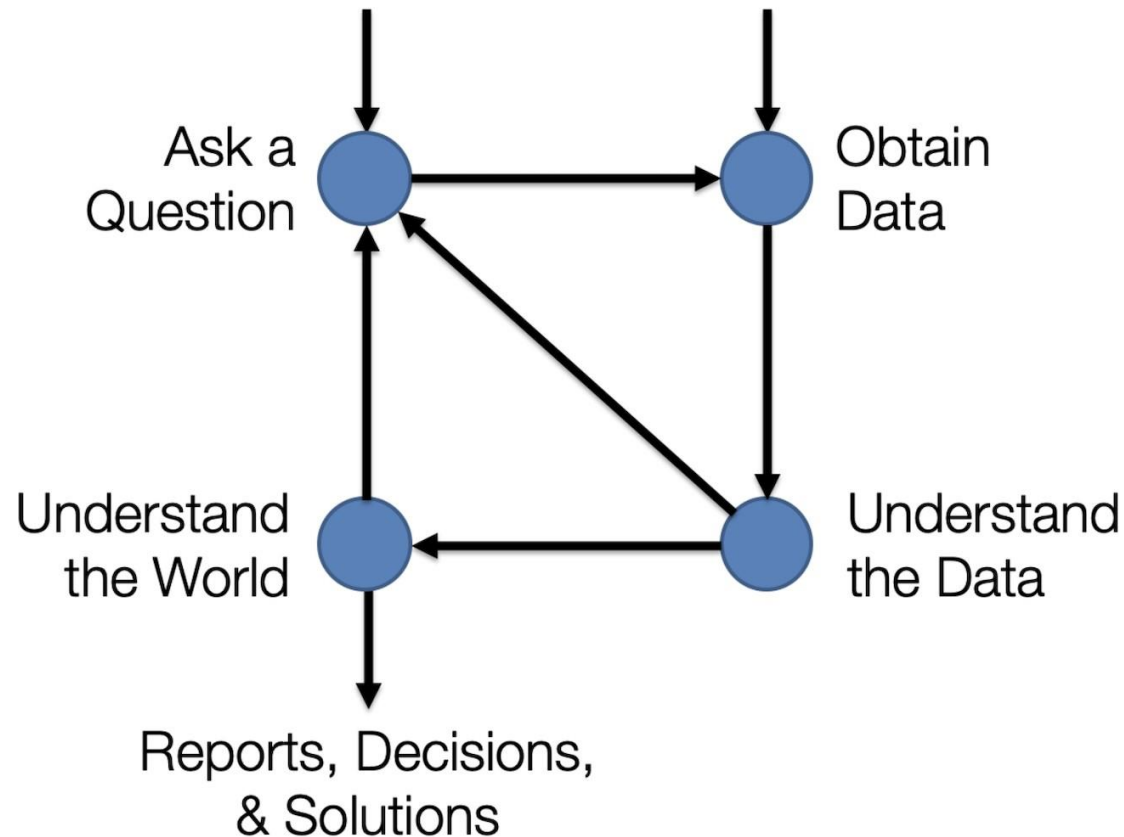


# Modeling

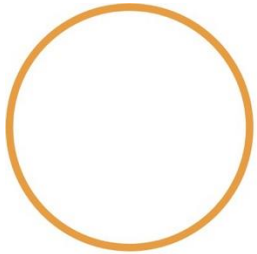
CHOOSING THE MODEL, CHOOSING THE OBJECTIVE,  
FITTING THE MODEL

Sean Kang

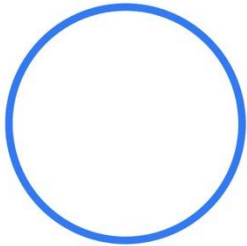
# The Data Science Lifecycle



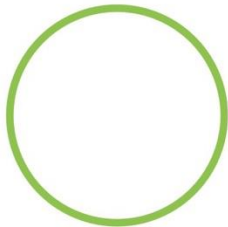
# Sampling



**Population:** the set of all units of interest, size  $N$ .



**Sampling frame:** the set of all possible units that can be drawn into the sample

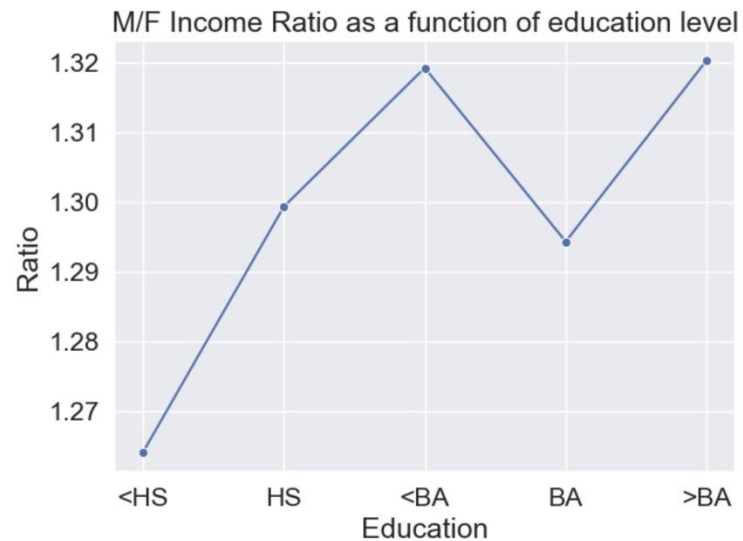
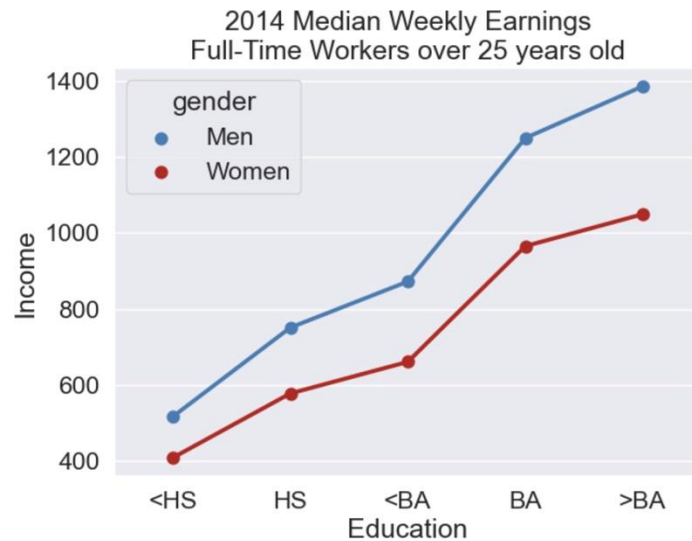


**Sample:** a subset of the sampling frame, size  $n$ .

# Data Wrangling

- Filtering rows
- Selecting columns
- Aggregation
- Pivot Tables
- String methods
  - Regular expressions
- Joins

# Data Visualization



1. Think about your Data

2. Think about your Model

# What is a model?

**Definition:** A model is a useful simplification of reality.

## Example

We can model the fall of an object to earth as subject to a constant acceleration due to gravity at  $9.81 \text{ m/s}^2$ .

The model ignores:

- 
- local variation in gravity
- air resistance
- non-linear dynamics in trajectory

# How can a model be *useful*?

## DESCRIPTION

To understand the world we live in.

- What factors play a role in the spread of COVID-19?
- How do an object's velocity and acceleration impact how far it travels?

## PREDICTION

The predict the value of unseen data.

- Is this email spam?
- Is this shape a pedestrian?



# Two classes of models

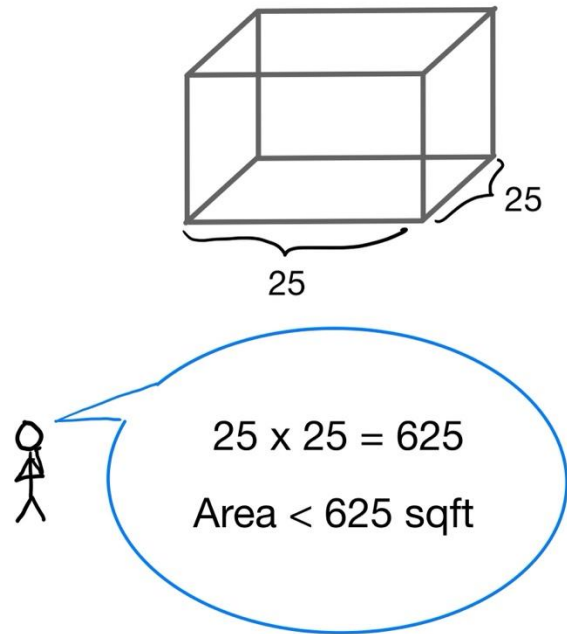
## PHYSICAL MODELS

Based upon well-established theories of how the world works.

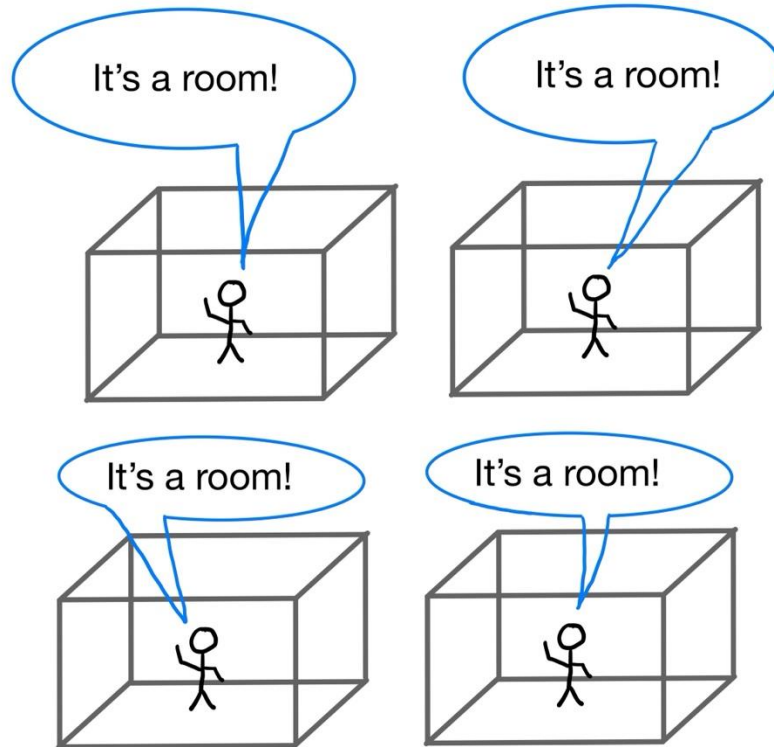
## STATISTICAL MODELS

Based upon observation and data.

## PHYSICAL MODEL

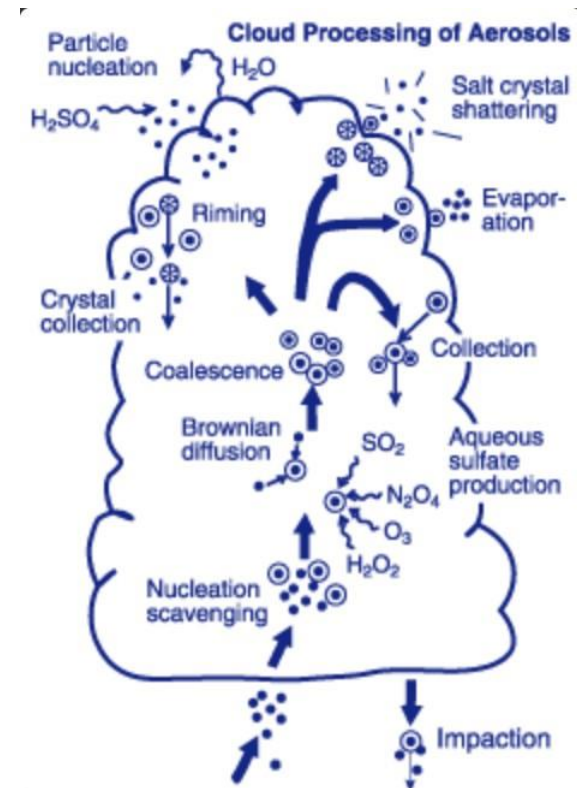
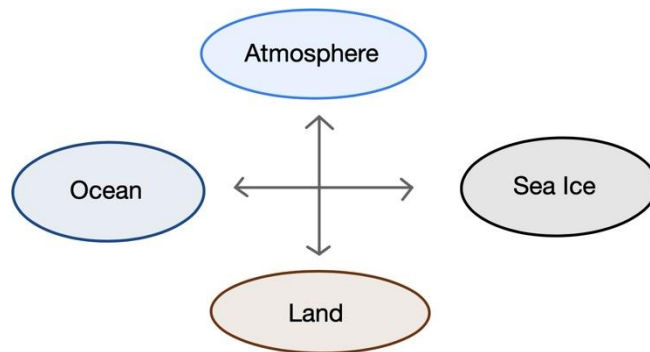


## STATISTICAL MODEL

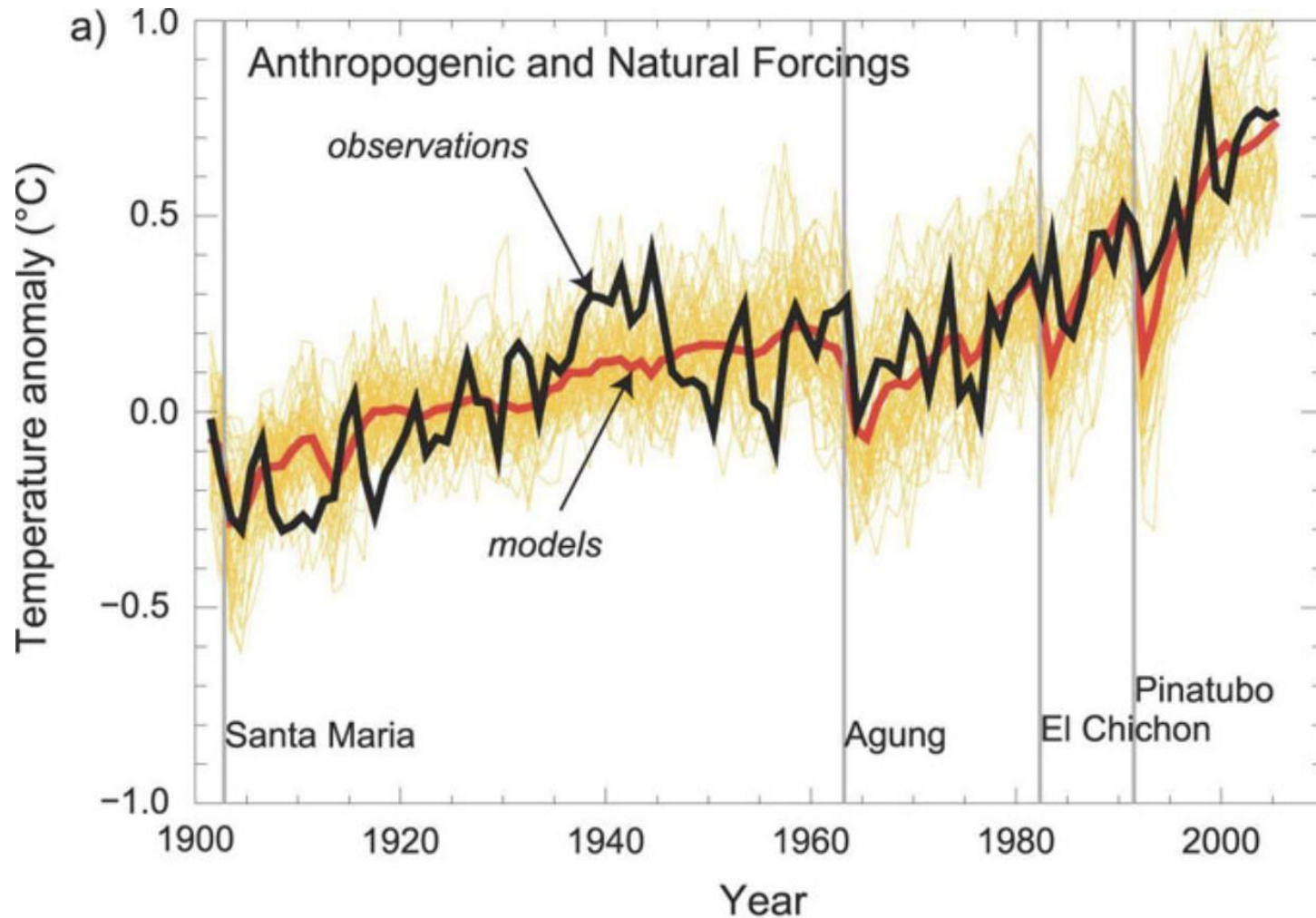


# A physical model: GCM

Global Climate Model (GCM)



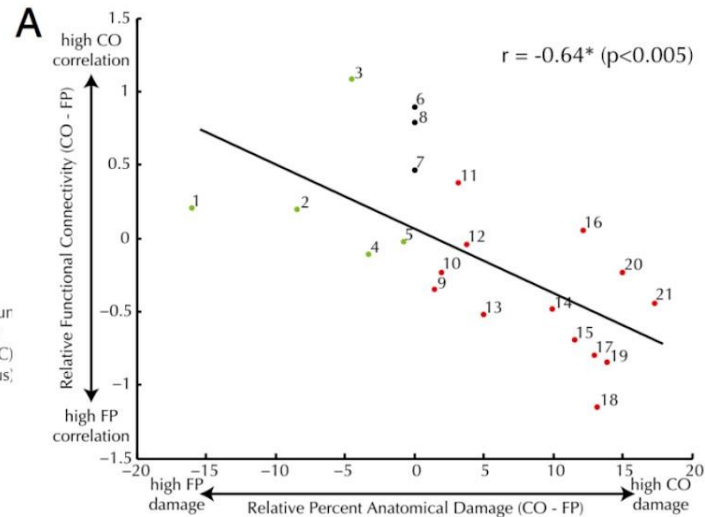
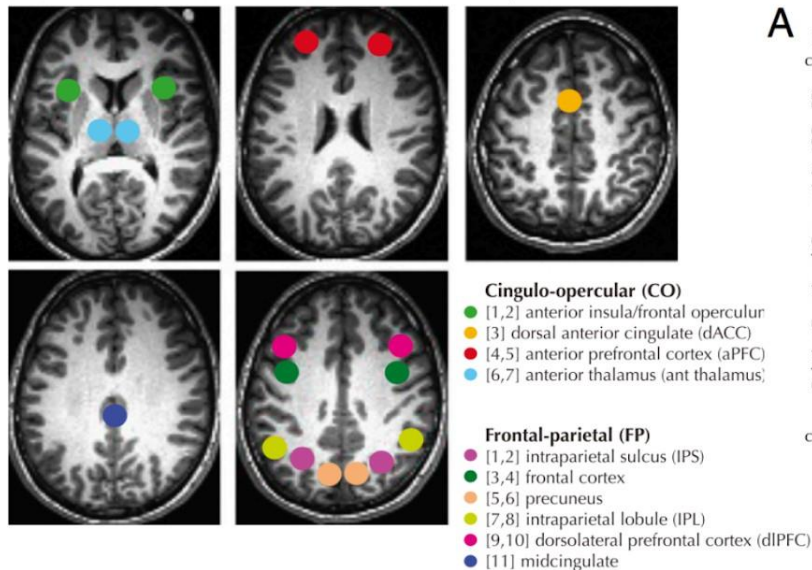
# A physical model: GCM



# A statistical model: FMRI

Other times, we don't have such a precise understanding of some natural relationship. In such cases, we collect data and use statistical tools to learn more about the relationships between variables.

**A** ROI coordinates from Dosenbach et al, 2007



# Choosing the Model

# The modeling process: 3 steps

I. CHOOSE A MODEL

II. CHOOSE AN OBJECTIVE  
FUNCTION

III. FIT THE MODEL BY  
OPTIMIZING YOUR  
OBJECTIVE FUNCTION

# Choose a Model

Choose a simple model

A function with a constant.

A constant model predicts the same number regardless of the circumstances, ignoring all other information.

Example: Tips are 15%

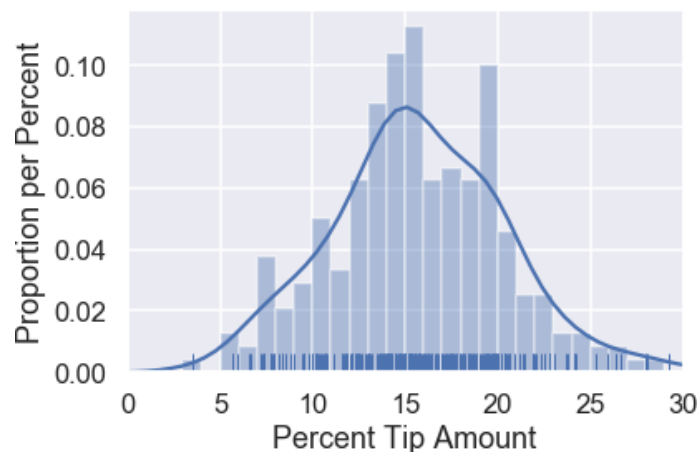


Useful? Descriptive and predictive

Simple? Ignores bill price, time of day, customer type



# The Tips Dataset



Tip rate at a restaurant  
across 244 bills.

Which constant best  
models these tips?

- 15% seems better than 25%
- Is 15% better than 14%
- We need a more precise formulation of this process.

# Notation

$y$  observations (data on tip %)

- $y_i$  individual observations
- $y_1, y_2, \dots, y_n$  data set of size  $n$

$\hat{y}$  Predicted observations (predicted tip %)

- $y_1$  individual prediction

$\theta$  Model parameter(s) ( “true” constant tip %)

$\hat{\theta}$  Fitted, or optimal, parameter(s) (est constant tip %)

# Notation

The constant model can be stated as:

$$\hat{y} = \theta$$

- Parameters define the model
  - Some models are *nonparametric* (e.g. KDEs)
- A constant model ignores any input
- Models can have many parameters

$$\hat{y} = \theta_0 + \theta_1 x \qquad \hat{y} = \frac{1}{1 + \exp(-x^\top \theta)}$$

- Goal: find "best" value of the parameter, denoted  $\hat{\theta}$

# Estimation

Using data to determine model parameters

$$\hat{\theta} = f_1(y, x)$$

# Prediction

Using the fitted model parameters to predict outputs for unseen data

$$\hat{y}_i = f_2(\hat{\theta}, x_i)$$

# The modeling process: 3 steps

## I. CHOOSE A MODEL

- Constant model
- Linear model
- Non-linear model

## II. CHOOSE AN OBJECTIVE FUNCTION

- Prediction: Loss function
- Description: e.g. Likelihood function

## III. FIT THE MODEL BY OPTIMIZING YOUR OBJECTIVE FUNCTION

- Analytical approach (calculus, algebra)
- Numerical approach (optimization, gradient descent)

# Loss functions

# The cost of doing business (making predictions)

We need some metric of how “good” or “bad” our predictions are. This is what loss functions provide us with. **Loss functions quantify how bad a prediction is for a single observation.**

- If our prediction is **close** to the actual value, we want **low loss**.
- If our prediction is **far** from the actual value, we want **high loss**.

A natural choice of loss function is **actual - predicted**, or  $y_i - \hat{y}_i$ . We call this the **error** for a single prediction.

- But, this treats “negative” predictions and “positive” predictions differently.
  - Predicting 16 when the true value is 15 should be penalized the same as predicting 14.
- This leads to two natural loss functions.

# Squared and absolute loss

The most common loss function you'll see is the **squared loss**, also known as L2 loss.

$$L_2(y, \hat{y}) = (y - \hat{y})^2$$

- For a single data point in general, this is  $(y_i - \hat{y}_i)^2$ .
- For our constant model, since  $\hat{y} = \theta$ , this is  $(y_i - \theta)^2$ .

If our prediction is equal to the actual observation, in both cases, our **loss is 0**.

Low loss means a good fit!

Another common loss function is the **absolute loss**, also known as L1 loss.

$$L_1(y, \hat{y}) = |y - \hat{y}|$$

- For our constant model, for a single point, this is  $|y_i - \theta|$ .

There are benefits and drawbacks to both of the above loss functions. We will examine those shortly. **These are also not the only possible loss functions; we will see more later.**



# Loss functions and empirical risk

We care about how bad our model's predictions are for our entire data set, not just for one point. A natural measure, then, is of the **average loss across all points**. Assuming  $n$  points:

$$\frac{1}{n} \sum_{i=1}^n L(y_i, \hat{y}_i)$$

Other names for **average loss** include **empirical risk** and an **objective function**.

**The average loss of a model tells us how well it fits the given data.** If our model has a low average loss across our dataset, that means it is good at making predictions. As such, we want to **find the parameter(s) that minimize average loss**, in order to make our model as good at making predictions as it can be.

# MSE and MAE

If we choose squared loss as our loss function, then average squared loss is typically referred to as **mean squared error (MSE)**, and is of the following form:

$$MSE(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

If we choose absolute loss as our loss function, then average absolute loss is typically referred to as **mean absolute error (MAE)**, and is of the following form:

$$MAE(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

These definitions hold true, regardless of our model. We want to **minimize** these quantities.

# Exploring MSE

Average loss is typically written as a function of  $\theta$ , since  $\theta$  defines what our model is (and hence what our predictions are). For example, with squared loss and the constant model, our average loss (and hence, the function we want to minimize) is

$$R(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2$$

**Average loss is also a function of our data.** But unlike theta, we can't change our data: it is given to us (i.e. it is fixed).

# Exploring MSE

When our model is the constant model, and we choose to use L2 loss, again, our average loss looks like:

$$R(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2$$

Let's examine a toy example. Suppose we have 5 observations, **[20, 21, 22, 29, 33]**.

$$L_2(20, \theta) = (20 - \theta)^2$$

The loss for the first observation ( $y_1$ ).

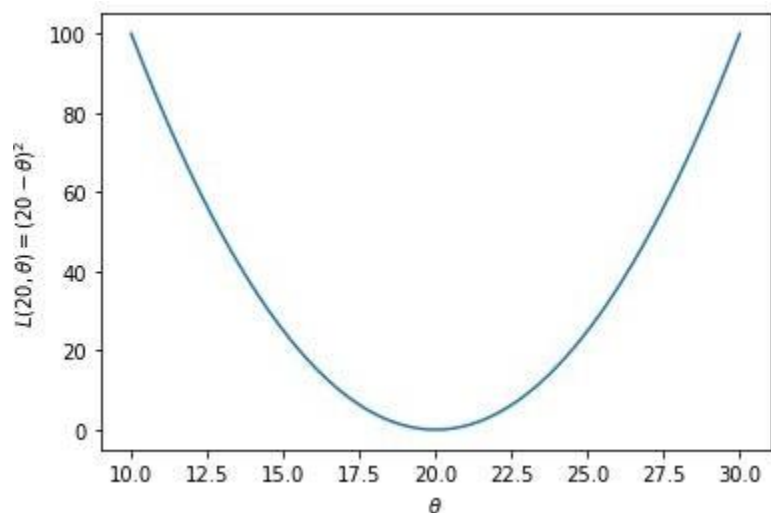
$$R(\theta) = \frac{1}{5} ((20 - \theta)^2 + (21 - \theta)^2 + (22 - \theta)^2 + (29 - \theta)^2 + (33 - \theta)^2)$$

The average loss across all observations (the MSE).

# Exploring MSE

$$L_2(20, \theta) = (20 - \theta)^2$$

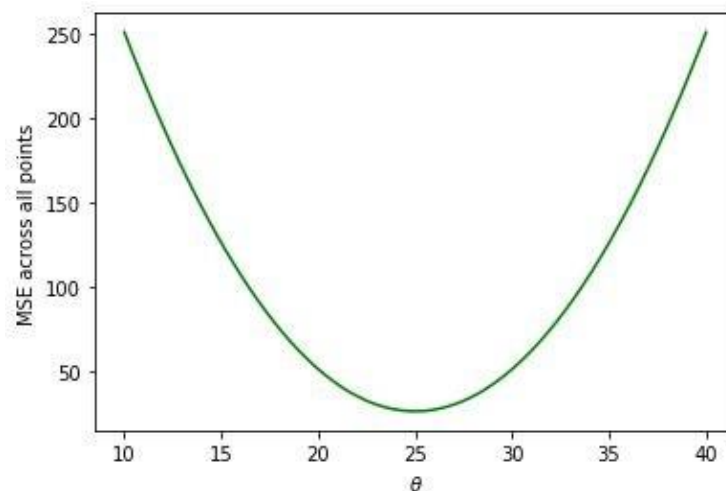
The loss for the first observation ( $y_1$ ).



A parabola, minimized at  $\theta = 20$ .

$$R(\theta) = \frac{1}{5} ((20 - \theta)^2 + (21 - \theta)^2 + (22 - \theta)^2 + (29 - \theta)^2 + (33 - \theta)^2)$$

The average loss across all observations (the MSE).



Also a parabola! Minimized at  $\theta = 25$ .

# Minimizing mean squared error (MSE)

for the constant model

# Minimizing MSE

We saw with the toy example of [20, 21, 22, 29, 33] that the value that minimizes the MSE of the constant model was 25, which was the **mean of our observations**.

We can try other examples if we want to, and we'll end up with the same result. It can be proven using mathematics. There are two ways we'll go about doing this:

- Using calculus.
- Using algebraic approach.

# Minimizing mean absolute error (MAE)

for the constant model



# Exploring MAE

When we use absolute (or L1) loss, we call the average loss **mean absolute error**. For the constant model, our MAE looks like:

$$R(\theta) = \frac{1}{n} \sum_{i=1}^n |y_i - \theta|$$

Let's again re-visit our toy example of 5 observations, **[20, 21, 22, 29, 33]**.

$$L_1(20, \theta) = |20 - \theta|$$

The loss for the first observation ( $y_1$ ).

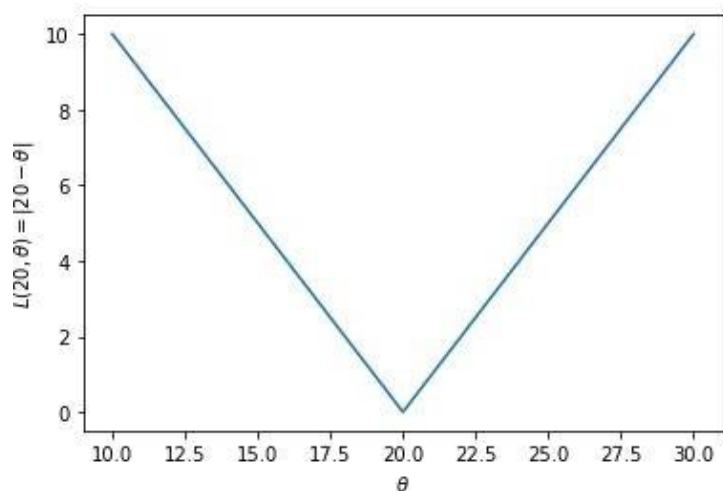
$$R(\theta) = \frac{1}{5} (|20 - \theta| + |21 - \theta| + |22 - \theta| + |29 - \theta| + |33 - \theta|)$$

The average loss across all observations (the MAE).

# Exploring MAE

$$L_1(20, \theta) = |20 - \theta|$$

The loss for the first observation ( $y_1$ ).



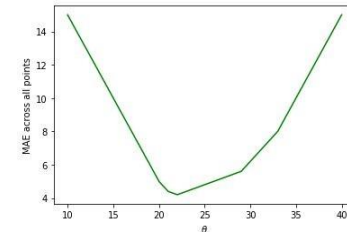
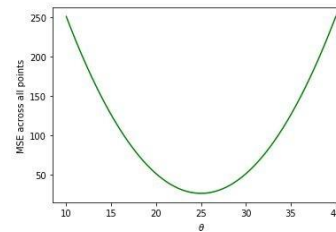
An absolute value curve, centered at  $\theta = 20$ .

$$R(\theta) = \frac{1}{5} (|20 - \theta| + |21 - \theta| + |22 - \theta| + |29 - \theta| + |33 - \theta|)$$

The average loss across all observations (the MAE).

# Comparing loss functions

# MSE vs. MAE



What else is different about squared loss (MSE) and absolute loss (MAE)?

**Mean squared error** (optimal parameter for the constant model is the [sample mean](#))

- **Very smooth**. Easy to minimize using numerical methods (coming later in the course).
- **Very sensitive to outliers**, e.g. if we added 1000 to our largest observation, the optimal  $\theta$  would become 225 instead of 25.

**Mean absolute error** (optimal parameter for the constant model is the [sample median](#))

- **Not as smooth** – at each of the “kinks,” it’s not differentiable. Harder to minimize.
- **Robust to outliers!** E.g, adding 1000 to our largest observation doesn’t change the median.

It’s not clear that one is “better” than the other.

In practice, **we get to choose our loss function!**

# The modeling process

We've implicitly introduced this three-step process, which we will use constantly throughout the rest of the course.

Choose a model

Choose a loss function

Fit the model by minimizing average loss

# The modeling process

We've implicitly introduced this three-step process, which we will use constantly throughout the rest of the course.

Choose a model

Choose a loss function

Fit the model by minimizing average loss

In this lecture, we focused exclusively on the **constant model**, which has a single **parameter**.

**Parameters define our model.** They tell us the relationship between the variables involved in our model. (Not all models have parameters, though!)

In the coming lectures, we will look at more sophisticated models.

# The modeling process

We've implicitly introduced this three-step process, which we will use constantly throughout the rest of the course.

Choose a model

Choose a loss function

Fit the model by minimizing average loss

We introduced two loss functions here: L2 (**squared**) loss and L1 (**absolute**) loss. There also exist others.

Both have their benefits and drawbacks. **We get to choose** which loss function we use, for any modeling task.

# The modeling process

We've implicitly introduced this three-step process, which we will use constantly throughout the rest of the course.

Choose a model

Choose a loss function

Fit the model by minimizing average loss

Lastly, we choose the **optimal parameters** by determining the parameters that **minimize average loss** across our entire dataset. **Different loss functions lead to different optimal parameters.**

This process is called **fitting the model to the data**. We did it by hand here, but in the future we will rely on computerized techniques.



# Summary

# Vocabulary review

- When we use squared (L2) loss as our loss function, the average loss across our dataset is called **mean squared error**.
  - “Squared loss” and “mean squared error” are not the exact same thing – one is for a single observation, and one is for an entire dataset.
  - But they are closely related.
- A similar relationship holds true between absolute (L1) loss and **mean absolute error**.
- Loss functions and summary statistics you already knew:
  - The **sample mean** is the value of  $\theta$  that minimizes the **mean squared error**.
  - The **sample median** is the value of  $\theta$  that minimizes the **mean absolute error**.
- “Average loss” and “empirical risk” mean the same thing for our purposes.
  - So far, our empirical risk was either mean squared error, or mean absolute error.
  - But generally, average loss / empirical risk could be the mean of any loss function across our dataset.

# What's next...

- **Changing the model.**
  - Next, we'll introduce the simple linear regression model that you saw in Data 8.
  - We'll also look at multiple regression, logistic regression, decision trees, and random forests, all of which are different types of models.
- **Changing the loss function.**
  - L2 loss (and, hence, mean squared error) will appear a lot.
  - But we'll also introduce new loss functions, like cross-entropy loss.
- **Changing how we fit the model to the data.**
  - We did this largely by hand in this lecture.
  - But shortly, we'll run into combinations of models and loss functions for which the optimal parameters can't be determined by hand.
  - As such, we'll learn about techniques like gradient descent.