# Lecture 1: Introduction

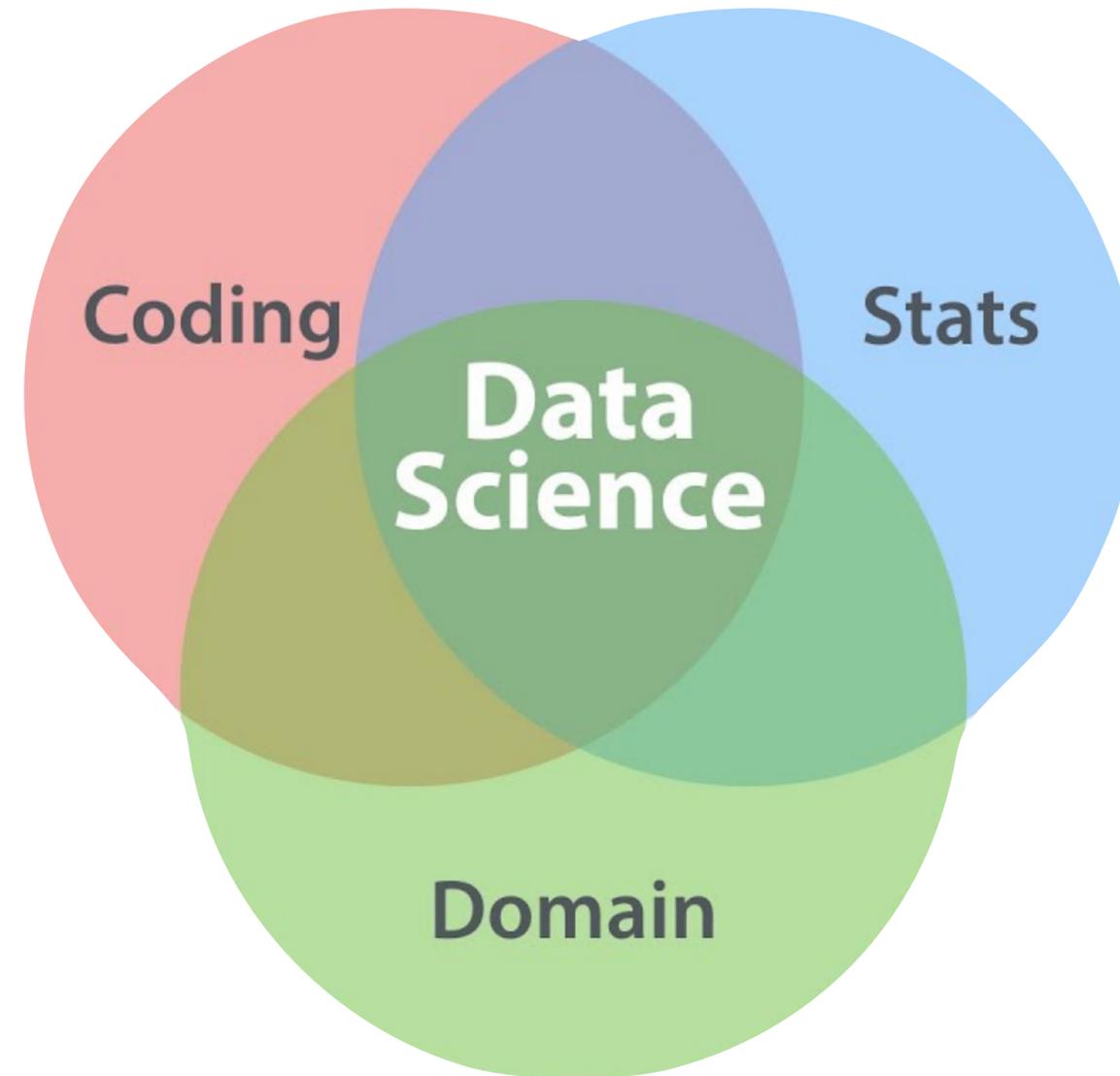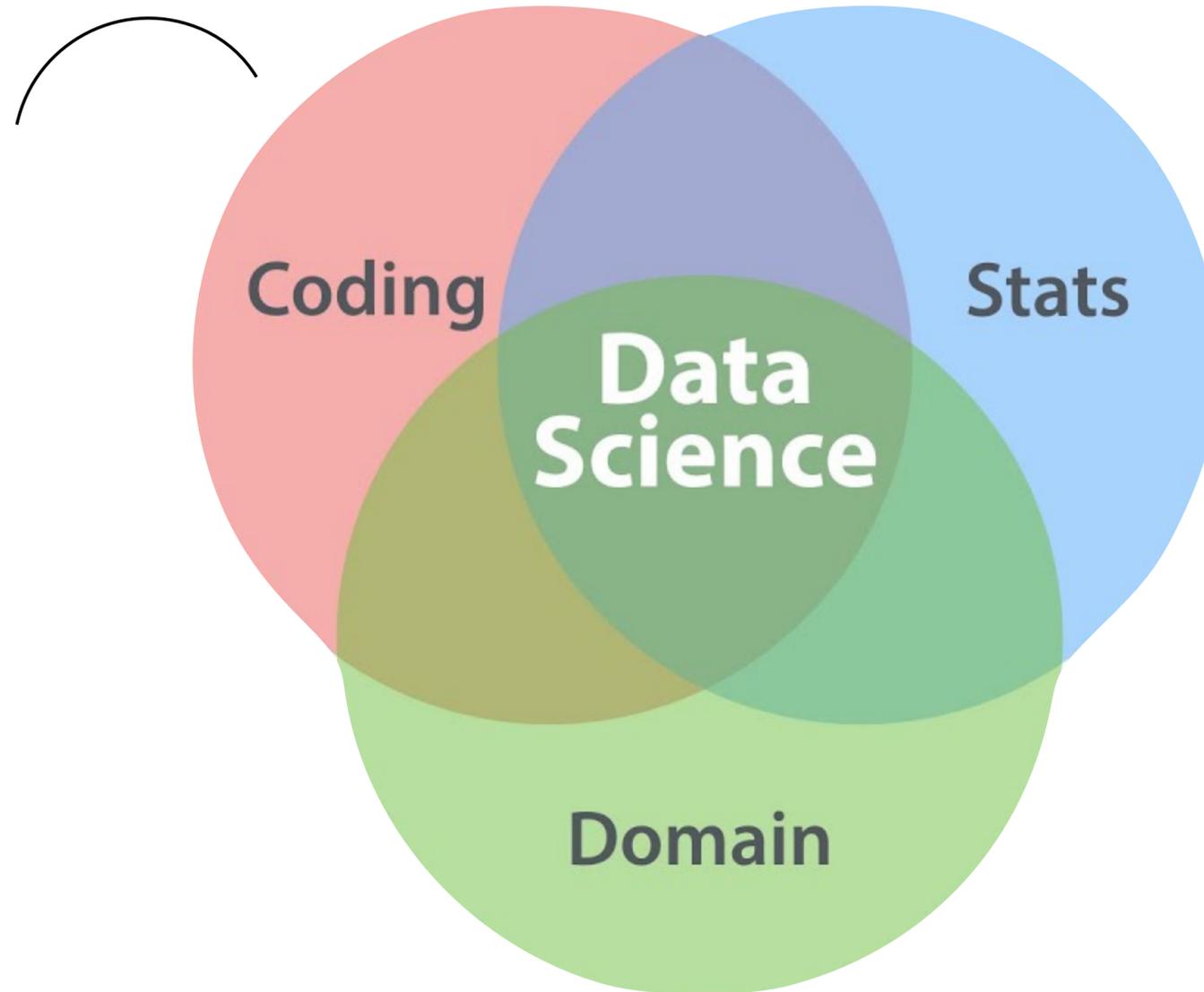# Instructor

Eysa Lee
eylee@barnard.edu

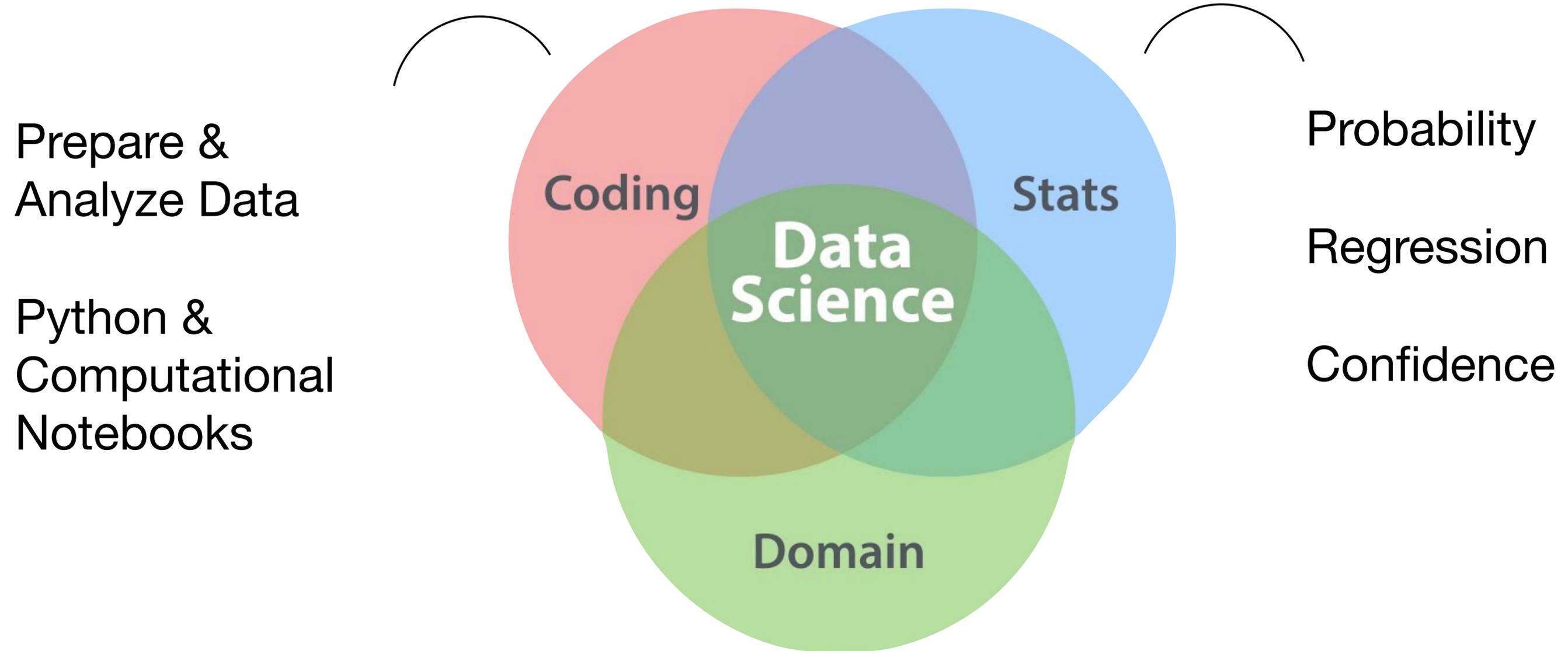# What is Data Science?

# What is Data Science?

Prepare &
Analyze Data

Python &
Computational
Notebooks

# What is Data Science?

Prepare &
Analyze Data

Python &
Computational
Notebooks

Probability

Regression

Confidence

# What is Data Science?



Prepare &
Analyze Data

Python &
Computational
Notebooks

Machine
Learning

Coding

Stats

Data
Science

Domain

Probability

Regression

Confidence

# What is Data Science?



Prepare & Analyze Data

Python & Computational Notebooks

Machine Learning

Coding

Data Science

Stats

Domain

Probability

Regression

Confidence

Identifying relevant & meaningful applications

# What is Data Science?



Prepare & Analyze Data

Python & Computational Notebooks

Probability

Regression

Confidence

Identifying relevant & meaningful applications

# What is Data Science?



Prepare & Analyze Data

Python & Computational Notebooks

Probability

Regression

Confidence

Identifying relevant & meaningful applications

# What is Data Science?

Data science is about drawing useful conclusions from large and diverse data sets through…
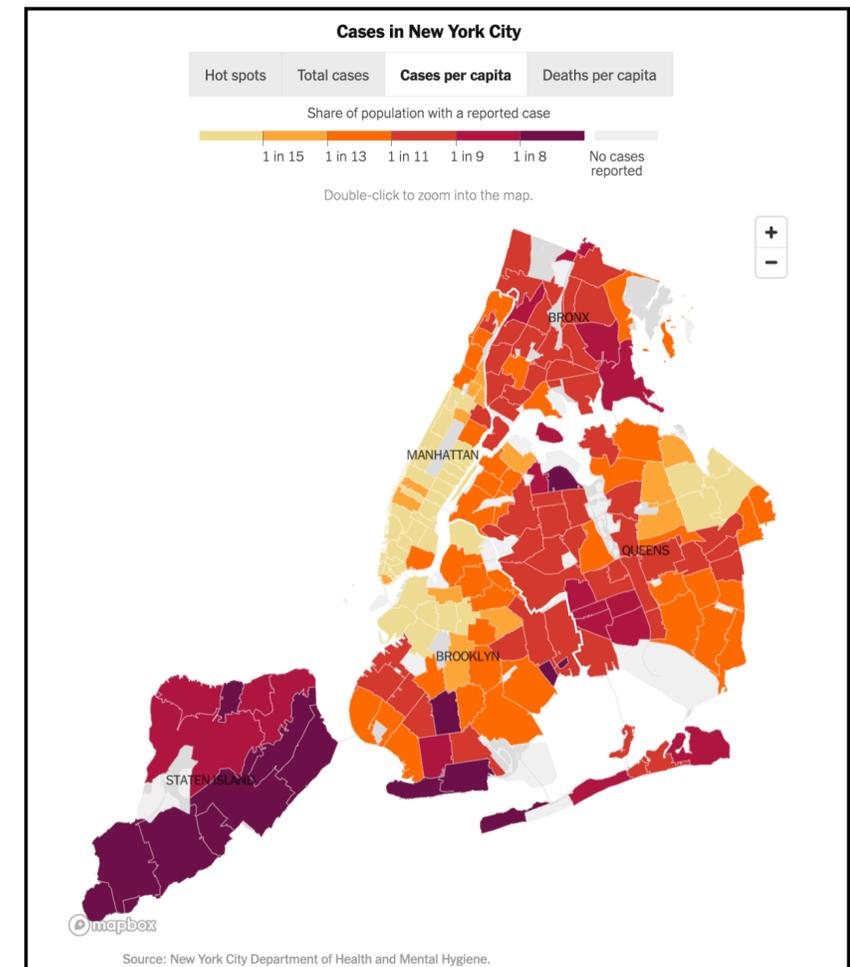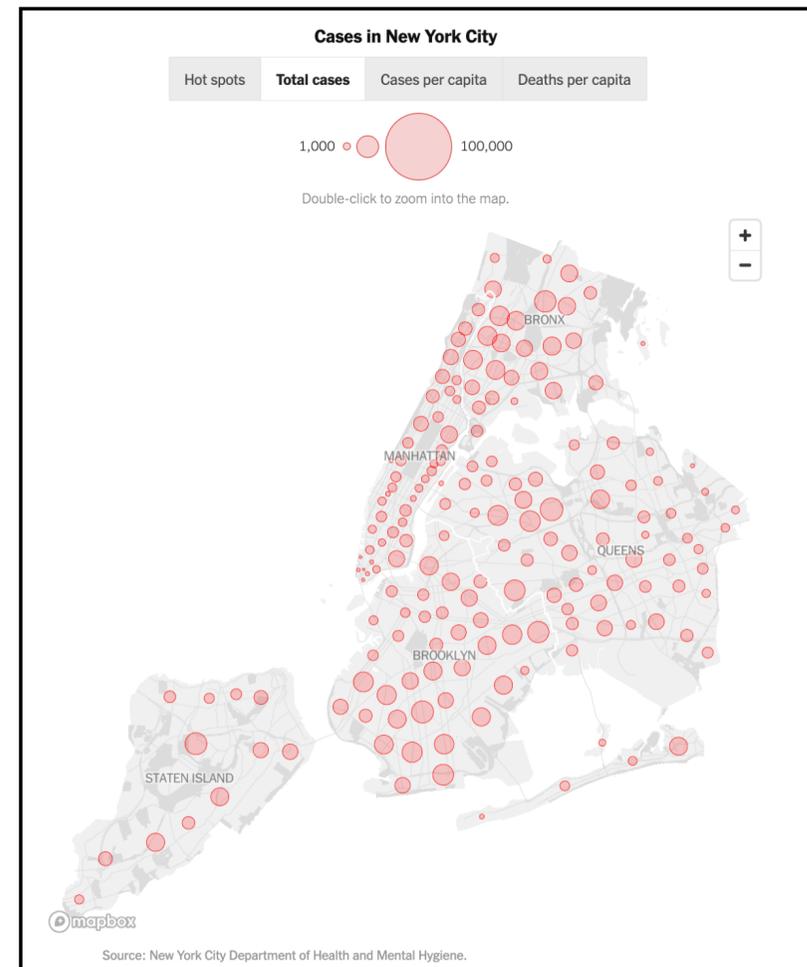
# What is Data Science?

Data science is about drawing useful conclusions from large and diverse data sets through…

- Exploration: Identifying patterns and trends using data (e.g., through visualization)



**Source: https://www.nytimes.com/interactive/2020/nyregion/new-york-city-coronavirus-cases.html**
**Data as of May 25, 2021**

# What is Data Science?

Data science is about drawing useful conclusions from large and diverse data sets through…

- Exploration: Identifying patterns and trends using data (e.g., through visualization)

- Inference: Drawing reliable conclusions using statistics



**Source: Murad Megjhani and his AI image generator of choice**

# What is Data Science?

Data science is about drawing useful conclusions from large and diverse data sets through…

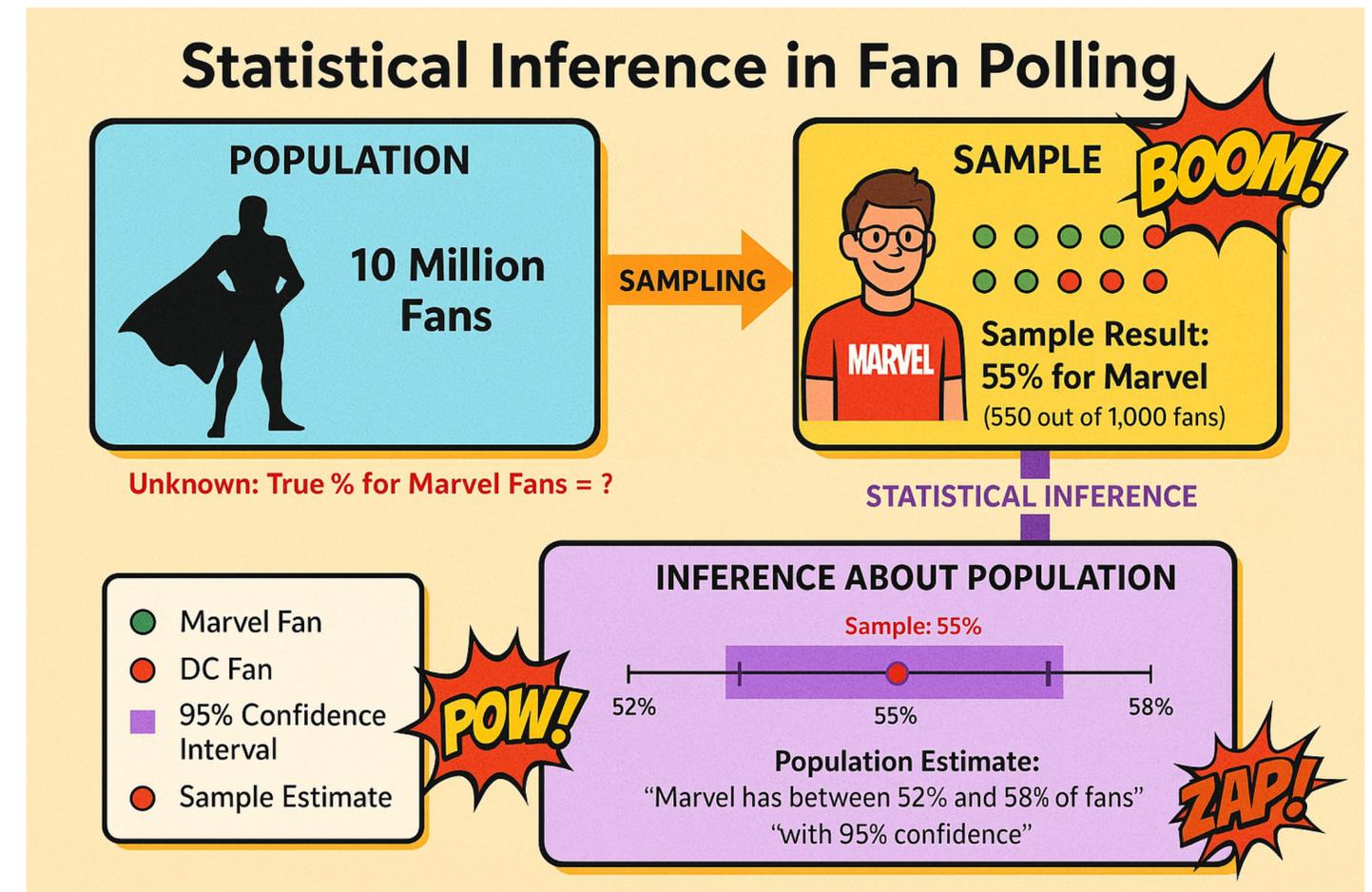- Exploration: Identifying patterns and trends using data (e.g., through visualization)

- Inference: Drawing reliable conclusions using statistics
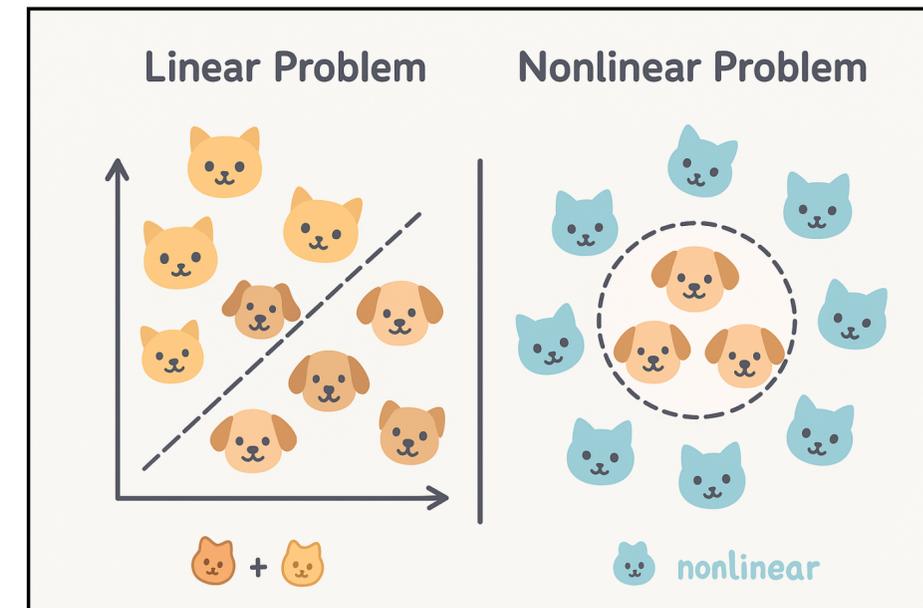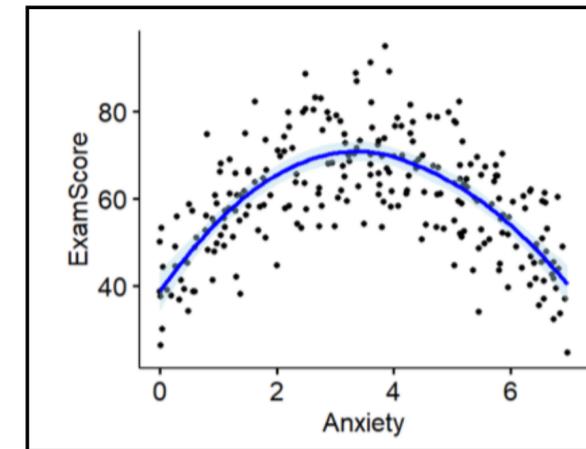
- Prediction: Making informed guesses about patterns using models

**Source: Murad Megjhani and his AI image generator of choice**

# Course Topics

# Course Topics

**Programming**

**Midterm Exam**

**Statistics**

**Final Project**

# Course Topics

**Programming**

Data Types

Iteration

Manipulating Arrays & Tables

Conditionals

Functions

Building Visualizations

**Midterm Exam**

**Statistics**

**Final Project**

# Course Topics

**Programming**

Data Types

Iteration

Manipulating Arrays & Tables

Conditionals

Functions

Building Visualizations

**Midterm Exam**

**Statistics**

Probabilities

Confidence Intervals

Correlation

Linear Regression

P-value & Statistical Significance

Residuals

**Final Project**

# Datasets You'll Explore

Climate Data

Vaccinations

Restaurant Reviews

Unemployment

Sports Records

Birth Rates

Movie Reviews

Compensation / Salaries
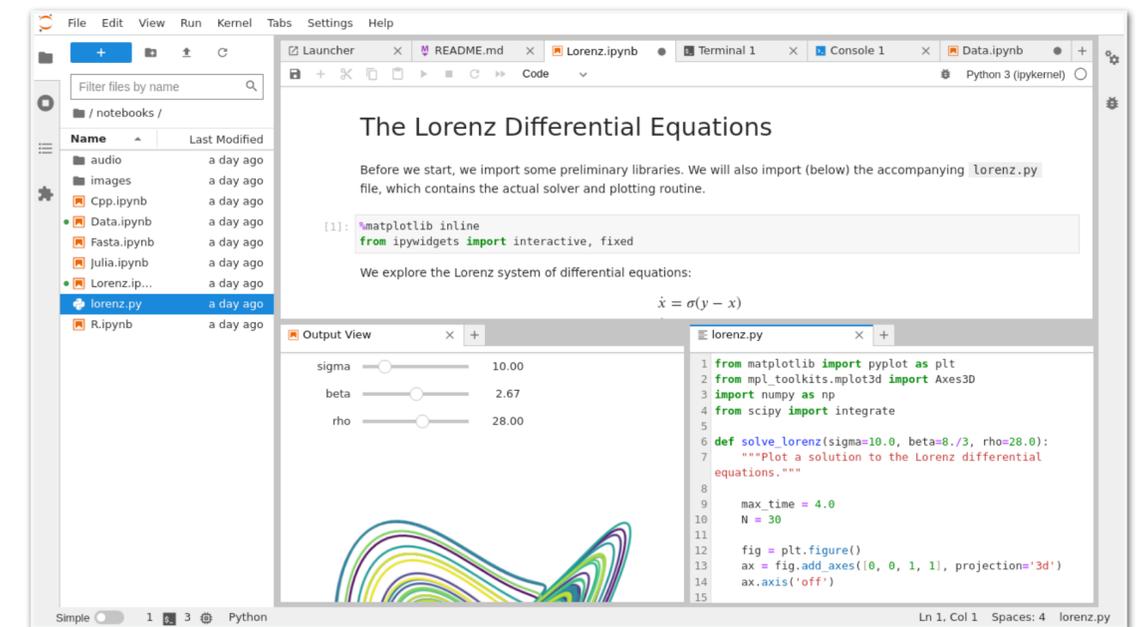
Happiness Scores

Ride Share Data

# What does Data look like?

- Tabular data typically in the form of a CSV

- Header row with clear field names

- You will use Jupyter Notebooks to read tabular data and perform analyses on it

| Num | Name | Type1 | Type2 | HP | Attack | Defense | SpAtk | SpDef | Speed |
|-----|------|-------|-------|-----|--------|---------|-------|-------|-------|
| 1 | Bulbasaur | Grass | Poison | 45 | 49 | 49 | 65 | 65 | 45 |
| 2 | Ivysaur | Grass | Poison | 60 | 62 | 63 | 80 | 80 | 60 |
| 3 | Venusaur | Grass | Poison | 80 | 82 | 83 | 100 | 100 | 80 |
| 3 | VenusaurMega Venusaur | Grass | Poison | 80 | 100 | 123 | 122 | 120 | 80 |
| 4 | Charmander | Fire | | 39 | 52 | 43 | 60 | 50 | 65 |
| 5 | Charmeleon | Fire | | 58 | 64 | 58 | 80 | 65 | 80 |
| 6 | Charizard | Fire | Flying | 78 | 84 | 78 | 109 | 85 | 100 |
| 6 | CharizardMega Charizard X | Fire | Dragon | 78 | 130 | 111 | 130 | 85 | 100 |
| 6 | CharizardMega Charizard Y | Fire | Flying | 78 | 104 | 78 | 159 | 115 | 100 |
| 7 | Squirtle | Water | | 44 | 48 | 65 | 50 | 64 | 43 |
| 8 | Wartortle | Water | | 59 | 63 | 80 | 65 | 80 | 58 |
| 9 | Blastoise | Water | | 79 | 83 | 100 | 85 | 105 | 78 |
| 9 | BlastoiseMega Blastoise | Water | | 79 | 103 | 120 | 135 | 115 | 78 |
| 10 | Caterpie | Bug | | 45 | 30 | 35 | 20 | 20 | 45 |
| 11 | Metapod | Bug | | 50 | 20 | 55 | 25 | 25 | 30 |
| 12 | Butterfree | Bug | Flying | 60 | 45 | 50 | 90 | 80 | 70 |
| 13 | Weedle | Bug | Poison | 40 | 35 | 30 | 20 | 20 | 50 |

# What are notebooks?

- Jupyter notebooks are environments for creating and sharing computational documents

  - Combination of notes (text and comments), code, data, and figures

- Data science is typically done in Jupyter notebooks using Python

  - Python has a rich developer community & set of libraries made for data science



https://docs.jupyter.org/en/latest/#what-is-a-notebook

# What will you learn to do?

- Take a dataset and explore it with visualizations

- Write Python code to support your experimentation

- Uncover interesting patterns and insights that might help you understand the data in a new way

# What is this course?

- Introduction to Computational Thinking and Data Science!

  - There is a lab section (BC 1017) associated with this course

  - You can register for any lab section, but you must be registered for one to take this course!

# Course Acknowledgments

- Builds on top of Data 8 (Berkeley Data Science course)

  - Lots of universities build on top of it (UW, NYU, UCSD, McGill, Cornell, etc.)

- Their textbook is great!

  - https://inferentialthinking.com/chapters/intro.html

# TAs, Computing Fellows, & Lab Sections

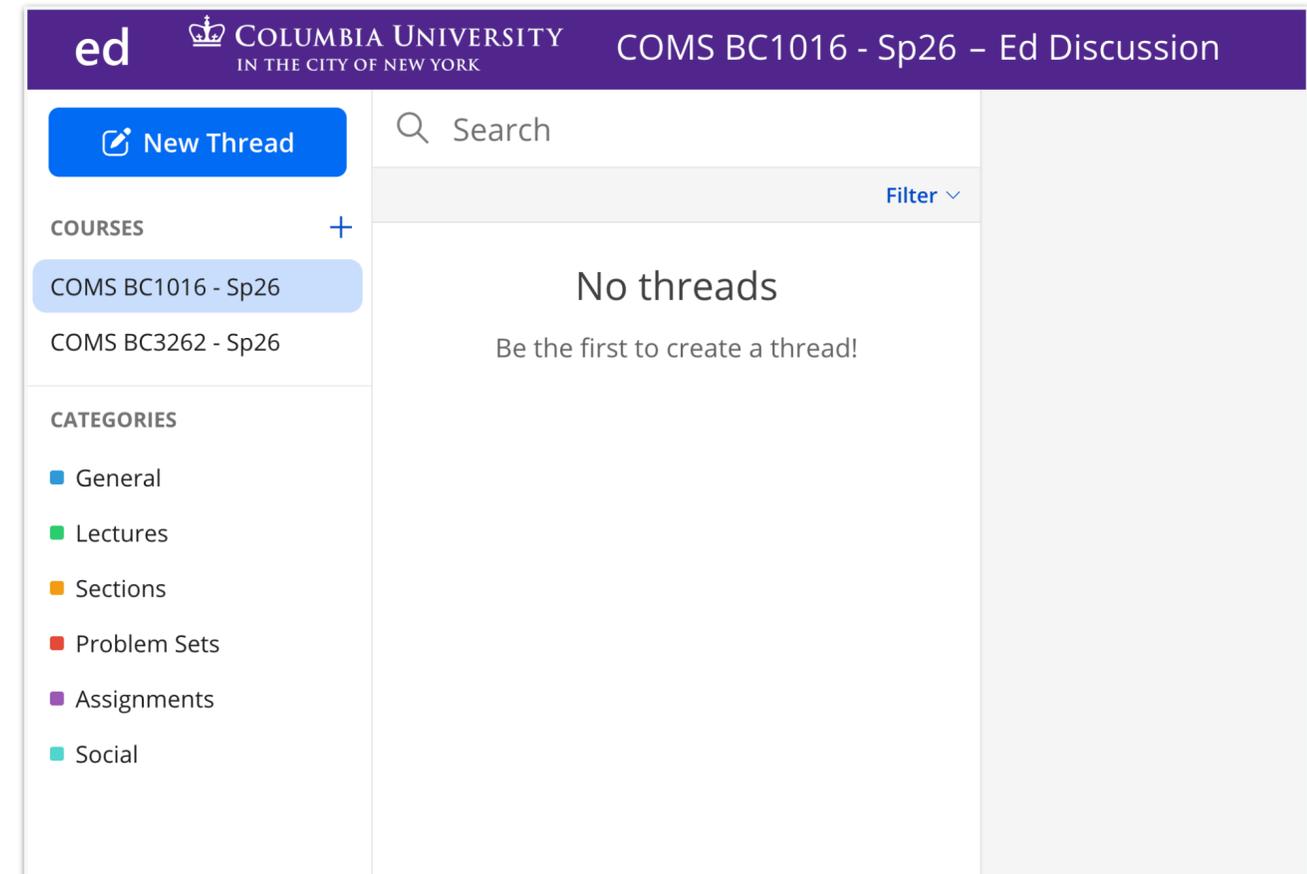| | TA | Computing Fellow |
|---|---|---|
| **Wednesday 2:10-3:40pm** | Nami Jain nbj2115@columbia.edu | Elena Lukac |
| **Thursday 9:00-10:30am** | Sathya Raman sr4213@columbia.edu | Madeline Gutierrez |

# Course Office Hours

- Office hours start week of Jan 26

- Professor Lee: Mondays 3:00pm-5:00pm (Milstein 512)

- TAs and Computing Fellows will offer 1.5 hours of OHs each week

  - Times TBA

# EdStem

- We also have a class discussion forum: https://edstem.org/us/courses/89029/discussion

    - The EdStem will help us not lose track of your questions

- The TAs, Computing Fellows, and I will be monitoring and answering questions

    - You are also able to post privately to instructors if needed

Note: labs start next week
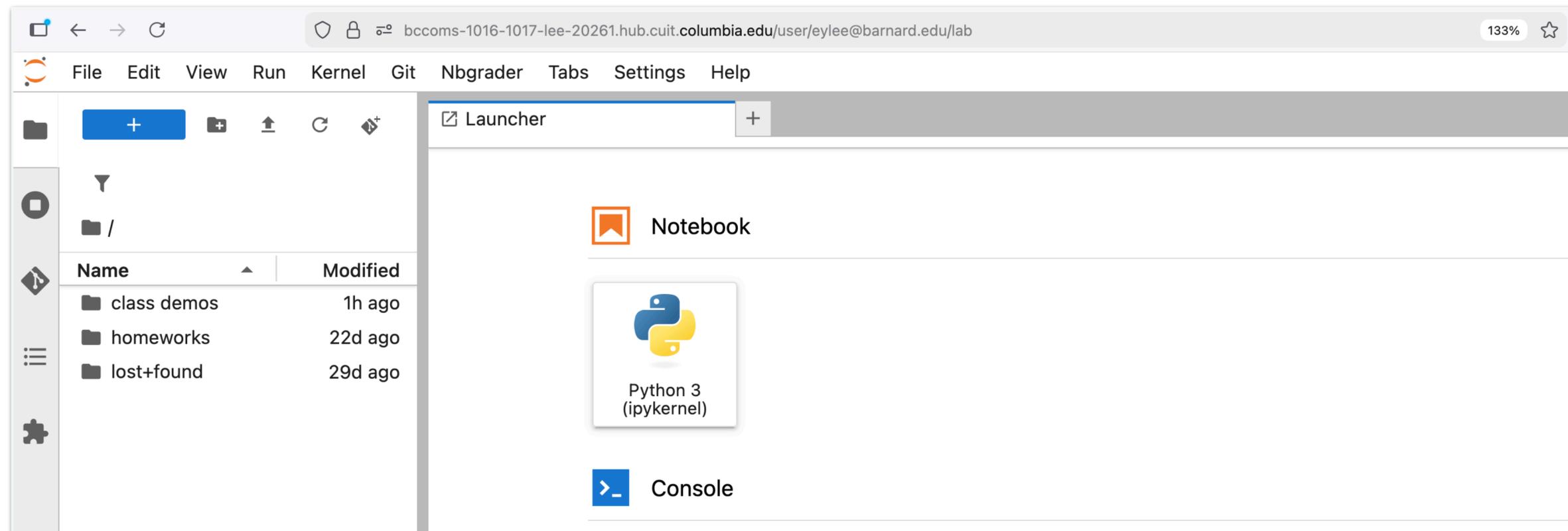**(no lab this week)**

# Course Expectations: Assignments

- Weekly Lab Assignments are due **Fridays 11:59pm** via Courseworks

  - Intended to be finished and submitted during lab itself

- (Mostly) weekly Homeworks, due **Wednesdays 11:59pm** through Gradescope (via Courseworks)

  - These may take a little bit more time… Start early

- The lowest lab and lowest homework grade will be dropped

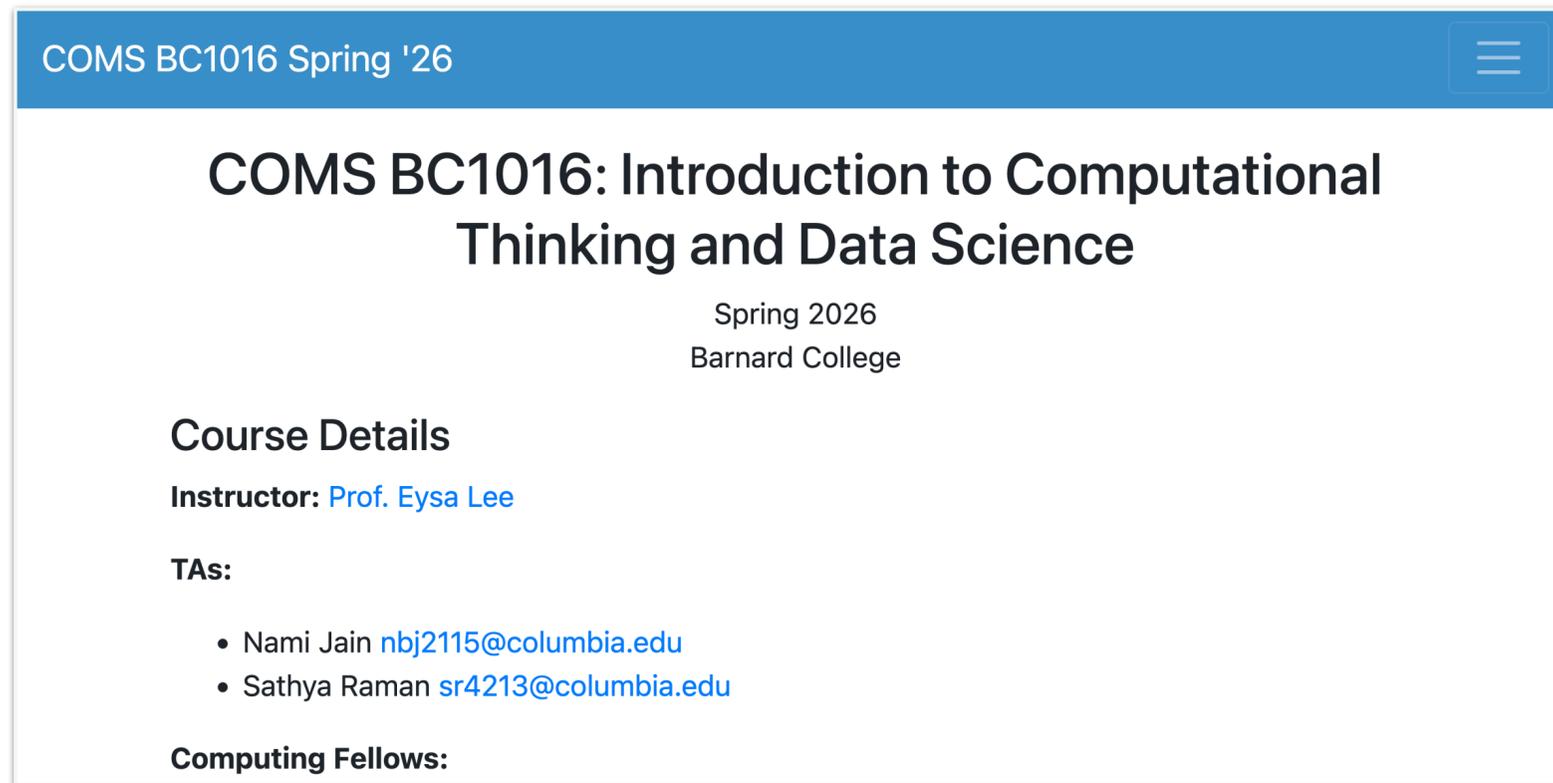| Monday | Tuesday | Wednesday | Thursday | Friday |
|---|---|---|---|---|
| Lecture<br><br>Homework Released | | Lecture<br><br>Lab<br><br>Lab Assignment Released | Lab | |
| | | HW due 11:59pm | | Lab due 11:59pm |

# Course Expectations: Assignments

- All assignments will be completed using cloud-based Jupyter notebooks

- You can access our course Jupyter Hub at: https://bccoms-1016-1017-lee-20261.hub.cuit.columbia.edu/

    - All you need is a web browser, no special software

# Course Website

https://www.eysalee.com/courses/s26/bc1016.html

- Contains syllabus, schedule (lectures, assignments, labs), and links to handy resources

- Slides and class demos will be uploaded or linked to from here

# Course Expectations: Lecture Grading

Your grade will be determined based on the following breakdown:

- 35%: Homework Assignments

- 25%: Midterm Exam

- 40%: Final Project

The lowest homework grade will be dropped.

# Course Expectations: Lab Grading

Labs are graded out of 10 points:

- 5 points: Lab Assignments

- 5 points: Attendance

**If you are going to be late or are unable to attend, email your lab TA in advance or you will receive 0 points for attendance.**

- You are permitted one unexcused absence from lab during the semester.

- The lowest lab grade will be dropped.

# Course Expectations: Lab Grading

Labs are graded out of 10 points:

- 5 points: Lab Assignments

- 5 points: Attendance

Assignment Grading:

- Complete and correct lab notebooks receive 5 points. Partially complete lab notebooks receive 3 points.

- Submit notebooks as a PDF via Courseworks.

# Course Expectations: Regrade Requests

- TAs will grade all assignments within one week of submission.

- Any regrade requests must be submitted **within 1 week** of your grade being received

  - We will not consider any regrades after this timeframe

- If you request a regrade, we reserve the right to lower your grade if the original grading was found to be too generous.

# Course Expectations: Late Policy

- Any late assignment (submitted after the due date) will be docked **10% of the total possible points per late day** for that assignment **up to five days**.

  - Any assignment submitted more than five days after the original due date will receive a no credit.

- This policy does not apply to the final project, which cannot be accepted after the due date except in exceptional circumstances.

# Course Expectations: Generative AI Policy

- This course is meant to build your programming skills, so it is not advised to use generative AI tools.

    - We want you to build intuition about how to write code and fix common bugs!

    - Be aware generated code does not always represent best practices and may be verbose (or potentially incorrect!).

- **AI generated code or final report text is not permitted**

# Midterm Exam

- Paper exam happens during class **Wednesday, March 11, 2026**

  - This is the week before Spring Recess

- You will be allowed a note sheet to use as a reference during the exam

  - It will be submitted along with your exam

- If you need particular accommodations, please contact CARDS

# Final Project

- Groups of 2-3

- We will provide datasets to select from

- You will use the dataset to perform analyses using a combination of visualization and statistical analyses

- Final project report & Jupyter Notebook will be submitted during Finals week

- More info will be released after the midterm

# Syllabus Recap

- Info can be found on the course website: https://www.eysalee.com/courses/s26/bc1016.html

- Assignments done through the course JupyterHub: https://bccoms-1016-1017-lee-20261.hub.cuit.columbia.edu/

- **Labs start next week**

  - Email your TA if you will be late or absent

  - Attendance is 50% of your lab grade, 1 excused absence
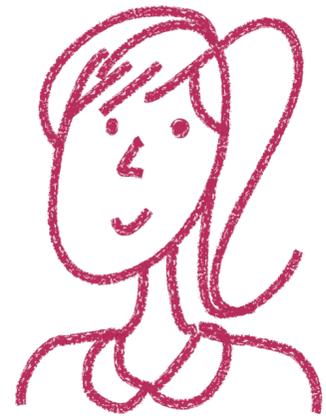
# Let's think about data

# What is Data Science?

Data science allows us to draw useful conclusions from large and diverse data sets through:

- Exploration: Identifying patterns and trends using data

- Inference: Drawing reliable conclusions using statistics

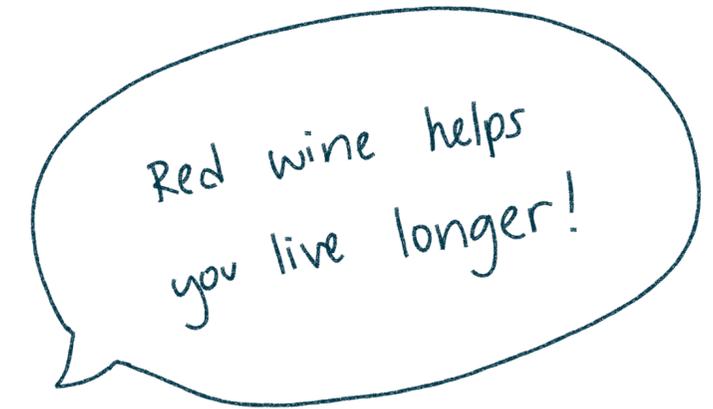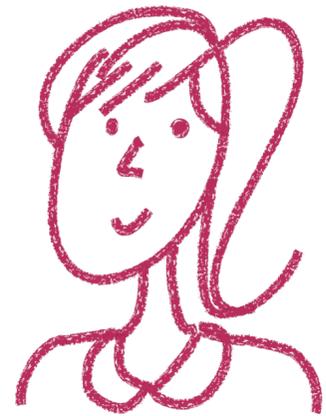- Prediction: Making informed guesses about patterns using models

# Cause and Effect

# A link between chocolate and health

# A link between red wine and health

# A link between cats and grades

# Observational Studies

Definition: A study in which researchers / scientists make conclusions based on *observed data* in which they had no hand in generating

- Individuals: study subjects, participants, units

- Treatment: factor of interest

- Outcome: result of treatment

# Observational Studies

Definition: A study in which researchers / scientists make conclusions based on *observed data* in which they had no hand in generating

→ people

- Individuals: study subjects, participants, units

→ chocolate consumption

- Treatment: factor of interest

- Outcome: result of treatment

→ heart health

Cardiac risk factors and prevention
Original article

Habitual chocolate consumption and risk of cardiovascular disease among healthy men and women

Chun Shing Kwok [1,2], S Matthijs Boekholdt [3], Marleen A H Lentjes [4], Yoon K Loke [5], Robert N Luben [4], Jessica K Yeong [6], Nicholas J Wareham [7], [iD] Phyo K Myint [1], Kay-Tee Khaw [4]

Correspondence to Dr Chun Shing Kwok, School of Medicine & Dentistry, University of Aberdeen, c/o Professor Phyo Kyaw Myint, Room 4:013, Polwarth Building, Foresterhill, Aberdeen AB25 2ZD, UK; phyo.myint@abdn.ac.uk

## Abstract

**Objective** To examine the association between chocolate intake and the risk of future cardiovascular events.

Source: https://heart.bmj.com/content/101/16/1279

# Association and Causality

Association: Any relation between the treatment and the outcome

Causality: if the treatment causes the outcome to occur

Cardiac risk factors and prevention
Original article

## Habitual chocolate consumption and risk of cardiovascular disease among healthy men and women

Chun Shing Kwok [1,2], S Matthijs Boekholdt [3], Marleen A H Lentjes [4], Yoon K Loke [5], Robert N Luben [4], Jessica K Yeong [6], Nicholas J Wareham [7], Phyo K Myint [1], Kay-Tee Khaw [4]

Correspondence to Dr Chun Shing Kwok, School of Medicine & Dentistry, University of Aberdeen, c/o Professor Phyo Kyaw Myint, Room 4:013, Polwarth Building, Foresterhill, Aberdeen AB25 2ZD, UK; phyo.myint@abdn.ac.uk
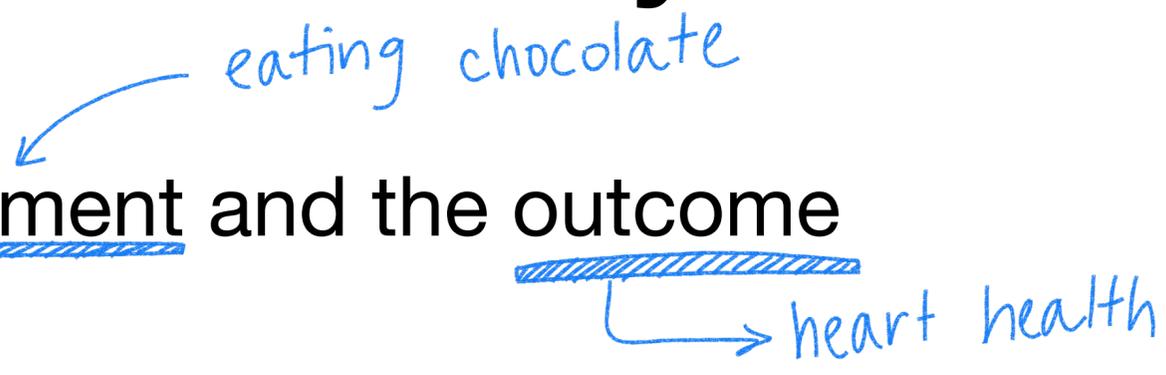
## Abstract

**Objective** To examine the association between chocolate intake and the risk of future cardiovascular events.

# Association and Causality

*eating chocolate*

Association: Any relation between the treatment and the outcome

*heart health*

Causality: if the treatment causes the outcome to occur

Habitual chocolate consumption and risk of cardiovascular disease among healthy men and women

Chun Shing Kwok [1,2], S Matthijs Boekholdt [3], Marleen A H Lentjes [4], Yoon K Loke [5], Robert N Luben [4], Jessica K Yeong [6],

Nicholas J Wareham [7], Phyo K Myint [1], Kay-Tee Khaw [4]

Correspondence to Dr Chun Shing Kwok, School of Medicine & Dentistry, University of Aberdeen, c/o Professor Phyo Kyaw Myint, Room 4:013, Polwarth Building, Foresterhill, Aberdeen AB25 2ZD, UK; phyo.myint@abdn.ac.uk

## Abstract

**Objective** To examine the association between chocolate intake and the risk of future cardiovascular events.

# Association and Causality

*eating chocolate*

Association: Any relation between the treatment and the outcome

*heart health*

Causality: if the treatment causes the outcome to occur

Is there an association between chocolate consumption and heart disease?

According to this study, yes ☺

Habitual chocolate consumption and risk of cardiovascular disease among healthy men and women

Chun Shing Kwok [1, 2], S Matthijs Boekholdt [3], Marleen A H Lentjes [4], Yoon K Loke [5], Robert N Luben [4], Jessica K Yeong [6], Nicholas J Wareham [7], (iD) Phyo K Myint [1], Kay-Tee Khaw [4]

Correspondence to Dr Chun Shing Kwok, School of Medicine & Dentistry, University of Aberdeen, c/o Professor Phyo Kyaw Myint, Room 4:013, Polwarth Building, Foresterhill, Aberdeen AB25 2ZD, UK; phyo.myint@abdn.ac.uk

## Abstract

**Objective** To examine the association between chocolate intake and the risk of future cardiovascular events.

**Conclusions** Cumulative evidence suggests that higher chocolate intake is associated with a lower risk of future cardiovascular events, although residual confounding cannot be excluded. There does not appear to be any evidence to say that chocolate should be avoided in those who are concerned about cardiovascular risk.

https://doi.org/10.1136/heartjnl-2014-307050

# Association and Causality

*eating chocolate*

Association: Any relation between the treatment and the outcome

*heart health*

Causality: if the treatment causes the outcome to occur

Is there an association between chocolate consumption and heart disease?

According to this study, yes ☺

Does eating chocolate lead to a reduction in heart disease?

Well... that's harder to say...

Cardiac risk factors and prevention
Original article

Habitual chocolate consumption and risk of cardiovascular disease among healthy men and women

Chun Shing Kwok [1, 2], S Matthijs Boekholdt [3], Marleen A H Lentjes [4], Yoon K Loke [5], Robert N Luben [4], Jessica K Yeong [6], Nicholas J Wareham [7], (iD) Phyo K Myint [1], Kay-Tee Khaw [4]

Correspondence to Dr Chun Shing Kwok, School of Medicine & Dentistry, University of Aberdeen, c/o Professor Phyo Kyaw Myint, Room 4:013, Polwarth Building, Foresterhill, Aberdeen AB25 2ZD, UK; phyo.myint@abdn.ac.uk

## Abstract

**Objective** To examine the association between chocolate intake and the risk of future cardiovascular events.

**Conclusions** Cumulative evidence suggests that higher chocolate intake is associated with a lower risk of future cardiovascular events, although residual confounding cannot be excluded. There does not appear to be any evidence to say that chocolate should be avoided in those who are concerned about cardiovascular risk.

https://doi.org/10.1136/heartjnl-2014-307050

# Confounding Factors

Definition: In observational studies, underlying difference(s) between two groups other than the treatment (that make it difficult to identify causality)

It's possible that people who like to eat chocolate do something else that offers heart protection, like eat a wide variety of healthful foods. One of the interesting things about this research is that participants in the non-chocolate group had higher average weight, more artery-damaging inflammation, more diabetes, were less physically active and had diets with the least amount of fat compared to chocolate eaters.

By **Howard E. LeWine, MD**, Chief Medical Editor, Harvard Health Publishing; Editorial Advisory Board Member, Harvard Health Publishing

In general, observational studies are not enough to determine causation

# Association and Causation

# London 1850s: Cholera Epidemic

- Large migration into London in the 1700-1800s, leading to overcrowding

- What's causing cholera to spread?
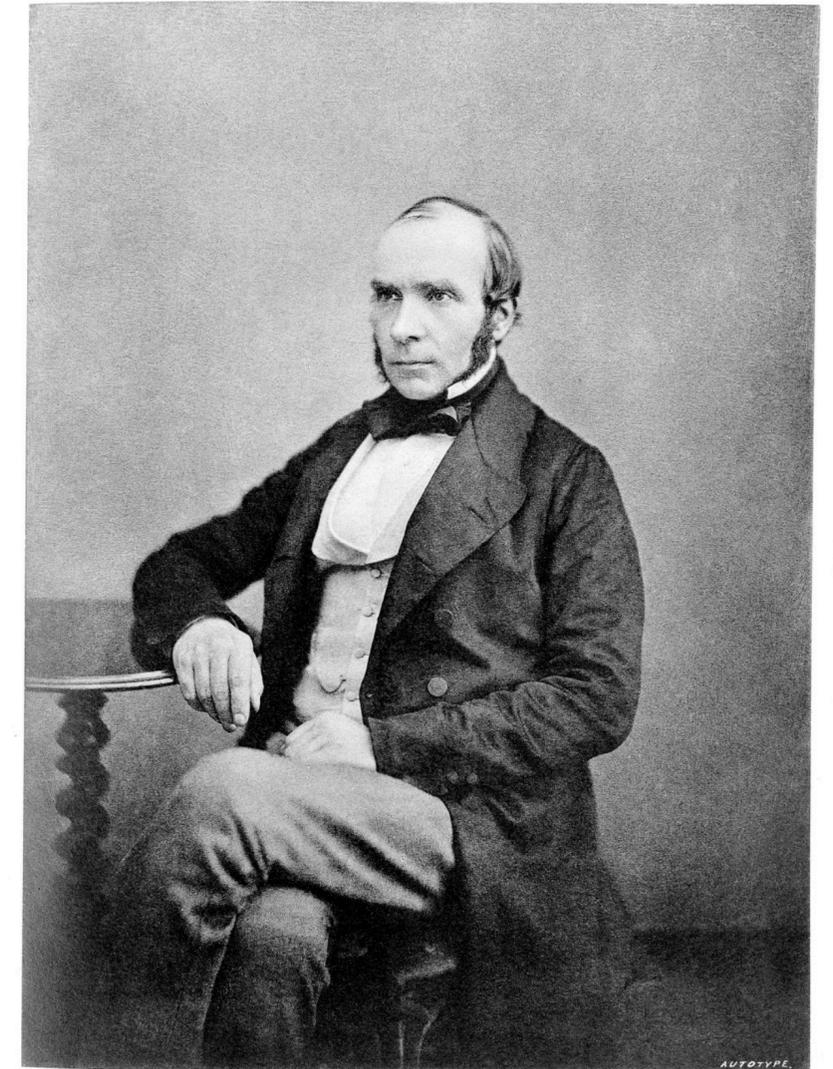


A COURT FOR KING CHOLERA.

https://www.sciencemuseum.org.uk/objects-and-stories/medicine/cholera-victorian-london

# Miasma

- 'Bad air' given off by waste and rotting matter

- Believed to the main source of how infectious disease spread

- Potential remedies: "Fly to clean air", "pocket full of posies", "fire off barrels of gunpowder"

- Popular medical theory of the time

  - Florence Nightingale (founder of modern nursing)

  - Edwin Chadwick (Commissioner of the Board of Health)

# John Snow (1813 - 1858)

- English physician

- Used data and visualizations to understand why cholera was spreading in the way that it was

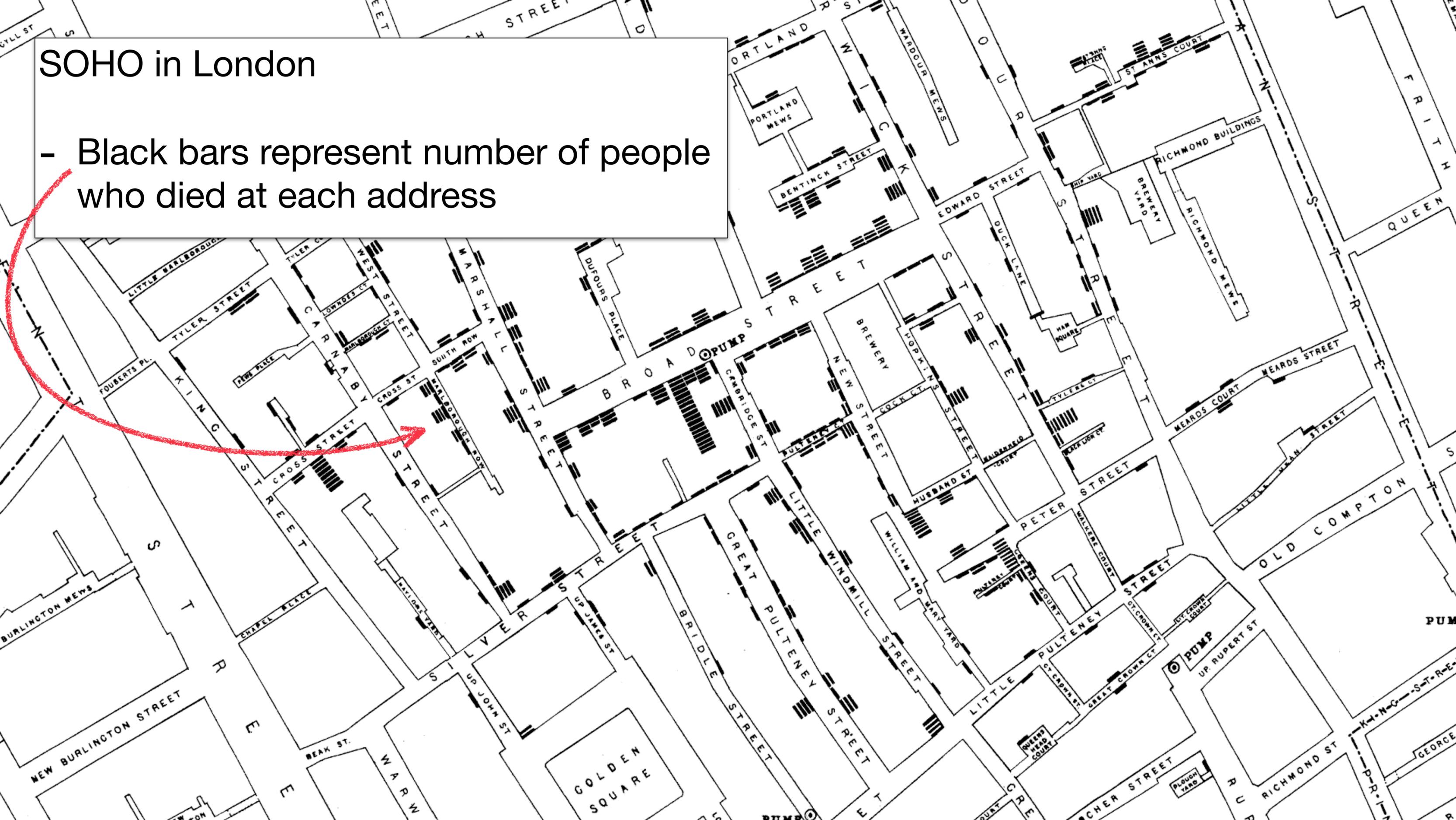  - Noticed people could be in close proximity without the same outcomes

not this dude

SOHO in London

– Black bars represent number of people who died at each address

GYLL ST

GREAT MARLBOROUGH STREET

PORTLAND

STREET

WARDOUR MEWS

ST ANNS COURT

STREET

PLACE

STREET

WORK
HOUSE

PORTLAND
MEWS

W
I
C
K

RICHMOND BUILDINGS

RICHMOND MEWS

ARGYLL

PLACE

PUMP

GREEN
DRAGON
YARD

LITTLE MARLBOROUGH ST.

MARSHALL ST.

TYLER COURT

WEST STREET

CARNABY

STREET

BENTINCK STREET

DUFOURS PLACE

STREET

EDWARD STREET

DUCK LANE

QUEEN

STRE

BREWERY
YARD

TYLER STREET

LOWNDES CT.

MARLBOROUGH ROW

SOUTH ROW

CROSS ST

MARLBOROUGH ROW

STREET

BROAD

PUMP

STREET

BROAD

PUMP

CAMBRIDGE ST

BREWERY

NEW STREET

HOPKINS STREET

MAIDENHEAD COURT

MAN SQUARE

MEARDS STREET

MEARDS COURT

FOUBERTS PL.

KING STREET

PEW PLACE

CROSS STREET

SILVER STREET

PULTENEY

COCK CT.

BLACK LION CT.

TYLERS CT.

PETER STREET

HUSBAND ST

WILLIAM AND MARY YARD

CASTLE COURT

OLD COMPTON

BURLINGTON MEWS

CHAPEL PLACE

NAYLORS YARD

UP JAMES ST

UP JOHN ST

GREAT PULTENEY STREET

LITTLE WINDMILL STREET

LITTLE PULTENEY STREET

CT. CROWN CT.

GT. CROWN CT.

LIT CROWN COURT

PUMP

UP. RUPERT ST

STREET

NEW BURLINGTON STREET

BEAK ST.

WARW

BRIDLE STREET

GOLDEN
SQUARE

QUEENS HEAD COURT

GREAT CROWN CT.

PLOUGH YARD

ACHER STREET

RICHMOND ST.

GEORGE

PUMP

K—I—N—G—S—T—R—E—E

GYLL ST

STREET

GREAT MARLBOROUGH STREET

PORTLAND STREET

WORK HOUSE

PORTLAND MEWS

WARDOUR MEWS

ST ANNS COURT

RICHMOND BUILDINGS

WICK STREET

BENTINCK STREET

HIP YARD

BREWERY YARD

RICHMOND MEWS

QUEEN

FRITH

ARGYLL PLACE

PUMP

GREEN DRAGON YARD

LITTLE MARLBOROUGH ST

TYLER COURT

MARSHALL ST.

WEST STREET

LOWNDES CT

MARLBOROUGH ROW

MARSHALL STREET

DUFOURS PLACE

EDWARD STREET

DUCK LANE

NAR SQUARE

CARNABY STREET

SOUTH ROW

CROSS ST

BROAD STREET

PUMP

CAMBRIDGE ST

BREWERY

HOPKINS STREET

MEARDS STREET

MEARDS COURT

TYLERS CT

FOUBERTS PL.

PEW PLACE

KING STREET

SILVER STREET

MARLBOROUGH

NEW STREET

PULTENEY

COCK CT

MAIDENHEAD COURT

BLACK LION CT

WALKERS COURT

PETER STREET

WILLIAM AND MARY YARD

HUSBAND ST

BURLINGTON MEWS

CHAPEL PLACE

NAYLORS TERRY

UP JAMES ST

UP JOHN ST

BRIDLE STREET

GREAT PULTENEY STREET

LITTLE WINDMILL STREET

GREAT

CT CROWN

CT CROWN

LITTLE PULTENEY

PUMP

UP RUPERT ST

NEW BURLINGTON STREET

BEAK ST.

WARW

GOLDEN SQUARE

QUEENS HEAD COURT

GREAT CROWN CT

UP CROWN COURT

ARCHER STREET

PLOUGH YARD

RICHMOND ST

GEORGE

PUMP

K-I-N-G--S-T-R-E

S-T-R-E
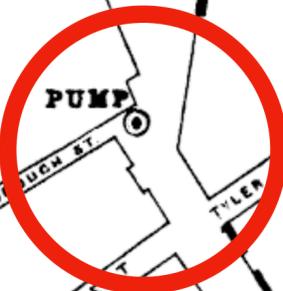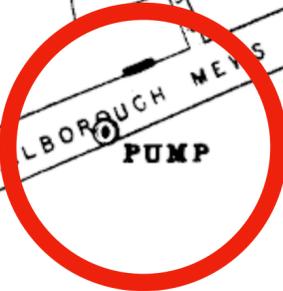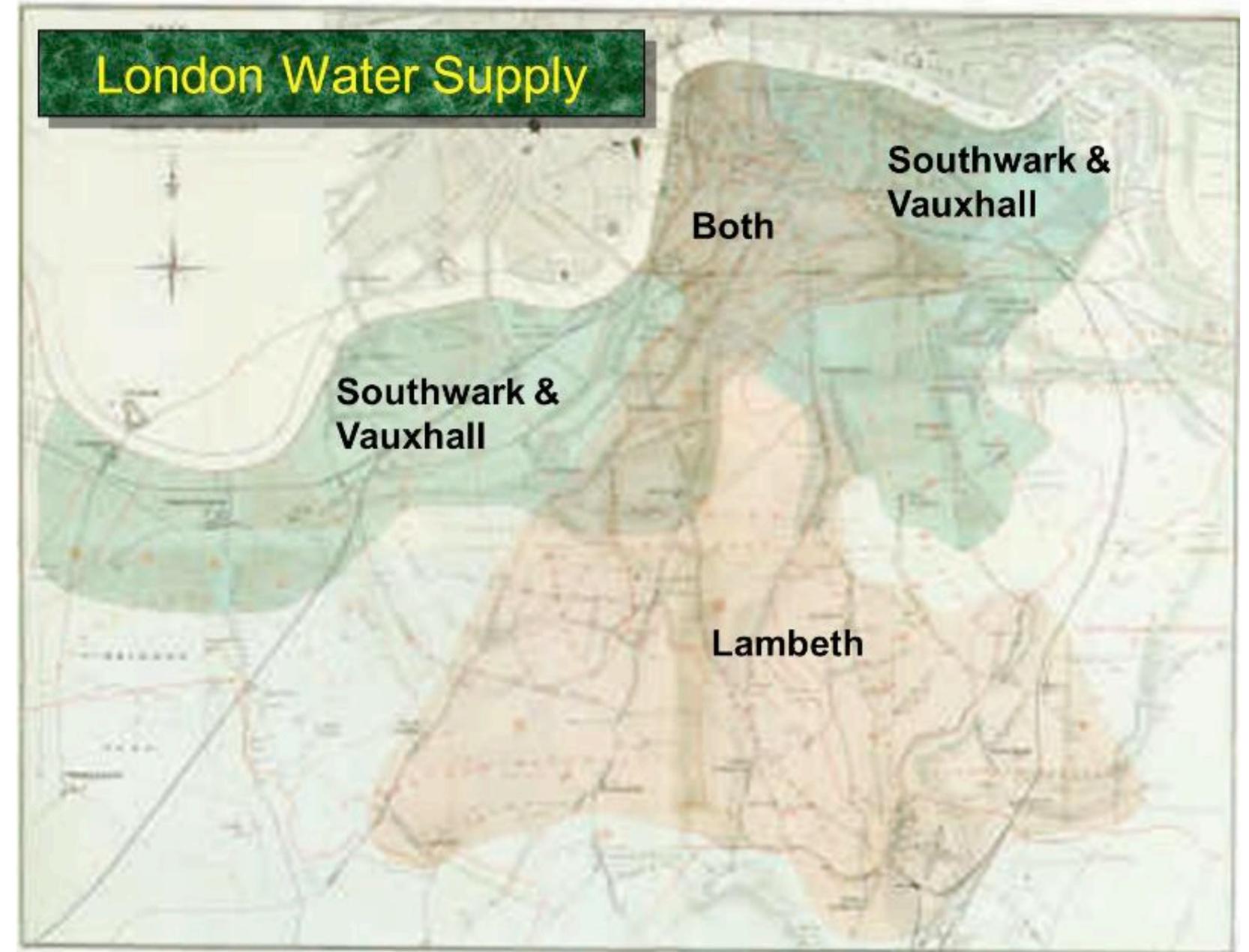
PUMP

Water Pumps

# Broad Street Pump

- Snow got pump handle removed

    - It was later discovered a cesspit had been leaking into the well

- Observational study

    - Snow did not control the pump or the people

    - Association was strong

# Showing Causation

Snow investigated cases in an area served by two water companies

- Southwark and Vauxhall's intake was in heavily polluted area of the Thames

- Lambeth's water intake was in a less polluted area upstream
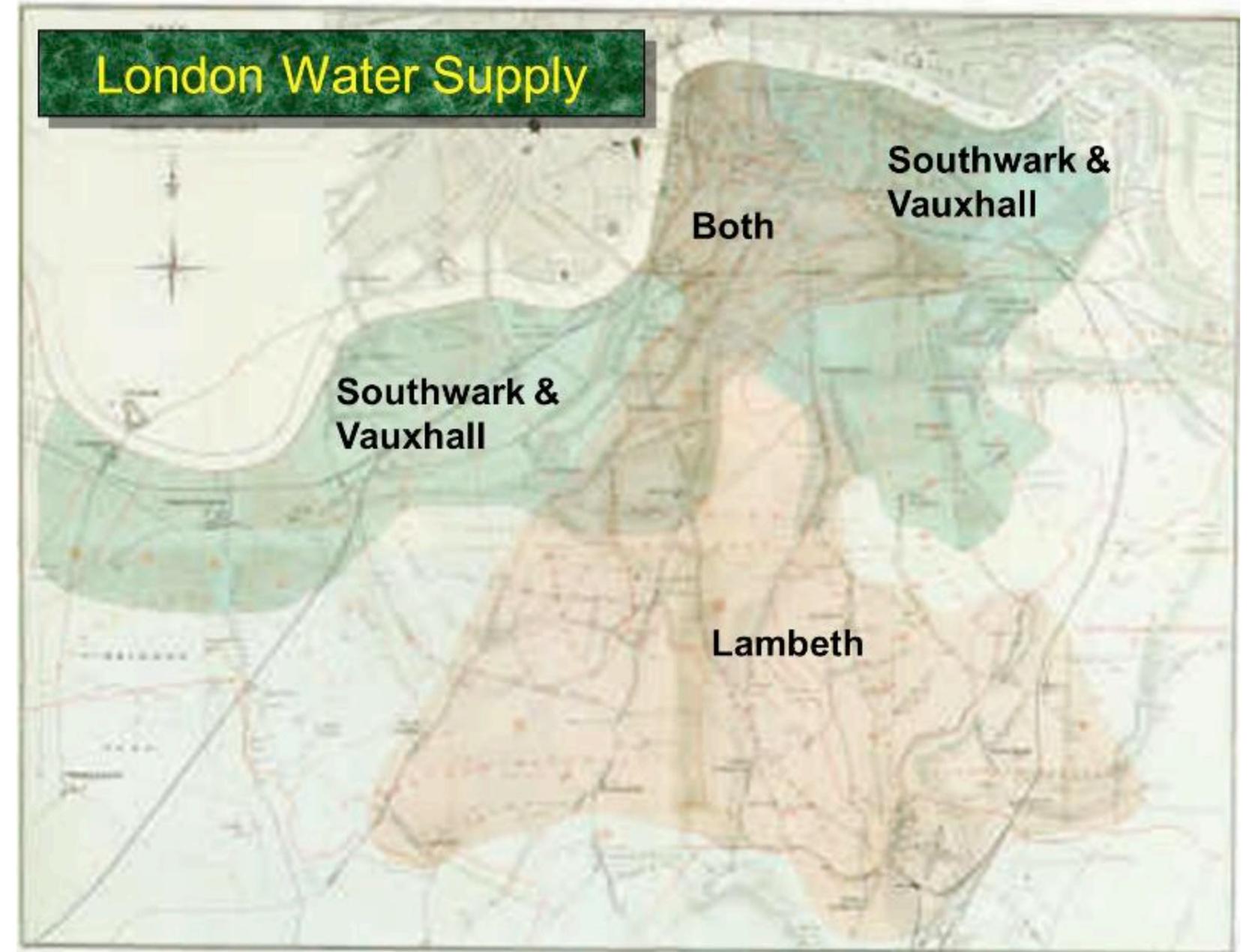
# Snow's "Grand Experiment"

Treatment group: S&V

Control group: Lambeth

"… there is no difference whatever in the houses or the people receiving the supply of the two Water Companies, or in any of the physical conditions with which they are surrounded…"

*Two groups are similar except for the treatment*


London Water Supply

# Snow's "Grand Experiment"

Deaths From Cholera Epidemic in Districts of London Supplied by Two Water Companies Over 7 Weeks, 1854

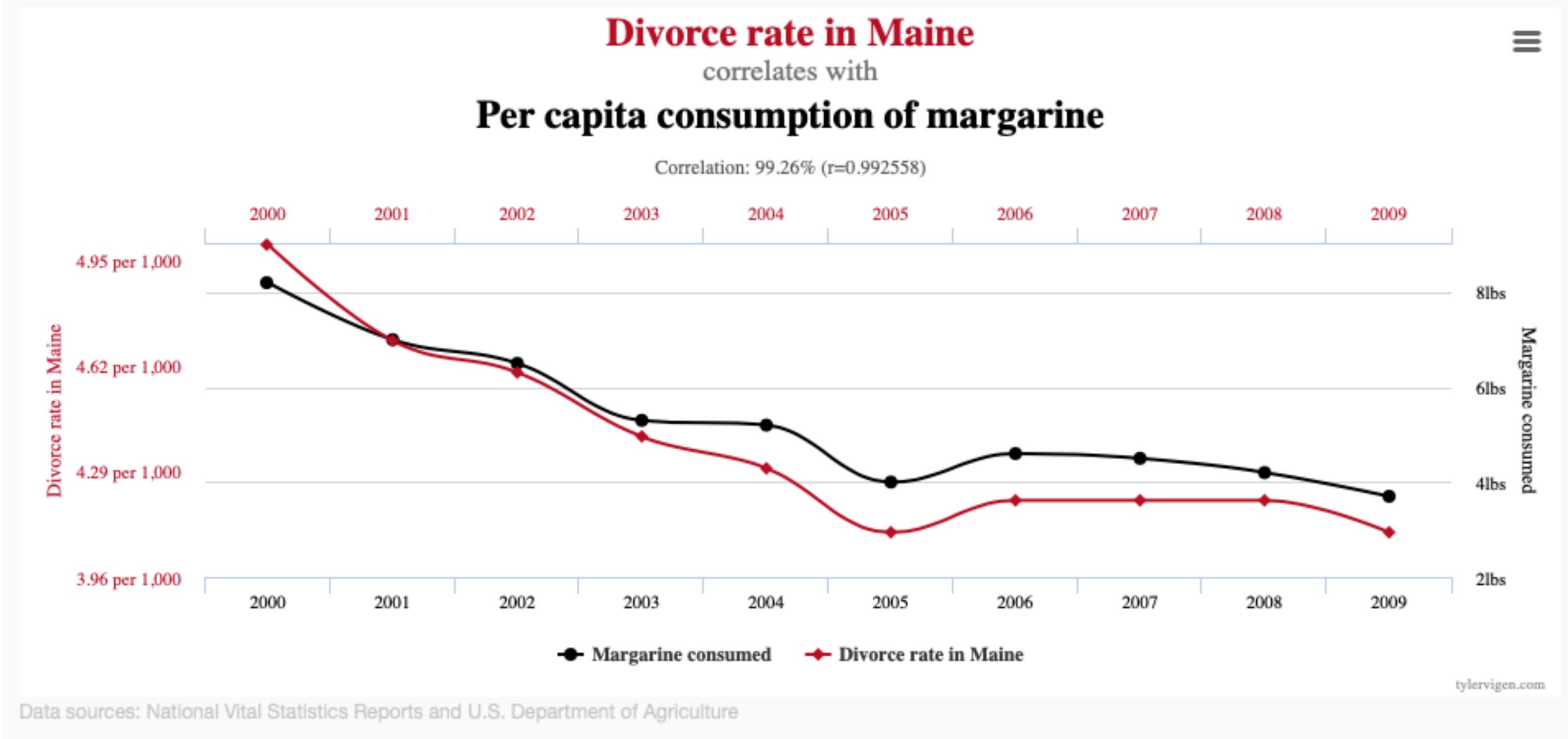| Water Supply Company | Number of Houses | Deaths From Cholera | Cholera Deaths per 10,000 Houses |
|---|---|---|---|
| Southwark and Vauxhall | 40,046 | 1,263 | 315 |
| Lambeth | 26,107 | 98 | 37 |
| Rest of London | 256,423 | 1,422 | 59 |

# Key to Establishing Causality

- <u>Treatment group</u>: Receives the treatment
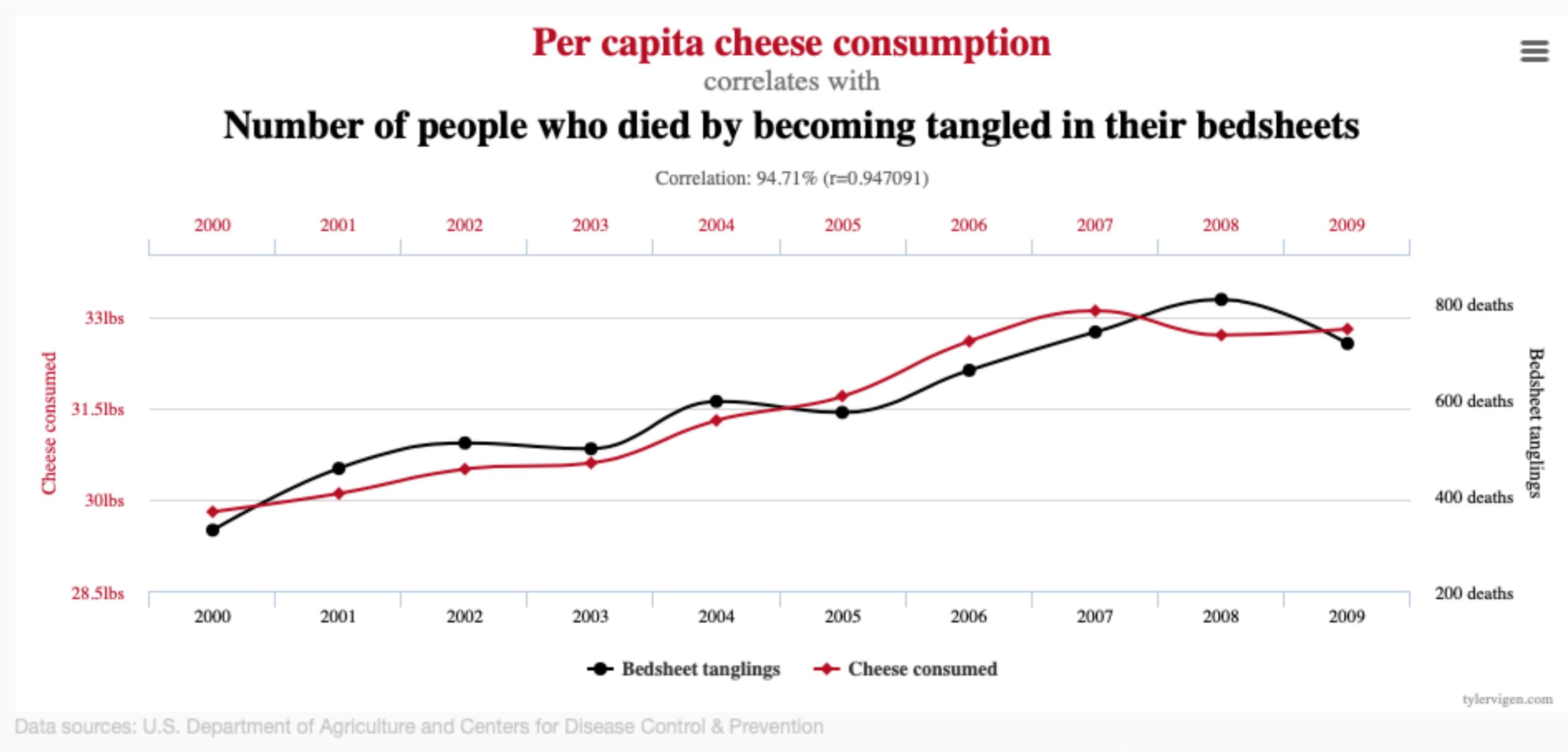
- <u>Control group</u>: Does not receive the treatment

If the treatment and control groups are *similar apart from the treatment*, then differences between the outcomes in the two groups can be ascribed to the treatment

# Confounding Factors / Spurious Correlations



**Divorce rate in Maine**
correlates with
**Per capita consumption of margarine**

Correlation: 99.26% (r=0.992558)

Data sources: National Vital Statistics Reports and U.S. Department of Agriculture

Source: https://www.tylervigen.com/spurious-correlations

# Confounding Factors / Spurious Correlations

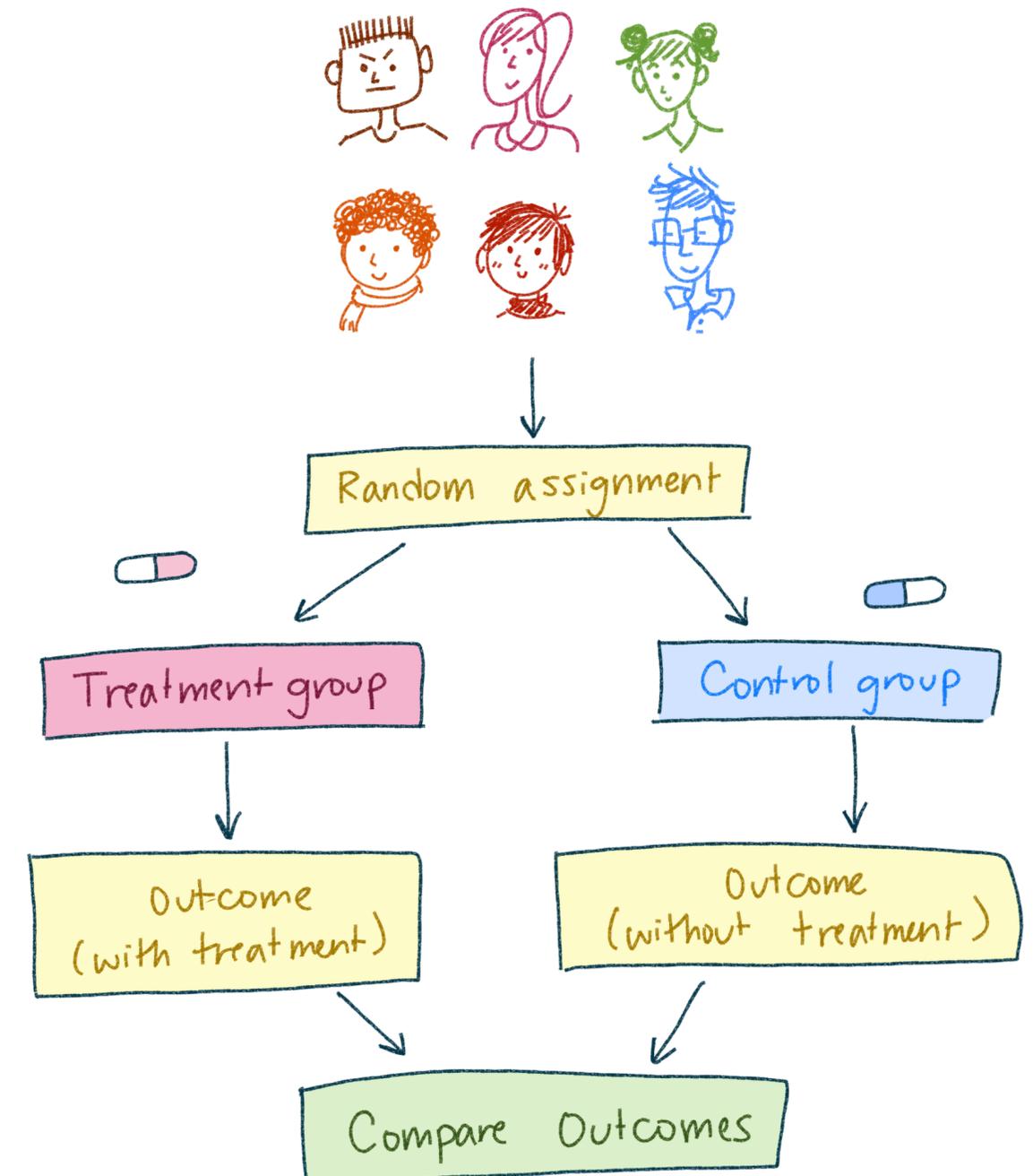

Source: https://www.tylervigen.com/spurious-correlations

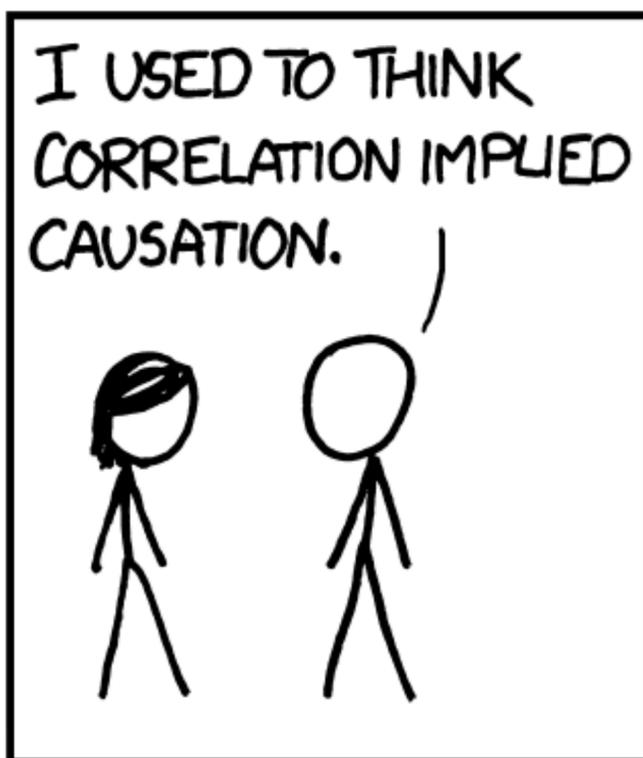# How to determine causation vs random correlation?

**Randomization**!

- If you assign individuals to treatment and control at random, then the two groups are likely to be similar apart from the treatment.

- Known as a **Randomized Controlled Experiment** or **Randomized Controlled Trial (RCT)**

- You can account (mathematically) for variability in the assignment

# CORRELATION

https://xkcd.com/552/

# Next Class

- Intro to Python: Expressions and Data Types

- Remember:

  - Labs start next week!!