

# Correlation vs Regression

Oliver Snellman  
oliver.snellman@gmail.com

May 2022

## Abstract

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

**Pearson's correlation coefficient  $\rho$**  expresses how much two variables, X and Y, co-vary *linearly*. **Linear Regression coefficient  $\beta_1$**  expresses how one variable Y depends on another variable X *linearly*.

$$\rho = \frac{\mathbb{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} \in [-1, 1] \qquad Y = \beta_0 + \beta_1 X + \varepsilon \qquad (1)$$

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. where  $\beta_0, \beta_1 \in \mathbb{R}$ ,  $\beta_0$  is an intercept and  $\varepsilon$  is an iid error term.

Correlation coefficient  **$\rho$  describes** how tightly the points (X,Y) align on a plane. Regression coefficient  **$\beta_1$  predicts** the increase in Y, when X increases by one.

**$\rho$  does not** tell us about the relative sizes of X and Y, that is what  $\beta_1$  does.  **$\beta_1$  does not** tell us how close the values of Y are from the prediction line, that is expressed by  **$\rho$**  and  $\varepsilon$ .

Next I visualize how variables can perfectly correlate  $|\rho| = 1$ , while the regression coefficient  $\beta_1$  differs. Likewise, the same value of  $\beta_1$  can occur for differing  $\rho$ .

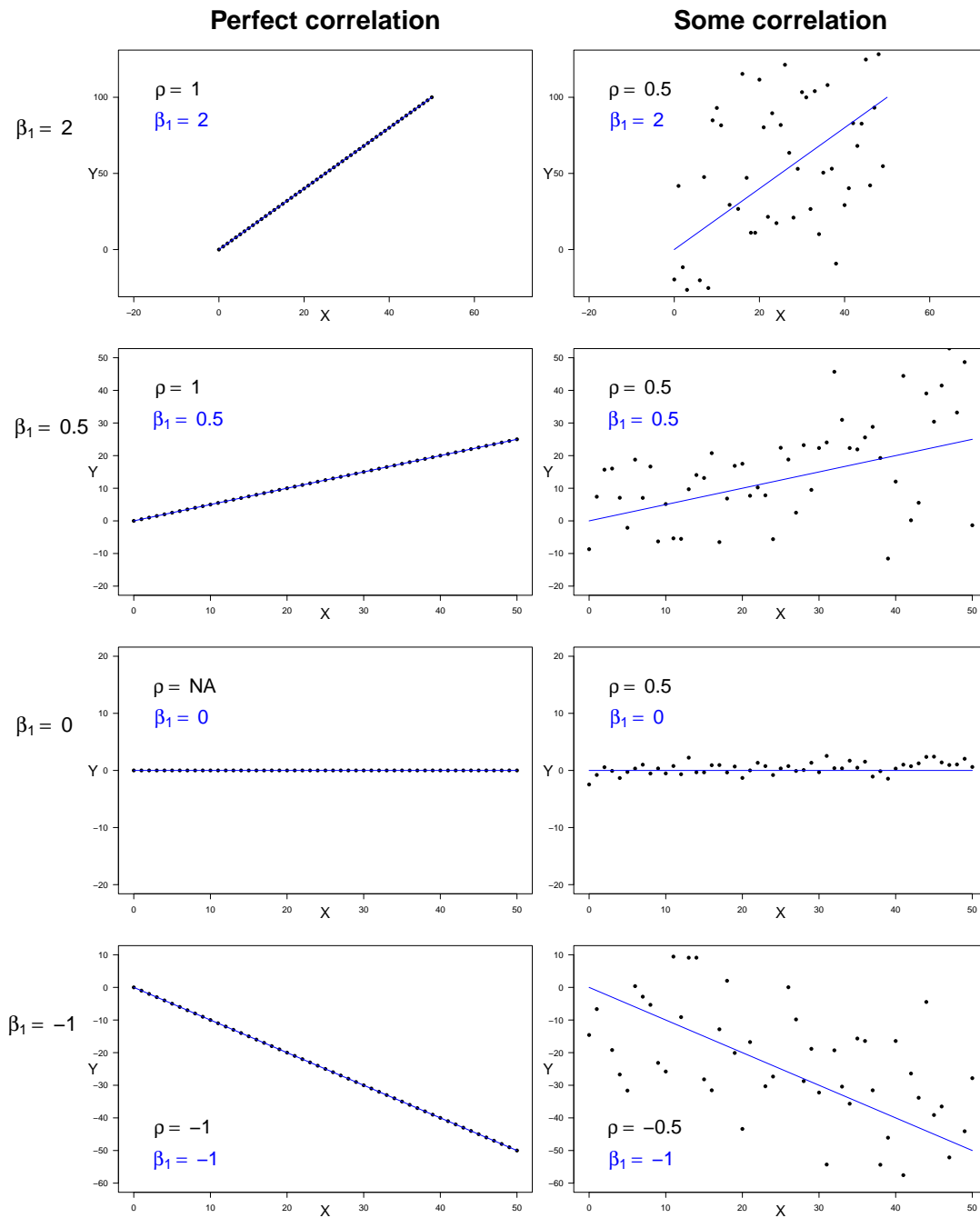


Figure 1: Each scatterplot has 51 observations with  $X$  and  $Y$  values, simulated from a linear model in Eq 1 without and with an error term. In the **left column** (no error terms) we can see that the correlation  $\rho$  equals one for points in all straight ascending lines, it is undefined for the horizontal line because of a zero in the denominator, and it equals minus one for all descending straight lines. **The rows** illustrate that the regression coefficient  $\beta_1$  (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.