

DATA SCIENCE

Introduction to data science / DSC 105 / Fall 2024

Marcus Birkenkrahe

August 26, 2024

You today



Pythagoras of Samos (ca. 570-495 BC)

- You are like Pythagoras (570-495 B.C.)
- He saw nature as a structured system of numbers
- Data science is (still) the "sexiest" job of the century

WHAT WILL YOU LEARN?



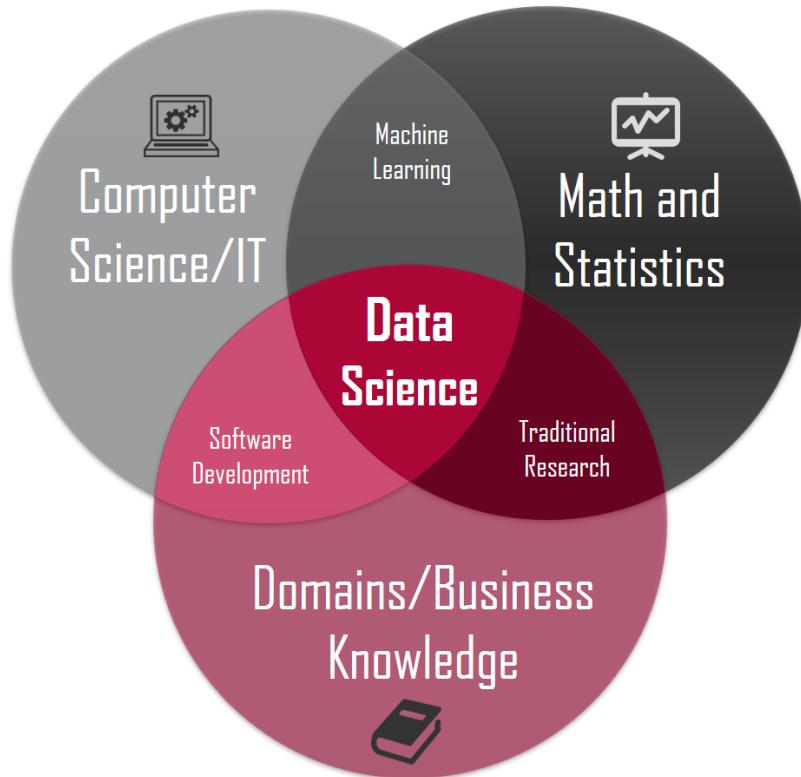
- How and why data science is so **popular**
- What **skills** you need to do data science
- Which **problems** data science can solve
- What data scientists **do** all the time

TOP DATA SCIENCE SKILLS



1. Data analysis (Python, R)
2. Data visualization (R, Power BI, Tableau)
3. Database management (SQL)

WHAT IS IT?



- Vast, new, hard to define, interdisciplinary
- The science of "that what is given" (data)
- Short: **data + code + stats = story**

In this lesson, I am giving an overview of several aspects of data science. Though young as a declared field (2012), it is a field both ill-defined and (or perhaps because of it) vast and hard to pin down. This outline will be applied rather than scholarly, focusing on applications and practice rather than concepts or theory.

In its name, "data science" carries both aspects of science and craft: the 'science' part is responsible for the modus operandi, which is informed by statistics and math, systematic and logically rigorous. The 'data' part relates mostly to craft: the ability to extract insights from data using computing tools. Most data scientists are more occupied by and with the craft part than with the science part (cp. Kozyrkov 2018).

Hence, data science so far is a typical support science. It supports other, more established disciplines in the natural and in the social sciences. Prominent examples are: economics, genomics, and epidemiology.

The need to use the data "to tell a story" sets data science apart from both traditional data craft and science. It is the reason why visualization techniques and theory ("grammar of graphics", cp. Sarkar 2018) play such an important role.

I would argue that data science is most successful when supporting fields that themselves are interdisciplinary and therefore need a higher degree of communication across different cultures of science and practice. This is the quasi-definition that I came up with while preparing these notes:

[RAW] DATA + [LITERATE] CODE + [APPLIED] Stats = [DECISION] STORY

Why? Because data always come in "raw" form and have to be wrestled with. To do this, you need to be able to code (a little anyway). But in order to achieve the main goal, namely add value, process-oriented science has to come in, most importantly through systematic methods and the accompanying processes, to validate insights and help communicate results. Well, so far, so good.

In the following lecture, I will focus on four aspects of data science: the popularity it currently enjoys (and has enjoyed for the past 10 years), the skills required to "do data science", and the processes or activities involved in doing it. We will look at each of these with some examples.

How popular is data science?



How would you try to find out how popular data science is?
Image: Selfie by Cristina Zaragoza (Unsplash)

Exploring popularity

..../img/2_4th_july.gif

- **Search** (how? where?)
- Find relevant **models** (how?)
- Generate **primary** data (how?)
- Use **secondary** data (how?)

Question: Can you think of any issues with these methods?

1. Search - where? How?

- Google (Scholar) - disadvantage of Google searches?
- arxiv.org
- data science blogs (R-Blogger, Towards Data Science, Analytics Vidhya, R Weekly, DataCamp)

2. Find relevant models - what is that?

- Metaphors are models
- Mathematical model may not exist
- Example for models?

3. Generate primary data

- Which measures are used?
- Which methods are used?

4. Look secondary data

- public?
- Valid?
- How do you validate?

Example: social networking analysis - Predicting Tie Strength (2009).

Paper: <https://1drv.ms/b/s!AhEvK3qWokrvqz6uRFC1uk1LE0W5>

This paper uses a model to distinguish between weak and strong ties (with over 85% accuracy) based on a parametrization (= features to establish splitting the data) and a linear model (= assumption that the predictive variables are linearly correlated). Data science is used to address questions hidden in the data, such as how users relate to one another in social media, how they behave, perhaps even why they do what they do (= statistical inference).

Image: Google doodle 4th July 2022

"The sexiest job"

The image consists of two side-by-side screenshots from the Harvard Business Review website. Both screenshots show articles related to data scientists.

Left Screenshot (Original Article):

- Title:** Data Scientist: The Sexiest Job of the 21st Century
- Author:** Thomas H. Davenport and DJ Patil
- Date:** From the Magazine (October 2010)
- Image:** A colorful graphic showing a network of circles and lines, representing data connections.
- Text:** Andrew J. Blaustein, silk screen on a page from a high school yearbook, 8.5" x 12", 2011. Tamar Cohen.
- Summary:** Back in the 1990s, computer engineer and Wall Street "quant" were the hot occupations in business. Today data scientists are the hires firms are competing to make. As companies wrestle with unprecedented volumes... [more](#)

Right Screenshot (Follow-up Article):

- Title:** Is Data Scientist Still the Sexiest Job of the 21st Century?
- Author:** Thomas H. Davenport and DJ Patil
- Date:** July 15, 2022
- Image:** A graphic featuring a pair of glasses and a pie chart.
- Text:** July 15, 2022. [more](#)
- Summary:** Ten years ago, the authors posited that being a data scientist was the "sexiest job of the 21st century." A decade later, does the claim stand up? The job has grown in popularity and is generally well-paid, and the field is projected to experience more growth than almost... [more](#)

What do you think has changed since 2012? What has changed since 2022?

Image: Harvard Business Review covers 2012 and 2022

In the graph from trends.google.com, "numbers represent search interest relative to the highest point on the chart for the given region [worldwide] and time [since logging trends in 2004]." The trend increased is noticeable. It peaked in March 2022 (Source: Google Trends).

In October 2012, almost 10 years ago, Davenport and Patil published "Data Scientist: The Sexiest Job of the 21st Century" and put the term on the map.

What has changed since 2012?

1. (According to Davenport/Patil, 2022)

- **Demand** in 2012 restricted to a few cities, startups, tech firms
- Data scientists in 2012 were **science PhDs**, exceptional at math, who knew how to code
- Data scientists now need to develop **AI models**
- By 2019, postings on career site Indeed had risen by 256%
- Projected 15% increase from 2019 to 2029
- Lack of "data-driven cultures" (no use for data insights)
- Turnover is high (data scientists often don't stay long)
- Data science is better institutionalized (= widely accepted)
- Diversification and proliferation of roles (many skills needed)
- Changes in technology (like AutoML, MLOps tools)
- Need for an ethical dimension widely acknowledged (politicized)

2. Other changes that might have affected data science:

- COVID-19 pandemic (2020-2022)
- Rise of cloud computing, quantum computing, deep learning
- Political divide deepened (immigration, abortion, gun laws)

3. Since 2022:

- Generative AI (with Large Language Models) became popular
- ChatGPT the leading application, has over 180 mio active users
- An endless number of AI apps are built into other applications
- Generative AI is a subfield of data science: deep [machine] learning

The definition of sexy (for scientists)



»The best data scientists are product and process innovators and sometimes, developers of new data-discovery tools. That is the definition of sexy.«
-Gil Press (Forbes, 09/27/12)

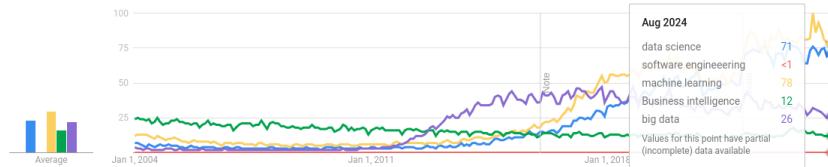
Image: Richard Feynman with drums (ca 1964)

Popularity contest

What do you think: which of these terms is most searched?

1. Big data?
2. Business intelligence?
3. Software engineering?
4. Data science?
5. Machine learning?

Most searched

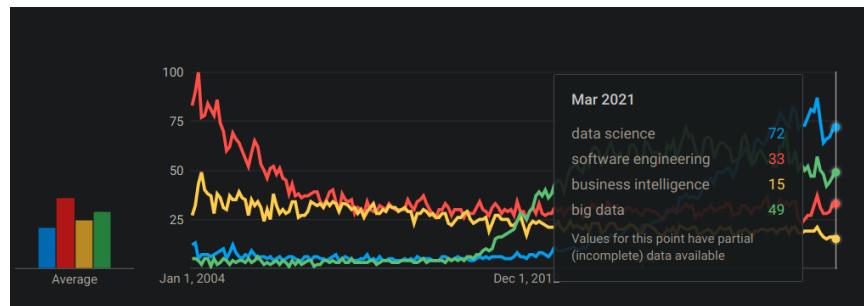


How do you like the visualization?

#+begin_notes Image Google Trends, August 2024

1. Bar chart (averages) difficult to read (percentages are missing).
2. List follows the search order, not the results
3. Grid lines (vertical lines) could improve reading

Three years ago...

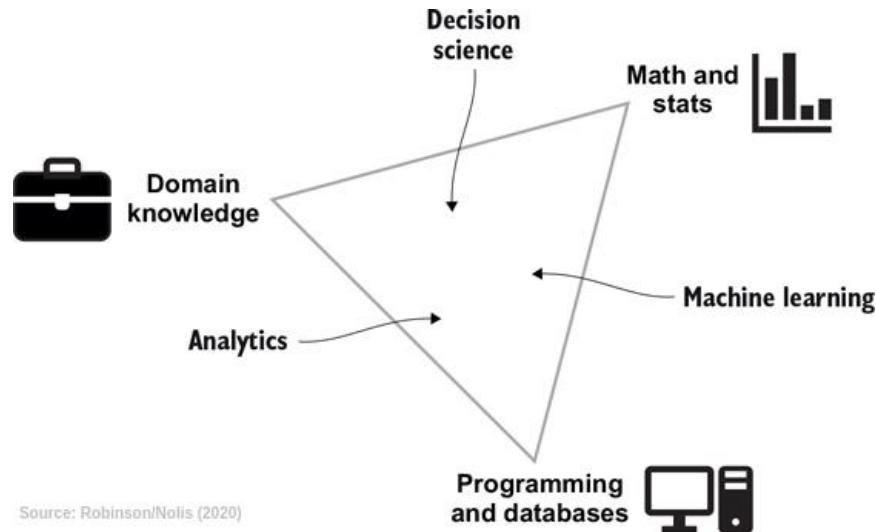


- "Software engineering" way down
- "Big data" down
- "Data science" constant

Image: Google Trends, March 2021

What changed: Web searches in "software engineering" went way down, interest in "big data" waned, relative position of "data science" constant.

WHICH SKILLS DO YOU NEED?



Give some examples for any of these skills!

- What do you know for example if you have "domain knowledge"?
- Which professional activities correspond to "math and stats"?
- What kind of "programming" would you have to do?

The three skill areas in the figure (from Robinson/Nolis (2020)) give rise to different tasks and problem settings:

Skill	Sample area	Sample activity	Sample analysis
Domain knowledge	Marketing Education Finance	Analyze customer data Learner data Investment data	What do customers like? How did students learn? Which stock performed?
Coding & databases	R, Python, SQL Cloud computing RStudio, Emacs Package creation	Analyze/automate/query Share data and code Improve your workflow Write new functions	Count customers by type Work in virtual teams Create a notebook Distribute package
Maths & stats	Data structure Model building Distribution	Data wrangling Linear regression Check significance	Check data tidiness Fit line graph to data Apply t-test

Between two of these areas each are application areas:

1. Domain knowledge and statistics support **decision science**. See info-graphic (source: Bobriakov 2019).
2. Data analytics are the result of applying **database programming** (e.g. with SQL) to domain knowledge problems(this is also sometimes called '**business intelligence**' or BI).
3. Programming, maths and statistics give rise to various machine learning (ML) techniques concerned in particular with **prediction** and automatic pattern recognition.

What are your skills?

URL: tinyurl.com/data-science-skills

1. In which **domain** do you have knowledge?
2. Which **decisions** have you made?
3. What do you know in **maths** and **stats**
4. Which programming **languages** do you know?
5. Which analytics **tools** have you used?
6. What are your skills in **machine learning**?

Compare: "My IT Skills Stack"

1. Problem solving skills:
 - Understand the problem: the conditions, the unknowns, the data. Of these, I am particularly good with data.
 - Design a plan of attack (e.g. by modeling - abstracting from the details to identify one or more routes or options)
 - Carry out the plan of attack: this is execution. Probably my least favorite part (often, when I see the solution path, I get bored). But I can do it, and it's satisfying to finish something.
 - Look back, review and discuss your solution. I am especially good at this type of postmortem analysis - it's probably what I use most when it comes to teaching stuff.
2. Computational thinking skills

- 10 programming languages - recommended: SQL and R

3. Data literacy skills

- Wikipedia definition is not bad: "Ability to understand, create, and communicate data as information." (I.e. structured data)
- Use of visualization and storytelling techniques
- Business process modeling

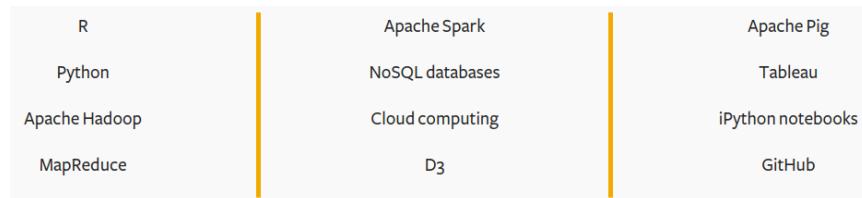
4. Communication skills

- team / leadership experience

5. Tool skills

- I love tools
- In my courses usually use about 20 different IT tools

What are technical data science skills?

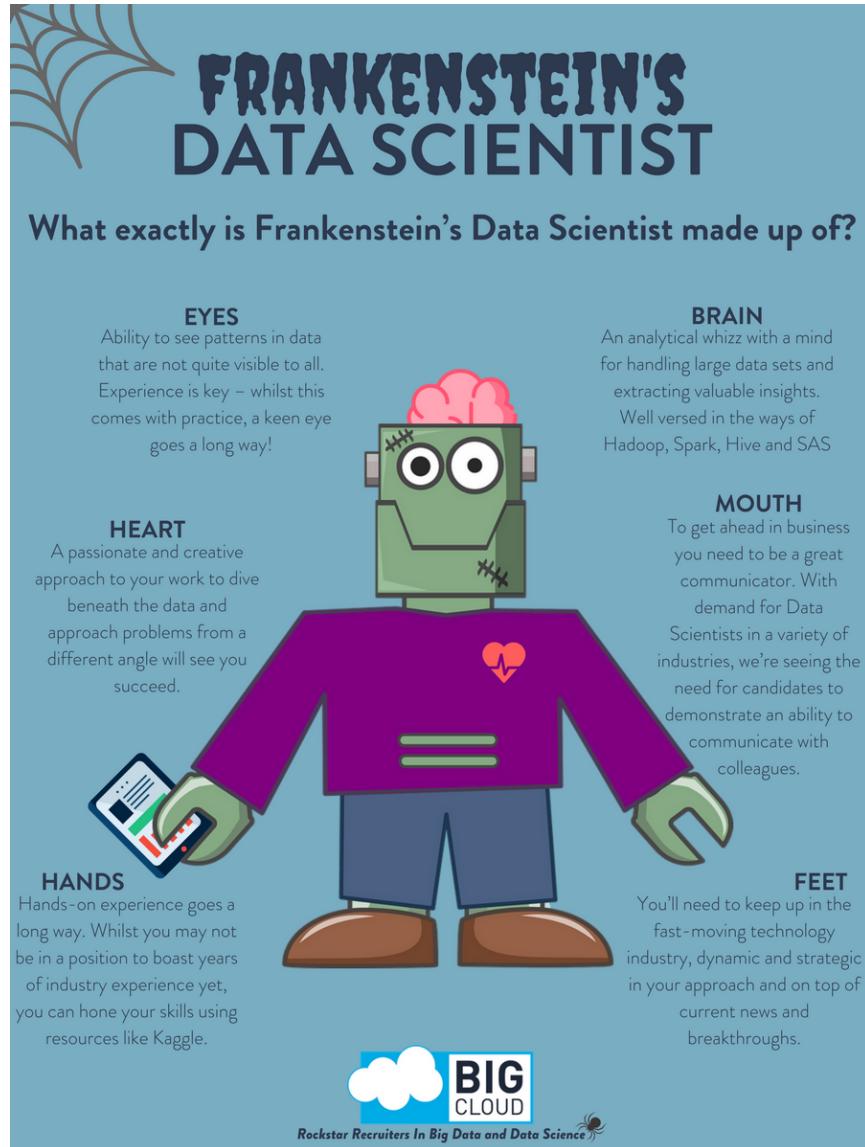


Have you heard of any of these?

Tip: when you come across products you don't know, make it a habit to look them up - knowing the names and what they stand for will help you anchor yourself in anything you read, and the most important products, which are most talked about, are often talked about for a reason - e.g. because they represent an innovation and/or an advantage. By knowing the products, you can also learn something about the innovation. This dependency on products also shows that both computer and data science are crafts.

TOOL	PURPOSE	TOOL	PURPOSE
D3.js	Visualization	Apache Hadoop	distributed computing
Apache Spark	Analytics engine	MapReduce	Google scalability
Apache Pig	Analytics platform	NoSQL	Unstructured big data
Tableau	Visualization	iPython nb	Literate Programming
GitHub	Version control		

All of these are also potential project topics!

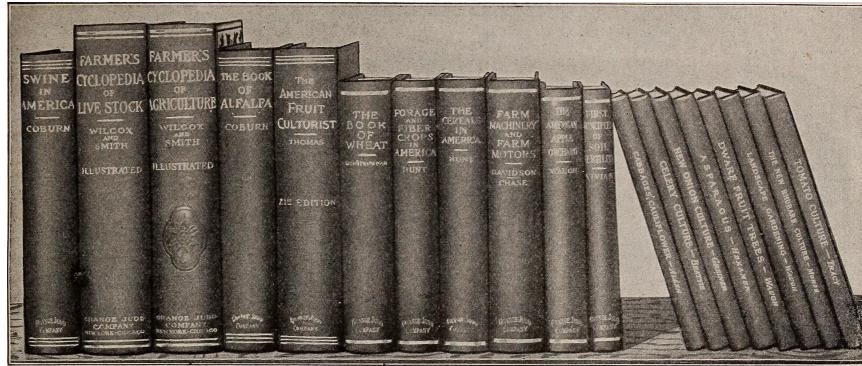


Source: datasciencecentral.com

"Frankenstein's monster" (based on the novel by "Frankenstein, or The Modern Prometheus", by Mary Shelley, 1818) is used in the figure as a metaphor for a working data scientist. it is a rich metaphor with many connotations.

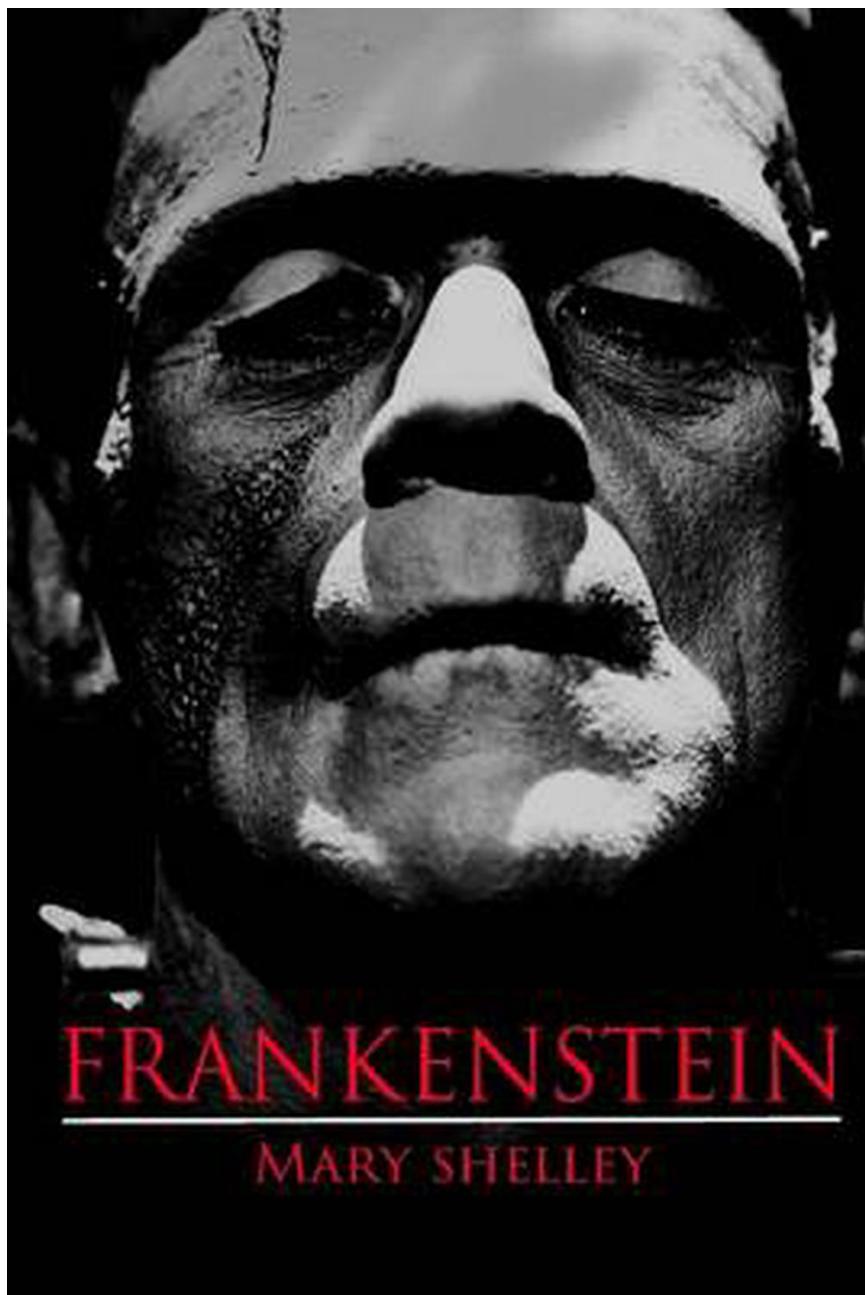
- "Eyes": experience with detecting data patterns. to do this actually with your eyes is unlikely - you need some tools for that, but you also need experience to know which tools will work. example: `head(dataset)` only prints the first 6 rows of a dataset giving you an idea of the type of data in the dataset.
- "Heart": passion for and creativity with data. "passion" is perhaps more relevant for the data's origin and for what you can do with well interpreted data - namely change the world! example: hans rosling's gapminder animations (and his passionate storytelling, demonstrated e.g. in Hans Rosling's TED videos).
- "Hands": domain knowledge gained by working in an industry for years, supported by activity in communities like kaggle (owned by google since 2017), which hosts datasets, notebooks and ml competitions.
- "Brain": analytical mindset and knowledge of analysis tools (none of the tools mentioned here, hadoop, spark, hive - a data warehouse - or sas - another statistical analysis workbench - are necessary - they are merely nice to know). how do you know that you have this kind of brain? e.g. if you enjoy getting quantitative (number-based) answers and if you like visualizations of complex or complicated data (like the gapminder data). also, if you like programming or maths, you've likely got such a brain.
- "Mouth": communication with colleagues - but not only. in fact, especially being able to communicate with people who are not your colleagues (so they are perhaps very different from you) is key. this is another way of saying that you need to be able to "tell a story" after data analysis (e.g. Prevos 2020).
- "Feet": data science is a very fast-moving technology field, especially its "machine learning" offshoot (which is not part of this course) - cp. Kozyrkov 2019. you need to keep on top of the available information. at the same time, there is too much to take in and digest - this means that it is very important to have a sound understanding of the foundations of data science.

A brain for numbers



- What if you don't have a "brain for numbers"?
 - What if graphs scare you because of the underlying math?
 - What if you like novels but hate manuals?
 - What if you actually hate computers and machines?

Metaphors



- What are the connotations of "Frankenstein's Data Scientist"?

- Do you find this metaphor apt or not?
- Which metaphor would you have chosen?

...youR tuRn: What are the connotations of using "Frankenstein's monster" as a metaphor for "data scientist"? Metaphors are especially important when definitions are not easily forthcoming, are confused or not standardized (all of which is the case for data science). Metaphors are a type of model.

WHAT'S THE DEMAND?

Is There a Demand for Data Scientists in 2024?

A few years ago, data science was a fancy word with little meaning to the average person. Now, the world has a better understanding of what the field encompasses. With a growing awareness of the significance of data in all areas of work and life and the commercialization of AI solutions, the demand for jobs in data science is on the rise. But the job requirements are changing.

According to the [U.S. Bureau of Labor Statistics](#), data scientist positions will continue to be among the fastest-growing jobs in 2024. The projected increase in job openings from 2022 to 2032 is 35%.

In addition, the [World Economic Forum's Future of Jobs 2023](#) report estimates that by 2027, the demand for AI and machine learning specialists will increase by 40%, and for data analysts, scientists, engineers, BI analysts, and other big data and database professionals will grow by 30%–35%.

Search a job portal for "data scientist" jobs. Find anything?

The value of statistics depends on the exact definitions of the job, on the ability of business to recruit exactly for what they want etc. However, as a rule, you can never go wrong with growing your skill stack, especially with regard to STEM skills, and within these especially with regard to your ability to analyse data quantitatively - which is what data science boils down to. For more details on "data science careers", see Robinson/Nolis (2020).

Mathematics, especially statistics, programming and databases are the skill-based disciplines that you need to master. Having said that: "mastering" could easily take not one, but several life times, and you need to begin somewhere. If you do this in earnest, you'll soon find that you start learning faster and faster the more connections with what you already know you can make.] Here is a (free) book called, incidentally, "Foundations of Data Science" (Blum et al 2015, 466 p.). It includes some geometry, graph theory,

linear algebra, markov chains, and a variety of algorithms for "massive data problems" like streaming, sketching and sampling.

Source: 365datascience.com/career-advice/data-scientist-job-market/

Job profiles (according to datacamp)



Data Engineer	Data Analyst	Data Scientist	Machine Learning Scientist
Store and maintain data	Visualize and describe data	Gain insights from data	Predict with data
SQL + Java/Scala/Python	SQL + BI Tools + Spreadsheets	Python/R	Python/R

- Who would you rather be?
- Why?
- Which job is most in demand?

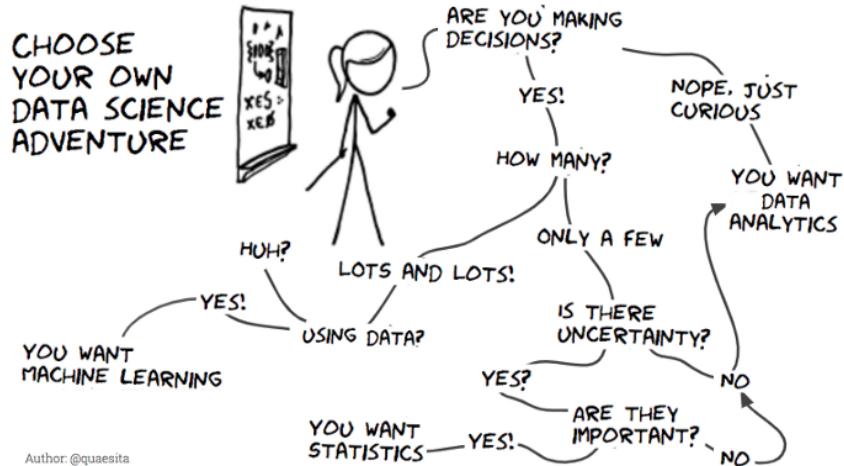
Image: Source: DataCamp, Understanding Data Science

Introductory DataCamp courses on data science "for everyone" (that is, without being tied to one of the three dominant languages - Python, R, or SQL), contain a job profile section to help users find their professional data science niche.

The figure shows four such profiles from a 2020 course. What is notably missing here is the maths and/or CS or software engineering knowledge required or desirable to fill these roles. But there are also people who say that you best come to a firm as a general-purpose computer scientist and then learn any of these on the job depending on the needs and the available experience.

In 2024, you'll probably find most jobs in the "Data analysts" category but I have not checked that.

Data science problems



- **Stats:** few important decisions, high uncertainty
- **Analytics:** explorative or explanatory
- **ML:** many decisions involving big data

Image: Cartoon by Cassie Kozyrkov (@quaesita)

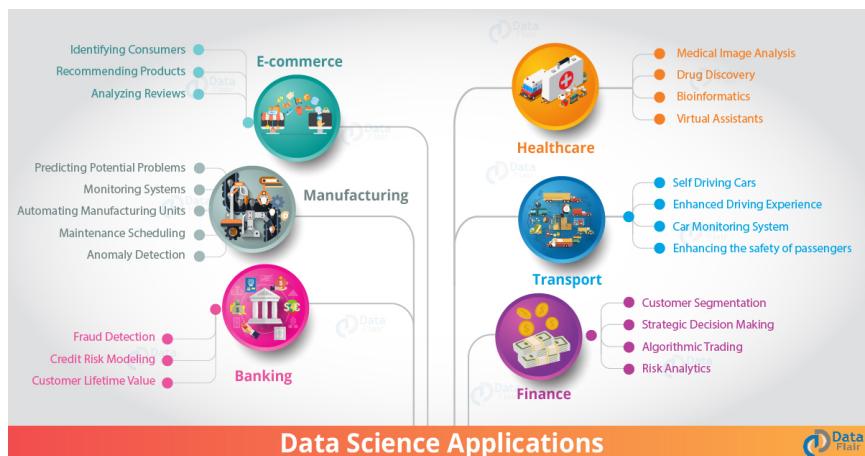
The cartoon in the figure is by Google's (former) head of "decision intelligence", Cassie Kozyrkov (2018) (in the style of xkcd). She has a specific, business- and decision-oriented idea of the purpose of data science, which I share: data science is there to help you make decisions. The option tree shown distinguishes three sub-fields of data science: data analytics, statistics and machine learning. It asks if you're "making decisions" at the start (many, few, hardly any), it quickly focuses on the type of data (few vs big) and the 'uncertainty' and 'importance' of the decisions. This is still a data-centric, not a decision-centric taxonomy. A focus on the latter would allow for many more options (e.g. strategic vs. tactical, organizational vs. managerial, routine vs. exceptional decisions etc.) Hence, for decision science, this kind of breakdown is not very useful.

The dominance of "big data" has also been doubted, especially when it comes to making (business) decisions. "Small [not big] data" (Saklani, 2017) and "thick [qualitative, descriptive] data" may be just as good depending on what you want to know. The article by Chiu (2020) is a bit of a history hack (in the scholarly sense) but it raises some good points.

Brandon Rohrer, [then] a data scientist at Microsoft, has addressed this question in a 3-part series of short articles (Rohrer, 2015a, 2015b, 2015c). His examples are a more specific, especially because he also says which family of algorithms match which type of data-related question. It is too early for us to discuss his taxonomy but at the end of the course, you should have a better idea about what you can do with data science tools.

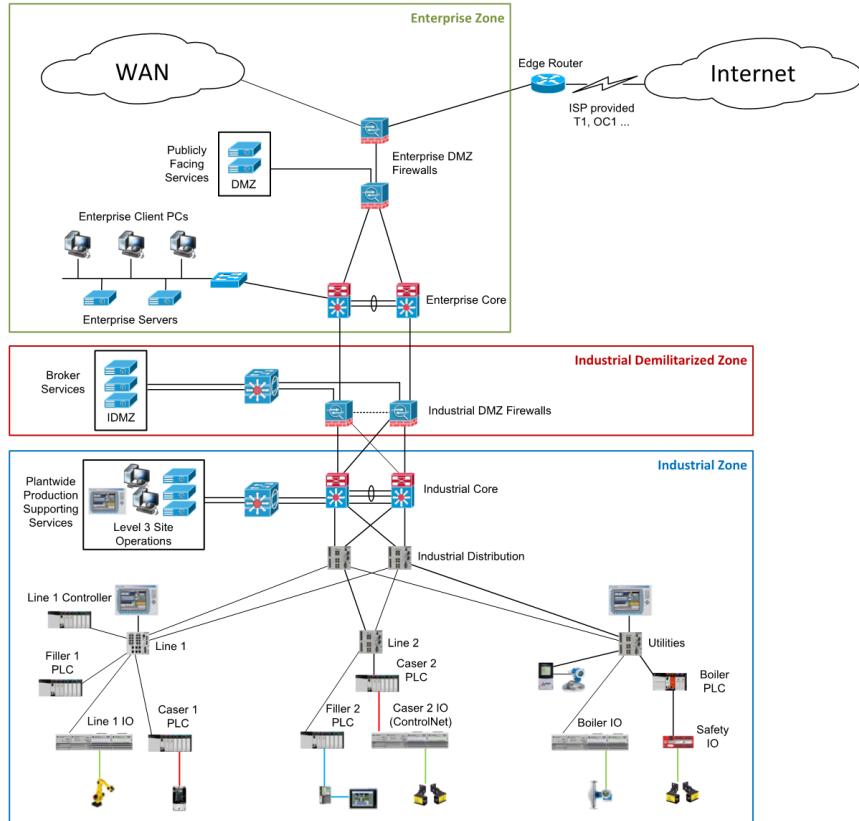
... youR tuRN: Think about any decision you make - what are the steps you go through? Do they amount to a "data science adventure" as shown in the figure - why (or why not)?

Data science applications



Source: data-flair.training

Example 1: cybersecurity



- Problem: how to secure critical digital infrastructure
 - Solution: Industrial Control System
 - DS: EDA (user data), simulation (sample data)
- Source: Industrial Cybersecurity (2017)

Example 2: time series analysis & text mining

```
Jul 16 09:10:11 linux systemd-timesyncd[1219]: Synchronized to time server [2001:67c:1560:8003::c8]:123 (ntp.ubuntu.com).
Jul 16 09:11:41 linux kdeconnectd.desktop[6764]: kdeconnect.core: TCP connection done (I'm the existing device)
Jul 16 09:11:41 linux kdeconnectd.desktop[6764]: kdeconnect.core: Starting server ssl (I'm the client TCP socket)
Jul 16 09:11:41 linux kdeconnectd.desktop[6764]: kdeconnect.core: TCP connection done (I'm the existing device)
Jul 16 09:11:41 linux kdeconnectd.desktop[6764]: kdeconnect.core: Starting server ssl (I'm the client TCP socket)
Jul 16 09:11:41 linux kdeconnectd.desktop[6764]: kdeconnect.core: Socket successfully established an SSL connection
Jul 16 09:11:41 linux kdeconnectd.desktop[6764]: kdeconnect.core: It is a known device "Xperia L2"
```

- Data: Linux /var/log/syslog event log
- Problem: Textual time series data

- Solution: Text or process mining of the event log data

All system components continuously write data protocols in the form of simple event logs, which you can view easily on Linux systems e.g. on Ubuntu. Check available system logs with `ls -la /var/log/`. The figure shows a sample section from my computer's system log in `/var/log/syslog`.

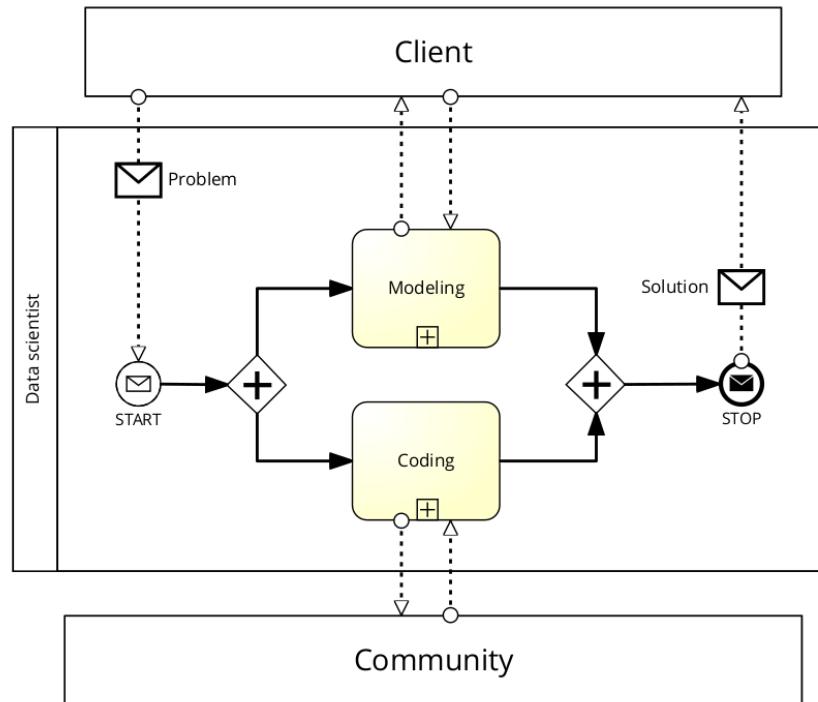
This excerpt shows how and when the computer clock was set remotely, and the starting of various servers and one socket where my mobile phone ("Xperia L2") was connected.

The language we're about to use in this course (and in the follow up course on machine learning), R, is well suited for rapid interactive exploration of datasets such as this one. The two immediately relevant problem areas are "text mining" (notice that all system files are human-readable to aid debugging), and "time series analysis" (event logs are time series).

Text mining is considered a part of "Natural Language Processing", and Time Series Analysis is also really important in finance, e.g. when analysing portfolio performance.

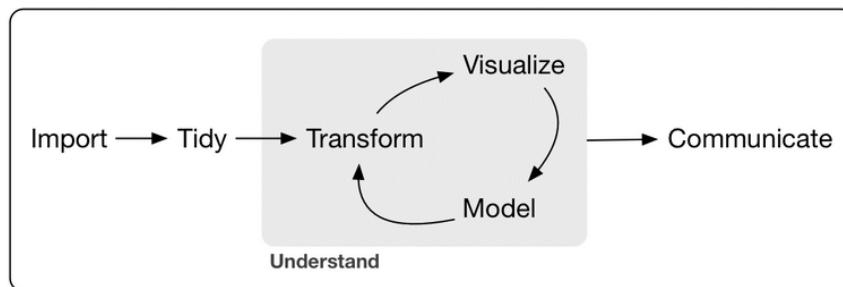
A separate technique (not immediately part of an R programming course) is "process mining". I'll be talking about that in data visualization.

WHAT IS THE DATA SCIENCE PROCESS?



Source: Birkenkrahe (2021)

Exploratory data analysis (eda) process model



Here is my interactive BPMN version.

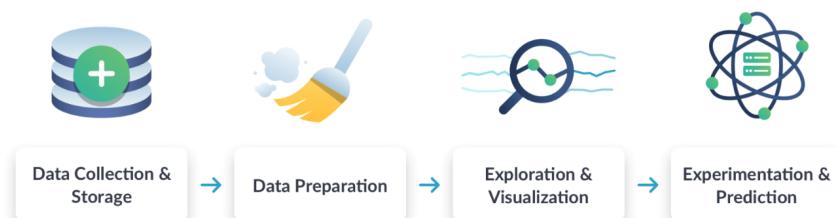
Image source: Wickham/Grolemund (2017)

The figure shows a process that begins with raw data. Such data are usually not formatted as "tidy" data, i.e. "each row represents one observation and columns represent the different variables available for each of these

observations" (Irizarry 2020). This is also the tabular format, which is usual for storing data in relational databases for analysis with SQL.

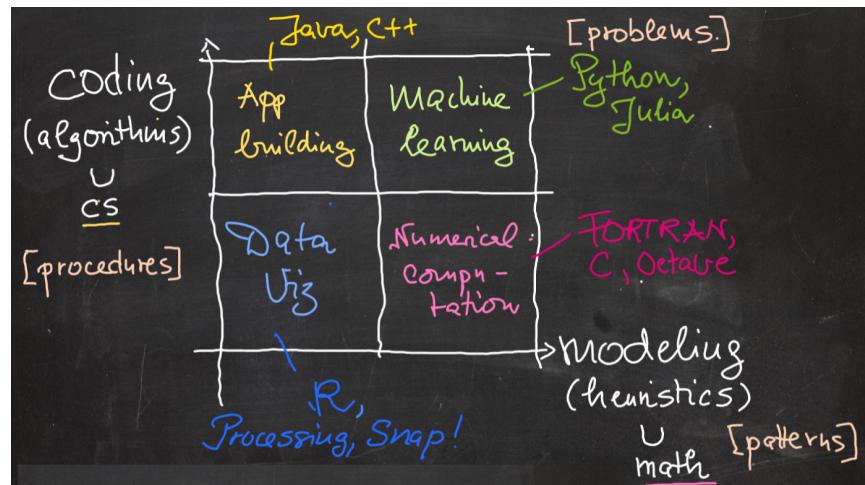
Once we have tidy data, an (often repeated) sub-process begins: "transform" refers to any operation on the dataset that helps us understand the data better. Depending on the size of the data tables, we will use different methods of visualization to make underlying structure visible. But visualization does not always have to be graphical - it could also be making a table, or creating a metaphor.

Data science workflow



Source: Understanding Data Science (DataCamp)

A model for learning data science



- Algorithmic vs heuristic
- Coding vs modeling

- Dashboards vs. Prediction

Image: Talk@Lyon College (Birkenkrahe, 2021)

CONCEPT SUMMARY

- Data science is used for **decision support, process analytics** and **machine learning**.
- Data science makes use of **domain knowledge** - experience in a particular field of business.
- The **job market** (for data science) is great
- The data science **process** includes modeling, visualizing, and communicating data analysis results.

REFERENCES

1. Blum A/Hopcroft J/Kannan R (4 Jan 2018). Foundations of Data Science - Cornell U. Online: cornell.edu.
2. Bobriakov I (16 Apr 2020). Data Science vs. Decision Science [Info-graphic]. Online: medium.com/@bobriakov.
3. Bolles R and Brooks K (2021). What color is your parachute? Online: <https://www.parachutebook.com/>
4. Chiu J (17 Aug 2020). Why Data Doesn't Have to Be That Big. Online: datacamp.com.
5. Davenport TH/Patil DJ (2012). Data Scientist: The Sexiest Job of the 21st Century. Online: hbr.org.
6. Davenport TH/Patil DJ (July 15, 2022). Is Data Scientist Still the Sexiest Job of the 21st Century? Online: hbr.org.
7. Devlin K (1 Jan 2017). Number Sense: the most important mathematical concept in 21st Century K-12 education. Online: huffpost.com.
8. Gapminder Foundation (15 Dec 2014). DON'T PANIC - Hans Rosling showing the facts about population. Online: youtube.com

9. Grolemund G/Wickham H (2017). R for Data Science. O'Reilly.
10. Irizarry R (2020). Introduction to Data Science. CRC Press.
11. Kozlykov C (10 Aug 2018). What on earth is data science? Online: [hackernoon.com](https://hackernoon.com/what-on-earth-is-data-science).
12. Kozlykov C (22 May 2019). Automated Inspiration. Online: [Forbes.com](https://www.forbes.com/sites/kozlykov/2019/05/22/automated-inspiration/#:~:text=Data%20science%20is%20the%20process,of%20data%20and%20information%20to%20inform%20decisions%20and%20actions.)].
13. Knuth D (1992). Literate Programming. Stanford, Center for the Study of Language and Information Lecture Notes 27.
14. Myers A (28 Apr 2020). Data Science Notebooks - A Primer. Online: [medium.com/memory-leak](https://medium.com/memory-leak/data-science-notebooks-a-primer-5a2a2a2a2a2a).
15. Porras E M (18 Jul 2018). Linear Regression in R. Online: [datacamp.com](https://www.datacamp.com/courses/linear-regression-in-r).
16. Prevost P (14 Aug 2020). Storytelling with Data: Visualising the Receding Sea Ice Sheets. Online: lucidmanager.org]].
17. Robinson E/Nolis, J (2020). Build a Career in Data Science. Manning.
18. Rohrer B (2015a). What Can Data Science Do For Me? Online: [microsoft.com](https://www.microsoft.com).
19. Rohrer B (2015b). What Types of Questions Can Data Science Answer? Online: [microsoft.com](https://www.microsoft.com).
20. Rohrer B (2015c). Which Algorithm Family Can Answer My Question? Online: [microsoft.com](https://www.microsoft.com).
21. Saklani P (19 Jul 2017). Sometimes “Small Data” Is Enough to Create Smart Products. Online: hbr.org.
22. Sarkar DJ (12 Sept 2018). A Comprehensive Guide to the Grammar of Graphics for Effective Visualization of Multi-dimensional Data. Online: towardsdatascience.com
23. Scherpereel CM (2006). Decision orders: A decision taxonomy. In: Management Decision 44(1):123-136.
24. Wing JM (2 Jul 2019). The data life cycle. Harvard Data Science Review. Online: hdsr.mitpress.mit.edu.