

# Simple Linear Regression

ST552 Lecture 2

Charlotte Wickham

January 6, 2015



# Review

## High level review

Since simple linear regression is a special case of multiple linear regression, we'll leave the "whys?" to when we cover multiple linear regression.

- ▶ the simple linear regression model →
- ▶ interpretation ✓
- ▶ assumptions ~
- ▶ how the estimates are found ✓
- ▶ properties of the estimates ✓
- ▶ F-tests ~

$$y = ax + b + \text{error}$$
$$X\beta + \varepsilon$$

# The simple linear regression model

$n$  observations are collected in pairs,  $(x_i, y_i), i = 1, \dots, n$  where the  $y_i$  are generated according to the model,

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

Annotations for the equation above:  
- A bracket above  $\beta_0 + \beta_1 x_i$  is labeled "straight line".  
- An arrow points from  $\beta_0$  to the word "intercept".  
- An arrow points from  $\beta_1$  to the word "slope".  
- An arrow points from  $\epsilon_i$  to the words "noise or error term random".  
- An arrow points from the entire equation to the word "assumptions".  
- An arrow points from the word "assumptions" to the text below.

## What is random?

where  $\epsilon_i$  are independent and identically distributed with expected value zero, and variance  $\sigma^2$ .

For inference, we often also assume the  $\epsilon_i$  are Normally distributed,

$$\epsilon_i \stackrel{i.i.d}{\sim} N(0, \sigma^2)$$

$\beta_0, \beta_1$  fixed parameters  
 $\hat{\beta}_0, \hat{\beta}_1$  random variables

$x_i$  fixed

# The assumptions in words

- ▶ **Linearity:** the mean response is a straight line function of the explanatory variable
- ▶ **Constant spread:** the standard deviation around the mean response is the same at all values of the explanatory variable
- ▶ **Normality:** the deviations from the mean response, the errors, are Normally distributed.
- ▶ **Independence:** the deviations from the mean response are independent.

# Interpretation of the parameters

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$
$$E[y_i] = \beta_0 + \beta_1 x_i \quad \cdot \quad E(\varepsilon_i) = 0$$

**Intercept,  $\beta_0$ ,**

When the explanatory variable is zero, the mean response is  $\beta_0$ .

**Slope,  $\beta_1$ ,**

An increase in the explanatory variable of one unit is associated with a change in mean response of  $\beta_1$ .

(Careful with causal language... is it justified?)

But we don't know  $\beta_0$  and  $\beta_1$ ...



# Properties of the least squares estimates

$$\hat{\beta}_0 \quad \hat{\beta}_1$$

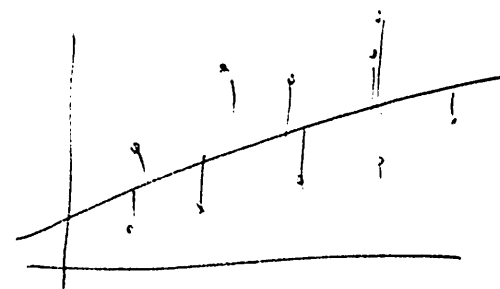
$$E(\epsilon_i) = 0$$

$$E(\hat{\beta}_0) = \beta_0$$

Using the moment assumptions of  $\epsilon_i$ , the least squares estimates can be shown to be unbiased. You can derive their variances,  $\text{Var}(\hat{\beta}_0)$ ,  $\text{Var}(\hat{\beta}_1)$ ,  $\text{Cov}(\hat{\beta}_0, \hat{\beta}_1)$ , but they depend on the unknown  $\underline{\sigma}$ .

An unbiased estimate of  $\sigma$  is

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2$$



Intuition:  $\frac{1}{n} \sum_{i=1}^n e_i^2$  seems a reasonable place to start to estimate the variance of the errors, but this tends to underestimate the variance because we picked our estimates to make the sum of squared errors as small as possible.