# Machine Learning with Python:
# A Hands-On Introduction

**https://github.com/cbrownley/2023MLWEEK_MLWITHPYTHON**

## Clinton Brownley, PhD

https://cbrownley.github.io/

https://www.linkedin.com/in/clintonbrownley/

# Agenda (too much…we'll take our time : )

8:30-8:45 – Setup and Overview

8:45-9:30 – **Data preprocessing**

9:30-10:00 – Hands-on: Data preprocessing

10:00-10:30 – **Cross-validation**

10:30-11:00 – Hands-on: Cross-validation

11:00-11:30 – Hands-on: K-fold cross-validation

11:30-12:00 – **Classification** (breast tumor diagnosis)

12:00-12:30 – Hands-on: Classification

12:30-1:00 – Hands-on: Decision Trees

1:00-1:30 – **Regression** (california housing)

1:30-2:00 – Hands-on: Regression

2:00-2:30 – Hands-on: **Shrinkage methods**

2:30-3:00 – Classification (credit card fraud)

3:00-3:30 – Regression (cycling counts)

3:30-4:00 – Hands-on: Classification (hotel bookings)

4:00-4:30 – Hands-on: **Explainable ML**

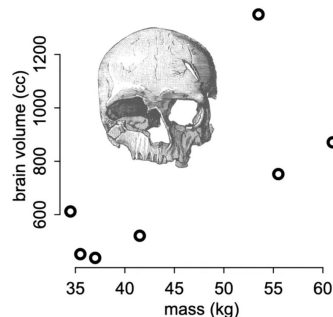# Prediction vs Causal Inference

## Problems of Prediction

What function describes these points? (fitting, compression)

What function explains these points? (causal inference)

What would happen if we changed a point's mass? (intervention)

What is the next observation from the same process? (prediction)



## Good & Bad Controls

**"Control" variable**: Variable introduced to an analysis so that a causal estimate is possible

Common **wrong** heuristics for choosing control variables
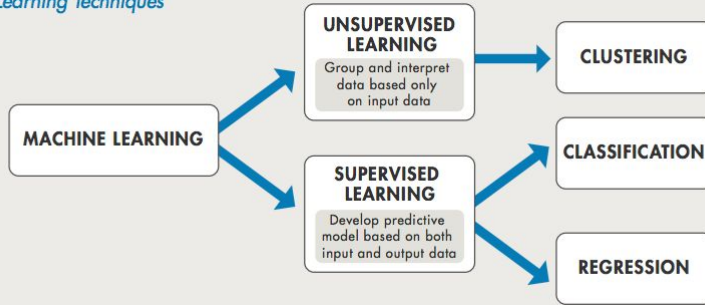
Anything in the spreadsheet **YOLO**!

Any variables not highly **collinear**

Any **pre-treatment** measurement (baseline)

# Supervised vs Unsupervised



Machine Learning Techniques

UNSUPERVISED LEARNING
Group and interpret data based only on input data

MACHINE LEARNING

SUPERVISED LEARNING
Develop predictive model based on both input and output data

CLUSTERING

CLASSIFICATION

REGRESSION

### Supervised Learning

| $X_1$ | $X_2$ | $X_3$ | $X_p$ | Y |
|-------|-------|-------|-------|---|
|       |       |       |       |   |
|       |       |       |       |   |
|       |       |       |       |   |
|       |       |       |       |   |

Target

### Un-Supervised Learning

| $X_1$ | $X_2$ | $X_3$ | $X_p$ | |
|-------|-------|-------|-------|---|
|       |       |       |       |   |
|       |       |       |       |   |
|       |       |       |       |   |
|       |       |       |       |   |

No Target

# Regression vs Classification

**Regression**
What is the temperature going to be tomorrow?

PREDICTION
84°

Fahrenheit °F
-50 -40 -30 -20 -10 0 10 20 30 40 50 60 70 80 90 100 110 120 130 140 150 160 170 180 190 200 210 220 230

**Classification**
Will it be Cold or Hot tomorrow?

PREDICTION
COLD    HOT

Fahrenheit °F
-50 -40 -30 -20 -10 0 10 20 30 40 50 60 70 80 90 100 110 120 130 140 150 160 170 180 190 200 210 220 230

**Regression**
What will house prices be like in my town next year?

$208K

Price in $    100K  120K  140K  160K  180K  200K  220K  240K  260K  280K  300K

**Classification**
Will houses be affordable in my town next year?

AFFORDABLE    EXPENSIVE

Price in $    100K  120K  140K  160K  180K  200K  220K  240K  260K  280K  300K

# The Bias-Variance Trade-off



$$\mathrm{E}\left[(y - \hat{f}(x))^2\right] = \mathrm{Bias}[\hat{f}(x)]^2 + \mathrm{Var}[\hat{f}(x)] + \sigma^2$$

Where:

$$\mathrm{Bias}[\hat{f}(x)] = \mathrm{E}[\hat{f}(x) - f(x)]$$

and

$$\mathrm{Var}[\hat{f}(x)] = \mathrm{E}[\hat{f}(x)^2] - \mathrm{E}[\hat{f}(x)]^2$$
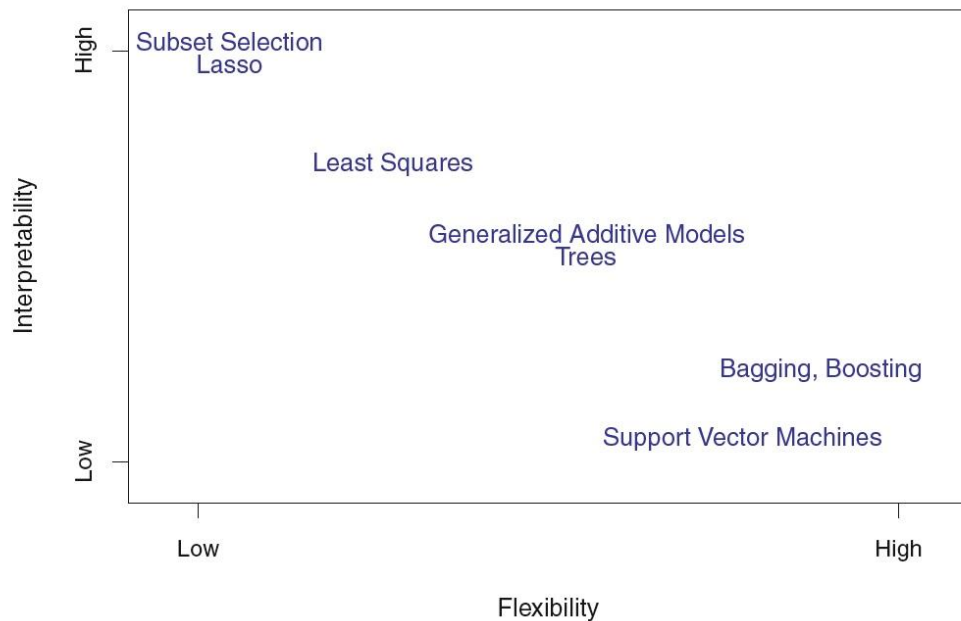
# Flexibility vs Interpretability



**FIGURE 2.7.** *A representation of the tradeoff between flexibility and interpretability, using different statistical learning methods. In general, as the flexibility of a method increases, its interpretability decreases.*