



190F Foundations of Data Science  
Fall 2018

# Lecture 2

---

Cause and Effect

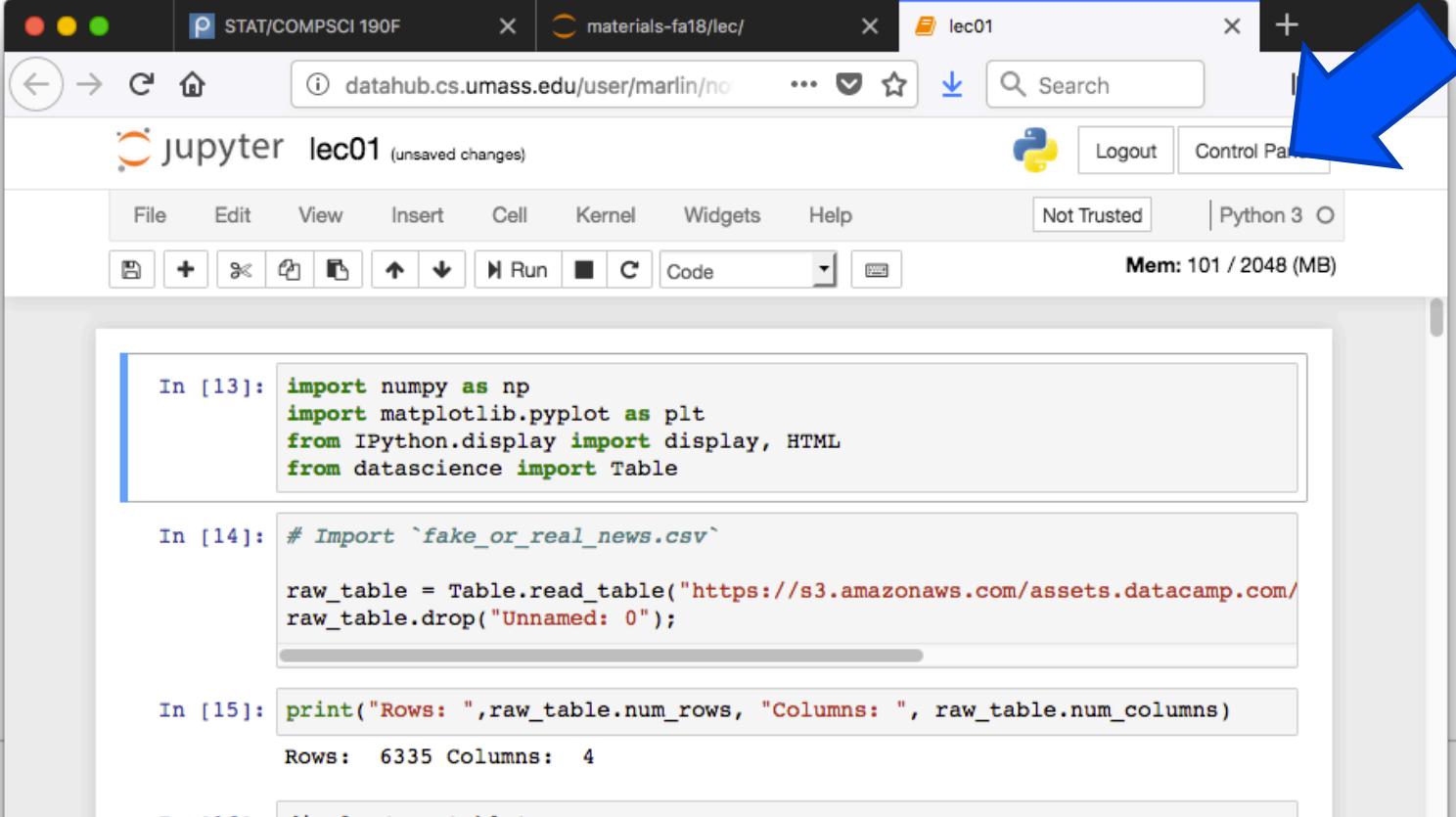
# Announcements

# Datahub

---

- If you have an @umass.edu email address, you should now be able to access the course's data hub.
  - Datahub uses UMass Google Apps authentication. Use your @umass.edu email address and Spire password to log in. It takes a minute to start up.
  - When you're done working with the Datahub, make sure to shut your datahub server down and then log out.
-

# Stopping Your Data Hub



The screenshot shows a Jupyter Notebook interface running in a web browser. The browser tabs are labeled "STAT/COMPSCI 190F", "materials-fa18/lec/", and "lec01". The main window title is "jupyter lec01 (unsaved changes)". The toolbar includes File, Edit, View, Insert, Cell, Kernel, Widgets, Help, Logout, and Control Panel. A status bar at the bottom indicates "Mem: 101 / 2048 (MB)". The notebook cells show the following code:

```
In [13]: import numpy as np
import matplotlib.pyplot as plt
from IPython.display import display, HTML
from datascience import Table

In [14]: # Import `fake_or_real_news.csv`

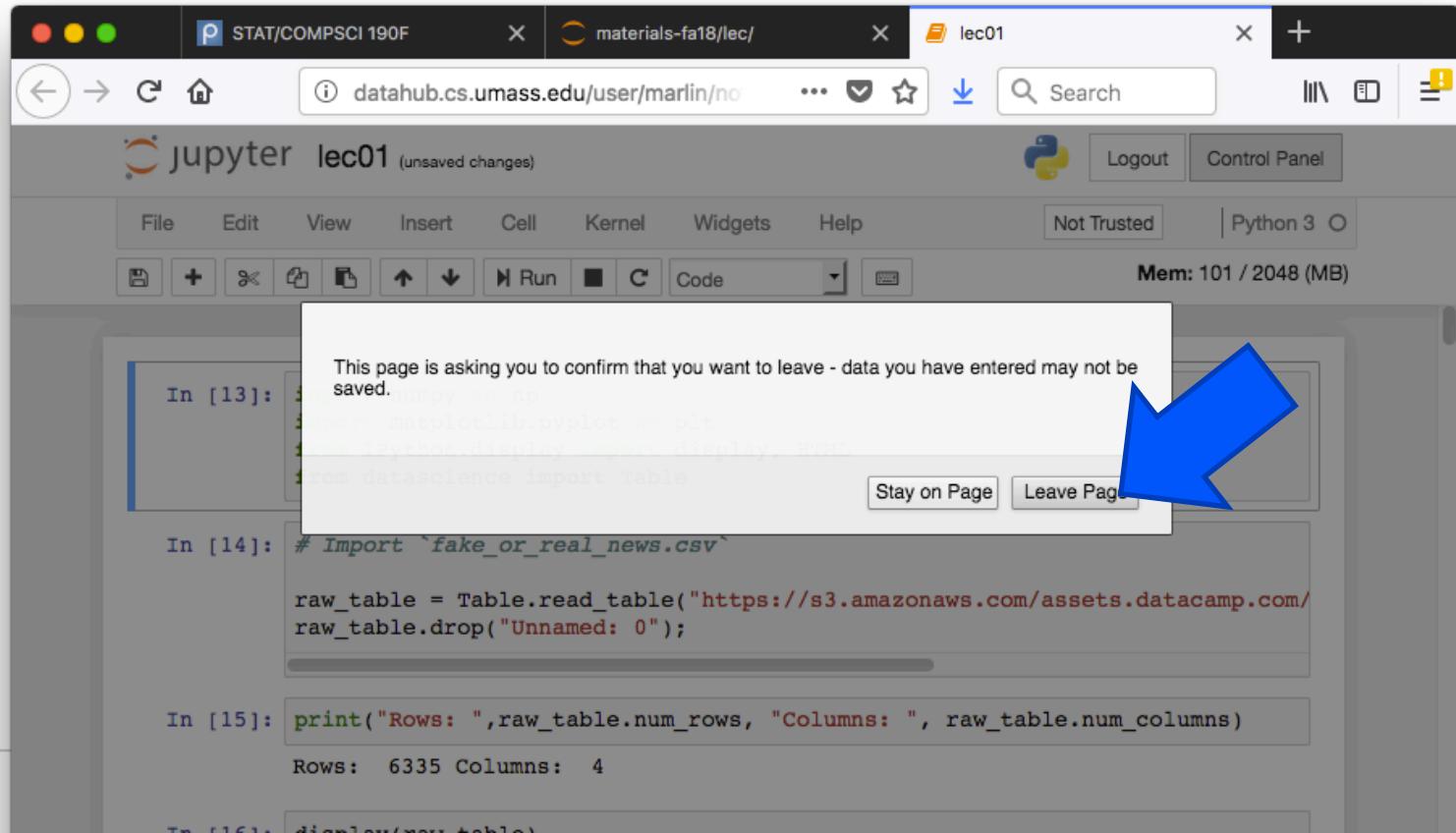
raw_table = Table.read_table("https://s3.amazonaws.com/assets.datacamp.com/
raw_table.drop("Unnamed: 0");

In [15]: print("Rows: ", raw_table.num_rows, "Columns: ", raw_table.num_columns)

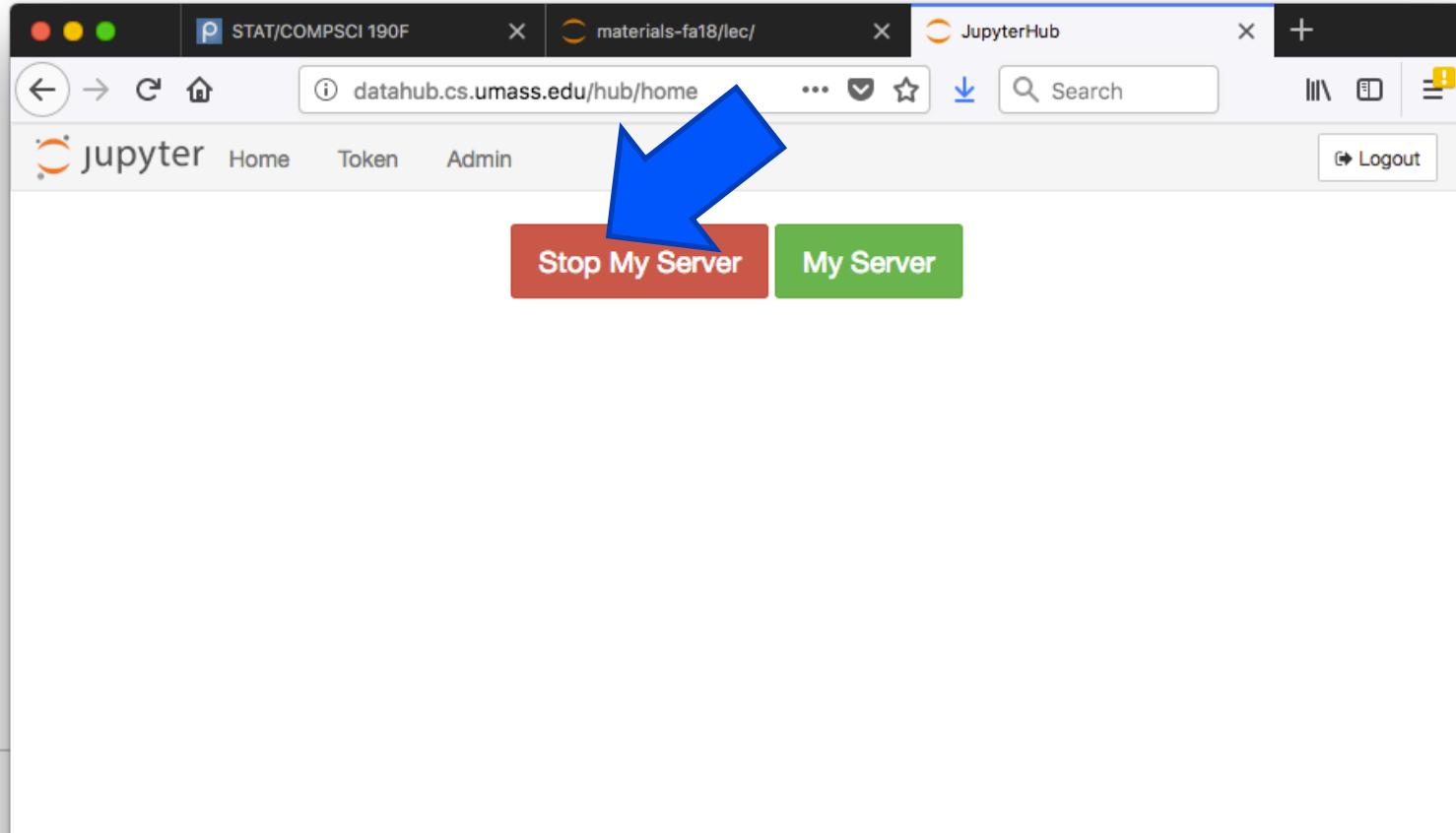
Rows: 6335 Columns: 4

In [16]: display(raw_table)
```

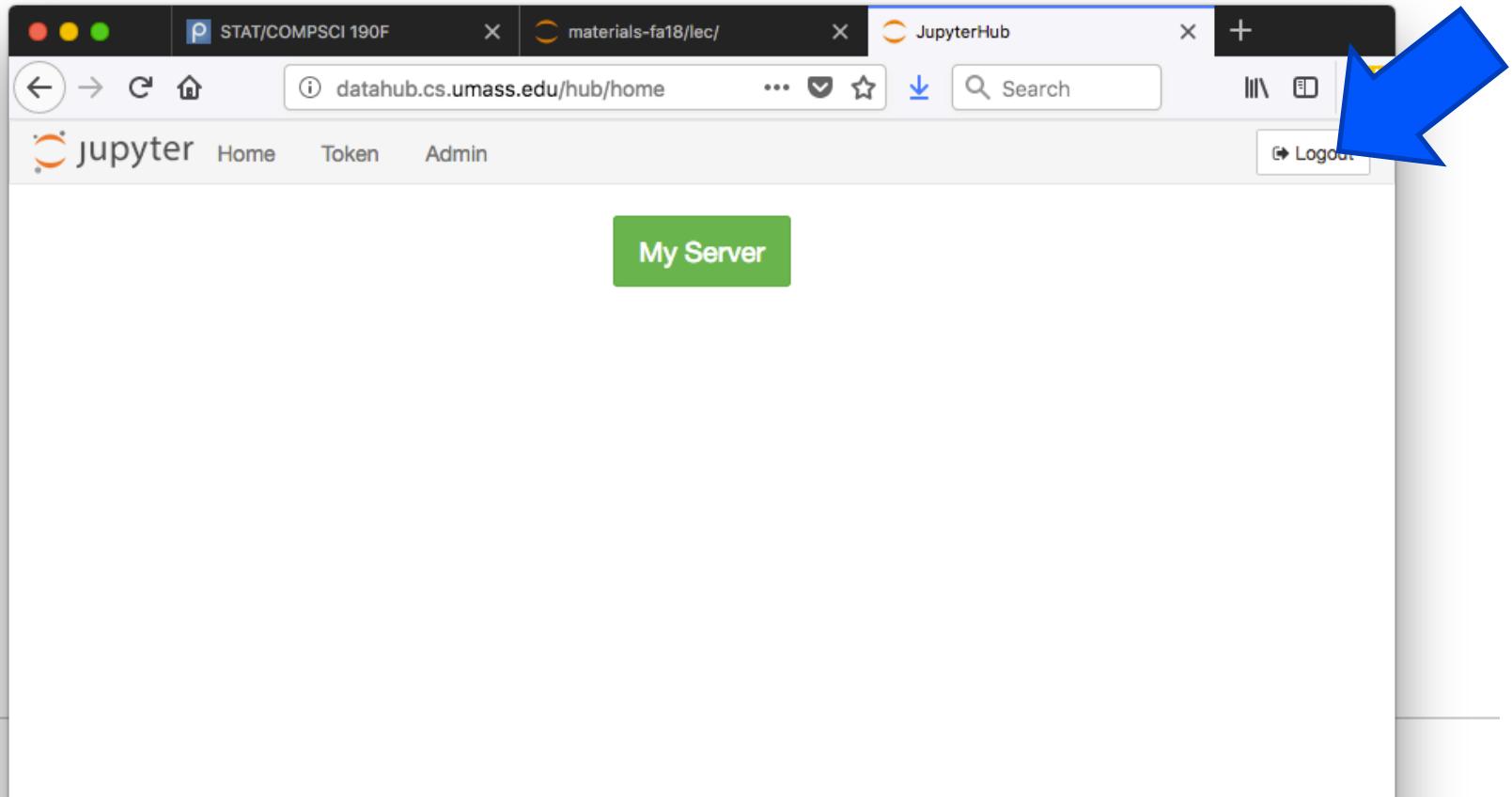
# Stopping Your Data Hub



# Stopping Your Data Hub



# Stopping Your Data Hub



# **Association vs Causation**

# Examples

---

No level of alcohol consumption is healthy, scientists say



Fo:

Dairy and meat 'beneficial for heart health and longevity'



Medic

Eating cheese may be associated with a lower risk of death –  
and that's good news for cheese lovers



that's good news for cheese lovers  
according to new study

Huddersfield Examiner • yesterday



# Chocolate and Heart Disease: Study

Chocolate, Chocolate, It's Good For Your Heart, Study Finds

June 19, 2015 · 5:03 AM ET  
Heard on [Morning Edition](#)



- **Population** (Individuals, study subjects, participants, units): 20K *European adults followed for 12 years.*
- **Treatment:** *chocolate consumption*
- **Outcome:** *heart disease*

# Chocolate and Heart Disease: Association

---

**Question 1:** Is there **any association** (any relationship) between chocolate consumption and heart disease?

- **Data:** “Among those in the top tier of chocolate consumption, 12 percent developed or died of cardiovascular disease during the study, compared to 17.4 percent of those who didn’t eat chocolate.”
  - **Answer:** Yes, this points to an **association**
-

# Chocolate and Heart Disease: Causation

---

**Question 2:** Does chocolate consumption lead to a reduction in heart disease?

- This question asks about **causality**
  - This question is often harder to answer.
  - “[The study] doesn’t prove a cause-and-effect relationship between chocolate and reduced risk of heart disease and stroke.” - *JoAnn Manson, chief of Preventive Medicine at Brigham and Women’s Hospital, Boston*
-

# Chocolate and Heart Disease: Alternatives

---

**Question 3:** Is the fact that people ate more chocolate the only possible cause for the observed effect of decreased heart disease risk?

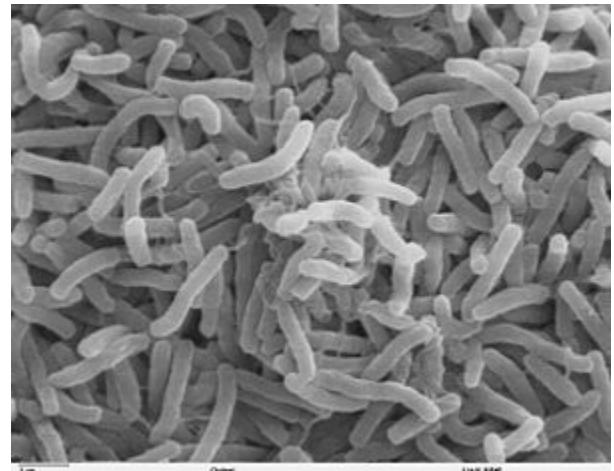
- For example, suppose the people who ate more chocolate tended to live in European countries with better health care?
  - What if wealthier people eat more chocolate and can also afford better health care?
-

# **Proving Cause and Effect**

# Broad Street Cholera Outbreak

---

- The Broad Street cholera outbreak was a severe outbreak of cholera that occurred in 1854 near Broad Street in the Soho district of London, England.
- This outbreak killed 616 people.



# Two Theories of Cholera

---

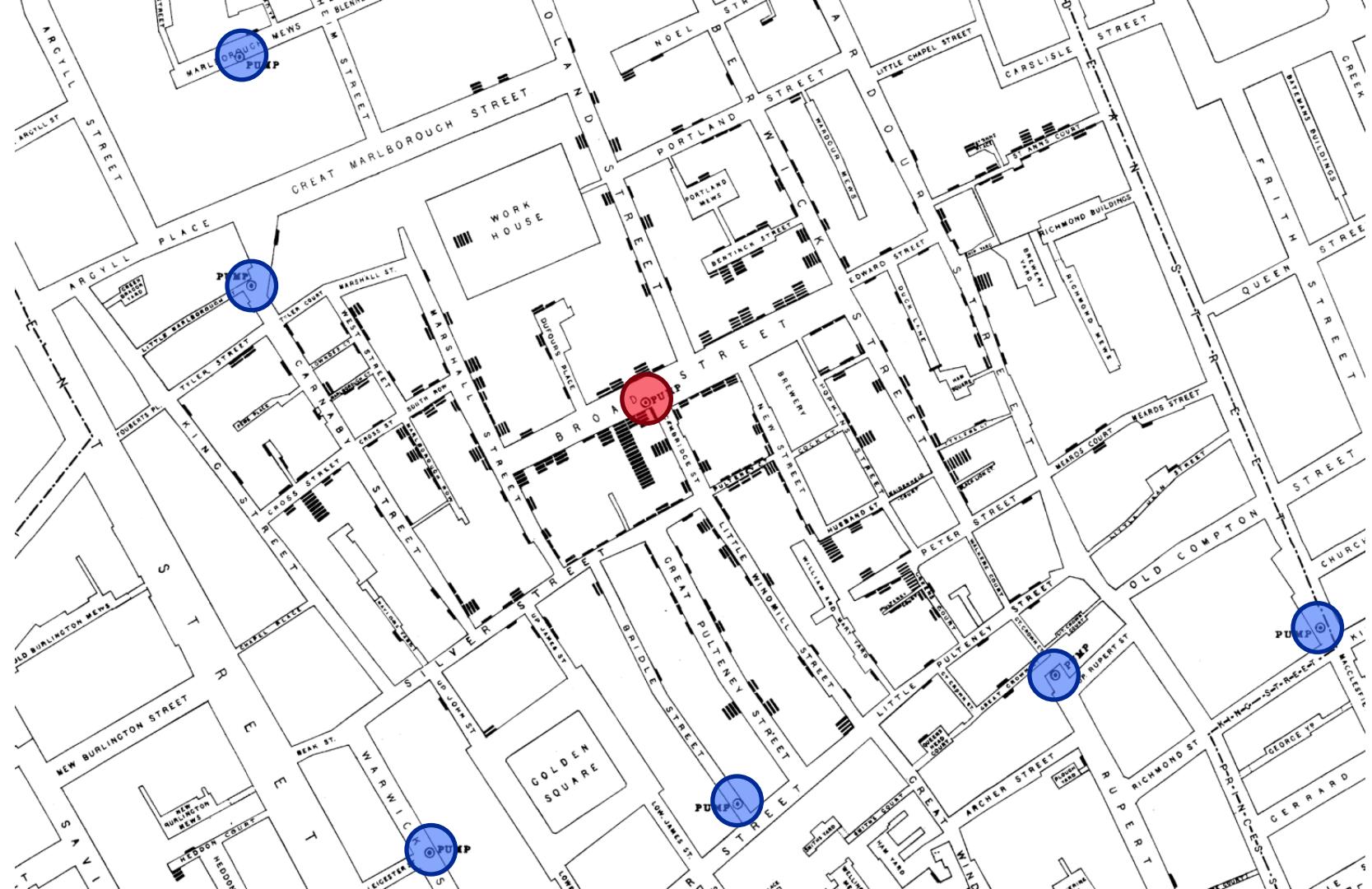
- **Miasma theory:** cholera was caused by **particles in the air**, or "miasmata", which arose from decomposing matter or other dirty organic sources.
- **Germ theory:** the principal cause of cholera was a germ cell that had not yet been identified, but was **transmitted through food or drink**.

# John Snow, 1813-1858

---



---





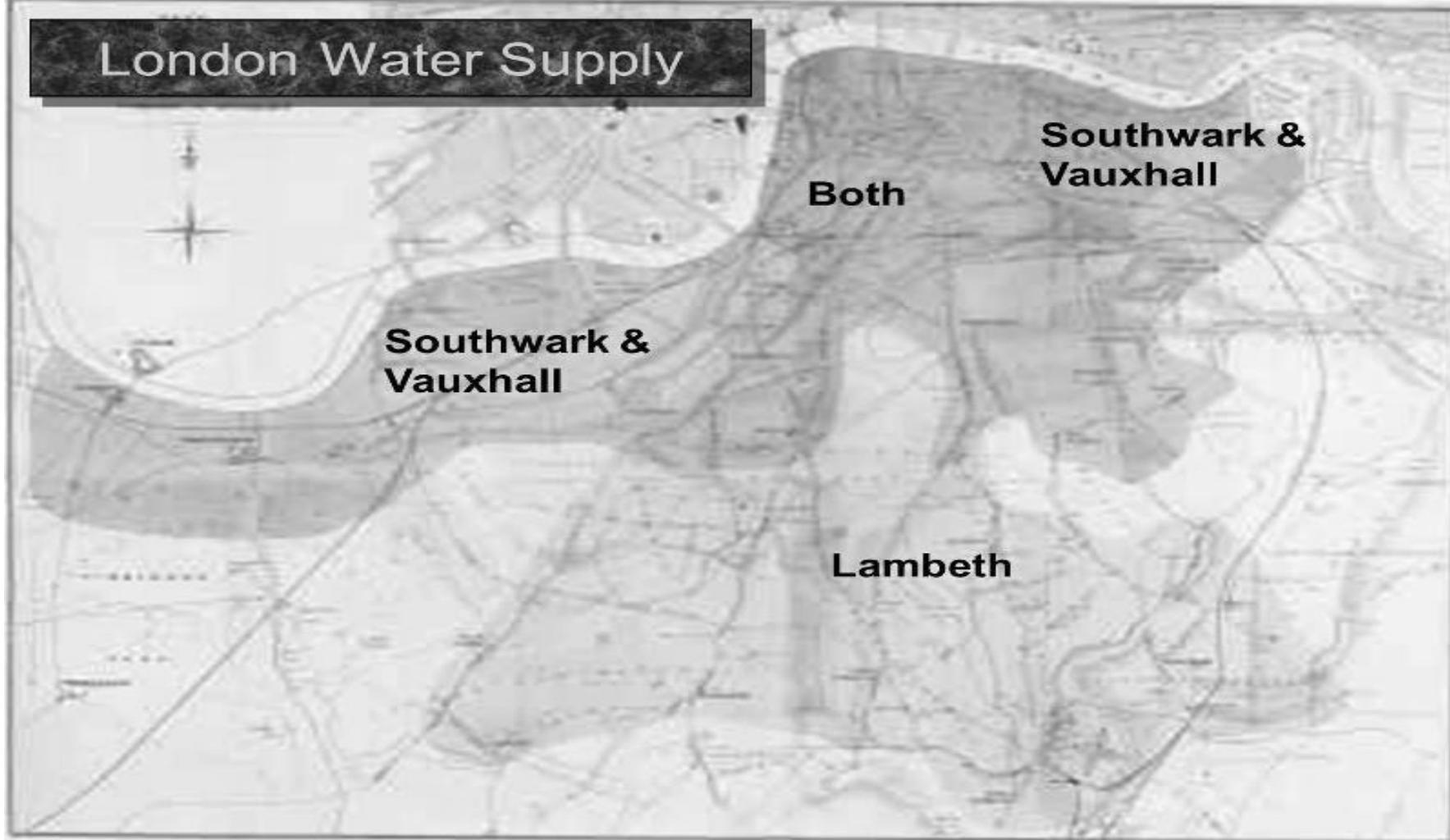
LEXINGTON  
STREET

JOHN SNOW

JOHN SNOW



# London Water Supply



# Comparison

---

- **Treatment group:** Do receive the treatment
- **Control group:** Do not receive the treatment

# Snow's “Grand Experiment”

---

“... 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 ...”

- The two groups were *similar except for the treatment.*

# Snow's table

---

Supply Area	Number of houses	Cholera deaths	Deaths per 10,000 houses
S&V	40,046	1,263	315
Lambeth	26,107	98	37

# Key to establishing causality

---

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

---

- If the treatment and control groups have **systematic differences other than the treatment**, then it might be difficult to identify causality.
- Such differences are often present in **observational studies**.
- When they lead researchers astray, they are called **confounding factors**.

# Randomization and Confounding

---

- If you assign individuals to treatment and control **at random**, then the two groups are likely to be similar apart from the treatment.
  - You can account – mathematically – for variability in the assignment.
  - **Randomized Controlled Experiments are the gold standard for establishing cause and effect.**
-

# RCEs vs Observational Studies

---

- **Question:** If randomized controlled experiments can establish causality while observational studies are subject to confounding, why are so many studies observational?