# Data as a Resource

David Bamman
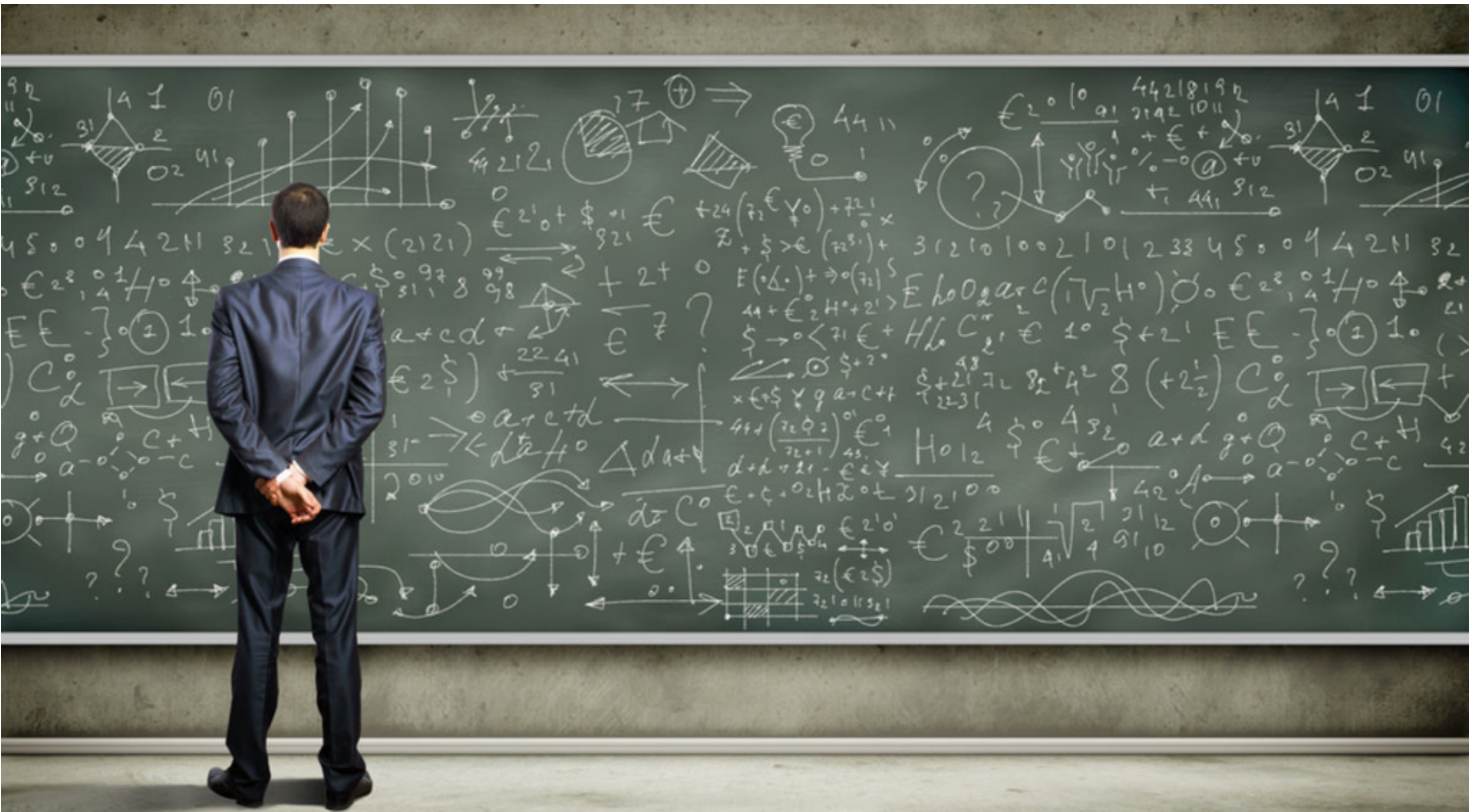Info 202: Information Organization and Retrieval

September 19, 2016
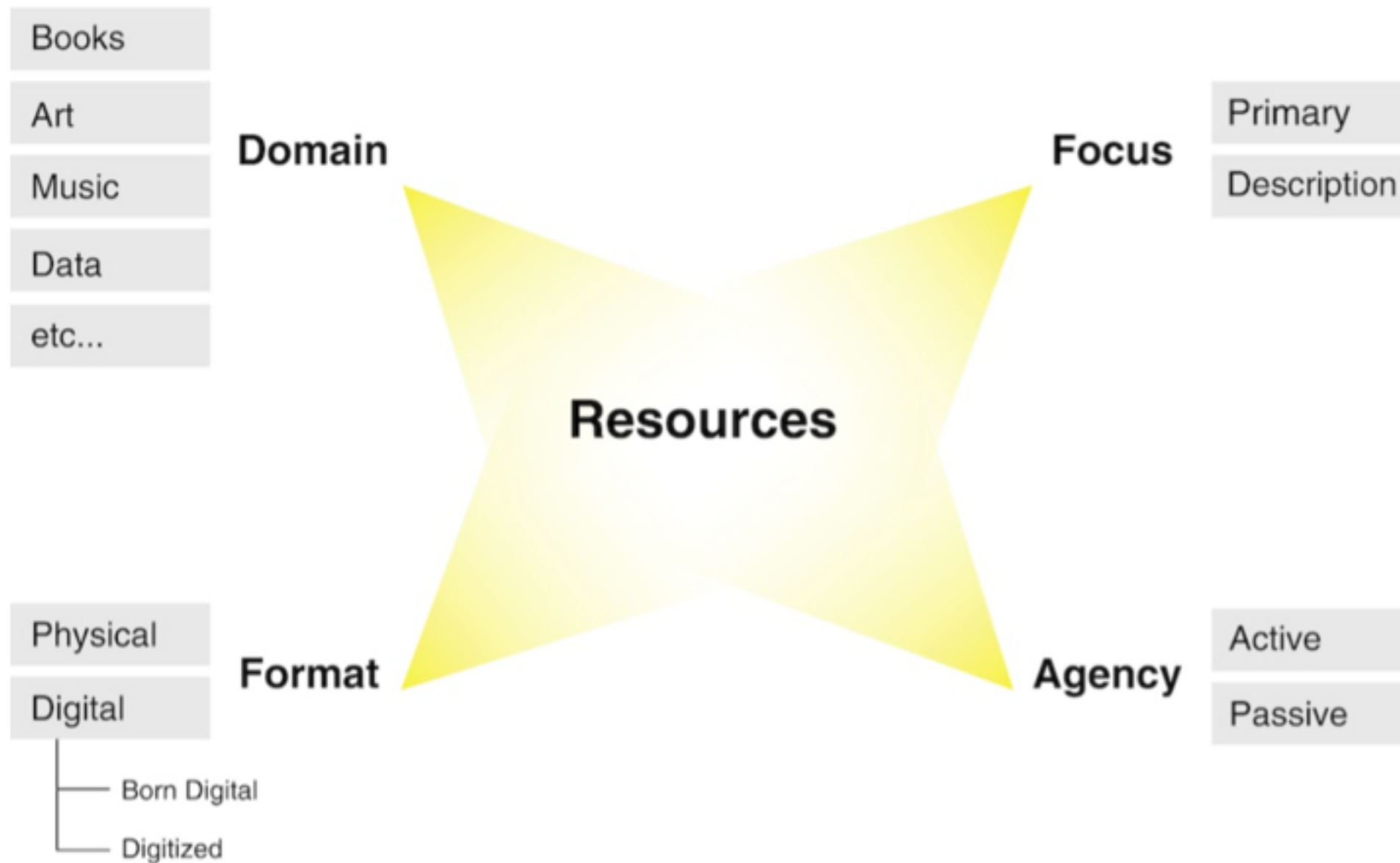
# Data as resource

- Data as a resource to be organized

  - data analysis
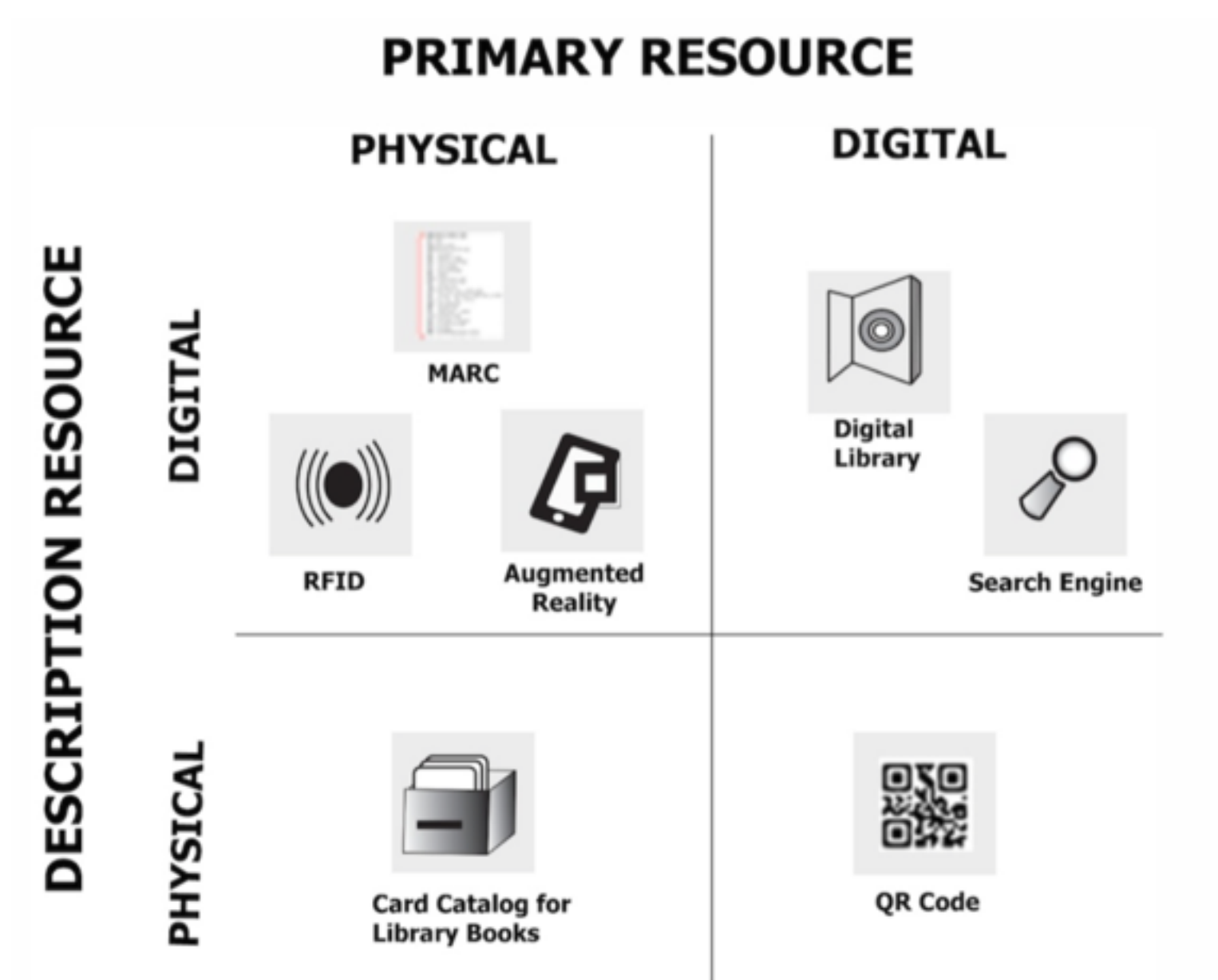
  - sensor data streams

  - internet of things

# "data" science

# Data as resource

# Format x Focus

# Data as a resource

- In many analyses, the data you have is not what you actually care about; it's the population they represent.

- What you organize is your knowledge about the population.

# Identity

- What is the unit of analysis?

- Population vs. sample

# Identity



Average age of Shamu?

# Population

# Population



population vs. hypothetical population

# Sample

- You can't always measure an entire population

  - expensive
  - can be destructive (measure how long light bulbs last)
  - impossible to measure (future users)

# How Big is Enough?

# How Big is Enough?



Average: 272 pages
Standard deviation: 252 pages

# Standard Deviation

A measure of the dispersion of a set of data points

| $x_i$ | $(x_i-272)^2$ |
|---|---|
| 3 | 72361 |
| 89 | 33489 |
| 273 | 1 |
| 501 | 52441 |
| 494 | 49284 |

Average: 41515.2

Variance is the average squared distance from the mean (41515.2)
Standard deviation = square root of the variance (203.8)

From the true population, let's take
samples of size n (= 5) and measure
the average of those samples
to see how much they vary

| | | | | | Average |
|---|---|---|---|---|---|
| 72 | 6 | 78 | 192 | 326 | 134.8 |
| 458 | 314 | 12 | 336 | 20 | 228.0 |
| 44 | 134 | 64 | 28 | 934 | 240.8 |
| 216 | 544 | 296 | 278 | 215 | 309.8 |
| 21 | 206 | 234 | 1024 | 330 | 363.0 |

n=5, standard deviation of samples = 111

n=10, standard deviation of samples = 79

n=100, standard deviation of samples = 25

n=1000, standard deviation of samples = 8

n=10000, standard deviation of samples = 2.5

# Standard error

- The standard deviation of the sample is known as the standard error.

$$se = \frac{\sigma}{\sqrt{n}}$$

# Margin of error

Under the assumption the sampling distribution is normally distributed (via the CLT), the standard error gives you confidence intervals for your measurement

| | | |
|---|---|---|
| 90% | measurement | $\pm 1.65$ x se |
| 95% | measurement | $\pm 1.96$ x se |
| 99% | measurement | $\pm 2.58$ x se |

# Bias in sampling

- Your knowledge about the data you really care about it is too uncertain.

    - Small samples

- Your knowledge about the data you have is not the same as the data you really care about.

    - Selection bias

    - Response bias

# Sampling bias

- Non-random process by which data points are selected to be in the sample and others are not

- Canvassers conducting in-person interviews on voting preferences, avoiding:

  - all houses with pit bulls chained out front

  ⇒ only "nice" looking houses are canvassed

Watt et al., "Populations and Samples," p. 52

# Sampling bias



Lascaux cave painting

# (Non-)Response bias

- Non-random process by which data points participate in the sample and others are not

- Survey companies about their organizational culture, only take measurements from those that let their employees respond

Watt et al., "Populations and Samples," p. 53

# (Non-)Response bias

1936 poll of presidential election (Alf Landon v. FDR)

# Example

- How many children are in your family (including you)?

# Organizing data for analysis

- What considerations guide our choice for selecting data? What's the granularity of our unit of analysis?

  - Predict opening box office for a movie

# Organizing data for analysis

- What considerations guide our choice for selecting data?  What's the granularity of our unit of analysis?

  - Automatically identifying plagiarism

# Sensor Data Streams

- Sensors for the analysis of human behavior

- Instrument a person

# Sensor Data Streams

- Sensors for the analysis of human behavior

- Instrument a person

# Sensor Data Streams

# Sensor Data Streams

- Instrument groups





Choudhury and Peatland, "The Sociometer"

# Sensor Data Streams

- Instrument a space

# Sensor Data Streams

- Instrument a space

- Deb Roy (MIT)

# Sensor Data Streams

- Where do people travel in a day?

- Who do they communicate with?

- What tools do they use during the day?

- What routines define a "typical" day?

- How healthy are their behaviors?

Voida et al. (2014), "Sensor Data Streams", p. 293

# Sensor Data Streams

- Formulate research question

- Get/build sensors

- Determine how frequently to collection samples

- Install sensors

- Store data

- Sense making

Voida et al. (2014), "Sensor Data Streams", p. 301

# Data as resource

- What information do we actually get from sensors?

# Granularity

→
- Heartbeat
- Sleep patterns
- Health

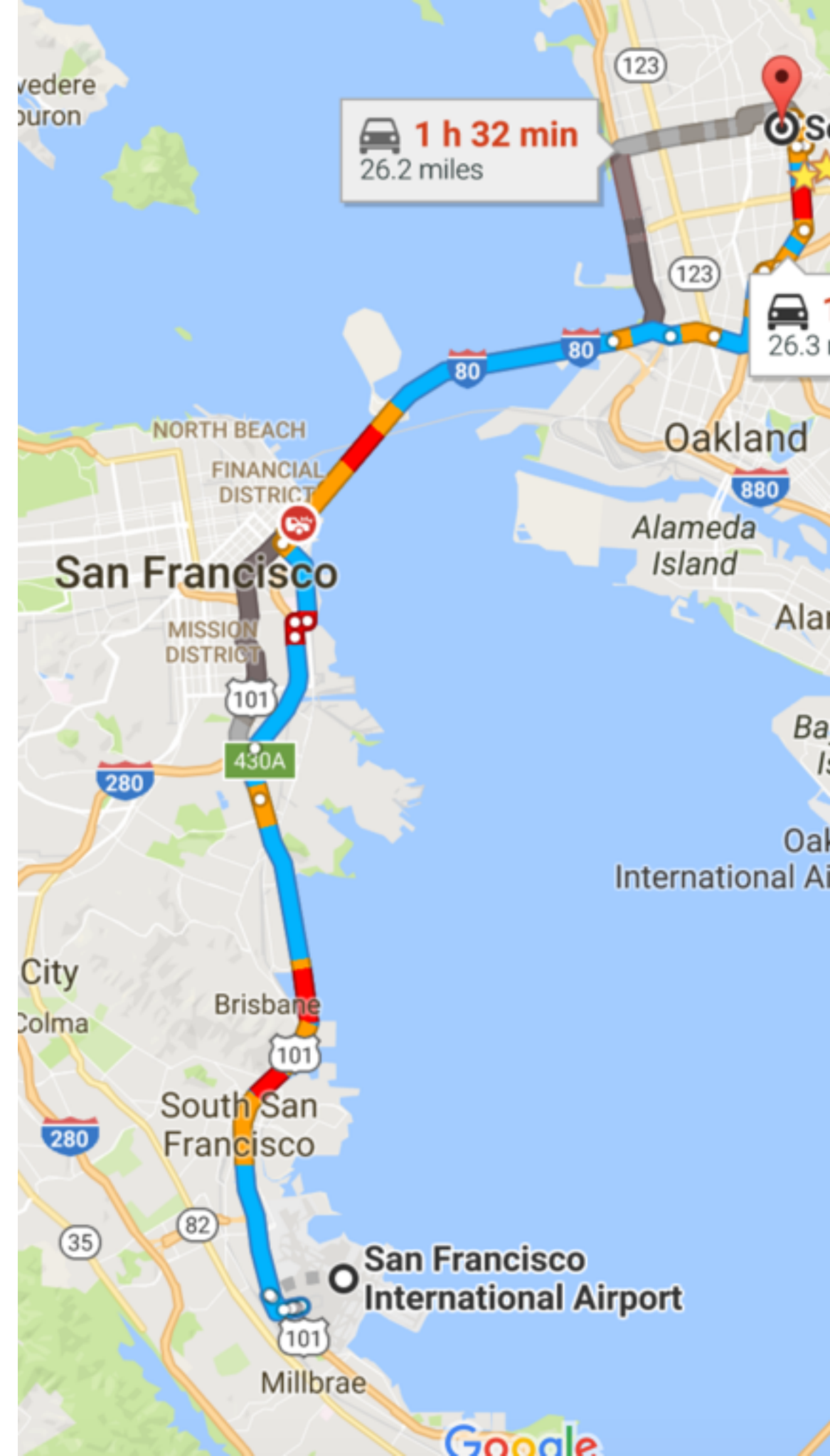→
- Location
- Most common transit routes
- Health

# Sensor data in an organizing system

- **what** is being organized?
- **why** is it being organized?
- **how much** is it being organized?
- **when** is it being organized?
- **how** (or by whom) is it being organized?
- **where** is it being organized?

# Granularity

Google Maps

- Low-level sensor data
- Route + timing
- Prediction of ETA

# Granularity

# Sensor data in an organizing system

- what is being organized?

- why is it being organized?

- how much is it being organized?

- when is it being organized?

- how (or by whom) is it being organized?

- where is it being organized?

# Midterm

- Answer 4 out of 6 questions

- Covers material through today

- You have 90 minutes (contiguous)

- Due Friday 9/23