

Lecture 5

# Simple Linear Regression

DSC 40A, Fall 2024

# Announcements

- Homework 1 is due **Friday night**.
- Look at the office hours schedule [here](#) and plan to start regularly attending!
- Remember to take a look at the supplementary readings linked on the course website.

# Agenda

- 0-1 loss
- Predictin rules using features
- Simple linear regression.
- Minimizing mean squared error for the simple linear model.

Question 🤔

Answer at [q.dsc40a.com](https://q.dsc40a.com)

**Remember, you can always ask questions at [q.dsc40a.com](https://q.dsc40a.com)!**

If the direct link doesn't work, click the "🤔 Lecture Questions"  
link in the top right corner of [dsc40a.com](https://dsc40a.com).

## Another example: 0-1 loss

Consider, for example, the **0-1 loss**:

$$L_{0,1}(y_i, h) = \begin{cases} 0 & y_i = h \\ 1 & y_i \neq h \end{cases}$$

The corresponding empirical risk is:

$$R_{0,1}(h) = \frac{1}{n} \sum_{i=1}^n L_{0,1}(y_i, h)$$

## Question 🤔

Answer at [q.dsc40a.com](https://q.dsc40a.com)

$$R_{0,1}(h) = \frac{1}{n} \sum_{i=1}^n \begin{cases} 0 & y_i = h \\ 1 & y_i \neq h \end{cases}$$

Suppose  $y_1, y_2, \dots, y_n$  are all unique. What is  $R_{0,1}(y_1)$ ?













- A. 0.
- B.  $\frac{1}{n}$ .
- C.  $\frac{n-1}{n}$ .
- D. 1.

## Minimizing empirical risk for 0-1 loss

$$R_{0,1}(h) = \frac{1}{n} \sum_{i=1}^n \begin{cases} 0 & y_i = h \\ 1 & y_i \neq h \end{cases}$$

## Summary: Choosing a loss function

Key idea: Different loss functions lead to different best predictions,  $h^*$ !

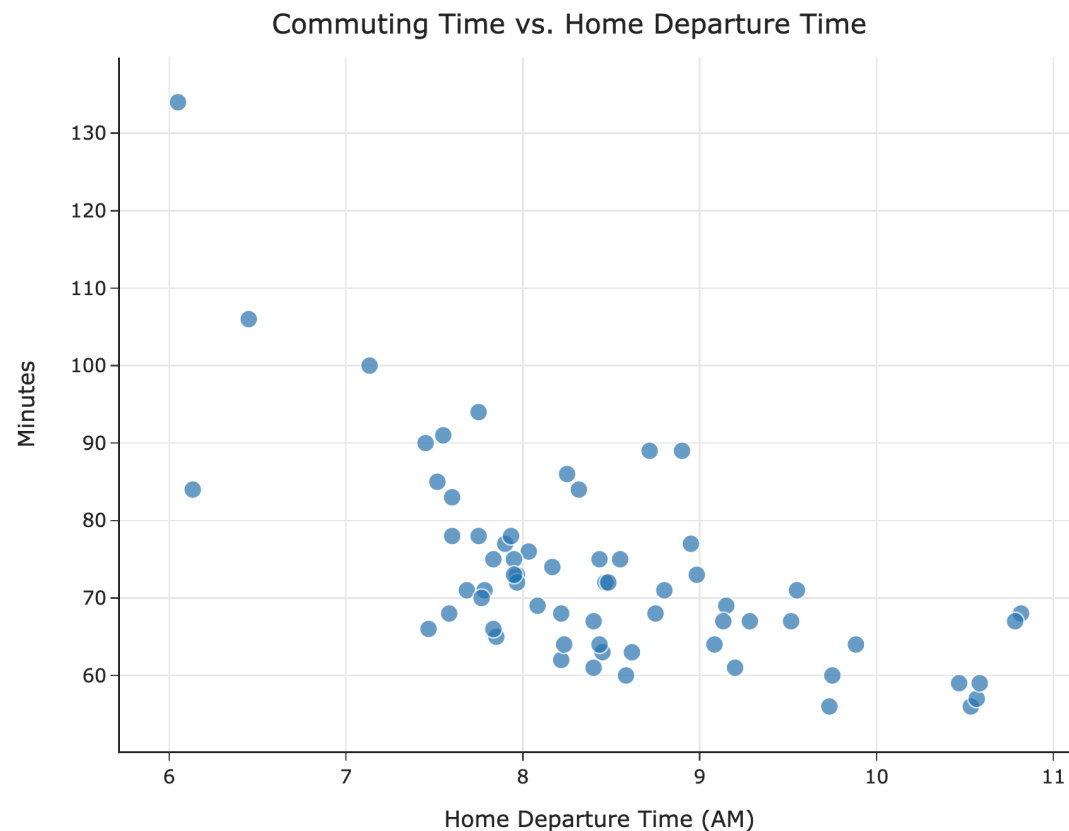
Loss	Minimizer	Always Unique?	Robust to Outliers?	Differentiable?
$L_{sq}$	mean	yes 	no 	yes 
$L_{abs}$	median	no 	yes 	no 
$L_{\infty}$	midrange	yes 	no 	no 
$L_{0,1}$	mode	no 	yes 	no 

The optimal predictions,  $h^*$ , are all **summary statistics** that measure the **center** of the dataset in different ways.



# Predictions with features

# Towards simple linear regression



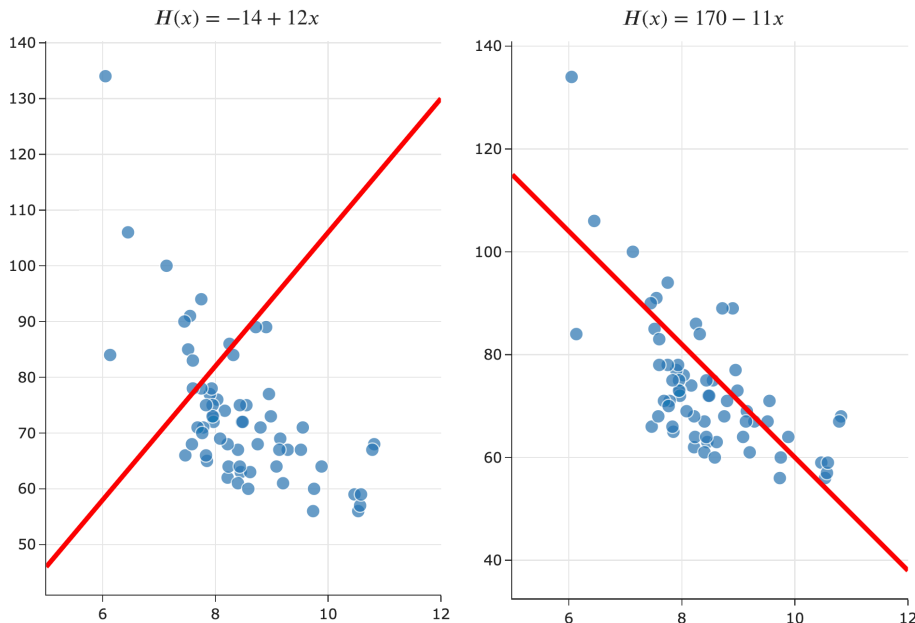
- In Lecture 1, we introduced the idea of a hypothesis function,  $H(x)$ .
- We've focused on finding the best **constant model**,  $H(x) = h$ .
- Now that we understand the modeling recipe, we can apply it to find the best **simple linear regression model**,  $H(x) = w_0 + w_1x$ .
- This will allow us to make predictions that aren't all the same for every data point.

## Recap: Hypothesis functions and parameters

A hypothesis function,  $H$ , takes in an  $x$  as input and returns a predicted  $y$ .

**Parameters** define the relationship between the input and output of a hypothesis function.

The simple linear regression model,  $H(x) = w_0 + w_1x$ , has two parameters:  $w_0$  and  $w_1$ .



# The modeling recipe

1. Choose a model.
2. Choose a loss function.
3. Minimize average loss to find optimal model parameters.

# Features

A **feature** is an attribute of the data – a piece of information.

- **Numerical**: maximum allowed speed, time of departure
- **Categorical**: day of week
- **Boolean**: was there a car accident on the road?

Think of features as columns in a DataFrame (i.e. table).

Departure time	Day of week	Accident on route	Commute time
7:05	Monday	yes	101
8:03	Tuesday	no	87
10:20	Wednesday	yes	79
8:30	Thursday	no	76

# Modeling

- We believe that commute time is a function of departure time.
- I.e., there is a function  $H$  so that:  
commute time  $\approx H(\text{departure time})$
- $H$  is called a **hypothesis function** or **prediction rule**.
- Our goal: find a good prediction rule  $H$ .

## Possible Hypothesis Functions

- $H_1(\text{departure time}) = 90 - 10 \cdot (\text{departure time} - 7)$
- $H_2(\text{departure time}) = 90 - (\text{departure time} - 8)^2$
- $H_3(\text{departure time}) = 20 + 6 \cdot \text{departure time}$

These are all valid prediction rules.

Some are better than others.

## Comparing predictions

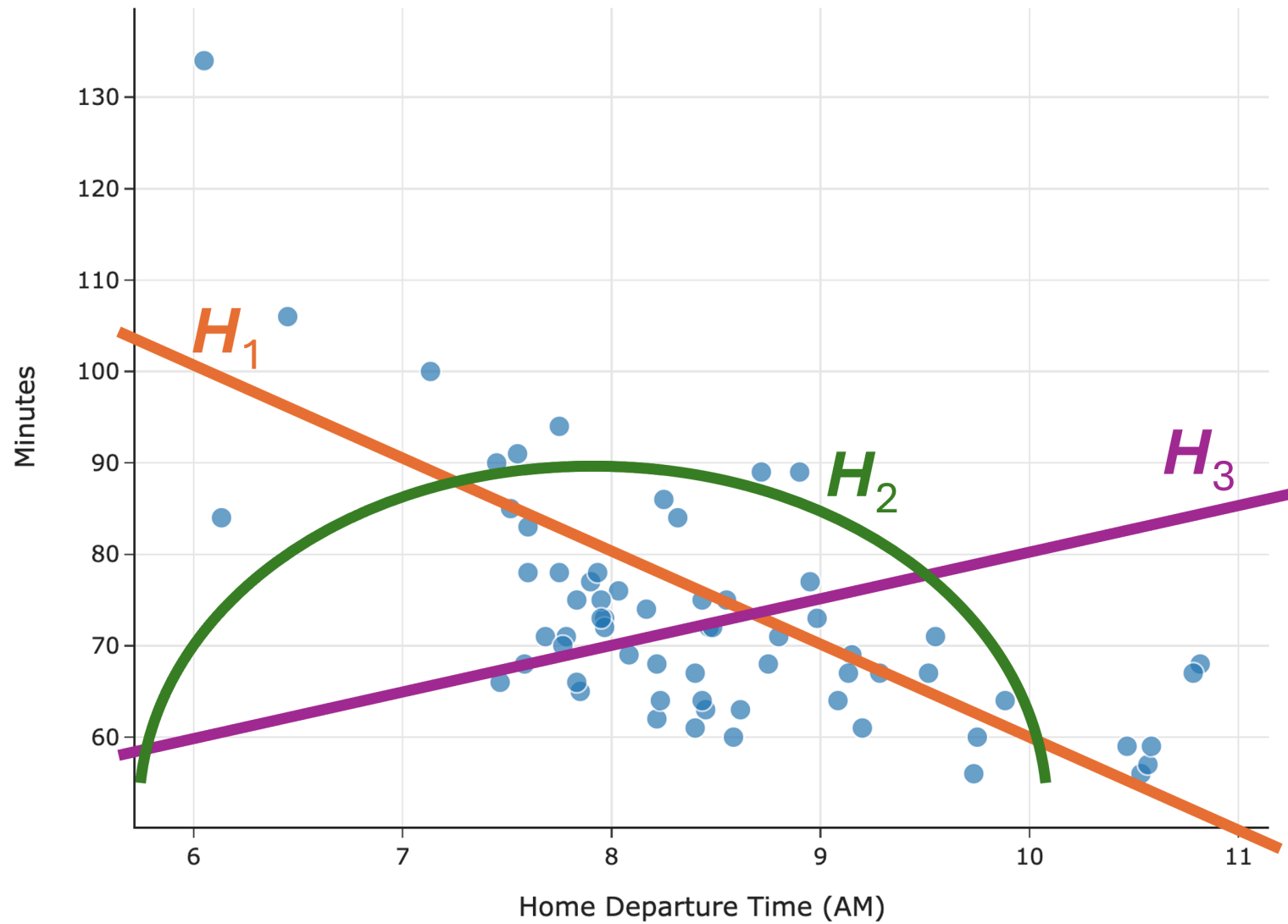
- How do we know which hypothesis is best:  $H_1$ ,  $H_2$ ,  $H_3$ ?
- We gather data from  $n$  days of commute. Let  $x_i$  be departure time,  $y_i$  be commute time:

(departure time <sub>1</sub> , commute time <sub>1</sub> )	$(x_1, y_1)$
(departure time <sub>2</sub> , commute time <sub>2</sub> )	$(x_2, y_2)$
...	$\rightarrow$
(departure time <sub><math>n</math></sub> , commute time <sub><math>n</math></sub> )	$(x_n, y_n)$

- See which rule works better on data.



Commuting Time vs. Home Departure Time



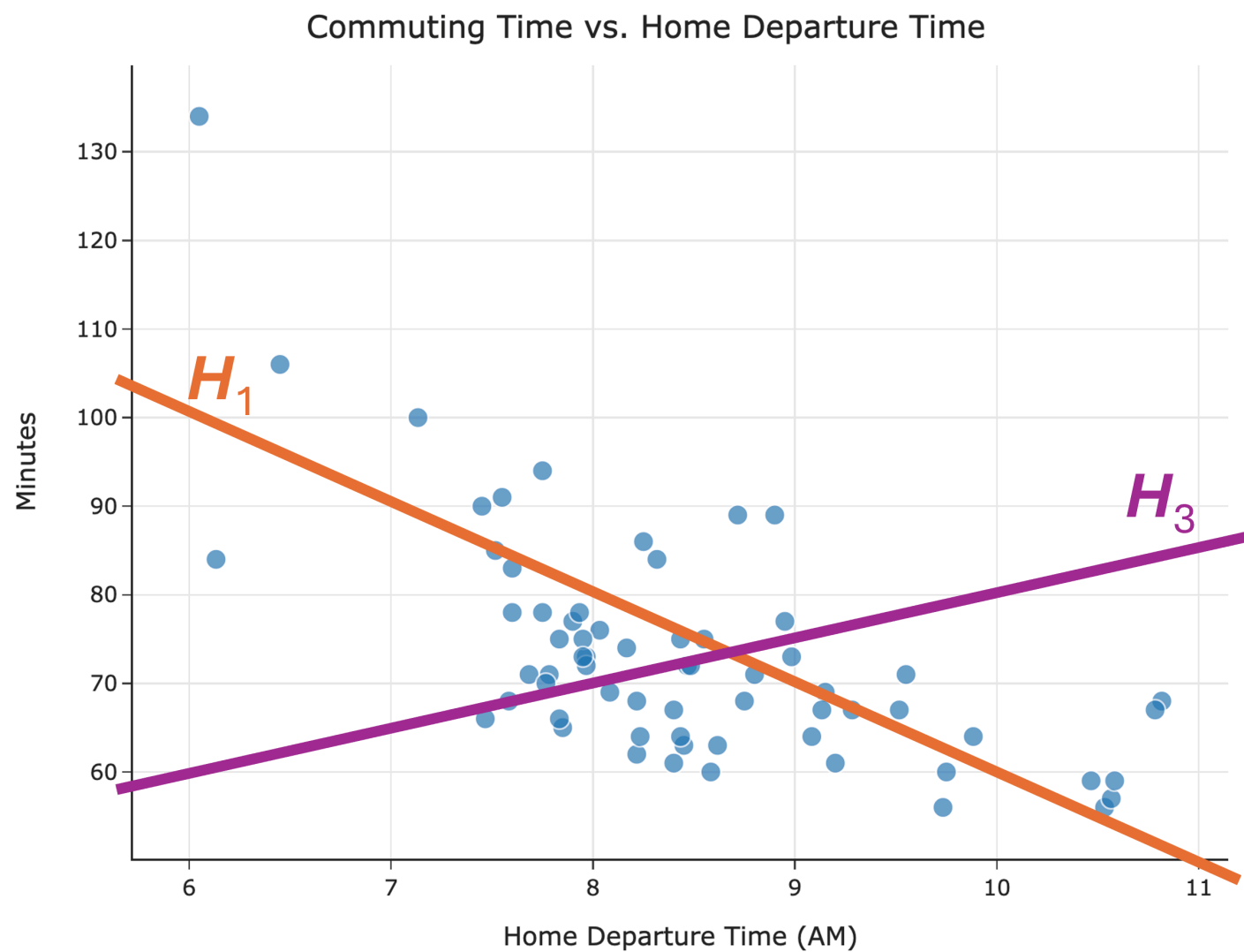
# Quantifying the performance of a model

- Reminder: one loss function, which measures how far  $H(x_i)$  is from  $y_i$ , is **absolute loss**.
- The mean absolute error of  $H(x)$  is

$$R_{\text{abs}}(h) = \frac{1}{n} \sum_{i=1}^n |y_i - H(x_i)|$$

- We want the **best** prediction,  $H^*(x)$ .
- The smaller  $R_{\text{abs}}(h)$  is, the better the hypothesis.

# Mean absolute error



## Finding the best hypothesis $H(x)$

- Goal: out of all functions  $\mathbb{R} \rightarrow \mathbb{R}$ , find the function  $H$  with the smallest mean absolute error.
- That is,  $H^*$  should be the function that minimizes

$$R_{\text{abs}}(h) = \frac{1}{n} \sum_{i=1}^n |y_i - H(x_i)|$$

## Finding the best hypothesis $H(x)$

- Goal: out of all functions  $\mathbb{R} \rightarrow \mathbb{R}$ , find the function  $H$  with the smallest mean absolute error.
- That is,  $H^*$  should be the function that minimizes

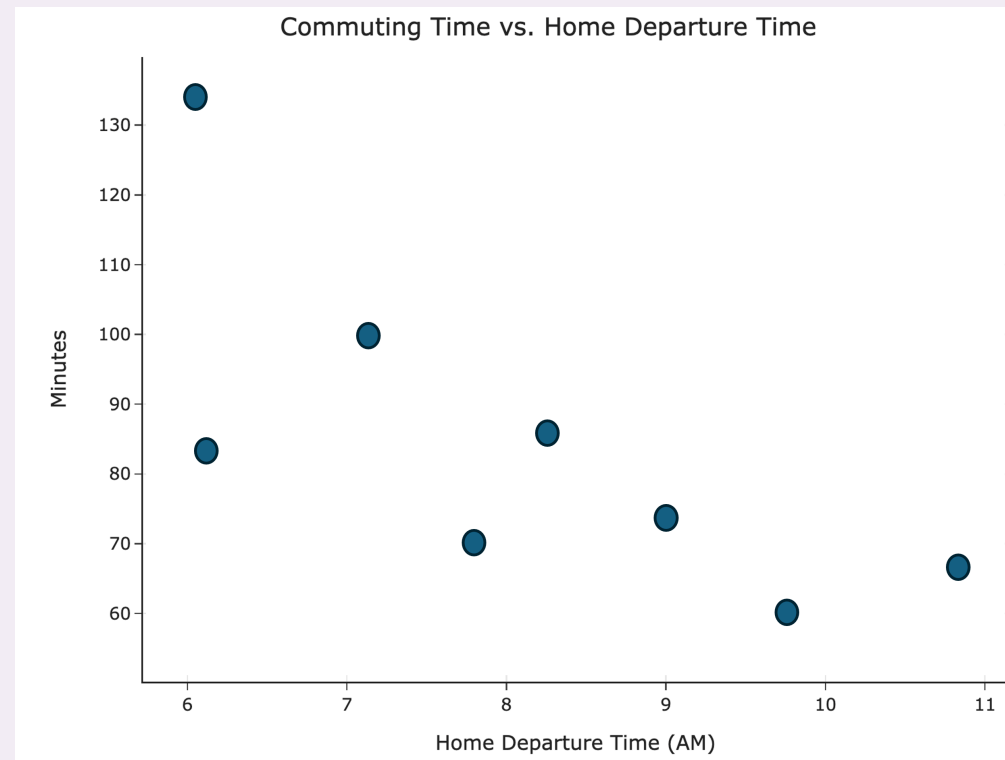
$$R_{\text{abs}}(h) = \frac{1}{n} \sum_{i=1}^n |y_i - H(x_i)|$$

- There are two problems with this.

## Question 🤔 Answer at [q.dsc40a.com](http://q.dsc40a.com)

Given the data below, is there a prediction rule  $H$  which has zero mean absolute error?

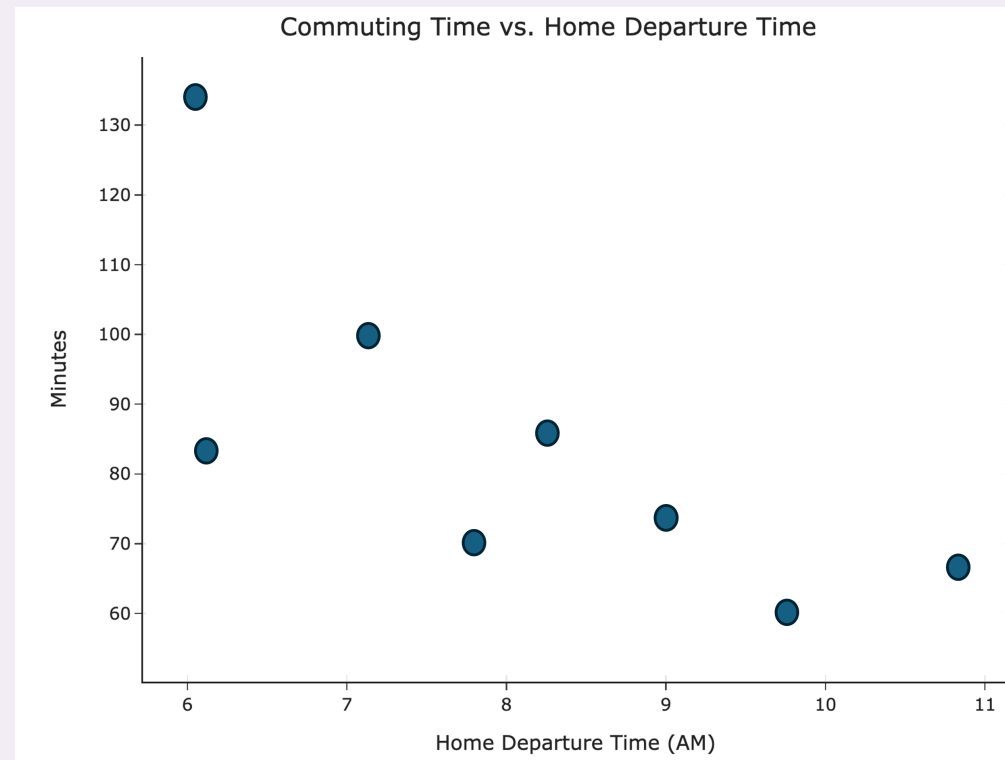
- A. yes
- B. no



## Question 🤔 Answer at [q.dsc40a.com](http://q.dsc40a.com)

Given the data below, is there a prediction rule  $H$  which has zero mean absolute error?

- A. yes
- B. no



# Problem

- We can make mean squared error very small, even zero!
- But the function will be weird.
- This is called **overfitting**.
- Remember our real goal: make good predictions on data **we haven't seen**.



## Solution

- Don't allow  $H$  to be just any function.
- Require that it has a certain form.
- Examples:
  - Linear:  $H(x) = w_0 + w_1x$ .
  - Quadratic:  $H(x) = w_0 + w_1x_1 + w_2x^2$ .
  - Exponential:  $H(x) = w_0e^{w_1x}$ .
  - Constant:  $H(x) = w_0$ .

## Finding the best linear model

- **Goal:** Out of all linear functions  $\mathbb{R} \rightarrow \mathbb{R}$ , find the function  $H^*$  with the smallest mean squared error.
  - Linear functions are of the form  $H(x) = w_0 + w_1x$ .
  - They are defined by a slope ( $w_1$ ) and intercept ( $w_0$ ).
- That is,  $H^*$  should be the linear function that minimizes

$$R_{abs}(H) = \frac{1}{n} \sum_{i=1}^n |y_i - H(x_i)|$$

## Finding the best linear model

- **Goal:** Out of all linear functions  $\mathbb{R} \rightarrow \mathbb{R}$ , find the function  $H^*$  with the smallest mean squared error.
  - Linear functions are of the form  $H(x) = w_0 + w_1x$ .
  - They are defined by a slope ( $w_1$ ) and intercept ( $w_0$ ).
- That is,  $H^*$  should be the linear function that minimizes

$$R_{abs}(H) = \frac{1}{n} \sum_{i=1}^n |y_i - H(x_i)|$$

- ◦ **There is still a problem with this.**

## Problem #2

It is hard to minimize the mean absolute error:

$$R_{abs}(H) = \frac{1}{n} \sum_{i=1}^n |y_i - H(x_i)|$$

- Not differentiable!
- What can we do?

# Minimizing mean squared error for the simple linear model

- We'll choose squared loss, since it's the easiest to minimize.
- Our goal, then, is to find the linear hypothesis function  $H^*(x)$  that minimizes empirical risk:

$$R_{\text{sq}}(H) = \frac{1}{n} \sum_{i=1}^n (y_i - H(x_i))^2$$

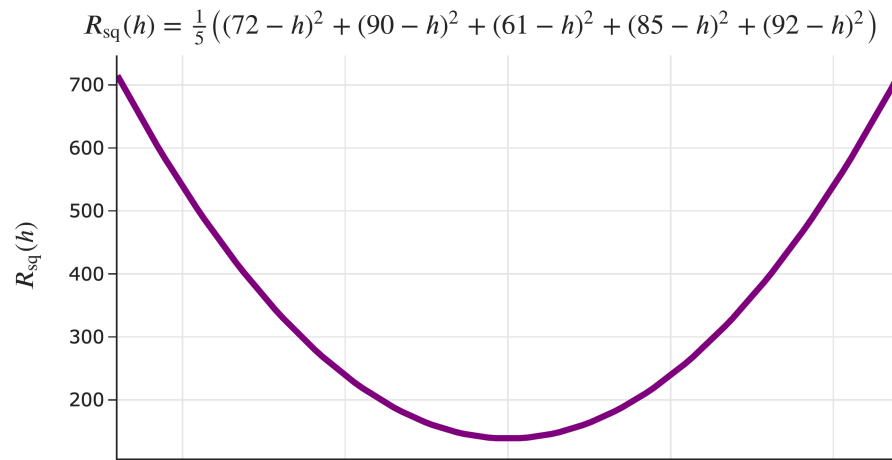
- Since linear hypothesis functions are of the form  $H(x) = w_0 + w_1x$ , we can rewrite  $R_{\text{sq}}$  as a function of  $w_0$  and  $w_1$ :

$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1x_i))^2$$

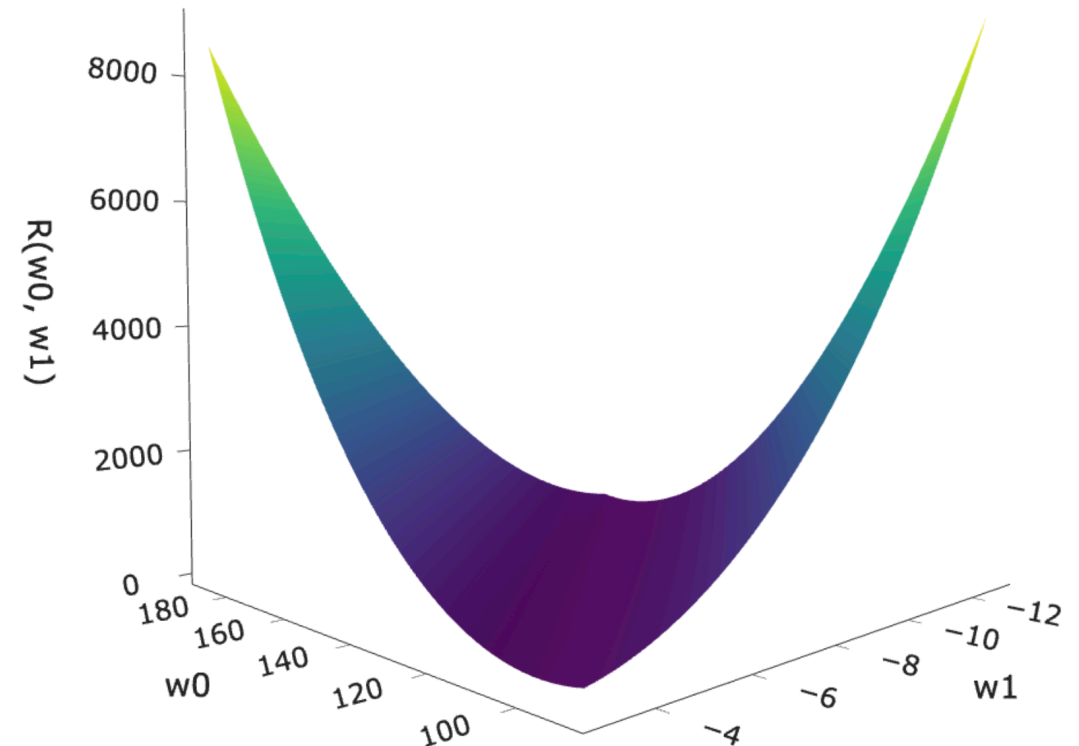
- How do we find the parameters  $w_0^*$  and  $w_1^*$  that minimize  $R_{\text{sq}}(w_0, w_1)$ ?

## Loss surface

For the constant model, the graph of  $R_{\text{sq}}(h)$  looked like a parabola.



What does the graph of  $R_{\text{sq}}(w_0, w_1)$  look like for the simple linear regression model?



# Minimizing mean squared error for the simple linear model

# Minimizing multivariate functions

- Our goal is to find the parameters  $w_0^*$  and  $w_1^*$  that minimize mean squared error:

$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

- $R_{\text{sq}}$  is a function of two variables:  $w_0$  and  $w_1$ .
- To minimize a function of multiple variables:
  - Take partial derivatives with respect to each variable.
  - Set all partial derivatives to 0.
  - Solve the resulting system of equations.
  - Ensure that you've found a minimum, rather than a maximum or saddle point (using the [second derivative test](#) for multivariate functions).



## Example

Find the point  $(x, y, z)$  at which the following function is minimized.

$$f(x, y) = x^2 - 8x + y^2 + 6y - 7$$

## Minimizing mean squared error

$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

To find the  $w_0^*$  and  $w_1^*$  that minimize  $R_{\text{sq}}(w_0, w_1)$ , we'll:

1. Find  $\frac{\partial R_{\text{sq}}}{\partial w_0}$  and set it equal to 0.
2. Find  $\frac{\partial R_{\text{sq}}}{\partial w_1}$  and set it equal to 0.
3. Solve the resulting system of equations.

## Question 🤔

Answer at [q.dsc40a.com](https://q.dsc40a.com)

$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

Which of the following is equal to  $\frac{\partial R_{\text{sq}}}{\partial w_0}$ ?

- A.  $\frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))$
- B.  $-\frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))$
- C.  $-\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) x_i$
- D.  $-\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))$

$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

$$\frac{\partial R_{\text{sq}}}{\partial w_0} =$$

$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

$$\frac{\partial R_{\text{sq}}}{\partial w_1} =$$

## Strategy

We have a system of two equations and two unknowns ( $w_0$  and  $w_1$ ):

$$-\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) = 0 \quad -\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) x_i = 0$$

To proceed, we'll:

1. Solve for  $w_0$  in the first equation.

The result becomes  $w_0^*$ , because it's the "best intercept."

2. Plug  $w_0^*$  into the second equation and solve for  $w_1$ .

The result becomes  $w_1^*$ , because it's the "best slope."

**Solving for  $w_0^*$**

$$-\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) = 0$$

Solving for  $w_1^*$

$$-\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) x_i = 0$$



## Least squares solutions

We've found that the values  $w_0^*$  and  $w_1^*$  that minimize  $R_{\text{sq}}$  are:

$$w_1^* = \frac{\sum_{i=1}^n (y_i - \bar{y})x_i}{\sum_{i=1}^n (x_i - \bar{x})x_i} \qquad w_0^* = \bar{y} - w_1^* \bar{x}$$

where:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \qquad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

These formulas work, but let's re-write  $w_1^*$  to be a little more symmetric.

## An equivalent formula for $w_1^*$

Claim:

$$w_1^* = \frac{\sum_{i=1}^n (y_i - \bar{y})x_i}{\sum_{i=1}^n (x_i - \bar{x})x_i} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Proof:

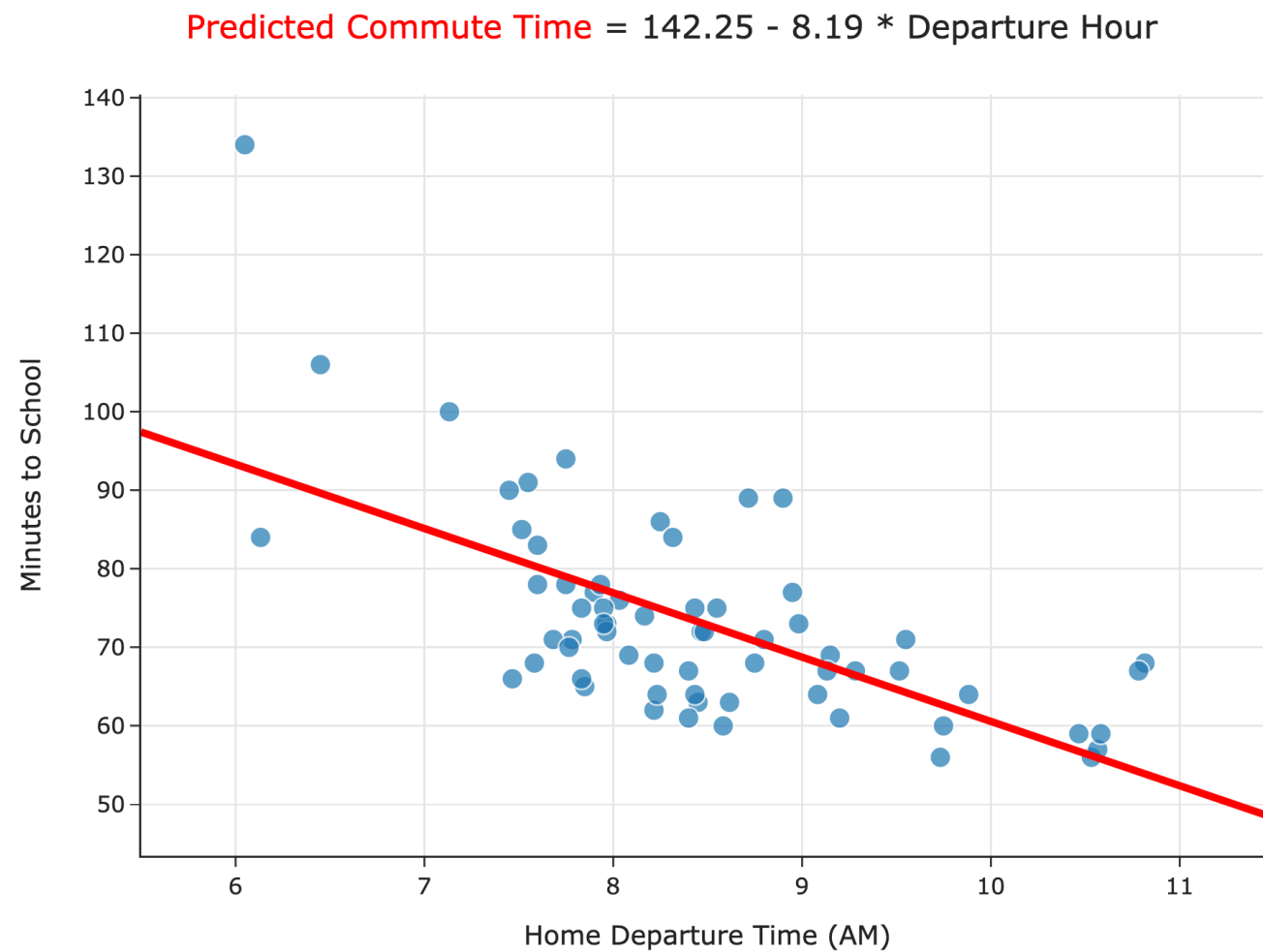
## Least squares solutions

- The **least squares solutions** for the intercept  $w_0$  and slope  $w_1$  are:

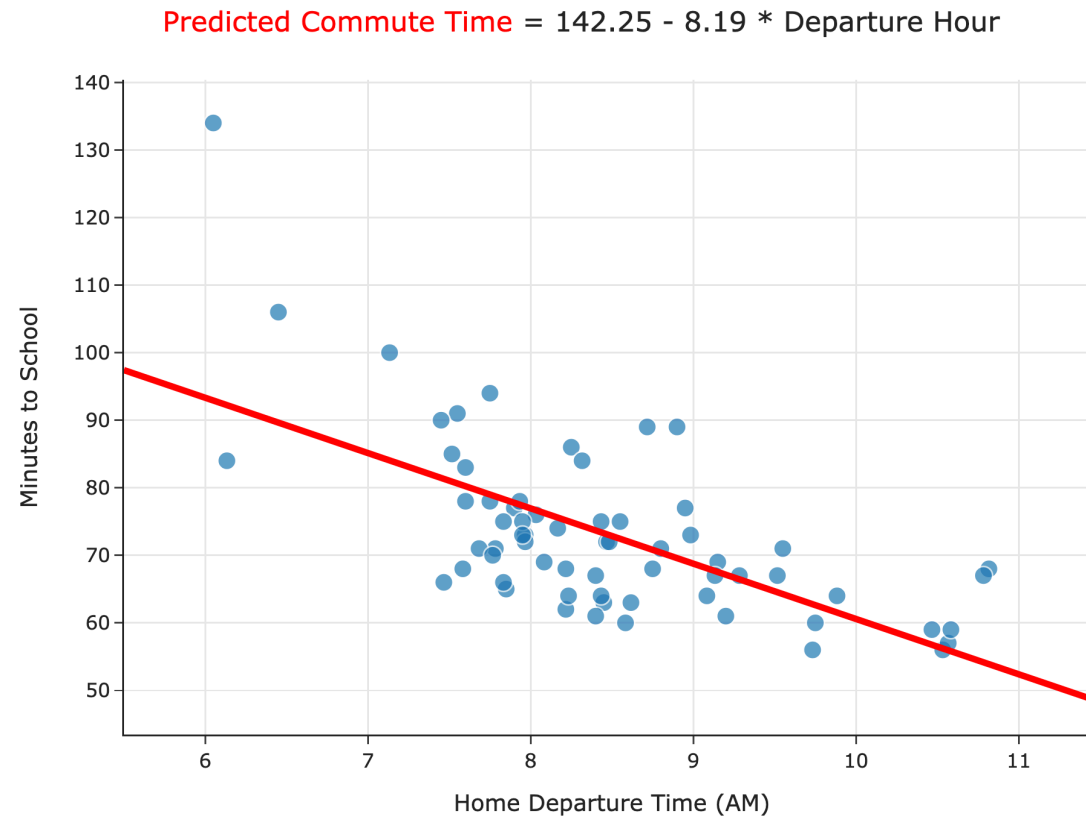
$$w_1^* = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \qquad w_0^* = \bar{y} - w_1^* \bar{x}$$

- We say  $w_0^*$  and  $w_1^*$  are **optimal parameters**, and the resulting line is called the **regression line**.
- The process of minimizing empirical risk to find optimal parameters is also called "**fitting to the data**."
- To make predictions about the future, we use  $H^*(x) = w_0^* + w_1^*x$ .

Let's test these formulas out in code! Follow along [here](#).



# Causality



Can we conclude that leaving later **causes** you to get to school quicker?

## What's next?

We now know how to find the optimal slope and intercept for linear hypothesis functions. Next, we'll:

- See how the formulas we just derived connect to the formulas for the slope and intercept of the regression line we saw in DSC 10.
  - They're the same, but we need to do a bit of work to prove that.
- Learn how to interpret the slope of the regression line.
- Discuss *causality*.
- Learn how to build regression models with **multiple inputs**.
  - To do this, we'll need linear algebra!