

# INTRODUCTION TO MODELING

*Matt Brems*

*General Assembly, D.C.*

---

# AGENDA

---

1. Data Science Process & Modeling
2. Linear Regression

---

# DATA SCIENCE PROCESS

---

1. Define problem.
2. Gather data.
3. Explore data.
4. Model with data.
5. Evaluate model.
6. Answer problem.

---

# MODELING

---

- Modeling is something that we naturally do.

---

# MODELING

---

- Modeling is something that we naturally do.
- A **model** is a simplification of reality.

---

# MODELING

---

- Modeling is something that we naturally do.
- A **model** is a simplification of reality.
  - How do we simplify?
    - Making assumptions about how things behave.
    - Taking into account only really important factors.

---

# MODELING

---

“Essentially, all models are wrong, but some are useful.”  
– George Box, 1987

---

# MODELING

---





---

# WHY DO WE MODEL?

---

- Prediction
  - How long does it take me to get to work?
  - How much money is a 29-year-old DSI alum expected to make?
- Inference
  - What is the effect of sex on income?
  - How much more money can I be expected to make in a year?

---

# MACHINE LEARNING ALGORITHMS

---

- **Machine learning** is a term we use to describe getting computers (machines) to learn without needing to be explicitly programmed.
- There are many different machine learning algorithms we'll cover in the class - from linear regression to neural networks!

---

# MACHINE LEARNING ALGORITHMS

---

---

# TERMINOLOGY

---

- $\mathbf{X}$ : our data, the independent/explanatory variables we use to predict  $\mathbf{Y}$ .
- $\mathbf{Y}$ : our data, the dependent variable we want to predict.
- $\hat{\mathbf{Y}}$ : our predicted values of  $\mathbf{Y}$ .

---

# MODELING GOALS

---

1. Use observed values of  $\mathbf{X}$  and  $\mathbf{Y}$  to model relationship between them.
2. Build model that makes  $\mathbf{Y}$  and  $\hat{\mathbf{Y}}$  as close as possible.
3. Use observed values of  $\mathbf{X}$  and existing model to make predictions  $\hat{\mathbf{Y}}$ .