

Introduction to Data Science



Introduction to Data Science

Instructor: Patrick Boily
Slides: P. Boily (IACS, DAL, uOttawa), M. Kashef (datascience2go), J. Schellinck (Sysabee, DAL, AI Guides)

uOttawa
Institut de développement professionnel
Professional Development Institute

1

Introduction to Data Science

Instructor: Patrick Boily

Slides: P. Boily (IACS, DAL, uOttawa), M. Kashef (datascience2go), J. Schellinck (Sysabee, DAL, AI Guides)

2

Outline

Module 1
Data Insight Fundamentals

Module 2
Data Collection and Data Management

Module 3
Data Visualization and Data Communication

Module 4
Data Processing and Data Cleaning

Module 5
Data Exploration and Data Analysis

Module 6
Data Mining and Machine Learning

3

Instructor

Bio

- Prof. Math/Stat ['19 – now, uOttawa]
- Manager and Senior Consultant ['12 – '19, CQADS, Carleton]
- Lecturer ['99 – '19, uOttawa | UQO | Carleton]
- Public Service ['08 – '12, ASFC | StatCan | TC | TPSGC]
- 60+ uni course; 250+ workshop days

Projects

- GAC; NWMO; CATSA; etc.
- 40+ projects

Specialization

- Data visualization; data cleaning (... unfortunately)
- Application of wide breadth of techniques to all kinds of data



4

Instructor: Patrick Boily

Introduction to Data Science

Suggested References

Data Understanding, Data Analysis, and Data Science
P. Boily
idlewyldanalytics.com

Data Science Basics (suggested exercise: #4)
Data Preparation (suggested exercise: #4)
Data Visualization & Data Exploration (sugg. ex: #7)
Machine Learning 101 (suggested exercise: #18)



@ IACS (2022)

Roundtable

?

Quick Intro

Experience

Why this course?

5

6

Module 1

Data Insight Fundamentals



7

"Reports that say that something hasn't happened are always interesting to me, because as we know, there are **known knowns**; there are things we know that we know. There are **known unknowns**; that is to say, there are things that we now know we don't know. But there are also **unknown unknowns** – there are things we do not know we don't know."

Donald Rumsfeld, US Department of Defense News Briefing, 2002

8

<https://archive.ics.uci.edu/ml/datasets/Mushroom>

Poisonous Mushroom Dataset

Amanita muscaria

Habitat: woods
Gill Size: narrow
Odor: none
Spores: white
Cap Colour: red

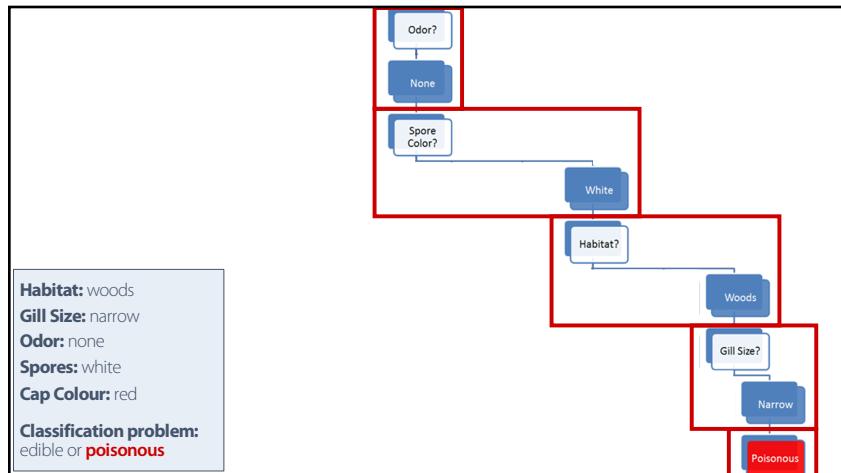
Classification problem:
Is Amanita muscaria edible, or poisonous?



9



10



11

Discussion

Would you have trusted an “**edible**” prediction?
Where is the model coming from?
What would you need to know to trust the model?
What’s the cost of making a classification mistake, in this case?

12

Solve Problems Create Meaningful Change Support 'Gut Checks'

Asking the Right Questions

The diagram illustrates various types of questions framed in rounded rectangles. Some are highlighted with red borders, while others have black borders. Arrows point from one question to another, suggesting a flow or relationship between them. The questions include:

- Is this an image of a cat or a dog?
- Will the customer click this link?
- What topics are described in this article?
- What's the sentiment of this tweet?
- Is this credit card transaction suspicious?
- Is this insulin reading unusual?
- What will the temperature be next Friday?
- What will sales for next quarter be?

datascience2go

13

Source: kdnuggets

Roadmap to Framing Questions

Understand the problem (opportunity vs problem)
What initial assumptions do I have about the situation?
How will the results be used?
What are the risks and/or benefits of answering this question?
What stakeholder questions might arise based on the answer(s)?
Do I have access to the data necessary to answering this question?
How will I measure my 'success' criteria?

datascience2go

14

Exercise: Roadmap to Framing Questions

Possible Initial Question 1: Should I buy a house? (vague)

Possible Initial Question 2: Should I buy a single house in Scotland?

datascience2go

15

Additional Rules

Avoid **glazing over the data** before you settle on the question.
You can be **blinded by love**; you can be **blinded by solutions**.
Do you **fully understand** what you're asking?

Source: kdnuggets

datascience2go

16

Source: Healthy Families BC

Yes/No Trap

Examples of **bad** questions:

- Are our revenues **increasing** over time? **Has it** increased year-over-year?
- Are most of our customers from **this demographic**?
- **Does this project have** valuable ambitions to the broader department?
- **How great** is our hard-working customer success team?
- How often do you **triple check** your work?

Examples of **good** questions:

- What's the **distribution** of our revenues over the past three months?
- Where are our **top 5** high-spending cohorts from?
- What are the **different benefits** of pursuing this project?
- What are **three good** and **bad traits** of our customer success team?
- Do you **tend to** do quality assurance testing on your deliverables?

datascience2go

17

Question Audit Checklist

1. Did I avoid creating any **yes/no** questions?
2. Would **everyone** in my team/department understand the question, regardless of their backgrounds?
3. Does the question need more than one sentence to express?
4. Is the question '**balanced**' – is the scope **so broad** that the question will never truly be answered; **so narrow** that the resulting impact is minimal?
5. Is the question being **skewed to what may be easier to answer** for my team's particular skillset(s)?

Source: The Head Game

datascience2go

18

Are these good questions?

Question	Specific?	What's the range in answers to this question?
How does rain affect goal percentage at a soccer match?	No, could be any soccer field	Could completely vary based on location, teams, level of players
Did the Toronto Maple Leafs beat the Edmonton Oilers ?		
Did you like watching the Tokyo Olympics ?		
What types of recovery drinks do hockey players drink?		
How many medals will Canada achieve at the Paris 2024 Olympics ?		
Should we fund the Canadian Basketball team <i>more</i> than the Canadian Hockey team?		

datascience2go

19

What is Data Science?

Data science is the collection of processes by which we extract useful and **actionable** insights from data.

T. Kwartler (paraphrased)

Data science is the **working intersection** of statistics, engineering, computer science, domain expertise, and "hacking." It involves two main thrusts: **analytics** (counting things) and **inventing new techniques** to draw insights from data.

H. Mason (paraphrased)

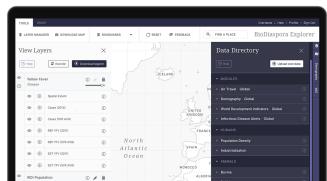
20

Case Study: BlueDot

Digital health company that tracks the spread of infectious diseases **globally** and assesses their risk of spread and impact worldwide.

Using advanced data science techniques, they've built a **global early warning system for infectious diseases**:

- Mapping out 200+ diseases 24/7, processing 100,000+ articles a day in over 65 languages
- Understanding impact of the spread
- Alerting clients to better inform policy decisions



Source: BlueDot

datascience2go

21

Analytics Modes

Analytics can be broken down into four core **key buckets**:

Descriptive



Show **what** happened

Diagnostic



Explain **why** something happened

Predictive



Guess **what will** happen

Prescriptive



Suggest **what should** happen

datascience2go

22

Data Science Ecosystem

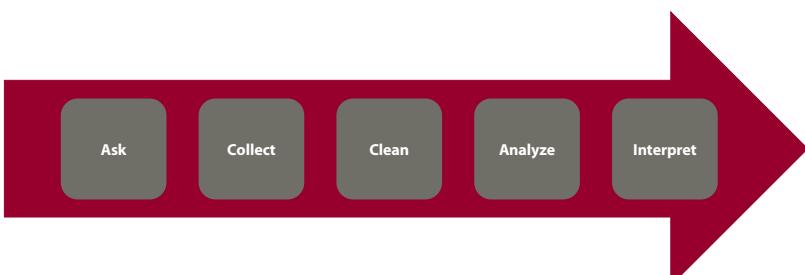
Data analysis is a **team sport**, with team members needing a good understanding of both **data** and **context**

- data management
- data preparation
- analysis
- communications

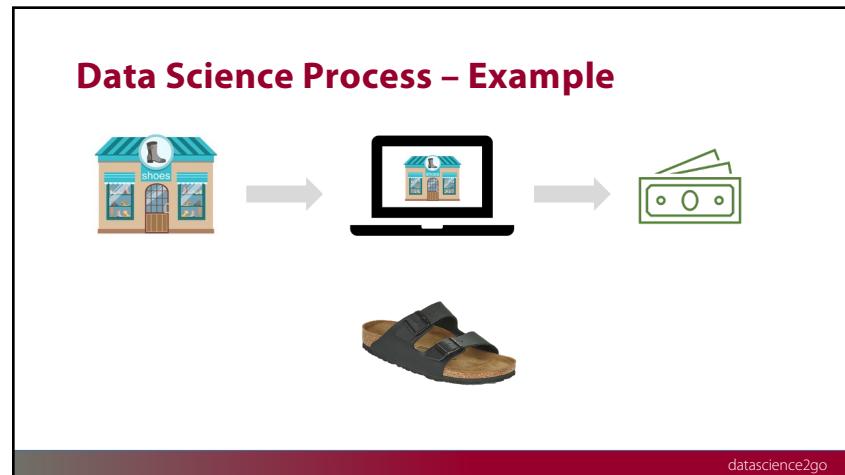
Even **slight improvements** over a current approach can find a useful place in an organization – data science is not solely about **Big Data**, disruption, the “Cloud”, etc!

23

Data Science Workflow



24



25



26

Example: Data Science Process

COLLECT

This slide shows the 'COLLECT' phase, displaying a list of transactions from a database. The transactions are categorized by status (Refund or Payment) and include details like customer names, payment methods, and creation dates.

27

Example: Data Science Process

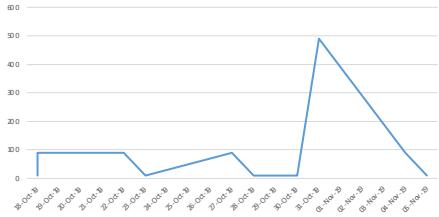
CLEAN

This slide shows the 'CLEAN' phase, where the raw transaction data has been processed and presented in a clean, structured table format. The original color-coding is removed, and the data is organized by transaction status.

28

Example: Data Science Process

ANALYZE



datascience2go

29

Example: Data Science Process

INTERPRET



datascience2go

30

Representations

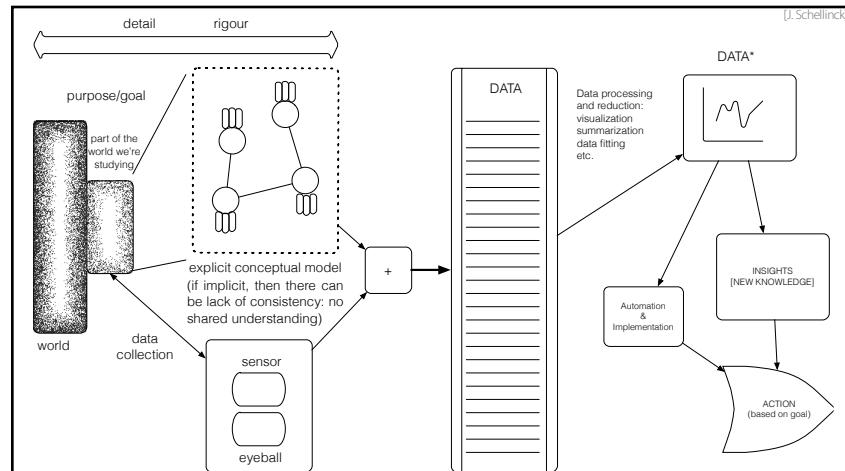
A **representation** is an object that stands in for another object.

A representation may or may not physically resemble the object it represents.

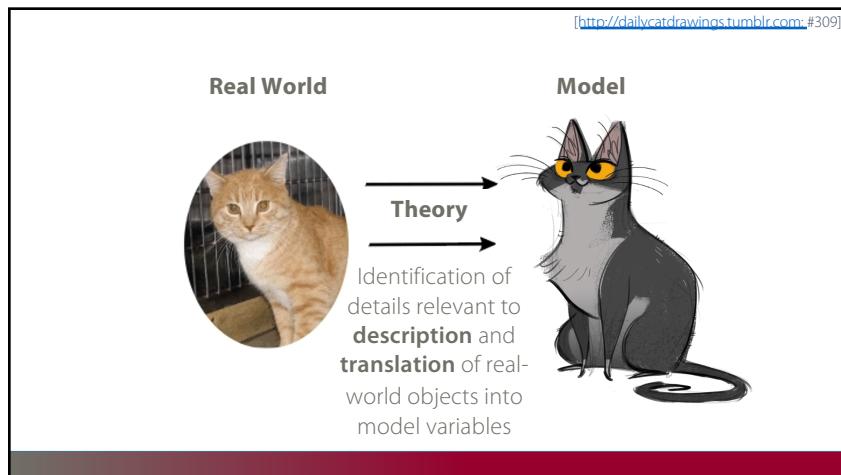
Representations of the world help us to **understand**, **navigate**, and **manipulate** the world.



31



32



33

Systemic Thinking Take-Aways

Systems can approximate certain aspects of the Universe.

System models provide the basis under which data is identified and collected, but data itself is **approximate** and **selective**.

Knowledge gaps happen – be ready to re-visit your set-up regularly.

Implicit conceptual modeling can lead to problematic situations.

If the data, the system, and the world are **out of alignment**, data analysis insights might ultimately prove useless.

34

What are Ethics?

"Ethics" refers to the **study** and **definition** of **right** and **wrong** conducts:

- "not [...] social convention, religious beliefs, or laws". (R.W. Paul, L. Elder)

Influential ethical theories:

- Kant's **golden rule** (do onto others...), **consequentialism** (the ends justify the means), **utilitarianism** (act in order to maximize positive effect), etc.
- **Confucianism, Taoism, Buddhism** (?), etc.
- **Ubuntu, Maori, OCAP**, etc.

Discussion: What harm can come from data?

35

Ethics in the Data Context

Data ethics questions:

- **Who**, if anyone, owns data?
- Are there **limits** to how data can be used?
- Are there **value-biases** built into certain analytics?
- Are there categories that should **not** be used in analyzing personal data?
- Should some data be **publicly available** to **all** researchers?

Analytically, the **general** is preferred to the anecdotal – decisions made based on machine learning and AI. (security, financial, marketing, etc.) may affect real beings in **unpredictable ways**.

36

Best Practices

"Do No Harm": data collected from an individual **should not be used to harm** the individual.

Informed Consent:

- Individuals must **agree to the collection and use** of their data
- Individuals must have a **real understanding of what they are consenting to**, and of **possible consequences** for them and others

Respect "Privacy": excessively hard to maintain in the age of constant trawling of the Internet for personal data.

37

Best Practices

Keep Data Public: data should be kept **public** (all? most? any?).

Opt-In/Opt-Out: Informed consent requires the ability to **opt out**.

Anonymize Data: removal of id fields from data prior to analysis.

"Let the Data Speak":

- no cherry picking
- importance of validation (more on this later)
- correlation and causation (more on this later, too)
- repeatability

38

Gapminder Exercises

We will conduct the exercises using Gapminder Tools.

The online version is available at <https://www.gapminder.org/tools/> [there is also an offline version].

Take some time to explore the tool. In the online version, the default starting point is a bubble chart of 2020 life expectancy vs. income, per country (with bubble size associated with total population). In the offline version, select the "Bubbles" option.

Do the exercises for Module 1.

39

Module 2 Data Collection and Data Management



40

What is Data?

4,529 'red' 25.782 'Y'

41

Objects and Attributes



Object: apple
Shape: spherical
Colour: red
Function: food
Location: fridge
Owner: Jen

A person or an object is **not simply the sum of its attributes!**

42

From Attributes to Datasets

Attributes are **fields** (columns) in a database; objects are **instances** (rows).

Objects are described by their **feature vector**, the collection of attributes associated with value(s) of interest.

ID#	Shape	Colour	Function	Location	Owner
1	spherical	red	food	fridge	Jen
2	rectangle	brown	food	office	Pat
3	round	white	tell time	lounge	School
...

43

Data is Real



Data is a representation, but data is still **physical**.

It has physical properties.

Physical space and energy are required to process and work with it.

44

Data Decay

Data ages over time – it has a **shelf life**.

We use the phrase “rotten data” or “decaying data”

- **literally** – the data storage medium might decay
- **metaphorically** – when the data no longer accurately represents the relevant objects and relationships or even when those objects no longer exist in the same way

Data must be kept ‘fresh’ and ‘current’, not ‘stale’
(context and model dependent!)



45

[G. Smith, [The Exaggerated Promise of So-Called Unbiased Data Mining](#)]

“A Dartmouth graduate student used an MRI machine to study the brain activity of a salmon as it was shown photographs and asked questions. The most interesting thing about the study was not that a salmon was studied, but that **the salmon was dead**. Yep, a dead salmon purchased at a local market was put into the MRI machine, and some patterns were discovered. There were inevitably patterns—and they were invariably meaningless.”

46

What's a Sample?

A **sample** is a portion of a ‘population’ from which the data is collected

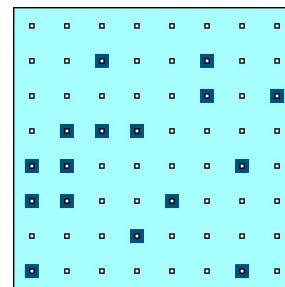
Biased	Unbiased
One or more parts of the population are favoured over others	Everyone has an equal chance of being chosen
Does not accurately represent the population	Accurately represents the population
Leads to invalid conclusions	Provides a valid conclusion

Source: Making Big Data Work

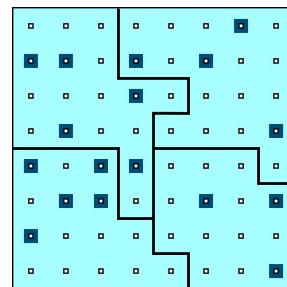
datascience2go

47

Sampling Designs



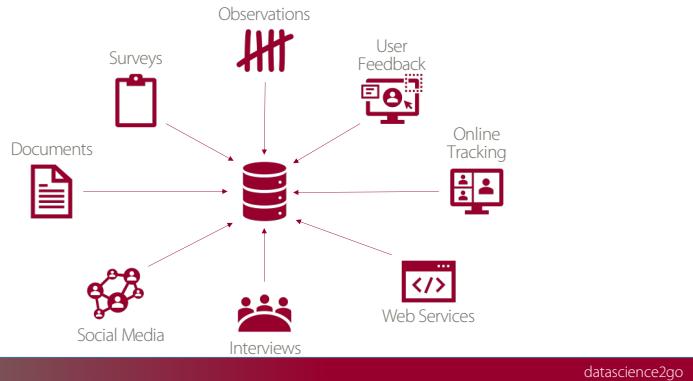
Simple Random Sampling (SRS)



Stratified Random Sampling (STS)

48

Collect/Create Data



49

Web Scraping – Example

Let's say you want to know what people think of a new phone.

Standard approach: market research (e.g. telephone survey, reward system, etc.).

Pitfalls:

- unrepresentative sample: the selected sample might not represent the intended population
- systematic non-response: people who don't like phone surveys might be less (or more) likely to dislike the new phone
- coverage error: people without a landline can't be reached, say
- measurement error: are the survey questions providing suitable info for the problem at hand?

50

Web Scraping – Example

These solutions can be **costly, time-consuming, ineffective**.

Proxies are indicators that are strongly related to the information of interest, without measuring it directly.

If **popularity** is defined as large groups of people preferring one product over a competitor, then sales statistics on a commercial website may provide a proxy for popularity.

Rankings on Amazon could provide a **more comprehensive** view of the phone market than a traditional survey.

51

Web Scraping – Example

Representativeness of the listed products

- are all phones listed?
- if not, is it because that website doesn't sell them?
- is there some other reason?

Representativeness of the customers

- are there specific groups buying/not-buying online products?
- are there specific groups buying from specific sites?
- are there specific groups leaving/not-leaving reviews?

Truthfulness of customers and **reliability** of reviews.

52

Scraping Dos and Don'ts

1. Stay identifiable
 2. Reduce traffic
 3. Do not bother server with multiple requests
 4. Write modest scrapers (efficient and polite)

Use **application programming interface** (APIs) as much as possible!

53

Fundamental Concepts

It is important to structure **data** and **knowledge** so that it can be:

- stored and accessible
 - added to/amended
 - usefully and efficiently extracted from that store (extract – transform – load)
 - operated over by **humans** and **computers** (programs, bots, A.I.)

Different options are used in terms of fundamental **data and knowledge** modeling or structuring strategies:

- key-value pairs (e.g., JSON)
 - triples (e.g., RDF – resource description framework)
 - graph databases
 - relational databases

55

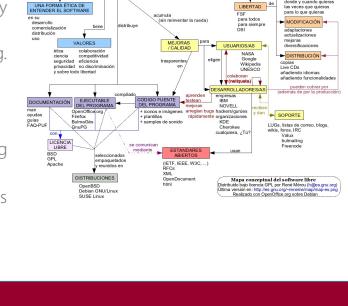
Conceptual Model

A **conceptual model** is, roughly speaking:

- a model that is not implemented, which exists only conceptually
 - a diagram or verbal description of a system (e.g. boxes and arrows, mind maps, lists, definitions)

Focus is :

- not on capturing specific behaviors but emphasizing **possible states**
 - on object types, not on specific instances; the goal is **abstraction**.



54

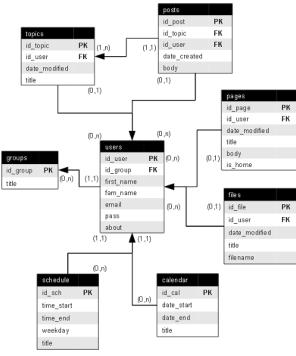
Data Modeling

Data models are **abstract/logical** descriptions of a system, constructed in terms that can then be implemented as the structure of a type of data management software.

This is half-way between a conceptual model and a database implementation.

The data itself is about **instances** – the model is about the **object types**.

Another option to consider: **ontologies**.

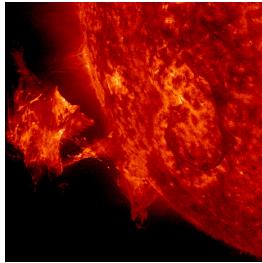


56

Structured/Unstructured Data

A major motivator for new developments in database types and data storing strategies is the increasing availability of **unstructured** data and 'blob' data

- **structured data:** labeled, organized, discrete structure is constrained and pre-defined
- **unstructured data:** not organized, no specific pre-defined structure data model (text)
- **blob data:** Binary Large Object (BLOB) – images, audio, multi-media



57

Flat Files and Spreadsheets

What about keeping data in a single giant table (spreadsheet)?

Or multiple spreadsheets?

How bad can it be?

Wayne Eckerson coined the term 'spreadmart' to describe a situation with many (ad hoc) spreadsheets as a data strategy.

Date	Con	Lab	Lzn	SMP	Ukp	Grosses	Con av	Lab av	Lzn av	SMP av	Ukp av	Gross av	
15 September 2017	41	41	5	4	5	3	40.7	41.4	6.8	3.3	4	2.7	
15 September 2017	38	38	8	3	6	4	40.7	41.7	7	3.2	3.6	2.6	
15 September 2017	42	42	7	3	4	4	40.9	42.2	6.4	3.3	3.5	2.6	
15 September 2017	38	40	7	3	4	4	40.9	42.3	7	3.2	3.5	2.6	
1 September 2017	38	40	7	3	1	4	40.9	42.3	7	3.2	3.4	2.3	
31 August 2017	23	23	8	3	4	3	40.7	41.4	6.8	3.3	4	2.7	
22 August 2017	19	19	10	4	5	3	40.7	41.4	6.8	3.3	4	2.7	
18 August 2017	15	15	30	8	3	4	40.7	41.7	7	3.2	3.8	2.6	
11 August 2017	15	15	30	8	3	4	40.7	41.7	7	3.2	3.8	2.6	
6 July 2017	15	15	30	8	3	4	40.9	42.2	7	3.2	3.5	2.6	
6 July 2017	15	15	30	8	3	4	40.9	42.2	7	3.2	3.5	2.6	
5 July 2017	15	15	30	8	3	4	40.9	42.2	7	3.2	3.5	2.6	
5 July 2017	15	15	30	8	3	4	40.9	42.3	7	3.2	3.4	2.6	
30 June 2017	15	15	30	8	3	4	40.9	42.3	7	3.2	3.4	2.6	
29 June 2017	15	15	30	8	3	4	40.9	42.3	7	3.2	3.4	2.6	
16 July 2017	20	20	20	7	2	3	3	41	42.2	7	3.1	4	2.6
20 August 2017	41	42	8	4	4	1	40.8	42.5	7	3.3	3.9	2.6	
15 August 2017	41	42	8	4	4	2	40.5	42.8	6.8	3.3	3.9	2.6	
11 July 2017	41	42	8	4	4	2	40.5	42.8	6.8	3.3	3.9	2.6	
6 July 2017	41	42	7	3	4	3	40.5	43.0	6.9	3.2	3.4	2.6	
5 July 2017	41	42	7	3	4	3	40.5	43.0	6.9	3.2	3.4	2.6	
30 June 2017	41	42	7	3	3	2	40.5	43.1	6.7	3.2	3.6	2.6	
29 June 2017	41	42	7	3	3	2	40.5	43.1	6.7	3.2	3.6	2.6	
16 July 2017	20	20	20	7	2	3	3	41	42.2	7	3.1	4	2.6
15 August 2017	41	42	8	4	4	1	40.8	42.5	7	3.3	3.9	2.6	
11 July 2017	41	42	8	4	4	2	40.5	42.8	6.8	3.3	3.9	2.6	
6 July 2017	41	42	7	3	3	2	40.5	43.0	6.9	3.2	3.4	2.6	
5 July 2017	41	42	7	3	3	2	40.5	43.0	6.9	3.2	3.4	2.6	
30 June 2017	41	42	7	3	3	2	40.5	43.1	6.7	3.2	3.6	2.6	
29 June 2017	41	42	7	3	3	2	40.5	43.1	6.7	3.2	3.6	2.6	
16 July 2017	20	20	20	7	2	3	3	41	42.2	7	3.1	4	2.6
15 August 2017	39	41	8	3	6	1	40.5	43.0	6.4	3.1	3.4	2.6	
11 July 2017	40	40	7	3	2	1	40.5	43.1	6.4	3.1	3.4	2.6	
6 July 2017	40	40	7	3	2	1	40.5	43.1	6.4	3.1	3.4	2.6	
5 July 2017	40	40	7	3	2	1	40.5	43.1	6.4	3.1	3.4	2.6	
30 June 2017	40	40	7	3	2	1	40.5	43.1	6.4	3.1	3.4	2.6	
29 June 2017	40	40	7	3	2	1	40.5	43.1	6.4	3.1	3.4	2.6	
16 July 2017	20	20	20	7	2	3	3	41	42.2	7	3.1	4	2.6
15 August 2017	39	41	8	3	6	1	40.5	43.0	6.4	3.1	3.4	2.6	
11 July 2017	40	40	7	3	2	1	40.5	43.1	6.4	3.1	3.4	2.6	
6 July 2017	40	40	7	3	2	1	40.5	43.1	6.4	3.1	3.4	2.6	
5 July 2017	40	40	7	3	2	1	40.5	43.1	6.4	3.1	3.4	2.6	
30 June 2017	40	40	7	3	2	1	40.5	43.1	6.4	3.1	3.4	2.6	
29 June 2017	40	40	7	3	2	1	40.5	43.1	6.4	3.1	3.4	2.6	

58

Database Management

Once data has been collected, it must also be **managed**.

Fundamentally, this means that the database must be maintained, so that the data is

- accurate,
- precise,
- consistent
- complete

Don't let your data lake turn into a data swamp!

59

Tools and Buzzwords

SQL, SQLite, MySQL, NoSQL

MongoDB, ArangoDB

Document store

JSON, YAML

API, GraphQL

Linked Data

Semantic Web

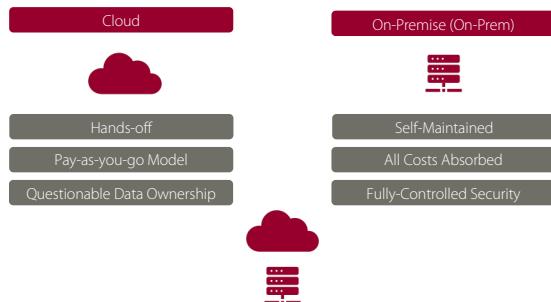
Ontology Web Language (OWL)

Protégé

etc.

60

Cloud vs. On-Premise



61

Roundtable: About Your Data



Does it exist?

Where does it live?

How is it structured and accessed?

datascience2go

62

Module 3 Data Visualization and Data Communication

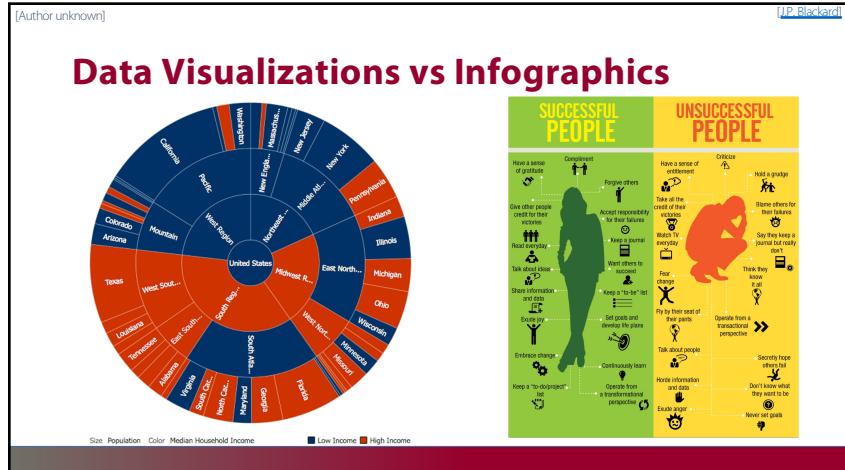


63

Gapminder Exercises

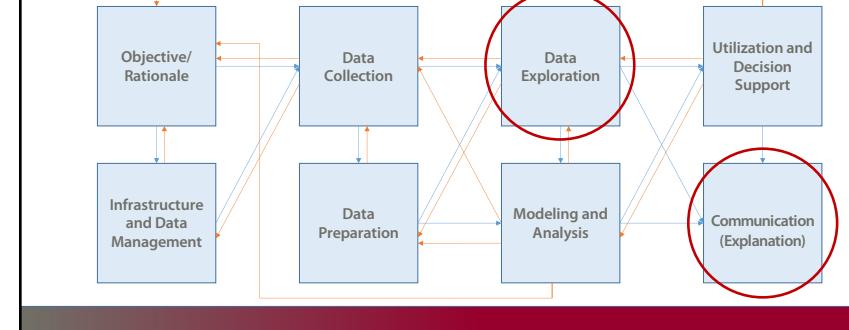
Do the exercises for Module 2.

64



65

The (Messy) Analysis Process



66

Pre-Analysis Uses

Data visualization can be used to set the stage for analysis:

- **detecting anomalous entries**
invalid entries, missing values, outliers
- **shaping the data transformations**
binning, standardization, Box-Cox transformations, PCA-like transformations
- **getting a sense for the data**
data analysis as an art form, exploratory analysis
- **identifying hidden data structure**
clustering, associations, patterns informing the next stage of analysis

67

Fundamental Principles of Data Viz

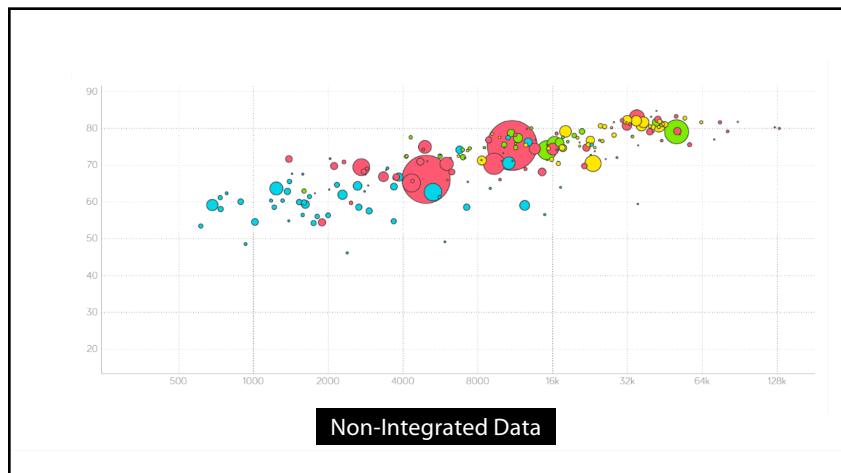
There is a **symmetry** to visual displays of evidence. Consumers should be seeking exactly what producers should be providing, namely:

- meaningful comparisons
- potential causal networks and underlying structure
- multivariate links
- integrated and relevant data
- honest documentation
- primary focus on content

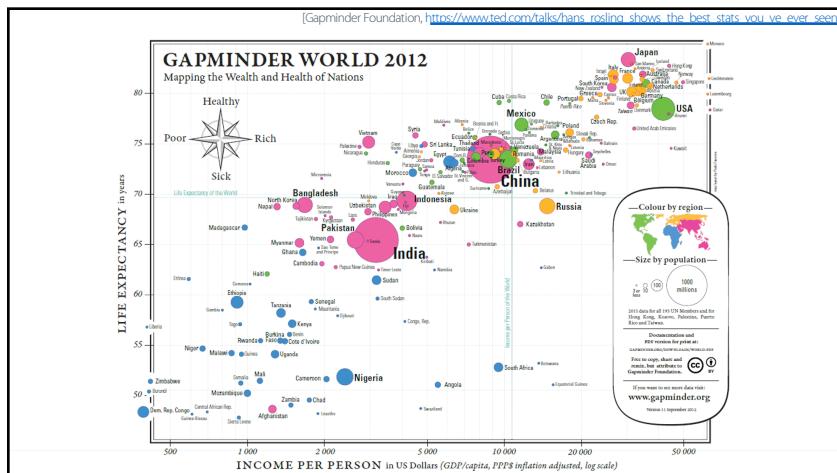
Graphics should be **clear** and **engaging**.

Don't be afraid to try something new if it helps **convey the message**.

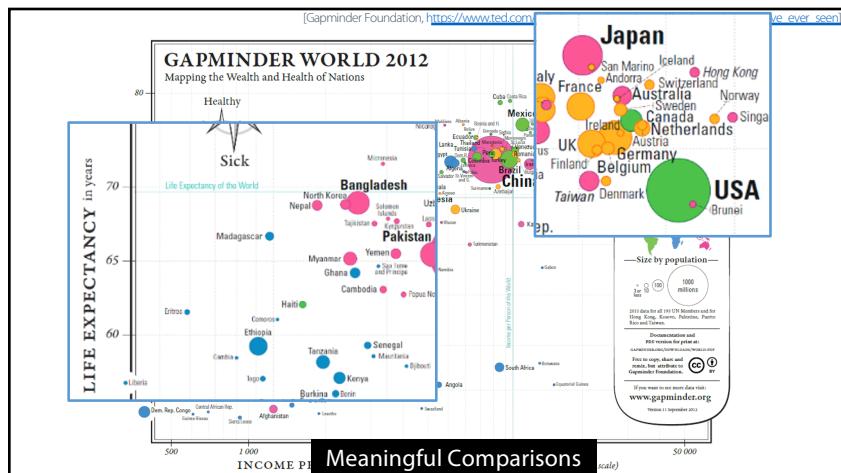
Introduction to Data Science



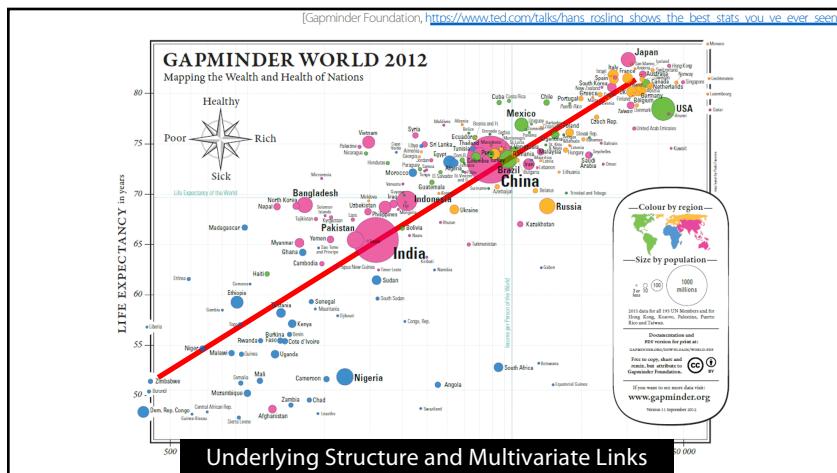
69



70

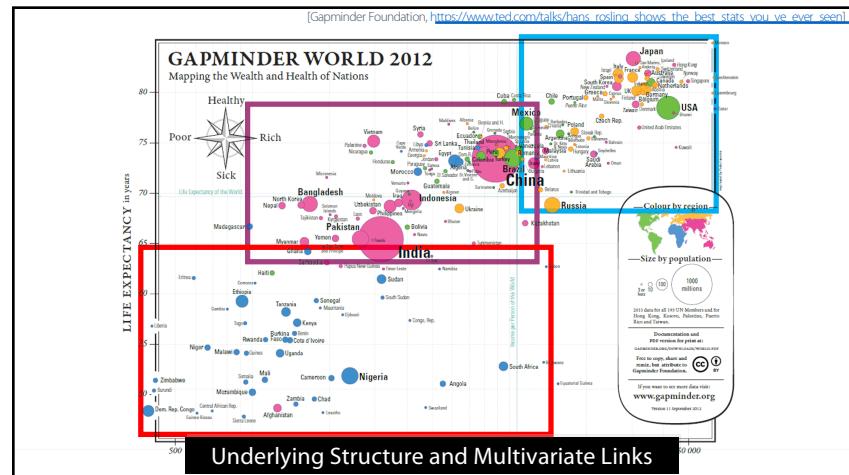


71

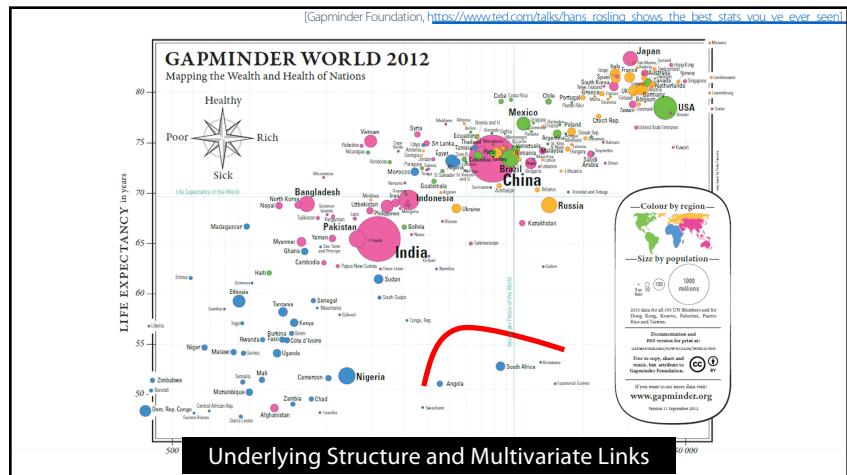


72

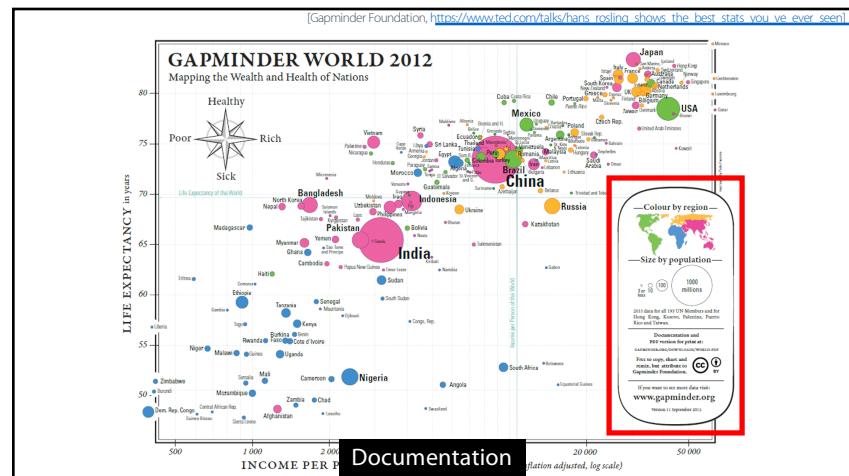
Introduction to Data Science



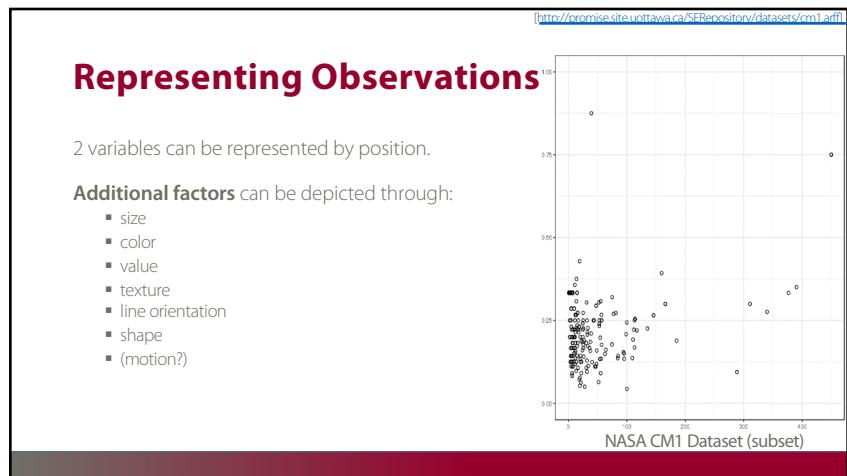
73



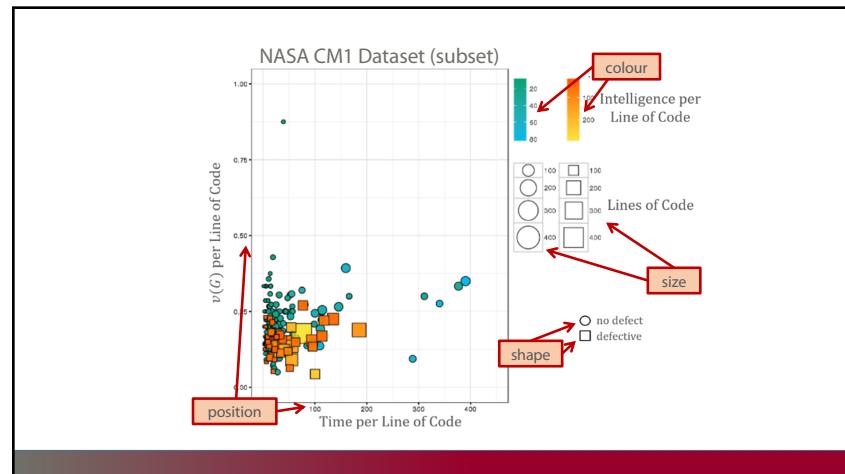
74



75



76



77

<http://dataphys.org/list/>

A Word About Accessibility

Charts cannot usually be translated to Braille. Describing the features and emerging structures in a visualization is a possible solution... **if they can be spotted.**

Analysts must produce clear and meaningful visualizations, but they must also describe them and their features in a fashion that allows all to "see" the insights. This requires analysts to have "seen" all the insights, which is not always possible.

Conditions: colourblindness, low vision, motor impairment, cognitive disability, ADHD, etc.

Best Practices: high contrast elements, zoom/magnifications, keyboard navigation, assistive design, short summaries, un/re-do functionality, text-to-voice, etc. [Elavsky]

78

A Word About Accessibility

Data Perception:

- texture-based representations
- text-to-speech
- sound/music
- odor-based or taste-based representations (?!?)

Sonifications:

- [TRAPPIST Sounds : TRAPPIST-1 Planetary System Translated Directly Into Music](#)
- [Listening to data from the Large Hadron Collider, L. Asquith](#)

79

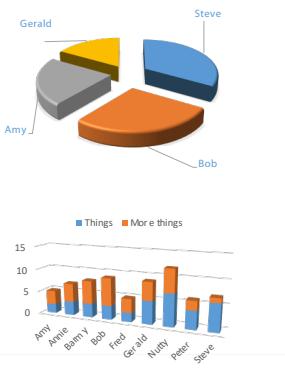
Chart Types

Simple text and tables	Geographical maps
Scatterplot	Parallel coordinates
Line chart	Chernoff faces
Bar charts	Word clouds
Stacked bar charts	Network diagrams
100% bar charts	Dendograms and trees
Area charts	Sparklines
Treemaps	Interactive charts
Gauge charts	Small multiples
Heatmaps and choropleth maps	etc.

80

Charts to Avoid

AVOID (?) anything with an arc (except gauge charts): pie, donut, etc. Human brains have a hard time **comparing arcs** – which is larger, Steve or Bob?



81

Decluttering

CLUTTER IS THE ENEMY!

- every element on a page adds **cognitive load**
- identify anything that isn't adding value and **remove**
- think of cognitive load as mental effort required to process information (lower is better)
- Tufte refers to the **data-to-ink** ratio – “the larger the share of a graphic's ink devoted to data, the better”
- in Resonate, Duarte refers to this as **“maximizing the signal-to-noise ratio”** where the signal is the information or the story we want to communicate.

82

Decluttering

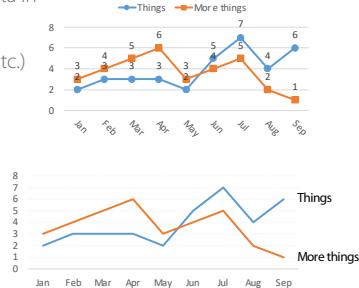
Use **Gestalt Principles** to organize/highlight data in a chart.

Align all the elements (graphs, text, lines, titles, etc.)

- DONT rely on eye, use position boxes and values

Charts:

- remove border, gridlines, data markers
- clean up axis labels
- label data directly



83

Decluttering

Use **consistent** font, font size, colour and alignment.

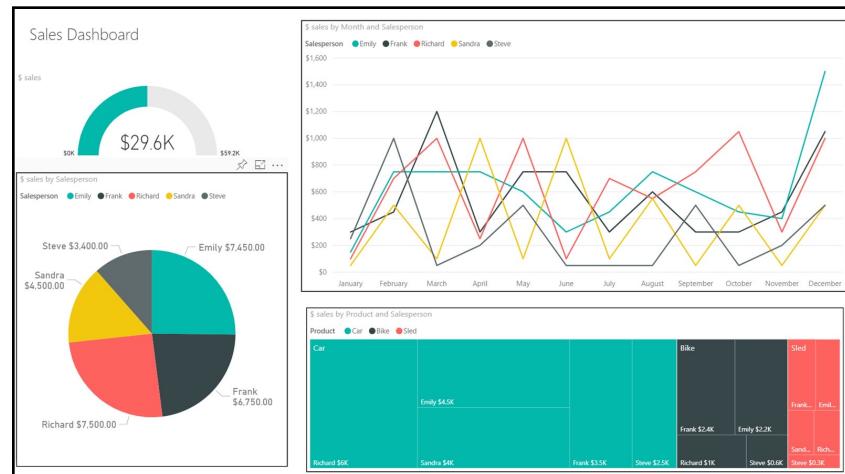
Don't rotate text to anything other than 0 or 90 degrees.

Use white space

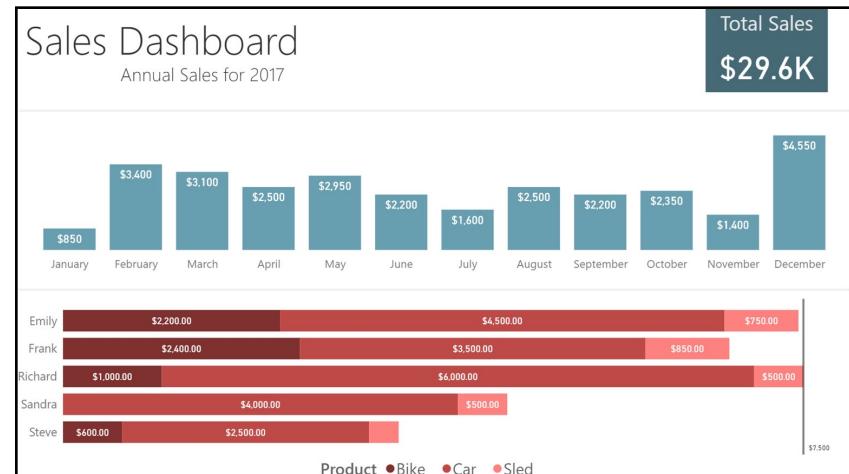
- margins should remain free of text and visuals
- don't stretch visuals to edge of page or too close to other visuals
- think of white space as a border

84

Introduction to Data Science



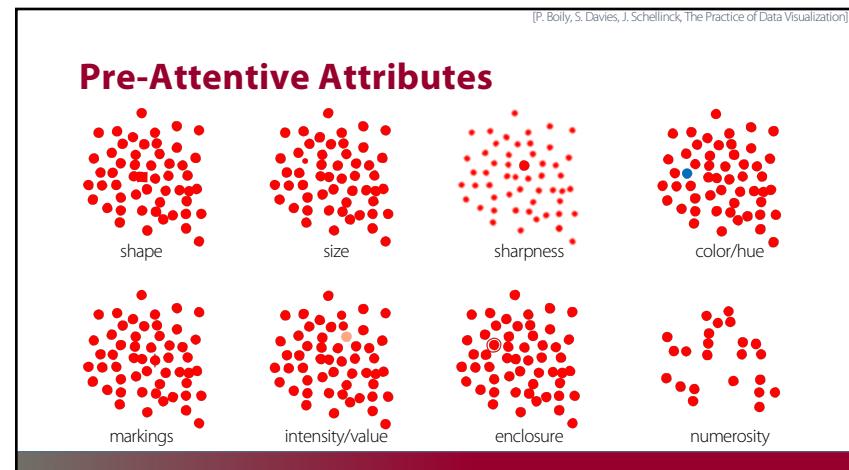
85



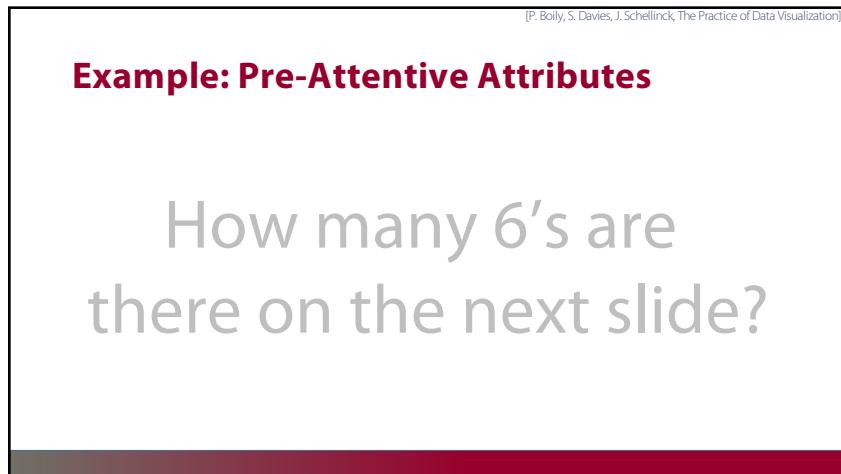
86



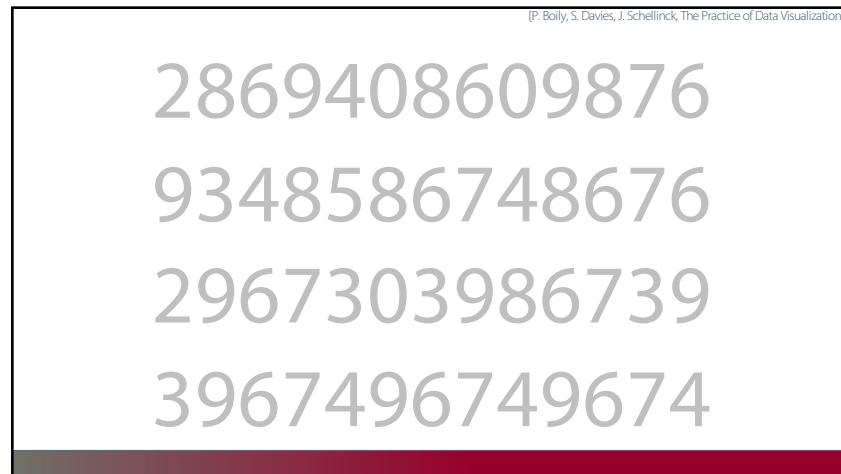
87



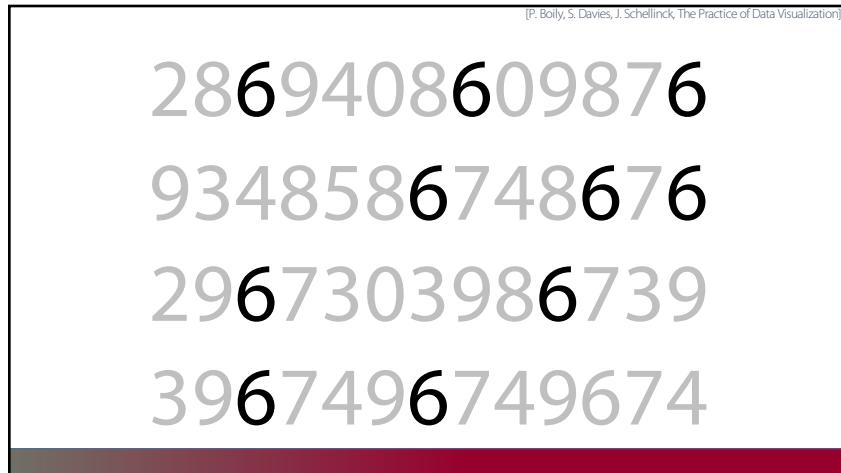
88



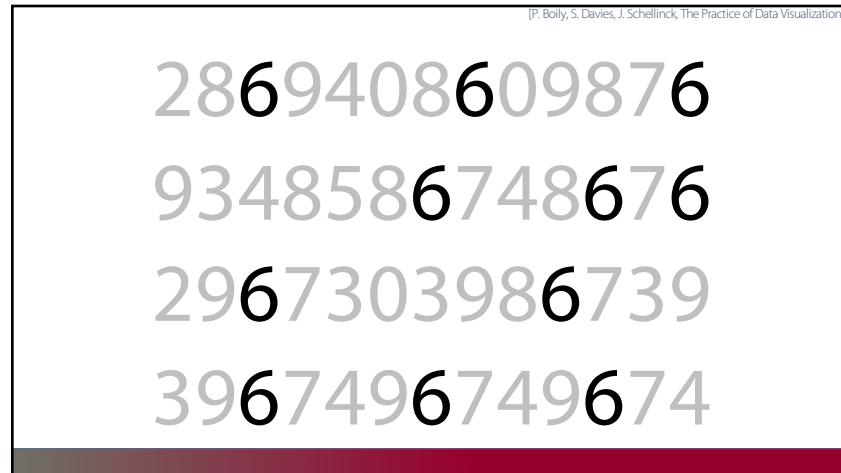
89



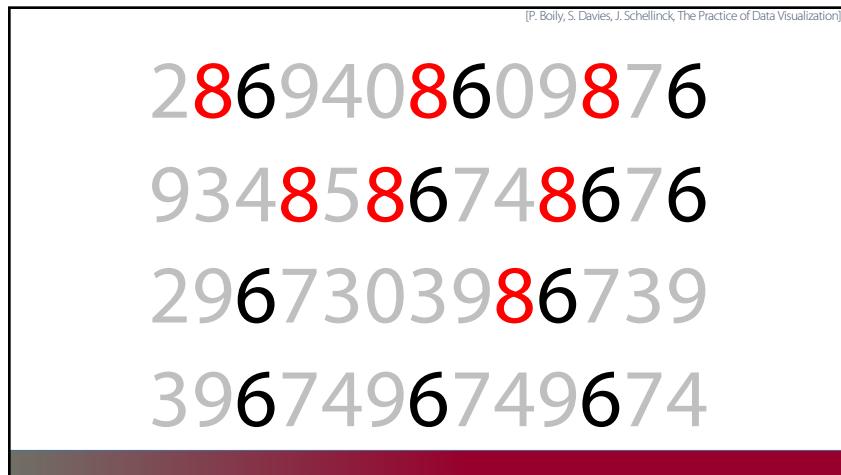
90



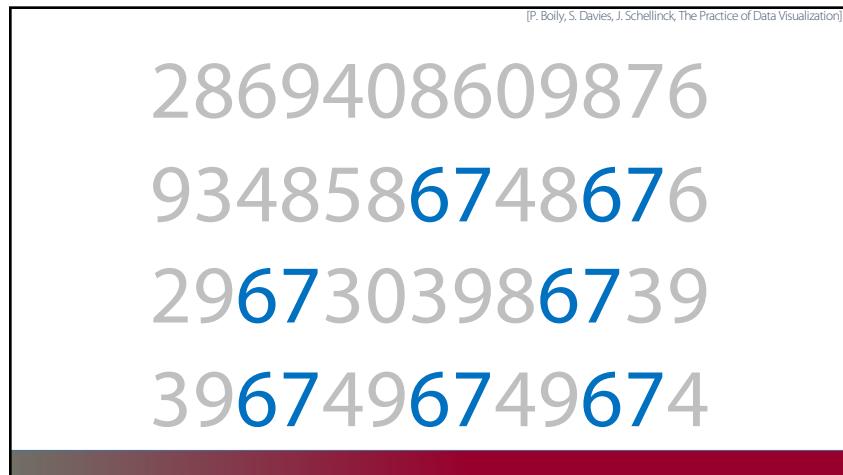
91



92



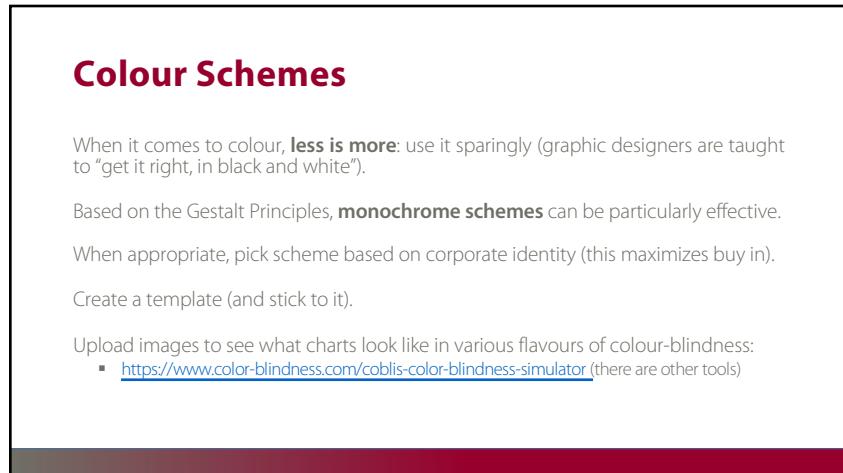
93



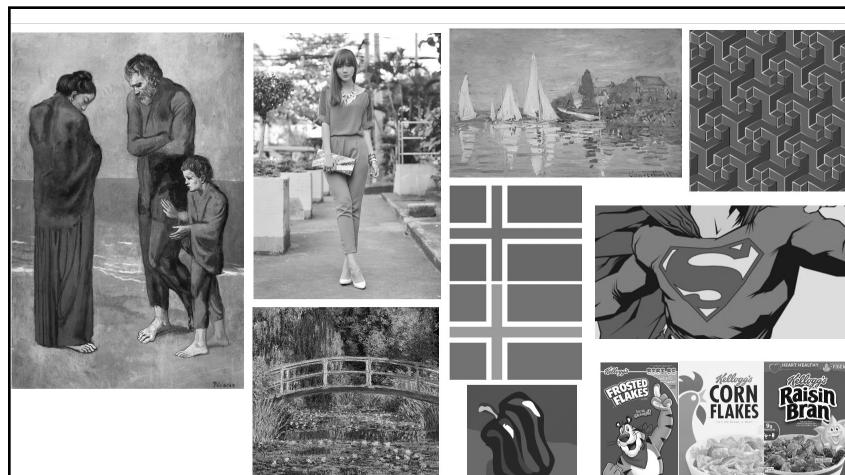
94



95



96



97

Mastering Colour

Country Level Sales Rank Top 5 Drugs

Rainbow distribution in color indicates sales rank in given country from #1 (red) to #10 or higher (dark purple)

Country	A	B	C	D	E
AUS	1	2	3	6	7
BRA	1	3	4	5	6
CAN	2	3	6	7	8
CHI	1	2	3	4	7
FRA	3	2	4	5	6
GER	3	1	6	5	4
IND	4	5	3	10	9
ITA	2	4	9	6	5
MEX	1	5	4	6	3
PER	4	3	7	9	12
SPA	2	3	4	5	11
TUR	7	2	3	4	8
UK	1	2	3	6	7
USA	1	2	4	3	5

Top 5 drugs: country-level sales rank

RANK	1	2	3	4	5+
COUNTRY DRUG	A	B	C	D	E
Australia Paracetamol	1	2	3	6	7
Brazil Aspirin	1	3	4	5	6
Canada Tylenol	2	3	6	12	8
China Ibuprofen	1	2	8	4	7
France Aspirin	3	2	4	8	10
Germany Paracetamol	3	1	6	5	4
India Paracetamol	4	1	8	10	5
Italy Aspirin	2	4	10	9	8
Mexico Paracetamol	1	5	6	6	5
Russia Paracetamol	1	3	7	9	12
Spain Paracetamol	2	3	4	5	11
Turkey Paracetamol	7	2	3	4	8
United Kingdom Paracetamol	1	2	3	6	7
United States Paracetamol	1	2	4	3	5

Source: Storytelling with Data

98

Evolving a Visualization



Source: Storytelling with Data

99

Evolving a Visualization

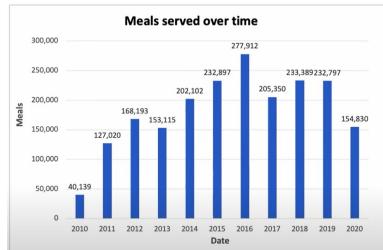
Meals served over time

Campaign Year	Meals Served
2010	40,139
2011	127,020
2012	168,192
2013	153,115
2014	202,102
2015	232,897
2016	277,912
2017	205,350
2018	233,389
2019	232,797
2020	154,830

Source: Storytelling with Data

100

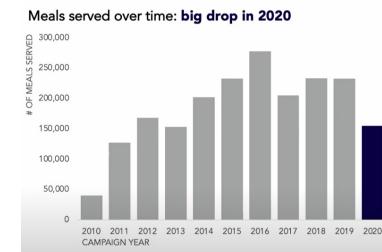
Evolving a Visualization



Source: Storytelling with Data

101

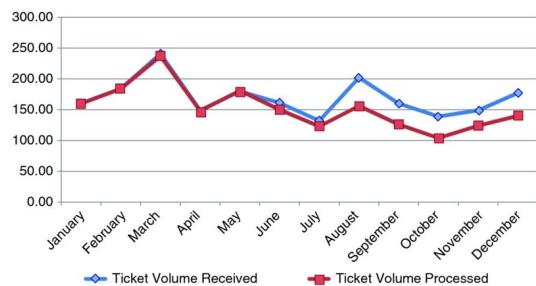
Evolving a Visualization



Source: Storytelling with Data

102

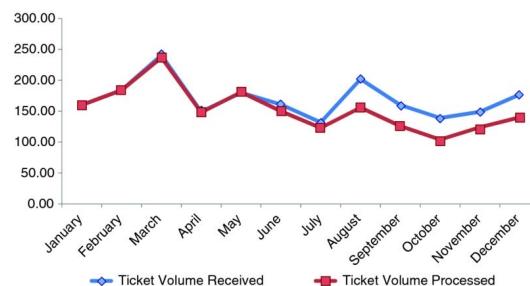
Decluttering – Step-by-Step Example



Source: Storytelling with Data

103

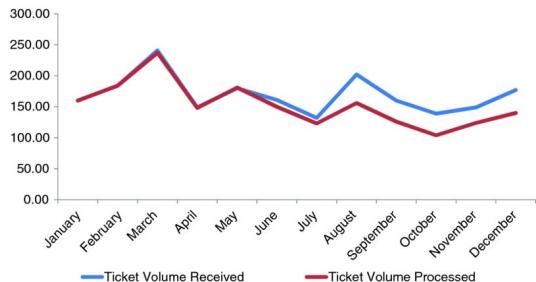
1. Remove Chart Border & Gridlines



Source: Storytelling with Data

104

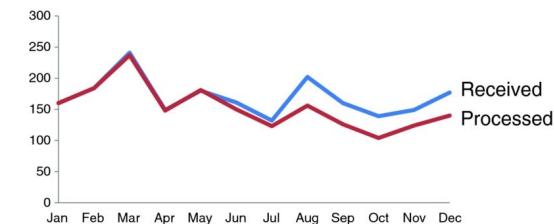
2. Remove Data Markers



Source: Storytelling with Data

105

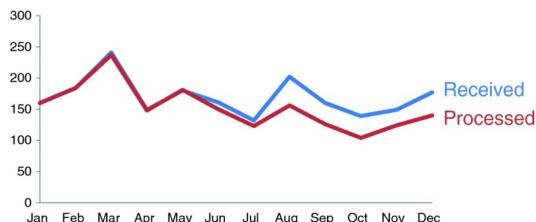
3. Clean Up Axis Labels



Source: Storytelling with Data

106

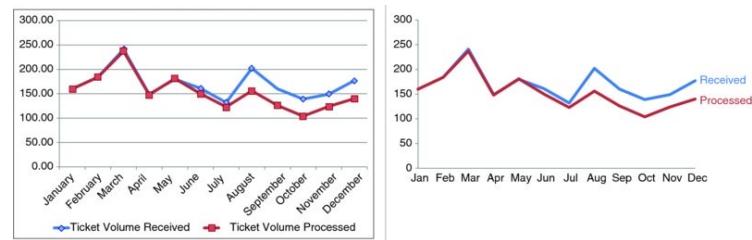
4. Colour Code the Lines



Source: Storytelling with Data

107

5. Before & After



Source: Storytelling with Data

108

Dashboards

A **dashboard** is any visual display of data used to monitor conditions and/or facilitate understanding.

In a car's dashboard, a small number of **key indicators** (speed, gasoline level, lights, etc.) need to be understood **immediately**.

A dashboard design that does not take these two characteristics under consideration can have **catastrophic consequences**. The same is true for data dashboards.



109

Dashboards Best Practices

The most amount of time someone will spend on a dashboard is **10-15 minutes**

- no more than 7-8 pages per dashboard (fewer is better)

Short-term memory makes it difficult to see **more than 4 visual chunks** at once

- no more than 4-5 charts on a single page (fewer is better)

Pre-attentive features can help **direct the eye**

- each chart should have 1 iconic memory trigger

Long-term memory is more easily triggered by a **combination** of words and visuals

- explain: tell us, in a few words, what we are supposed to be seeing

110

Exercise

Consider the following dashboards.

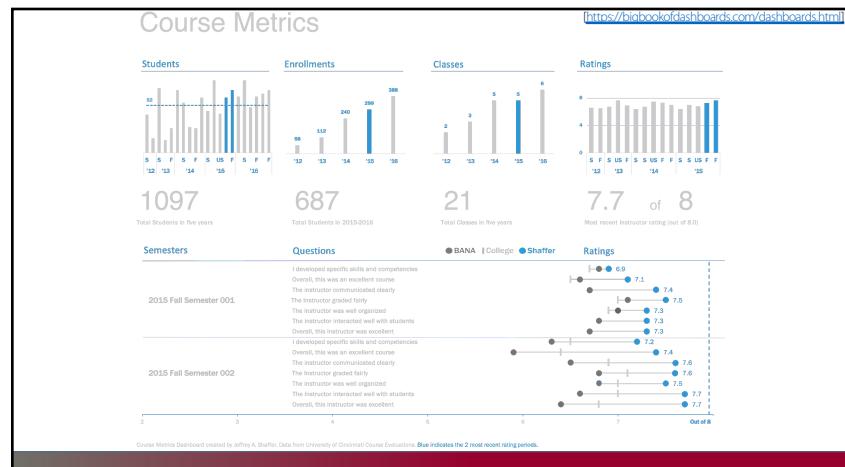
Can you figure out, at a glance, who their audience is?

What are their strengths?

What are their limitations?

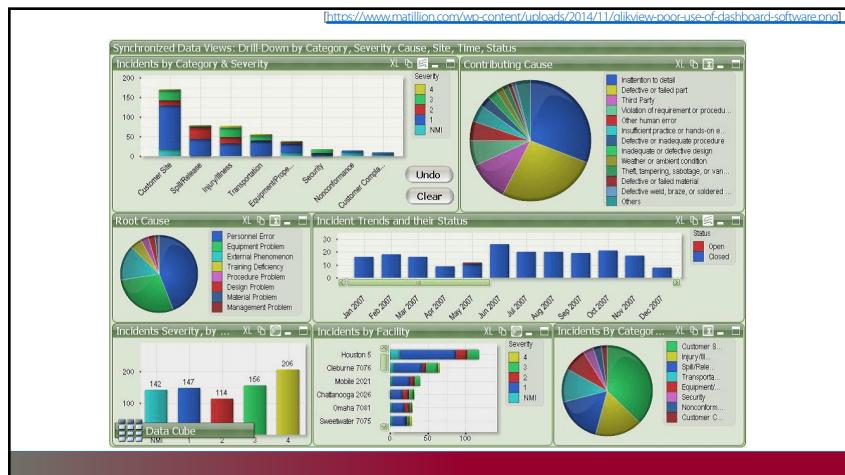
How would you improve them?

111

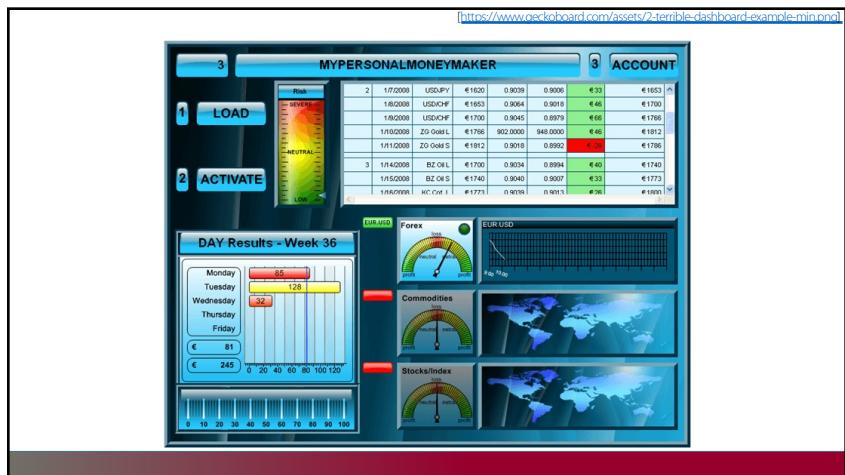


112

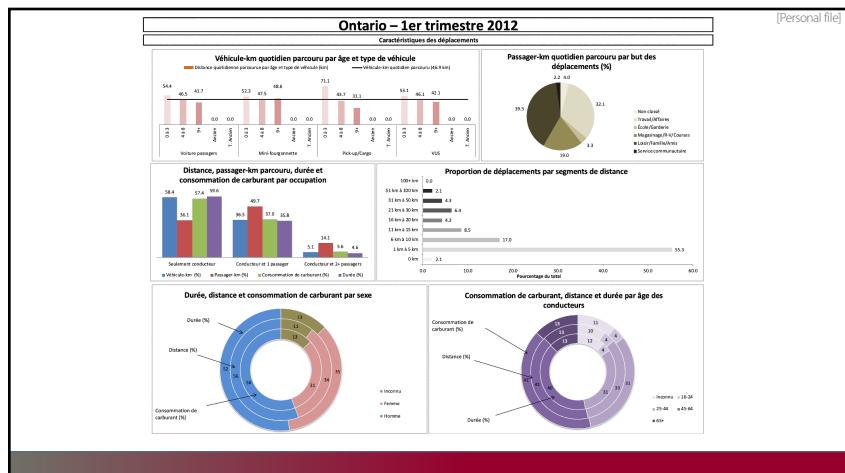
Introduction to Data Science



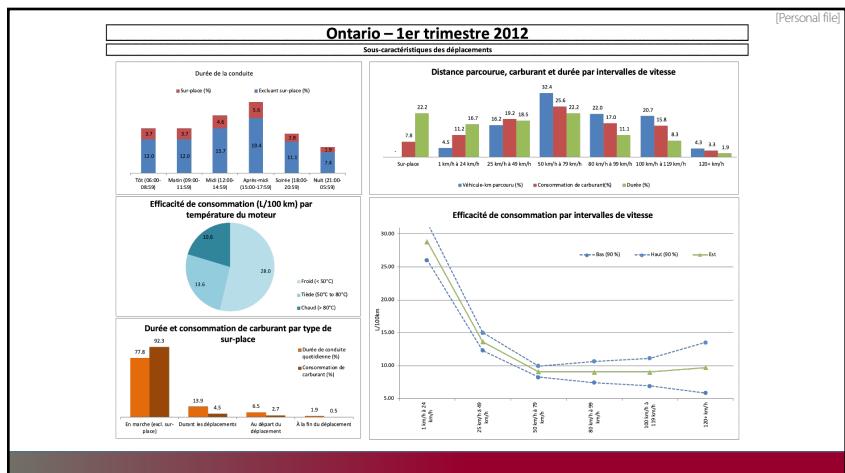
113



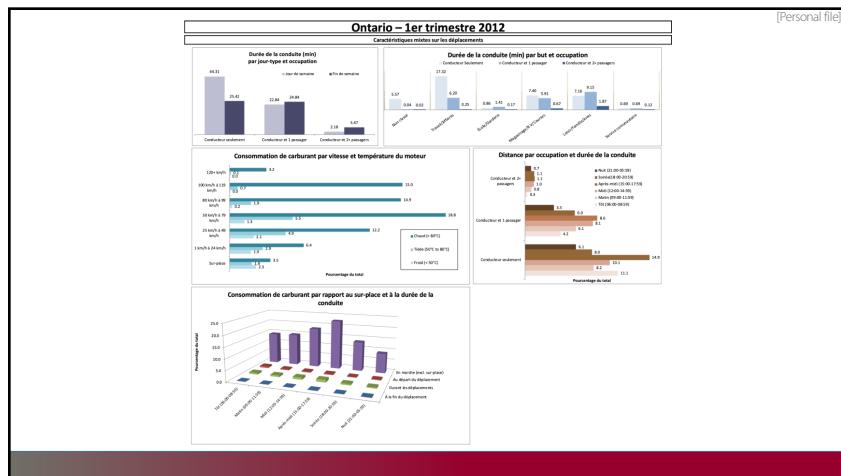
114



115



116



117

Exercise

In teams or individually, identify a scenario for which a dashboard could prove useful.

Determine specific questions that the dashboard could help answer or insights that it could provide.

Identify data sources and data elements that could be fed into your dashboard.

Design a display (with pen and paper) with mock charts.

What are the strengths and limitations of your dashboard? Is it functional? Elegant?

118

Roundtable: Data Viz Posters



Hits?

Misses?

Thoughts?

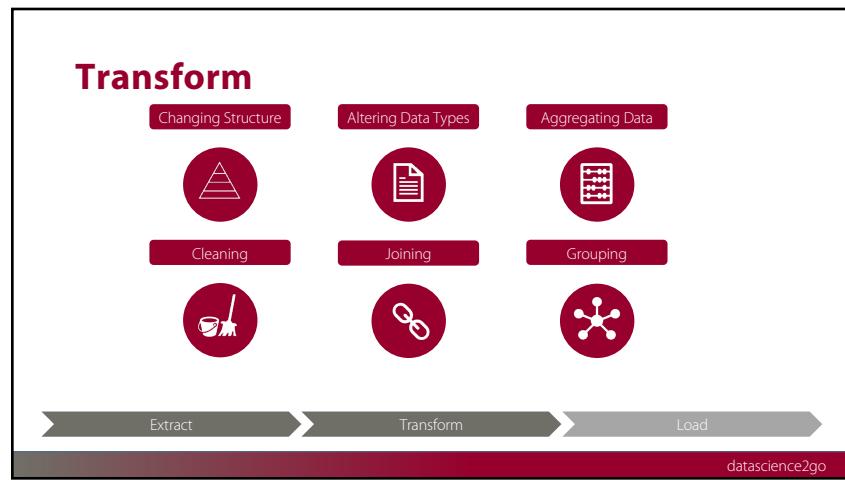
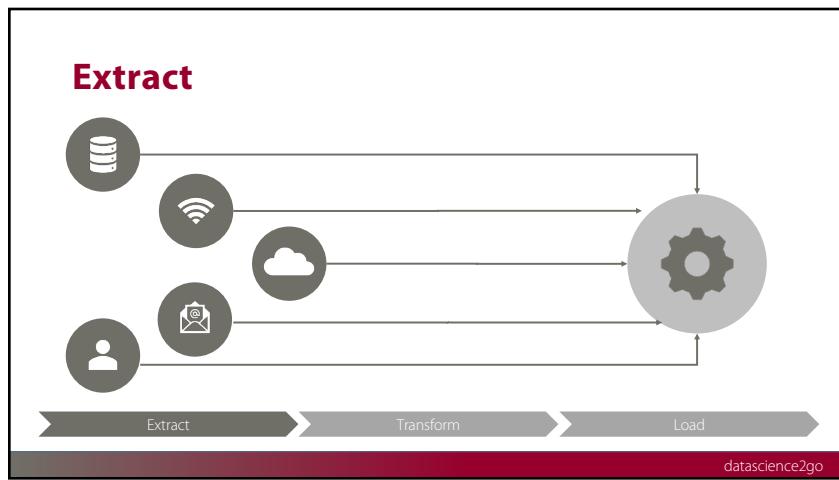
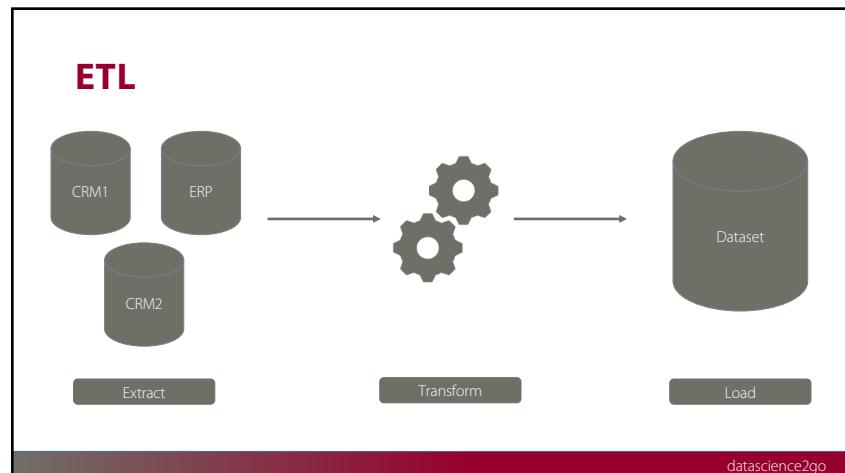
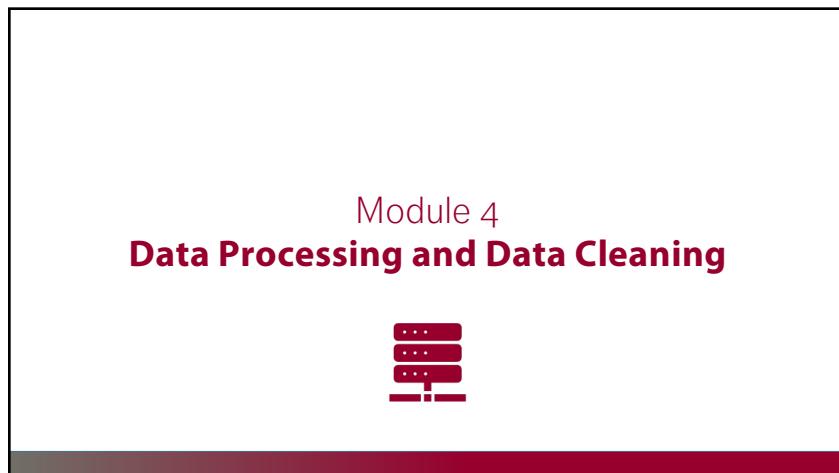
datascience2go

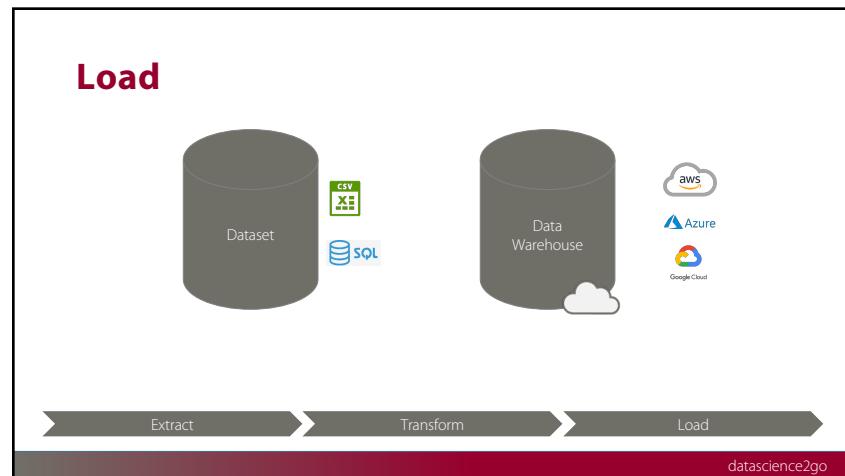
119

Gapminder Exercises

Do the exercises for Module 3.

120





125

Approaches to Data Cleaning

There are two **philosophical** approaches to data cleaning and validation:

- methodical
- narrative

The **methodical** approach consists of running through a **check list** of potential issues and flagging those that apply to the data.

The **narrative** approach consists of **exploring** the dataset and trying to spot unlikely and irregular patterns.

126

Pros and Cons

Methodical (syntax)

- Pros: checklist is **context-independent**; pipelines **easy to implement**; common errors and invalid observations **easily identified**
- Cons: may prove **time-consuming**; cannot identify **new** types of errors

Narrative (semantics)

- Pros: process may simultaneously yield **data understanding**; false starts are (at most) as costly as switching to mechanical approach
- Cons: may miss important sources of errors and invalid observations for datasets with **high number of features**; domain knowledge may **bias the process** by neglecting uninteresting areas of the dataset

127

Tools and Methods

Methodical

- list of potential problems (Data Cleaning Bingo)
- code which can be re-used in different contexts

Narrative

- visualization
- data summary
- distribution tables
- small multiples
- data analysis

128

Data Cleaning Bingo				
random missing values	outliers	values outside of expected range - numeric	factors incorrectly/consistently coded	date/time values in multiple formats
impossible numeric values	leading or trailing white space	badly formatted date/time values	non-random missing values	logical inconsistencies across fields
characters in numeric field	values outside of expected range - date/time	DCB!	inconsistent or no distinction between null, 0, not available, not applicable, missing	possible factors missing
multiple symbols used for missing values	???	fields incorrectly separated in row	blank fields	logical inconsistencies within field
entire blank rows	character encoding issues	duplicate value in unique field	non-factor values in factor	numeric values in character field

129

Approaches to Data Cleaning

The narrative approach is akin to working out a crossword puzzle with a pen and putting down potentially wrong answers **occasionally**, to see where that takes you.

The mechanical approach is akin to working it out with a pencil, a dictionary, and never jotting down an answer unless you are certain it is correct.

You'll solve more puzzles (and it will be flashier) the first way, but you'll rarely be wrong the second way.

It's the same thing with data: analysts must be comfortable **with both approaches**.

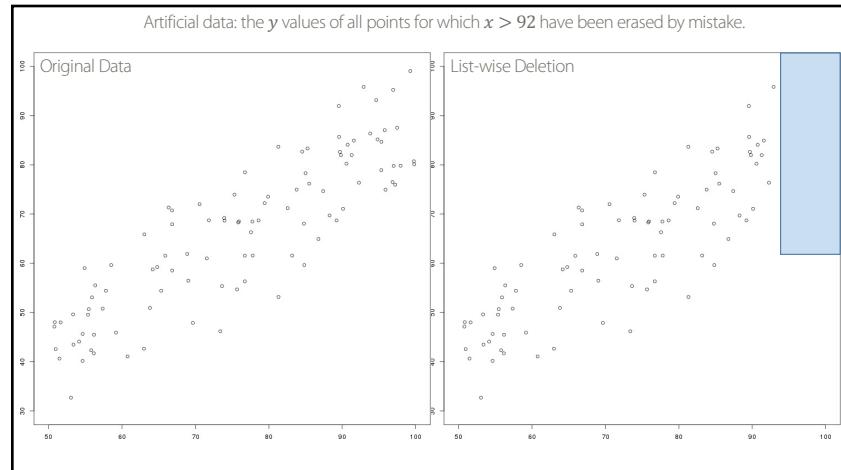
130

The Case for Imputation

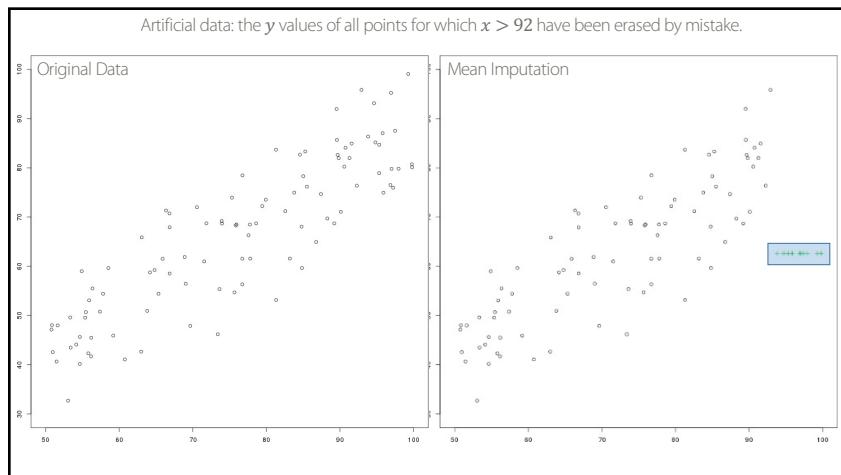
Not all analytical methods can easily accommodate missing observations – 2 options:

- **Discard** the missing observation
 - not recommended, unless the data is missing completely randomly in the dataset
 - acceptable in certain situations (small number of missing values in a large dataset)
- Establish a **replacement (imputation) value**
 - main drawback: we never know what the true value would have been
 - often the best available option

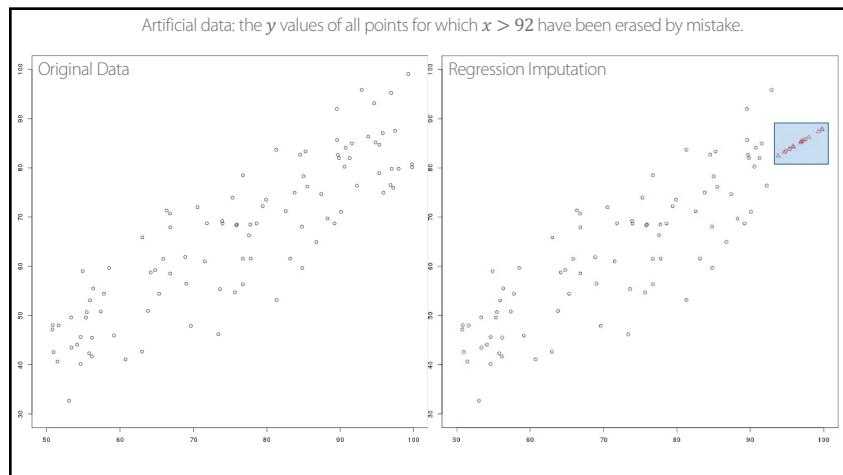
131



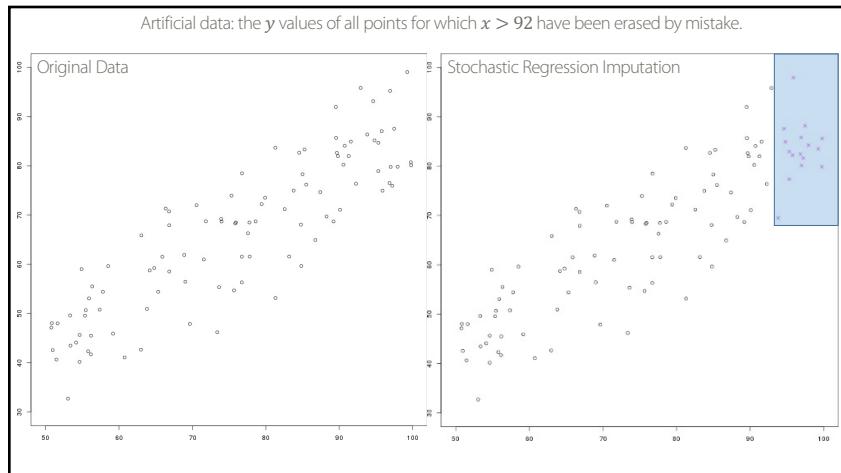
132



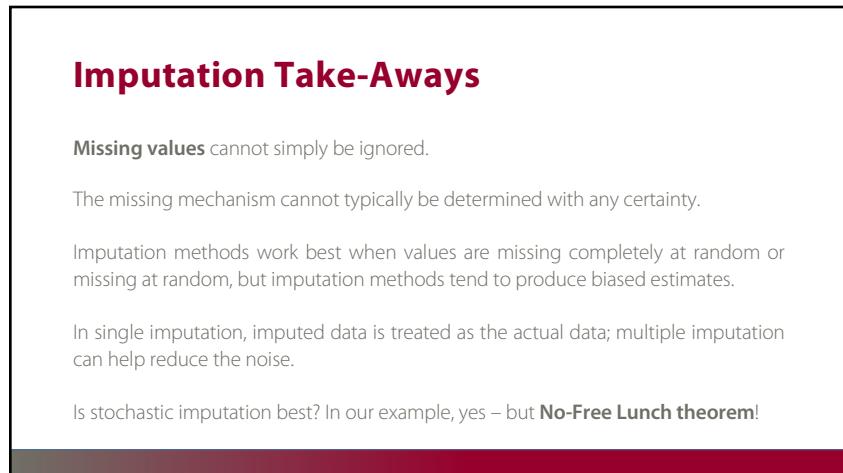
133



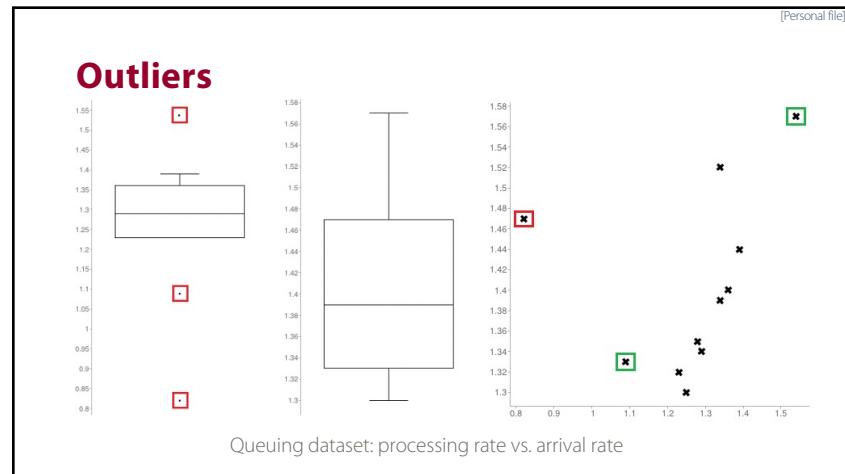
134



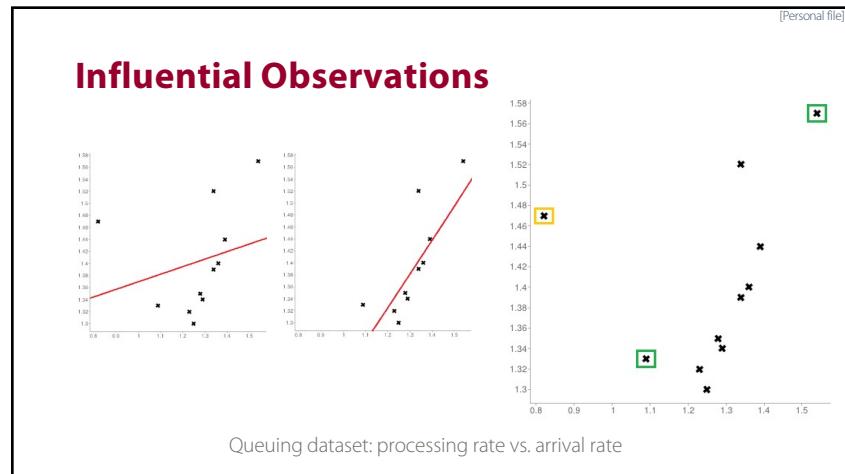
135



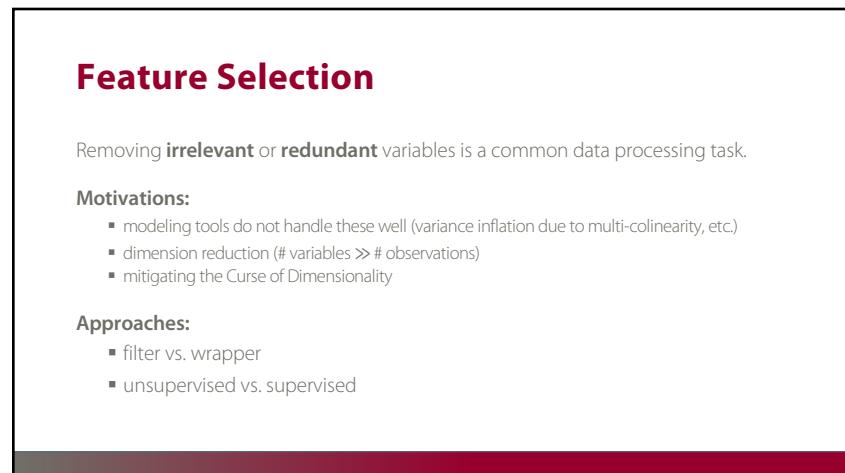
136



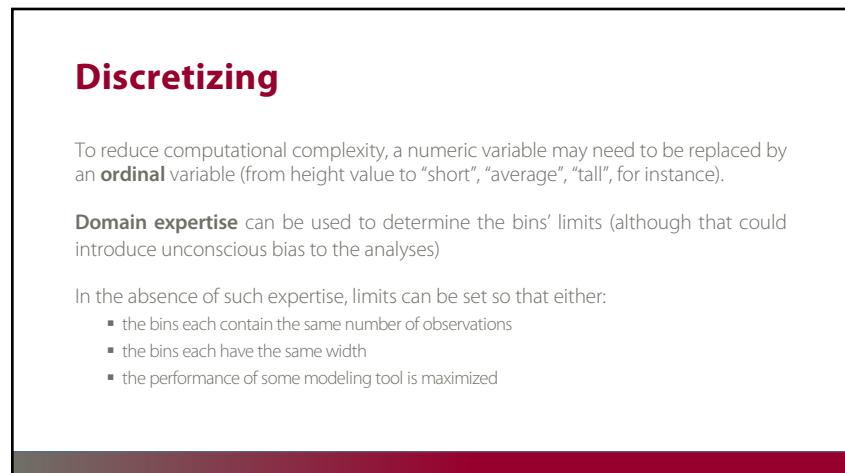
137



138



139



140

Sound Data

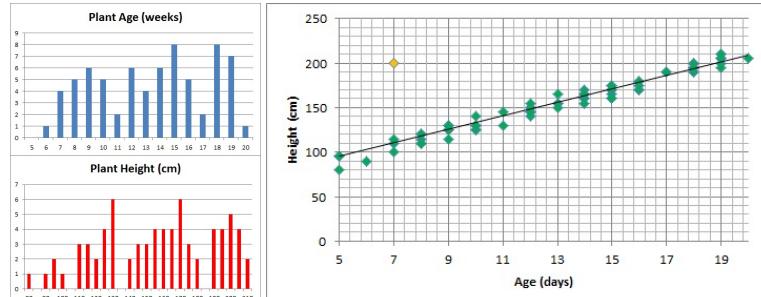
The ideal dataset will have as few issues as possible with:

- **validity:** data type, range, mandatory response, uniqueness, value, regular expressions
- **completeness:** missing observations
- **accuracy and precision:** related to measurement and/or data entry errors; target diagrams (accuracy as bias, precision as standard error)
- **consistency:** conflicting observations
- **uniformity:** are units used uniformly throughout?

Checking for data quality issues **early** can save headaches at a later analytical stage.

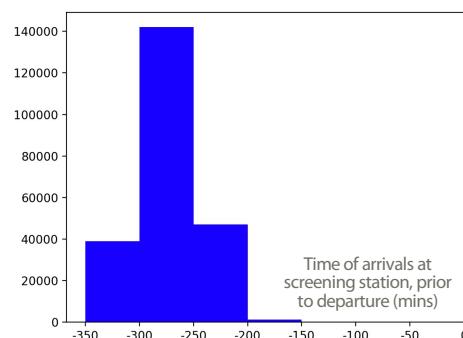
141

Detecting Invalid Entries



142

Detecting Invalid Entries



143

Data Quality Take-Aways

Don't wait until **after** the analysis to find out there was a problem with data quality.

Univariate tests don't always tell the whole story.

Visualizations can help.

Context is crucial – you may need more context about the data in order to make sense of what you see... but whatever the situation, you need to understand the data quality.

144

Putting it All Together

-  Iterating Constantly
-  Communicating the Data
-  Documenting the Process and Assumptions
-  Collecting, Creating, and Cleaning
-  Aligning All Projects
-  Planning and Designing the 'Plan of Attack'
-  Understanding Organizational Needs

datascience2go

145

Gapminder Exercises

Do the exercises for Module 4.

146

Module 5 Data Exploration and Data Analysis



147

Exploratory Data Analysis (Big Picture)

It is important to understand what the data looks like before conducting analyses

EDA = Visualize + Compute Basic Statistics



datascience2go

148

Intuition for Data Analysis



Why are some people 'poor'?

datascience2go

149

Intuition for Data Analysis



Surface level answer:
Because they don't have money

datascience2go

150

Intuition for Data Analysis



datascience2go

151

Intuition for Data Analysis

Assuming vs Analyzing

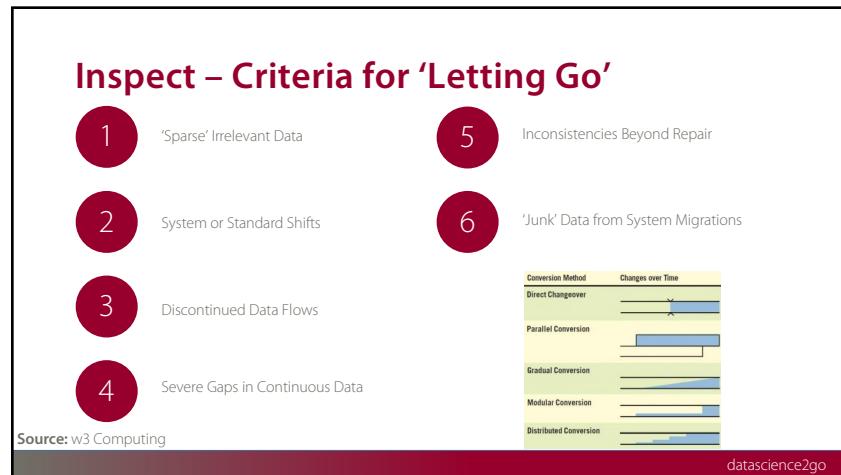


datascience2go

152



153



154



155

Verify – Previewing Numerical Data

The data preview pane shows the following statistics for Column1:

	A	B	C	D	E	F	G	H
1	Scores	83	Column1					
2		93	Mean	81.21428571				
3		91	Standard Deviation	14.64531624				
4		69	Median	85				
5		96	Mode	93				
6		61	Compressed Min	11.54400000				
7		60	Compressed Max	139.00000000				
8		58	Min	143.00000000				
9		59	Max	0.4921930000				
10		100	Range	42				
11		93	Minimum	58				
12		71	Maximum	100				
13		78	Sum	1137				
14		98	Count	14				
15								
16								

datascience2go

156

Verify – Data Scavenger Hunting

- 1 Use 'Look-alike' Data
- 2 Leverage 'Open' Datasets
- 3 Create 'Synthetic' Data
- 4 Extrapolate Data if Statistically Sig.

datascience2go

157

Report – Create a Data Dictionary

Field Name	Data Type	Data Format	Field Size	Description	Example
License ID	Integer	NNNNNN	6	Unique number ID for all drivers	12345
Surname	Text		20	Surname for Driver	Jones
First Name	Text		20	First Name for Driver	Arnold
Address	Text		50	First Name for Driver	11 Rocky st Como 2233
Phone No.	Text		10	License holders contact number	0400111222
D.O.B	Date / Time	DD/MM/YYYY	10	Drivers Date of Birth	08/05/1956

Source: Medium

datascience2go

158

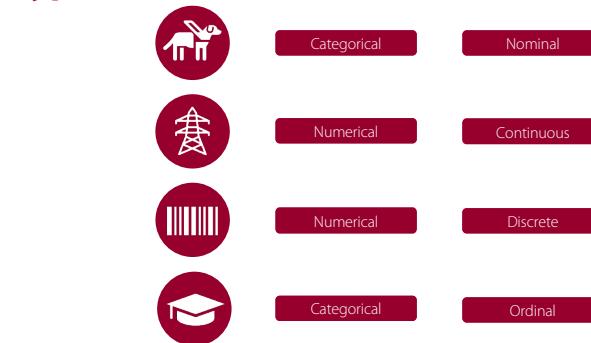
Types of Data



datascience2go

159

Types of Data



datascience2go

160

Types of Data



Categorical	Nominal
Numerical	Continuous
Numerical	Continuous
Categorical	Ordinal

datascience2go

161

Special Role of Categorical Data

Categorical data plays a special role:

- in **data science**, categorical variables come with a pre-defined set of values
- in **experimental science**, a factor is an independent variable with its levels being defined (it may also be viewed as a category of treatment)
- in **business analytics**, these are dimensions (with members) vs. measures

However they are labeled, they are used to subset or **roll up/summarize** the data.

162

Data Summarizing

Min: smallest value

Max: largest value

Median: "middle" value

Mode: most frequent value

Unique Values: list of unique values

etc.

Signal	Type
4.31	Blue
5.34	Orange
3.79	Blue
5.19	Blue
4.93	Green
5.76	Orange
3.25	Orange
7.12	Orange
2.85	Blue

163

Contingency/Pivot Tables

Contingency table: examines the relationship between two categorical variables via their relative (cross-tabulation).

Pivot table: a table generated by applying operations (sum, count, mean, etc.) to variables, possibly based on another (categorical) variable.

Contingency tables are **special cases** of pivot tables.

	Large	Medium	Small
Window	1	32	31
Door	14	11	0

Type	Count	Signal avg	Signal stdev
Blue	4	4.04	0.98
Green	1	4.93	N.A.
Orange	4	5.37	1.60

164

Analysis Through Visualization

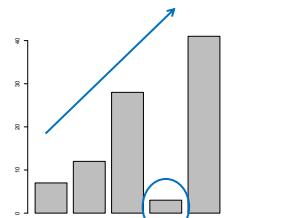
Analysis (broad definition):

- identifying patterns or structure
- adding meaning to these patterns or structure by interpreting them in the context of the system.

Option 1: use analytical methods to achieve this.

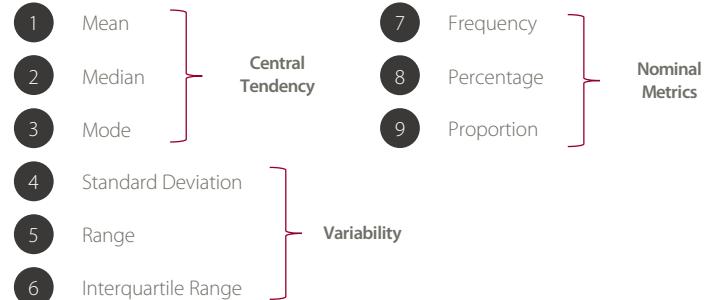
Option 2: visualize the data and use the brain's analytic power (perceptual) to reach meaningful conclusions about these patterns.

We will discuss further.



165

Descriptive Statistics



datascience2go

166

Nominal Data

Frequencies

Count number of events

Proportion

Divide frequency by total number of events

Percentage

Multiply proportion by 100

datascience2go

167

Descriptive Statistics – Central Tendency

Mean

Sum of values divided by the number of observations.

$$\text{Mean} = \frac{7 + 2 + 1 + 6 + 4 + 5 + 5 + 5}{8} = 4.25$$

Median

Number that divides the data set in half.

$$\text{Median} = \frac{2 + 1 + 6 + 4 + 5 + 5 + 5 + 6}{8} = 4.5$$

Mode

The most frequently occurring number.

$$\text{Mode} = \frac{2 + 1 + 6 + 4 + 5 + 5 + 5 + 6}{8} = 5$$

datascience2go

168

Descriptive Statistics – Variability

Standard Deviation

Amount of variation between the mean and rest of the data points.

Range

Difference between Min and Max values.

Interquartile Range

Middle fifty of the data – where the majority of the data lies.

datascience2go

169

Descriptive Statistics – Variability



Quartile	Result	Definition
0	31	Minimum Value
1	43.25	25 th Percentile
2	48.5	50 th Percentile (median)
3	52.25	75 th Percentile
4	65	Maximum Value

$$\begin{aligned} \text{Range} &= 65 - 31 = 34 \\ \text{IQR} &= 52.25 - 43.25 = 9 \\ \text{Std Dev} &= 10.36 \end{aligned}$$

datascience2go

170

Visual Summary - Boxplot

The boxplot is a **graphical summary** of a univariate distribution.

Draw a box along the observation axis, with **endpoints** at Q_1 and Q_3 , and with a “belt” at the median.

Plot a line extending from Q_1 to the smallest obs. less than $1.5 \times \text{IQR}$ below Q_1 .

Plot a line extending from Q_3 to the smallest obs. more than $1.5 \times \text{IQR}$ above Q_3 .

Any suspected outlier is plotted separately.

171

Visual Summary – Histogram

Histograms can also provide an indication of the distribution of a variable.

They should include/contain the following information:

- the **range** of the histogram is $r = Q_4 - Q_0$;
- the **number of bins** should approach $k = \sqrt{n}$, where n is the number of observations;
- the **bin width** should approach r/k , and
- the **frequency of observations** in each bin should be added to the chart.

172

Example

Consider the daily number of car accidents in Sydney over a 40-day period:

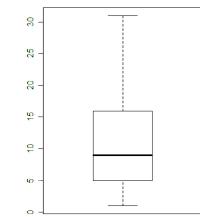
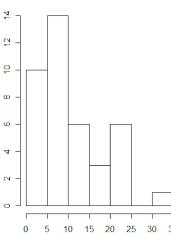
6, 3, 2, 24, 12, 3, 7, 14, 21, 9, 14, 22, 15, 2, 17, 10, 7, 7, 31, 7, 18, 6, 8, 2, 3, 2, 17, 7, 7, 21, 13, 23, 1, 11, 3, 9, 4, 9, 9, 25

The sorted values are:

1 2 2 2 2 3 3 3 3 4 6 6 7 7 7 7
7 8 9 9 9 10 11 12 13 14 14 15 17
17 18 21 21 22 23 24 25 31

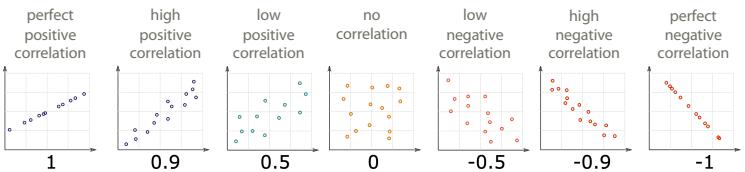
min	Q_1	med	Q_3	max
1	5.5	9	15.5	31

Is it more likely that we have between 5-15 accidents on a given day, or between 25-35?

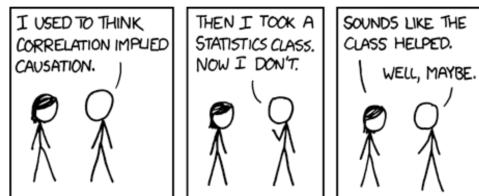


173

Correlation



174



Correlation doesn't imply causation, but it does waggle its eyebrows suggestively and gesture furtively while mouthing 'look over there'.

175

Regression Modeling

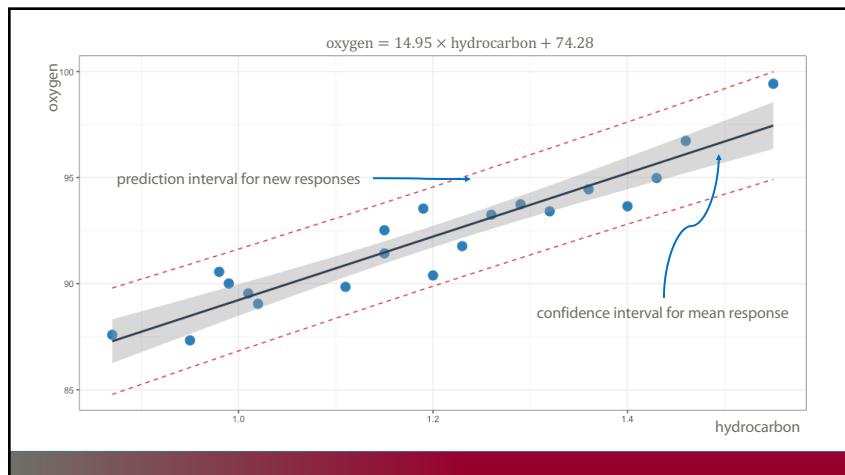
The most common data modeling methods are **regressions**, both linear and logistic.

About 80% of real data applications use a regression as their **final model**, typically after very careful **data preparation**, **encoding**, and **creation of variables**.

There are several reasons for their frequent use:

- generally straightforward to **understand** and to **train**
- mean square error (MSE) objective function has a closed-form linear solution
- system of equations can usually be solved through matrix inversion or linear manipulation

176



177

Other Analytical Approaches

Categorical analysis
Monte-Carlo simulations
Design of experiments
Bayesian data analysis
Times series analysis
Machine learning
Optimization
Queueing models
etc.

178

Gapminder Exercises

Do the exercises for Module 5.

179

Module 6 Data Mining and Machine Learning



What is Machine Learning?

Starting around the 1940s, researchers began the earnest study of how to **teach machines to learn**.

The goal of **machine learning** was (is?) to create machines that can **learn, adapt, and respond** to novel situations

A wide variety of techniques, accompanied by a great deal of theoretical underpinning, was created to achieve this goal.

181

What is Artificial/Augmented Intelligence?

Artificial Intelligence (A.I.) is non-human intelligence that has been engineered rather than one that has evolved naturally.

A.I. research is research carried out in pursuit of this goal.

Pragmatically speaking, A.I. is “computers carrying out tasks that only humans can usually do”.

Augmented Intelligence is human intelligence that is supported or enhanced by machine intelligence.

182

The Mining Analogy

What are we mining? data (**earth**)

What are we using to mine? data mining techniques (**digging tools**)

What are we mining for? looking for patterns/knowledge (**raw minerals**)

What do we do with the raw material? describe patterns/relationships (**refine minerals into something useful**)

What is the output, or product? models (**Ge, Ga, Si to build transistors**)

What do we do with the product? apply models to evidence-based decision support (**use transistor in electrical systems**)

183

Learning in General

Beyond “just taking a quick look,” humans learn through:

- answering questions
- testing hypotheses
- creating concepts
- making predictions
- creating categories and classifying objects
- grouping objects

The **central Data Science/Machine Learning problem** is:

can (**should**) we design algorithms that can learn?

184

Types of Learning

Supervised Learning (learning with a teacher)

- classification, regression, rankings, recommendations
- uses labeled training data (**student gives an answer to each test question based on what they learned from worked-out examples**)
- performance is evaluated using testing data (**teacher provides the correct answers**)

Unsupervised Learning (grouping similar exercises together as a study aid)

- clustering, association rules discovery, link profiling, anomaly detection
- uses unlabeled observations (**teacher is not involved**)
- accuracy cannot be evaluated (**students might not end up with the same groupings**)

185

Types of Learning

Semi-Supervised Learning (teacher providing worked-out examples and a list of unsolved problems)

Reinforcement Learning (embarking on a research project with an advisor?)

In supervised learning, there's a target against which to train the model.

In unsupervised learning, we don't know what the target is, or if there is one.

The distinction is crucial.

186

Learning Tasks

Classification and **class probability estimation**: which clients are likely to be repeat customers?

Clustering: do diplomatic missions form natural groups?

Association rule discovery: what books are commonly purchased together?

Others:

profiling and behaviour description; link prediction; value estimation (how much is a client likely to spend in a restaurant); similarity matching (which prospective clients are similar to a company's best clients?); data reduction; influence/causal modeling, etc.

187

Case Study: Association Rules Mining

The Danish National Patient Registry contains 68 million health observations on 6.2 million patients over a 15-year time span ('96 –'10).

Objectives:

- finding connections between different diagnoses
- determining how a diagnosis at some point in time might allow for the prediction of another diagnosis at a later point in time

Temporal disease trajectories condensed from population-wide registry data covering 6.2 million patients
Jensen, A.B., Moseley, P.L., Oprea, T.I., Ellesøe, S.G., Eriksson, R., Schmock, H., Jensen, P.B., Jensen, L.J., Brunak, S. [2014]. Nature Communications.

188

Methodology

1. Compute strength of correlation for **pairs of diagnoses** over a 5-year interval on a **representative subset** of the data
2. Test pairs for **directionality** (one repeatedly occurring before the other)
3. Determine reasonable **diagnosis trajectories** (thoroughfares) by combining smaller frequent trajectories with overlapping diagnoses
4. Validate the trajectories by **comparison with non-Danish data**
5. Cluster the thoroughfares to identify **central medical conditions** around which disease progression is organized

189

Results

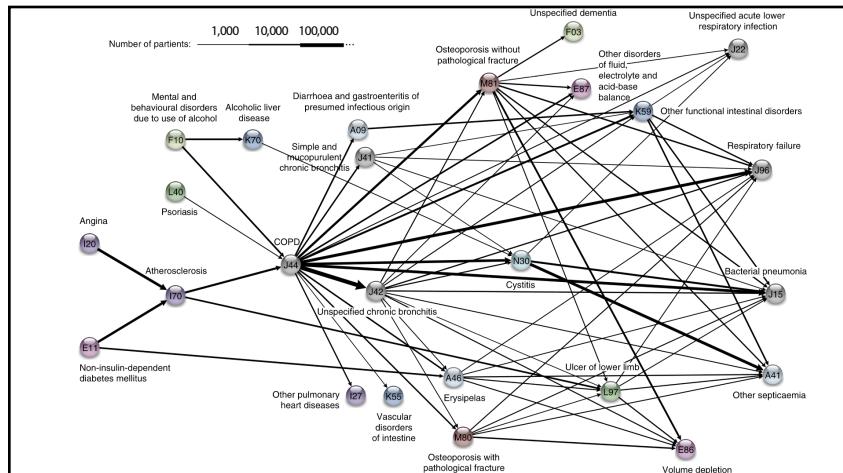
Data was reduced to **1,171 thoroughfares**, with **5 key diagnoses**:

- diabetes
- chronic obstructive pulmonary disease (COPD)
- cancer
- arthritis
- cardiovascular disease.

The data analysis showed, for example:

- diagnoses of **anemia** followed later by the discovery of **colon cancer**
- **gout** was identified as a step toward **cardiovascular disease**
- **COPD** is under-diagnosed and under-treated

190



191

Case Study Take-Aways

Data makes it possible to view diseases in a larger context, which could yield tangible **health benefits** beyond one-size-fits-all medicine.

The **sooner** a health risk pattern is identified, the **better** we can prevent and treat critical diseases.

Instead of looking at each disease in isolation, we can talk about a **complex system** with many different **interacting factors**.

The **order** in which different diseases appear can help find patterns and complex correlations outlining the direction for each individual person.

192

Association Rules Basics

Association Rule Discovery is unsupervised learning that finds connections among attributes (and combinations of attributes).

Example: we might analyze a dataset on the physical activities and purchasing habits of North Americans and discover that

- runners who are also triathletes (the **premise**) tend to drive Subarus, drink microbrews, and use smartphones (the **conclusion**), or
- individuals who have purchased home gym equipment are unlikely to be using it 1 year later (to name some fictitious possibilities)

193

Market Basket Analysis

Supermarkets record the contents of shopping carts at check-outs to determine items which are **frequently purchased together**.

Examples:

- **bread** and **milk** are often purchased together, but that's not so interesting given how often they are purchased individually
- **hot dog buns** and **wieners** are also often purchased as a pair, but more rarely purchased individually

A supermarket could then have a sale on hot dogs to drive in customers, while raising the price on condiments, to drive in sales.

194

Applications

Related Concepts

- looking for pairs (triplets, etc) of words that represent a joint concept
- {Ottawa, Senators}, {Michelle, Obama}, {veni, vidi, vici}, etc.

Plagiarism

- looking for sentences that appear in various documents
- looking for documents that share sentences

Bio-markers

- diseases that are frequently associated with a set of bio-markers

195

Applications

Making predictions and decisions based on these rules.

Alter circumstances or environment to take advantage of these correlations (often mis-used).

Use the connections to modify the likelihood of certain outcomes.

Imputing missing data.

Text autofill and autocorrect.

196

Causation and Correlation

Association rules can help automate **hypothesis discovery**, but we must remain correlation-savvy (which is less prevalent among analysts than one would hope...).

If attributes A and B are shown to be **correlated**, then the possibilities are:

- A and B are correlated entirely by chance in this dataset
- A is a relabeling of B
- A causes B and/or B causes A
- combinations of other attributes C_1, \dots, C_n (known or not) cause A & B
- etc?

197

Causation and Correlation

Insight	Organization
Pop-Tarts before a hurricane	Walmart
Higher crime, more Uber rides	Uber
Typing with proper capitalization indicates creditworthiness	A financial services startup company
Users of the Chrome and Firefox browsers make better employees	A human resources professional services firm, over employee data from Xerox and other firms
Men who skip breakfast get more coronary heart disease	Harvard University medical researchers
More engaged employees have fewer accidents	Shell
Smart people like curly fries	Researchers at the University of Cambridge and Microsoft Research
Female-named hurricanes are more deadly	University researchers
Higher status, less polite	Researchers examining Wikipedia behavior

198

Case Study: Minnesota Tax Audit

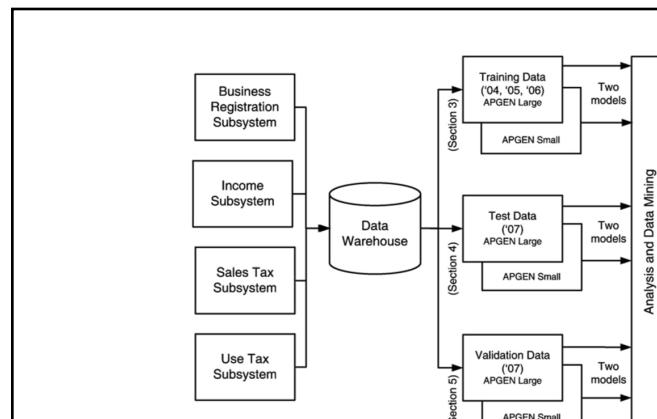
Large gaps between revenue owed (in theory) and revenue collected (in practice) are problematic for governments.

Revenue agencies implement various fraud detection strategies (such as audit reviews) to bridge that gap.

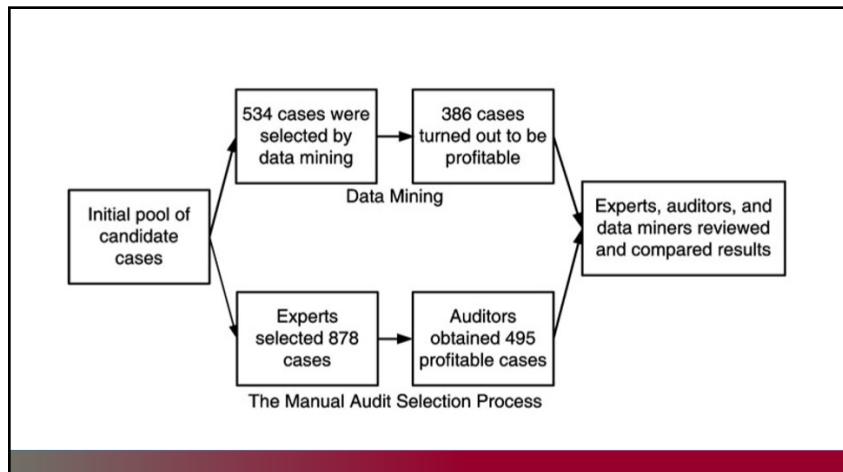
Business audits are costly – are there **algorithms that can predict whether an audit is likely to be successful or a waste of resources?**

Data mining-based tax audit selection: a case study of a pilot project at the Minnesota Department of Revenue
Hsu, W., Pathak, N., Srivatsava, J., Tschida, Bjorklund, E. [2013], Real World Data Mining Applications, Annals of Information Systems, v.17, Springer.

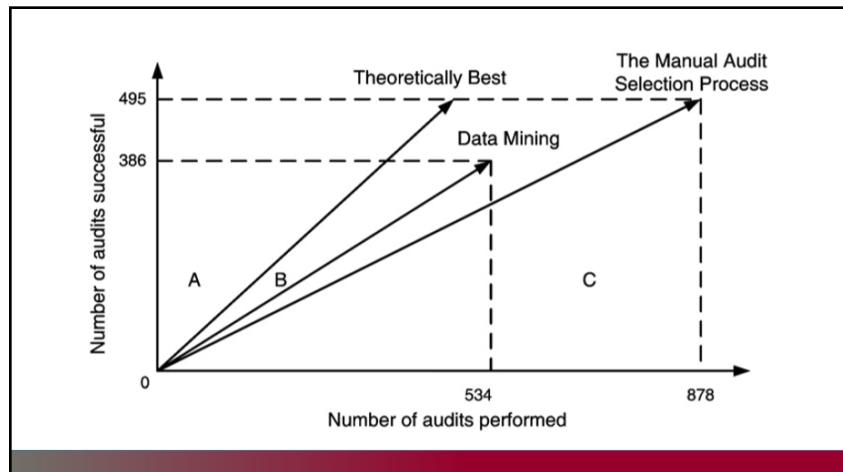
199



200



201



202

	Predicted as good	Predicted as bad
Actually good	386 (Use tax collected) R = \$5,577,431 (83.6 %) C = \$177,560 (44 %)	109 (Use tax lost) R = \$925,293 (13.9 %) C = \$50,140 (12.4 %)
Actually bad	148 (costs wasted) R = \$72,744 (1.1 %) C = \$68,080 (16.9 %)	235 (costs saved) R = \$98,105 (1.4 %) C = \$108,100 (26.7 %)

203

Classification Overview

In **classification**, a sample set of data (the **training set**) is used to determine **rules** and **patterns** that divide the data into pre-determined groups (**classes**).

Classification is a **supervised learning** task.

The training data usually consists of a **randomly selected** subset of the **labeled** (target) data.

Value estimation (regression) is akin to classification, but the target variable is **numerical**.

204

Classification Overview

In the **testing phase**, the model is used to assign a class to observations for which the **label is hidden**, but ultimately known (the **testing set**).

The **performance** of a classification model is evaluated on the testing set, never on the training set.

Technical challenges include:

- selecting the features to include in the model
- selecting the algorithm
- etc.

205

Applications

Medicine and Health Science

- predicting which patient is at risk of suffering a second, fatal heart attack within 30 days based on health factors (blood pressure, age, sinus problems, etc.)

Social Policies

- predicting the likelihood of requiring assisted housing in old age based on demographic information/survey answers

Marketing and Business

- predicting which customers are likely to switch to another cell phone company based on demographics and usage

206

Other Uses

Predicting that an object belongs to a particular class.

Organizing and grouping instances into categories.

Enhancing the detection of relevant objects

- avoidance: "this object is an incoming vehicle"
- pursuit: "this borrower is unlikely to default on her mortgage"
- degree: "this dog is 90% likely to live until it's 7 years old"

In the absence of testing data, classification may be **descriptive** but not predictive.

207

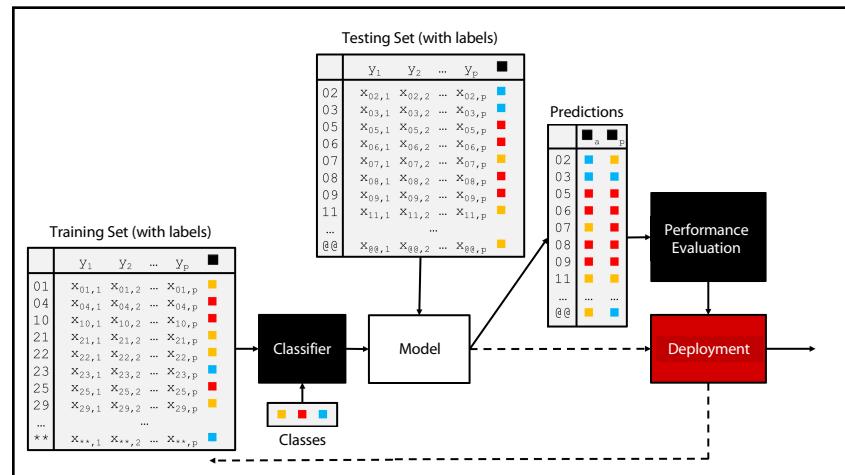
Example

Scenario: a motor insurance company has a fraud investigation dept. that studies up to 30% of all claims made, yet money is still getting lost on fraudulent claims.

Questions: can we predict

- whether a claim is likely to be fraudulent?
- whether a customer is likely to commit fraud in the near future?
- whether an application for a policy is likely to result in a fraudulent claim?
- the amount by which a claim will be reduced if it is fraudulent?

208



209

Classification Methods

Logistic Regression

- classical model
- affected by variance inflation and variable selection process

Neural Networks

- hard to interpret
- requires all variables to be of the same type
- easier to train since backpropagation (chain rule)

Decision Trees

- may overfit the data if not pruned correctly (manually?)

210

Classification Methods

Naïve Bayes Classifiers

- quite successful for text mining applications (spam filter)
- assumptions not often met in practice

Support Vector Machines

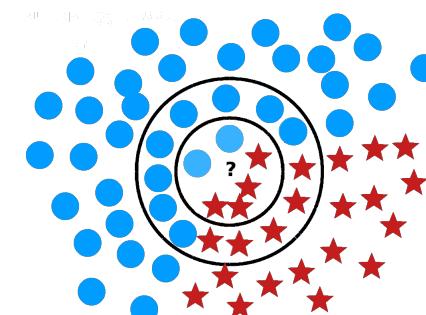
- may be difficult to interpret (non-linear boundaries)
- can help mitigate big data difficulties

Nearest Neighbours Classifiers

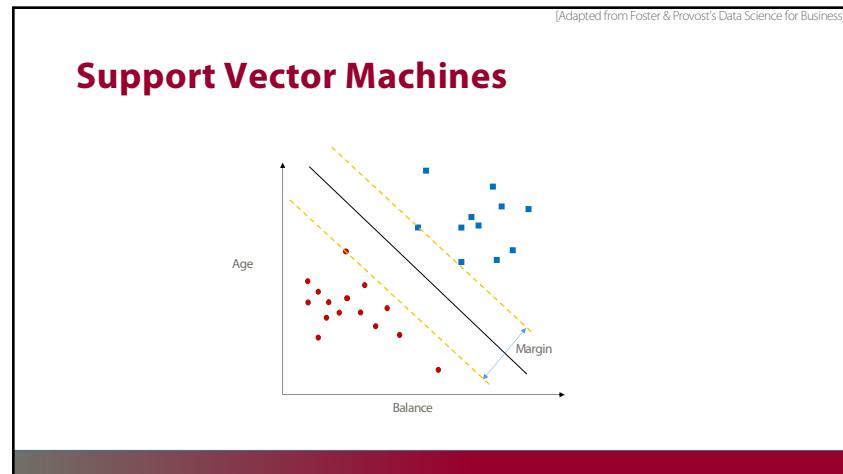
- require very little assumptions about the data
- not very stable (adding points may substantially modify the boundary)

211

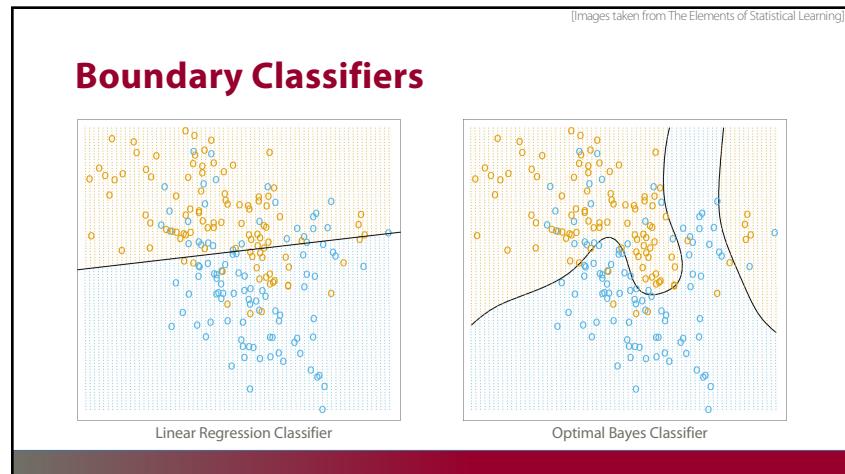
k – Nearest Neighbours



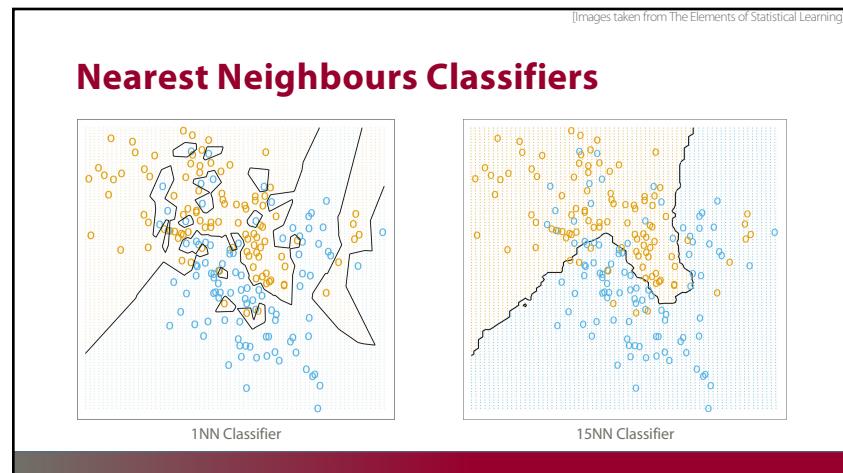
212



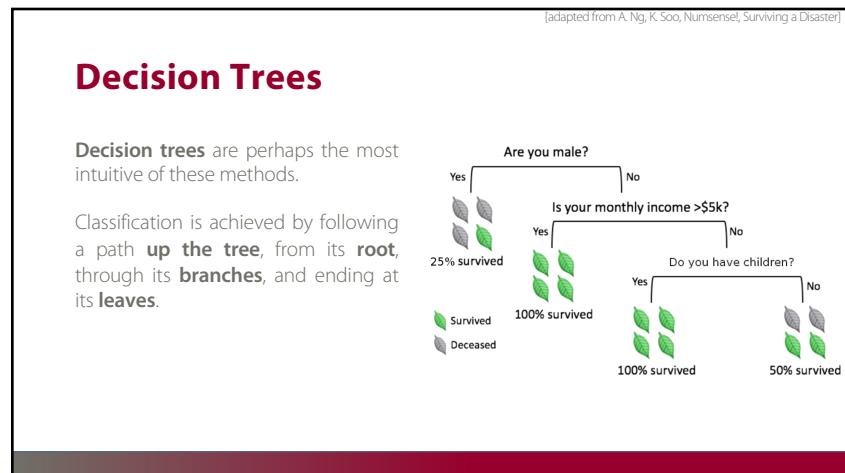
213



214



215



216

Performance Evaluation

Classifiers are evaluated on a **testing** set.

Ideally, a good classifier would have high rates of both **True Positives** (TP) and **True Negatives** (TN), and low rates of both **False Positives** (FP, Type I error) and **False Negatives** (FN, Type II error).

Evaluation metrics mean very little on their own: context requires comparison with other classifiers, and other evaluation metrics.

217

Performance Evaluation

$$\text{sensitivity} = \frac{TP}{TP + FN}$$

$$\text{specificity} = \frac{TN}{FP + TN}$$

$$\text{precision} = \frac{TP}{TP + FP}$$

$$\text{recall} = \frac{TP}{TP + FN}$$

$$\text{negative predictive value} = \frac{TN}{TN + FN}$$

$$\text{false positive rate} = \frac{FP}{FP + TN}$$

$$\text{false discovery rate} = \frac{FP}{FP + TP}$$

$$\text{false negative rate} = \frac{FN}{FN + TP}$$

$$\text{accuracy} = \frac{(TP + TN)}{T}$$

		Predicted		Actuals	Category I	Category II	Total
		TP	FN				
Actuals	Category I	TP	FN				
	Category II	FP	TN				
Total							T

Other metrics:

F_1 -score, ROC AUC, informedness, markedness, Matthews' Correlation Coefficient (MCC), etc.

218

Performance Evaluation

		Predicted		Actuals	Classification Rates		Performance Metrics	
		A	B		Sensitivity	Accuracy	F1-Score	
Actuals	A	54	10	64	0.84	0.80	0.87	79.0%
	B	6	11	17	0.65	0.65	0.67	21.0%
Total	60	21	81		0.90	0.79	0.83	
		74.1%	25.9%					
		Predicted		Actuals	Classification Rates		Performance Metrics	
		A	B		Sensitivity	Accuracy	F1-Score	
Actuals	A	54	0	54	1.00	0.80	0.87	66.7%
	B	16	11	27	0.41	0.41	0.41	33.3%
Total	70	11	81		0.77	0.77	0.77	
		86.4%	13.6%					

219

Case Study: OK Cupid

Chris McKinlay, a 35 year old UCLA Math PhD Student, was looking for a romantic partner online with little luck

- OK Cupid algorithms use only the questions that both potential matches decide to answer, and the questions he had chosen (more or less at random up to that point) were not popular

Between June 2012 and December 2013, he:

- used statistical sampling to find questions which mattered to the kind of partner he had in mind;
- constructed a new profile that answered only those questions;
- matched only with women in LA who might be right for him.

K. Poulsen, How a Math Genius Hacked OK Cupid to Find True Love, WIRED

220

Process

This story provides a **great** example of the data mining process, from start to finish:

1. **Collect** data
2. Collect **more** and **slightly better** and **different** data
3. Collect **still more** data
4. Figure out a data mining technique that would be **relevant** to what he wanted to know (clustering)
5. **Validate** the results of the analysis
6. **Investigate** the results, and narrow down which results were interesting
7. Analyze the interesting results **some more**, and use this to solve the original problem
8. Use the data to **improve other areas** of his profile as well
9. Sit back and reap the benefits of data mining?

221

Methodology and Results

Used k -mode to cluster 20,000 women into seven statistically distinct clusters based on their questions and answers.

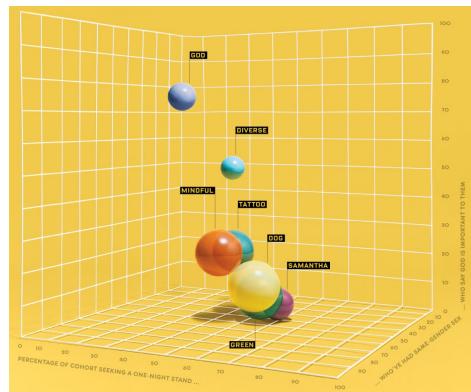
Validated the clustering with another 5,000 profiles from the site.

Analyzed the clusters to find two that interested him

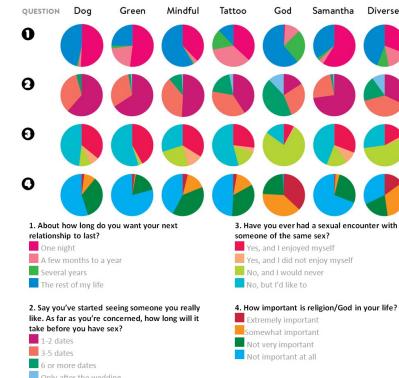
- women in their mid-twenties who looked like indie types, musicians and artists
- slightly older women who held professional creative jobs, like editors and designers.

Used results to derive **which questions he should answer** in his profile, leading to more matches based on his profile, to more first dates, to some second dates, and ... to a lone third date.

222



223

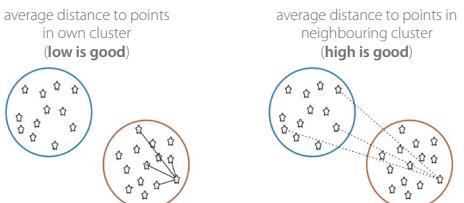


224

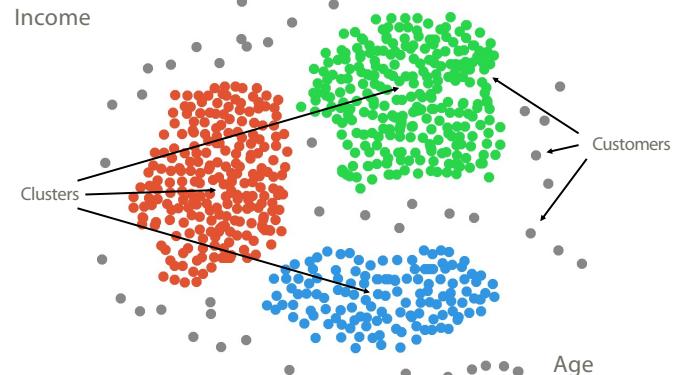
Clustering Overview

In **clustering**, the data is divided into **naturally occurring groups**. Within each group, the data points are **similar**; from group to group, they are **dissimilar**.

The grouping labels are not determined ahead of time, so clustering is an example of **unsupervised** learning.



225



226

Clustering Overview

Clustering is a relatively **intuitive** concept for human beings as our brains do it unconsciously

- facial recognition
- searching for patterns, etc.

In general, people are very good at **messy** data, but computers and algorithms have a harder time.

Part of the difficulty is that there is **no agreed-upon definition of what constitutes a cluster**

- "I may not be able to define what it is, but I know one when I see one"

227

Clustering Overview

Clustering algorithms can be **complex** and **non-intuitive**, based on varying notions of similarities between observations

- in spite of that, the temptation to explain clusters a posteriori is strong

They are also (typically) **non-deterministic**:

- the same algorithm, applied twice (or more) to the same dataset, can discover completely different clusters
- the order in which the data is presented can play a role
- so can starting configurations

228

Applications

Text Documents

- grouping similar documents according to their topics, based on the patterns of common and unusual words

Product Recommendations

- grouping online purchasers based on the products they have viewed, purchased, liked, or disliked
- grouping products based on customer reviews

Marketing and Business

- grouping client profiles based on their demographics and preferences

229

Other Uses

Dividing a larger group (or area, or category) into **smaller** groups, with members of the smaller groups guaranteed to have similarities of some kind

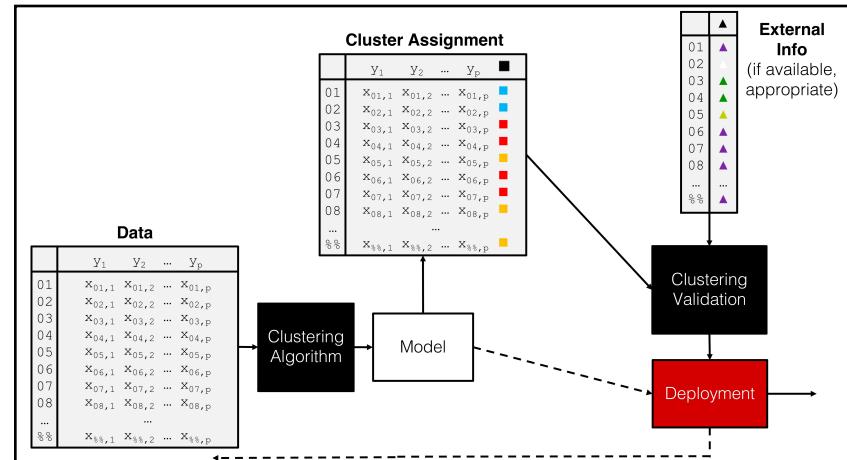
- tasks may then be solved separately for each of the smaller groups
- this may lead to increased accuracy once the separate results are aggregated

Creating (new) taxonomies **on the fly**, as new items are added to a group of items

- this would allow for easier product navigation on a website like Netflix, for instance.

See Spotlight on Clustering, in Data Understanding, Data Analysis, and Data Science for more examples.

230



231

Clustering Methods

k-Means

- classical (and over-used) model
- assumptions made about the shape of clusters

Hierarchical Clustering

- easy to interpret, deterministic

Latent Dirichlet Allocation

- used for topic modeling

Expectation-Maximization

232

Clustering Methods

Balanced Iterative Reducing and Clustering using Hierarchies

- aka BIRCH

Density-Based Spatial Clustering of Applications with Noise

- graph-based

Affinity Propagation

- selects the optimal number of clusters automatically

Spectral Clustering

- recognizes non-blob clusters

233

Clustering Validation

What does it mean for a clustering scheme to be **better** than another?

What does it mean for a clustering scheme to be **valid**?

What does it mean for a single cluster to be **good**?

How many clusters are there in the data, really?

Right vs. wrong is meaningless: seek **optimal vs. sub-optimal**.

234

Clustering Challenges

Automation

relatively intuitive for humans, but harder for machines

Lack of a clear-cut definition

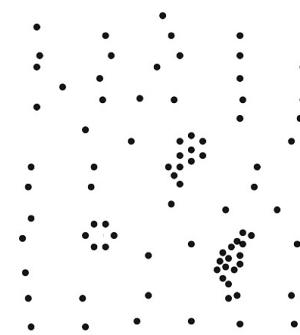
no universal agreement as to what constitutes a cluster

Lack of repeatability

non-deterministic: the same algorithm, applied twice to the same dataset can discover completely different clusters

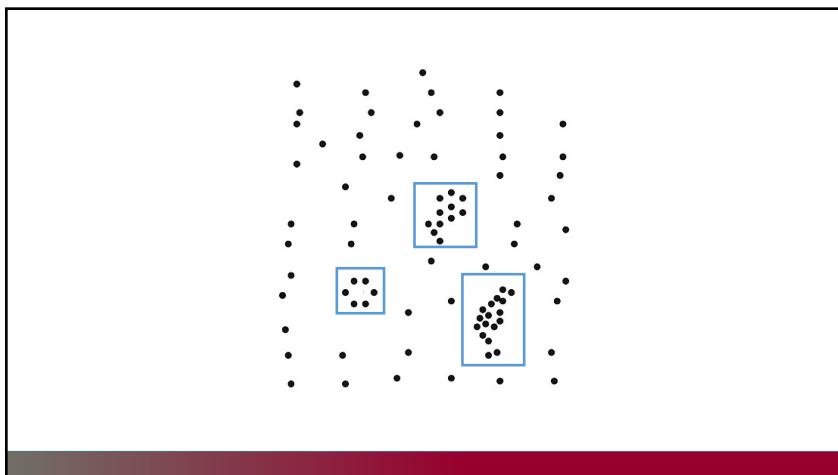
Number of clusters

optimal number of clusters difficult to determine

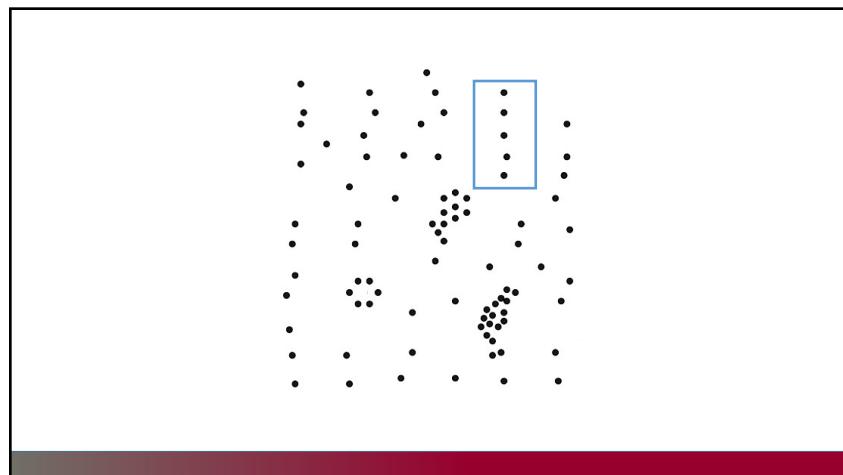


235

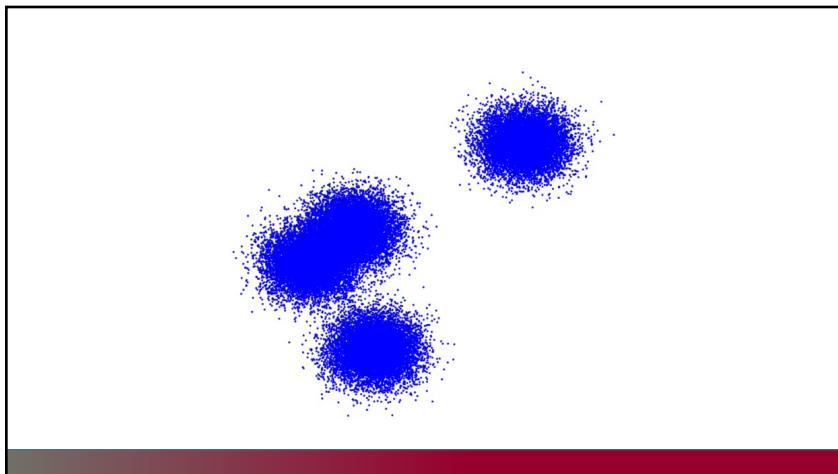
236



237



238



239

Clustering Challenges

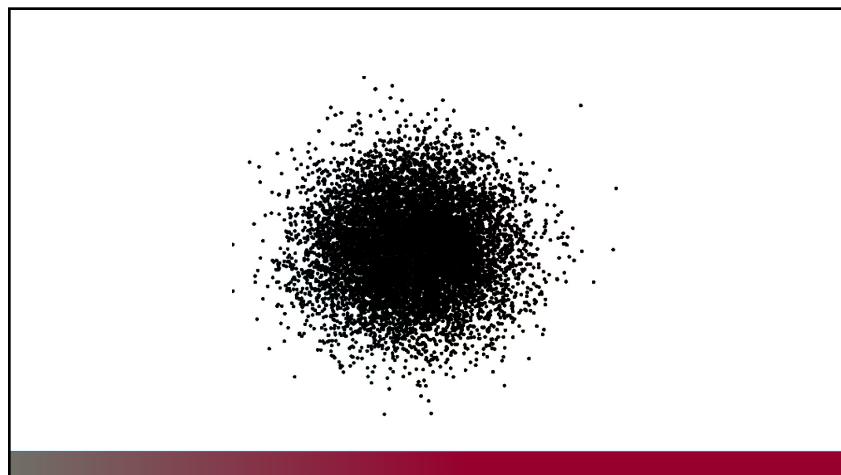
Cluster description
should clusters be described using representative instances or average values?

Model validation
no true clustering information against which to contrast the clustering scheme, so how do we determine if it is appropriate?

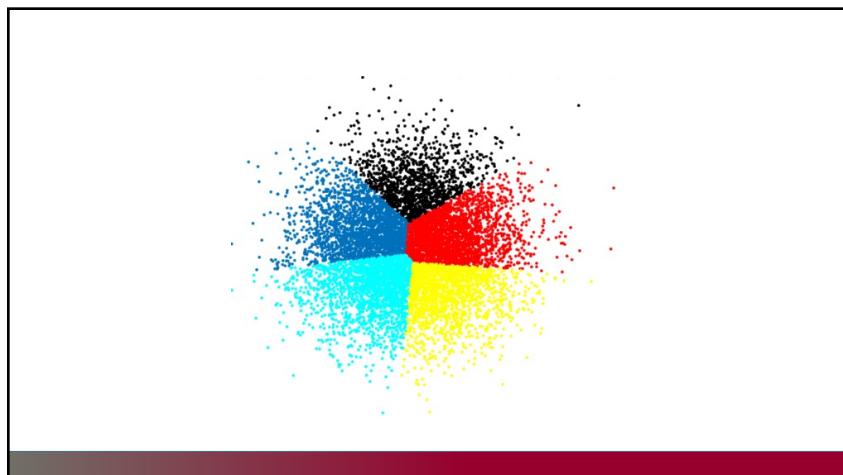
Ghost clustering
most methods will find clusters even if there are none in the data

A posteriori rationalization
once clusters have been found, it is tempting to try to "explain" them ...

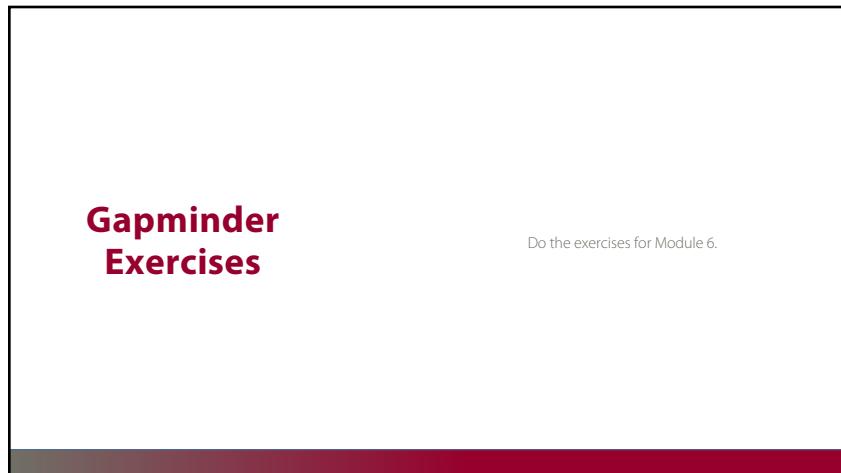
240



241



242



243

Bad Data

Does the dataset pass the **smell test**? (invalid entries, etc.)

Detecting **lies** and **mistakes** (reporting errors, use of polarizing language)

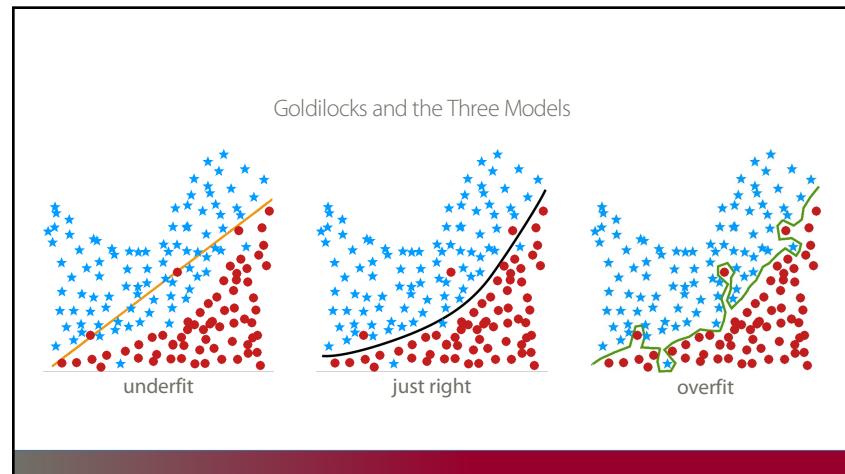
Is **close enough, good enough?**

Sources of **bias** and **errors**

Seeking **perfection** (academic, professional, government, service data)

Data science pitfalls: analysis without understanding, using only one tool (by choice/flat), analysis for the sake of analysis, unrealistic expectations of data science.

244



245

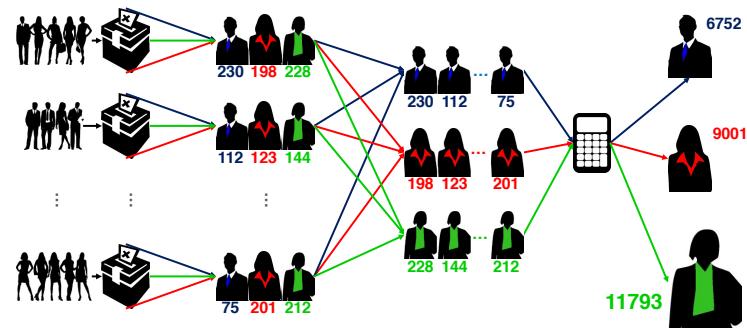
Big Data vs. Small Data

Parameters	Traditional Data	Big Data
Volume	GB	TB or PB
Generate	per hour, per day	every second or microsecond
Structure	structured	semi-structured or unstructured
Data Source	centralized	fully distributed
Data Integration	easy	difficult
Data Store	RDBMS	HDFS, NoSQL
Data Store	interactive	batch or near real time
Access Update Scenario	repeated read and write	write once repeated read
Data Structure	static schema	dynamic schema
Scaling Potential	non-linear	somewhat close to linear

datascience2go

246

Analogy: Election



247

Analogy: Pizzeria

The gains from parallelism depend on whether serial algorithms can be adapted to make use of parallel hardware.

Pizzeria analogy for limitations of parallelization/bottleneck:

- multiple cooks can prepare toppings in parallel
- but baking the crust can't be parallelized
- doubling oven space will increase the number of pizzas that can be made simultaneously but won't substantially speed up any one pizza
- sometimes bottlenecks prevent any gains from parallelism: people line up on both sides of a table to get some soup but there's only one ladle

248

Biases, Fallacies, and Interpretation

When consulting (or conducting) studies, you should try to determine how the following biases could have come into play:

- **Selection bias** (what data was included, how was it selected?)
- **Omitted-variable bias** (were relevant variables ignored?)
- **Detection bias** (did prior knowledge affect the results?)
- **Funding bias** (who's paying for this?)
- **Publication bias** (what's not being published?)
- **Data-snooping bias** (trying too hard?)
- **Analytical bias** (did the choice of specific method affect the results?)
- **Exclusion bias** (are specific observations/units being excluded?)
- etc. (there are tons)

249

Biases, Fallacies, and Interpretation

- | | |
|--------------------------------------------------|-----------------------------------------------------|
| Correlation is not causation (but it is a hint!) | Randomness plays a role |
| Extreme patterns can mislead | Human component to any analytical activity |
| Stay within a study's range | Small effects can be (statistically) significant |
| Keep the base rate in mind | Beware of sacrosanct statistics (p -value, etc.) |
| Odd results happen (Simpson's Paradox) | |

250

Data Science Myths & Mistakes

- [AK Maheshwari, Business Intelligence and Data Mining]
- Mistake #1** – Selecting the wrong problem.
 - Mistake #2** – Getting buried under tons of data without metadata understanding.
 - Mistake #3** – Not planning the data analysis process.
 - Mistake #4** – Insufficient business and domain knowledge.
 - Mistake #5** – Using incompatible data analysis tools.
 - Mistake #6** – Using tools that are too specific.
 - Mistake #7** – Ignoring individual predictions/records in favour of aggregated results.
 - Mistake #8** – Running out of time.
 - Mistake #9** – Measuring results differently than the sponsor.
 - Mistake #10** – Naïvely believing what one's told about the data.

251

What We Didn't Talk About

- Tons of other classification and clustering algorithms
- Recommender systems
- Data streams
- Natural language processing (in depth)
- Feature selection and dimension reduction (curse of dimensionality)
- Data engineering

... and much, much more!

252

The Future of DS/ML/AI

Self-driving vehicles
Machine translation and language understanding
Detection and prevention of climate and ecosystem disturbances
Automated data science (?)
Detection and prevention of astronomical catastrophic events
Explainable AI.

What else?

253

Future Trends

New questions
New tools
New data sources
Data science as job component
Augmented/swarm intelligence

What else?

254

Data Science Buffet

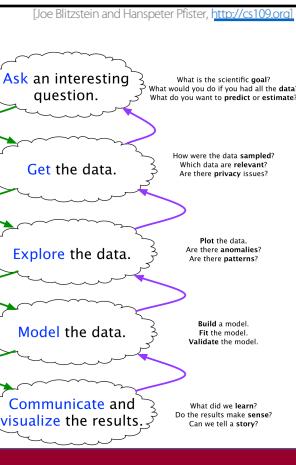


datascience2go

255

In Conclusion

Data science is a team activity, with subject matter experts.
Ethical considerations are paramount and need not conflict with profitability.
Let the data speak (but be careful).
Look for actionable insights.
Supervised vs. unsupervised vs...
Much time must be spent on data preparation.



256