**Lectures 6-7**

# Simple Linear Regression

**DSC 40A, Fall 2024**

# Agenda

- Simple linear regression.

  *Least squares solution*

- Correlation.

- Interpreting the formulas.

- Connections to related models.

− HW 1 due tonight, HW 2 released

− Submit regrade requests — no need for emails

# Least squares solutions

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

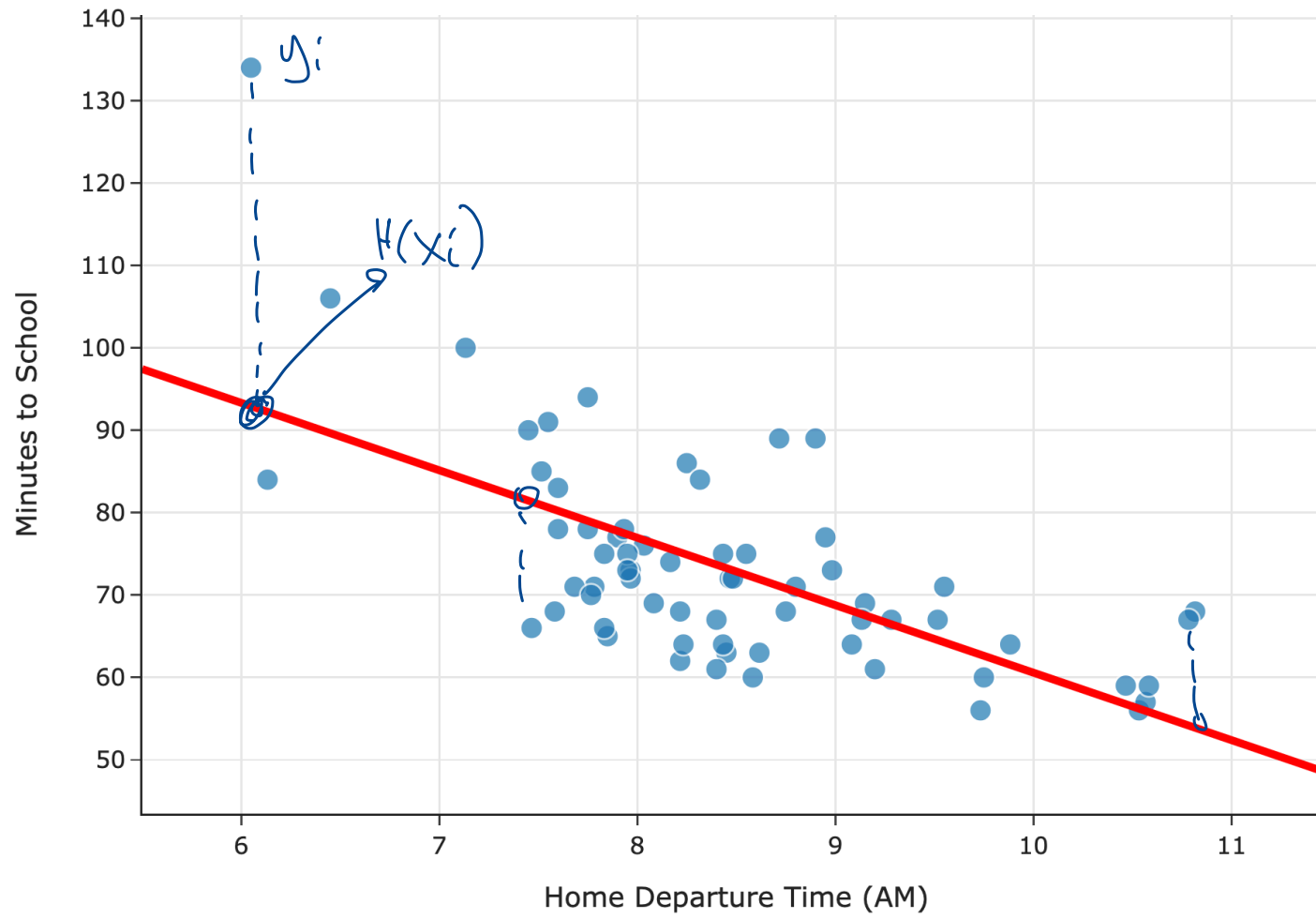$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^{n} (y_i - (w_0 + w_1 x_i))^2$$

_（handwritten annotation: $H(x_i')$）_

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

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

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

optimal intercept

optimal slope

Predicted Commute Time = 142.25 - 8.19 * Departure Hour

$y_i$

$H(x_i)$

There is no other line for this dataset with smaller MSE

$$W_0^*, W_1^* = \arg\min_{W_0, W_1} MSE(H)$$

24

# 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:
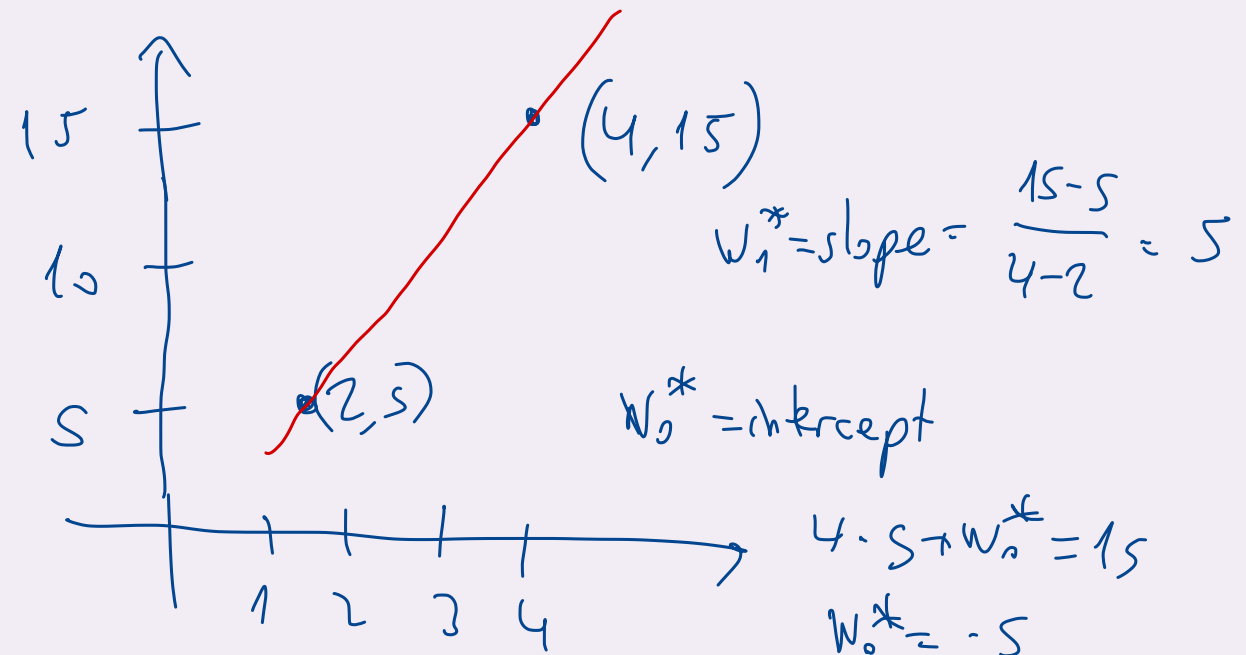
- See how the formulas we just derived connect to the formulas for the slope and intercept of the regression line we saw in DSC 10.
  - They're the same, but we need to do a bit of work to prove that.
- Learn how to interpret the slope of the regression line.
- Understand connections to other related 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 linear hypothesis function 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. $w_0^* = -5$, $w_1^* = 5$



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

$$w_0^* = intercept$$

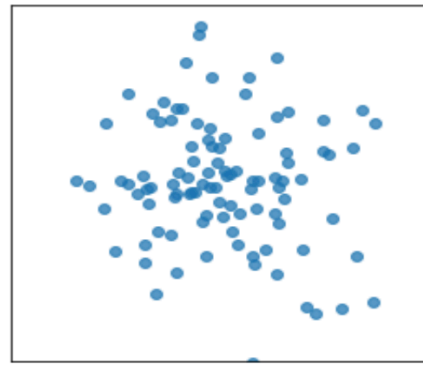$$4 \cdot 5 + w_0^* = 15$$

$$w_0^* = -5$$

26

# Correlation

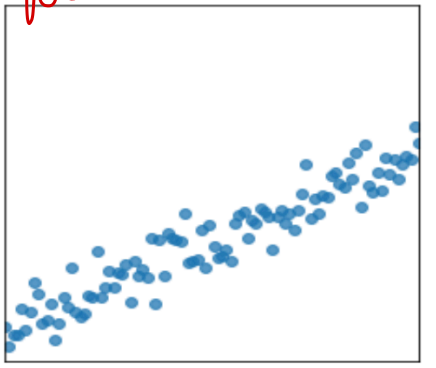# Quantifying patterns in scatter plots

*Correlation = linear association*
*does it look like a line?*

- 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.
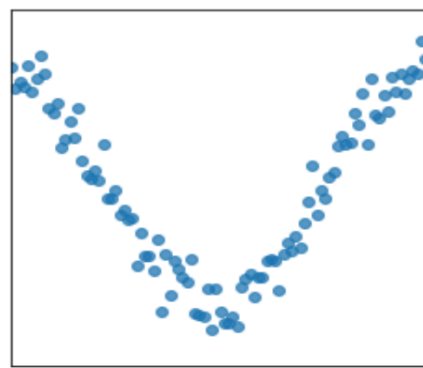
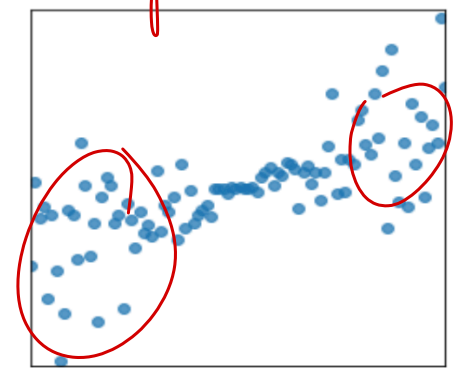- It ranges between -1 and 1.

*no association*

*strong positive correlation*

*non linear association*

*positive correlation*

*r negative : negative association*
*r positive : positive association*
*the closer r is to ±1, the stronger the correlation!*

# The correlation coefficient

*Pearson's correlation*

- 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 $\frac{x_i - \bar{x}}{\sigma_x}$.  *← mean centering*  *← std.*

- The correlation coefficient, then, is:

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

*average*  *x in standard units*  *y in standard units*

*covariance*

$$\sigma_{xy} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})$$

*alternative*
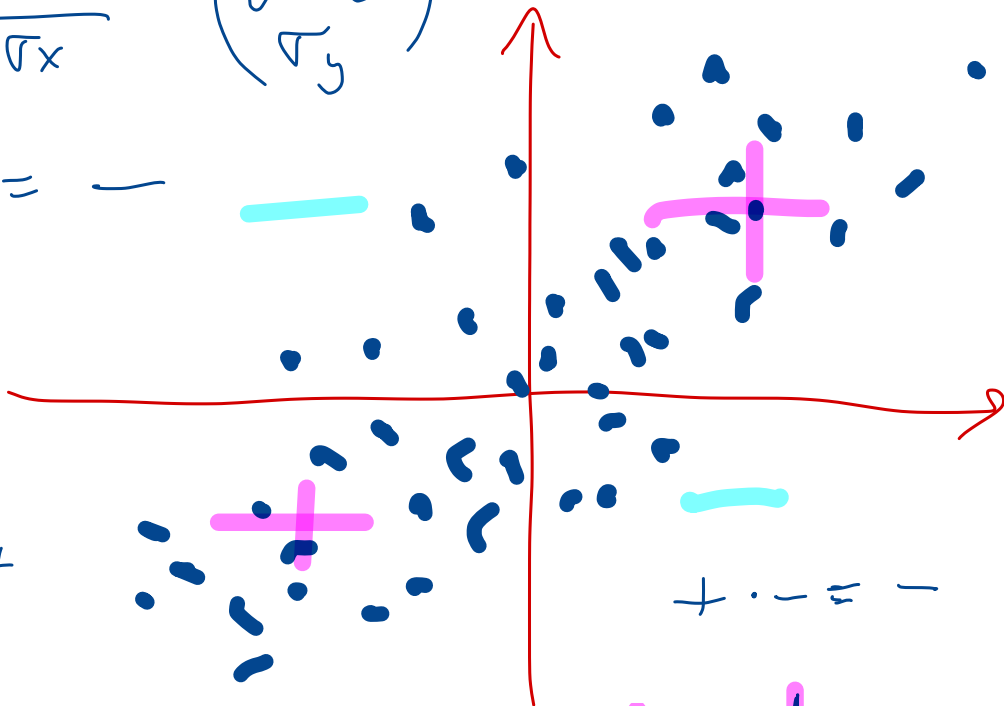
$$r = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

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

29

$$\frac{1}{n} \sum_i \frac{(X_i - \bar{X})}{\sigma_x} \cdot \left(\frac{y_i - \bar{y}}{\sigma_y}\right)$$

$$\left(\begin{array}{c} X_i \text{ in} \\ SD \end{array}\right)\left(\begin{array}{c} y_i \text{ in} \\ SD \end{array}\right)$$

$$+ \cdot +$$

$- \cdot + = -$

$- \cdot - = +$

$+ \cdot - = -$

$$r = \sum (\ ) = + \; + \; + \; + \; (-) \; + \; (-) \quad > 0$$

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

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

$$= \dfrac{n\, r\, \sigma_x\, \sigma_y}{n\, \sigma_x^2} =$$

$$= \dfrac{r\, \sigma_y}{\sigma_x}$$

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

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

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

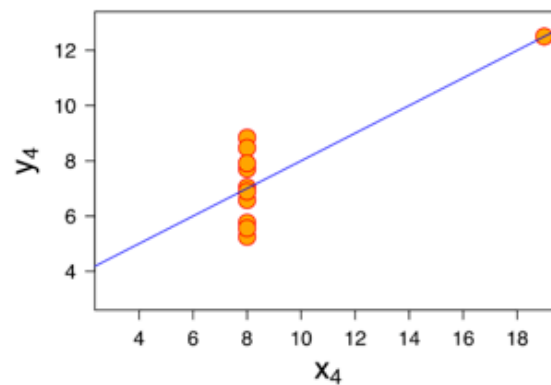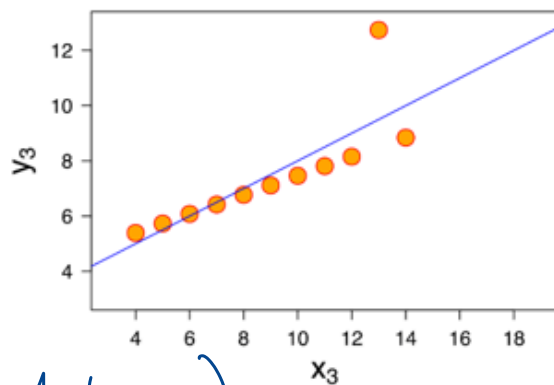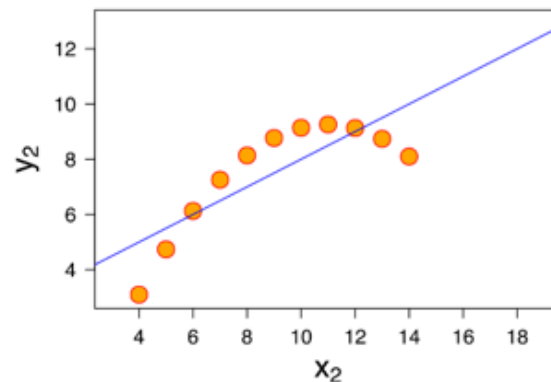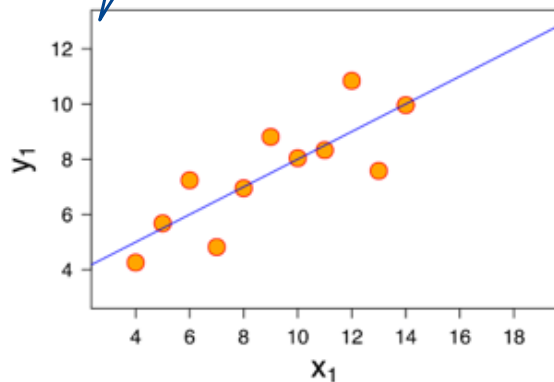$$\sigma_x = \sqrt{\dfrac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

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

$$\dfrac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})^2 = n\, \sigma_x^2$$
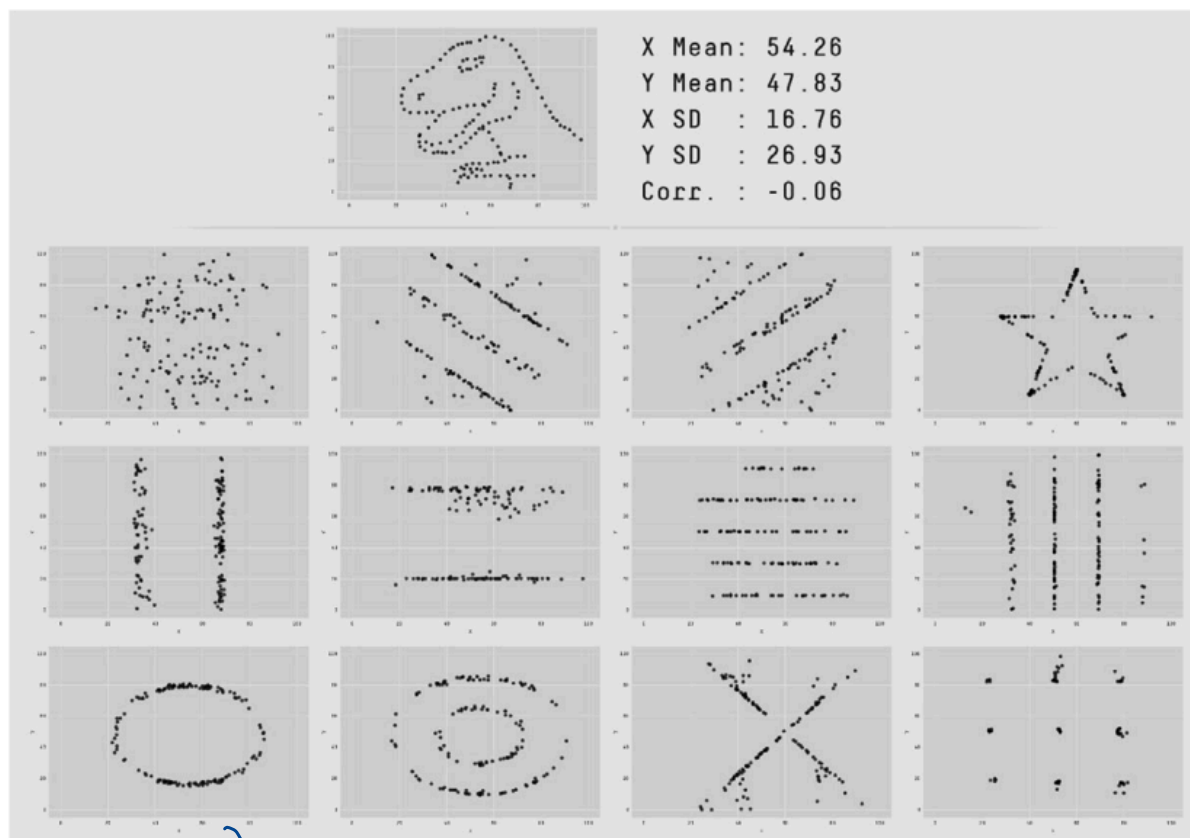
32

# Dangers of correlation

Anscombe's quartet · same mean, std, correlation



(Anscombe 1973)

# Dangers of correlation

Datasauras dozen



X Mean: 54.26
Y Mean: 47.83
X SD  : 16.76
Y SD  : 26.93
Corr. : -0.06

(Matejka et al. 2017)

# Interpreting the formulas

# Interpreting the slope

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

no units

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**.
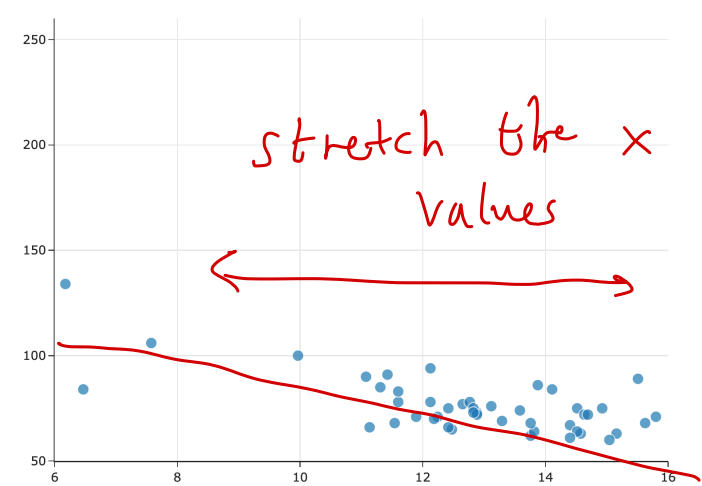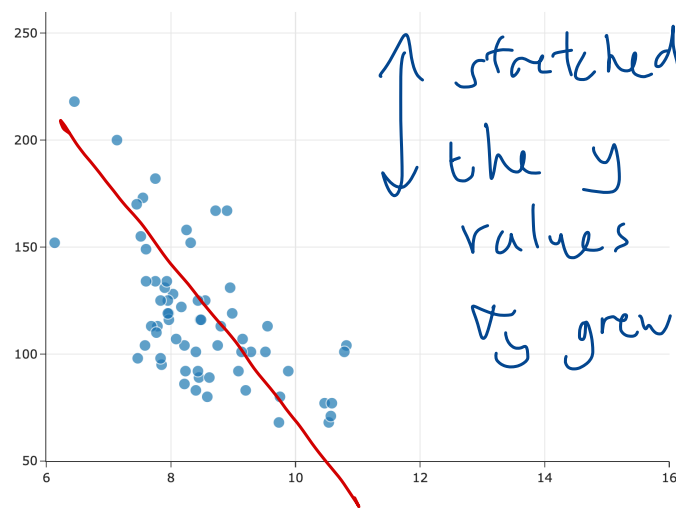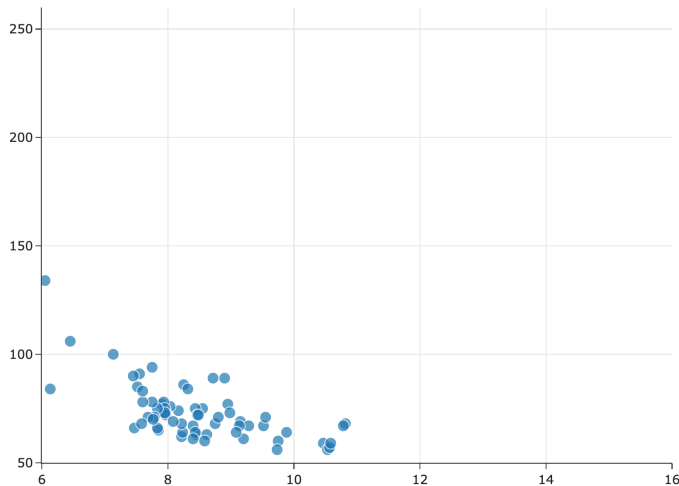
$x$    departure time   in hours

$y$   Commute time in minutes

# Interpreting the slope

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



stretched the y values → grew

stretch the x values

- 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 more spread out, $\sigma_x$ increases, so the slope gets shallower.

37

# Interpreting the intercept

$$\bar{x} = \text{mean}\{x\}$$
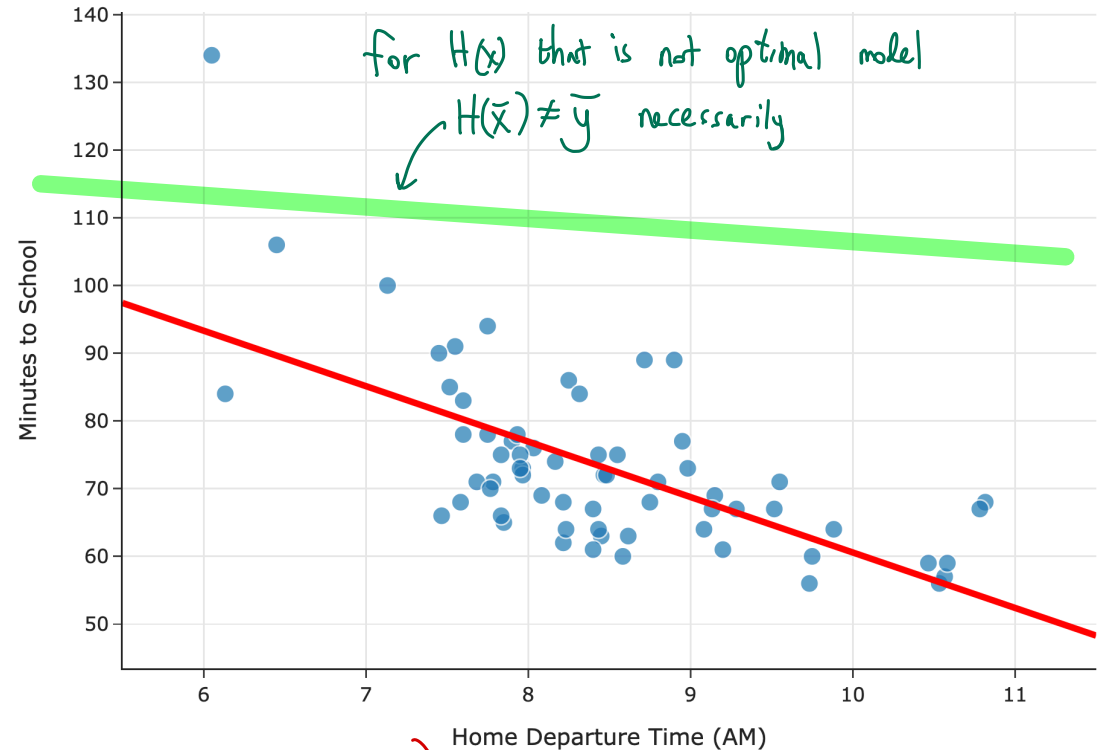$$= \frac{1}{n}\sum_i x_i$$

$$\bar{y} = \text{mean}\{y_i\}$$
$$= \frac{1}{n}\sum_i y_i$$

units of $y$ $\qquad \dfrac{\text{units of } y}{\text{units of } x} \cdot \left(\text{units of } x\right)$

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

$= \text{units of } y$

Predicted Commute Time = 142.25 - 8.19 * Departure Hour

for H(x) that is not optimal model

$H(\bar{x}) \neq \bar{y}$ necessarily

[scatter plot: Minutes to School (y-axis, 50–140) vs Home Departure Time (AM) (x-axis, 6–11), with red regression line and green line]

$H(x=0) = \text{intercept}$
$= \text{predicted commute time @ midnight}$

- What are the units of the intercept?

  units of $y$

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

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

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

$$\boxed{H^*(\bar{x}) = \bar{y}}$$

38

# 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.