

**SIMPLE LINEAR MODEL VIA ORDINARY LEAST SQUARES (OLS)**

Assume that on  $n$  statistical units (individuals) we observe  $(x_i, y_i)$ ,  $i=1, \dots, n$ .

Hence the data are  $\mathbf{y} = (y_1, \dots, y_n)$  and  $\mathbf{x} = (x_1, \dots, x_n)$

We consider that each  $y_i$  is realization of a r.v.  $Y_i$ ,  $i=1, \dots, n \rightarrow$  sample space  $\mathcal{S} = \mathcal{Y}^n = \mathcal{R}^n$

We do not specify a distribution for  $(Y_1, \dots, Y_n)$ : we only make assumptions about the first two moments  $E[Y_i]$  and  $var(Y_i)$ .

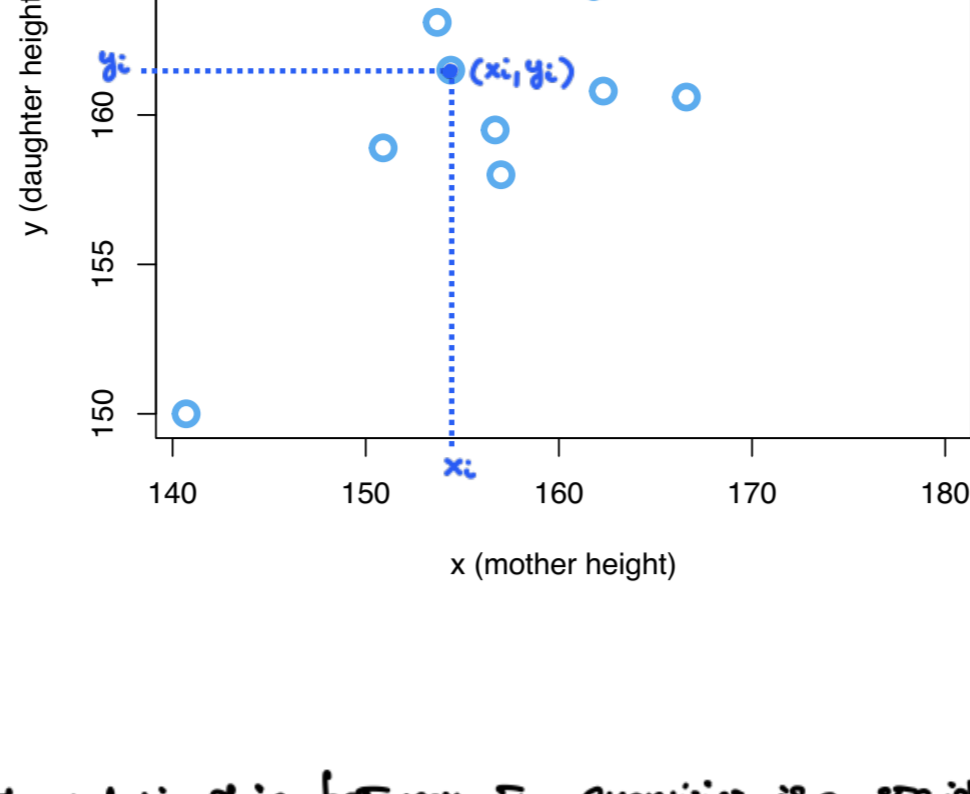
We specify a simple linear model (only 1 covariate)

We estimate the model parameters only through "intuitive" considerations and a simple optimization ("ordinary least squares" method)

We start with a simple example

relationship between the height of 11 mothers ( $x_i$ ) and the height of their daughters ( $y_i$ ).

$i$	$x$	$y$
1	153.7	163.1
2	156.7	159.5
3	173.5	169.4
4	157.0	158.0
5	161.8	164.3
6	140.7	150.0
7	179.8	170.3
8	150.9	158.9
9	154.4	161.5
10	162.3	160.8
11	166.6	160.6



Intuition:

the simplest way to describe the relationship between two quantities is a straight line:

$$Y_i = \beta_1 + \beta_2 x_i \quad i=1, \dots, n$$

However, such a relationship does not hold exactly: the points are not PERFECTLY ALIGNED.

hence we add an error term to take into account this discrepancy:

$$Y_i = \beta_1 + \beta_2 x_i + \epsilon_i \quad i=1, \dots, n$$

**1st step: MODEL SPECIFICATION**

Consider the model:

$$Y_i = \beta_1 + \beta_2 x_i + \epsilon_i \quad i=1, \dots, n$$

height of the  $i$ -th daughter      systematic component      ERROR TERM: the linear relationship is not exact

$(\beta_1, \beta_2)$  are the REGRESSION COEFFICIENTS

We specified a straight line with the intercept  $(\beta_1)$

We only observe 1 covariate, but we also introduce one additional "variable" taking value 1 for each individual.

The model matrix then is:

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

$\Rightarrow \beta_1$  is the INTERCEPT (coefficient of 1)

$\beta_2$  is the COEFFICIENT of  $x$  (slope)

**ASSUMPTIONS on the independent variables**

- $x_1, \dots, x_n$  fixed and non-stochastic
- the  $x_i$  can not be all equal (sample variance of  $(x_1, \dots, x_n)$  must be  $\neq 0$ )

The systematic component is now fully specified, we need to define the stochastic component  $(\epsilon)$ .

**ASSUMPTIONS on the STOCHASTIC COMPONENT**

- $E[\epsilon_i] = 0$  for  $i=1, \dots, n$
- $Var(\epsilon_i) = \sigma^2 > 0$   $i=1, \dots, n$  (common variance across subjects)
- $cov(\epsilon_i, \epsilon_k) = 0$  if  $i \neq k$ ,  $i, k=1, \dots, n$

- $E[\epsilon_i] = 0$   $i=1, \dots, n$  ABSENCE OF SYSTEMATIC ERROR

Implications for  $Y_i$ :  $E[\epsilon_i] = 0$

$$E[Y_i] = E[\beta_1 + \beta_2 x_i + \epsilon_i] = E[\underbrace{\beta_1 + \beta_2 x_i}_{\text{non-stochastic}}] + E[\underbrace{\epsilon_i}_0] = \beta_1 + \beta_2 x_i$$

What happens if there is a systematic error? i.e.  $E[\epsilon_i] = c \neq 0$

$$E[Y_i] = \beta_1 + \beta_2 x_i + c = (\beta_1 + c) + \beta_2 x_i$$

the systematic error  $c$  is embedded into the intercept (not a problem)

it is equivalent to a model

$$Y_i = \beta_1^* + \beta_2 x_i + \epsilon_i^* \quad \text{where } \beta_1^* = \beta_1 + c \quad \epsilon_i^* = \epsilon_i - c \Rightarrow E[\epsilon_i^*] = 0$$

- $Var(\epsilon_i) = \sigma^2 > 0$  for all  $i=1, \dots, n$  HOMOSEDNEITY OF THE ERRORS

Implications for  $Y_i$ :

$$var(Y_i) = var(\beta_1 + \beta_2 x_i + \epsilon_i) = var(\epsilon_i) = \sigma^2 \quad \forall i=1, \dots, n$$

$\Rightarrow$  homoscedasticity of the response

- $cov(\epsilon_i, \epsilon_k) = 0$  for  $i \neq k$  THE ERRORS ARE UNCORRELATED

Implication for  $Y_i$

$$cov(Y_i, Y_k) = cov(\beta_1 + \beta_2 x_i + \epsilon_i; \beta_1 + \beta_2 x_k + \epsilon_k) = cov(\epsilon_i, \epsilon_k) = 0$$

**2nd step: ESTIMATE**

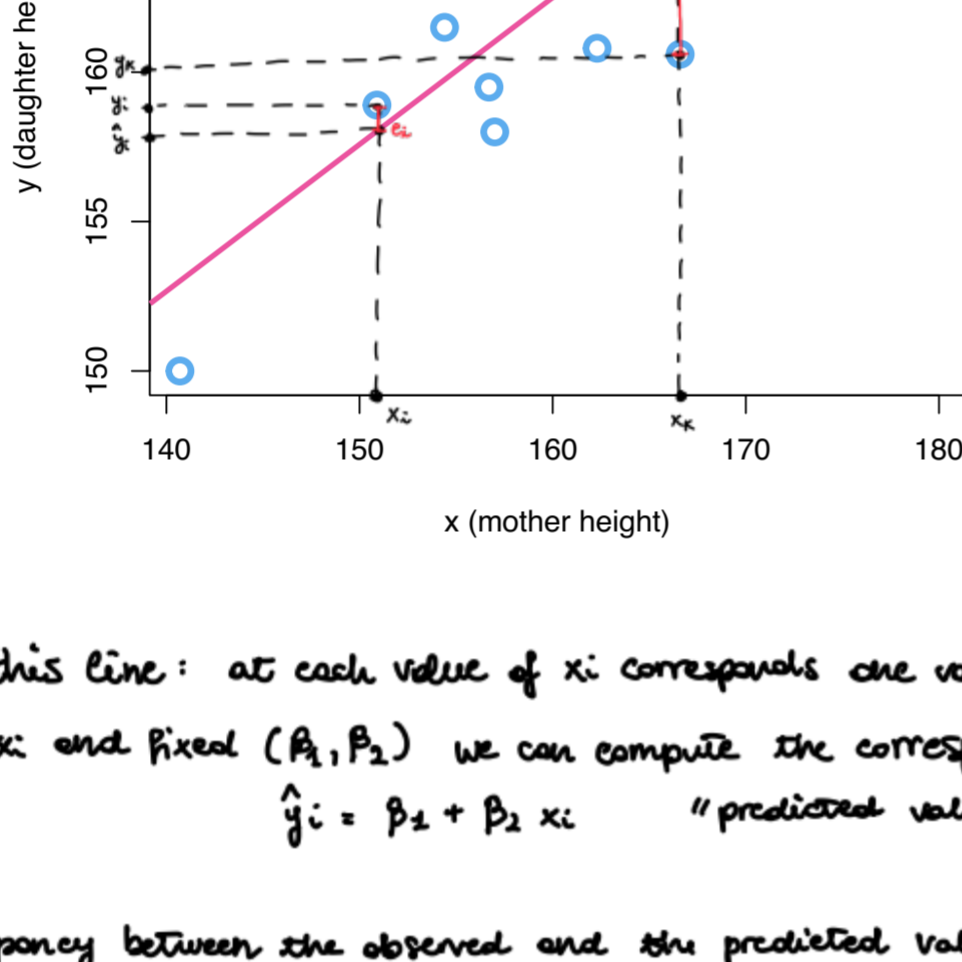
what do we need to estimate? Unknown quantities are  $(\beta_1, \beta_2, \sigma^2)$

Hence the PARAMETER SPACE is  $\Theta = \mathbb{R}^2 \times (0, +\infty)$

Every combination of  $(\beta_1, \beta_2)$  determines a specific line: how do we select the "best" one?

We need a criterion of what is a "good" line.

We want a line which is the closest to the observed points.



Consider this line: at each value of  $x_i$  corresponds one value of  $\hat{y}_i$  that lies on the line

$\Rightarrow$  given  $x_i$  and fixed  $(\hat{\beta}_1, \hat{\beta}_2)$  we can compute the corresponding value of  $\hat{y}_i$  (according to the line)

$$\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 x_i \quad \text{"predicted value"}$$

The discrepancy between the observed and the predicted value (at the observed locations  $x_i$ )

$$\text{RESIDUAL: } \epsilon_i = y_i - \hat{y}_i$$

A good line will have small residuals OVERALL.

- we could consider the sum of the residuals  $\sum_{i=1}^n \epsilon_i$  and select the  $(\beta_1, \beta_2)$  that minimize it

$\rightarrow$  not a good idea: positive and negative values cancel out.

- we could consider the sum of the absolute values  $\sum_{i=1}^n |\epsilon_i| \rightarrow$  mathematically not very practical

- we consider instead the sum of the squared residuals

$$\sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \beta_1 - \beta_2 x_i)^2 = S(\beta_1, \beta_2)$$

and take as an estimate of  $(\beta_1, \beta_2)$  the combination that minimizes it.

**DEF:** the LEAST SQUARES estimate of  $(\beta_1, \beta_2)$  is the combination of values  $(\hat{\beta}_1, \hat{\beta}_2)$  that minimizes  $S(\beta_1, \beta_2)$

$$\begin{aligned} (\hat{\beta}_1, \hat{\beta}_2) &= \underset{(\beta_1, \beta_2) \in \mathbb{R}^2}{\text{argmin}} S(\beta_1, \beta_2) \\ &= \underset{(\beta_1, \beta_2) \in \mathbb{R}^2}{\text{argmin}} \sum_{i=1}^n (y_i - \beta_1 - \beta_2 x_i)^2 \end{aligned}$$

We have hence turned a problem of estimation into an optimization.

**THEM:** The least squares estimate of  $(\beta_1, \beta_2)$  is

$$\begin{aligned} \hat{\beta}_1 &= \bar{y} - \hat{\beta}_2 \bar{x} \\ \hat{\beta}_2 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{aligned}$$

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

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \quad (\text{sample mean}).$$

Remark:

recall that the sample variance of  $(x_1, \dots, x_n)$  is  $s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

(and similarly for  $s_y^2$ )

the sample covariance is  $s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$

$$\text{Hence } \hat{\beta}_2 = \frac{s_{xy}}{s_x^2}$$

**Proof:** we want to show that  $\hat{\beta}_1, \hat{\beta}_2$  minimize  $S(\beta_1, \beta_2) = \sum_{i=1}^n (y_i - \beta_1 - \beta_2 x_i)^2$ .

We need to find the critical points ( $1^{st}$  derivative = 0)

and then check that  $(\hat{\beta}_1, \hat{\beta}_2)$  is a minimum ( $2^{nd}$  derivative  $> 0$ )

$$\begin{cases} \frac{\partial S(\beta_1, \beta_2)}{\partial \beta_1} = 0 \\ \frac{\partial S(\beta_1, \beta_2)}{\partial \beta_2} = 0 \end{cases} \Leftrightarrow \begin{cases} \sum_{i=1}^n 2(y_i - \beta_1 - \beta_2 x_i)(-1) = 0 \\ \sum_{i=1}^n 2(y_i - \beta_1 - \beta_2 x_i)(-x_i) = 0 \end{cases}$$

$$\begin{cases} \sum_{i=1}^n (y_i - \beta_1 - \beta_2 x_i) = 0 & \text{①} \\ \sum_{i=1}^n x_i (y_i - \beta_1 - \beta_2 x_i) = 0 & \text{②} \end{cases}$$

$$\text{① } n\bar{y} - n\beta_1 - \beta_2 \sum_{i=1}^n x_i = 0 \quad (\text{since } \sum_{i=1}^n y_i = n\bar{y})$$

$$\beta_1 = \bar{y} - \beta_2 \bar{x}$$

$$\text{② } \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} - \beta_2 \sum_{i=1}^n x_i^2 = 0$$

$$\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} - \beta_2 \sum_{i=1}^n x_i^2 = 0$$

$$\hat{\beta}_2 = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2}$$

$$\text{we obtain } \hat{\beta}_2 = \frac{s_{xy}}{s_x^2}$$

$$\text{and } \hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}$$

Is  $(\hat{\beta}_1, \hat{\beta}_2)$  a minimum? We compute the Hessian

$$H = \begin{bmatrix} \frac{\partial^2 S(\beta_1, \beta_2)}{\partial \beta_1^2} & \frac{\partial^2 S(\beta_1, \beta_2)}{\partial \beta_1 \partial \beta_2} \\ \frac{\partial^2 S(\beta_1, \beta_2)}{\partial \beta_1 \partial \beta_2} & \frac{\partial^2 S(\beta_1, \beta_2)}{\partial \beta_2^2} \end{bmatrix} = \begin{bmatrix} 2n & 2n\bar{x} \\ 2n\bar{x} & 2 \sum_{i=1}^n x_i^2 \end{bmatrix}$$

$$\det(H) = 4n \sum_{i=1}^n x_i^2 - 4n^2 \bar{x}^2 = 4n \sum_{i=1}^n (x_i - \bar{x})^2 > 0$$

since  $\det(H) > 0$  and  $H_{1,1} = 2n > 0$ ,  $(\hat{\beta}_1, \hat{\beta}_2)$  is a minimum of  $S(\beta_1, \beta_2)$

Moreover, it is the global minimum.

Remarks:

- we did not use the assumptions on  $\epsilon_i$

- we used the assumption on the  $x_i$ : what happens if  $x_i = x_0$  for all  $i=1, \dots, n$ ?

$$(x_i - \bar{x}) = 0 \quad \forall i \Rightarrow s_x^2 = 0 \quad \text{and } s_{xy} = 0 \Rightarrow \hat{\beta}_2 = \frac{0}{0} \quad \text{not defined}$$

- once we estimate  $(\hat{\beta}_1, \hat{\beta}_2)$ , we automatically obtain  $\hat{y} = \hat{\beta}_1 + \hat{\beta}_2 x$ , i.e. the estimated regression line.

-  $\hat{y}$  allows us to make predictions: given a generic value  $x_1$ , we predict the corresponding value of the response.

As usual, careful with extrapolation, i.e., estimating the response for a value of  $x$  outside of the observed range of  $(x_1, \dots, x_n)$ .

**INTERPRETATION of  $(\hat{\beta}_1, \hat{\beta}_2)$**

we have estimated a line  $\hat{y} = \hat{\beta}_1 + \hat{\beta}_2 x$

$\hat{\beta}_1$  is the intercept, i.e., the predicted value of  $y$  when  $x=0$ .

Not always interpretable! E.g. with the heights example: height = 0 is meaningless

Now consider two individuals observed at  $x_1 = x_0$  and  $x_2 = x_0 + 1$

The predicted values are

$$\hat{y}_1 = \hat{\beta}_1 + \hat{\beta}_2 x_0$$

$$\hat{y}_2 = \hat{\beta}_1 + \hat{\beta}_2 (x_0 + 1)$$

let's study the difference in their predicted values

$$\begin{aligned} \hat{y}_2 - \hat{y}_1 &= \hat{\beta}_1 + \hat{\beta}_2 (x_0 + 1) - \hat{\beta}_1 - \hat{\beta}_2 x_0 \\ &= \hat{\beta}_2 x_0 + \hat{\beta}_2 - \hat{\beta}_2 x_0 \\ &= \hat{\beta}_2 \end{aligned}$$

Hence  $\hat{\beta}_2$  is the expected change in  $y$  when I increase  $x$  of 1 unit

