

# Simple Linear Regression

Dustin Long, PhD

Department of Biostatistics  
University of Alabama at Birmingham

September 3, 2019

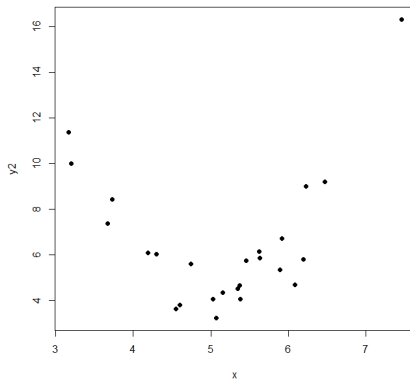
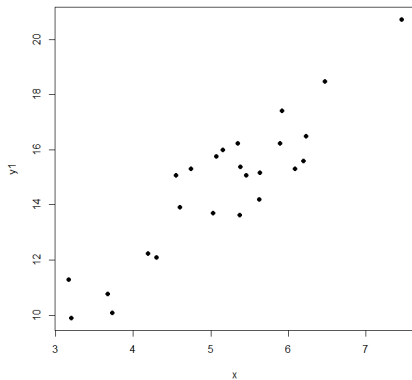
## Outline:

- Simple Linear Regression

## Outline:

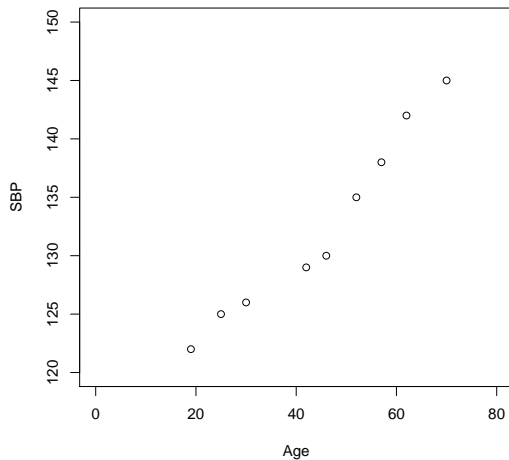
- Assumptions and properties of regression with one variable
- Determination and measures for the line of best fit
- Inference and Interpretations of parameters

- When do we use linear regression?
- What type of model do we use?
- What is the best fitting model, and what do we mean by “best fit”?



## Example: SBP and Age

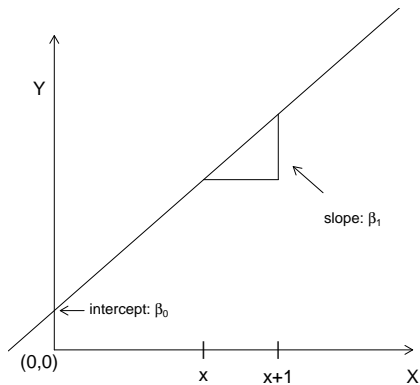
Obs	Age	SBP
1	19	122
2	25	125
3	30	126
4	42	129
5	46	130
6	52	135
7	57	138
8	62	142
9	70	145



## Mathematical Properties of a Straight Line

- A line is defined by an intercept and slope, i.e.,  $y = Mx + B$
- Recall properties of straight lines





## Definition of Linear Regression

- The model used for simple linear regression is:
- $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$
- where  $Y_i$  is a dependent (outcome) random variable,  $X_i$  is an independent random variable,  $i = 1 \dots n$ ,  $\beta_0$  is the population intercept,  $\beta_1$  is the population slope, and  $\epsilon_i \sim N(0, \sigma^2)$
- **$\mathbf{Y} = \mathbf{X}\beta + \mathbf{e}$**

## Statistical Assumptions of a Straight Line

- Assumption 1: Homoscedasticity,  $\sigma_x^2 = \sigma^2 \quad \forall x$
- Assumption 2: Independence,  $Y$  values are independent of one another
- Assumption 3: Linearity,  $E[Y|X = x] = \mu_{Y|X=x} = \beta_0 + \beta_1 x$  or  $Y = \beta_0 + \beta_1 x + E$
- Assumption 4: Existence,  $(Y|X = x \sim f(\mu, \sigma_x^2)$ , where  $\mu, \sigma_x^2 < \infty, \forall x$
- Assumption 5: Normality,  $Y|X = x \sim N(\mu_{Y|X=x}, \sigma^2)$
-

## Determining the Best-fitting Line

- Least Squares Method, OLS Example
- Another example
- Minimum-variance method, Best Linear Unbiased Estimators (BLUE)
- Under Assumptions 1-5, OLS estimates and BLUE are the same
- $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$
- $\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$

## Least Squares Estimation

- Least squares estimators are values of  $\beta_0$  and  $\beta_1$  that minimize

$$\sum_{i=1}^N (Y_i - \beta_0 - \beta_1 X_i)^2$$

- Set partial derivatives equal to 0, solve for  $\beta_0$  and  $\beta_1$
- Or can set  $\mathbf{Y} = \mathbf{X}\hat{\beta}$  and solve for  $\hat{\beta}$

- When using hats, all estimated quantities get a hat
- WRONG:  $Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$
- CORRECT:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$
- Regression lines with out hats must have  $\epsilon_i$  or it is just a line NOT regression
- $\hat{\epsilon}_i = \hat{e}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$  and  $\hat{\mathbf{e}}$  are all residuals
- $\sum_{i=1}^n \hat{e}_i = 0$ .

- $MSE$ , mean square error, is the primary estimate of  $\sigma^2$
- $SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$
- $SSE = \hat{\mathbf{e}}' \hat{\mathbf{e}}$
- If  $SSE = 0$  then the estimated regression line fits the data perfectly
- $MSE = SSE / (n - 2)$  ONLY FOR SIMPLE LINEAR REGRESSION
- $SSE$  and  $MSE$  are given directly from software

## Inference about the Slope and Intercept

- Under Assumptions 1-5,  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are Normally distributed
- To test  $H_0 : \beta_1 = \beta_1^{(0)}$ , the statistic is
- $T = \frac{\hat{\beta}_1 - \beta_1^{(0)}}{S_{\hat{\beta}_1}}$
- where  $S_{\hat{\beta}_1} = \frac{\sqrt{MSE}}{\sqrt{(n-1) \sum_{i=1}^n (x_i - \bar{x})^2}}$
- $T \sim T_{n-2}$  under  $H_0$
- C.I. for  $\beta_1$  is constructed by inverting the previous test statistic
- Test for  $\beta_0$  exists but rarely used



## Interpretations of Tests

- Most common test for slope,  $H_0 : \beta_1 = 0$
- If  $H_0$  not rejected, it does NOT mean there is no relationship between  $Y$  and  $X$ , just there is no evidence that the relationship is linear
- If  $H_0$  is rejected, there is at minimum a linear relationship between  $Y$  and  $X$  but that might not be the entire story

## Inferences about the regression line

- Recall that  $\mu_{Y|X=x} = \beta_0 + \beta_1 x$ , a particular point in the regression line
- Can test  $H_0 : \mu_{Y|X=x} = \mu_{Y|X=x}^{(0)}$
- $S_{\hat{Y}_x} = S_{Y|X} \sqrt{\frac{1}{n} + \frac{(x-\bar{x})^2}{(n-1)S_x^2}}$
- This tests the value of the line at a single point,  $x$
- Inverted C.I. called confidence bands when constructed for all observed values of  $X$

- For prediction of  $\hat{Y}_i$  at a particular value of  $X_i$  use
- $S_{\hat{Y}_x} = S_{Y|X} \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{(n-1)S_x^2}}$
- to create prediction intervals and bands

# ANOVA TABLE

Questions?