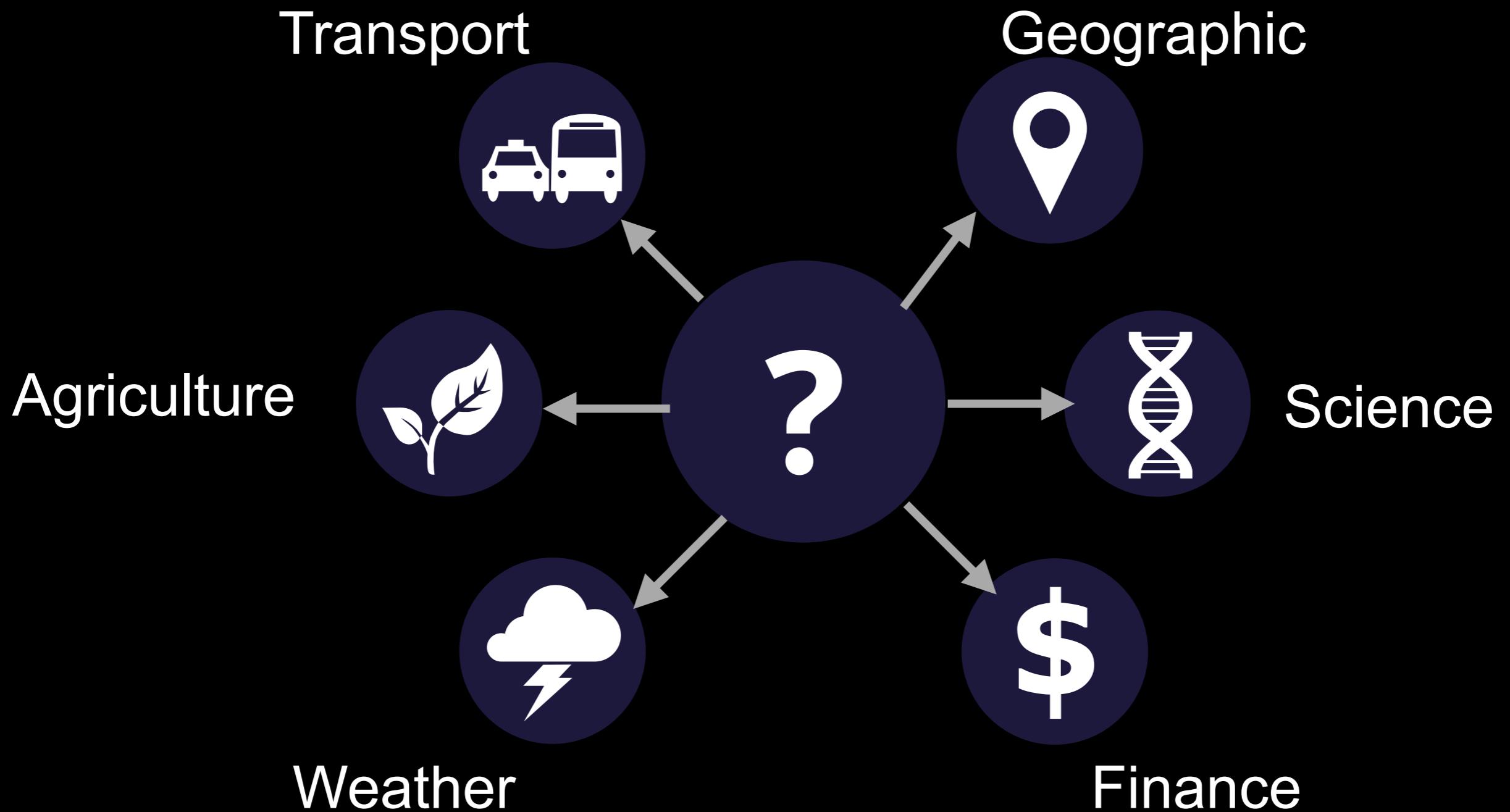


Foundations of Data Analysis

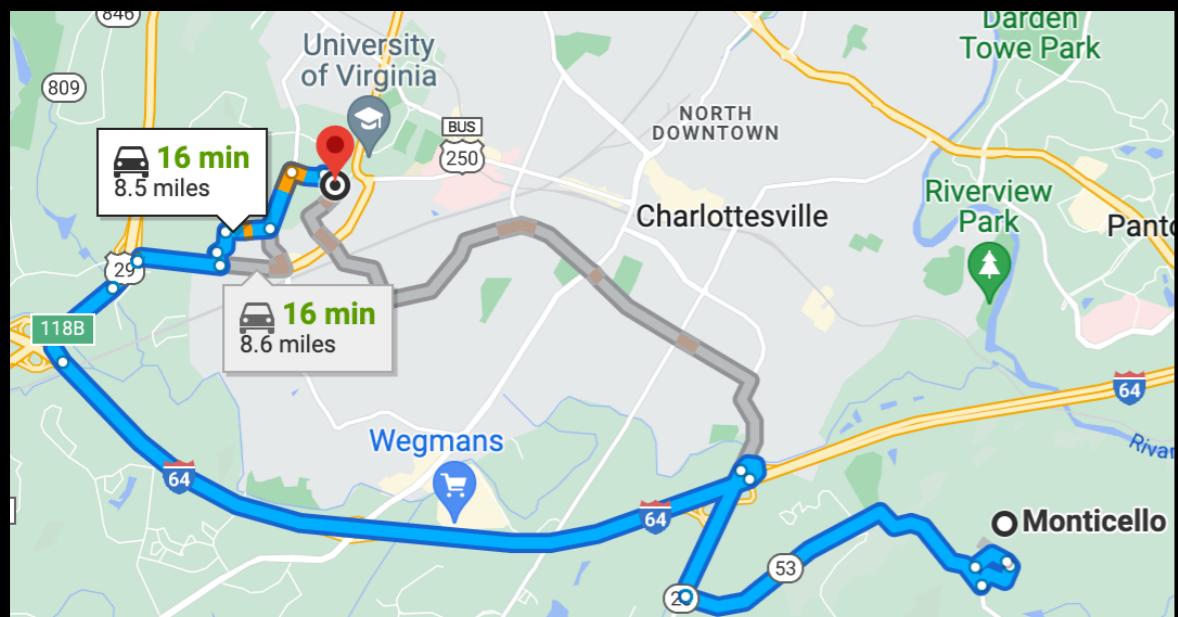
Lecture 1: Introduction

Foundations of Data Analysis



Foundations of Data Analysis

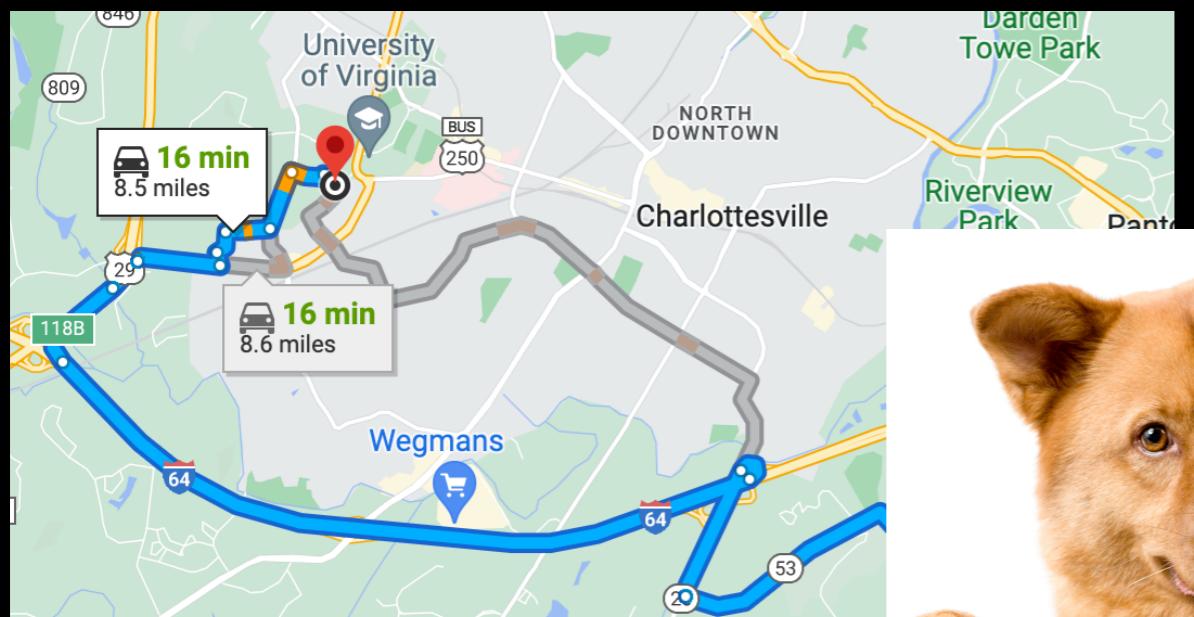
Analyzing data to answer questions.



Optimal path?

Foundations of Data Analysis

Analyzing data to answer questions.



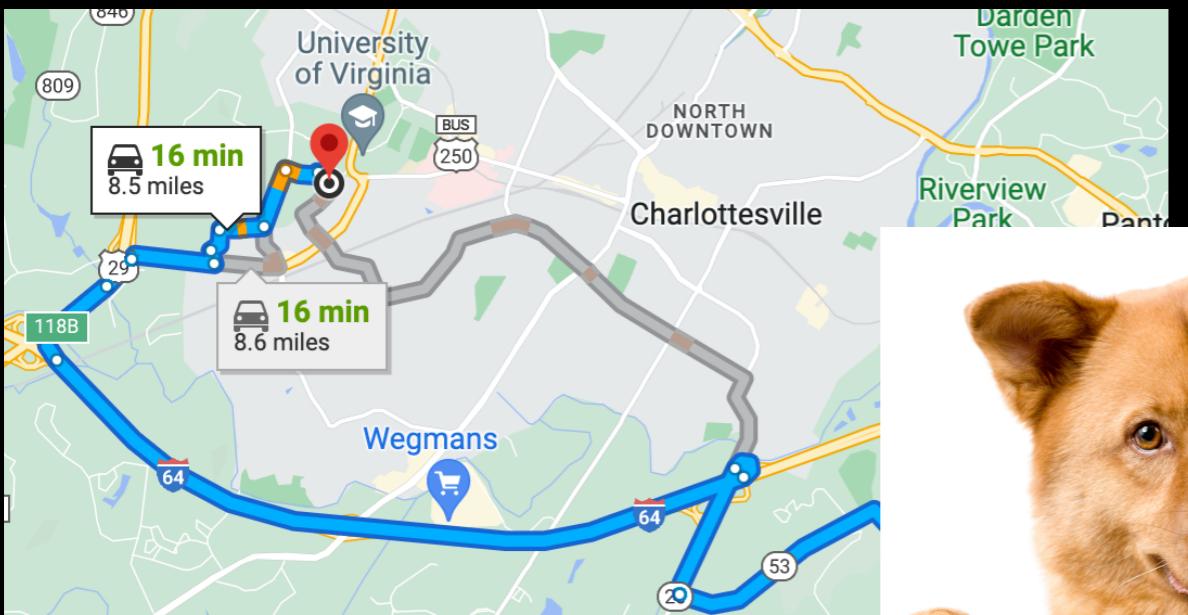
Optimal path?



Dog or cat?

Foundations of Data Analysis

Analyzing data to answer questions.



Dog or cat?



Voices?

Foundations of Data Analysis

Learning basic skills to

- Describe and visualize the data (software/tools)
- Model data (linear algebra, calculus, probability....)
- Manipulate, interpret, and process data
(programming)

Course Purpose

An introduction to the foundations behind modern data analysis and machine learning.

Prerequisites

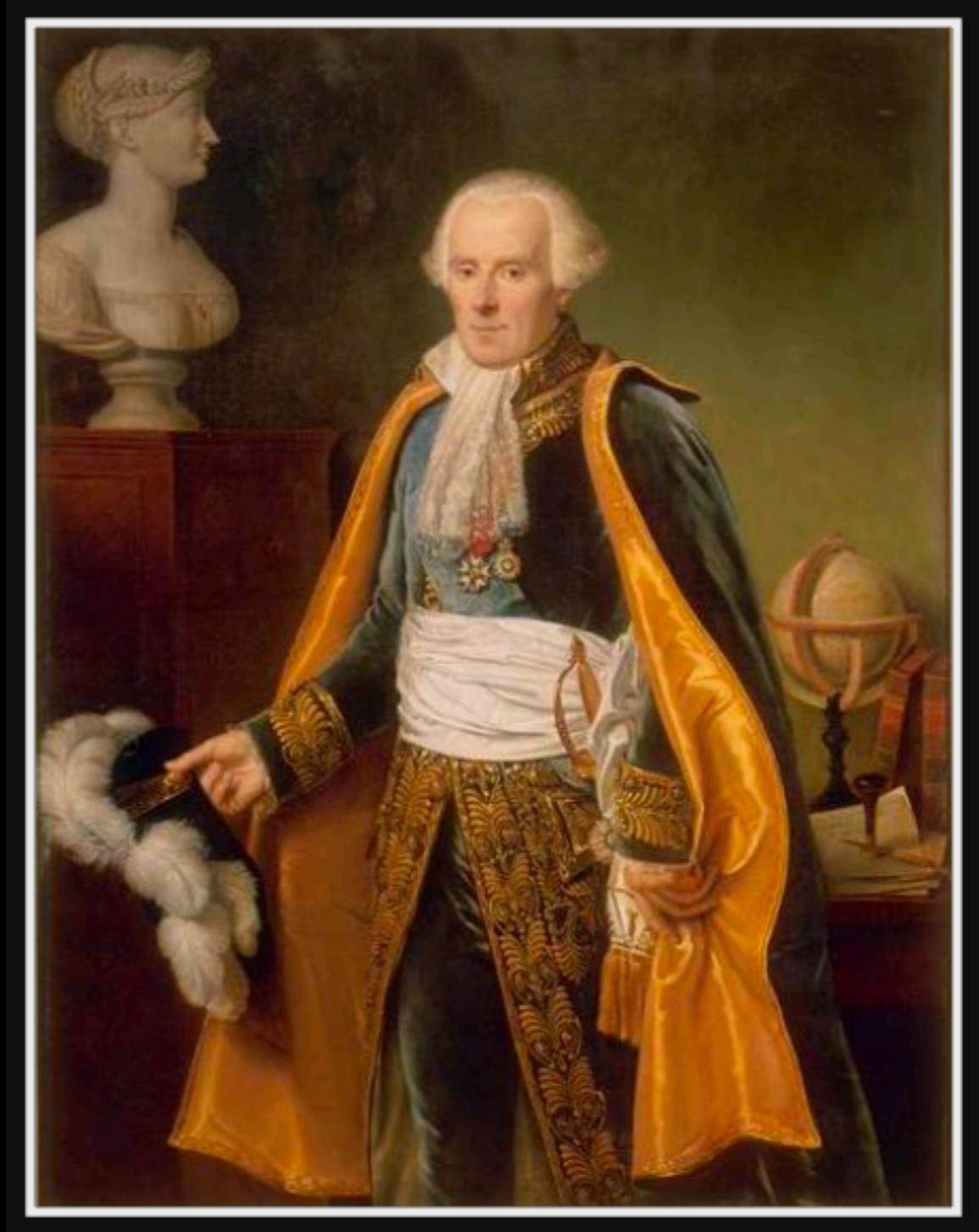
- Programming in Python (CS 2110 or equivalent)

Course Requirements

- Four homework assignments.
- Midterm exam (IN CLASS with one-page A4 cheat sheet).
- Final (Take home)
- NO bonus projects / assignments

Course Information

All course information will be distributed online:
<https://mz8rr.github.io/FoDA/>



Pierre-Simon Laplace (1749-1827)

Births in Paris

1745 - 1770

Births in Paris

1745 - 1770

251,527 Boys

241,945 Girls

Are males born at a higher rate than females?

251,527 Boys

241,945 Girls

Some possible descriptive statistics

251,527 Boys

241,945 Girls

Some possible descriptive statistics

251,527 Boys

241,945 Girls

- **Difference:** +9582 Boys

Some possible descriptive statistics

251,527 Boys

241,945 Girls

- **Difference:** +9582 Boys
- **Ratio:** 104 Boys to 100 Girls

Some possible descriptive statistics

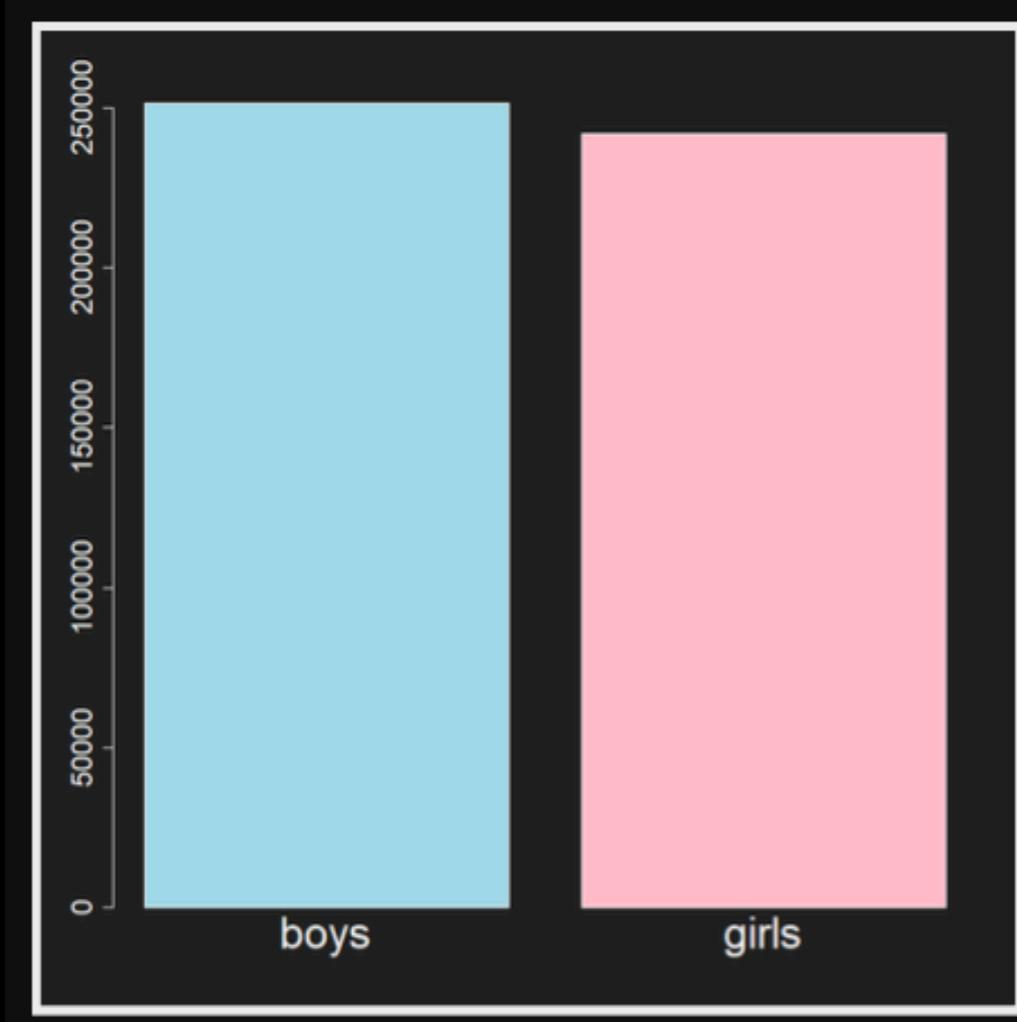
251,527 Boys

241,945 Girls

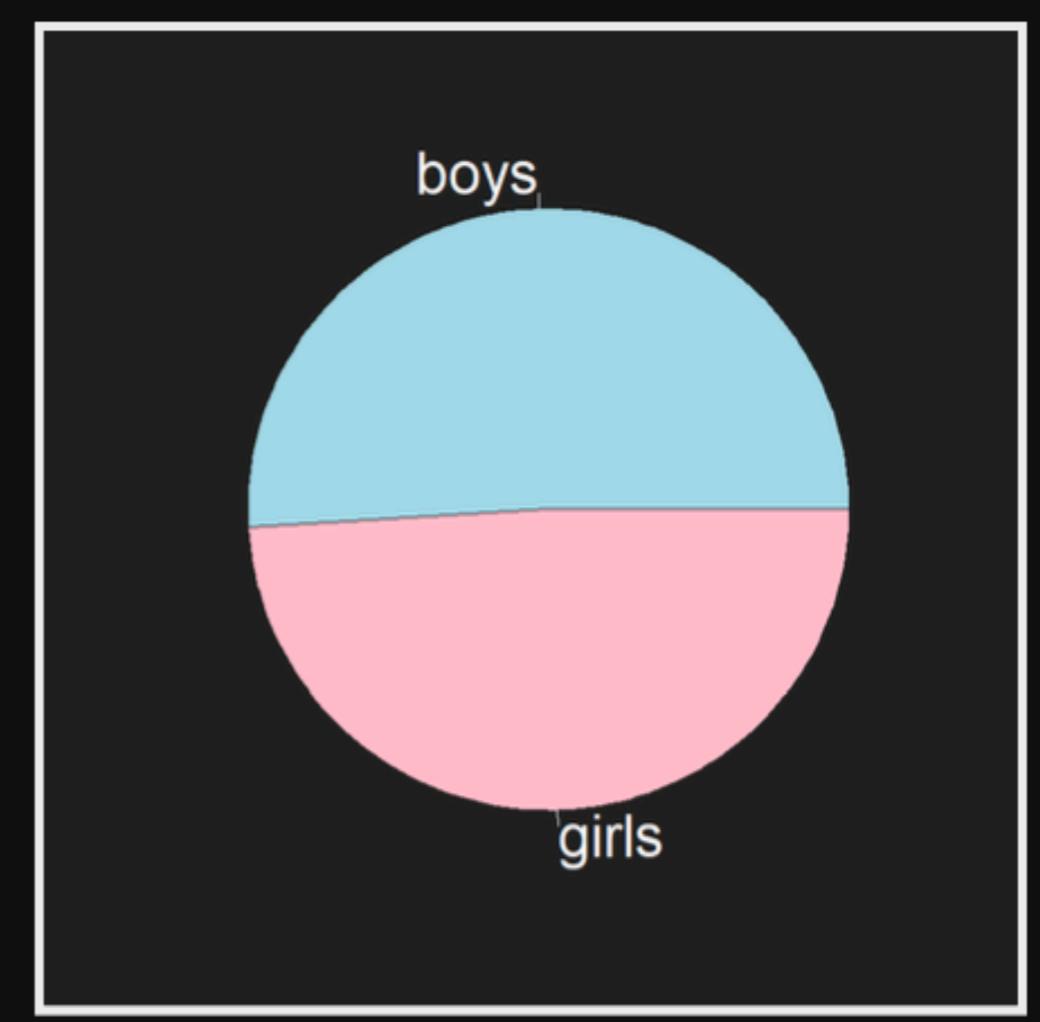
- **Difference:** +9582 Boys
- **Ratio:** 104 Boys to 100 Girls
- **Proportion:** 50.97% Boys

Some possible visualizations

251,527 Boys



241,945 Girls



How did Laplace solve this?

Answer?

How did Laplace solve this?

Conditional Probability that:

rate of boys, θ , is greater than girls,

given

observed data

What is probability ?

What is probability ?

Definition: *Probability* is the study of the mathematical rules that govern random events.

But what is randomness?

But what is randomness?

Informally, a *random event* is an event where we do not know the outcome without observing it.

But what is randomness?

Informally, a *random event* is an event where we do not know the outcome without observing it.

Probability tells us what we can say about such events, given our assumptions about the possible outcomes.

What is mathematical statistics?

What is statistics?

Definition: *Statistics* is the application of probability to the collection, analysis, and description of random data.

Statistics is used to:

Statistics is used to:

- **Design experiments**

Statistics is used to:

- **Design** experiments
- **Summarize** data

Statistics is used to:

- **Design experiments**
- **Summarize data**
- **Make conclusions about the world**

Statistics is used to:

- **Design** experiments
- **Summarize** data
- **Make conclusions** about the world
- **Explore** complex data

Machine learning

Statistics

Probability

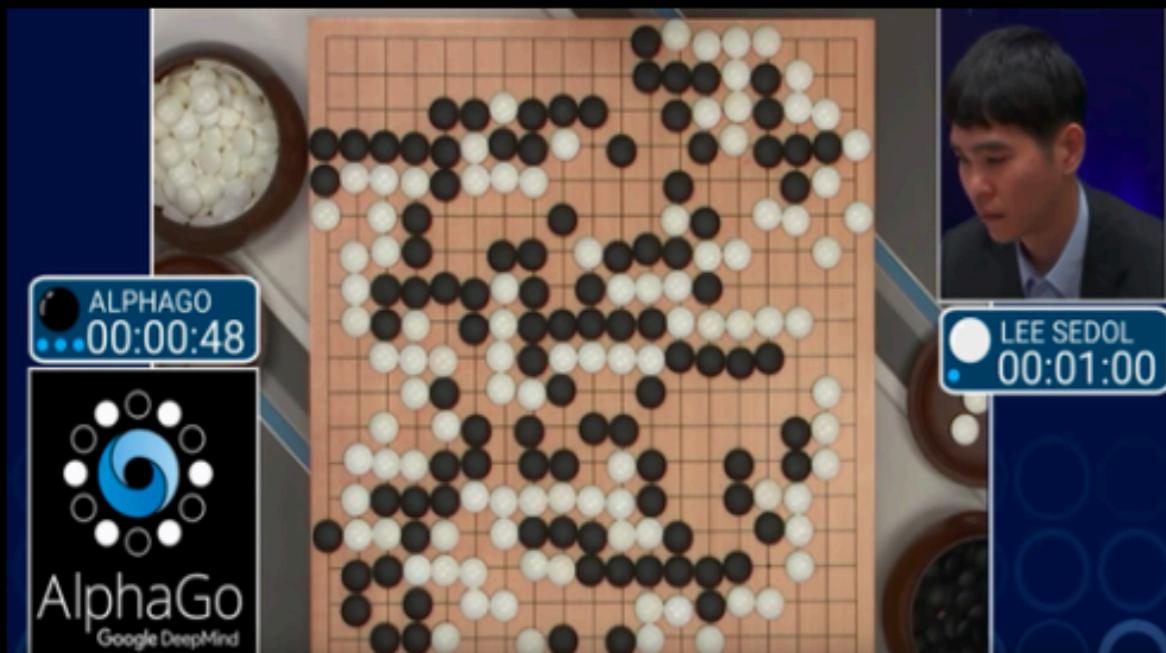
Linear Algebra

What is machine learning?

What is machine learning?

Definition: *Machine learning* builds statistical models of data in order to recognize complex patterns and to make decisions based on these observations.

Machine learning is Everywhere?



Games



Recommendation system



Assisted driving



Cancer diagnosis

Levels of data analysis expertise

Levels of data analysis expertise

- 0: what is data?

Levels of data analysis expertise

- 0: what is data?
- 1: I know how to run data analysis software

Levels of data analysis expertise

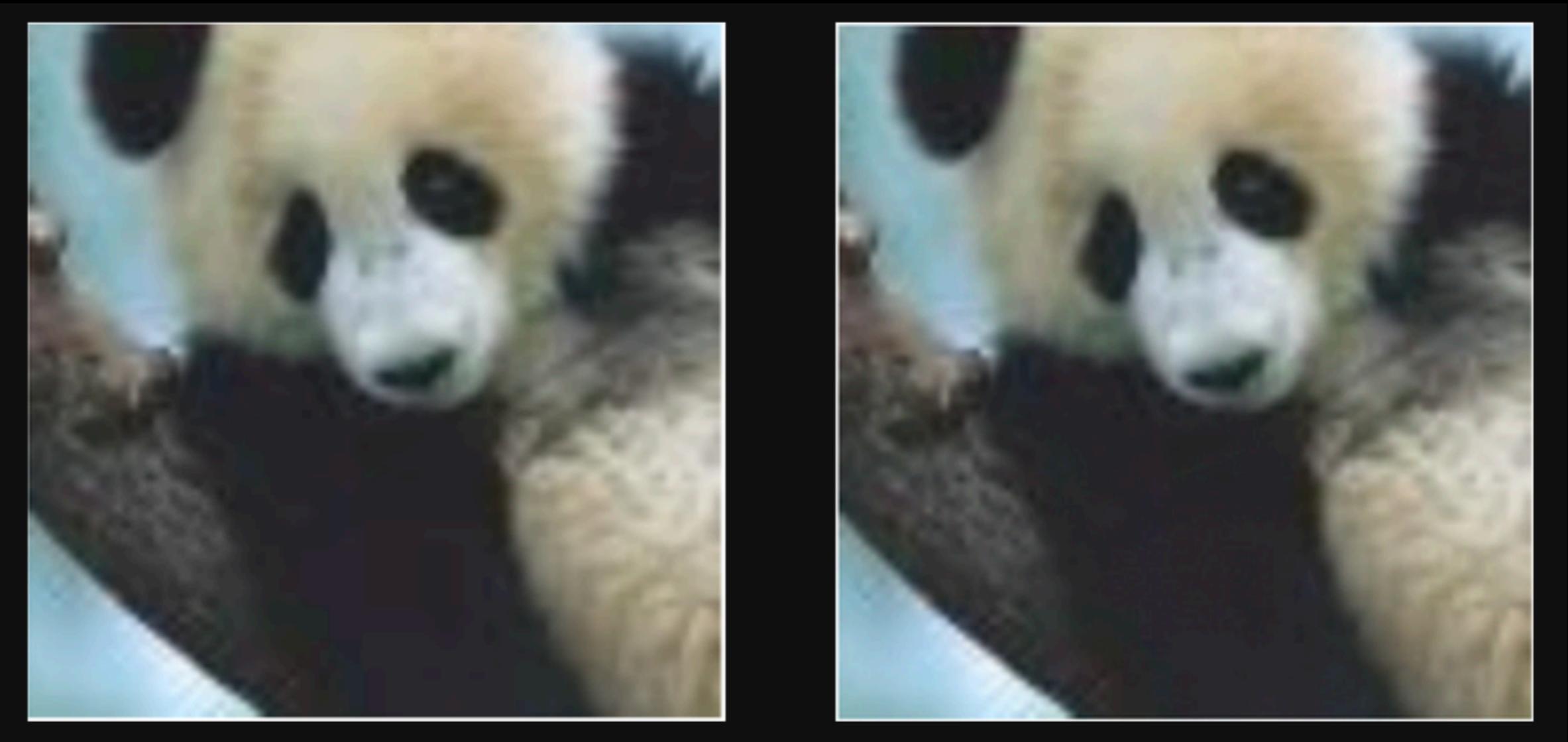
- 0: what is data?
- 1: I know how to run data analysis software
- 2: I understand the math behind the analysis

Levels of data analysis expertise

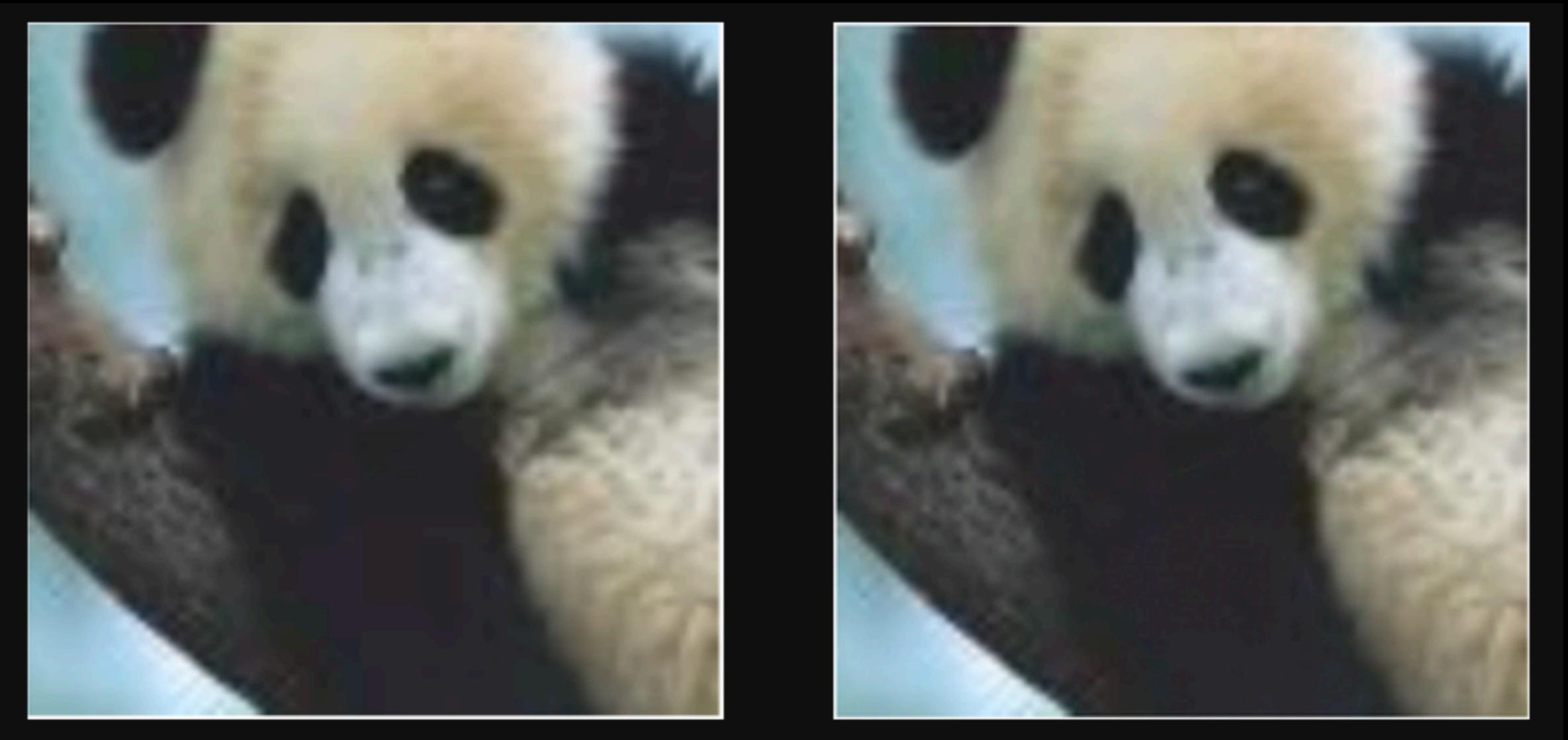
- 0: what is data?
- 1: I know how to run data analysis software
- 2: I understand the math behind the analysis
- 3: I'm able to invent new data analysis methods

**Why should you know the
mathematical foundations?**

When machine learning goes wrong

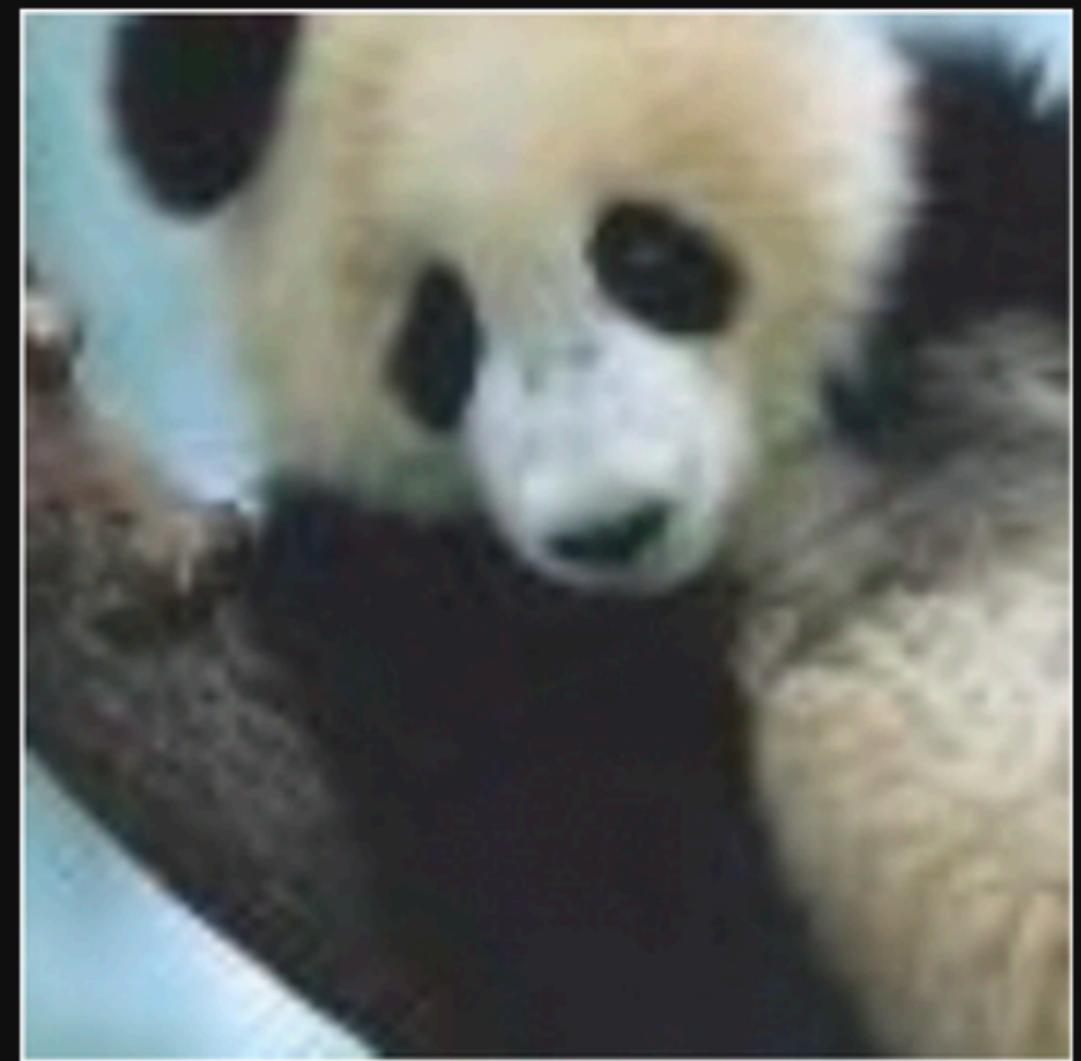
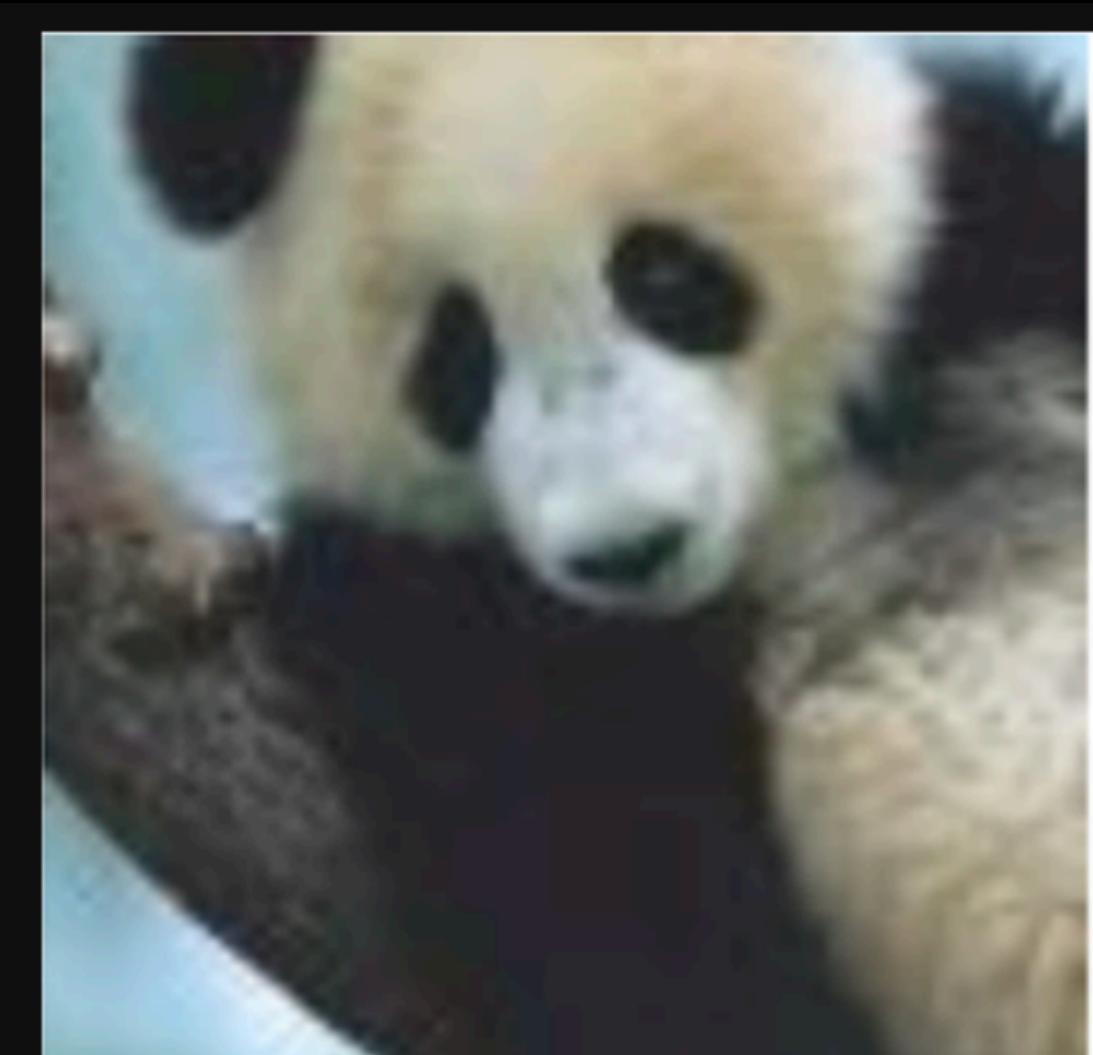


When machine learning goes wrong



Panda (57.7% confidence)

When machine learning goes wrong



Panda (57.7% confidence)

Gibbon (99.3% confidence)