

# SIMPLE LINEAR REGRESSION: OLS POINT ESTIMATORS OF $\hat{\beta}_0$ & $\hat{\beta}_1$ [DEVORE 12.2]

## SIMPLE LINEAR REGRESSION MODEL (DEFINITION):

The **simple linear regression model** is:  $Y_i = \beta_0 + \beta_1 x_i + E_i$  where  $E_1, \dots, E_n \stackrel{IID}{\sim} \text{Normal}(0, \sigma^2)$

- $Y_i$   $\equiv$  rv for measurement of  $i^{\text{th}}$  response
- $x_i$   $\equiv$  Actual measurement of  $i^{\text{th}}$  regressor
- $\beta_0$   $\equiv$  Expected value of response when regressor is zero
- $\beta_1$   $\equiv$  Expected change in response per unit increase in regressor
- $E$   $\equiv$  Effect of random error on  $i^{\text{th}}$  response

The **realized simple linear regression model** is:  $y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + e_i \quad \forall i = 1, \dots, n$

- $y_i$   $\equiv$  Actual measurement of  $i^{\text{th}}$  response
- $x_i$   $\equiv$  Actual measurement of  $i^{\text{th}}$  regressor
- $\hat{\beta}_0$   $\equiv$  Estimated value of response when regressor is zero
- $\hat{\beta}_1$   $\equiv$  Estimated change in response per unit increase in regressor
- $e_i$   $\equiv$  Actual error on  $i^{\text{th}}$  response

**SUMS (DEFINITION):** Let vectors  $\mathbf{x} = (x_1, \dots, x_n)^T$  and  $\mathbf{y} = (y_1, \dots, y_n)^T$ . Then:

$$S_x := \sum_i x_i \quad S_y := \sum_i y_i \quad S_{xx} := \sum_i x_i x_i \quad S_{yy} := \sum_i y_i y_i \quad S_{xy} := \sum_i x_i y_i$$

**CAUCHY-SCHWARZ INEQUALITY:** Let vectors  $\mathbf{x} = (x_1, \dots, x_n)^T$  and  $\mathbf{y} = (y_1, \dots, y_n)^T$ . Then:

$$(a) \quad (S_{xy})^2 \leq S_{xx} \cdot S_{yy} \quad (b) \quad (S_x)^2 \leq n \cdot S_{xx}$$

where equality in part (a) holds when either  $\mathbf{x} = \vec{0}$  or  $\mathbf{y} = \vec{0}$  or  $\mathbf{y} = c\mathbf{x}$  where  $c \neq 0$ .

where equality in part (b) holds when either  $\mathbf{x} = \vec{0}$  or  $\mathbf{x} = k \cdot \vec{1}$  where  $k \in \mathbb{R}$ .

**CENTERED SUMS (DEFINITION):** Let vectors  $\mathbf{x} = (x_1, \dots, x_n)^T$  and  $\mathbf{y} = (y_1, \dots, y_n)^T$ . Then:

$$SC_{xx} := \sum_i (x_i - \bar{x})^2 \quad SC_{yy} := \sum_i (y_i - \bar{y})^2 \quad SC_{xy} := \sum_i (x_i - \bar{x})(y_i - \bar{y})$$

## CENTERED SUMS LEMMA:

$$(a) \quad SC_{xx} = S_{xx} - \frac{1}{n} S_x S_x \quad (b) \quad SC_{yy} = S_{yy} - \frac{1}{n} S_y S_y \quad (c) \quad SC_{xy} = S_{xy} - \frac{1}{n} S_x S_y$$

## OLS ESTIMATORS:

$$\begin{cases} \hat{\beta}_1 &= \frac{SC_{xy}}{SC_{xx}} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})(x_i - \bar{x})} = \frac{S_{xy} - \frac{1}{n} S_x S_y}{S_{xx} - \frac{1}{n} S_x S_x} \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} = \frac{\sum_i y_i - \hat{\beta}_1 \sum_i x_i}{n} = \frac{S_y - \hat{\beta}_1 S_x}{n} \end{cases}$$

**GAUSS-MARKOV THEOREM:** Suppose the following conditions are all satisfied:

$$\begin{aligned} \mathbb{E}[E_i] &= 0 && \text{(errors are all centered at zero)} \\ \mathbb{V}[E_i] &= \sigma^2 && \text{(errors all have the same finite variance)} \\ \mathbb{C}[E_i, E_{i'}] &= 0 && \text{(errors are uncorrelated when } i \neq i') \end{aligned}$$

Then, the least-squares estimators (LSE's)  $\hat{\beta}_0, \hat{\beta}_1$  are both BLUE's.

# SIMPLE LINEAR REGRESSION: SUMS OF SQUARES [DEVORE 12.2]

## DEVIATION LEMMA:

$$(a) \sum_i (x_i - \bar{x}) = 0 \quad (b) \sum_i (x_i - \bar{x})x_i = SC_{xx} \quad (c) \sum_i (x_i - \bar{x})y_i = SC_{xy}$$

## RESIDUAL LEMMA:

$$(a) \sum_i (y_i - \hat{y}_i) = 0 \quad (b) \sum_i (y_i - \hat{y}_i)x_i = 0 \quad (c) \sum_i (y_i - \hat{y}_i)\hat{y}_i = 0$$

## SUMS OF SQUARES (DEFINITION):

$$SS_{total} := \sum_i (y_i - \bar{y})^2 \quad SS_{res} := \sum_i (y_i - \hat{y}_i)^2 \quad SS_{reg} := \sum_i (\hat{y}_i - \bar{y})^2$$

## VARIATION PARTITIONING:

$$\underbrace{\sum_i (y_i - \bar{y})^2}_{SS_{total}} = \underbrace{\sum_i (y_i - \hat{y}_i)^2}_{SS_{res}} + \underbrace{\sum_i (\hat{y}_i - \bar{y})^2}_{SS_{reg}}$$

## SUMS OF SQUARES LEMMA:

$$(a) SS_{total} = SC_{yy} \quad (b) SS_{res} = SC_{yy} - \hat{\beta}_1^2 \cdot SC_{xx} \quad (c) SS_{reg} = \hat{\beta}_1^2 \cdot SC_{xx}$$

## DEGREES OF FREEDOM:

$$\underbrace{SS_{total}}_{\text{Total Variation}} = \underbrace{SS_{reg}}_{\text{Variation due to Regression}} + \underbrace{SS_{res}}_{\text{Unexplained Variation}}$$

$$\sum_i (y_i - \hat{\mu})^2 = \sum_i (\hat{y}_i - \hat{\mu})^2 + \sum_i (y_i^{res})^2$$

$$\sum_i (y_i - \bar{y})^2 = \sum_i (\hat{y}_i - \bar{y})^2 + \sum_i (y_i - \hat{y}_i)^2$$

$$\sum_i (y_i - \hat{\mu})^2 = \sum_i [(\hat{\beta}_0 + \hat{\beta}_1 x_i) - \hat{\mu}]^2 + \sum_i [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2$$

$$\underbrace{\nu}_{\text{Total dof's}} = \underbrace{\nu_{reg}}_{\text{Regression dof's}} + \underbrace{\nu_{res}}_{\text{Residual dof's}}$$

$$\nu = n - 1 \quad \nu_{reg} = 2 - 1 = 1 \quad \nu_{res} = n - 2$$

$$\left( \begin{array}{c} \# \text{ dof's in} \\ \text{SS expr} \end{array} \right) = \left( \begin{array}{c} \# \text{ total responses} \\ \text{or parameter estimates} \\ \text{in left difference term} \end{array} \right) - \left( \begin{array}{c} \# \text{ parameter estimates} \\ \text{in right difference term} \end{array} \right)$$

“... we can think of degrees of freedom as a form of currency. In a broad sense, statistics deals with the marketplace of knowledge. Labor, in this framework, corresponds to the task of gathering information. More precisely, our job is to draw a random sample from a population. Each observation we obtain results in a degree of freedom, much like every few minutes of work at a job results in a dollar earned. In statistics, we spend our degrees of freedom to estimate parameters, to increase the probability of reaching correct decisions, or to form models of the way the world behaves... In short, degrees of freedom buy knowledge.”<sup>†</sup>

<sup>†</sup>R.S. Schulman, *Statistics in Plain English with Computer Applications*, 1992. (§2.7)

**SIMPLE LINEAR REGRESSION:  
POINT ESTIMATION OF  $\sigma^2$  [DEVORE 12.2]**

$\hat{\beta}_1$  AS LINEAR COMBINATION OF RESPONSES:

$$\hat{\beta}_1 = \sum_i \xi_i y_i \quad \text{where} \quad \xi_i := \frac{(x_i - \bar{x})}{SC_{xx}}$$

EXPECTATION & VARIANCE OF  $\hat{\beta}_1$ :

$$\mathbb{E}[\hat{\beta}_1] = \beta_1 \qquad \mathbb{V}[\hat{\beta}_1] = \frac{\sigma^2}{SC_{xx}}$$

EE LEMMA:

$$(a) \mathbb{E} [\sum_i (E_i - \bar{E})] = 0 \quad (b) \mathbb{E}[E_i \bar{E}] = \frac{\sigma^2}{n} \quad (c) \mathbb{E} [\sum_i (E_i - \bar{E})^2] = (n - 1)\sigma^2$$

EXPECTATION OF  $SS_{res}$ :

$$\mathbb{E} [SS_{res}] = (n - 2)\sigma^2$$

MEAN SQUARED RESIDUAL:

$$MS_{res} := \frac{SS_{res}}{n - 2}$$

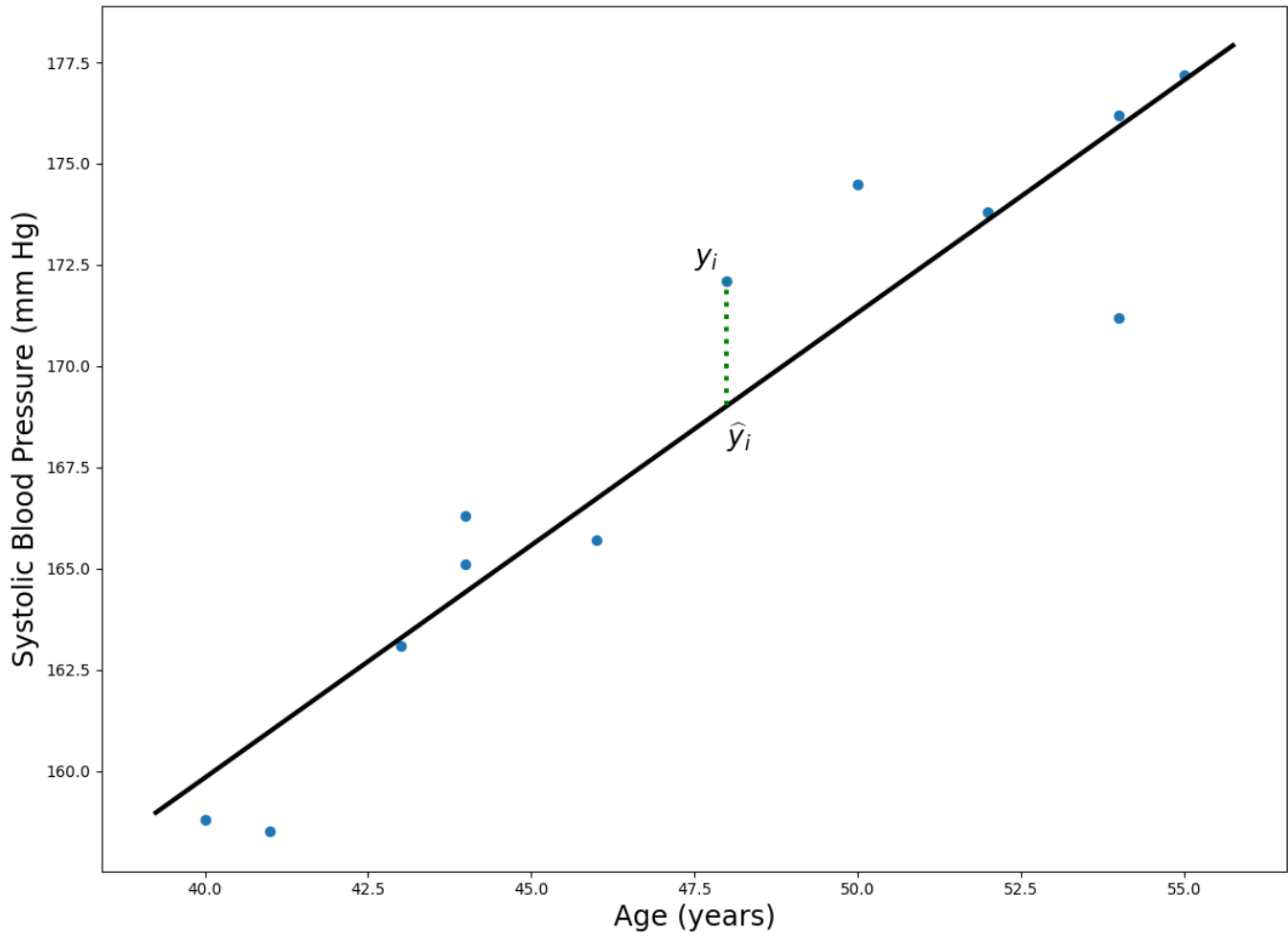
POINT ESTIMATOR OF  $\sigma^2$ :

$$\mathbb{E} [MS_{res}] = \sigma^2 \implies \hat{\sigma}^2 = MS_{res}$$

VARIATION EXPLANATION ( $R^2$ ):

$$\text{The coefficient of determination is: } R^2 := \frac{SS_{reg}}{SS_{total}} = \frac{SC_{xy} \cdot SC_{xy}}{SC_{xx} \cdot SC_{yy}}$$

SIMPLE LINEAR REGRESSION:  
VISUALIZATION [DEVORE 12.2]



The  $i^{th}$  residual ( $y_i^{res}$ ) is represented by the vertical dashed green line's length.

**EX 12.2.1:** The systolic blood pressures of twelve adults were measured as a single group without intervention<sup>†</sup>:

SYSTOLIC BLOOD PRESSURE (in mmHg) versus Age (in years)													TOTAL
Index ( $i$ )	1	2	3	4	5	6	7	8	9	10	11	12	$n = 12$
Age ( $x_i$ )	40	41	43	44	44	46	48	50	52	54	54	55	$\sum_i x_i = 571$
BP ( $y_i$ )	158.8	158.5	163.1	165.1	166.3	165.7	172.1	174.5	173.8	176.2	171.2	177.2	$\sum_i y_i = 2022.5$

<sup>†</sup>(Simplification of data) Y.H. Chan, “Biostatistics 201: Linear Regression Analysis”, *Singapore Med. J.*, **45** (2004), 55-61.

- (a) Identify the response and the sole regressor in this observational study.
- (b) Formulate this observational study as a linear model.
- (c) Compute the response sum and regressor sum:  $S_y$ ,  $S_x$ .
- (d) Compute the squared regressor sum,  $S_{xx}$ .
- (e) Compute the cross-termed regressor-response sum,  $S_{xy}$ .
- (f) Compute the cross-termed regressor-response centered sum,  $SC_{xy}$ .
- (g) Compute the squared regressor centered sum,  $SC_{xx}$ .
- (h) Compute the ordinary least-squares (OLS) estimates of the regression parameters  $\hat{\beta}_0, \hat{\beta}_1$ .
- (i) Determine the OLS best-fit line.
- (j) Compute the squared response sum,  $S_{yy}$ .
- (k) Compute the squared response centered sum,  $SC_{yy}$ .
- ( $\ell$ ) Compute the residual sum of squares,  $SS_{res}$ .
- (m) Compute the linear model's estimated variance,  $\hat{\sigma}^2$ .
- (n) Compute  $R^2$ .