# PLSC 197 / SODA 197N: Social Data, Technology, and Artificial Intelligence

Christopher Zorn

January 20, 2026

# Wherefore "data"?

<u>Linguistics</u>:

- "Data" in common usage: *mass noun* (like "rice," "rain," "music," etc.) $\rightarrow$ singular

- "Data" in scientific usage: plural of *datum* $\rightarrow$ plural

<u>Usage</u>:

- If: "data" $\equiv$ "information": *singular*
- If:
  - · "data" $\equiv$ "facts": *plural*
  - · "data" $\equiv$ "data points": *plural*

"*Data*: A reinterpretable representation of information in a formalized manner suitable for communication, interpretation, or processing."

- Borgman (2008)

# Data isn't...(only) numbers

Basic data types:

- Numeric data:
  - · Integers
  - · Real numbers
  - · Floating points, etc.
- Characters
- Booleans (true / false)
- Void / Null (empty)

Derived/user-defined data types:

- Functions
- Strings
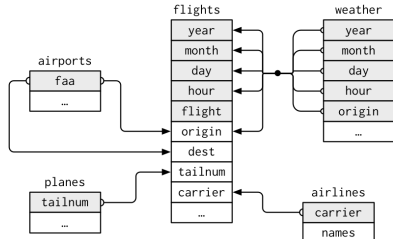- Lists
- Dates
- etc.

# Data isn't… rectangular

# Data isn't...(only) information

{celery, enchant, borscht, symmetry, jovial, daybreak}

{daisy, refuge, imagine, armadillo, futon, frolic}

"Data":

- ...requires *context*
- ...requires *purpose*
- ...implies *value*

# Data isn't...*natural*

"Data is the new oil."

- mathematician Clive Humby, 2006

Data:

- ...requires refining to be useful, *but*

- ...it's also renewable, non-rivalrous, non-excludable, non-fungible, etc. Moreover,

- ...data is a *specific asset*, which makes its value highly variable. And finally,

- ...it *doesn't exist at all without human intervention*

# Data Across (Academic) Disciplines

Adapted from Borgman (2008):

|                      | P&NS*             | S&BS**     | A&H***      |
|----------------------|-------------------|------------|-------------|
| Centrality           | Highest           | High       | Moderate    |
| Diversity            | High              | Higher     | Highest     |
| Structure[†]         | High              | Varies     | Varies/Low  |
| Data Generation      | Common            | Often      | Rarely      |
| Infrastructure       | Highly Developed  | Developed  | Developing  |

\* Physical and natural sciences.

\*\* Social and behavioral sciences.

\*\*\* Arts and humanities.

[†] More on this next week...

# Jones (2024): "AI is running out of data…"



**RUNNING OUT OF DATA**
The amount of text data used to train large language models (LLMs) is rapidly approaching a crisis point. An estimate suggests that, by 2028, developers will be using data sets that match the amount of text that is available on the Internet.

— Amount of available text on the Internet — Size of training data sets for LLMs
• Individual LLMs

*One token is about 0.8 words. †Technology Innovation Institute, Abu Dhabi.

©nature