**Lecture 5**

# More Simple Linear Regression

**DSC 40A, Spring 2024**

# Announcements

- Homework 2 is due on **Thursday**. Remember that using the Overleaf template is required for Homework 2 (and only Homework 2).

- Homework 1, Groupwork 1, and Groupwork 2 solutions are all available on Ed.

- Check out the new FAQs page and the tutor-created supplemental resources on the course website.

- If you asked for an alternate Final Exam and/or have OSD accommodations, you should've received an email from me a few days ago with the details of your Final Exam arrangement.

- You can access the Markdown source code for lectures here (potentially useful if you want to write your own notes).

# Agenda

- Recap: Simple linear regression.

- Correlation.

- Interpreting the formulas.

- Connections to related models.

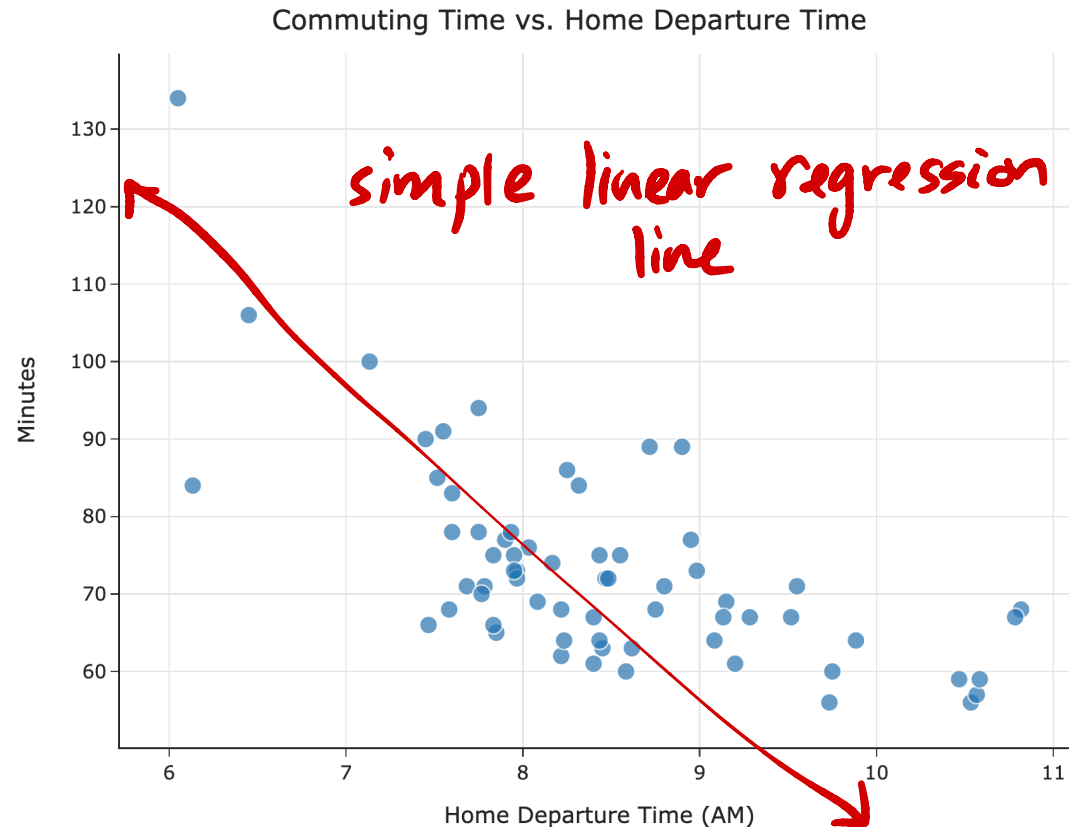- Introduction to linear algebra.

# 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.

# Recap: Simple linear regression

# Recap



Commuting Time vs. Home Departure Time

*simple linear regression line*

- In Lecture 4, our goal was to fit a **simple linear regression** model, $H(x) = w_0 + w_1 x$, to our commute times dataset.
  - $x_i$: The $i$th home departure time (e.g. 8.5, for 8:30 AM).
  - $y_i$: The $i$th actual commute time (e.g. 76 minutes).
  - $H(x_i)$: The $i$th predicted commute time.
- To do so, we used squared loss.

# The modeling recipe

1. Choose a model.

$$H(x) = w_0 + w_1 x$$

intercept

slope

2. Choose a loss function.

$$L_{sq}(y_i, H(x_i)) = (y_i - H(x_i))^2$$

$(actual - predicted)^2$

3. Minimize average loss to find optimal model parameters.

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

# Least squares solutions

- Our goal was to find the parameters $w_0^*$ and $w_1^*$ that minimized:

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

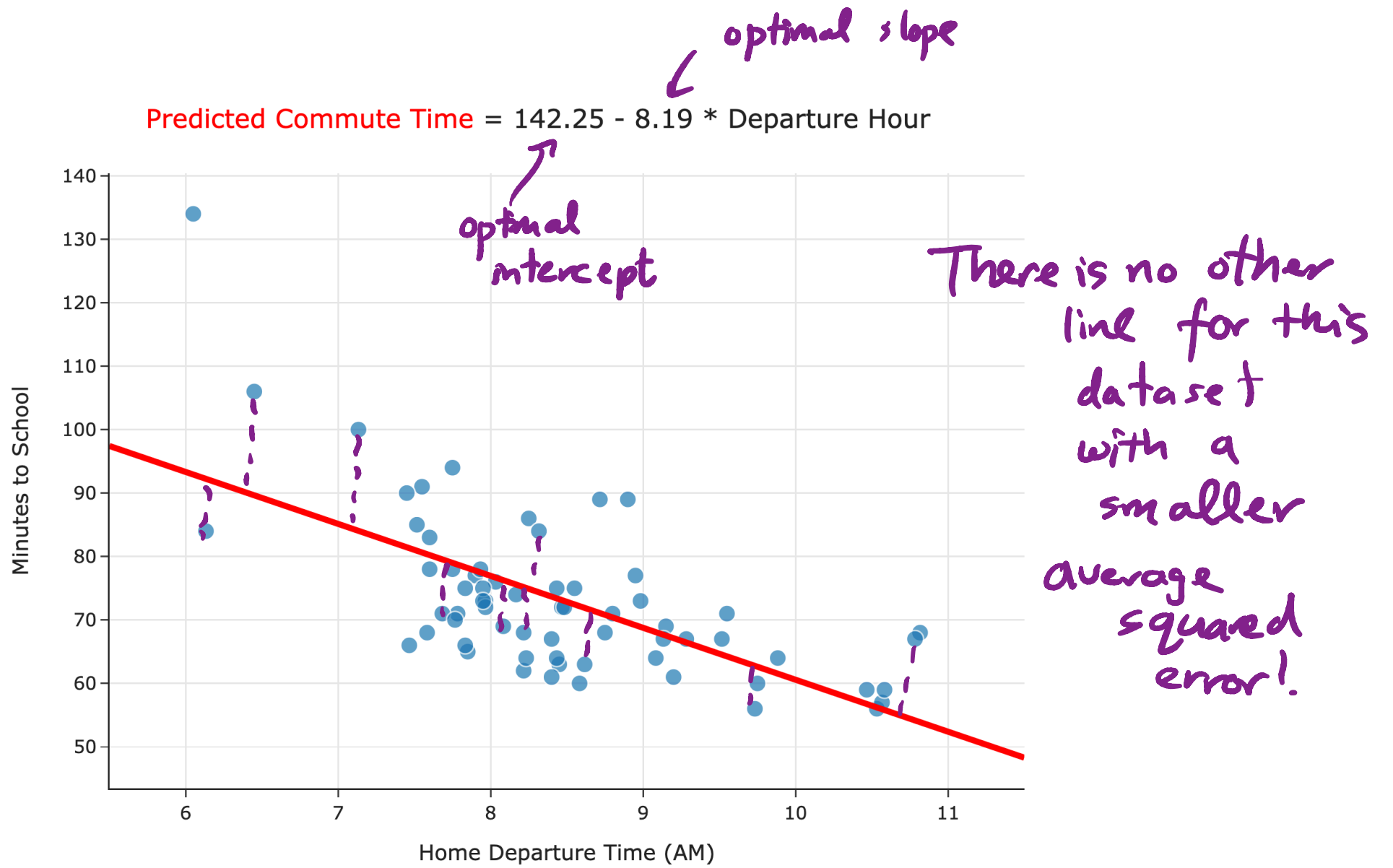- To do so, we used calculus, and we found that the minimizing values are:

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

*best slope* → $w_1^*$

*best intercept* → $w_0^*$

- We say $w_0^*$ and $w_1^*$ are **optimal parameters,** and the resulting line is called the **regression line**.

# Now what?

We've found the optimal slope and intercept for linear hypothesis functions using squared loss (i.e. for the regression line). Now, we'll:
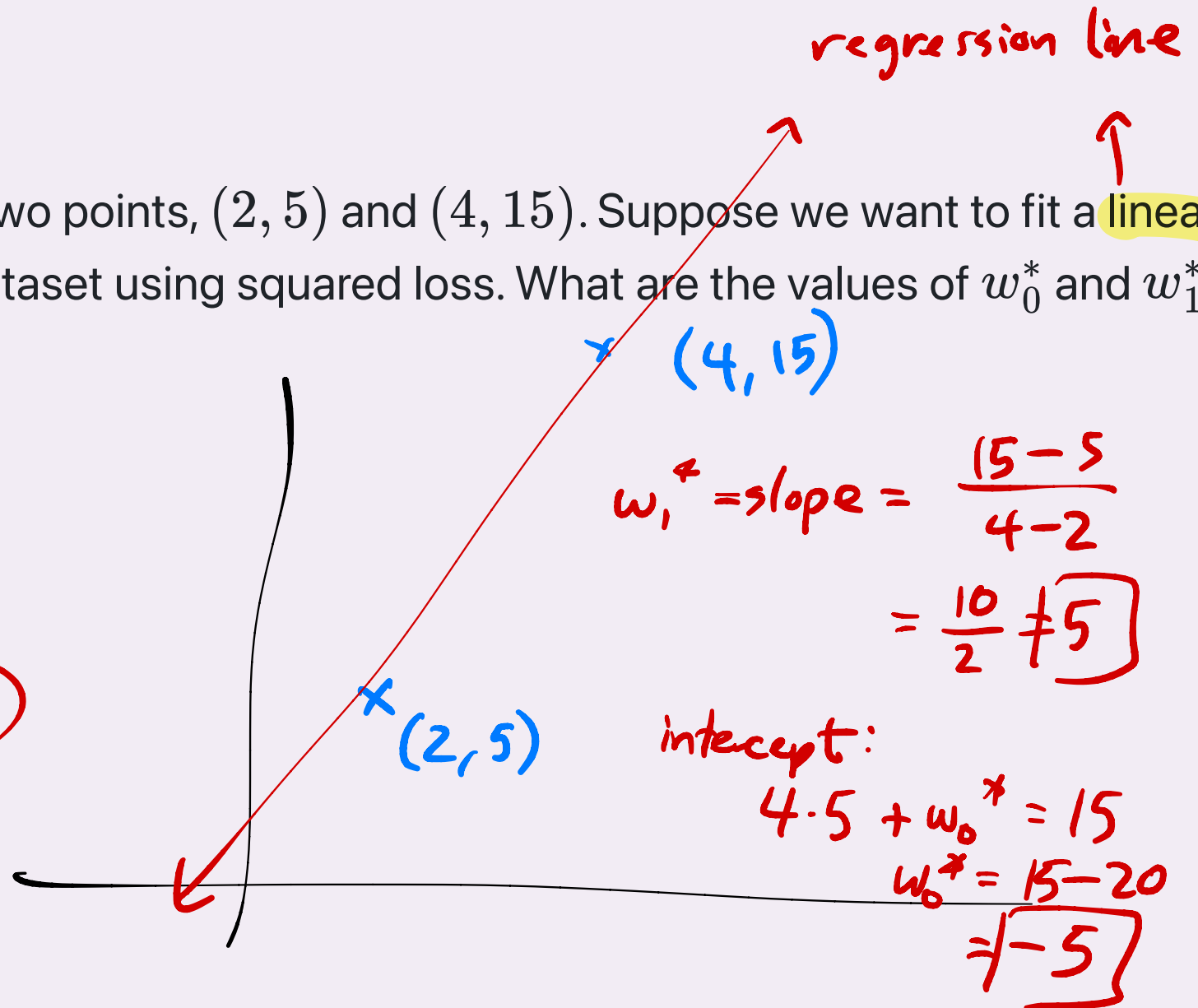
- See how the formulas we just derived connect to the formulas for the slope and intercept of the regression line we saw in DSC 10.
  - They're the same, but we need to do a bit of work to prove that.
- Learn how to interpret the slope of the regression line.
- Understand connections to other related models.
- Learn how to build regression models with **multiple inputs**.
  - To do this, we'll need linear algebra!

# Question 🤔

Consider a dataset with just two points, $(2, 5)$ and $(4, 15)$. Suppose we want to fit a <mark>linear hypothesis function</mark> to this dataset using squared loss. What are the values of $w_0^*$ and $w_1^*$ that minimize empirical risk?

- A. $w_0^* = 2, w_1^* = 5$
- B. ~~$w_0^* = 3, w_1^* = 10$~~
- C. $w_0^* = -2, w_1^* = 5$
- D. $w_0^* = -5, w_1^* = 5$

regression line

$(4, 15)$

$w_1^* = \text{slope} = \dfrac{15 - 5}{4 - 2}$

$= \dfrac{10}{2} = 5$

$(2, 5)$

intercept:

$4 \cdot 5 + w_0^* = 15$

$w_0^* = 15 - 20$

$= -5$

11

# Correlation
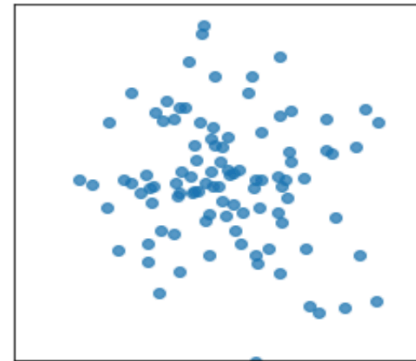
association : any pattern

correlation : <u>linear association</u>
pattern that looks like a line!

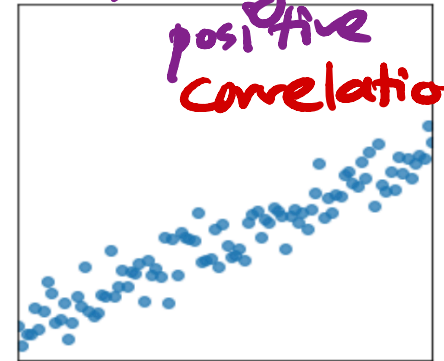## Quantifying patterns in scatter plots

- In DSC 10, you were introduced to the idea of the **correlation coefficient**, $r$.

- It is a measure of the strength of the **linear association** of two variables, $x$ and $y$.

- Intuitively, it measures how tightly clustered a scatter plot is around a straight line.
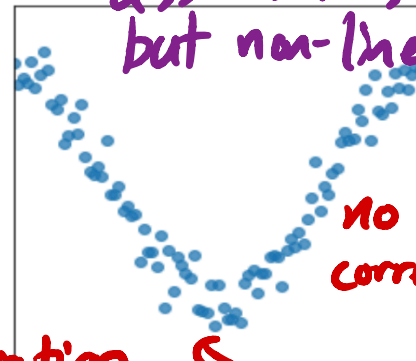
- It ranges between –1 and 1.

no association!

strong positive correlation

association, but non-linear

no correlation

positive correlation

weaker because of spread

- $r$ negative : negative linear association
- $r$ positive : positive linear association
- the closer $r$ is to $\pm 1$, the stronger the correlation!

13

## The correlation coefficient

*annotations:* Pearson's correlation coefficient — there are others

- The correlation coefficient, $r$, is defined as the **average of the product of** $x$ **and** $y$, **when both are in standard units**.

- Let $\sigma_x$ be the standard deviation of the $x_i$s, and $\bar{x}$ be the mean of the $x_i$s.

- $x_i$ in standard units is $\frac{x_i - \bar{x}}{\sigma_x}$.

  *annotations:* $\frac{\text{value} - \text{mean}}{\text{SD}}$ : measures the number of SDs above/below the mean

- The correlation coefficient, then, is:

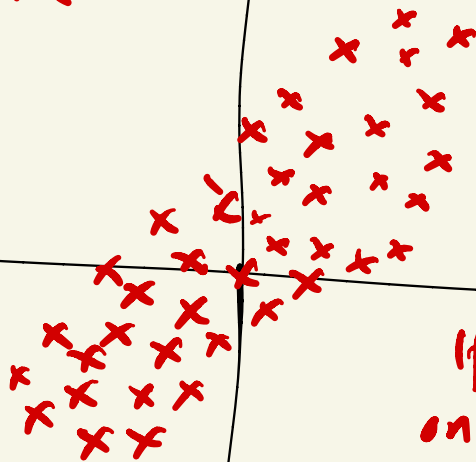$$r = \frac{1}{n}\sum_{i=1}^{n}\left(\frac{x_i - \bar{x}}{\sigma_x}\right)\times\left(\frac{y_i - \bar{y}}{\sigma_y}\right)$$

*annotations:* "sigma x"

*annotations:* average | $x_i$ in standard units | $y_i$ in standard units

# Question : Why multiply the SUs when calculating r ?

Suppose there's positive correlation.

Most points are in the top right and bottom left.

Top right :
$x_i$ (su) positive and
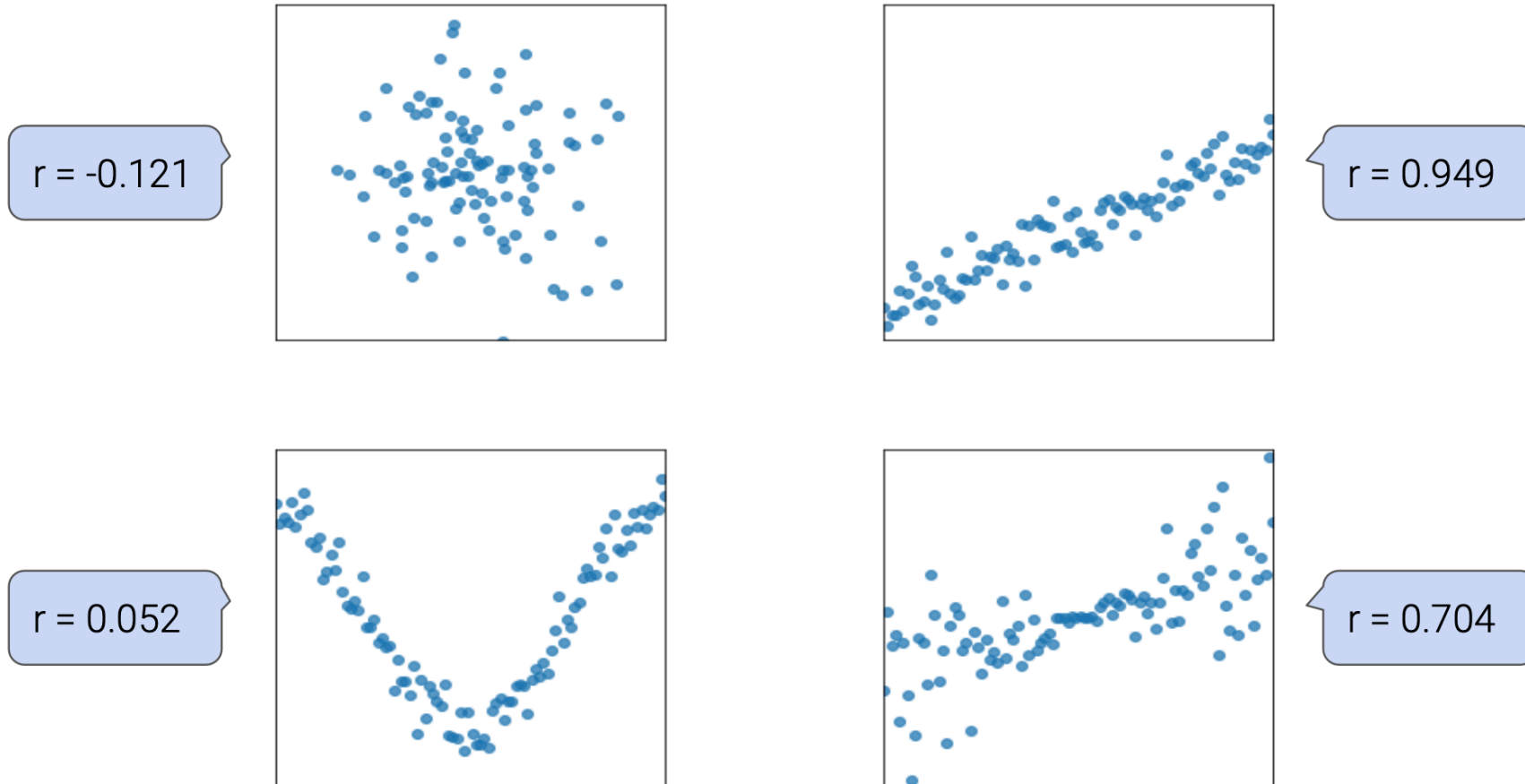$y_i$ (su) positive.

Bottom left :
$x_i$ (su) and
$y_i$ (su) both
negative



If there's positive <u>correlation</u>, on average, both
$\Rightarrow$ $x_i$ (su) and $y_i$ (su) will have
the same sign. $+ \cdot + = +$
$- \cdot - = +$
$\Rightarrow$ on average, $x_i$ (su) $\cdot$ $y_i$ (su) is positive!

# The correlation coefficient, visualized

r = -0.121

r = 0.949

r = 0.052

r = 0.704

# Another way to express $w_1^*$

- It turns out that $w_1^*$, the optimal slope for the linear hypothesis function when using squared loss (i.e. the regression line), can be written in terms of $r$!

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

- It's not surprising that $r$ is related to $w_1^*$, since $r$ is a measure of linear association.

- Concise way of writing $w_0^*$ and $w_1^*$:

$$w_1^* = r\frac{\sigma_y}{\sigma_x} \qquad w_0^* = \bar{y} - w_1^*\bar{x}$$

**Proof that** $w_1^* = r\dfrac{\sigma_y}{\sigma_x}$

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

$$= \frac{r \, n \, \sigma_x \, \sigma_y \quad ①}{n \, \sigma_x^2 \quad ②}$$

$$= \boxed{r \, \frac{\sigma_y}{\sigma_x}} \quad \text{done!}$$

$$r = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{x_i - \bar{x}}{\sigma_x} \right) \left( \frac{y_i - \bar{y}}{\sigma_y} \right)$$

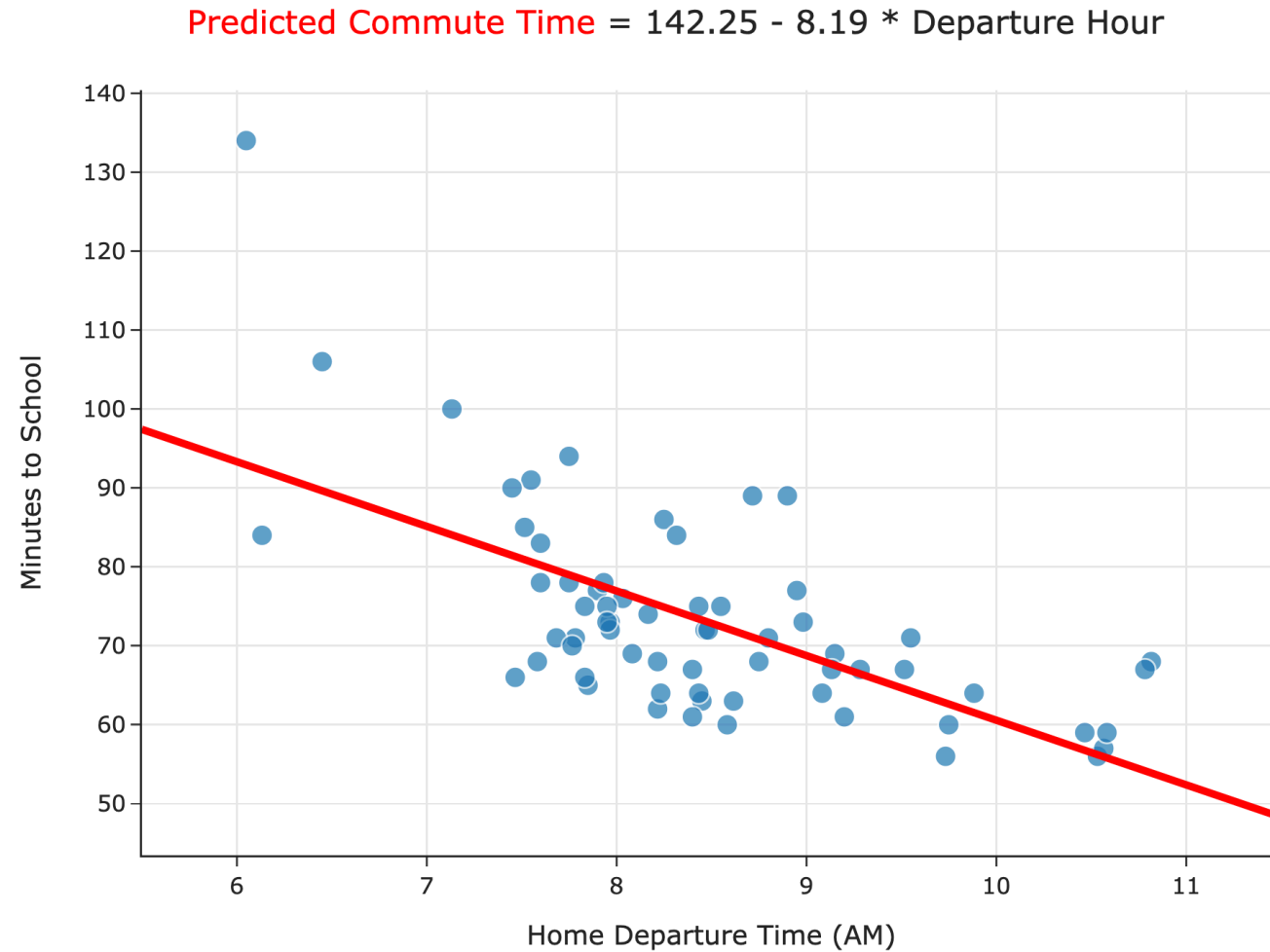$$r = \frac{1}{n \sigma_x \sigma_y} \sum_{i=1}^{\hat{}} (x_i - \bar{x})(y_i - \bar{y})$$

$$\Rightarrow \boxed{r \, n \, \sigma_x \sigma_y = \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})} \quad ①$$

$$\sigma_x = \sqrt{\frac{1}{n} \sum_{i=1}^{\hat{}} (x_i - \bar{x})^2}$$

$$\sigma_x^2 = \frac{1}{n} \sum_{i=1}^{\hat{}} (x_i - \bar{x})^2$$

$$\Rightarrow \boxed{n \sigma_x^2 = \sum_{i=1}^{\hat{}} (x_i - \bar{x})^2} \quad ②$$

17

Let's test these new formulas out in code! Follow along here.



Predicted Commute Time = 142.25 - 8.19 * Departure Hour

# Interpreting the formulas

# Interpreting the slope

no units!

$$w_1^* = r\frac{\sigma_y}{\sigma_x}$$

→ units of $y$

→ units of $x$

- The units of the slope are **units of $y$ per units of $x$**.

- In our commute times example, in $H(x) = 142.25 - 8.19x$, our predicted commute time **decreases by 8.19 minutes per hour**.
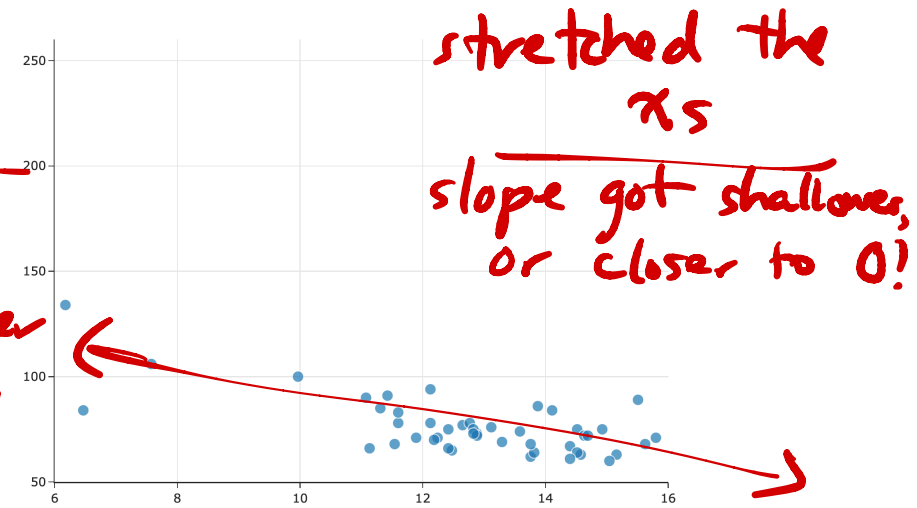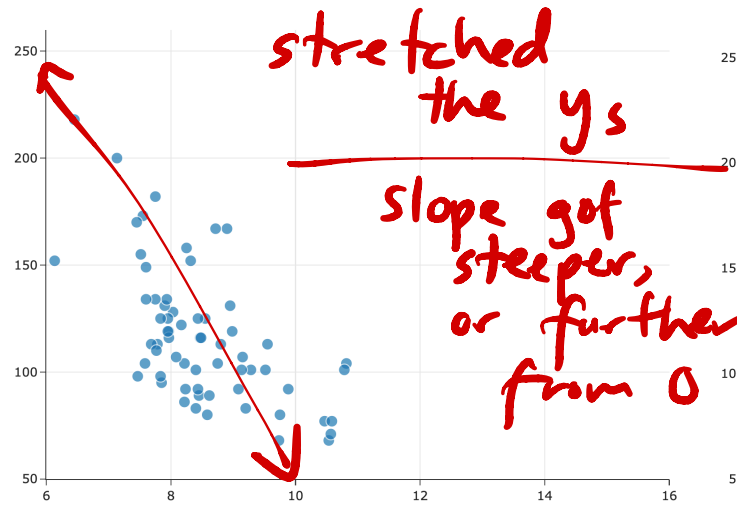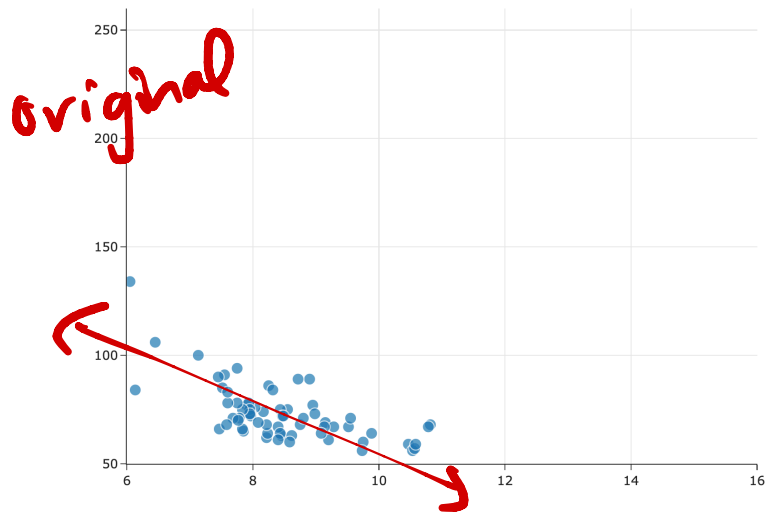
$x_i$ : departure time in **hours**

$y_i$ : commute time in **minutes**

## Interpreting the slope

$$w_1^* = r \frac{\sigma_y}{\sigma_x}$$



original

stretched the ys

slope got steeper, or further from 0

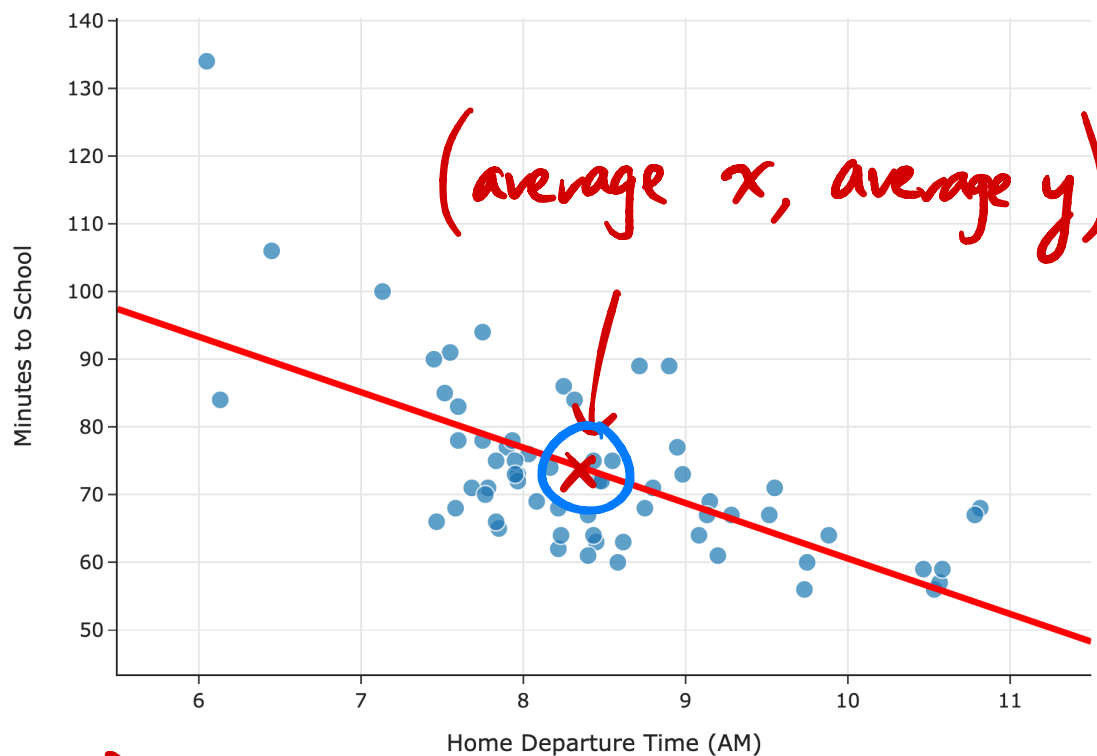stretched the xs

slope got shallower, or closer to 0!

- Since $\sigma_x \geq 0$ and $\sigma_y \geq 0$, the slope's sign is $r$'s sign.

- As the $y$ values get more spread out, $\sigma_y$ increases, so the slope gets steeper.

- As the $x$ values get ~~less~~ more spread out, $\sigma_x$ increases, so the slope gets shallower.

21

# Interpreting the intercept

$$w_0^* = \bar{y} - w_1^* \bar{x}$$

Predicted Commute Time = 142.25 - 8.19 * Departure Hour



(average x, average y)

$H(0)$ = intercept
= predicted commute time
@ midnight

- What are the units of the intercept?

units of $y$: minutes

- What is the value of $H^*(\bar{x})$?

$$H^*(x_i) = w_0^* + w_1^* x_i$$
$$= \bar{y} - w_1^* \bar{x} + w_1^* x_i$$
$$= \bar{y} + w_1^* (x_i - \bar{x})$$

$$H^*(\bar{x}) = \bar{y} + w_1^* (\bar{x} - \bar{x}) \rightarrow 0$$
$$= \bar{y}$$

22

# Question 🤔

We fit a regression line to predict commute times given departure hour. Then, we add 75 minutes to all commute times in our dataset. What happens to the resulting regression line?

- A. Slope increases, intercept increases.
- B. Slope decreases, intercept increases.
- C. Slope stays the same, intercept increases.
- D. Slope stays the same, intercept stays the same.

*intercept increases by 75!*

commute time

departure hour

23

# Correlation and mean squared error

- **Claim**: Suppose that $w_0^*$ and $w_1^*$ are the optimal intercept and slope for the regression line. Then,

$$R_{\text{sq}}(w_0^*, w_1^*) = \sigma_y^2(1 - r^2)$$

- That is, the mean squared error of the regression line's predictions and the correlation coefficient, $r$, always satisfy the relationship above.

- Even if it's true, why do we care?

  - In machine learning, we often use both the mean squared error and $r^2$ to compare the performances of different models.

  - If we can prove the above statement, we can show that **finding models that minimize mean squared error** is equivalent to **finding models that maximize $r^2$**
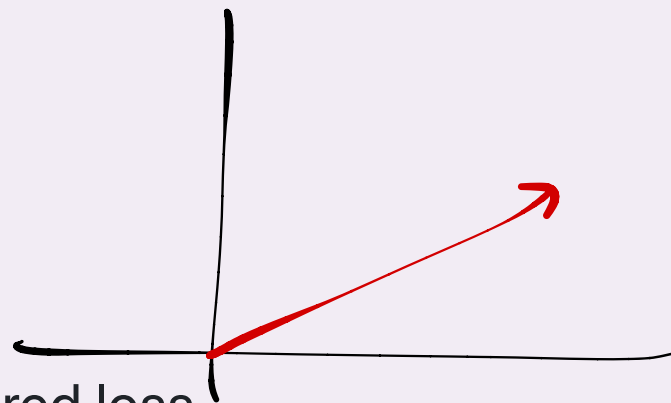
.

**Proof that** $R_{\mathrm{sq}}(w_0^*, w_1^*) = \sigma_y^2(1 - r^2)$

# Connections to related models

# Question 🤔

**Answer at q.dsc40a.com**

Suppose we chose the model $H(x) = w_1 x$ and squared loss.

What is the optimal model parameter, $w_1^*$?

*a line forced through (0,0)*

- A. $\dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$

- B. $\dfrac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2}$

- C. $\dfrac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2}$

- D. $\dfrac{\sum_{i=1}^{n} y_i}{\sum_{i=1}^{n} x_i}$

*the same! (and correct!)*

# Exercise

Suppose we chose the model $H(x) = w_1 x$ and squared loss.

What is the optimal model parameter, $w_1^*$?

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

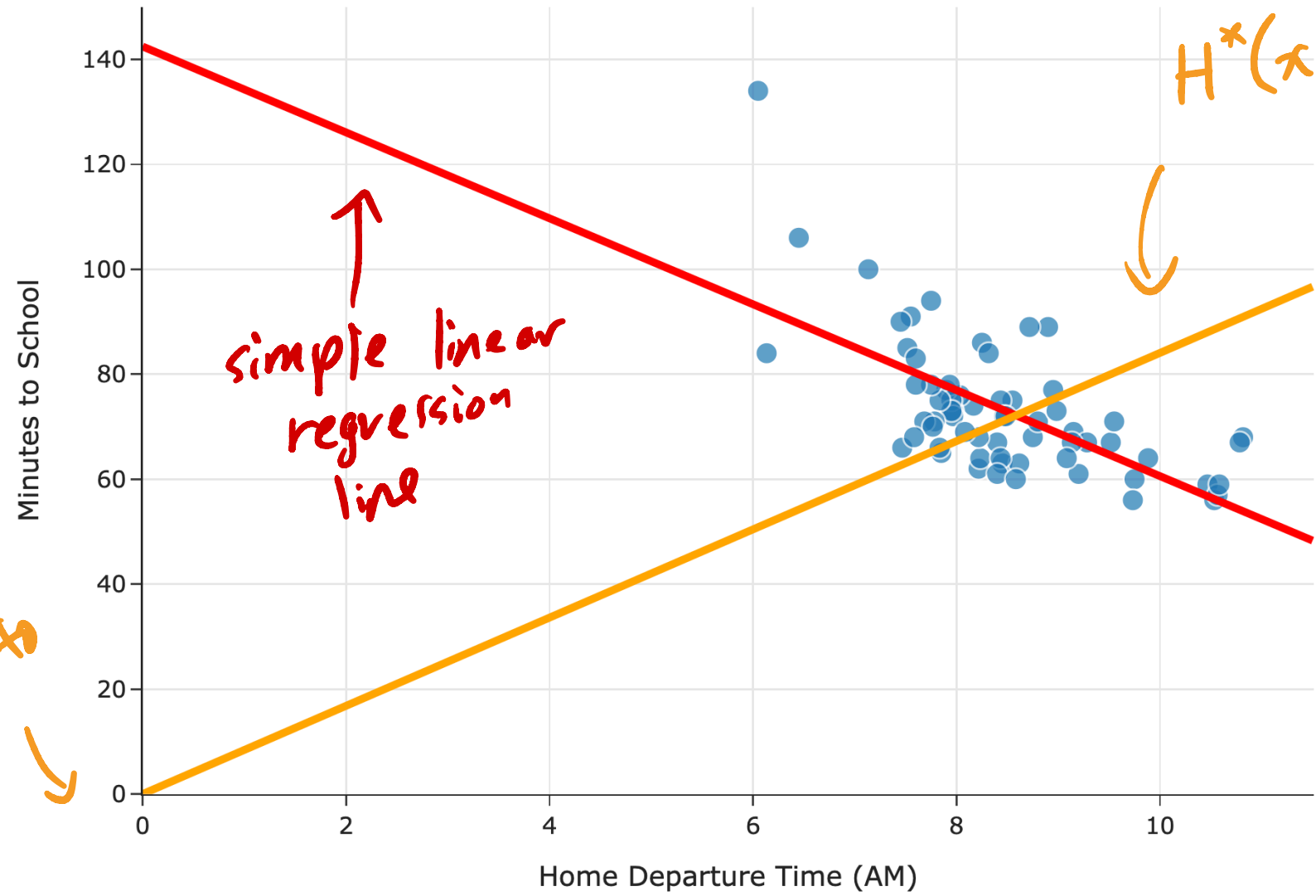$$\frac{dR_{sq}}{dw_1} = \frac{1}{n} \sum_{i=1}^{n} 2(y_i - w_1 x_i)(-x_i)$$

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

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

$$w_1^* = \frac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2}$$

Predicted Commute Time = 142.25 - 8.19 * Departure Hour
Predicted Commute Time = 8.41 * Departure Hour

$H^*(x) = W_1^* x$

simple linear regression line

intercept forced to be 0

# Exercise

Suppose we choose the model $H(x) = w_0$ and squared loss.

What is the optimal model parameter, $w_0^*$?

$$w_0^* = \text{Mean}(y_1, y_2, \ldots, y_n)$$

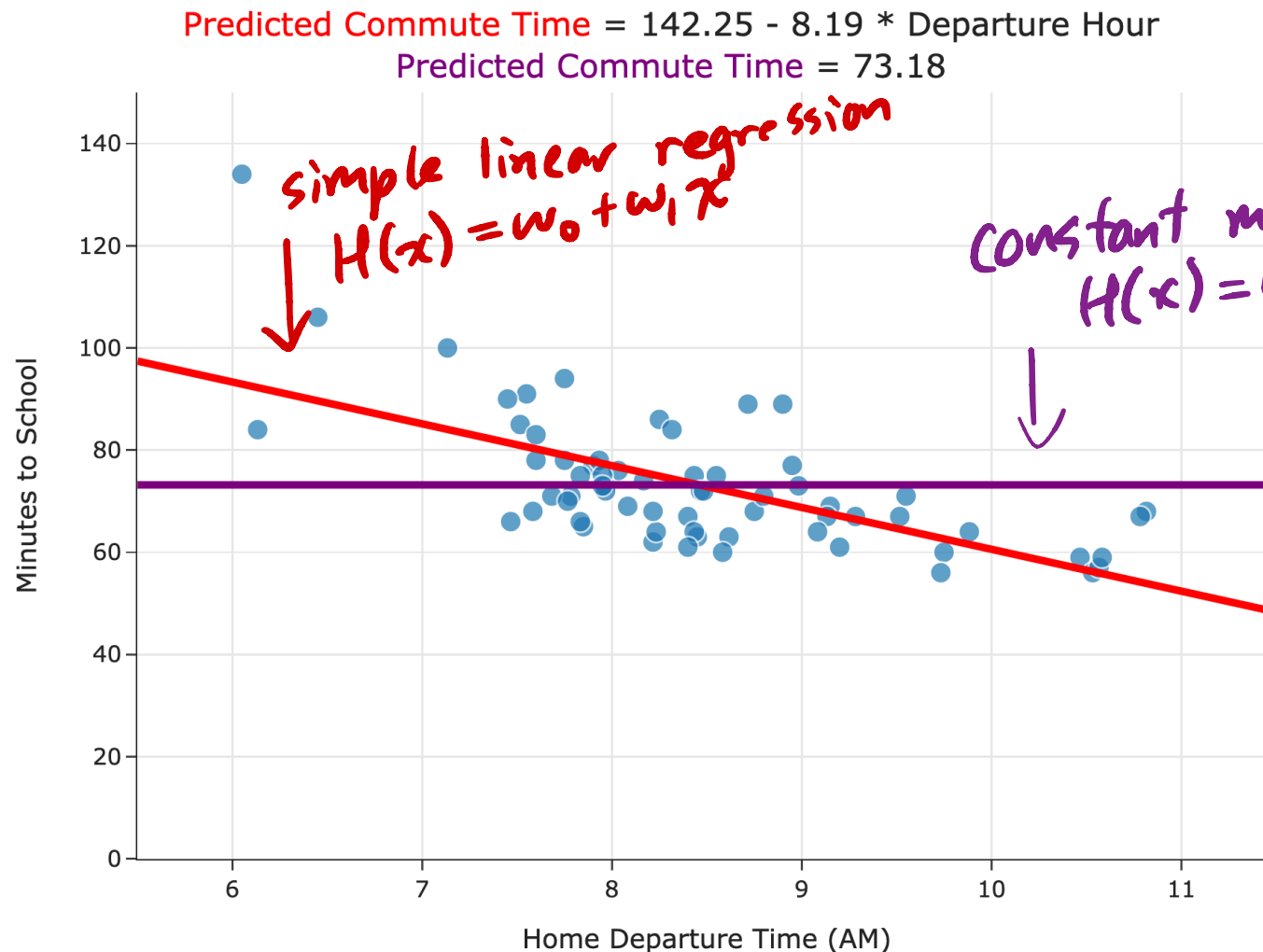# Comparing mean squared errors

- With both:

  - the constant model, $H(x) = h$, and

  - the simple linear regression model, $H(x) = w_0 + w_1 x$,

  when we chose squared loss, we minimized mean squared error to find optimal parameters:

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

- **Which model minimizes mean squared error more?**

# Comparing mean squared errors

Predicted Commute Time = 142.25 - 8.19 * Departure Hour
Predicted Commute Time = 73.18



$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - H(x_i) \right)^2$$

- The MSE of the best simple linear regression model is $\approx 97$.

- The MSE of the best constant model is $\approx 167$.

- The simple linear regression model is a more flexible version of the constant model.
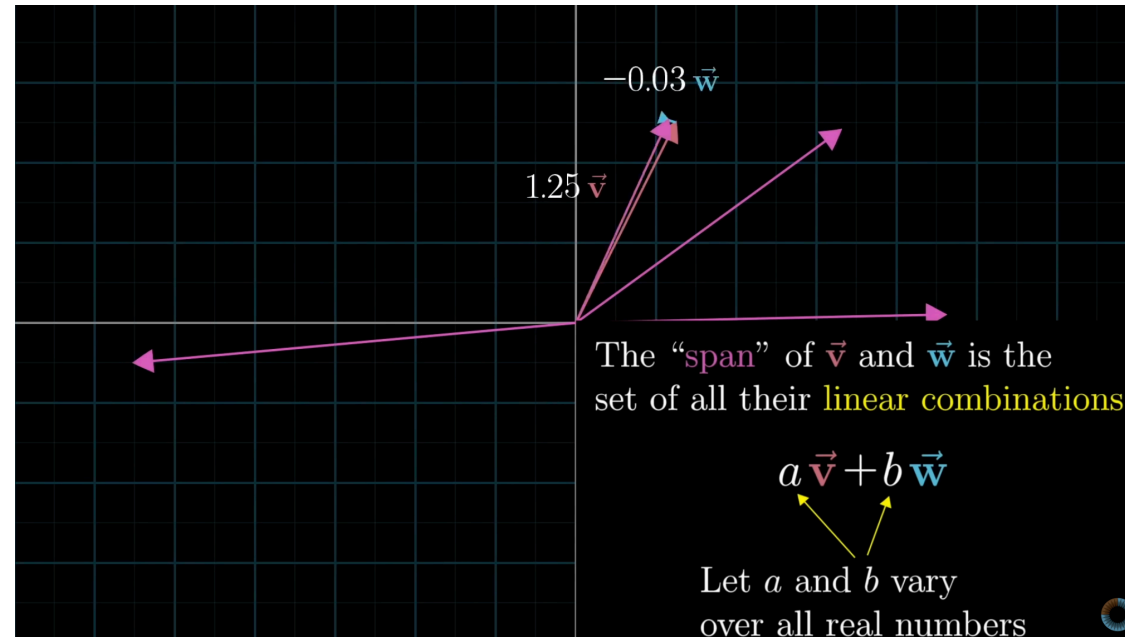
# Linear algebra review

# Wait... why do we need linear algebra?

- Soon, we'll want to make predictions using more than one feature.
  - Example: Predicting commute times using departure hour and temperature.
- Thinking about linear regression in terms of **matrices and vectors** will allow us to find hypothesis functions that:
  - Use multiple features (input variables).
  - Are non-linear, e.g. $H(x) = w_0 + w_1 x + w_2 x^2$.
- Before we dive in, let's review.

# Spans of vectors

- One of the most important ideas you'll need to remember from linear algebra is the concept of the **span** of two or more vectors.

- To jump start our review of linear algebra, let's start by watching 🎥 **this video by 3blue1brown**.

# Next time

- We'll review the necessary linear algebra prerequisites.

- We'll then start to formulate the problem of minimizing mean squared error for the simple linear regression model **using matrices and vectors**.

- We'll send some relevant linear algebra review videos on Ed.