

(Nov 17 2021)

## Two quantitative variables

Pairs of values  $(x_i, y_i)$ ,  $i = 1, \dots, n$

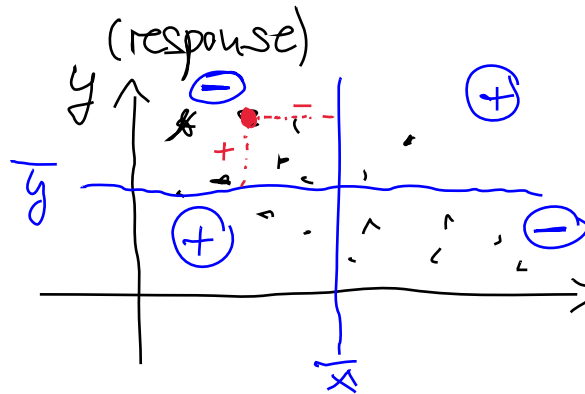
Arithmetic means:  $\bar{x}$ ,  $\bar{y}$

Variances:

$$S_{xx} = S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$S_{yy} = S_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

## Scatterplot



explanatory  
(if there is  
a direction)

## Covariance

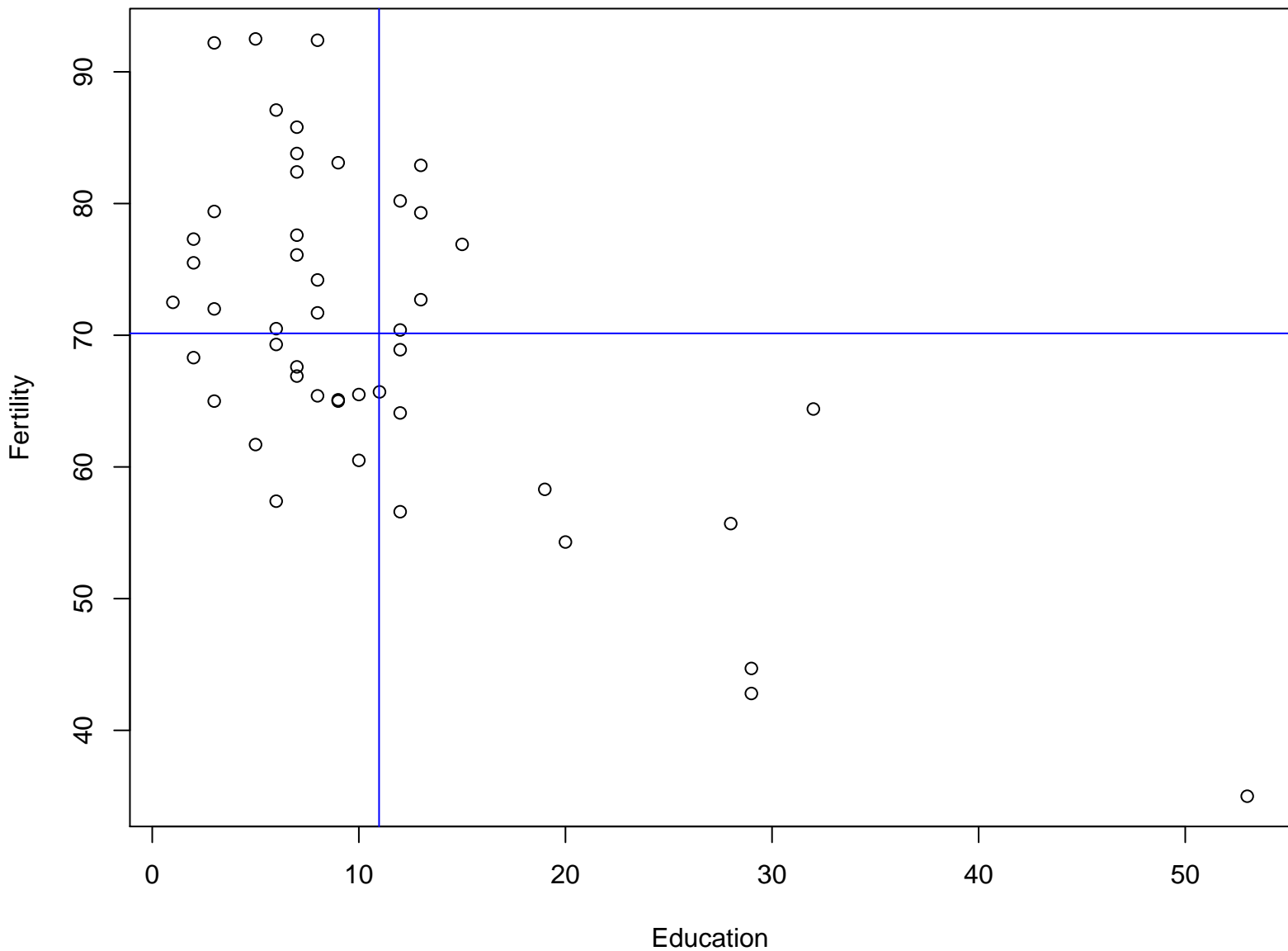
$$S_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- Sign shows direction of relation,  
e.g. Swiss data: higher Education  
comes with lower Fertility
- the absolute value is unscaled and  
thus hard to interpret
- covariances of a group of variables  
can be collected in a matrix:

$R$ : cov

$$\begin{array}{c}
 V_1 \quad \dots \quad V_p \\
 \begin{array}{cccc}
 V_1 & S_{11}^2 & S_{12} & \dots & S_{1p} \\
 \vdots & S_{21} & \ddots & \ddots & S_{2p} \\
 \vdots & \vdots & \ddots & \ddots & \vdots \\
 V_p & S_{p1} & S_{p2} & \dots & S_{pp}^2
 \end{array}
 \end{array}$$

symmetric,  
 $S_{ij} = S_{ji}$



Correlation: Normalized covariance

Bravais-Pearson correlation

$$r_{xy} = \frac{S_{xy}}{S_x S_y}$$

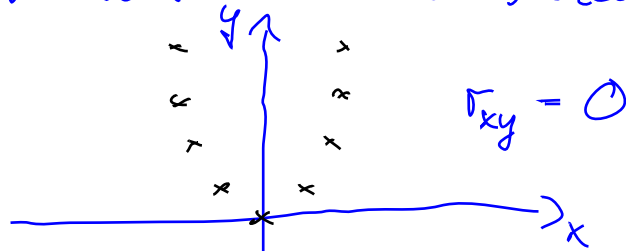
- between -1 and +1

-  $r_{xx} = r_{yy} = 1$

- -1: all values are on straight line with negative slope

+1: ... with positive slope

0: the relation between the two variables has no linear "portion"



- don't interpret correlation from small samples without a scatter plot

The Bravais-Pearson correlation measures the direction and strength of a linear relation between  $x$  and  $y$ .

The linear relation: (Simple) linear regression

Model:  $Y = \underbrace{\beta_0 + \beta_1 \cdot X}_{\substack{\text{straight line} \\ \text{intercept } \beta_0 \\ \text{slope } \beta_1}} + \underbrace{\varepsilon}_{\substack{\uparrow \text{random} \\ \text{variation}}}$

Data: determine suitable values ("estimates")  $\hat{\beta}_0$  and  $\hat{\beta}_1$  for  $\beta_0$  and  $\beta_1$   
 $\Rightarrow$  estimated straight line

The most common estimation method:  
Ordinary Least Squares (OLS, LS)

$y_i$   
 actual  
 $y_i$

$b_0 + b_1 x_i$   
 value of the  
 line based on  $x_i$  with  $b_0$  and  $b_1$

The OLS  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are obtained from minimizing

$$L(b_0, b_1) = \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2 \text{ w.r.t. } b_0 \text{ and } b_1$$

Minimization: take derivatives and set equal to 0

$$\frac{\partial L(b_0, b_1)}{\partial b_0} = -\sum_{i=1}^n 2(y_i - (b_0 + b_1 x_i))$$

$\vdots$   
 $\dots \rightarrow$  "Regression" for general  
 (2nd sem.) case

Result:  $\hat{\beta}_1 = \frac{S_{xy}}{S_x^2}$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

Because of the intercept:

$$\overline{\hat{y}} = \bar{y}$$

$x_i$	$y_i$	$x_i - \bar{x}$	$y_i - \bar{y}$	$\hat{y}_i$ (from R)	$y_i - \hat{y}_i$ (errors from R)
2	7	-2	4	3.09	3.90
4	1	0	-2	3	-2
3	2	-1	-1	3.045	-1.045
8	4	4	1	2.81	1.18
3	1	-1	-2	3.045	-2.045

$$\bar{x} = 4 \quad \bar{y} = 3$$

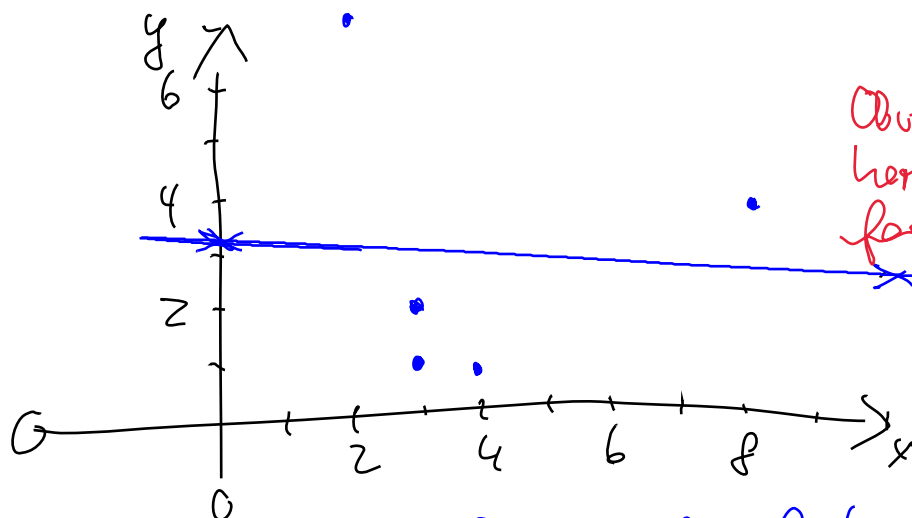
$$s_x^2 = \frac{1}{5-1} \cdot (4+0+1+16+1) = \frac{22}{4} = 5.5$$

$$s_y^2 = \frac{1}{5-1} (16+4+1+1+4) = \frac{26}{4} = 6.5$$

$$s_{xy} = \frac{1}{5-1} ((-2) \cdot 4 + 0(-2) + (-1) \cdot (-1) + 4 \cdot 1 + (-1) \cdot (-2)) = \frac{-1}{4} = -0.25$$

$$\hat{\beta}_1 = \frac{-0.25}{5.5} = -\frac{1}{22} = -0.045$$

$$\hat{\beta}_0 = 3 - \left(-\frac{1}{22}\right) \cdot 4 = 3 + \frac{4}{22} = 3.18$$



Obviously a  
horrible model  
for those fake  
data.

Draw line by calculating two points:  
 $(0, 3.18)$   
 $(10, 3.18 - 10 \cdot 0.045)$

## Quality of fit of a straight line

Decompose the sum of squares of  $y$  into "Model" and "Error":

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

Total
Model ( $\geq 0$ )
Error ( $\geq 0$ )

Model should be large  
Error " " small

Coefficient of determination  $R^2$ :

$$R^2 = \frac{\text{Model SS}}{\text{Total SS}}$$

$R^2 = 1 = 100\%$  : perfect

$R^2 = 0$  : no explanatory value  
at all

Convenient:  $R^2 = r_{xy}^2$

Correlation and  $R^2$  for example:

$$r_{xy} = \frac{s_{xy}}{s_x \cdot s_y} = \frac{-0,25}{\sqrt{5,5} \cdot \sqrt{6,5}} = -0,0418121$$

$$R^2 \approx 0.00175 = 0.175\% \text{ awful}$$

SS decomposition: (from R)

$$\begin{array}{rclcl} \text{Total SS} & = & \text{Model SS} & + & \text{Error SS} \\ 26 & = & 0.045 & + & 25.954 \end{array}$$