

# INTRODUCTION TO MODELING

*Matt Brems*

*General Assembly, D.C.*

---

# AGENDA

---

1. Data Science Process & Modeling
2. Linear Regression

---

# DATA SCIENCE PROCESS

---

1. Define problem. *Data Science Problem*
2. Gather data. *.csv, database*
3. Explore data. *EDA*
4. **Model with data.**
5. Evaluate model.
6. Answer problem.

*} 80%*

# MODELING

---

- Modeling is something that we naturally do.

Commute time

✓ 5 minutes to get to Metro  
✓ 20 minutes on Metro (40) (30)  
✓ 5 minutes from Metro to GA ] ~ 30

Lyft: 20

---

# 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.

---

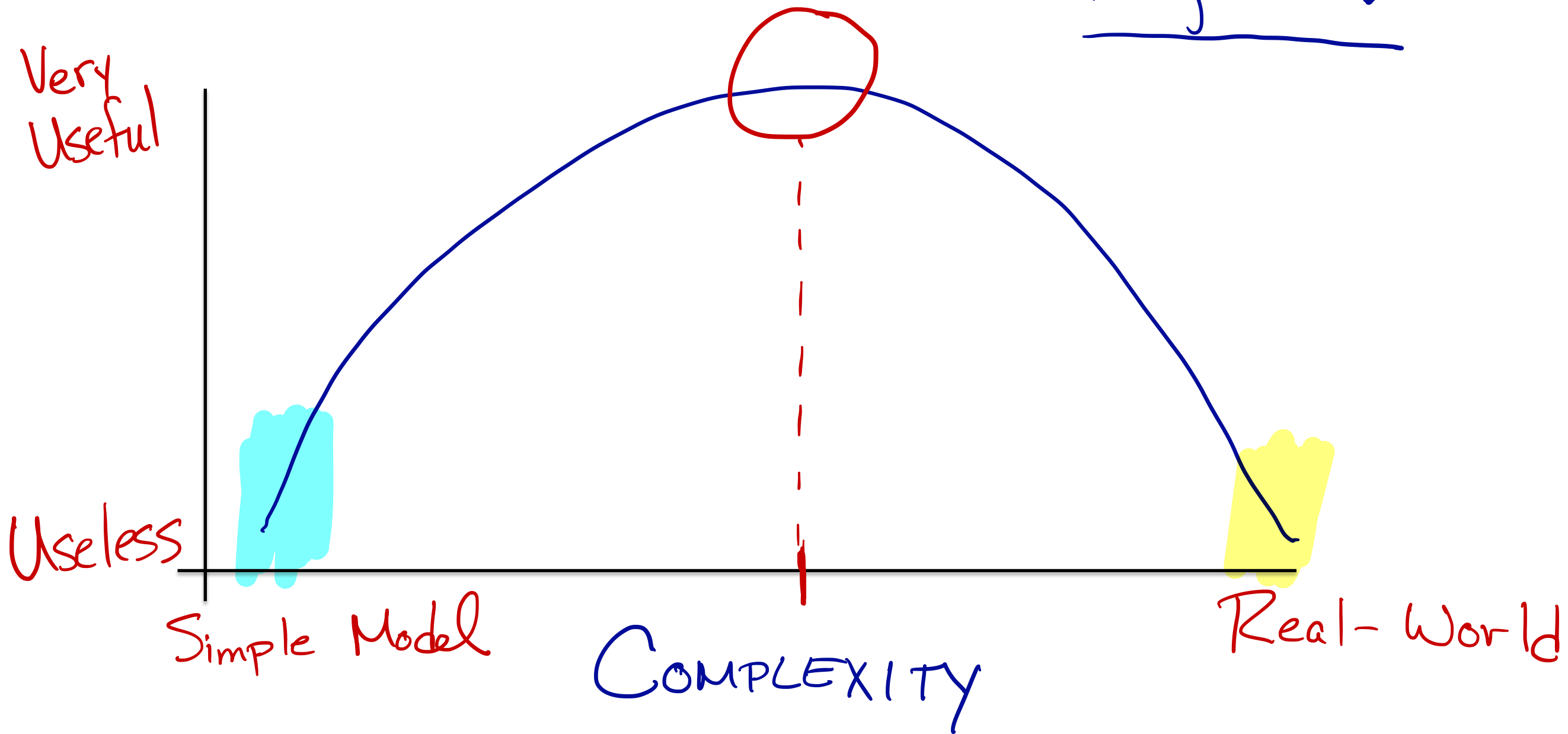
# MODELING

---

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

# MODELING

in general!





---

# WHY DO WE MODEL?

---

- Prediction *stock prices*
  - How long does it take me to get to work?
  - How much money is a 29-year-old DSI alum expected to make?
- Inference *medical studies*
  - 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

wk 8

wk 6, 7

## DATA SCIENCE PROBLEM



### Supervised Learning

↳ have access to our  $Y$  variable

wk 3

### Regression

↳  $Y$  is continuous

wk 4

### Classification

↳  $Y$  is discrete

↓  
output I want to predict

### Unsupervised Learning

↳ do not have access to output we want to predict.

↳ clustering

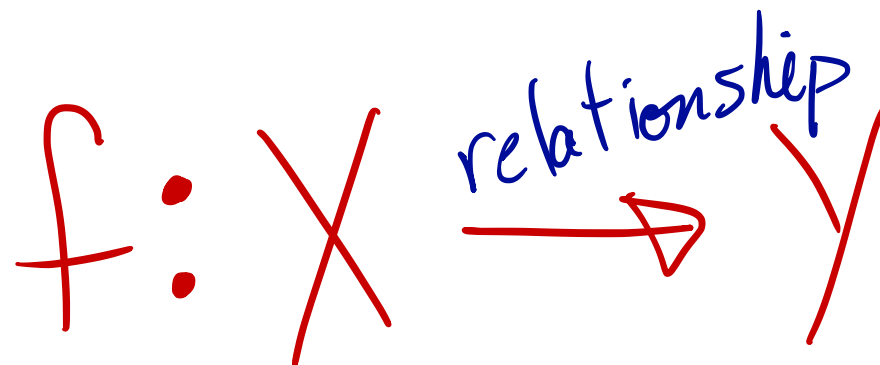
---

# TERMINOLOGY

---

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

↳ "y-hat"



# 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.

observed  $\rightarrow$  predicted

3. Use observed values of  $\mathbf{X}$  and existing model to make predictions  $\hat{\mathbf{Y}}$ .

$f$

$$f(X) = \hat{y}$$