# Data Visualization

- *Exploratory Data Analysis*
- *Graphing in Python*

James Balamuta

INMAS Statistical Methods Workshop Fall 2021

# Lecture Objectives

- **Apply** quantitative and visual exploratory techniques on data

- **Deduce** and **interpret patterns** revealed during exploratory data analysis (EDA)

- **Create** common statistical graphs using Seaborn.

# Exploratory Data Analysis

**Definition**:

*Exploratory Data Analysis (EDA)* is a philosophy for the beginning of an analysis that describes a variety of techniques that are *quantitative* and *visual* in nature to look for patterns in data.
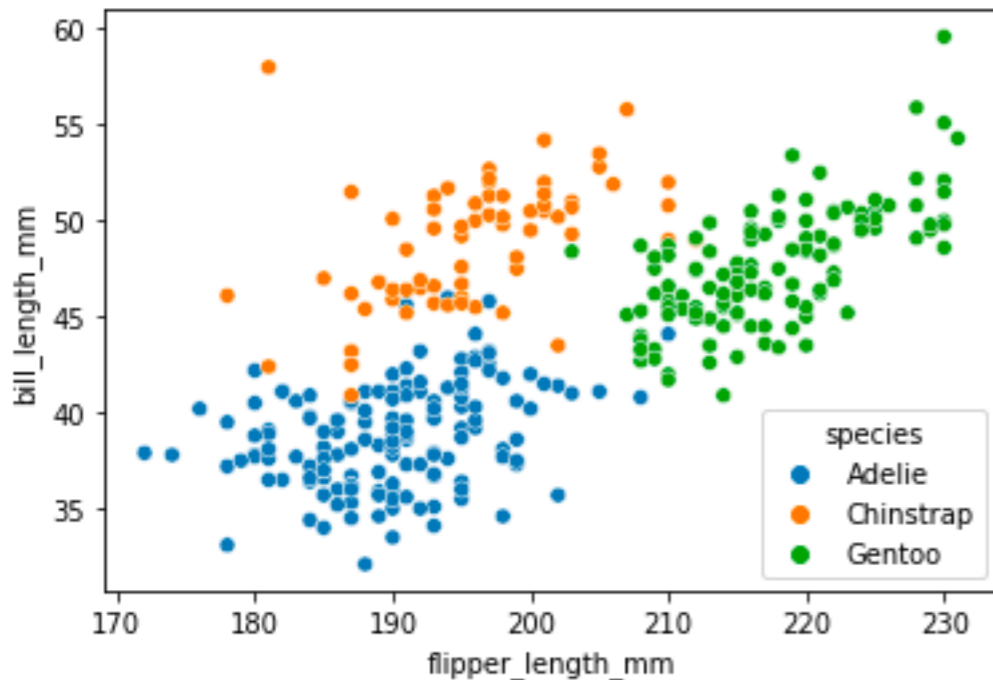
# Types of EDA

## Quantitative vs. Visual

|       | bill_length_mm | bill_depth_mm | flipper_length_mm | body_mass_g |
|-------|----------------|---------------|-------------------|-------------|
| count | 342.000000     | 342.000000    | 342.000000        | 342.000000  |
| mean  | 43.921930      | 17.151170     | 200.915205        | 4201.754386 |
| std   | 5.459584       | 1.974793      | 14.061714         | 801.954536  |
| min   | 32.100000      | 13.100000     | 172.000000        | 2700.000000 |
| 25%   | 39.225000      | 15.600000     | 190.000000        | 3550.000000 |
| 50%   | 44.450000      | 17.300000     | 197.000000        | 4050.000000 |
| 75%   | 48.500000      | 18.700000     | 213.000000        | 4750.000000 |
| max   | 59.600000      | 21.500000     | 231.000000        | 6300.000000 |

```python
# Obtain summary information
penguins.describe()


# Visualize information with
# a scatterplot
sns.scatterplot(
    data=penguins,
    x="flipper_length_mm",
    y="bill_length_mm",
    hue="species");
```
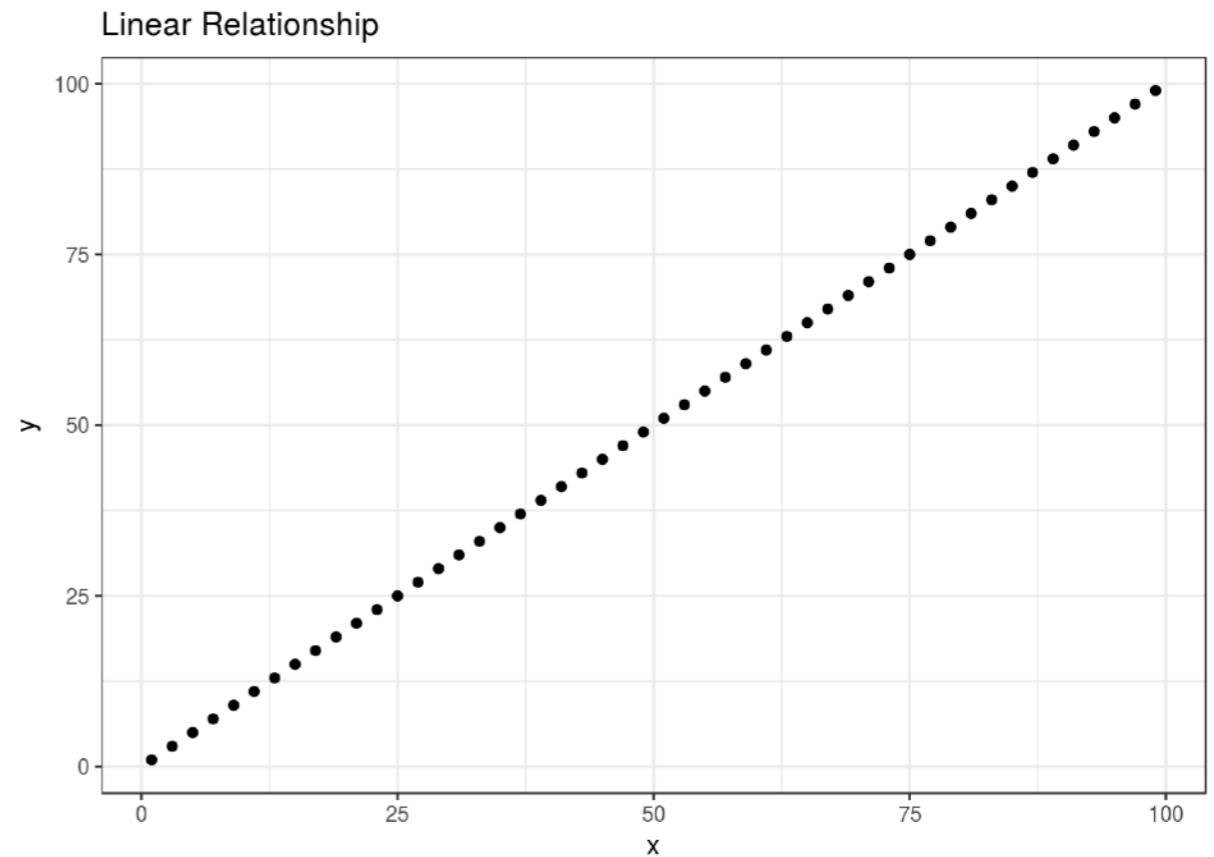
# Patterns

... detecting, analyzing, and communicating ...

1. What kind of **relationship exists** is present?
2. What's the **level of strength** of the relationship?
3. Are there any **confounding variables** that might be behind it?
4. How does the **pattern exist in subsets** of the data?



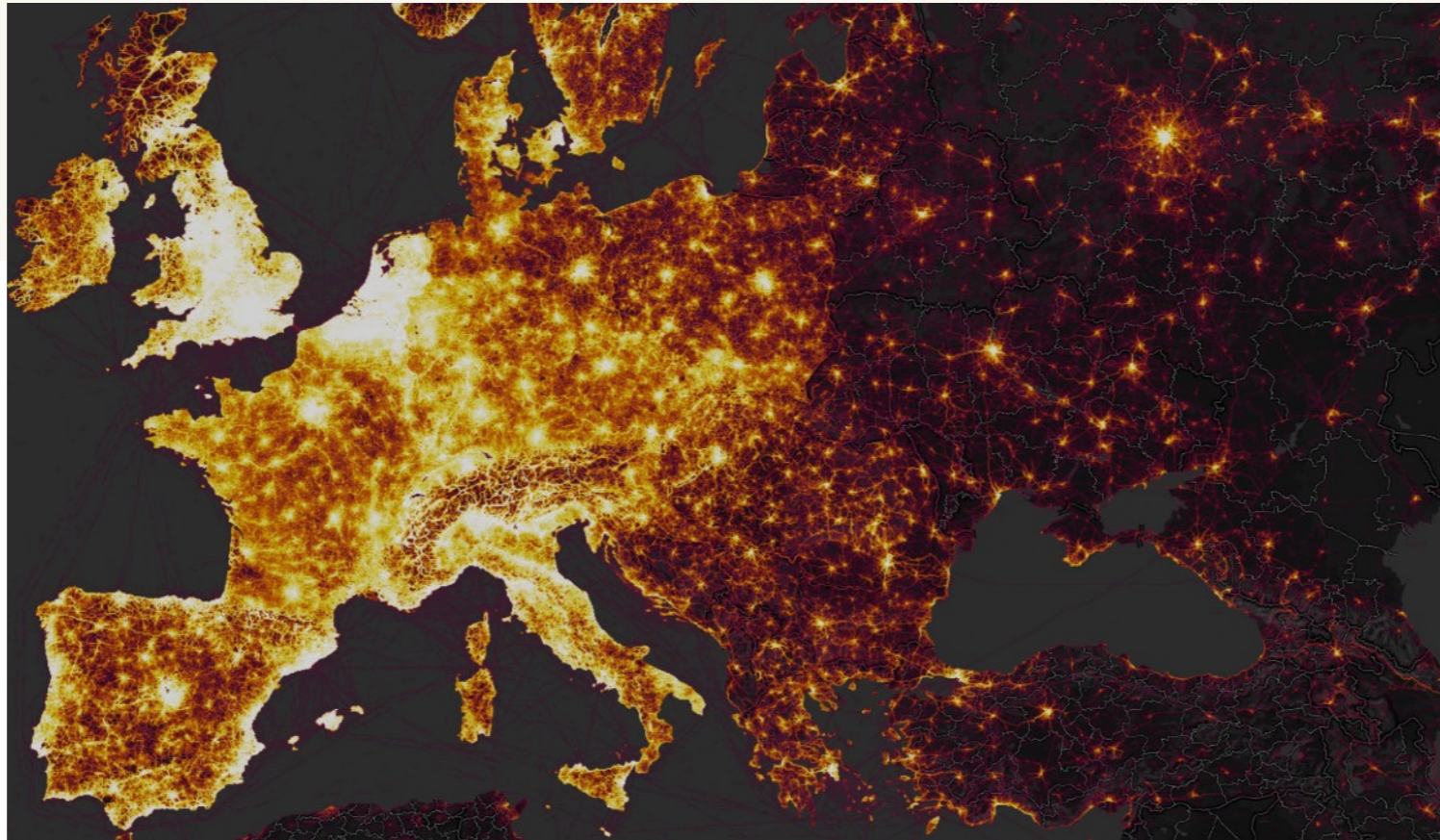Linear Relationship



Fuzzed Linear Relationship

**Definition**:
*Variation* is the difference between observations in **one variable**.

**Definition**:
*Covariation* is the difference among observations between **two or more** variables.

"…make **both** calculations **and** graphs. Both sorts of output should be studied; each will contribute to understanding."

– F. J. Anscombe, 1973

"The greatest value of a picture is when it forces us to notice what we never expected to see."

−John Turkey in Exploratory Data Analysis (1977)

# Datasaurus

## … the new anscombe's quartet …



X Mean: 54.26`59224`
Y Mean: 47.83`13999`
X SD  : 16.76`49829`
Y SD  : 26.93`42120`
Corr. : -0.06`42526`

# Pattern

## ... the new anscombe's quartet ...

# Graphing in Python

# Python Graphing Systems
## Comparison of Graphs on similar data



**matplotlib**

**seaborn**

**Bokeh**

**plotnine**

# Goldilocks Scenario

... what graph system to choose ???



**matplotlib**

Highly customizable

**seaborn**
**bokeh**

**plotnine**

# Grammar of a Graph

## Underlying Structure of Graphs



Data + Geometric Object + Coordinate System = Plot

Anatomy of a figure

# seaborn: statistical data visualization



Seaborn is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.

For a brief introduction to the ideas behind the library, you can read the introductory notes. Visit the installation page to see how you can download the package and get started with it. You can browse the example gallery to see what you can do with seaborn, and then check out the tutorial and API reference to find out how.

To see the code or report a bug, please visit the GitHub repository. General support questions are most at home on stackoverflow or discourse, which have dedicated channels for seaborn.

## Contents

- Introduction
- Release notes
- Installing
- Example gallery
- Tutorial
- API reference

## Features

- Relational: API | Tutorial
- Distribution: API | Tutorial
- Categorical: API | Tutorial
- Regression: API | Tutorial
- Multiples: API | Tutorial
- Style: API | Tutorial
- Color: API | Tutorial

https://seaborn.pydata.org/

# Graphing Outline
## … when to use a graph …

**Variables**

### One Variable (Variation)
x = <?>

**Continuous**

**Categorical or Discrete**

 Histogram

 Barplot

### Two Variables (Covariation)
x = <?> , y = <?>

Indep: Continuous
Dep: Numerical

Indep: Categorical
Dep: Numerical

Indep: Categorical
Dep: Categorical

 **Line Graph**

 **Scatterplot**

 **Boxplot**

 **Violin Graph**

 **Stacked Barplot**

 **Side-by-side Barplot**

**Indep**endent
**Dep**endent

# Histogram

## **Continuous** data for one variable

binwidth

Frequency of
Observations
that fall into
the bin

30

25

15

10

5

1   2   3   4

Continuous
variable

# Histograms
# imagined "bins" for continuous data



dataset: Geyser—272 records of delay (in seconds) between eruptions of Old Faithful

https://tinlizzie.org/histograms/

# Boxplot

**Categorical** paired with **numerical** to observe *covariation*

Interquartile range (**IQR** = Q3 - Q1)

**Minimum**
(excludes outlier)

**Maximum**
(excludes outlier)

**Outliers**
Upper:
Q3 + 1.5*IQR
Lower:
Q1 - 1.5*IQR

**1st Quantile**
(25%)

**Median**
(50%)

**3rd Quantile**
(75%)

# data-to-viz.com
## Flowchart for picking the best visualization given data

# Code Examples

# Sample Implementations
## Code to Generate Various Graphs in Python and R



Source

Source

# Summary

- Emphasized the benefits of graphing data.

- Discussed different approaches to graphing data with one variable vs. two variables.

This work is licensed under the