

CSCI 3022: Intro to Data Science

Lecture 20: Introduction to Regression

Rachel Cox
Department of Computer
Science



If she loves you more each and every day,
by linear regression she hated you before you met.

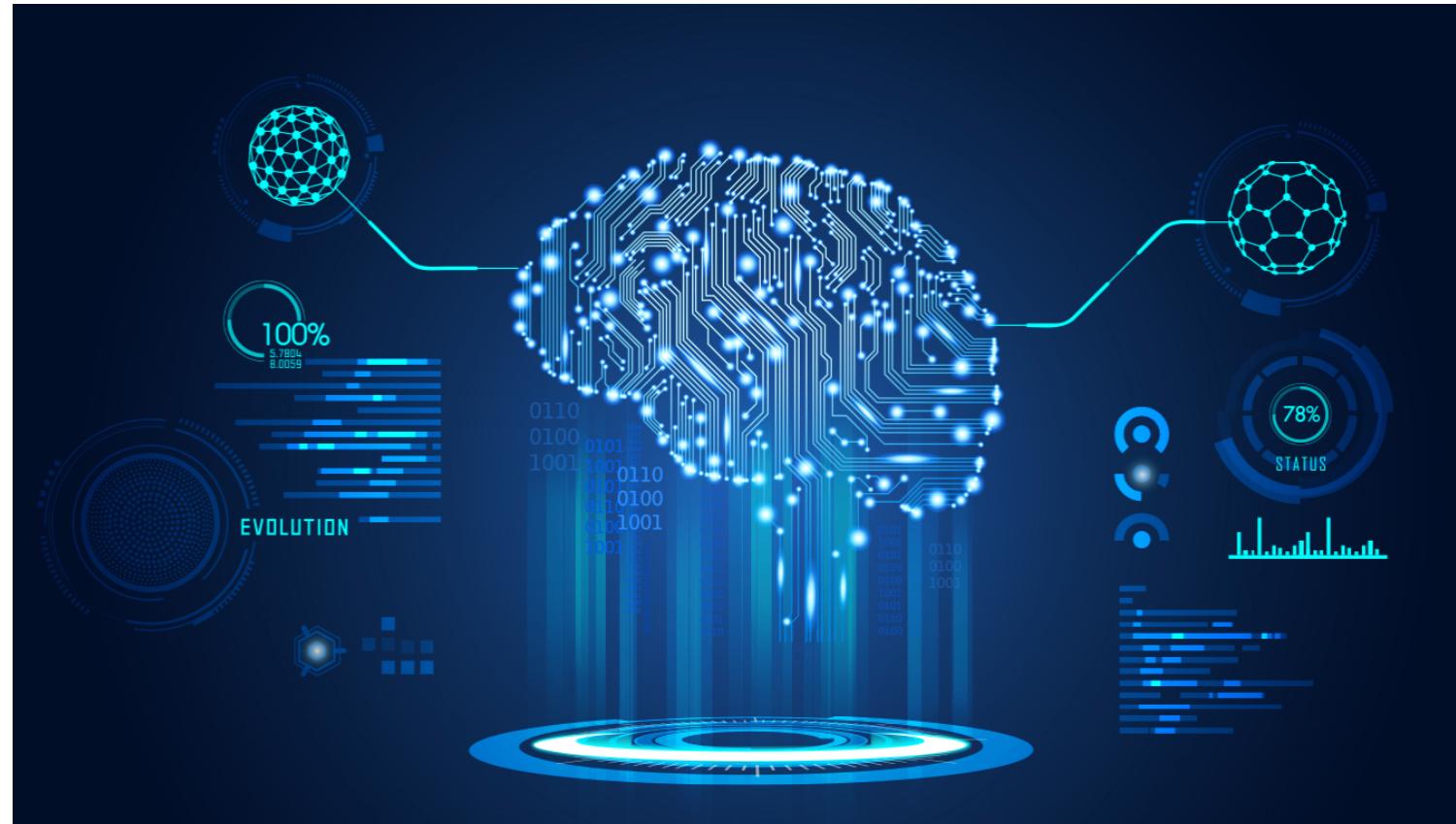
Announcements & Reminders

- HW4 Due Monday April 6th at 11:59pm
 - normal usage of late days

What will we learn today?

- Predictive statistics
- Simple linear regression (SLR) model
- Residuals
- Sum of squared-errors
- Fitted regression line
- Least-squares line

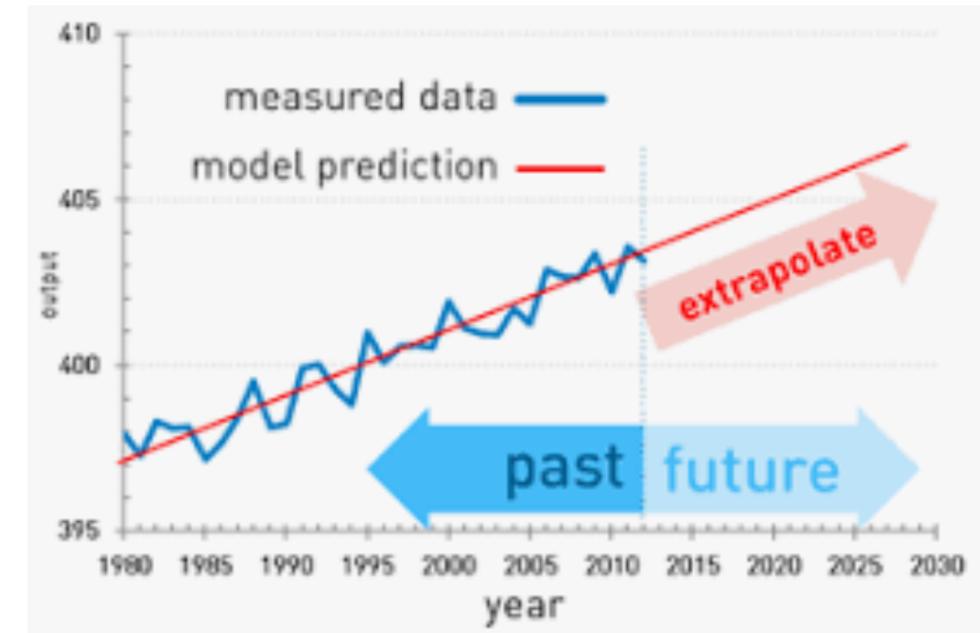
- A Modern Introduction to Probability and Statistics, Chapter 22*



Statistical Modeling

So far, we've talked about:

- Descriptive statistics: “this is the way my sample is”
- Inferential statistics: “This is what I can conclude from my sample [with $100(1 - \alpha)\%$ confidence]”
- Today: predictive statistics



Linear regression for prediction

Examples:

- Given a person's age and gender, predict their height
- Given the area of a house, predict its sale price
- Given unemployment, inflation, number of wars and economic growth, predict the president's approval rating
- Given a person's browser history, predict how long they'll stay on a product page
- Given the advertising budget expenditures in various media markets, predict the number of products they'll sell.

Simple Linear Regression

line: $y = mx + b$

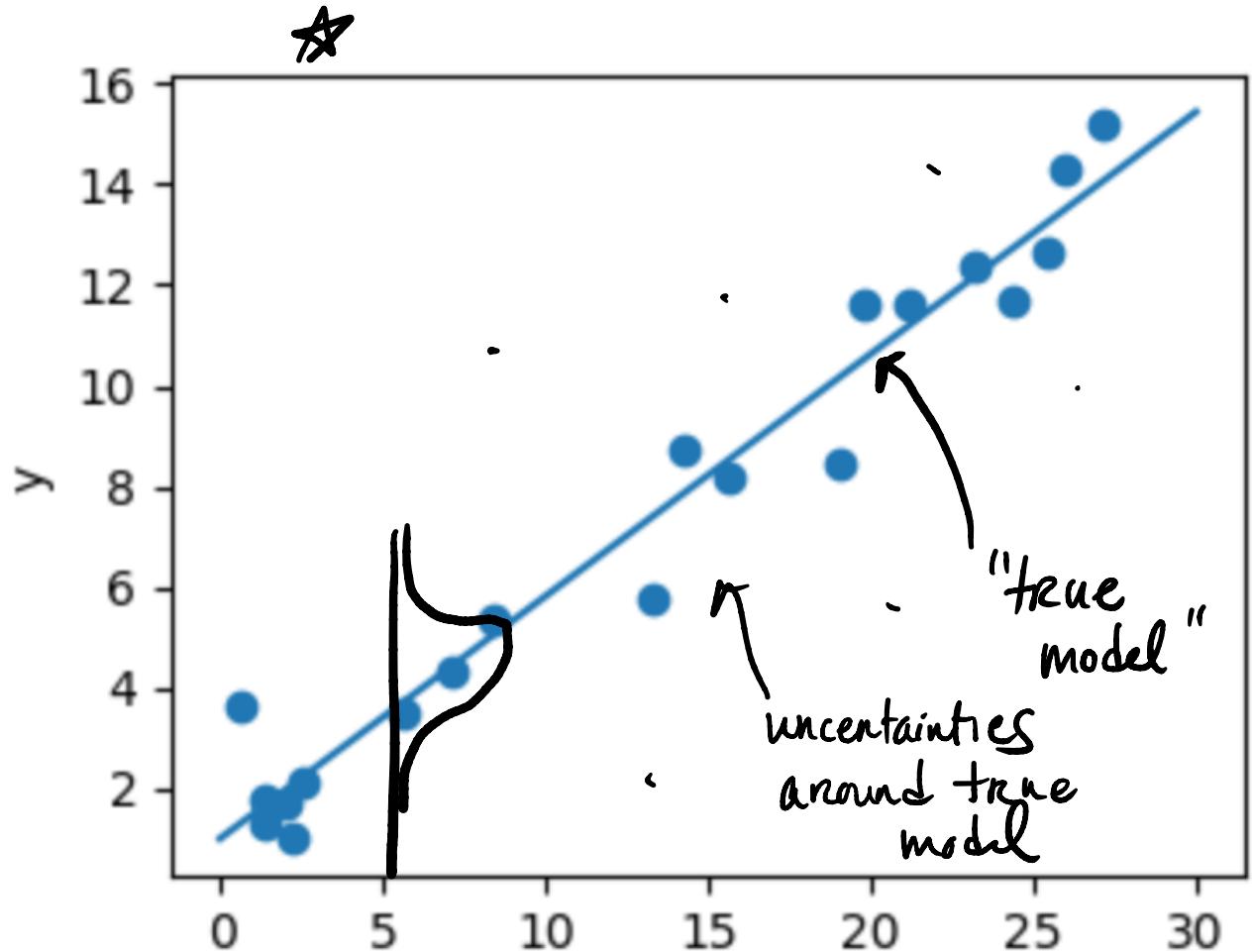
Our goal: figure out the equation of the line through the data.

Definitions and Assumptions:

1) $y_i = \alpha + \beta x_i + \epsilon_i$

2) Each of the ϵ_i are independent

3) $\epsilon_i \sim N(0, \sigma^2)$



Simple Linear Regression

$$Y = \alpha + \beta X + \epsilon$$

linear equation

noisy process

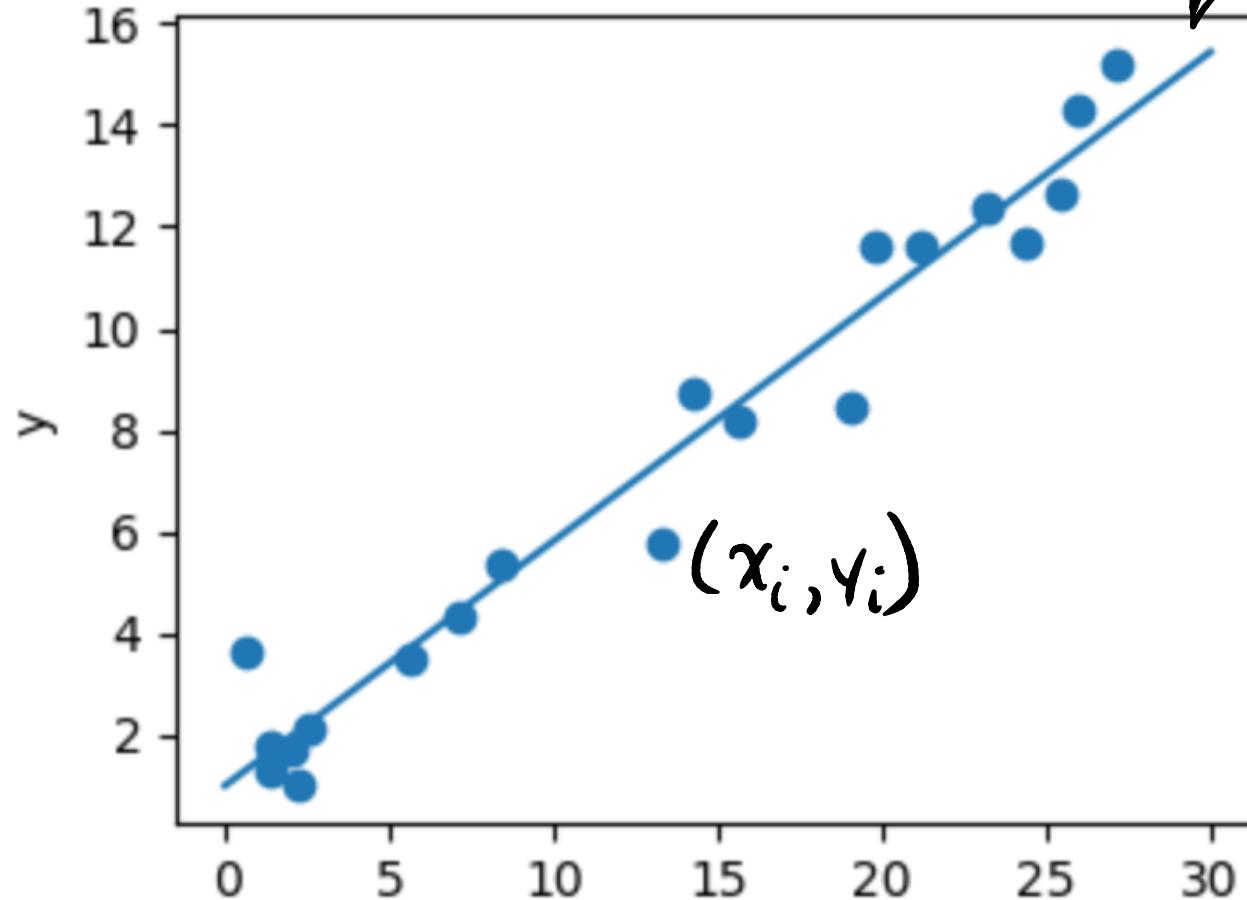
X: the independent variable, the predictor, the explanatory variable, the feature

Y: the dependent variable, the response variable

ϵ : the random deviation, random error – accounts for the fact that the world is uncertain and that there are random deviations around the true process.

Simple Linear Regression (SLR)

$$Y = \alpha + \beta X + \epsilon$$



The points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ resulting from n independent observations will be scattered about the true regression line.

Simple Linear Regression

$$\epsilon \sim N(0, \sigma^2)$$

Y is a random variable.

What is $E[Y]$?

$$E[Y] = E[\alpha + \beta X + \epsilon]$$

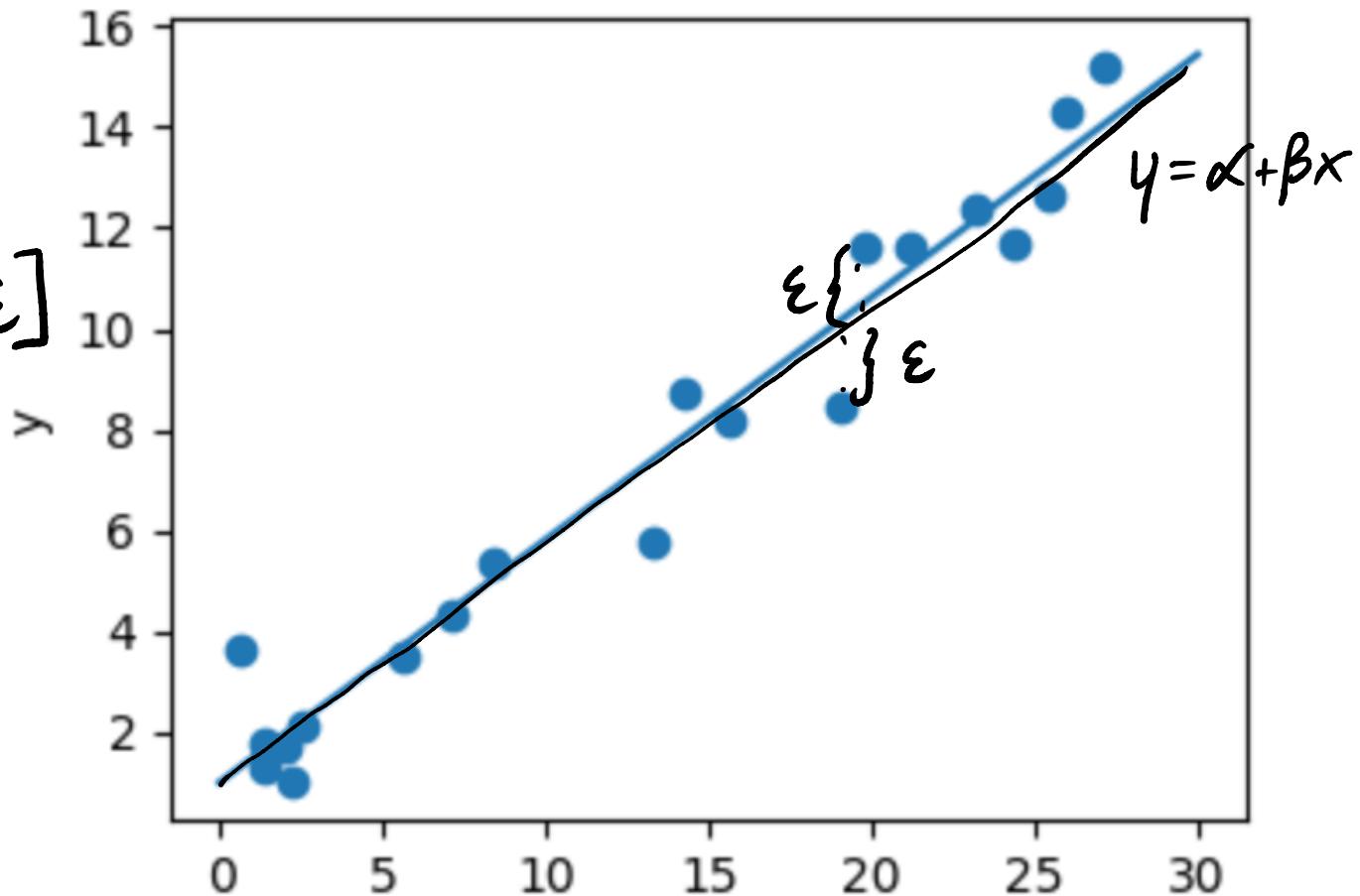
$$= E[\alpha] + E[\beta X] + E[\epsilon]$$

$$= \alpha + \beta x + 0$$

$$= \alpha + \beta x$$

$$Y = \alpha + \beta X + \epsilon$$

✓ model
for data set



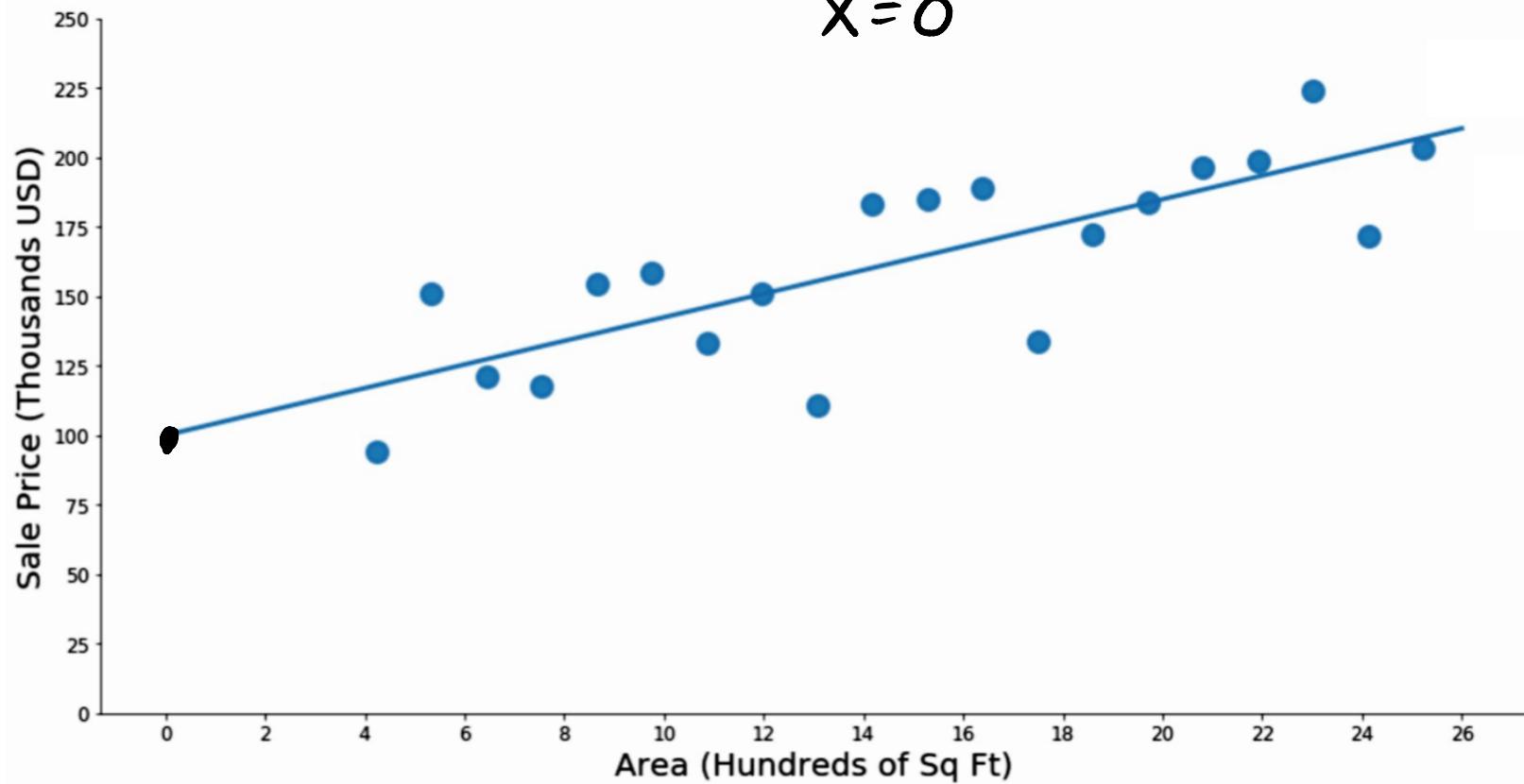
Simple Linear Regression – Interpreting parameters

true model

α is the intercept of the true regression line (aka the baseline average)

It's called the "baseline" because it's the response when the feature = 0

$$X=0$$



Simple Linear Regression – Interpreting parameters

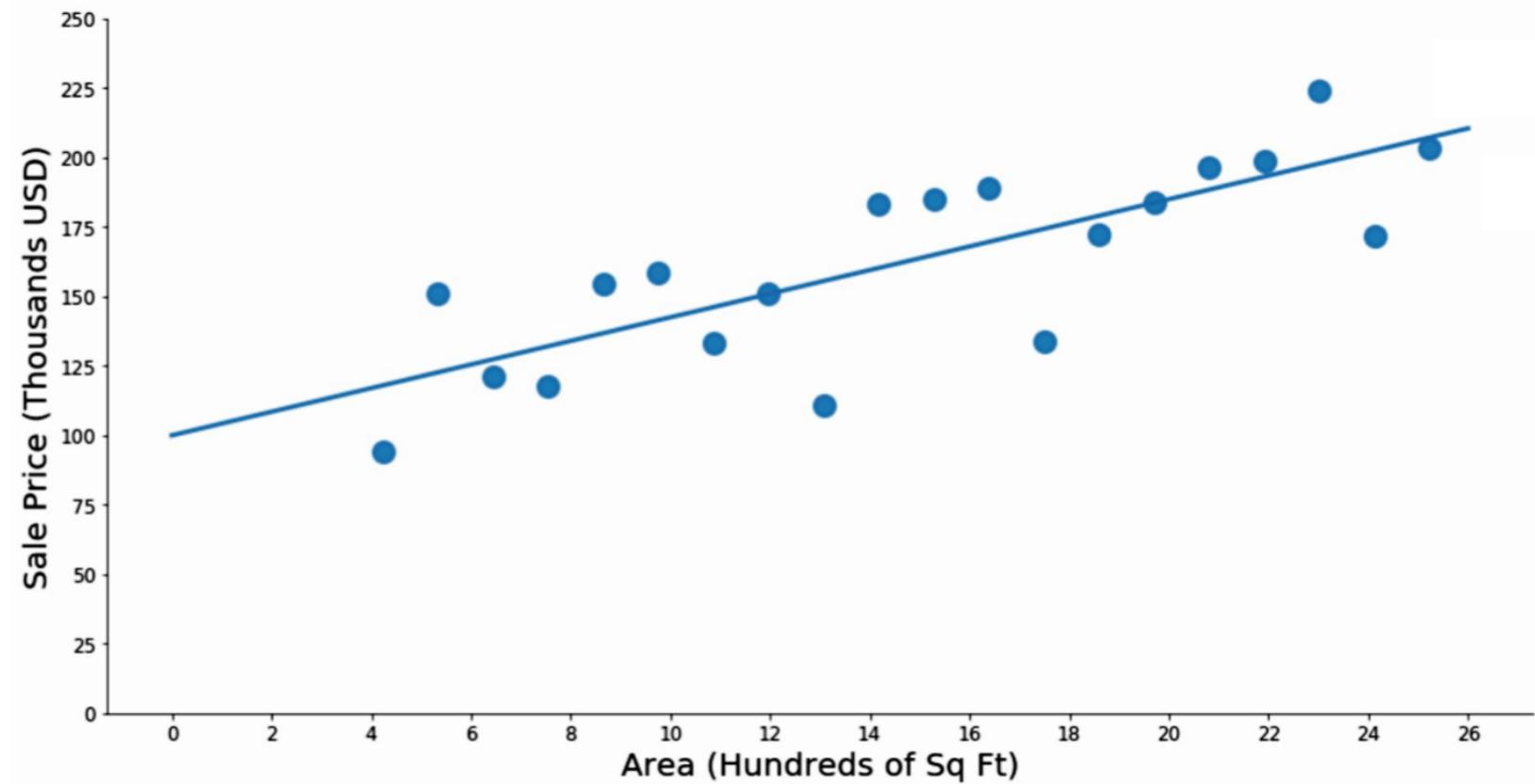
β is the slope of the true regression line

It's the increase in our response from a unit increase in our feature

$$Y(x+1) = \alpha + \beta(x+1) + \varepsilon$$

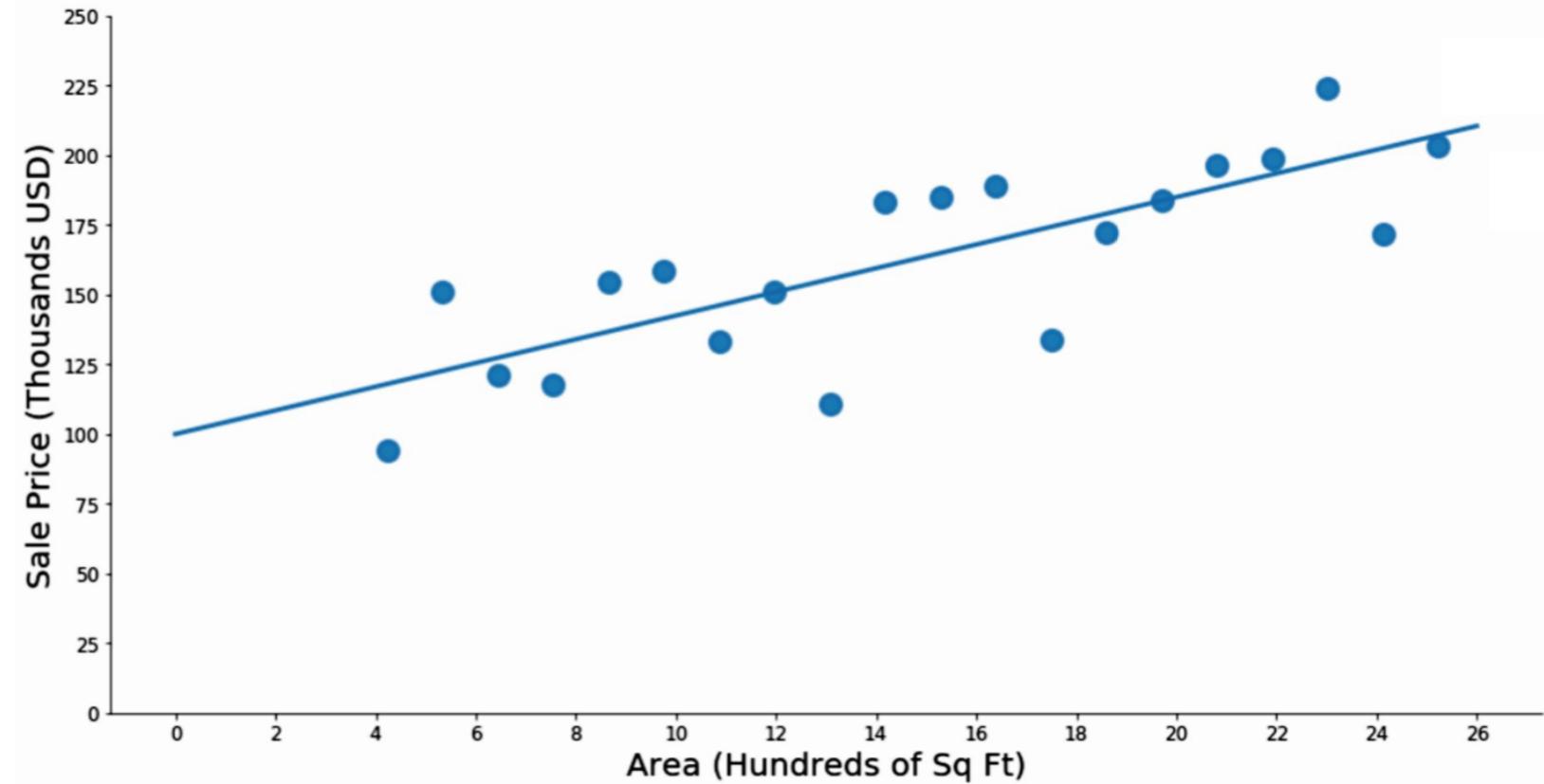
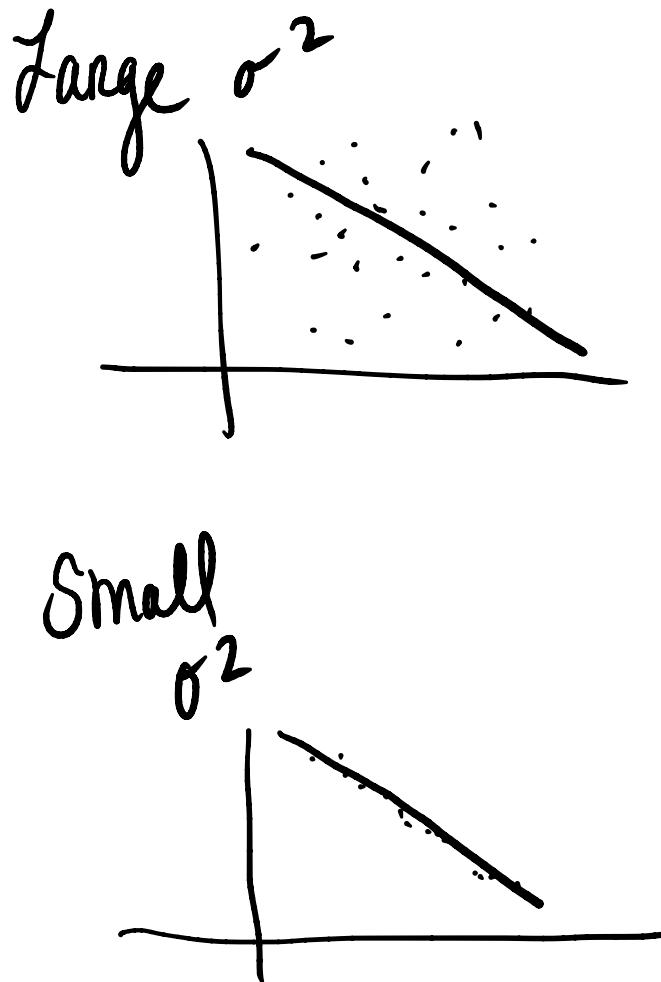
$$Y(x) = \alpha + \beta x + \varepsilon$$

$$\Delta y = \beta$$



Simple Linear Regression – Interpreting parameters

The variance parameter σ^2 determines the extent to which each normal curve spreads about the true regression line.

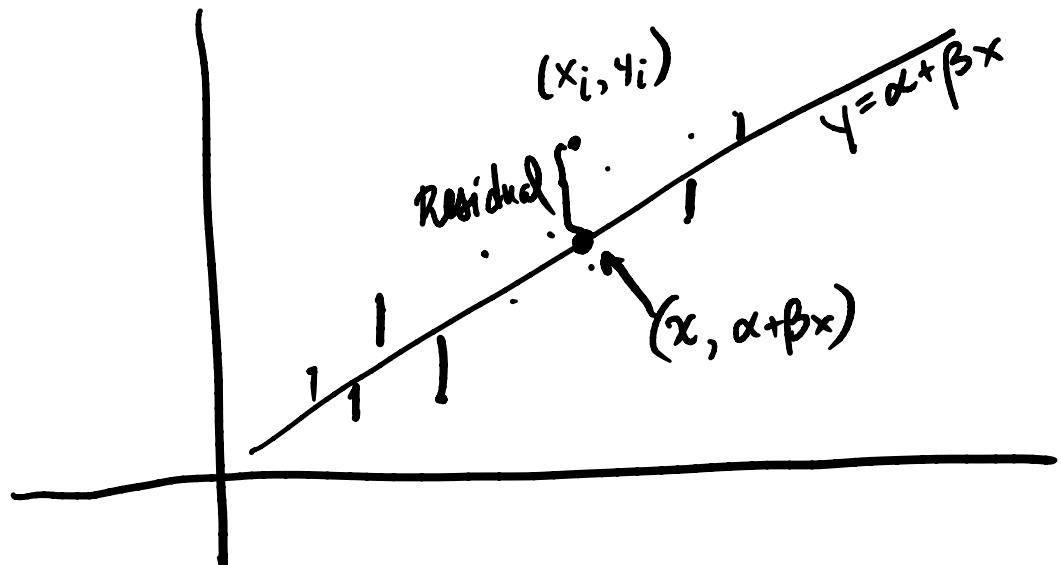


Simple Linear Regression – Estimating model parameters

We've got data: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

We want to figure out our regression line: $y = \alpha + \beta x$ vs. modeled line $\hat{y} = \hat{\alpha} + \hat{\beta} x$
 α, β - true model parameters $\hat{\alpha}, \hat{\beta}$ - estimated parameters

How can we minimize the residuals, $\epsilon_i = y_i - (\alpha + \beta x_i)$?

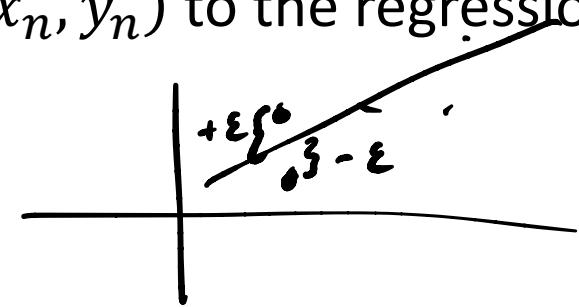


goal is to fit a line between the data so that we minimize the distance between our data and our line

Simple Linear Regression – Estimating model parameters

The sum of squared-errors for the points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ to the regression line is given by

$$\left\{ SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right.$$



The point-estimates (single value estimates from the data) of the slope and intercept parameters are called the least-squares estimates, and are defined to be the values that minimize the SSE.

our goal is to
minimize the SSE

Simple Linear Regression – Estimating model parameters

How do we actually find the parameter estimates?

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (\hat{\alpha} + \hat{\beta}x_i))^2$$

Finding

$$\frac{\partial SSE}{\partial \hat{\alpha}} = 0$$

$$f(x) = 3x^3$$

$$\begin{aligned} \frac{\partial SSE}{\partial \hat{\alpha}} &= \sum_{i=1}^n 2(y_i - (\hat{\alpha} + \hat{\beta}x_i))(-1) \\ &= -2 \sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta}x_i) = 0 \end{aligned}$$

set equal
to find minimum

$$\frac{\partial SSE}{\partial \hat{\beta}} = 0$$

$$\frac{df}{dx} = 9x^2$$

$$\sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta}x_i) = 0$$

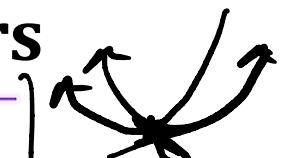
$$\sum_{i=1}^n y_i - \sum_{i=1}^n \hat{\alpha} - \sum_{i=1}^n \hat{\beta} x_i = 0$$

$$\begin{aligned} \bar{y} \cdot n - \hat{\alpha} \cdot n - \hat{\beta} \bar{x} \cdot n &= 0 \\ \bar{y} - \hat{\alpha} - \hat{\beta} \bar{x} &= 0 \end{aligned}$$

Recall:

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

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



Simple Linear Regression – Estimating model parameters

How do we actually find the parameter estimates?

$$\frac{\partial SSE}{\partial \beta} = \sum_{i=1}^n 2(y_i - \hat{\alpha} - \hat{\beta} x_i)(-x_i) = 0 \quad \begin{matrix} \text{set equal to 0} \\ \text{to find minimum} \end{matrix}$$

$$-2 \sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta} x_i)(x_i) = 0$$

$$\sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta} x_i)(x_i) = 0$$

$$\sum_{i=1}^n (x_i y_i - \hat{\alpha} x_i - \hat{\beta} x_i^2) = 0$$

$$\sum_{i=1}^n x_i y_i - \hat{\alpha} \sum_{i=1}^n x_i - \hat{\beta} \sum_{i=1}^n x_i^2 = 0$$

$$n \bar{xy} - \hat{\alpha} \cdot n \cdot \bar{x} - \hat{\beta} \cdot n \cdot \bar{x}^2 = 0$$

$$\frac{\sum_{i=1}^n x_i y_i}{n} = \bar{xy}$$

$$\frac{\sum_{i=1}^n x_i^2}{n} = \bar{x}^2$$

Simple Linear Regression – Estimating model parameters

How do we actually find the parameter estimates?

$$\bar{xy} - \hat{\alpha}\bar{x} - \hat{\beta}\bar{x^2} = 0$$

$$-\hat{\beta}\bar{x^2} = +\hat{\alpha}\bar{x} - \bar{xy}$$

$$\hat{\beta} = \frac{\hat{\alpha}\bar{x}}{-\bar{x^2}} - \frac{\bar{xy}}{-\bar{x^2}}$$

$$\hat{\beta} = \frac{\bar{xy}}{\bar{x^2}} - \alpha \frac{\bar{x}}{\bar{x^2}}$$

Two equations, two unknowns : $\alpha = \bar{y} - \hat{\beta}\bar{x}$

$$\hat{\beta} = \frac{\bar{xy}}{\bar{x^2}} - \frac{\bar{x}}{\bar{x^2}}(\bar{y} - \hat{\beta}\bar{x})$$

$$\hat{\beta} = \frac{\bar{xy}}{\bar{x^2}} - \frac{\bar{x}\bar{y}}{\bar{x^2}} + \frac{\hat{\beta}(\bar{x})^2}{\bar{x^2}}$$

$$\rightarrow \hat{\beta} - \hat{\beta} \cdot \frac{(\bar{x})^2}{\bar{x^2}} = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2}}$$

$$\hat{\beta} \left(1 - \frac{(\bar{x})^2}{\bar{x^2}} \right) = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2}}$$

$$\hat{\beta} \left(\frac{\bar{x^2} - (\bar{x})^2}{\bar{x^2}} \right) = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2}}$$

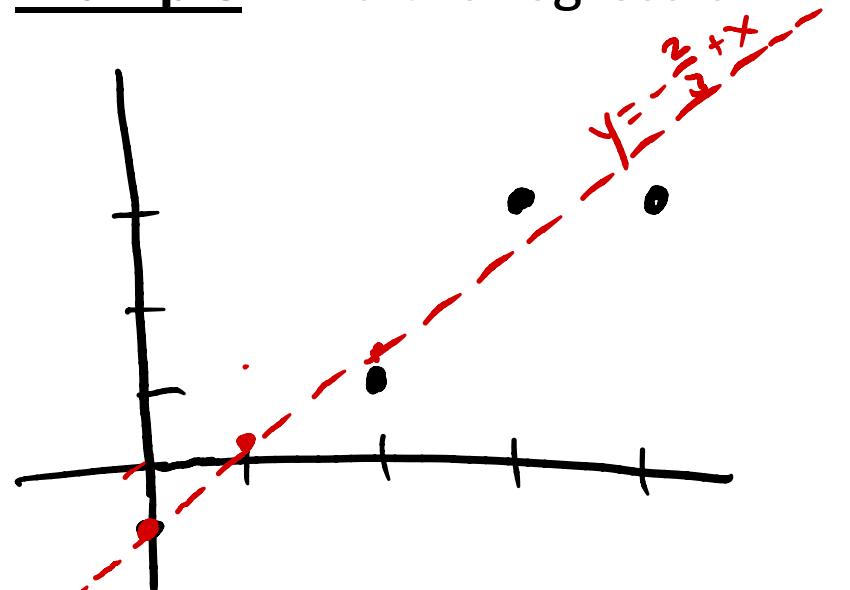
$$\hat{\beta} = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2}} \cdot \frac{\bar{x^2}}{\bar{x^2} - (\bar{x})^2}$$

$$\boxed{\hat{\beta} = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2} - (\bar{x})^2}},$$

$$\boxed{\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}}.$$

Simple Linear Regression – Estimating model parameters

Example: Find the regression line for the following data: (2,1), (3,3), (4,3)



$$\hat{\beta} = \frac{3 \cdot \frac{7}{3} - \frac{23}{3}}{9 - \frac{27}{3}} = \frac{\frac{21}{3} - \frac{23}{3}}{\frac{27}{3} - \frac{27}{3}} = \frac{-\frac{2}{3}}{-\frac{2}{3}} = 1$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} = \frac{7}{3} - 1 \cdot 3 = \frac{7}{3} - \frac{9}{3} = -\frac{2}{3}$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

$$\hat{\beta} = \frac{\bar{x} \bar{y} - \bar{x} \bar{y}}{(\bar{x})^2 - \bar{x}^2}$$

$$\bar{x} = \frac{2+3+4}{3} = \frac{9}{3} = 3$$

$$\bar{y} = \frac{1+3+3}{3} = \frac{7}{3}$$

$$\bar{x} \bar{y} = \frac{2 \cdot 1 + 3 \cdot 3 + 4 \cdot 3}{3} = \frac{23}{3}$$

$$\bar{x}^2 = \frac{2^2 + 3^2 + 4^2}{3} = \frac{29}{3}$$

$$(\bar{x})^2 = 3^2 = 9$$

⇒ Simple Regression Line

$$\boxed{\hat{y} = -\frac{2}{3} + x}$$

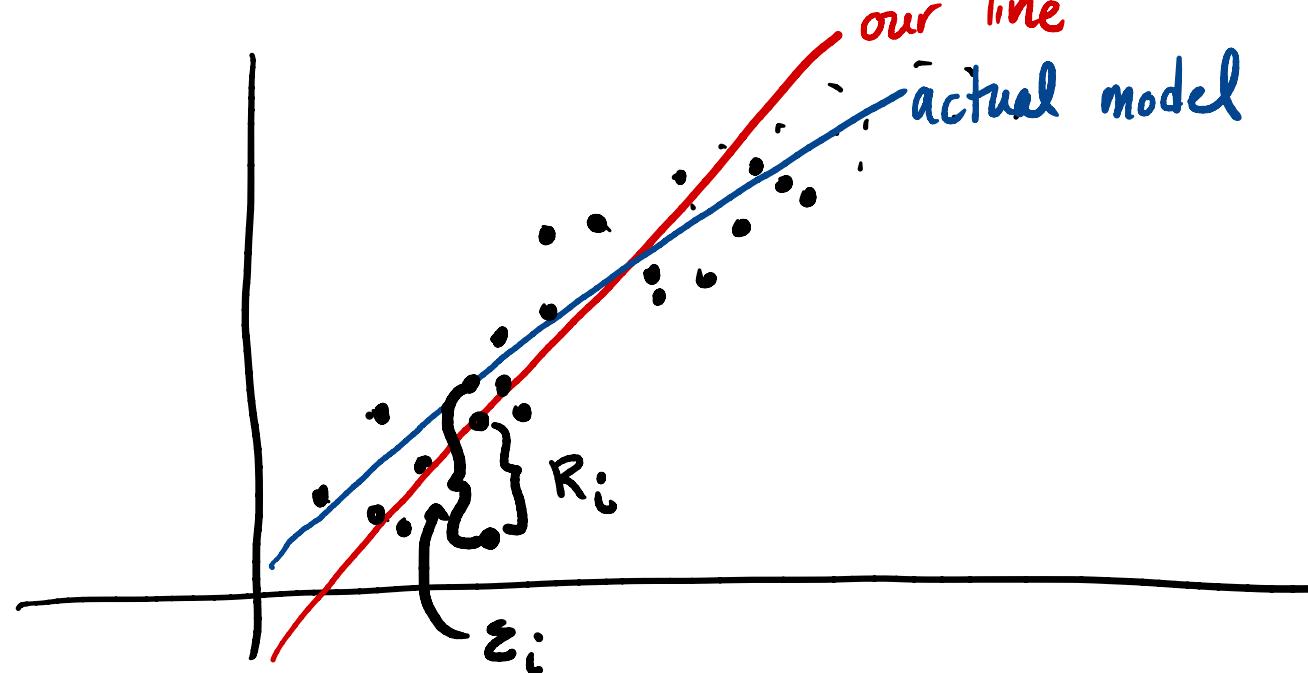
Residuals

The residuals are the difference between the observed and the predicted responses:

$$r_i = y_i - \hat{y}_i$$

The residuals r_i are estimates of the unknown true error ϵ_i

Differences between error and residuals



Next Time:

- ❖ Inference in Simple Linear Regression