

Correlation vs Regression

Oliver Snellman
oliver.snellman@gmail.com

May 2022

Abstract

Correlation and regression coefficients are both tools to describe a certain type of relation between two variables. Here I provide a visual intuition on how they differ.

Pearson's correlation coefficient ρ measures, how much two datasets, X and Y, co-vary linearly. **Linear Regression coefficient β_1** measures, how one variable Y linearly depends on another variable X.

$$\rho = \frac{\mathbb{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} \in [-1, 1]$$

where $\mu_X = \mathbb{E}[X]$ is the expected value and $\sigma_X = \mathbb{E}[(X - \mu_X)^2]^{\frac{1}{2}}$ is the standard deviation of X, and similarly for Y.

Correlation coefficient ρ **describes**, how tightly the points (X,Y) align to a single line, when plotted on a 2D plane.

ρ **does not** tell anything about the relative sizes, that the values of X and Y might have (that is what regression does).

$$Y = \beta_0 + \beta_1 X + \varepsilon$$
$$\beta_0, \beta_1 \in \mathbb{R}$$

where β_1 is the regression coefficient, β_0 is an intercept and ε is an error term.

Regression coefficient β_1 **predicts**, how much the value of Y is expected to increase, when X increases by one.

β_1 **does not** tell anything about the accuracy of the prediction; how close the values of Y in a dataset are from the prediction (this is expressed by ρ and ε).

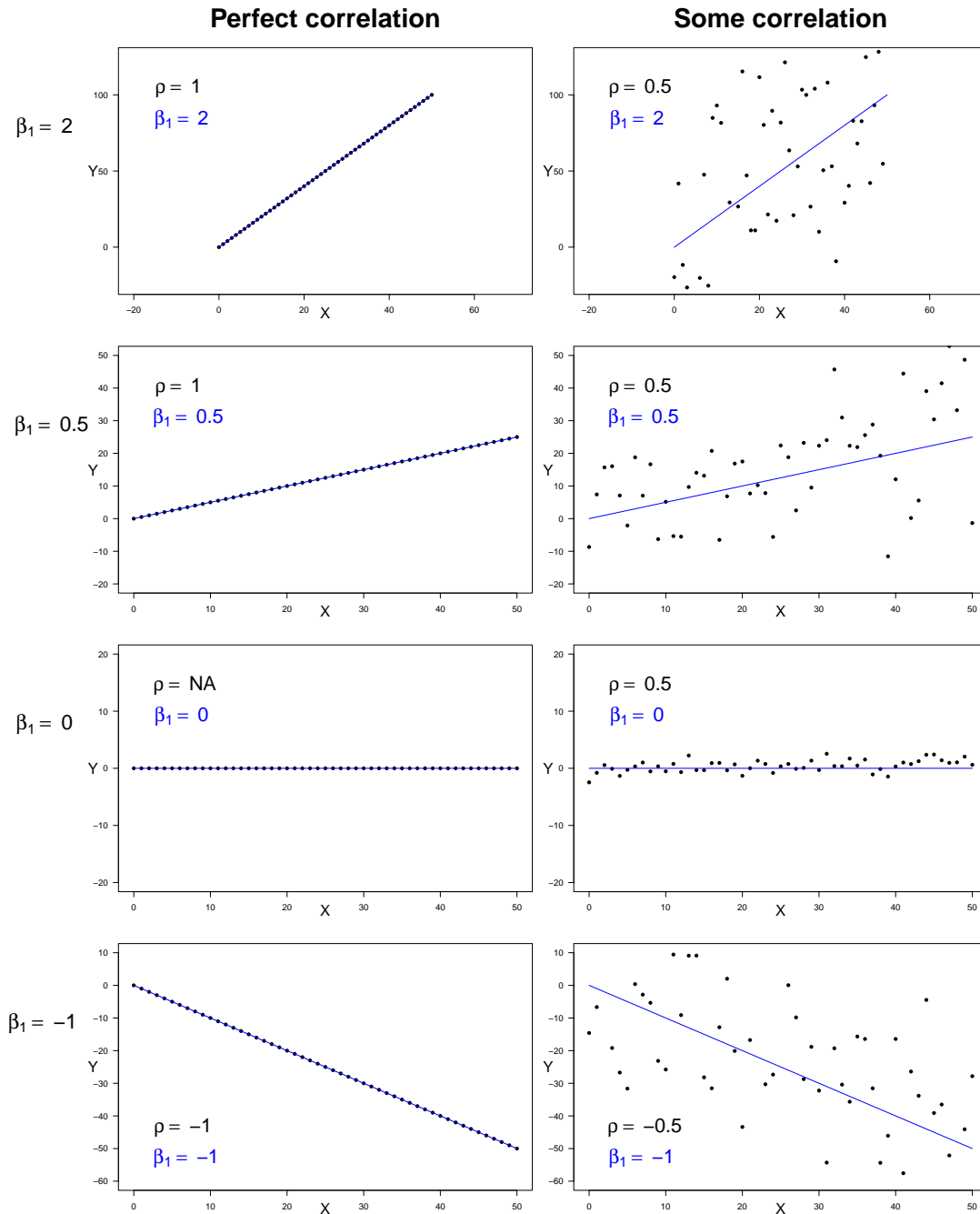


Figure 1: Each scatterplot shows 51 observations with X and Y values, simulated with R. In the **left column** the correlation ρ equals one for all points in any straight ascending line, it is undefined for the horizontal line because of a zero in the denominator, and it equals minus one for all descending straight lines. **Each row** shows, that the regression coefficient β_1 (the slope of the blue line) can have the same value for differently correlated variables. I left out the case of a straight vertical line, because it would require more explaining.