

Spatial Data Science

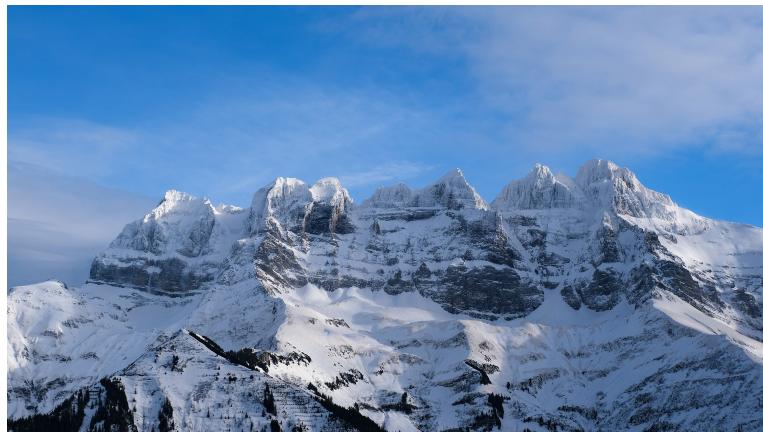
Introduction

(YMS31303)

Lecture 1

Hans Hoogenboom





tem3d_library.py

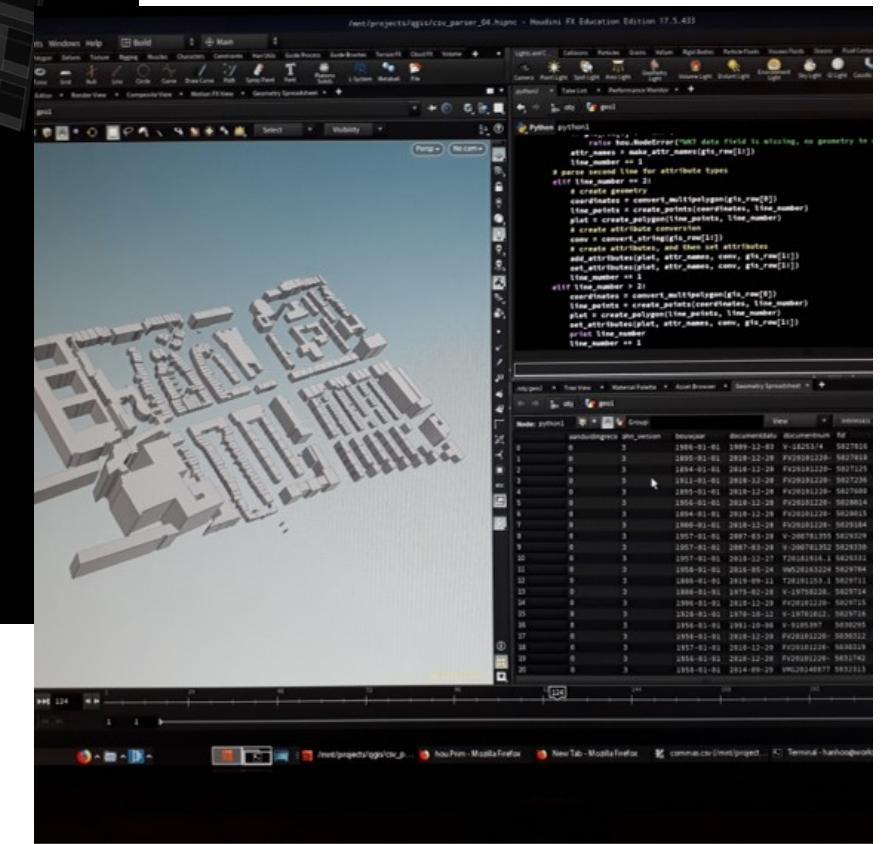
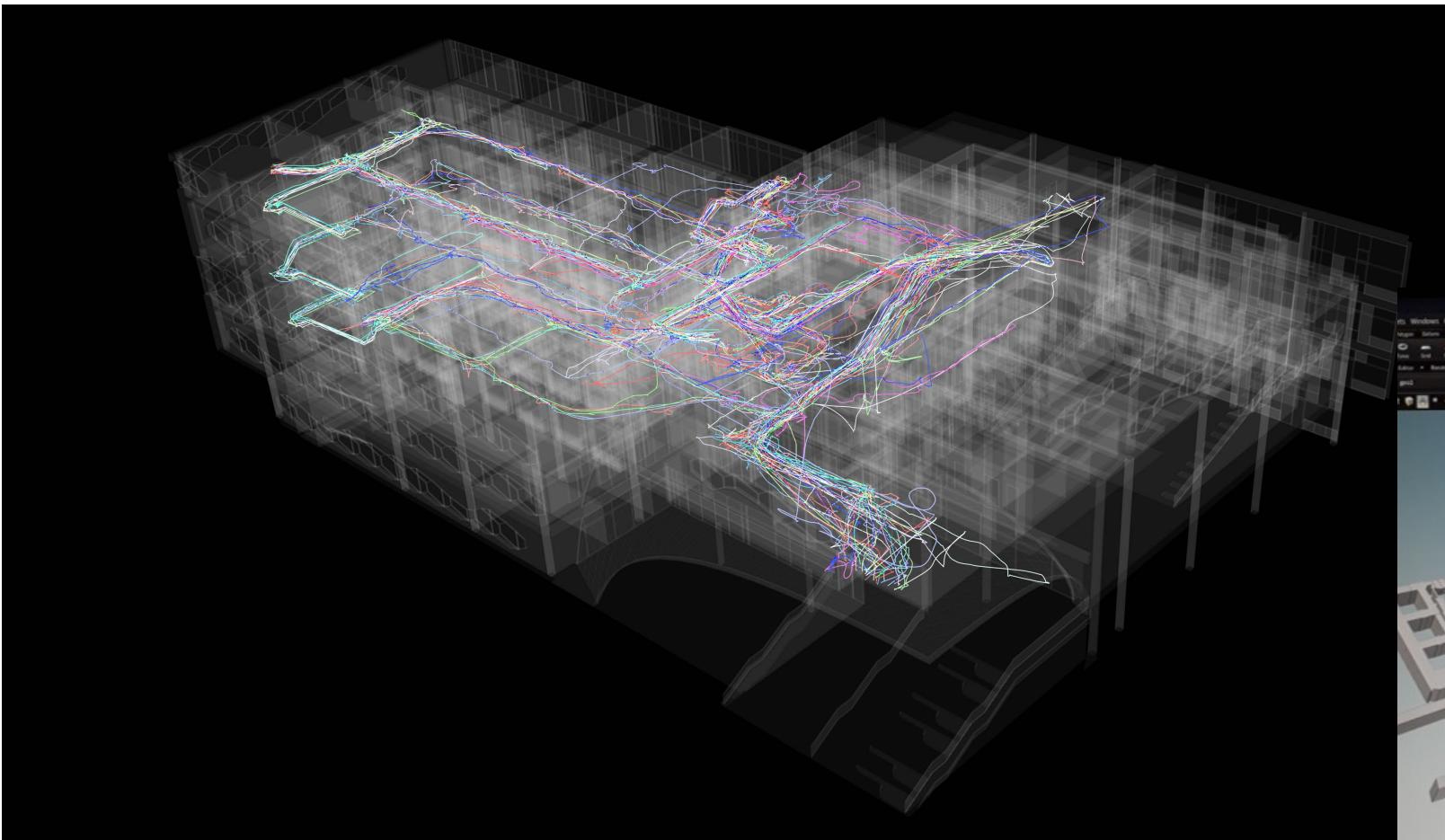
Finite Element Method solver for Houdini 12

Source master cfb1aa8 Full commit

femmes / hip / fem3d_library.py

```
1098     res[i] = sums / matx[i,i]
1099     #start the backsubstitution
1100     res[n] = res[n] / matx[n,n]
1101     for i in range(n-1,-1,-1):
1102         sums = 0.0
1103         #j is the column
1104         for j in range(n, i, -1):
1105             sums = sums + matx[j,i] * res[j]
1106         sums = res[i] - sums
1107         res[i] = sums / matx[i,i]
1108     return res
1109
1110 #below two iterative methods
1111
1112 ##### GaussSeidel iterative method (add preconditioning pivoting?) #####
1113 # GaussSeidl iterative method (add preconditioning pivoting?) #
1114 #####
1115 def GaussSeidel(matx, sol, res0, max_itter, tol):
1116     size = len(sol)
1117     res1 = Vector(size, 0.0)
1118     status = "normal"
1119     for k in range(size):
1120         if abs(matx[k,k]) < tol:
1121             print "Error"
1122             status = "singular"
1123     iter_num = 0
1124     if status == "normal":
1125         while (status == "normal") and (iter_num < max_itter):
1126             eps = 0.0
1127             for i in range(size):
1128                 sum = 0.0
1129                 for j in range(0, i):
```







Centre for Urban Science & Policy

We are a transdisciplinary research group working to advance urban research, planning and policy in a way that strives for just and equitable outcomes for communities. We use a mix of computational spatial science and qualitative participatory methods to investigate how social, economic, environmental and political processes shape cities. Our goal is to develop a body of computational approaches and curate evidence that facilitates an integrated systems-based approach for urban planning.

[EXPLORE OUR WORK](#)

Space and Place

Cities & Social Justice

Understand the complex nature of urban spaces in transformation

Understand how inequalities intersect with space

Propose guidelines and develop tools to support public participation and citizen action

Identify inequalities associated with lack of recognition and legitimacy

Intersectionality

Democratisation of urban design, planning & policy

Teaching Support

Ka Yi Chua
2nd year AMS

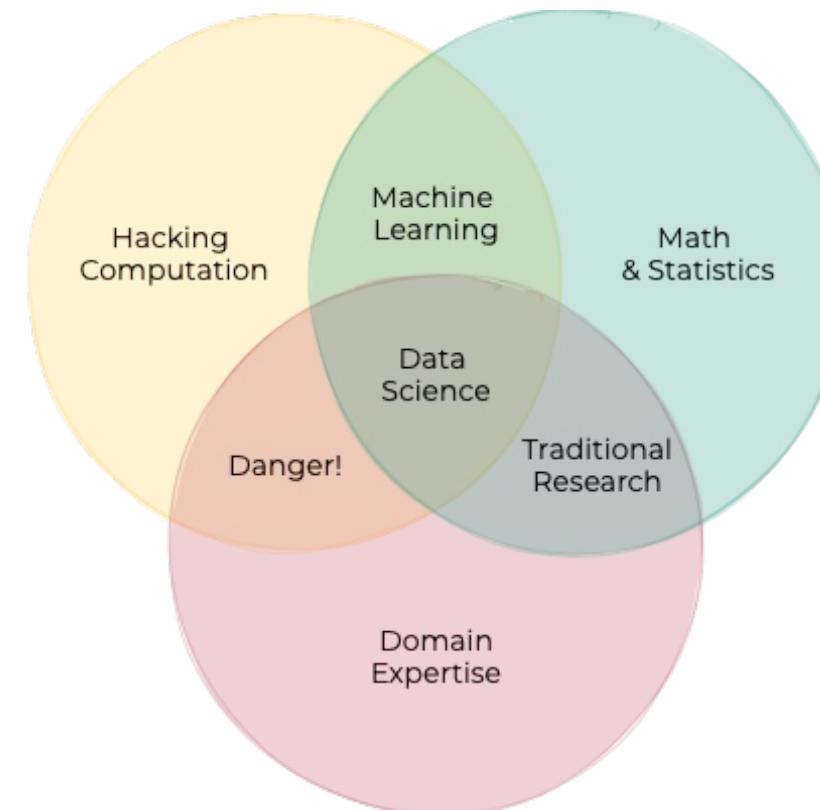
Today

- Introduction to the Course
- Tools - Python (and Conda)
- Post break, Intro part II

Introduction to the Course

More stats than a GIS course... more GIS than a stats course

With a few twists!



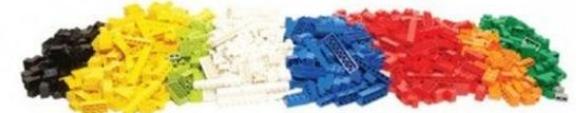
After this course

You will be able to...

- **Obtain**: Obtaining data from multiple **open** data sources.
- **Scrub**: Data cleaning, munging, sampling to consolidate all information into a dataset that is manageable, informative and relates to your problem.
- **Explore**: Exploratory data analysis to make sense of what your data is trying to say.
- **Model**: Estimation and modelling based on statistical tools such as regression and clustering.
- **Interpret**: Communicating results and reflections through visualisation, storytelling and interpretable summaries.



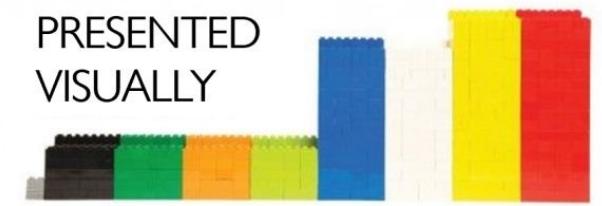
SORTED



ARRANGED



PRESENTED VISUALLY



EXPLAINED WITH A STORY



Can you find the source for me?

Philosophy

- (Lots of) **methods** and techniques
 - General overview
 - Intuition
 - Very little math
 - Lots of ways to continue learning more
- Emphasis on the **application** and **use**
- Close connection to “**real world**” applications

Format – Ways of Working

Six weeks of:

- **Prep. Materials**: videos, podcasts, articles... 1h. approx. (most recommended!)
- **2x 1h. Lecture**: interaction, concepts, methods, examples
- **1x 2h. Computer labs**: hands-on, application of concepts, Python (*highly employable*)
- **Extra material**: how to go beyond the minimum*

*(not necessary for the course but useful in life)

Content

- **Weeks 1-3:** “big picture” lectures + introduction to computational tools (learning curve) + lots and lots of data
- **Weeks 4-6:** lots of spatial, network and machine learning concepts
- **Weeks 7:** prepare an awesome final project

An overview of the course

- **Weeks 1-3:** “big picture” lectures + introduction to computational tools (learning curve) + lots and lots of data
-
- **Weeks 4-6:** lots of spatial, network and machine learning concepts + visualisations
- **Weeks 7:** prepare an awesome final project

Week	Date	Morning		Afternoon	
		10.00 - 11.00	11.00 - 12.00	13.00 - 14.00	14.00 - 15.00
1	04-Sep	Welcome /Introduction		Python Basics	
2	11-Sep	Data Grammar for Spatial & Urban Data		Tabular Data (Numpy & Pandas)	
3	18-Sep	Data Engineering		Numpy & Seaborn	
4	25-Sep	Exploratory Spatial Data Analysis		Plotting and Vizualization of Data	
5	02-Oct	Networks		Shortest Paths (Networkx)	

Week	Date	Morning		Afternoon	
		10.00 - 12.00	13.00 - 15.00		
6	09-Oct	Machine Learning for Everyone		Basic Python ML	
7	16-Oct			Hackathon (Final Project)	
8	23-Oct				
9	30-Oct				
10	06-Nov				

Logistics - Website

Course Material

- https://jhoogenboom.github.io/spatial-data-science/_index.html

Communication Channels

- With instructors and TA: Email
- Announcements on Brightspace

Submissions, Groups, Grades and Feedback

- Brightspace

The screenshot shows the homepage of the course website. At the top right are navigation icons: a magnifying glass for search, a gear for settings, a double arrow for refresh, a downward arrow for download, and a square with a dot for full screen. To the left of these is a menu icon (three horizontal lines). On the far right is a vertical sidebar with a blue header "Contents" and links to "Overview", "Learning Outcomes", "Career Prospects", and "Ethical Considerations". The main content area has a light gray background. At the top left of the content area is the QUSP logo (a stylized 'Q' and 'U' in a hexagon shape) next to the text "Centre for Urban Science & Policy". Below the logo is the course title "YMS31303 AMS Metropolitan Data 1". To the right of the title is a large, bold, dark gray section header "Overview". Under "Overview" is a detailed description of the course's purpose and content. Below the overview is a section titled "Learning Outcomes" with a list of learning objectives. At the bottom is a section titled "Career Prospects" with a list of career-related bullet points. On the left side of the content area, there is a sidebar with a blue header "YMS31303 Metropolitan Data 1" and a list of course modules: "Introduction", "Lectures", "Labs", "Assessment", "Software", "Helping Material", and "FAQ". Each module has a small icon to its left.

YMS31303 AMS Metropolitan Data 1

Overview

Urban planners, policymakers, and key decision-making stakeholders use data and data-based infrastructures to govern various urban systems from operations to planning, optimization, and distribution of resources. This course will introduce you to practices in data analysis and computational methods in the context of urban planning. It will illustrate how data can be used and misused, and how to critically evaluate datasets, models, and questions that arise from them. While learning how to collect, transform, and analyze data using machine learning techniques for understanding urban phenomena, you will learn about the process of data science and its positive and negative impacts on people and places.

Learning Outcomes

After successful completion of this course, you will be able to:

- Interpret and discuss spatial data sources that are usable and relatable for a problem presented.
- Transform spatial data and consolidate all information into a dataset that is manageable, informative, and relates to your problem.
- Describe and analyze the consolidated spatial datasets to support your problem with evidence.
- Apply models using statistical techniques and machine learning to infer results in the process of turning spatial data into valuable information.
- Report results and reflections through visualization, mapping, storytelling, and interpretable summaries, especially when faced with a new dataset.

Career Prospects

- Hopefully get great data-driven policy, governance, or civil society jobs in the future.
- Go on adventurous and sustainable journeys with an open mind.

Self-directed learning

This course is much more about “**learning to learn**” and problem solving rather than acquiring specific programming tricks or stats wizardry

- **Prepare** for the labs
- TAs will be present for abundant help and feedback.
- **Go over the notebooks** before the lecture and the computer lab
- If the first time you see a notebook is at the lab, you may find it difficult to follow on. The best thing to do is to prepare a set of questions to ask us.
- **Bring** questions, comments, feedback, (informed) rants to class/labs. The more you bring, the more we all learn.
- **Collaborate** (it’s **NOT** a zero-sum win!)

Python

- **General purpose** programming language
- “Sweet spot” between “*proof-of-concept*” and “*production-ready*”
- Industry standard: **GIS** (Esri, QGIS) and **Data Science** (510, World Bank, OECD, The Atlantic, Gemeente Den Haag...)

How many of you have written a line of computer code before?

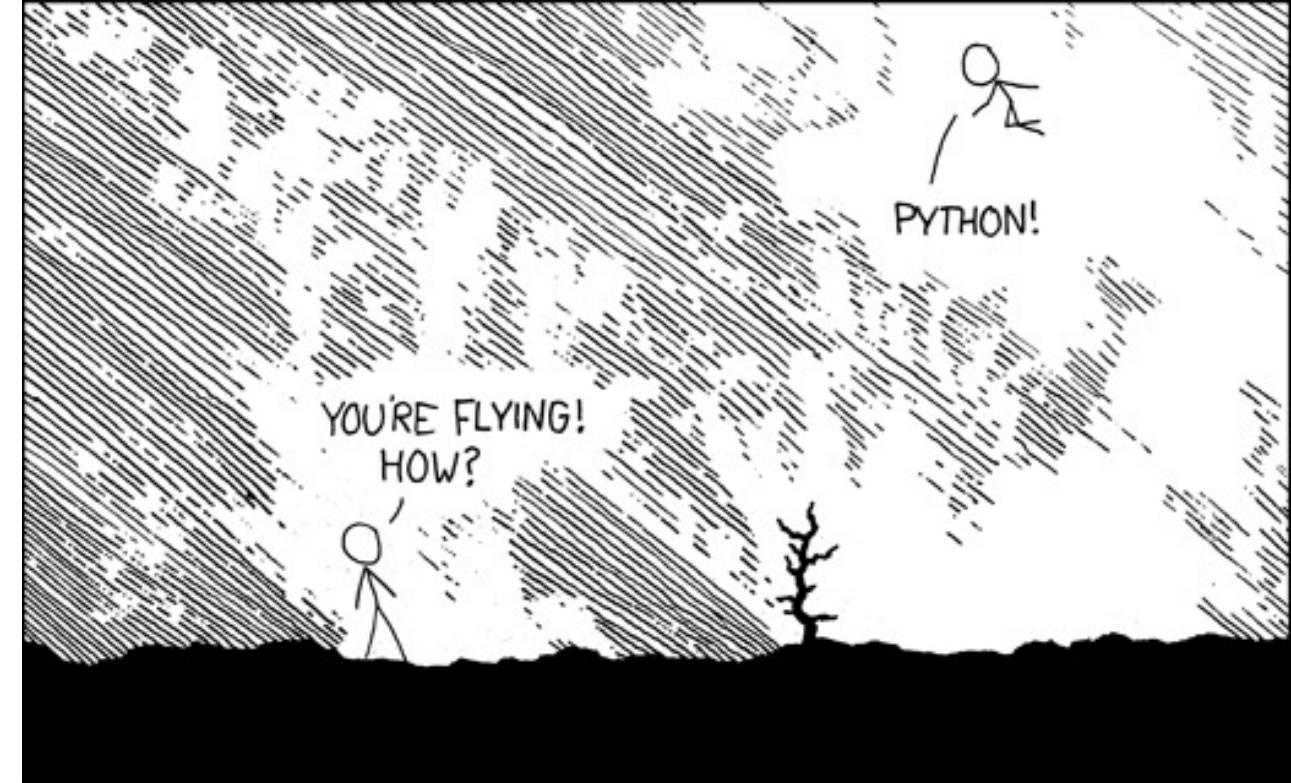


Figure under Creative Commons Attribution Non-Commercial 2.5 License



Assessments

Formative : These are ungraded

- 4 Group assignments
- 1 Hackathon

Discuss the completed homework with your peers using a list of Do's and Don'ts to evaluate each other's work

Assessments

Summative : There are 5 graded components that contribute to the final mark for the course as follows:

40% group assignments

- Assignment 1: Getting up and Running, Github Pages + homework (only once) (10%)
- Assignment 2: ? (10%)
- Assignment 3: ? (10%)
- Assignment 4: ? (10%)

60% Group assessment include

- Hackathon

Rubric for Assignments

- Assignments are graded based on **five criteria** if applicable
- The criteria have different weights that add up to 100%

Aspect (weight)	Level A- (0-5)	Level A (6)	Level B (7)	Level C (8)	Level D (9)
Analysis 25%	Incoherent analysis, lacking data sets and jumping to conclusions.	Analysis partially supported by data, incoherent conclusions.	Sound analysis supported by correct data set. No pulled out of the blue conclusions.	Extensive analysis using several data sets to support conclusion.	Extensive and in-depth analysis supported by literature.
Presentation 10%	One text file, not supported by any images or graphs.	Basic NB, minimal formatting.	Jupyter notebooks with graphs and images, formatted	Advanced use of NB or a formatted IO with graph and images.	Clear layout and properly navigable IO with images and graphs (animated).
Project documentation 20%	No to minimal documentation. No GH landing page. No folder structure.	Documentation is incoherent and not complete. GH landing page non descriptive.	Documentation is incomplete or excessive.	Documentation is present.	All functionality developed documented, very clear documentation.
Code readability 25%	No code or at basic workshop 2 level.	Lack of interpunction, function and variable names not descriptive.	Mixed coding styles. Interpunction taken into consideration.	Coding style according to PEP 8.	Coding style according to PEP. Descriptive variable/function/class names. Code reads as English.
Reproducibility 20%	Code developed for single use case.	The code developed is for a specific or limited data set. Results differ on different runs of the software. Adjustments in code are needed.	The code developed is for a specific or limited data set. Adjustments in code are needed.	Consistent results on a limited data set. Limited part of the code needs to be adjusted.	The code is general enough to run on different data sets. Consistent results.

Do's

- Finish the corresponding labs before starting an assignment
- Think about the objective of the assignment: data exploration vs data analysis
- Think about the method and its limitations
- Make sure that your code runs and produces the expected output
- Use clear and interpretable variable names: '**SP.DYN.TFRT.IN**' vs '**average_fertility_rate**'
- Think about the data quality: are there missing values? How will you deal with them? What are the implications?
- Use headers to structure your notebook
- Use markdown to explain the code and interpret the findings
- Please do not use any language other than **Python and English**

Don'ts

- Don't copy and paste the code from the labs without understanding what the added value is
- Don't print huge matrices or use enormous font. Make sure the notebook can be read as a report
- Don't print empty cells, import useless packages or import packages multiple times
- Don't print figures without labels or titles
- Don't use colours that are difficult for colour-blind people to distinguish
- Don't hardcode, try instead to use functions to make the code easily reusable later
- Don't overwork, focus on what is asked from you in the assignment

Rubric for Final Project

Same as assignments...

Aspect (weight)	Level A- (0-5)	Level A (6)	Level B (7)	Level C (8)	Level D (9)
Analysis 25%	Incoherent analysis, lacking data sets and jumping to conclusions.	Analysis partially supported by data, incoherent conclusions.	Sound analysis supported by correct data set. No pulled out of the blue conclusions.	Extensive analysis using several data sets to support conclusion.	Extensive and in-depth analysis supported by literature.
Presentation 10%	One text file, not supported by any images or graphs.	Basic NB, minimal formatting.	Jupyter notebooks with graphs and images, formatted	Advanced use of NB or a formatted IO with graph and images.	Clear layout and properly navigable IO with images and graphs (animated).
Project documentation 20%	No to minimal documentation. No GH landing page. No folder structure.	Documentation is incoherent and not complete. GH landing page non descriptive.	Documentation is incomplete or excessive.	Documentation is present.	All functionality developed documented, very clear documentation.
Code readability 25%	No code or at basic workshop 2 level.	Lack of interpunction, function and variable names not descriptive.	Mixed coding styles. Interpunction taken into consideration.	Coding style according to PEP 8.	Coding style according to PEP. Descriptive variable/function/class names. Code reads as English.
Reproducibility 20%	Code developed for single use case.	The code developed is for a specific or limited data set. Results differ on different runs of the software. Adjustments in code are needed.	The code developed is for a specific or limited data set. Adjustments in code are needed.	Consistent results on a limited data set. Limited part of the code needs to be adjusted.	The code is general enough to run on different data sets. Consistent results.

Support!

Centre for Urban Science & Policy 

EPA 122A Spatial Data Science

-  Introduction
-  Lectures
-  Labs
-  Assessment
-  Software
-  **Helping Material**
-  FAQ



Helping Material

Our teaching support team has prepared and updated two wonderful resources for the students. There is one resource on programming support. You can find links to library documentation, data science practices, analysis and much more. Another resources is on debugging. Everytime your program fails to do what you expected of it, go to this resource first.

Programming Help Sheet

This [document](#) is your go-to guide when you encounter challenges in your coding journey. It emphasizes a structured approach, beginning with checking official documentation to ensure correct usage, followed by exploring general Python tutorials for broader understanding.

Programming Help Sheet
Prepared by Ludovica Bindi and Dorukhan Yeşilli

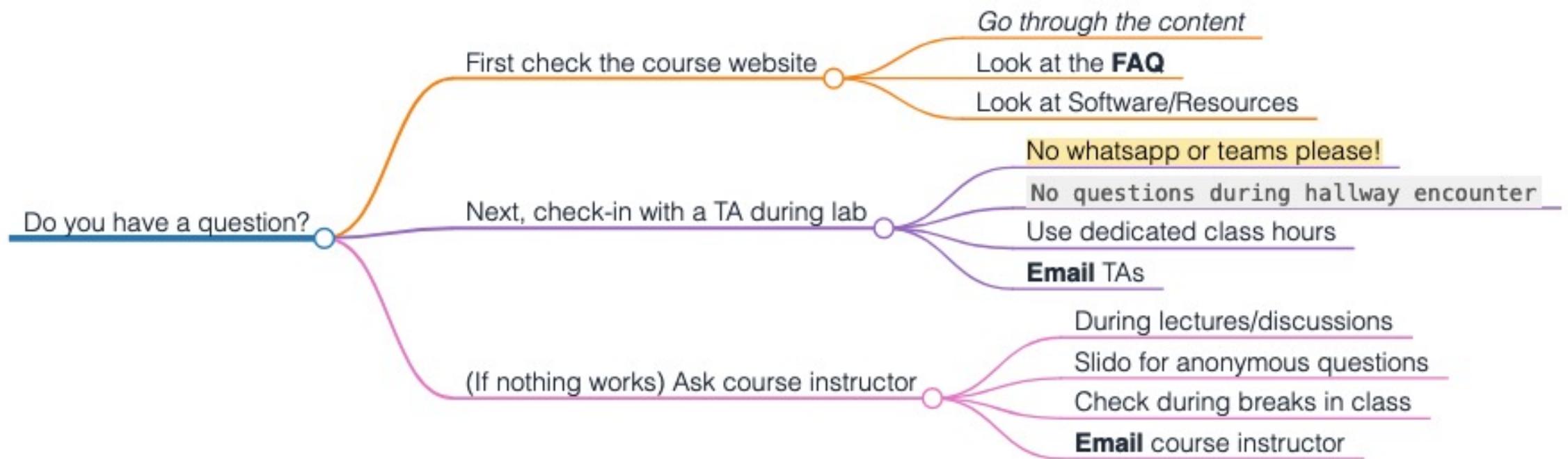
When facing a problem with your code, a library function, etc., to increase the chances of learning from your mistakes, the first thing you should do is to check the official documentation to see if you are doing things correctly (e.g., properly using the library function by passing the right kinds of arguments), then check more structured websites that contain general Python tutorials and read their explanations; afterwards, if you are still struggling check websites like StackOverflow (but beware: sometimes you can find there very specific solutions to the specific problems there, solutions that are too complicated and that you don't need, and the language on the website can become too technical). On this help sheet, you can find links to official documentation, helpful websites, and cheat sheets.

More Support!!!

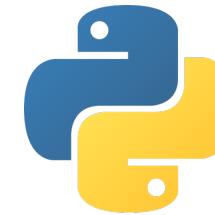
- Ask questions
- Help others as much as you can (the best way to learn is to share perspectives)
- Search heavily on browser + Stack Overflow (learn to troubleshoot)
- Bring questions, comments, feedback, (informed) rants to class
- Collaborate with each other



General Support Structure



Installing Python



EPA 122A Spatial Data Science

Introduction

Lectures

Labs

Assessment

Software

Standard Installation

Minimalist Installation

Comprehensive Installation

Virtual Environments

Helping Material

FAQ

Software

This course is best followed if you can reproduce the examples and tutorials provided with it. To do so, you will need to install in your machine a series of software packages. These are all open-source and available for free to download.

There are three main pathways to install required Python libraries on your machine.

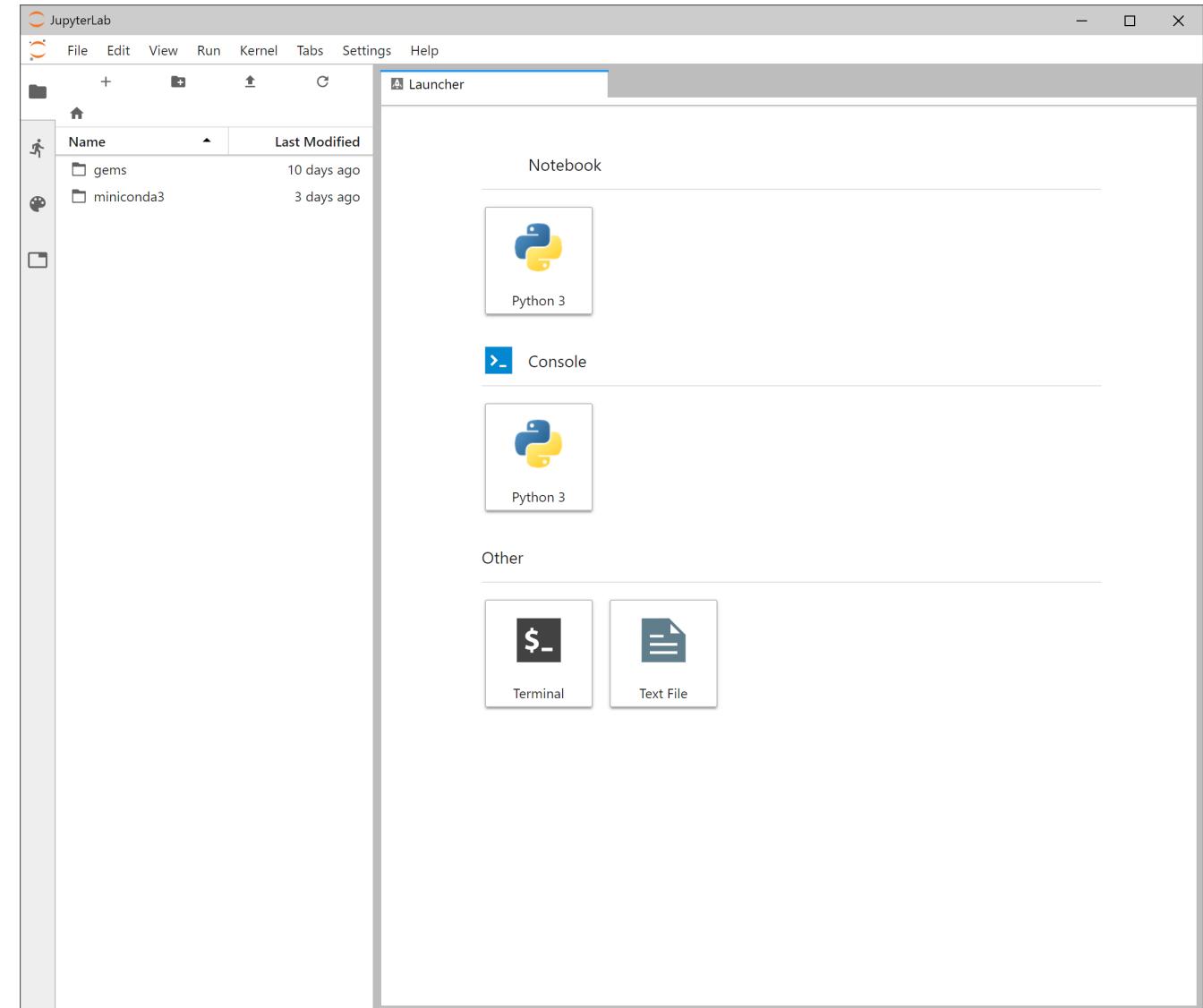
- A [\(1\) standard](#) one is installing the software on your operating system without using the command line interface.
- A [\(2\) minimalist](#) one will provide basic Python resources and the ability to expand them.
- A [\(3\) comprehensive](#) one will install not only a Python stack but also several useful libraries (including some from the programming language, R).

If you want to learn to explore Python and its capabilities, while going beyond this course, I recommend option 2. If all else fails, option 3 is the last resort for this course. It is guaranteed to work and very powerful, so you will not be limited in any way. But it does not allow you to install new libraries, which means you are limited by what it offers.

Hint

The difference in these options can be explained through the illustration of a living place. If you own a house, you might be able to expand it, paint the walls, add new furniture, even keep a dog. This is akin to the **minimalist** approach which gives you everything you need and the freedom to build upon it. Instead, if you rent a house, in most cases you will not be allowed to make any changes. A **comprehensive** approach gives you everything too, but no freedom to experiment with new python libraries. The **standard** option is like visiting a hotel where others service you for a bit without you having to do the heavy lifting. You choose what works for you!

Installing Python





Create Groups

Turn to your neighbors and chat for three minutes about,

- Name
- What did you study before coming here?
- What do you expect from this course?
- How do you see this course helping you with your future ambitions?
- Can we learn something from each other?

Community

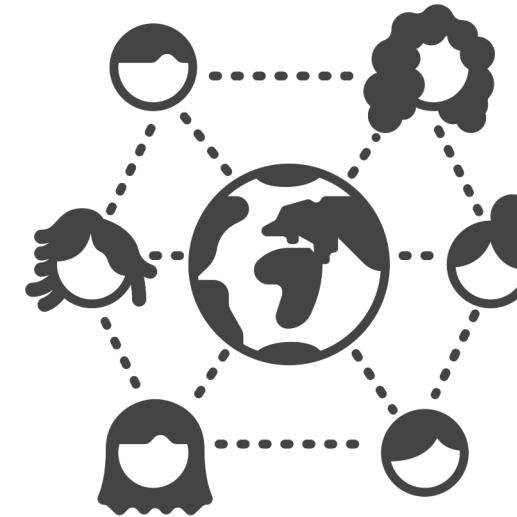
- We are from different parts of the world
- We have different educational backgrounds
- We have different experiences in life



Adapted from the work of Sean Perez

Our students bring DIVERSITY

Socioeconomic Backgrounds
Education, Culture, Wealth, Identity



Lived Experiences
Gender, Ethnicity, Family, Country

Theories, Concepts & Challenges
Transitions, Development, Poverty, Politics

Through **recognizing and analysing the cultures** in which we are **positioned**, and that therefore cannot help but mould our worldviews, we take steps to become more aware and even more objective. We **come to know the world more fully by knowing how we know the world**.

– Strong Objectivity, Sandra Harding

Quality of Education

What would you change if you were the course manager?

What are strong points you would definitely keep?

Share it...

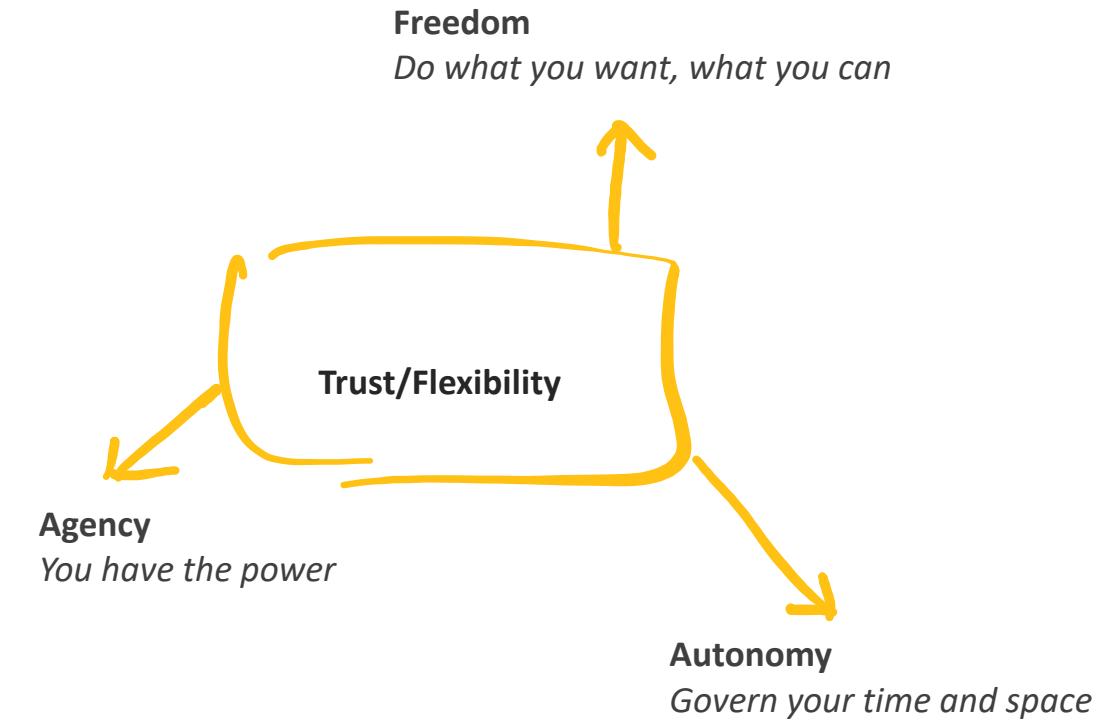
- With us, the lecturers
- With Management: Merel Hijzelendoorn.
- With your fellow students who join the student council
- Fill in the questionnaire at the end of this course

Flexibility in learning



Flexibility

in everything



You do what fits your lives!

Flexibility

in everything

- You do not owe us “productivity” or “efficiency”. We just want your participation, so we can all learn something new.
- You do not owe us any information concerning your personal situation or mental or physical health condition.
- If you want to, you are welcome to talk to us about anything you are going through. **Just drop us an email**, and we will figure it out from there.
- This course is just one small part of life. If you have to work around it to figure life out, go ahead and do it. We trust you will reach out to us if you need support, and we will be here to offer it.
- Exams only focus on your ability to regurgitate knowledge in a 3-hour window. There is no exam in this course. The final project will provide you an opportunity to learn from each other and create something awesome. *But we do have a 4-hour Hackathon...*

Deadline Philosophy

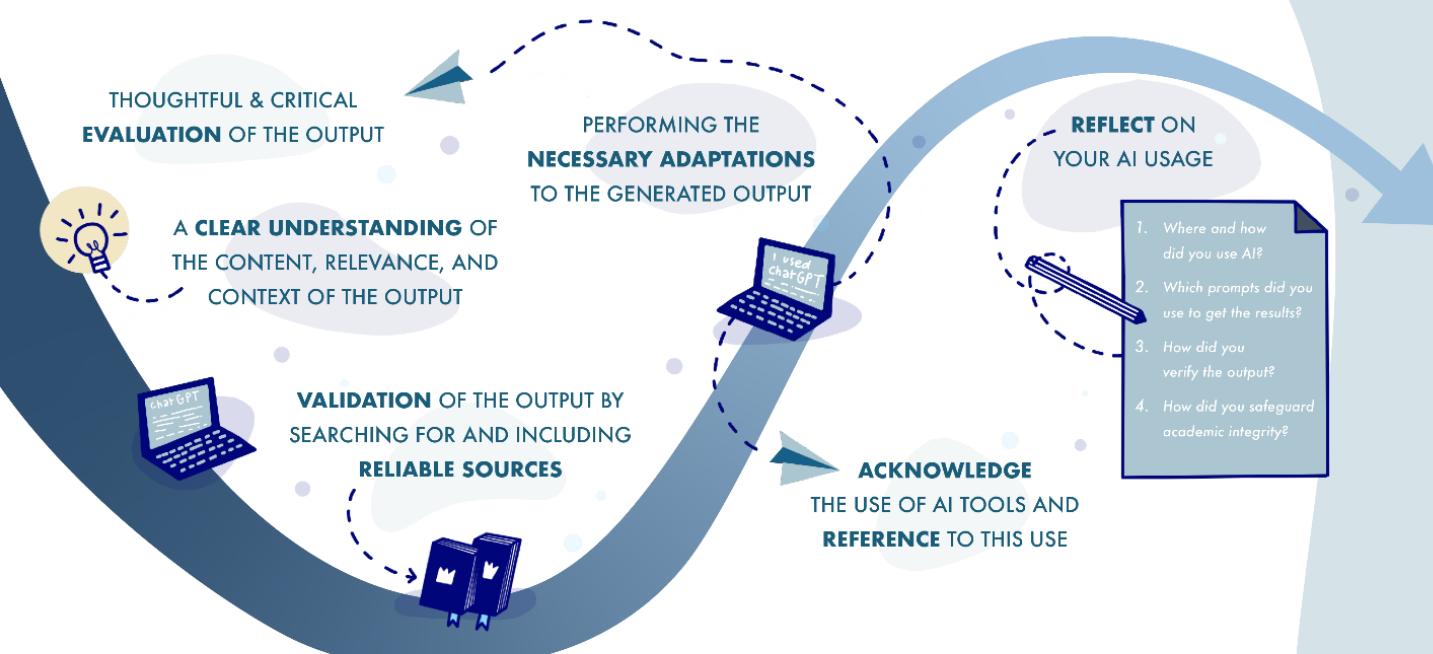
Deadlines are there for a reason!

But if you need more time because... well life, you need to send us an **email** explaining **why** you need more time and send us that email in good time.

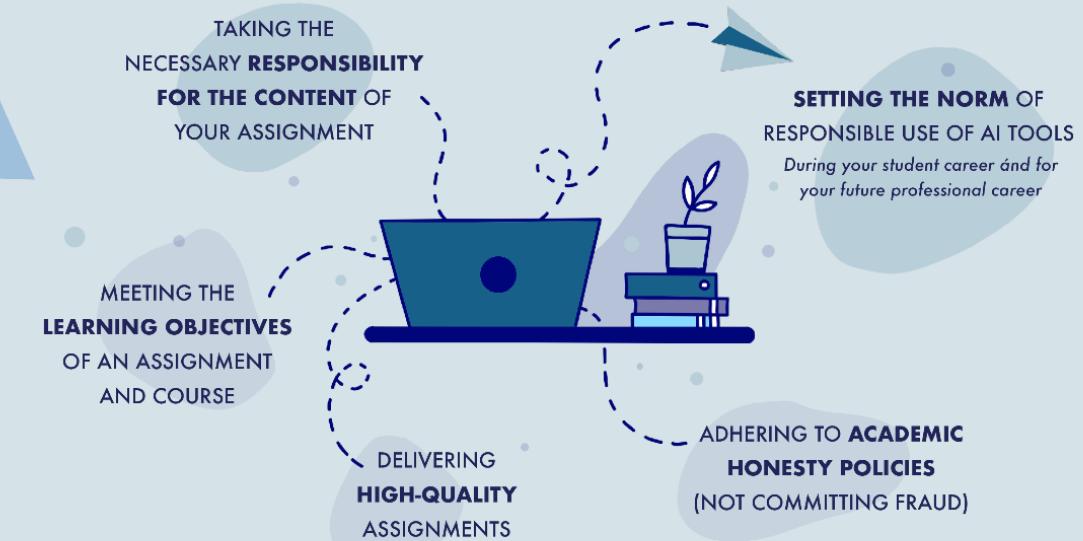


Many colossal screwups exist in the world in leading political positions.
We are not here to police you.
We are not here to filter those who can hack their way into differentiating AI vs non-AI.
You are all here because you pay for this degree. We are here so you can learn along the way.

RULES WHEN USING AI TOOLS



THE NEED FOR COMPLIANCE TO THESE RULES



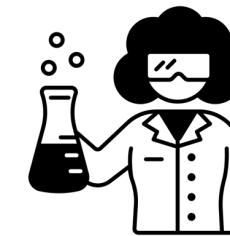
Break



CHILL



WALK



COFFEE OR TEA



MAKE FRIENDS

Spatial Data Science

Introduction-II

(YMS31303)

Lecture 1



Adapted from the work of Sean Perez

Just before the break

- Introduction to the Course
- Tools - Python and Conda

Now..

- The Data Revolution
- (Geo-)Data Science
- Why Data Science?
- What is Data Science?

The data revolution

Exciting times to be a:

- Data Scientist
- Urban Planner
- Policymaker

The world is producing a lot of “**data**”...

Massive Data Revolution

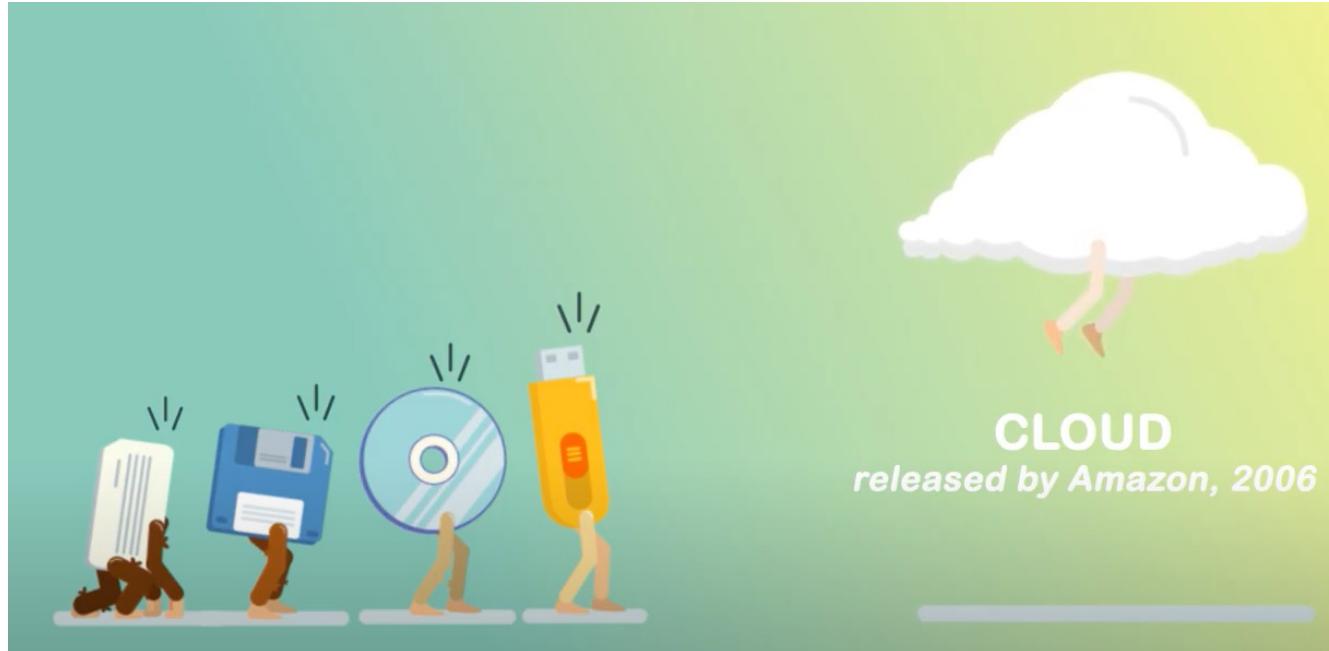
Quantification of phenomena through the systematic recording of data, “taking all aspects of life and turning them into data” ([Cukier & Mayer-Schoenberger](#))

Examples: credit transactions, public transit, tweets, facebook likes, spotify songs, etc.

Implications

- **Window** into human behaviour (this course)
- Opportunities for optimization of systems (Industrial IoT, planning systems...)
- Issues with **representation** and **privacy**
- ...

Why now?



Statistics

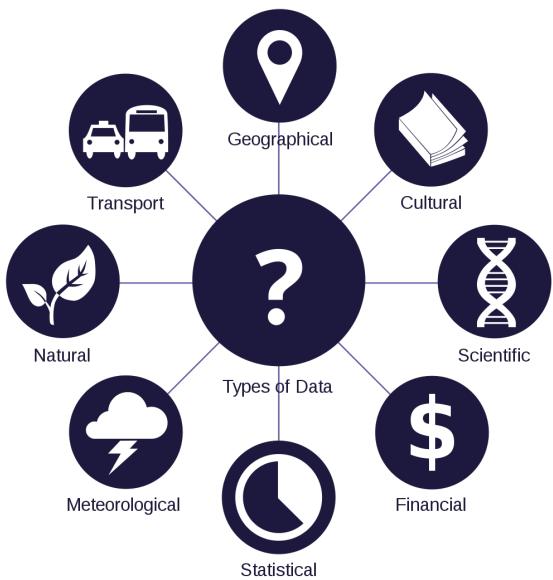
Machine
learning

What's
next?
o

- o Massive Data generation
- o Computing power
- o R + Python
- o Visualisation

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

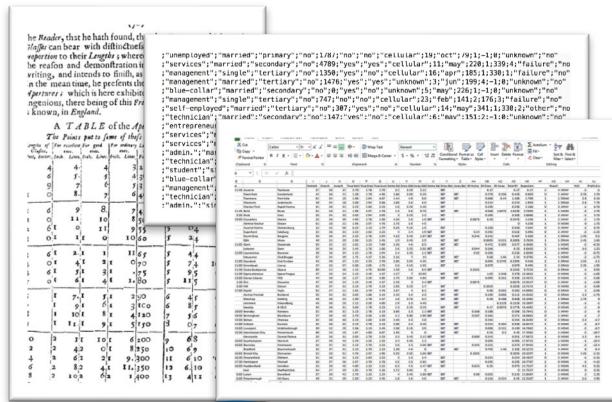
At this point many people have said it..



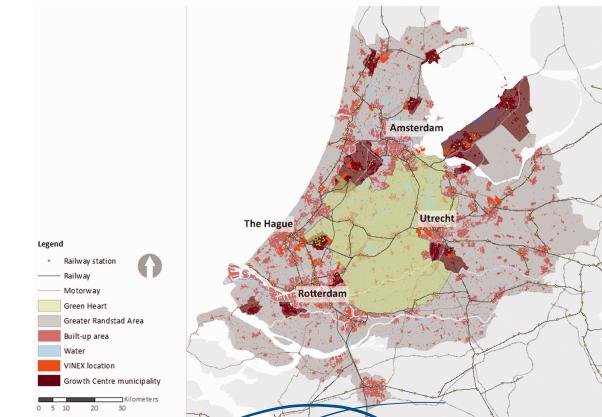
Examples: credit transactions, public transit usage, tweets, census, mobility and migration, etc.

Formats: CSV, Excel, JSON, Shapefiles

Now, data alone is not very valuable



Data



Action

Methods, tools and techniques to turn data into
actionable knowledge (hence this lesson!)

Class Quiz

Can you think of a real-world context where data and statistics are being used to make a difference? And how?

- Turn to your neighbour and discuss for two minutes
- Then I may ask you to summarise your discussion



Data Science

Statistics + ...

- **Computational** tools → Programming (hence this course's labs and homework!)
- **Communication** skills → “Story telling ” (hence this course's assignments)
- **Domain** expertise → Theories about why the data are the way they are (hence the rest of your degree)

Some examples...

Emmy-winning US TV Shows



Police Detective TV Dramas



Critically Acclaimed Witty TV Shows



Free Online Dating | OkCupid - Mozilla Firefox (Private Browsing)

Free Online D... https://www.okcupid.com

Have an account? Sign in

okcupid

Join the best free dating site on Earth.

I am a

Woman

Continue



Signing up takes two minutes and is totally free.



Our matching algorithm helps you find the right people.



iOS or Android?
You can take us to go.

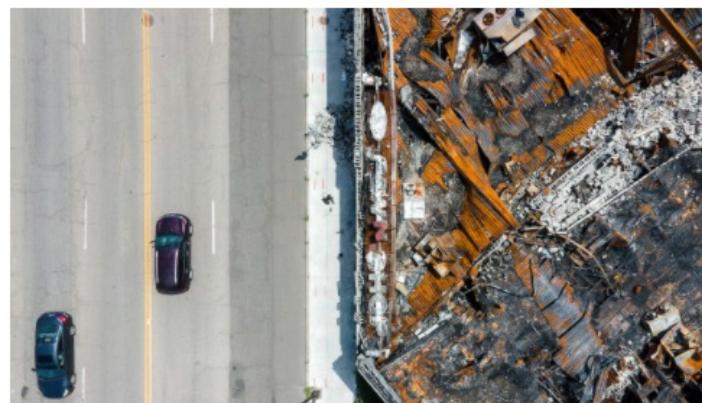
The (Geo-)Data Revolution

The Global Picture: Urban Inequalities

Rising Seas



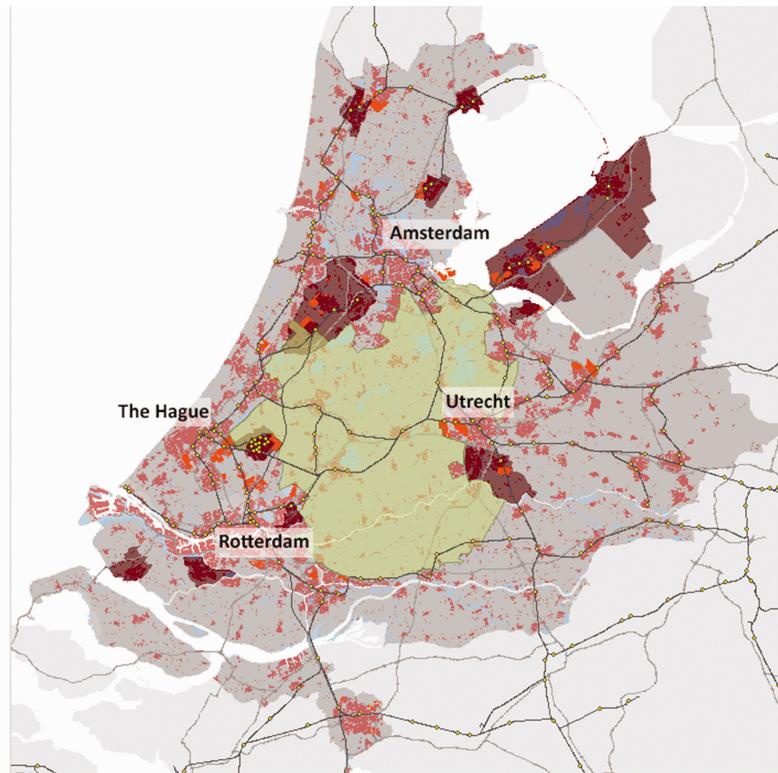
Economic
Inefficiencies



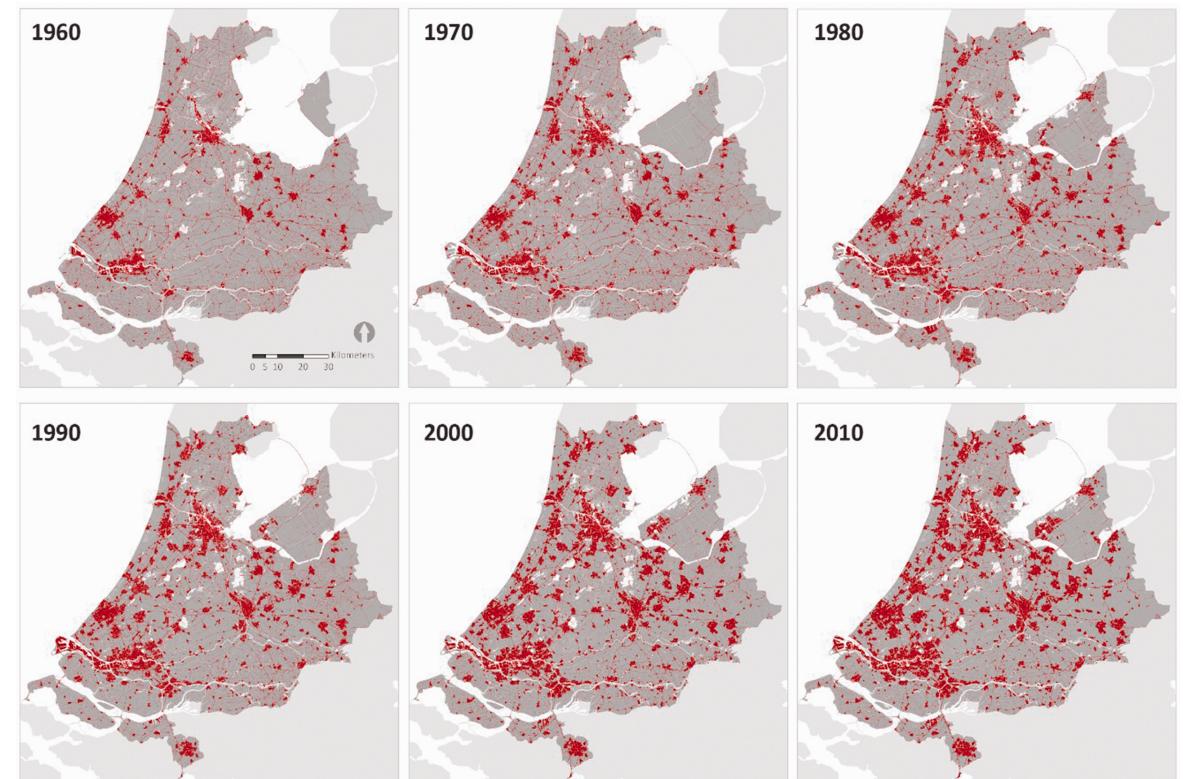
Growing Energy
Demand



The Local Picture: Randstad



Urbanisation of the Greater Randstad Area 1960-2010



Space is important!

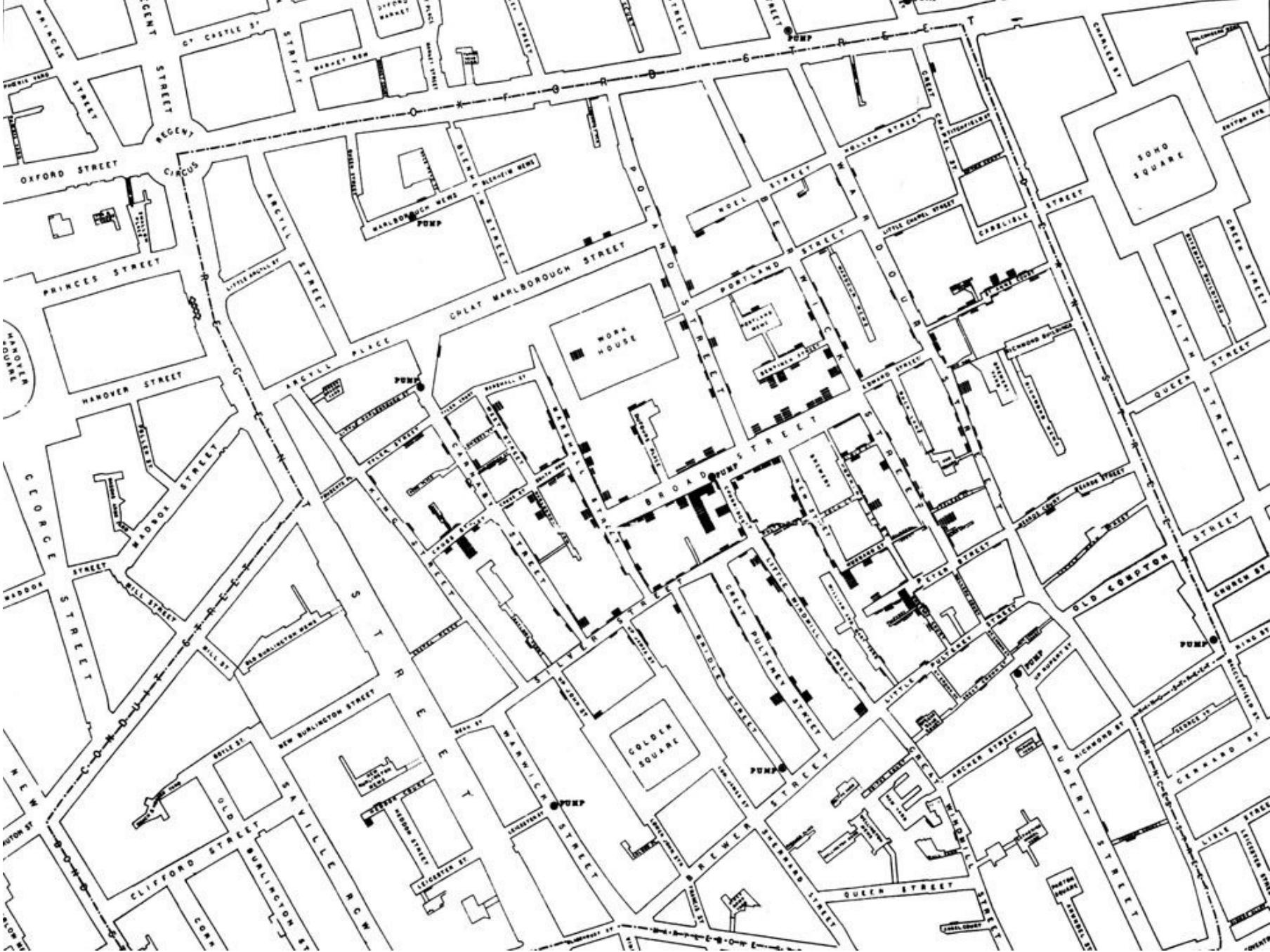
- The space around us
- Thus, the geography of our location
- Geolocated data

The world is producing a lot of geo-located “**data**”...

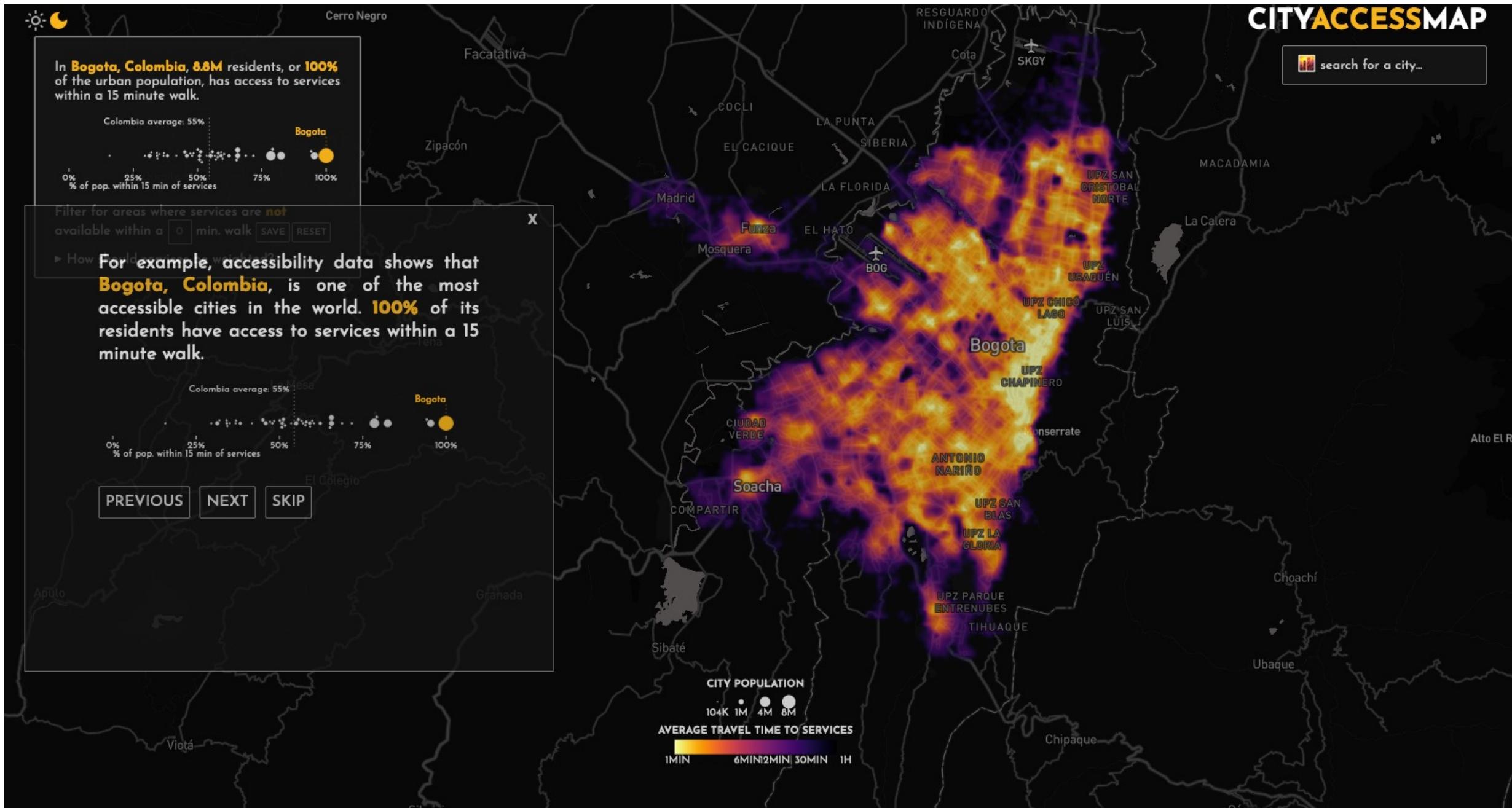
(Geo-)Data Science

(Geo-)Data Science

- A (very) large portion of all these new data are inherently **geographic** or can be traced back to some location over space.
- Spatial is special.
- Some of the methods require an explicitly spatial treatment -> (Geo-)Data Science
- Some examples...



Map of the book "On the Mode of Communication of Cholera" by John Snow, originally published in 1854 by C.F. Cheffins, Lith., Southampton Buildings, London, England.



To do before class [Takes about 1 hour of prep at home]

As a way to whet your appetite about the content of the first class, I recommend you:

- Listen to [this interview](#) with Hilary Mason, Max Shron, and Alex Pentland about the power of data.
- Watch [this video](#) by Mike Flowers, Chief Analytics Officer, at the City of New York about how data is used to influence policy decisions.
- Read [What New Yorkers are complaining about](#) and reflect on [if the cost of running such data systems worth the price of knowing?](#)

ARCHIVE

Is the Cost of 311 Systems Worth the Price of Knowing?

311 systems have revolutionized the way cities gather information, allowing them to tackle small problems before they get too big. But running them can be extremely costly.

February 24, 2014 • Tod Newcombe



Minneapolis, Minn. FlickrCC/Photo Phiend

Why Data Science?

History

Long time ago (thousands of years) science was only empirical, and people counted stars



© Trivik Verma. All rights reserved.

History (cont)

Long time ago (thousands of years) science was only empirical, and people counted stars or crops

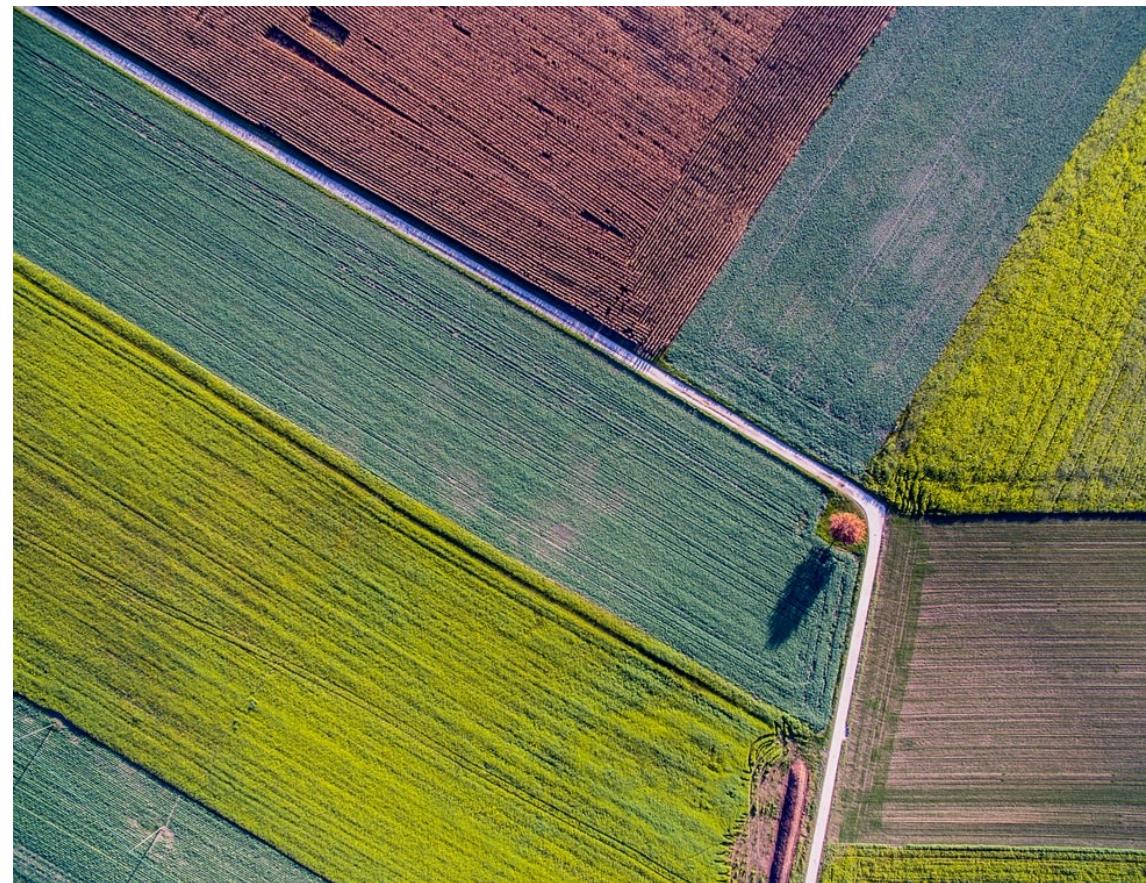


Photo by [jean wimmerlin](#) on [Unsplash](#)

History (cont)

Long time ago (thousands of years) science was only empirical, and people counted stars or crops and used the data to create machines to describe the phenomena



Photo by [Frank Chou](#) on [Unsplash](#)

History (cont)

Few hundred years: theoretical approaches, try to derive equations to describe general phenomena.

Maxwell's Equations	$\nabla \cdot \mathbf{E} = 0$ $\nabla \times \mathbf{E} = -\frac{1}{c} \frac{\partial \mathbf{H}}{\partial t}$	$\nabla \cdot \mathbf{H} = 0$ $\nabla \times \mathbf{H} = \frac{1}{c} \frac{\partial \mathbf{E}}{\partial t}$	J.C. Maxwell, 1865
Second Law of Thermodynamics	$dS \geq 0$		L. Boltzmann, 1874
Relativity	$E = mc^2$		Einstein, 1905
Schrodinger's Equation	$i\hbar \frac{\partial}{\partial t} \Psi = H\Psi$		E. Schrodinger, 1927
Information Theory	$H = - \sum p(x) \log p(x)$		C. Shannon, 1949
Chaos Theory	$x_{t+1} = kx_t(1 - x_t)$		Robert May, 1975
Black-Scholes Equation	$\frac{1}{2} \sigma^2 S^2 \frac{\partial^2 V}{\partial S^2} + rS \frac{\partial V}{\partial S} + \frac{\partial V}{\partial t} - rV = 0$		F. Black, M. Scholes, 1990

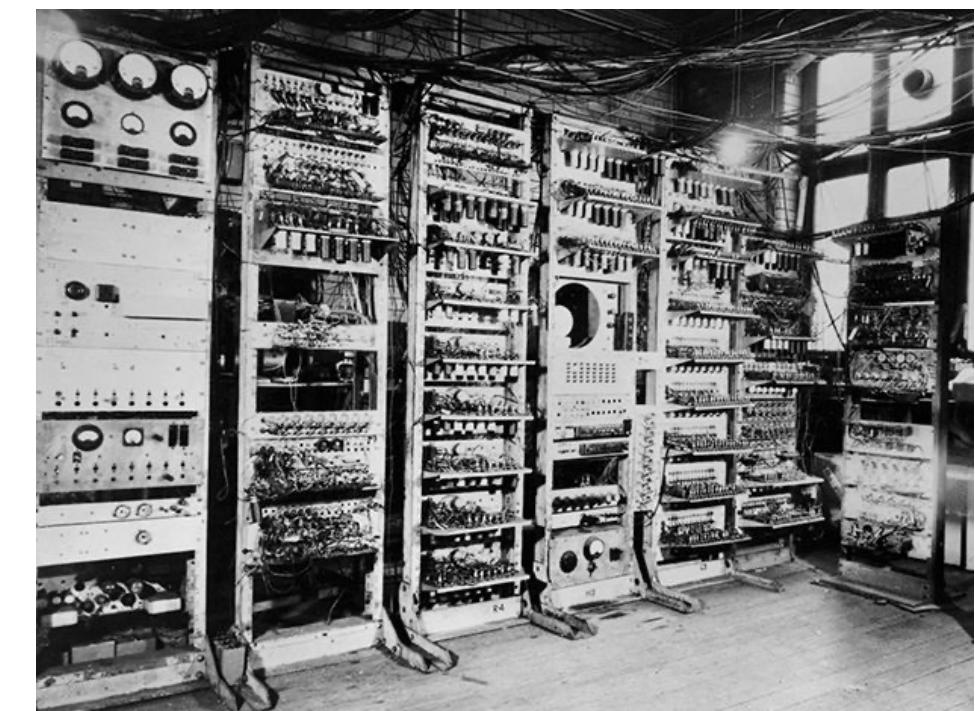
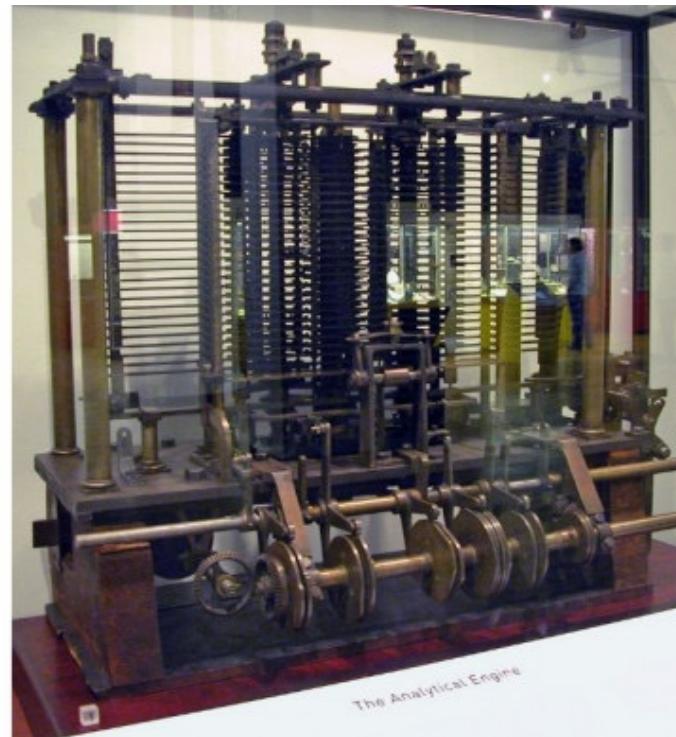
Stewart, I. (2012). *In pursuit of the unknown: 17 equations that changed the world*. Basic Books.

History (cont)

About a hundred years ago: computational approaches



Scanned from *The Calculating Passion of Ada Byron* by Joan Baum.
Analytical Machine [Wikimedia Commons](#)



SSPL/Getty Images The Manchester Mark I at Manchester University's Computer Machine Laboratory.

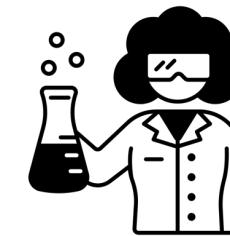
Break



CHILL



WALK



COFFEE OR TEA



MAKE FRIENDS

Spatial Data Science

Introduction-III

(YMS31303)

Lecture 1



Adapted from the work of Sean Perez



What is Data Science?

what my friends think I do



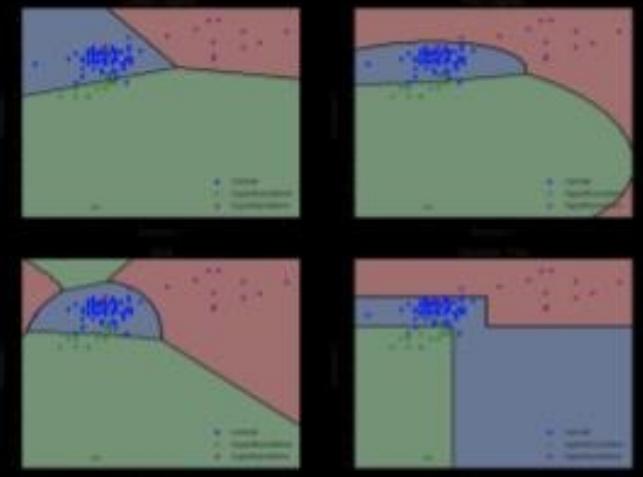
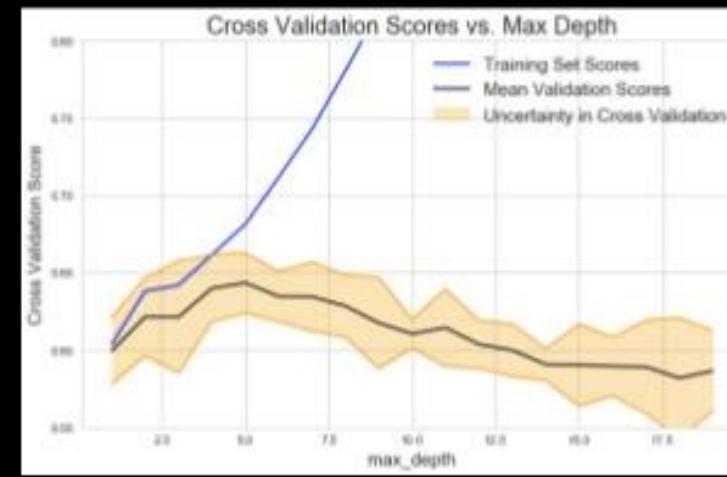
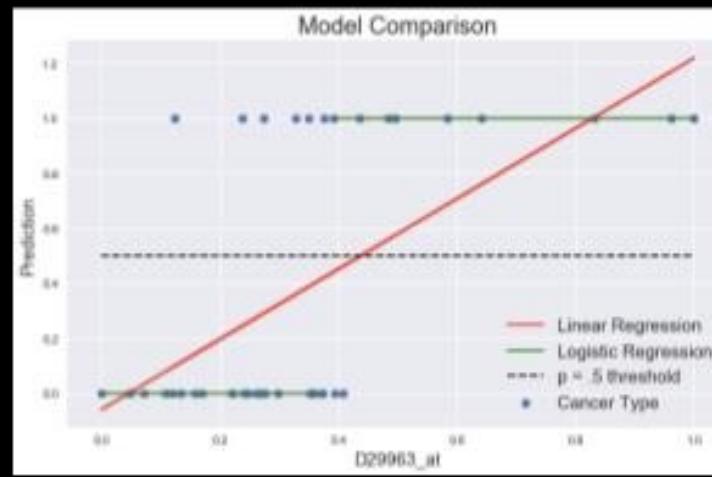
what my family thinks I do



what society thinks I do



what I actually (will) do in Data Science 1

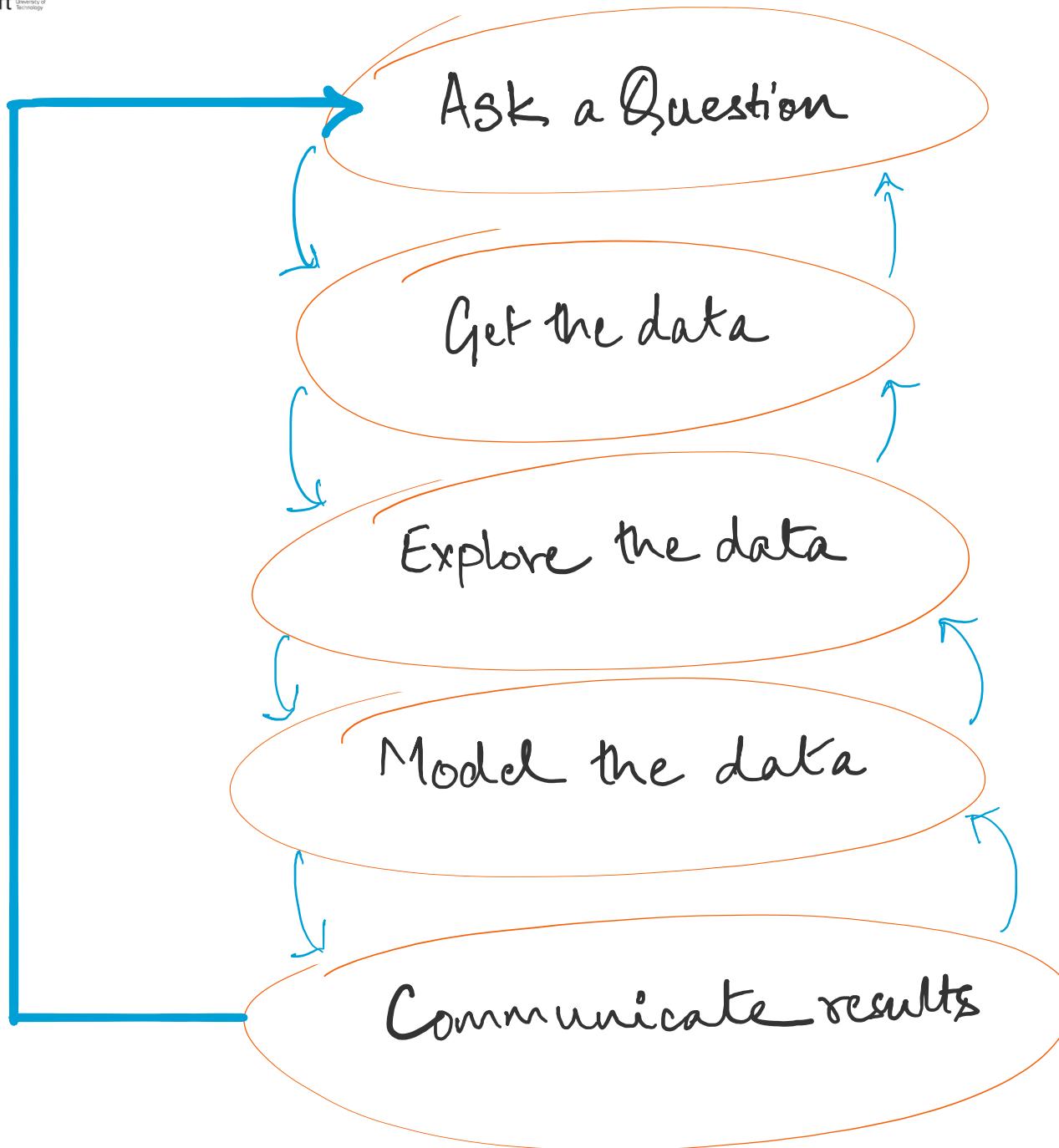


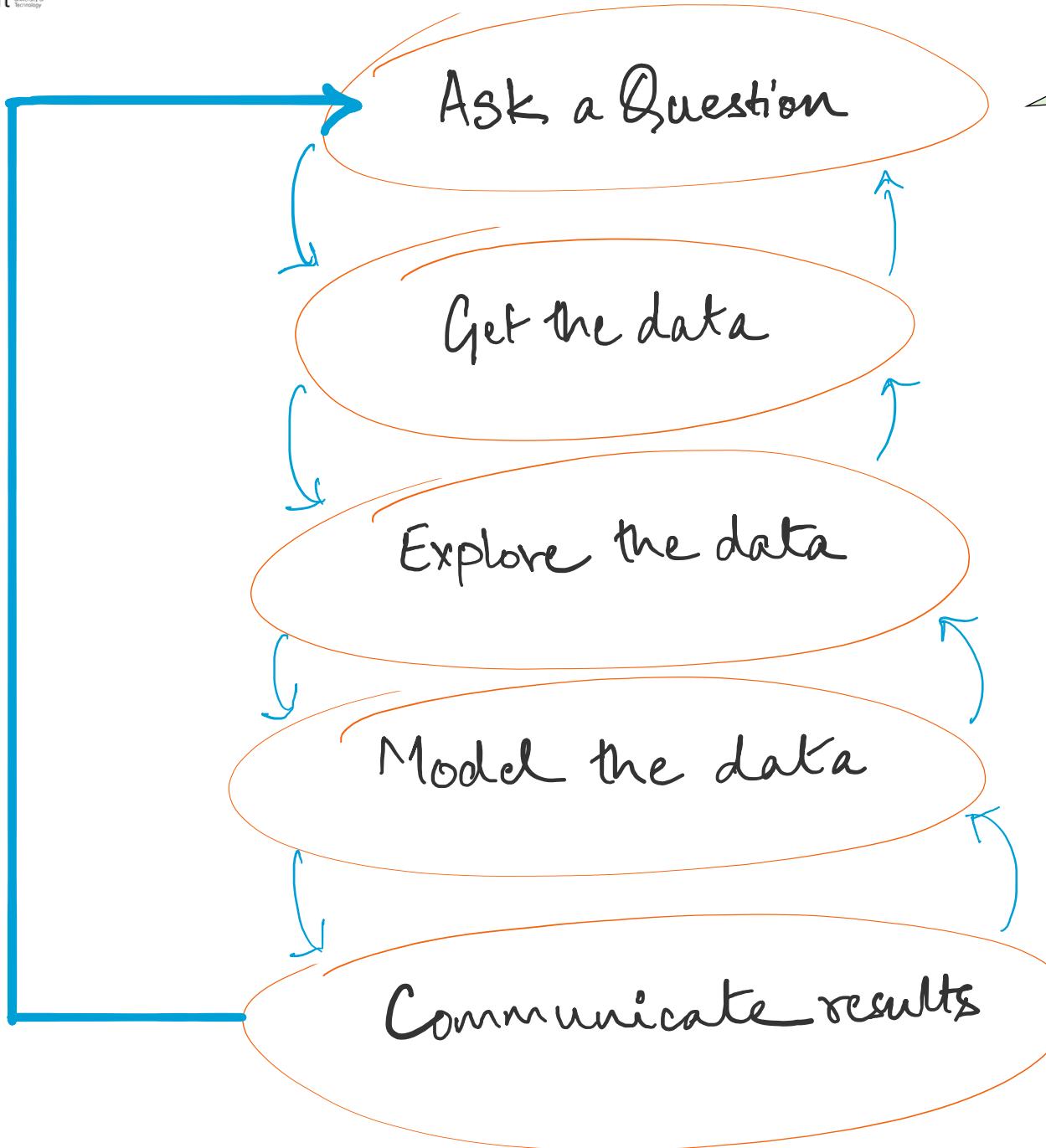
The Data Science Process

The Data Science Process is like the scientific process - one of observation, model building, analysis and conclusion:

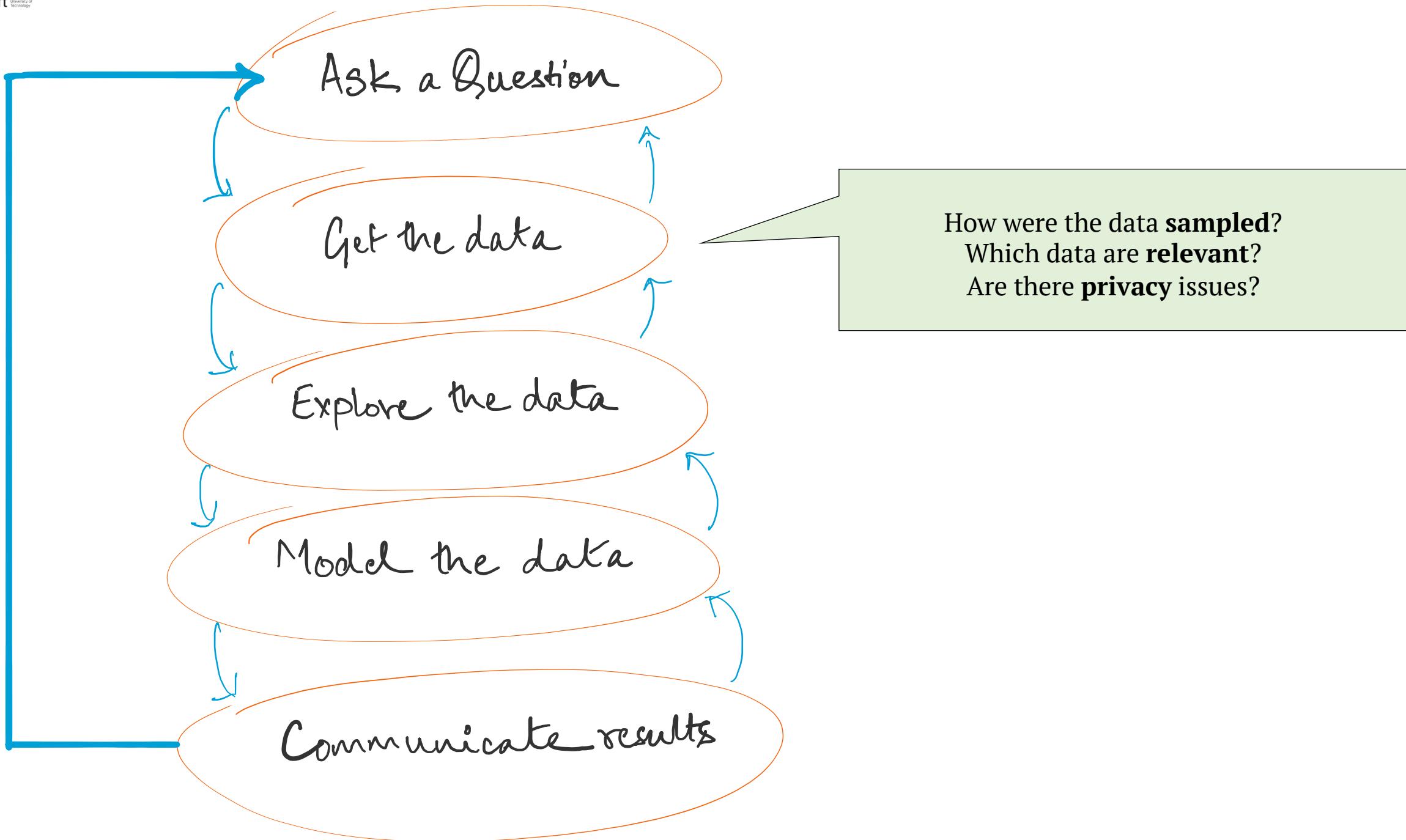
- Ask questions
- Data Collection
- Data Exploration
- Data Modeling
- Data Analysis
- Visualization and Presentation of Results

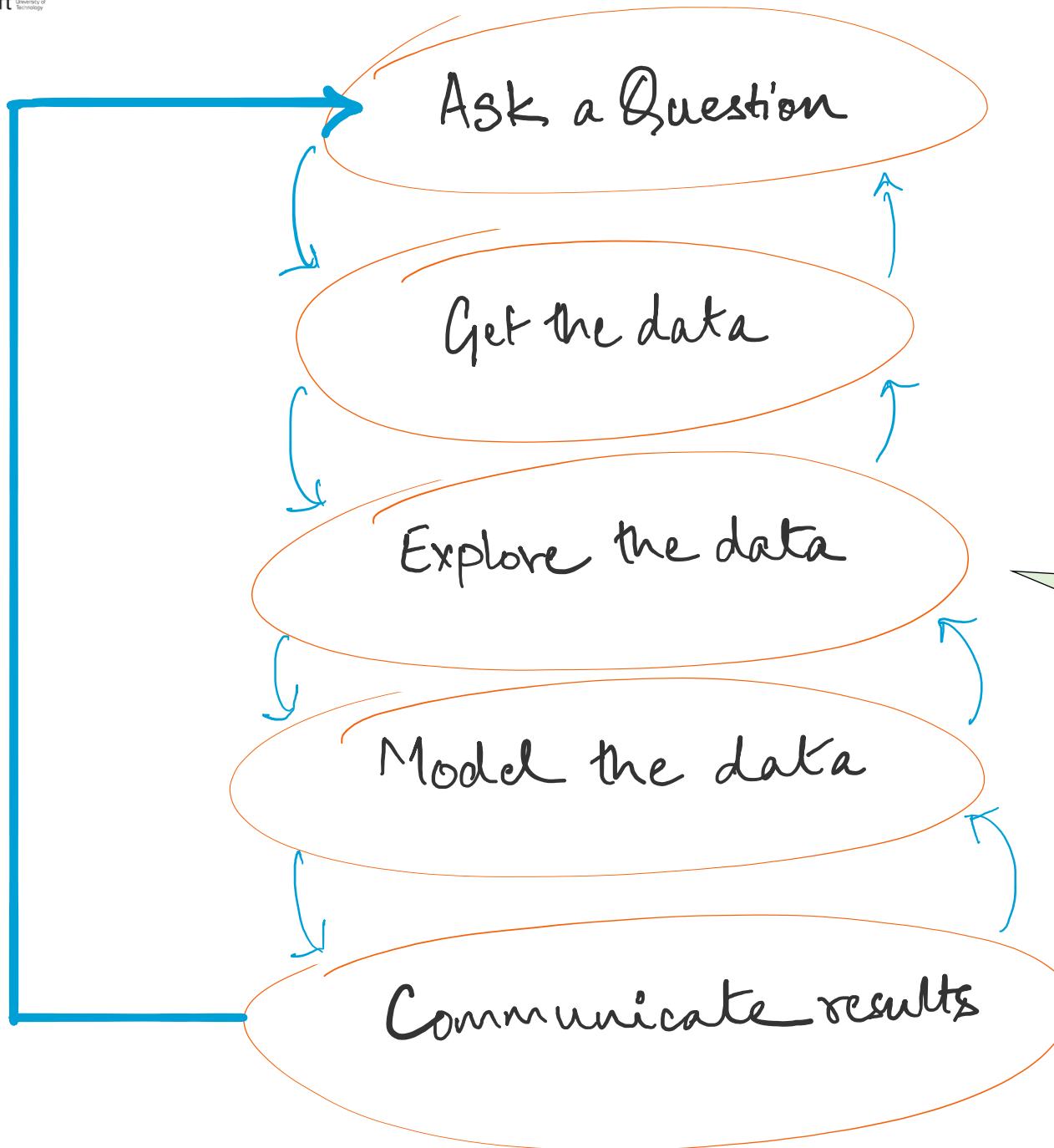
Note: This process is by no means linear!



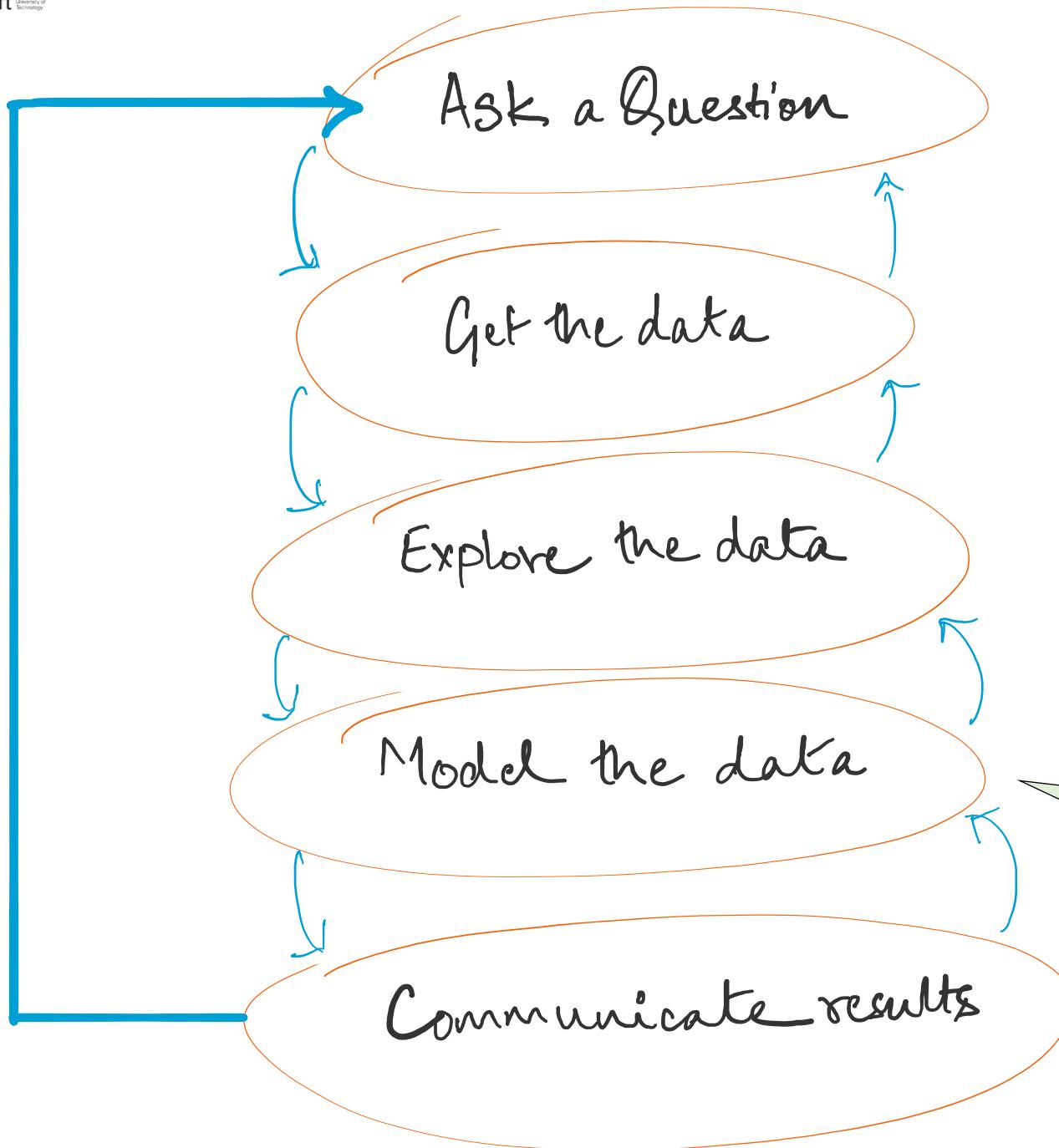


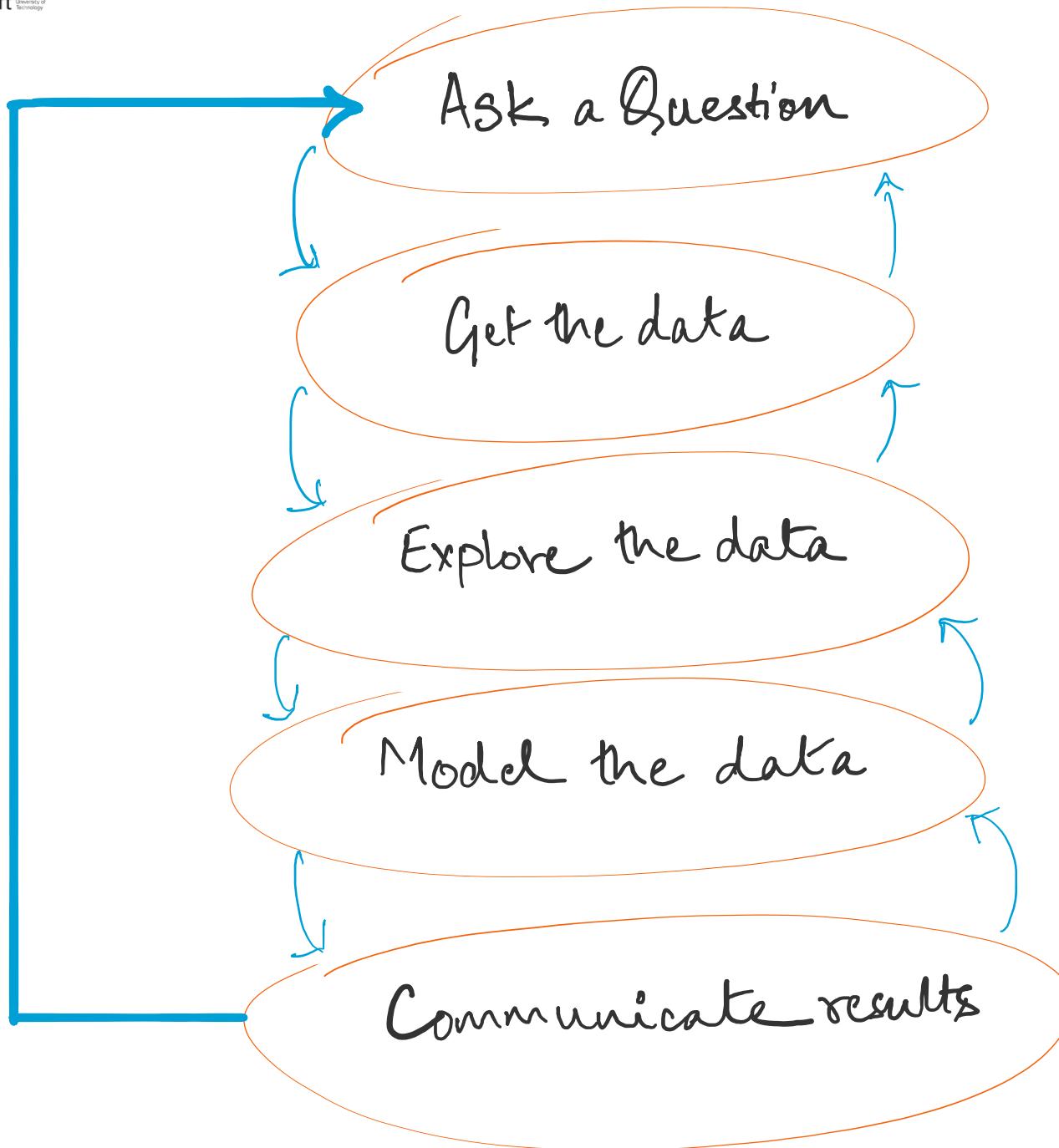
What is the scientific **goal**?
What would you do if you had all the **data**?
What do you want to **predict** or **estimate**?



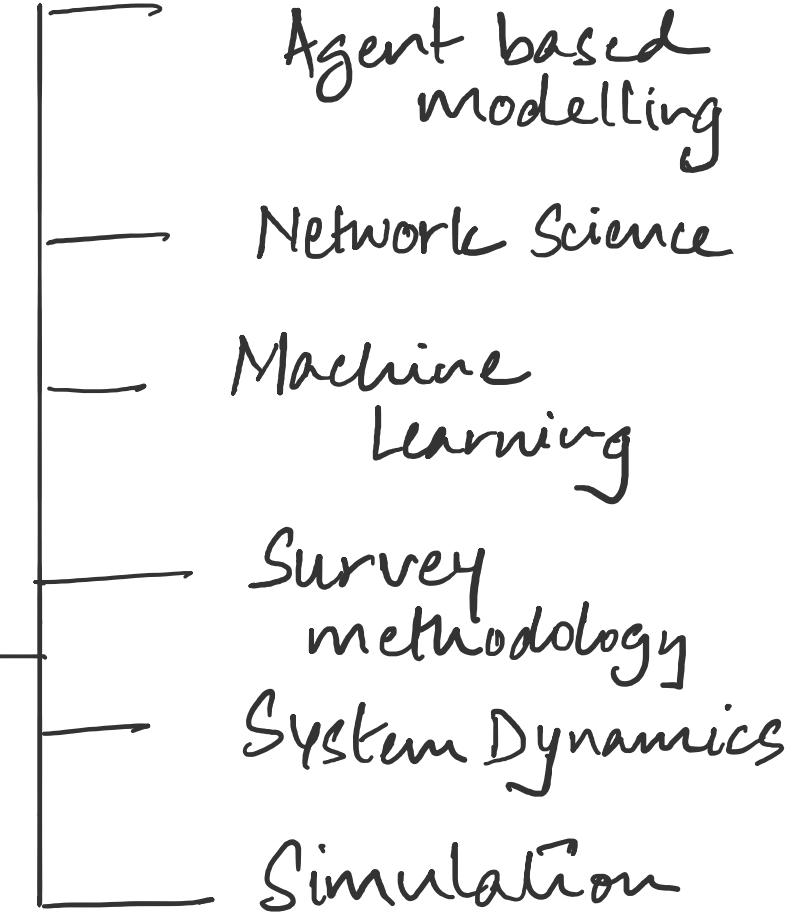


Plot the data.
Are there **anomalies**?
Are there **patterns**?





What did we **learn**?
Do the results make **sense**?
Can we tell a **story**?



Critical Data Science

- A student handbook led by Laura van Geene

Data Science Process	Inclusion <i>Who is (not) included in the data?</i>	Inequality <i>What role does inequality play in data science methods?</i>	Participation <i>Who is (not) involved in the data science process?</i>	Power <i>How does the data reflect existing power dynamics?</i>	Positionality <i>What is your own positionality with the research?</i>
Focus of Analysis <i>Theories, processes & stakeholders that drive the analysis</i>	<p>Investigation of exclusive practices of past and present relating to the research focus, and how these affect the diversity of the people represented in the data (Boyd, 2021a; Lee et al., 2022).</p>	<p>What tools can be reliably used to explore the research topic? Research the limitations of the methods, particularly their influence to structural inequalities (Boyd, 2021a).</p>	<p>Use a participatory modelling approach and include stakeholders that may not have otherwise been involved in the design of the research process and discuss how to include topics/perspectives that are not commonly researched (Lee et al., 2022). Discuss the possibility of multiple framings of the research topic (Delbos, 2023).</p>	<p>Investigation of where and with whom power was distributed in the situations referenced with the data (Lee et al., 2022). Also investigate potential histories of injustices and oppression of the sampling population (Harrington et al., 2021).</p>	<p>Critically reflect on your own position to the research (Boyd, 2021a):</p> <ol style="list-style-type: none"> 1. Why are <u>you</u> doing research about this specific topic? Why are you specifically involved in this research? What makes you suitable for this research? 2. What is the story that you are trying to tell with this research? Consider biases: do you already have ideas about how this story should go? 3. Is there potential that you cause harm or erasure with your research about this topic?

Before you start Lab 0..

Why do we use Functional programming

- **Organization** -- As programs grow in complexity, having all the code live inside the main() function becomes increasingly complicated. A function is almost like a mini-program that we can write separately from the main program, without having to think about the rest of the program while we write it. This allows us to reduce a complicated program into smaller, more manageable chunks, which reduces the overall complexity of our program.
- **Reusability** -- Once a function is written, it can be called multiple times from within the program. This avoids duplicated code (“Don’t Repeat Yourself”) and minimizes the probability of copy/paste errors. Functions can also be shared with other programs, reducing the amount of code that must be written from scratch (and retested) each time.
- **Testing** -- Because functions reduce code redundancy, there’s less code to test in the first place. Also, because functions are self-contained, once we’ve tested a function to ensure it works, we don’t need to test it again unless we change it. This reduces the amount of code we must test at one time, making it much easier to find bugs (or avoid them in the first place).
- **Extensibility** -- When we need to extend our program to handle a case it didn’t handle before; functions allow us to make the change in one place and have that change take effect every time the function is called.
- **Abstraction** -- In order to use a function, you only need to know its name, inputs, outputs, and where it lives. You don’t need to know how it works, or what other code it’s dependent upon to use it. This lowers the amount of knowledge required to use other people’s code (including everything in the standard library).

For next class..



Finish Labs to practice
programming



Complete Homework for
more practice



Check Assignment
contents and due date



See “To do before class”
for next lecture (~ 1 hour
of self-study)