**Lectures 5-7**

# Simple Linear Regression

**DSC 40A, Fall 2025**

# 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

- Prediction rules using features

- Simple linear regression.

- Minimizing mean squared error for the simple linear model.

# Question 🤔

Answer at q.dsc40a.com

**Remember, you can always ask questions at q.dsc40a.com!**

If the direct link doesn't work, click the "🤔 Lecture Questions"

link in the top right corner of 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 🤔

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

Suppose $y_1, y_2, \ldots, y_n$ 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? |
|------|-----------|----------------|---------------------|-----------------|
| $L_{\text{sq}}$ | mean | yes ✅ | no ❌ | yes ✅ |
| $L_{\text{abs}}$ | 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



Commuting Time vs. Home Departure Time

- 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**, $H(x) = w_0 + w_1 x$.
- 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, $H(x) = w_0 + w_1 x$, 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 |

# Variables

- The features, $x$, that we base our predictions on are called predictor variables.

- The quantity, $y$, that we're trying to predict based on these features is called the response variable, dependent variable or target.

- We are trying to predict our commute time as a function of departure time.

# Modeling

- We believe that commute time is a function of departure time.
- I.e., there is a function $H$ so that:

  commute time $\approx H$(departure time)
- $H$ is called a hypothesis function or prediction rule.
- Our goal: find a good prediction rule, $H$.

# Possible Hypothesis Functions

- $H_1$(departure time) = 90 - 10 ·(departure time-7)
- $H_2$(departure time) = 90 - (departure time-8)$^2$
- $H_3$(departure time) = 20 + 6·departure time

These are all valid prediction rules.
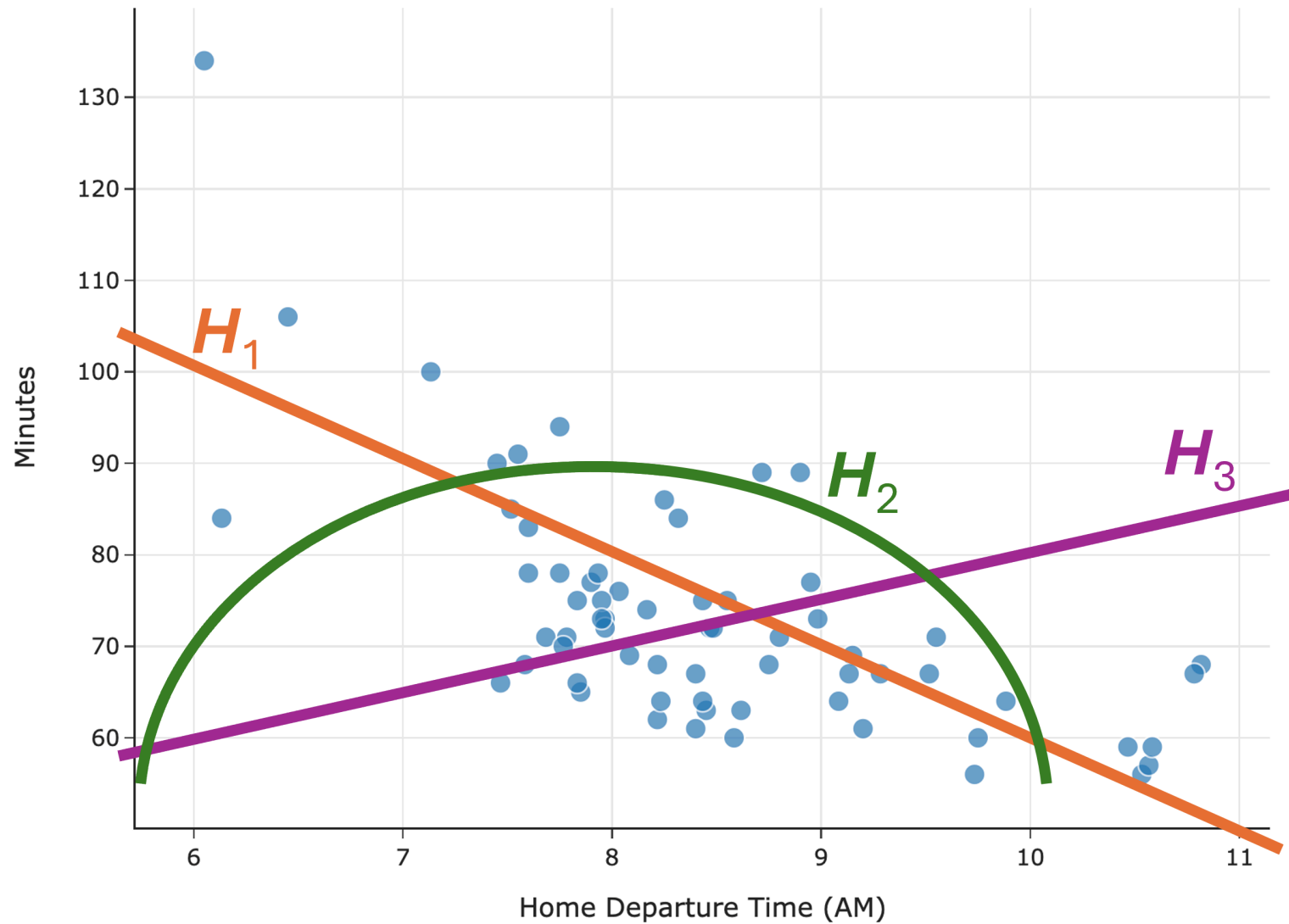
Some are better than others.

# Comparing predictions

- How do we know which is best: $H_1,\ H_2,\ H_3$?

- We gather data from $n$ days of commute. Let xi be experience, yi be salary:

  (departure time$_1$ , commute time$_1$)        $(x_1, y_1)$
  (departure time$_2$ , commute time$_2$)        $(x_2, y_2)$
  $\ldots$        $\rightarrow$
  (departure time$_n$ , commute time$_n$)        $(x_n, y_n)$

- See which rule works better on data.

Commuting Time vs. Home Departure Time

# Question 🤔 Answer at q.dsc40a.com

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

- A. yes

- B. no
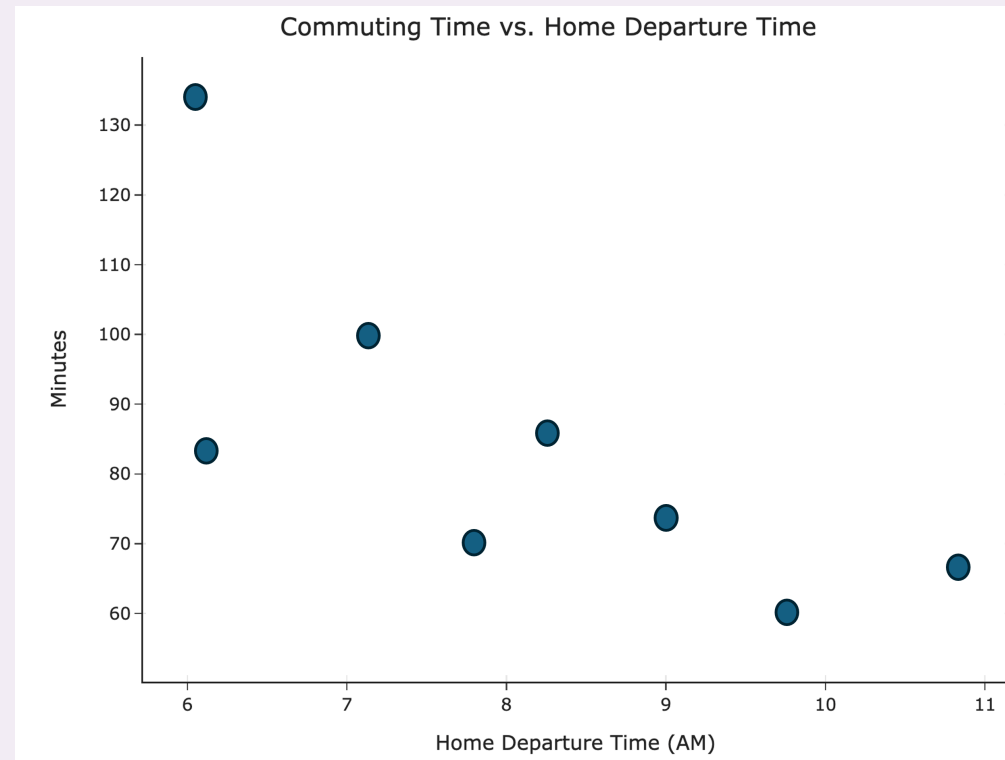


Commuting Time vs. Home Departure Time

# Question 🤔 Answer at q.dsc40a.com

Given the data below, is there a prediction rule H which has zero mean absolute error?
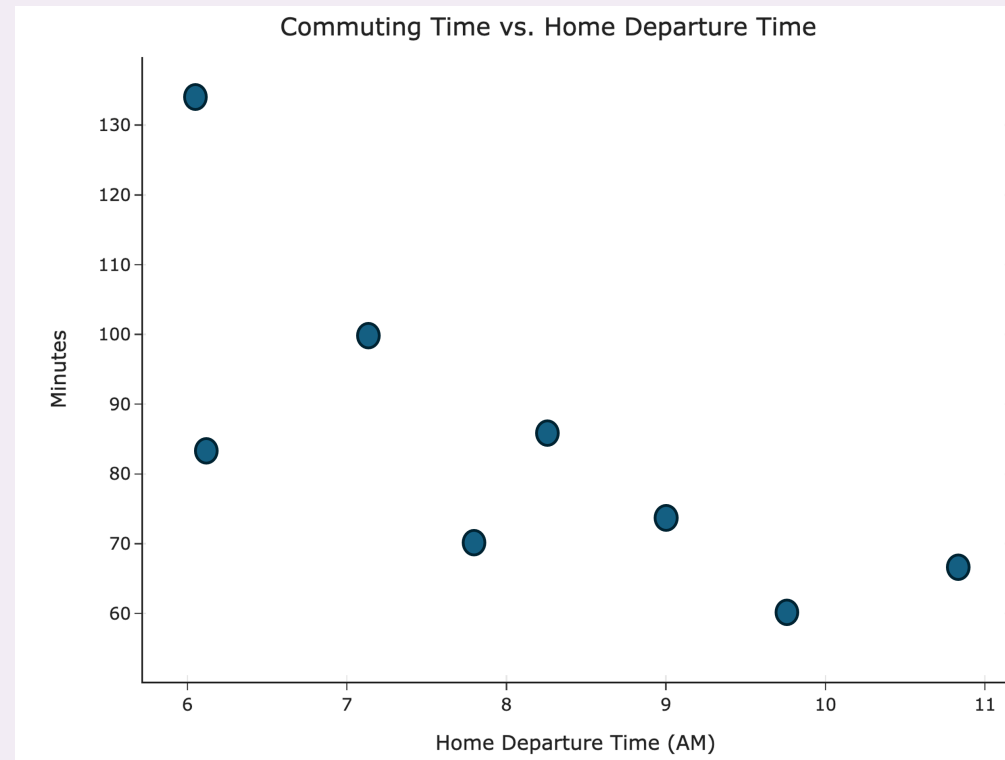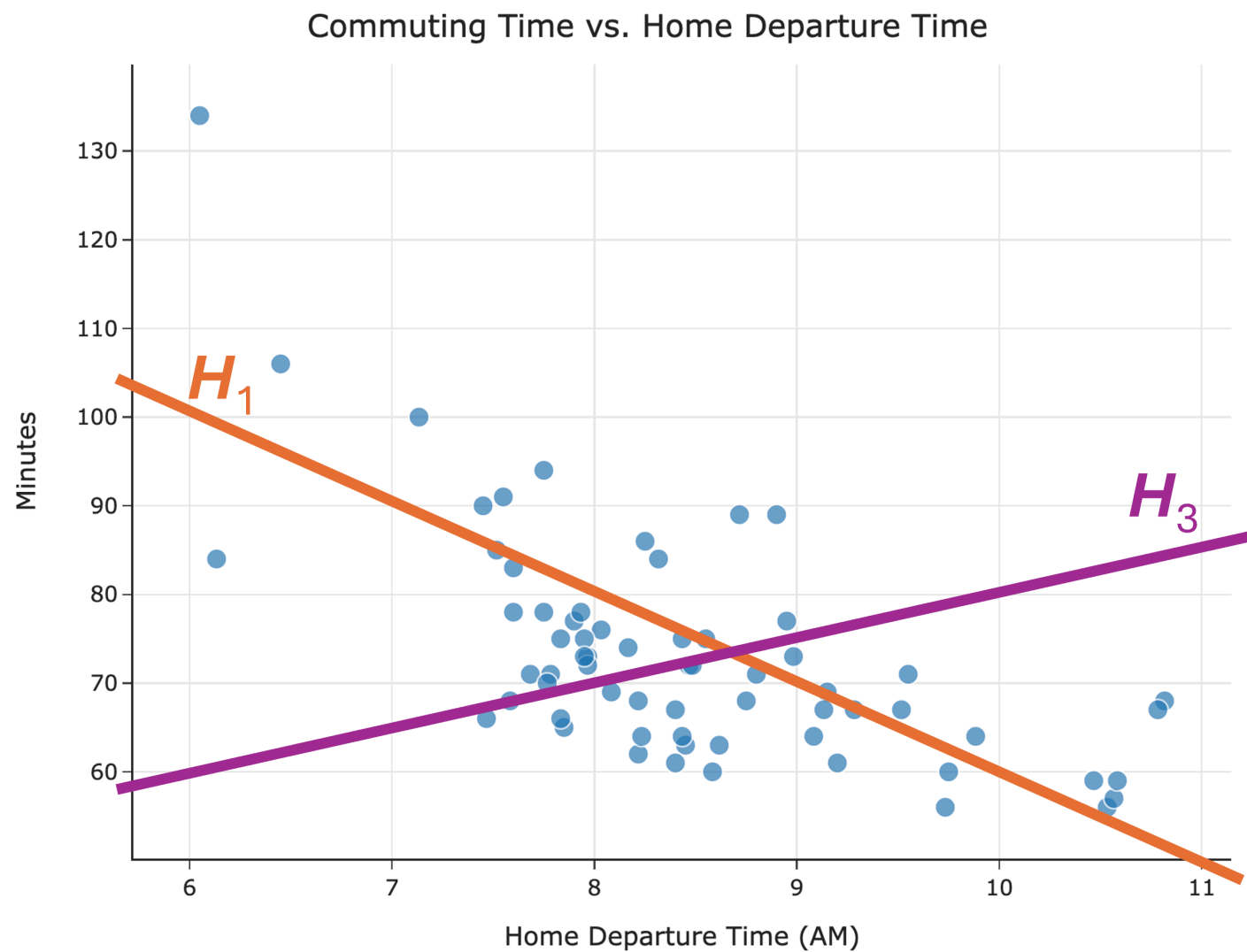
- A. yes

- B. no

# Problem

- We can make mean absolute 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: $H(x) = w_0 + w_1 x$.
  - Quadratic: $H(x) = w_0 + w_1 x_1 + w_2 x^2$.
  - Exponential: $H(x) = w_0 e^{w_1 x}$.
  - Constant: $H(x) = w_0$.

# Comparing predictions


Commuting Time vs. Home Departure Time

# 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 $H(x) = w_0 + w_1 x$, we can re-write $R_{\text{sq}}$ as a function of $w_0$ and $w_1$:

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

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