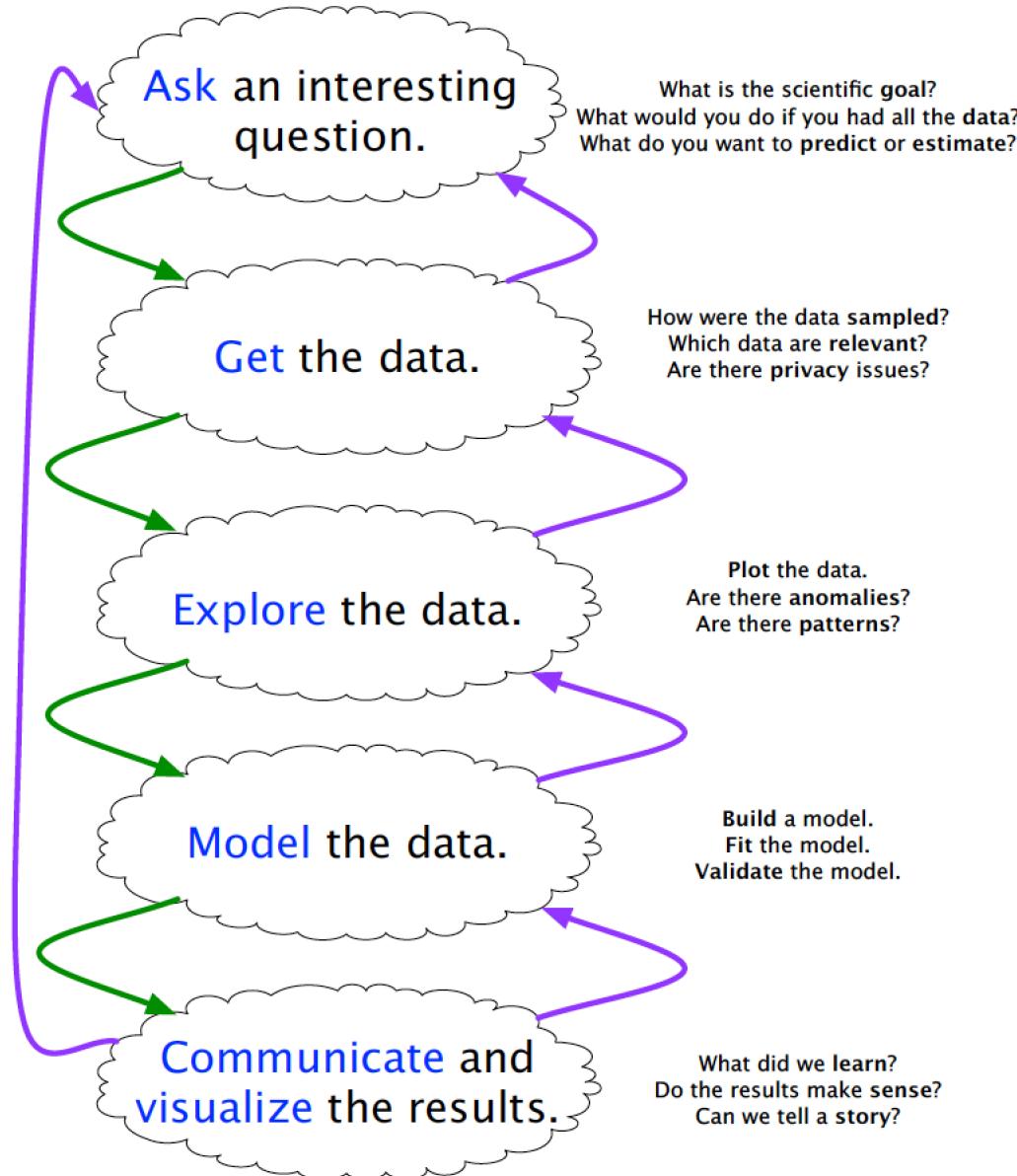
Science Paradigms

- Thousand years ago:
science was empirical
describing natural phenomena
- Last few hundred years:
theoretical branch
using models, generalizations
- Last few decades:
a computational branch
simulating complex phenomena
- Today: **data exploration (eScience)**
unify theory, experiment, and simulation
 - Data captured by instruments
or generated by simulator
 - Processed by software
 - Information/knowledge stored in computer
 - Scientist analyzes database/files
using data management and statistics

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G p}{3} - K \frac{c^2}{a^2}$$



Jim Gray, Microsoft



- *data munging/scraping/sampling/cleaning* in order to get an informative, manageable data set;
- *data storage and management* in order to be able to access data - especially big data - quickly and reliably during subsequent analysis;
- *exploratory data analysis* to generate hypotheses and intuition about the data;
- *prediction* based on statistical tools such as regression, classification, and clustering; and
- *communication* of results through visualization, stories, and interpretable summaries.

What do analysts do?

Enterprise Data Analysis and Visualization: An Interview Study

Sean Kandel, Andreas Paepcke, Joseph M. Hellerstein, and Jeffrey Heer

Abstract—Organizations rely on data analysts to model customer engagement, streamline operations, improve production, inform business decisions, and combat fraud. Though numerous analysis and visualization tools have been built to improve the scale and efficiency at which analysts can work, there has been little research on how analysis takes place within the social and organizational context of companies. To better understand the enterprise analysts’ ecosystem, we conducted semi-structured interviews with 35 data analysts from 25 organizations across a variety of sectors, including healthcare, retail, marketing and finance. Based on our interview data, we characterize the process of industrial data analysis and document how organizational features of an enterprise impact it. We describe recurring pain points, outstanding challenges, and barriers to adoption for visual analytic tools. Finally, we discuss design implications and opportunities for visual analysis research.

Index Terms—Data, analysis, visualization, enterprise.

1 INTRODUCTION

Organizations gather increasingly large and complex data sets each year. These organizations rely on data analysis to model customer engagement, streamline operations, improve production, inform sales and business decisions, and combat fraud. Within organizations, an increasing number of individuals—with varied titles such as “business analyst”, “data analyst” and “data scientist”—perform such analyses. These analysts constitute an important and rapidly growing user population for analysis and visualization tools.

Enterprise analysts perform their work within the context of a larger organization. Analysts often work as a part of an analysis team or business unit. Little research has observed how existing infrastructure, available data and tools, and administrative and social conventions within an organization impact the analysis process within the enterprise. Understanding how these issues shape analytic workflows can inform the design of future tools.

ery and wrangling, often the most tedious and time-consuming aspects of an analysis, are underserved by existing visualization and analysis tools. We discuss recurring pain points within each task as well as difficulties in managing workflows across these tasks. Example pain points include integrating data from distributed data sources, visualizing data at scale and operationalizing workflows. These challenges are typically more acute within large organizations with a diverse and distributed set of data sources.

We conclude with a discussion of future trends and the implications of our interviews for future visualization and analysis tools. We argue that future visual analysis tools should leverage existing infrastructures for data processing to enable scale and limit data migration. One avenue for achieving better interoperability is through systems that specify analysis or data processing operations in a high-level language, enabling retargeting across tools or platforms. We also note

What do analysts do?

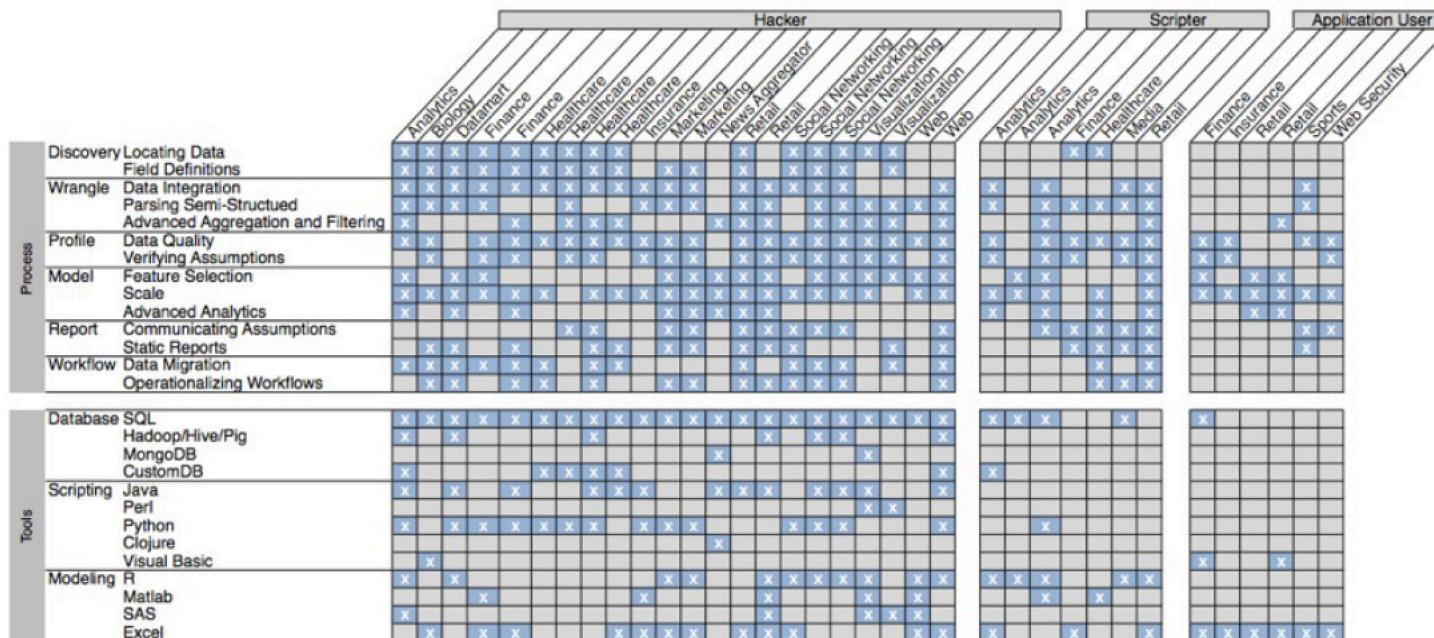


Fig. 1. Respondents, Challenges and Tools. The matrix displays interviewees (grouped by archetype and sector) and their corresponding challenges and tools. *Hackers* faced the most diverse set of challenges, corresponding to the diversity of their workflows and toolset. *Application users* and *scripters* typically relied on the IT team to perform certain tasks and therefore did not perceive them as challenges.

I spend more than half of my time integrating, cleansing and transforming data without doing any actual analysis. Most of the time I'm lucky if I get to do any analysis. Most of the time once you transform the data you just do an average... the insights can be scarily obvious. It's fun when you get to do something somewhat analytical

BE CAREFUL WITH STATISTICS

Anscombe's Quartet

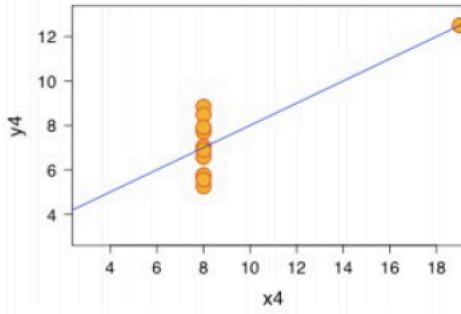
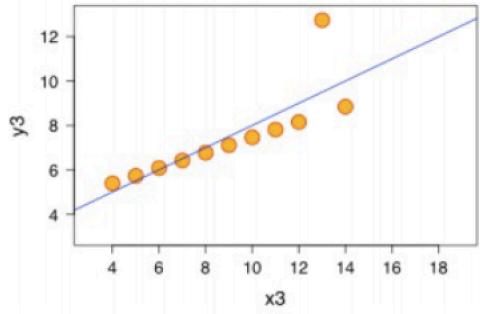
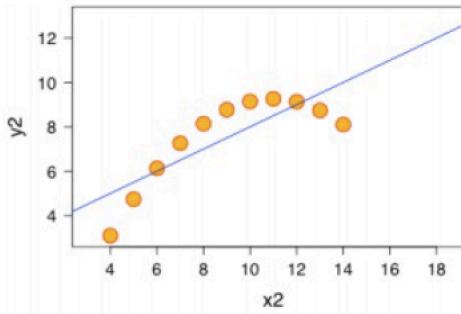
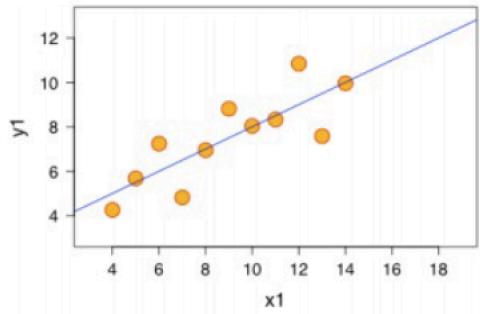
Same mean, variance, correlation, and linear regression line

Anscombe's Quartet: Raw Data								
	I		II		III		IV	
	x	y	x	y	x	y	x	y
mean	9.0		9.0		9.0		9.0	
var.	10.0		3.75		3.75		3.75	
corr.	0.816		0.816		0.816		0.816	
10.0	8.04		10.0	9.14		10.0	7.46	
8.0	6.95		8.0	8.14		8.0	6.77	
13.0	7.58		13.0	8.74		13.0	12.74	
9.0	8.81		9.0	8.77		9.0	7.11	
11.0	8.33		11.0	9.26		11.0	7.81	
14.0	9.96		14.0	8.10		14.0	8.84	
6.0	7.24		6.0	6.13		6.0	6.08	
4.0	4.26		4.0	3.10		4.0	5.39	
12.0	10.84		12.0	9.13		12.0	8.15	
7.0	4.82		7.0	7.26		7.0	6.42	
5.0	5.68		5.0	4.74		5.0	5.73	

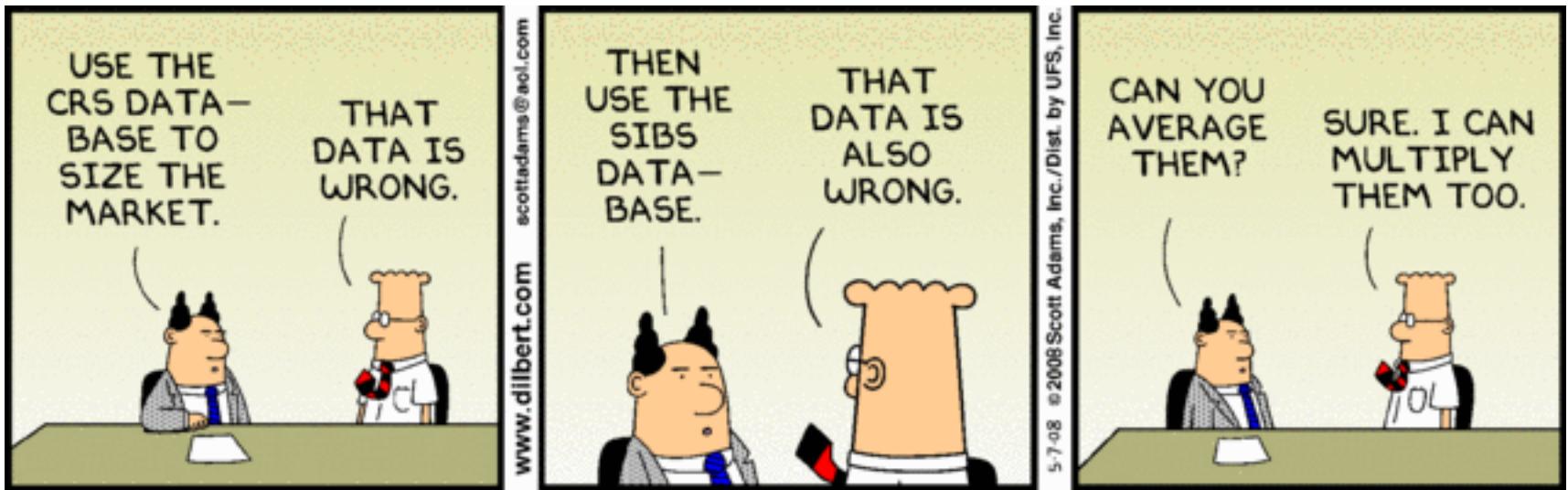
Anscombe '73

Anscombe's Quartet

Same mean, variance, correlation, and linear regression line



Anscombe '73



Example: Antibiotics (Will Burton, 1951)

TELLING A STORY

Data & Questions

What are the data types?
What are possible questions?

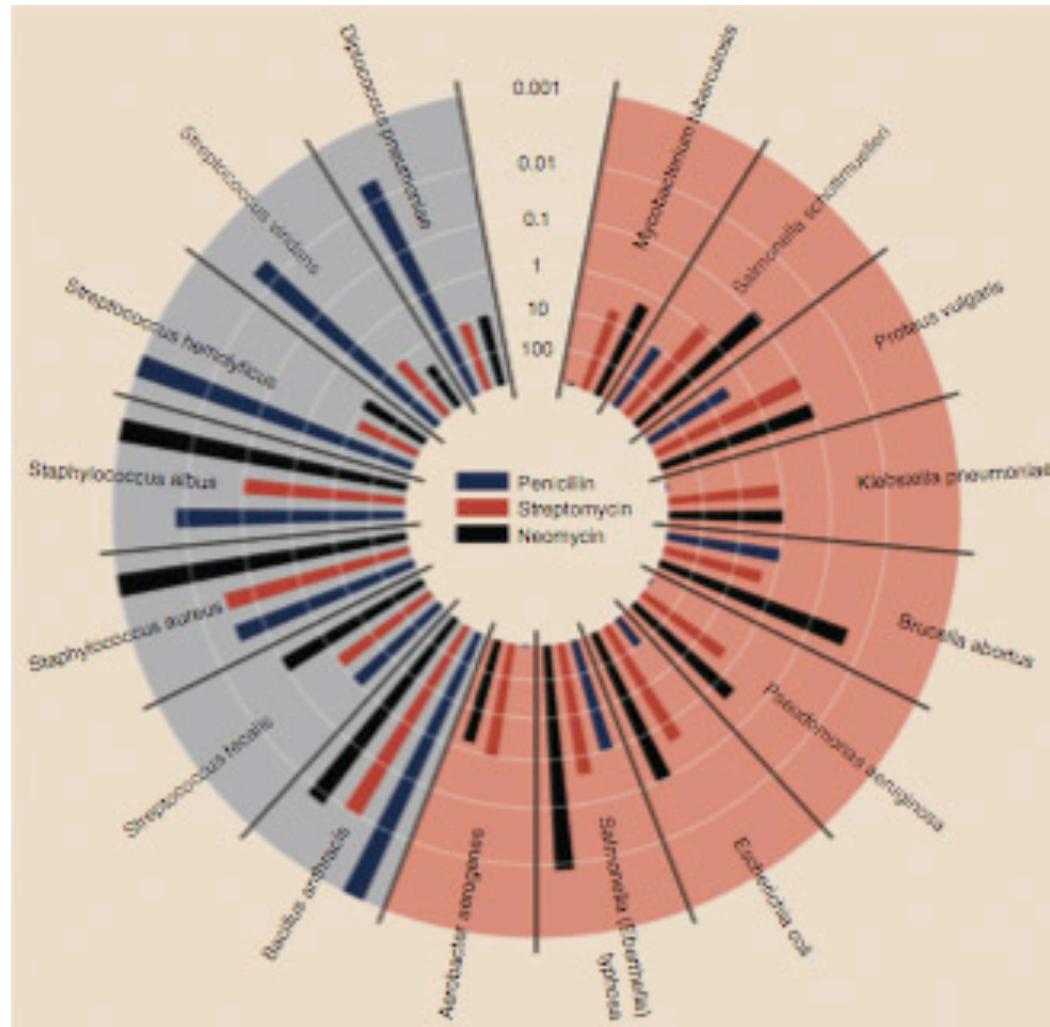
Bacteria	Penicillin	Antibiotic Streptomycin	Neomycin	Gram stain
<i>Aerobacter aerogenes</i>	870	1	1.6	–
<i>Brucella abortus</i>	1	2	0.02	–
<i>Bacillus anthracis</i>	0.001	0.01	0.007	+
<i>Diplococcus pneumoniae</i>	0.005	11	10	+
<i>Escherichia coli</i>	100	0.4	0.1	–
<i>Klebsiella pneumoniae</i>	850	1.2	1	–
<i>Mycobacterium tuberculosis</i>	800	5	2	–
<i>Proteus vulgaris</i>	3	0.1	0.1	–
<i>Pseudomonas aeruginosa</i>	850	2	0.4	–
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	–
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	–
<i>Staphylococcus albus</i>	0.007	0.1	0.001	+
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	+
<i>Streptococcus fecalis</i>	1	1	0.1	+
<i>Streptococcus hemolyticus</i>	0.001	14	10	+
<i>Streptococcus viridans</i>	0.005	10	40	+

<https://www.americanscientist.org/issues/pub/2009/4/thats-funny>

- Genus & species of bacteria [string]
- Antibiotic name [string]
- Gram staining? [pos/neg]
- Minimum inhibitory concentration (mg/ml) [float]
(lower == more effective)

Bacteria	Penicillin	Antibiotic Streptomycin	Neomycin	Gram stain
<i>Aerobacter aerogenes</i>	870	1	1.6	–
<i>Brucella abortus</i>	1	2	0.02	–
<i>Bacillus anthracis</i>	0.001	0.01	0.007	+
<i>Diplococcus pneumoniae</i>	0.005	11	10	+
<i>Escherichia coli</i>	100	0.4	0.1	–
<i>Klebsiella pneumoniae</i>	850	1.2	1	–
<i>Mycobacterium tuberculosis</i>	800	5	2	–
<i>Proteus vulgaris</i>	3	0.1	0.1	–
<i>Pseudomonas aeruginosa</i>	850	2	0.4	–
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	–
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	–
<i>Staphylococcus albus</i>	0.007	0.1	0.001	+
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	+
<i>Streptococcus fecalis</i>	1	1	0.1	+
<i>Streptococcus hemolyticus</i>	0.001	14	10	+
<i>Streptococcus viridans</i>	0.005	10	40	+

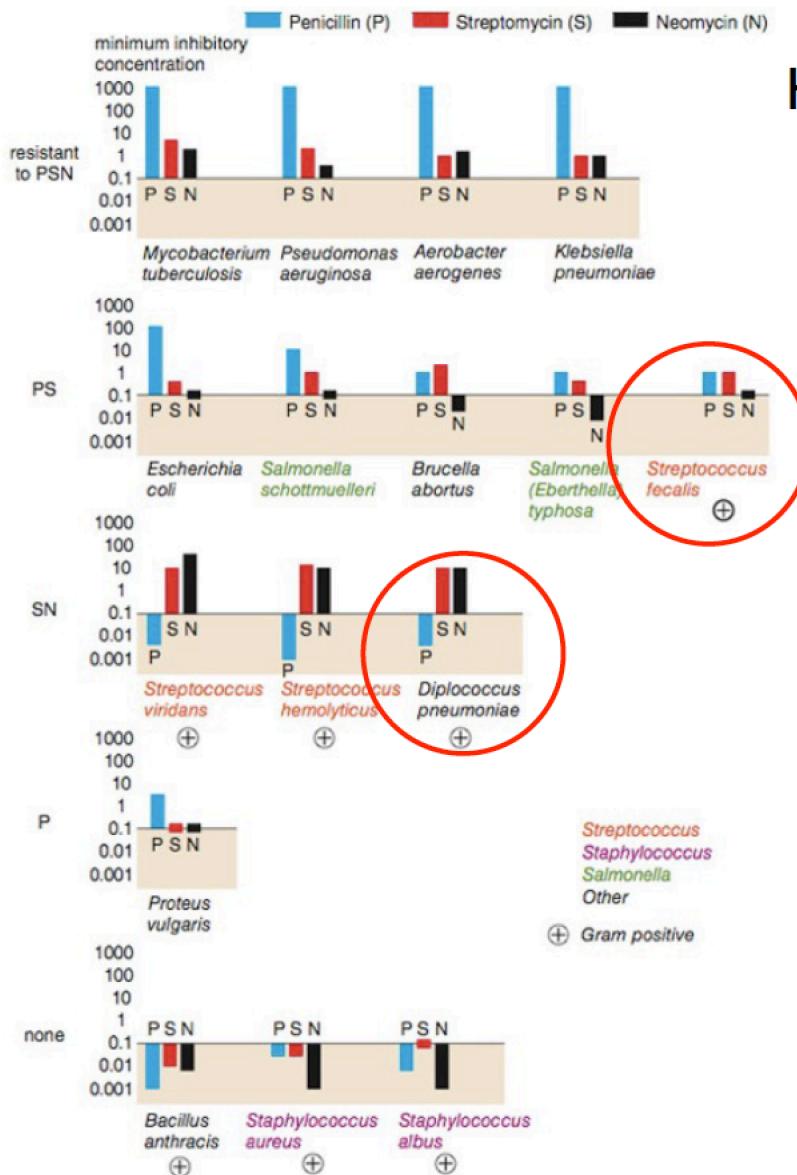
How effective are the drugs?



Not answered:
Which bacteria
behave likewise?

<https://www.americanscientist.org/issues/pub/2009/4/thats-funny>

How do the bacteria compare?

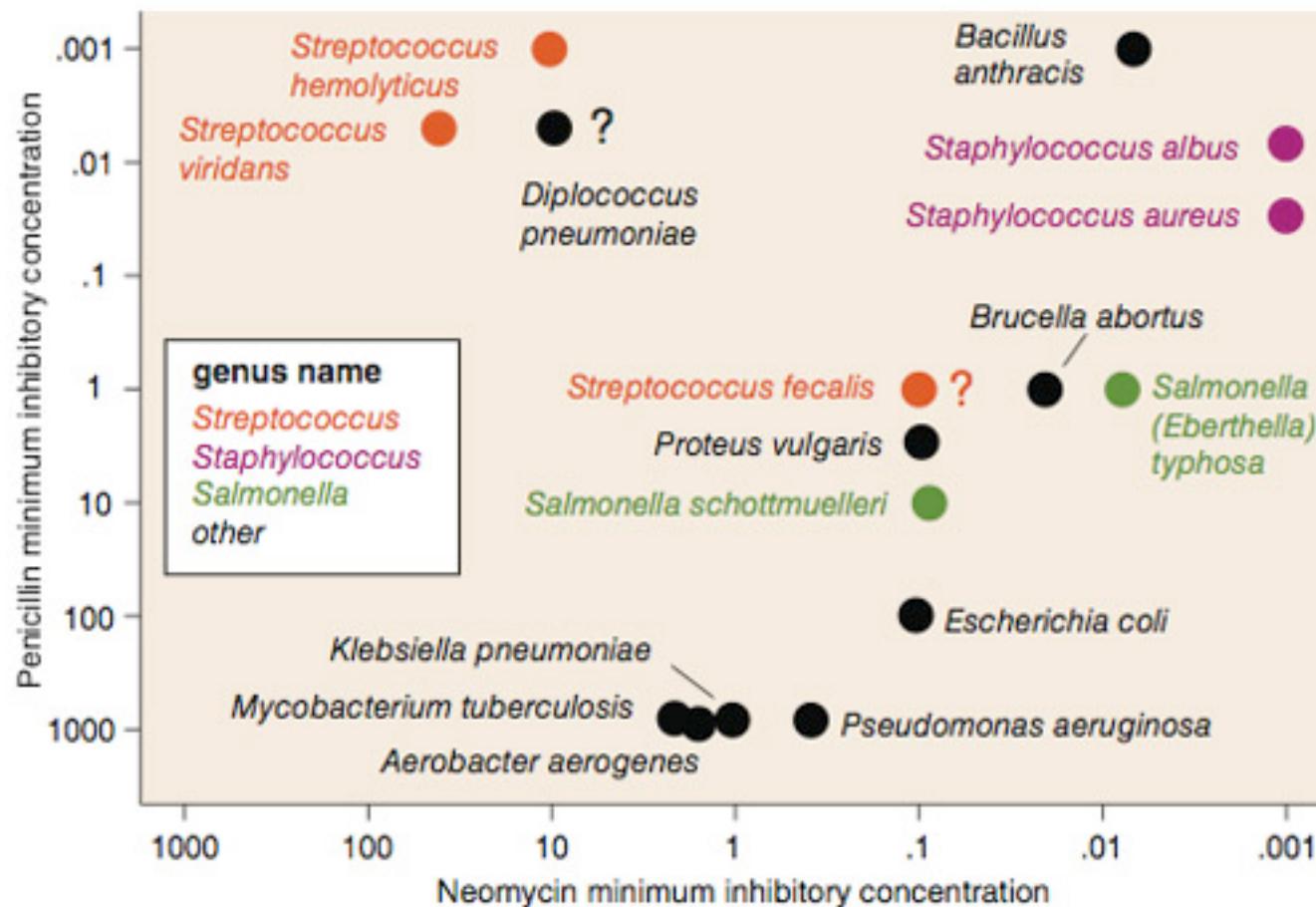


How do the bacteria compare?

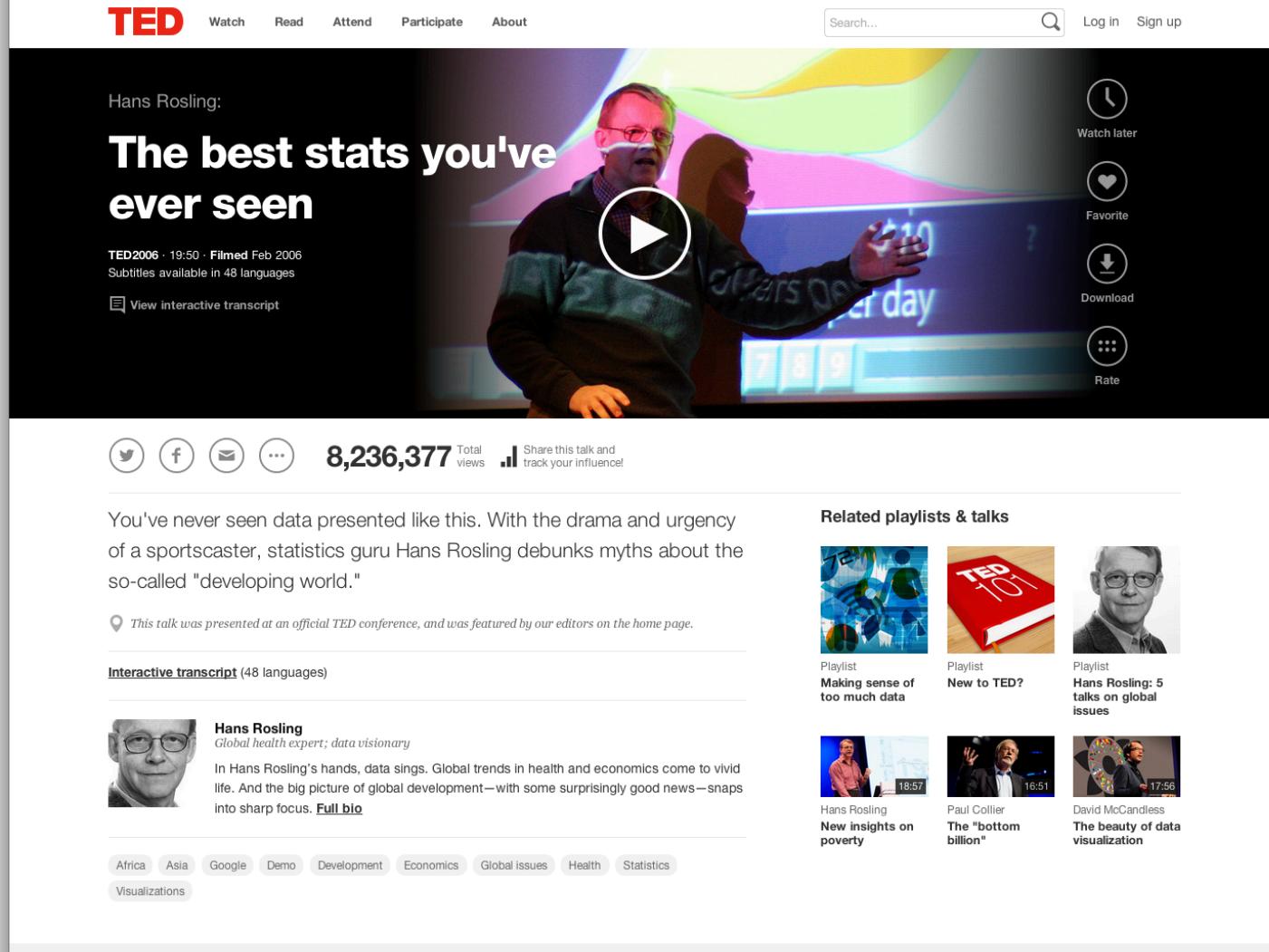
Not a streptococcus!
(realized ~30 years later)

Really a streptococcus!
(realized ~20 years later)

Wainer & Lysen, "That's funny..."
American Scientist, 2009
Adapted from Brian Schmotzer



http://www.ted.com/talks/hans_rosling_shows_the_best_stats_you_ve_ever_seen



The screenshot shows the TED talk page for Hans Rosling's presentation titled "The best stats you've ever seen". The video player displays Hans Rosling speaking in front of a colorful data visualization. The video is from TED2006, 19:50, and was filmed in February 2006. It has 8,236,377 total views. Below the video, a description reads: "You've never seen data presented like this. With the drama and urgency of a sportscaster, statistics guru Hans Rosling debunks myths about the so-called 'developing world.'". A note indicates the talk was presented at an official TED conference. An "Interactive transcript" is available in 48 languages. A bio for Hans Rosling, described as a "Global health expert; data visionary", is provided, along with links to his full bio and various TED talks. A navigation bar at the bottom includes categories like Africa, Asia, Google, Demo, Development, Economics, Global issues, Health, Statistics, and Visualizations. On the right, there are sections for "Related playlists & talks" featuring other TED talks by Hans Rosling, Paul Collier, and David McCandless.

LOOKING AT DATA

Tamara Munzner, 2013

DATASET TYPES

tables

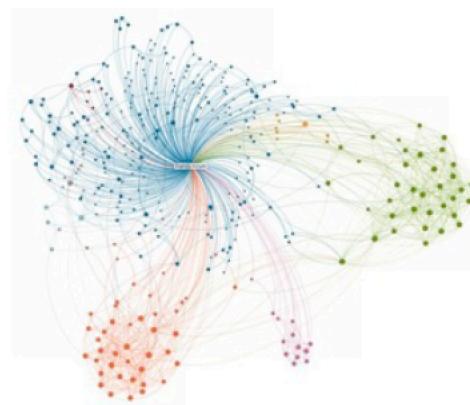
networks

text/logs

Google Docs

FriendFeed Audience

Site Name	Category	Compositio	Unique Use	Country	Re Page	Views Googl
friendfeed.com	/Online Cor	160000	150000	0.1	2000000	
twhrl.org	/Computers	47000	43000	0	74000	
tweetscan.com	/Online Cor	43000	18000	0	120000	
christbrogan.com	/Online Cor	39000	29000	0	74000	
brightblue.com	/Electronics	29000	64000	0	93000	
twitpic.com	/Home & G	24000	71000	0	340000	TRUE
web-strategist.c	/Online Cor	24000	32000	0	86000	
summize.com	/Arts & Hui	20000	54000	0	570000	



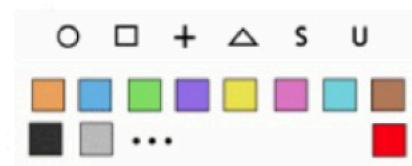
- **Data Semantics:** The real-world meaning
 - ◆ e.g., company name, day of the month, person height, etc.
- Data Type: Interpretation in terms of scales of measurements
 - ◆ e.g., quantity or category, sensible mathematical operations, data structure, etc

Data Types

- Nominal (Categorical) (N)

Are = or \neq to other values

Apples, Oranges, Bananas,...



- Ordinal (O)

Obey a $<$ relationship

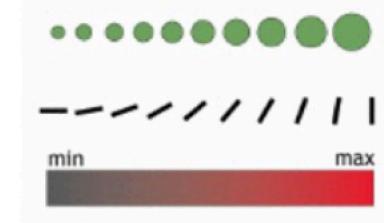
Small, medium, large



- Quantitative (Q)

Can do arithmetic on them

10 inches, 23 inches, etc.



On the theory of scales and measurements [S. Stevens, 46]

Data Types

- Q - Interval (location of zero arbitrary)
 - Dates: Jan 19; Location: (Lat, Long)
 - Like a geometric point. Cannot compare directly.
 - Only differences (i.e., intervals) can be compared
- Q - Ratio (zero fixed)
 - Measurements: Length, Mass, Temp, ...
 - Origin is meaningful, can measure ratios & proportions
 - Like a geometric vector, origin is meaningful

On the theory of scales and measurements [S. Stevens, 46]

Data Types

- N - Nominal (labels)
Operations: $=, \neq$
- O - Ordinal (ordered)
Operations: $=, \neq, >, <$
- Q - Interval (location of zero arbitrary)
Operations: $=, \neq, >, <, +, -$ (distance)
- Q - Ratio (zero fixed)
Operations: $=, \neq, >, <, +, -, \times, \div$ (proportions)

On the theory of scales and measurements [S. Stevens, 46]

Nominal

Ordinal

Interval

Ration

- Nominals: $x=y$, $x \neq y$
- Ordinals: $x=y$, $x \neq y$, $x < y$, $x > y$
- Interval scales:

Ration scales (also):

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$s_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\Delta x_t = x_t - x_{t-1}$$

$$\frac{\Delta y}{\Delta x} = \frac{y_{11} - y_{21}}{x_{11} - x_{21}}$$

Scale types with their properties according to Stanley Smith Stevens				
	Nominal scale	Ordinal scale	Interval scale	Ratio scale
Logical/ math operations	×	✗	✗	✓
	÷			
	+	✗	✗	✓
	-	✗		✓
	<	✗	✓	✓
	>		✓	✓
Examples: <i>Dichotomous and non-dichotomous</i>	<i>Dichotomous:</i> Gender (male vs. female)	<i>Dichotomous:</i> Health (healthy vs. sick),	Date (from 1457 BC to AD 2013)	Age (from 0 to 99 years)
			Latitude (from +90° to -90°)	Temperature in Kelvin (from 10K to 20K)
	<i>Variable name (data values)</i>	<i>Non-dichotomous:</i> Nationality (American/Chinese/etc)	Temperature in Celsius degrees (from 10°C to 20°C)	
		Beauty (beautiful vs. ugly)		
		<i>Non-dichotomous:</i> Opinion ('completely agree'/'mostly agree'/'mostly disagree'/'completely disagree')		
Measure of central tendency	Mode	Median Mode	Arithmetic Mean Median Mode	Geometric Mean Arithmetic Mean Median Mode
Qualitative or Quantitative	Qualitative	Qualitative	Quantitative	Quantitative

https://en.wikipedia.org/wiki/Level_of_measurement

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack		2/22/08
32	7/16/07	2-High	Small Pack		7/17/07
32	7/16/07	2-High	Jumbo Box	0.72	7/17/07
32	7/16/07	2-High	Medium Box	0.6	7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified	Wrap Bag	0.6	6/6/05
70	12/18/06	5-Low	Small Box	0.59	12/23/06
70	12/18/06	5-Low	Wrap Bag	0.82	12/23/06
96	4/17/05	2-High	Small Box	0.55	4/19/05
97	1/29/06	3-Medium	Small Box	0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08

Semantics

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box		7/17/07
32	7/16/07	2-High	Medium Box		7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified	Wrap Bag	0.6	6/6/05
70	12/18/06	5-Low	Small Box	0.59	12/23/06
70	12/18/06	5-Low	Wrap Bag	0.82	12/23/06
96	4/17/05	2-High	Small Box	0.55	4/19/05
97	1/29/06	3-Medium	Small Box	0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08

Item

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack		2/22/08
32	7/16/07	2-High	Small Pack		7/17/07
32	7/16/07	2-High	Jumbo Box		7/17/07
32	7/16/07	2-High	Medium Box		7/18/07
32	7/16/07	2-High	Medium Box		7/18/07
35	10/23/07	4-Not Specified	Wrap Bag		10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified	Wrap Bag	0.6	6/6/05
70	12/18/06	5-Low	Small Box	0.59	12/23/06
70	12/18/06	5-Low	Wrap Bag	0.82	12/23/06
96	4/17/05	2-High	Small Box	0.55	4/19/05
97	1/29/06	3-Medium	Small Box	0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08

**Attribute
aka
Feature**

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box	0.72	7/17/07
32	7/16/07	2-High	Medium Box	0.6	7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified		0.6	6/6/05
70	12/18/06	5-Low		0.59	12/23/06
70	12/18/06	5-Low		0.82	12/23/06
96	4/17/05	2-High		0.55	4/19/05
97	1/29/06	3-Medium		0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08

1 = Quantitative

2 = Nominal

3 = Ordinal

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box	0.72	7/17/07
32	7/16/07	2-High	Medium Box	0.6	7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified		0.6	6/6/05
70	12/18/06	5-Low		0.59	12/23/06
70	12/18/06	5-Low		0.82	12/23/06
96	4/17/05	2-High		0.55	4/19/05
97	1/29/06	3-Medium		0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08
194	4/5/08	3-Medium	Wrap Bag	0.84	4/7/08

I = Quantitative
2 = Nominal
3 = Ordinal

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box	0.72	7/17/07
32	7/16/07	2-High	Medium Box	0.6	7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05			44	6/6/05
69	6/4/05			0.6	6/6/05
70	12/18/06			59	12/23/06
70	12/18/06			82	12/23/06
96	4/17/05			55	4/19/05
97	1/29/06			38	1/30/06
129	11/19/08			37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08
194	4/5/08	3-Medium	Wrap Bag	0.84	4/7/08

Nominal /Ordinal = Dimensions

Describe the data, independent variables

Quantitative = Measures

Numbers to be analyzed, dependent variables

Data vs. Conceptual Model

- Data Model: Low-level description of the data
Set with operations, e.g., floats with +, -, /, *
- Conceptual Model: Mental construction
Includes semantics, supports reasoning

Data	Conceptual
1D floats	temperature
3D vector of floats	space

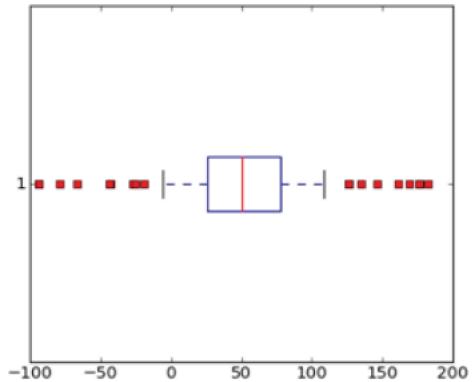
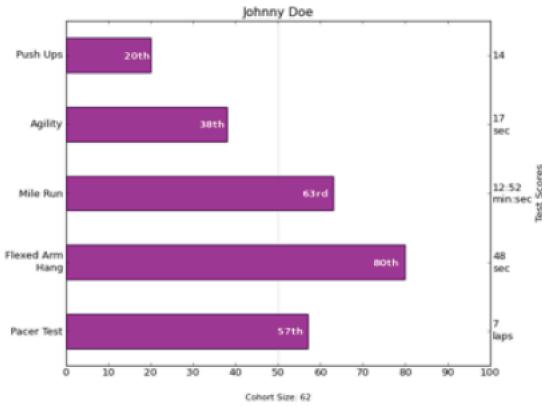
Data vs. Conceptual Model

- From data model...
32.5, 54.0, -17.3, ... (floats)
- using conceptual model...
Temperature
- to data type
Continuous to 4 significant figures (Q)
Hot, warm, cold (O)
Burned vs. Not burned (N)

Based on slide from Munzner

DATA DIMENSIONS

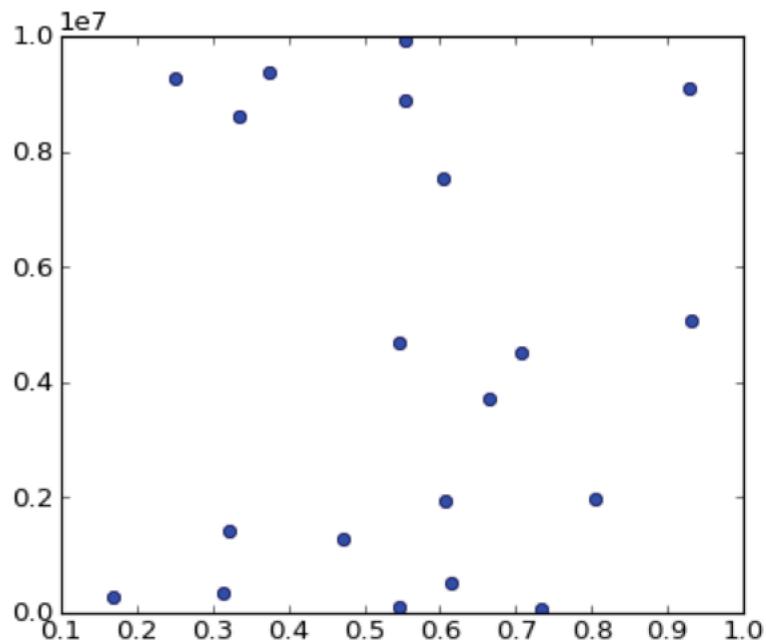
Univariate Data



Based on slide from M.Agrawala

Bivariate Data

Scatterplot is common

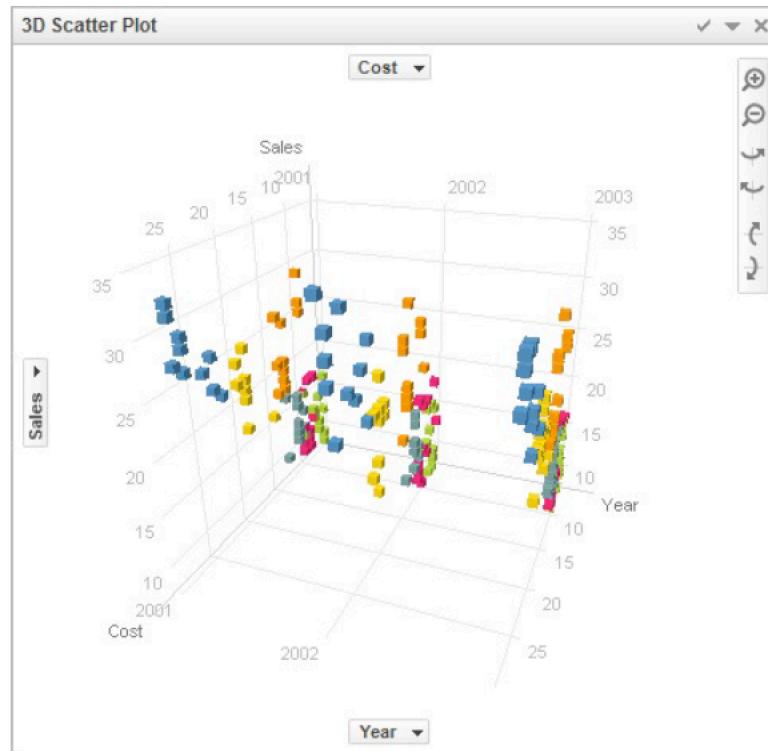


Not a function in general!

Based on slide from M.Agrawala

Trivariate Data

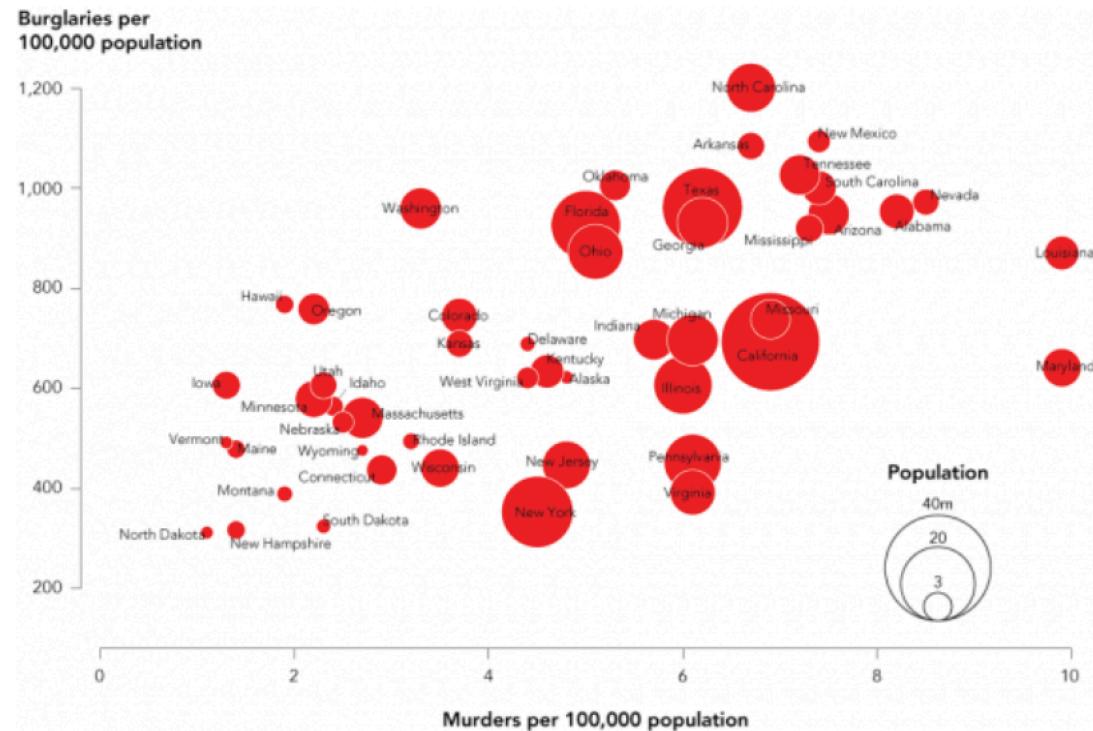
Do NOT use 3D scatterplots!



Based on slide from M.Agrawala

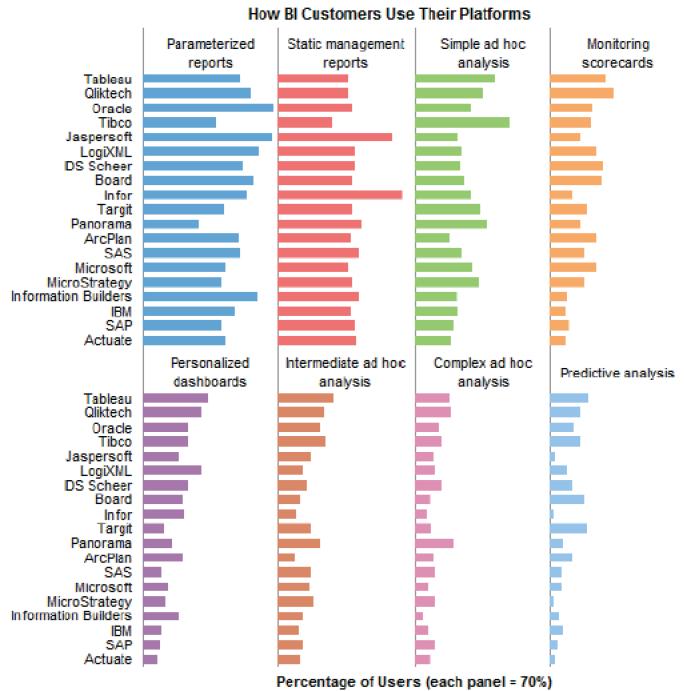
Trivariate Data

Map the third dimension to some other visual attribute



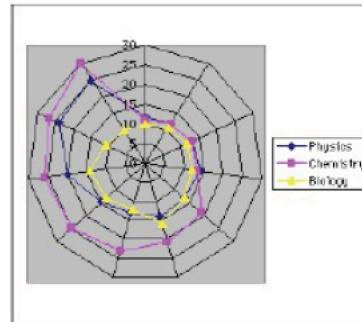
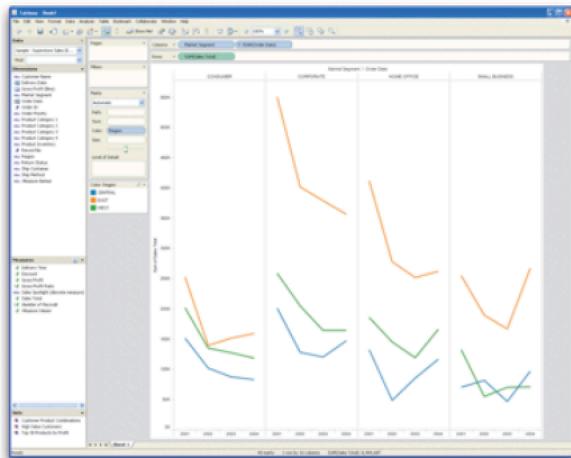
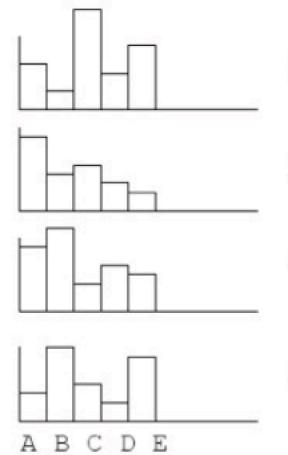
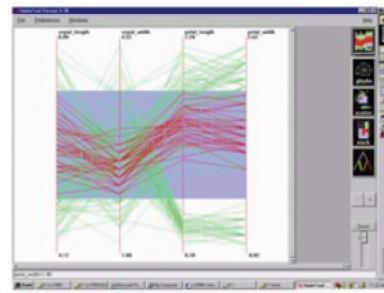
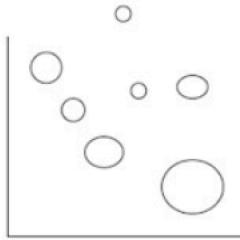
Multivariate Data

Give each attribute its own display (small multiples)



Based on slide from M.Agrawala

Multivariate Data Representations



- Filtering: Eliminate some items or attributes
 - ◆ e.g., select range of interest, zoom in, remove outliers, etc.
- Aggregation: Represent a group of elements by a new derived element
 - ◆ e.g., take average, min, max, count, sum
 - ◆ Attribute aggregation a.k.a. dimensionality reduction