

Introduction to Linear Modelling

NGSchool 2022

Nuno Sepúlveda, 16.09.2022
nuno.Sepulveda@mini.pw.edu.pl
<http://www.immune-stats.net>

Your thoughts about linear regression

Why learning linear regression?

Why learning linear regression?

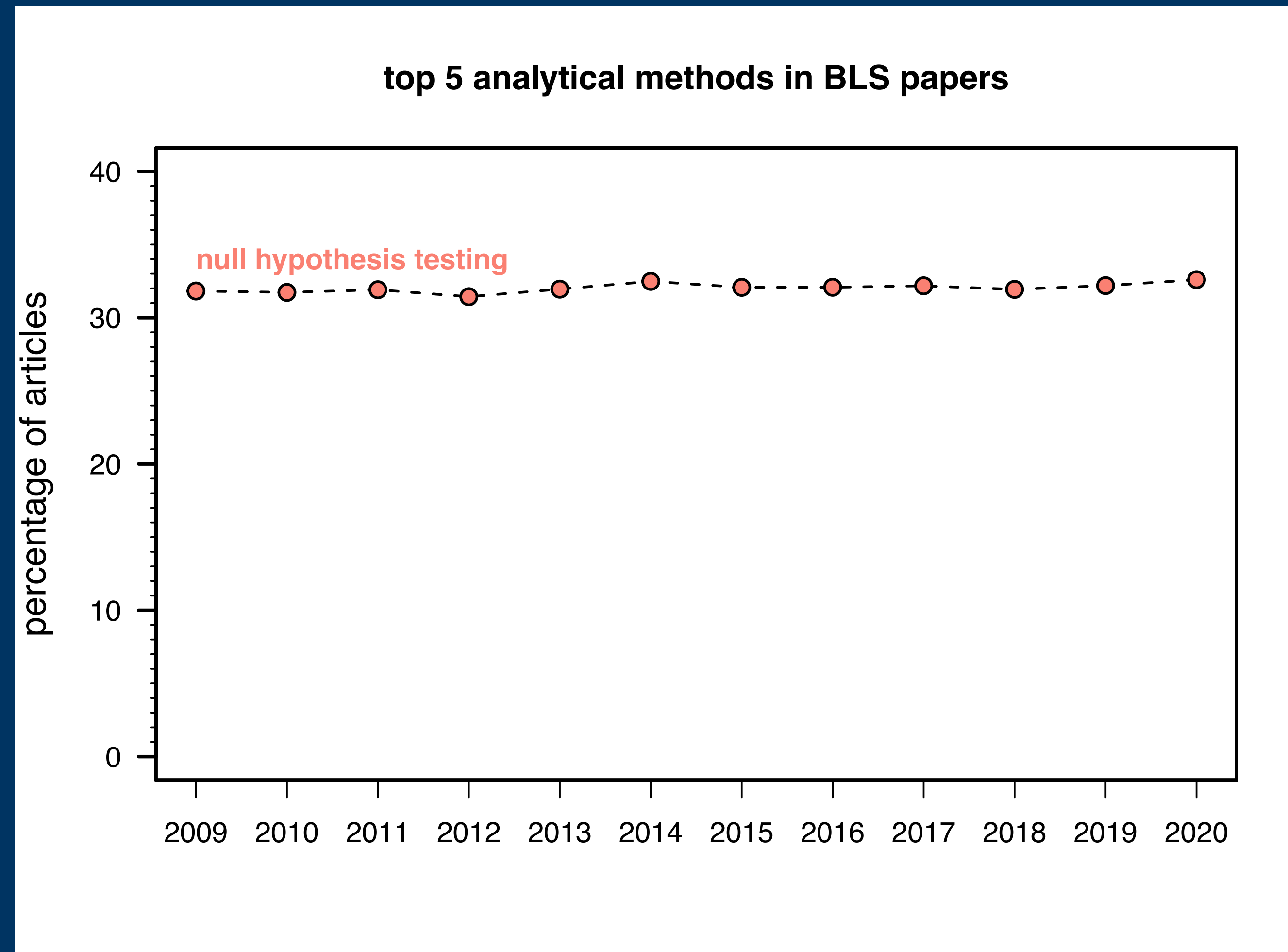
PLOS BIOLOGY

META-RESEARCH ARTICLE

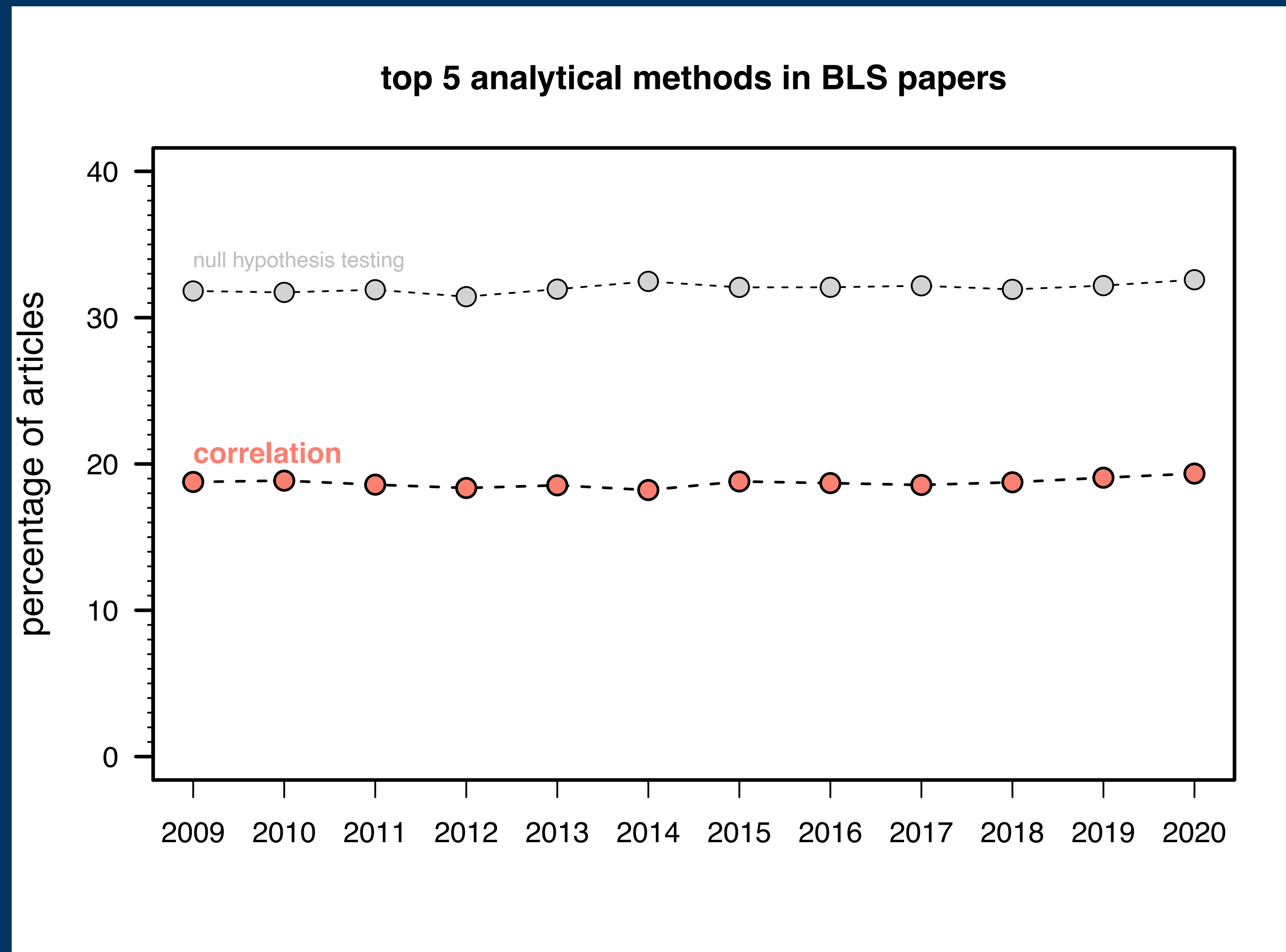
Educating the future generation of researchers: A cross-disciplinary survey of trends in analysis methods

Taylor Bolt^{1*}, Jason S. Nomi¹, Danilo Bzdok^{2,3}, Lucina Q. Uddin^{1,4}

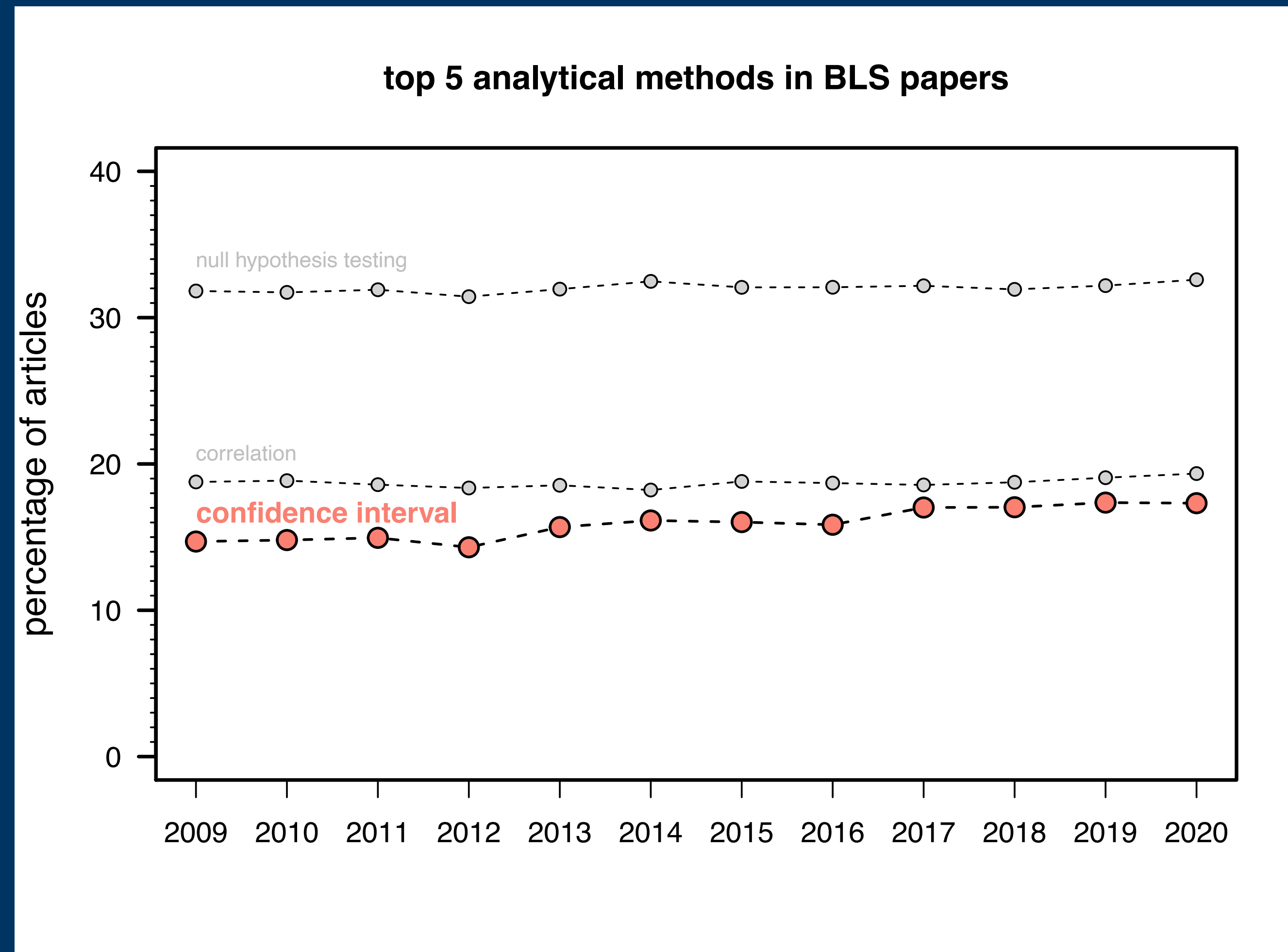
Why learning linear regression?



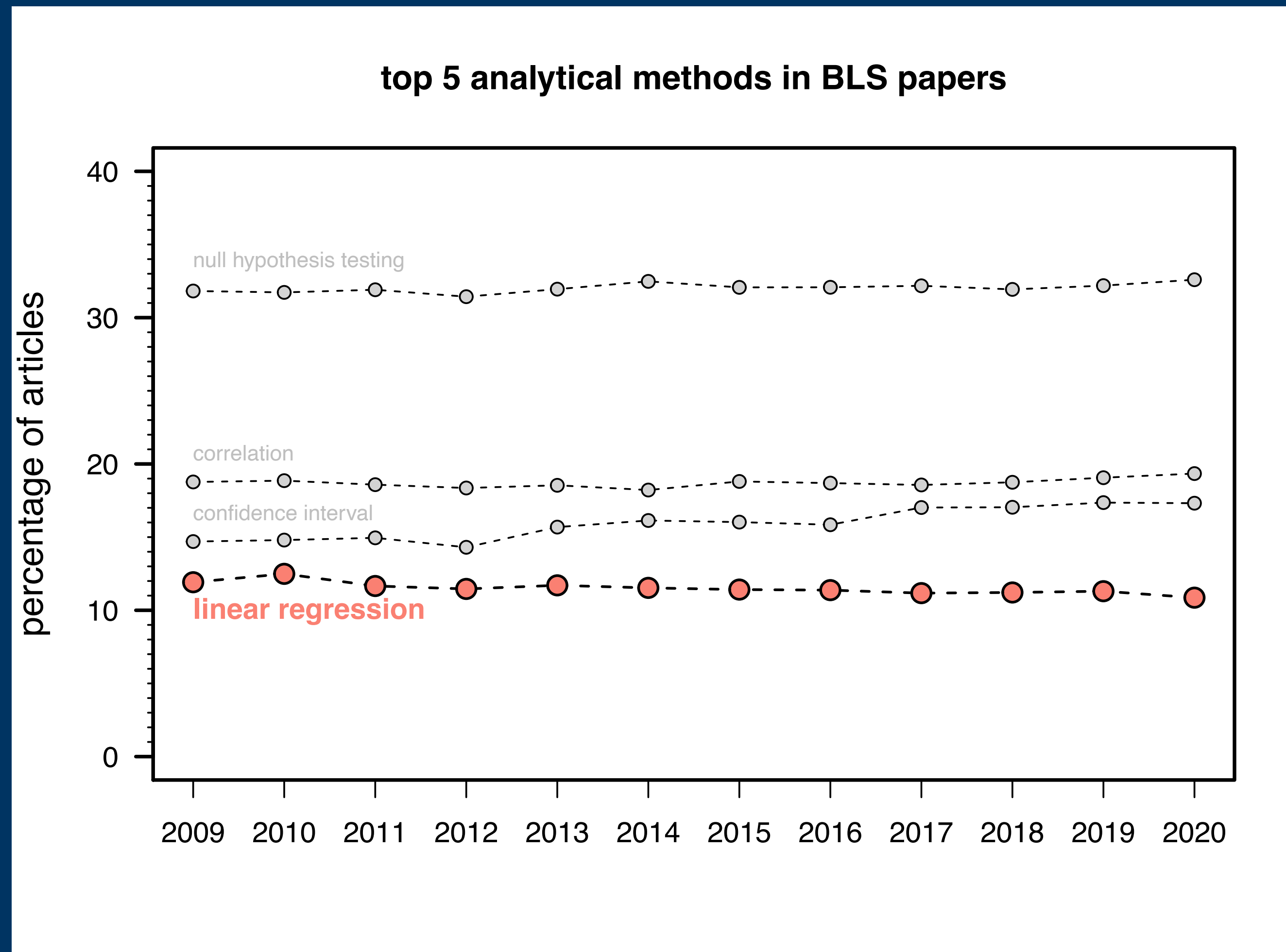
Why learning linear regression?



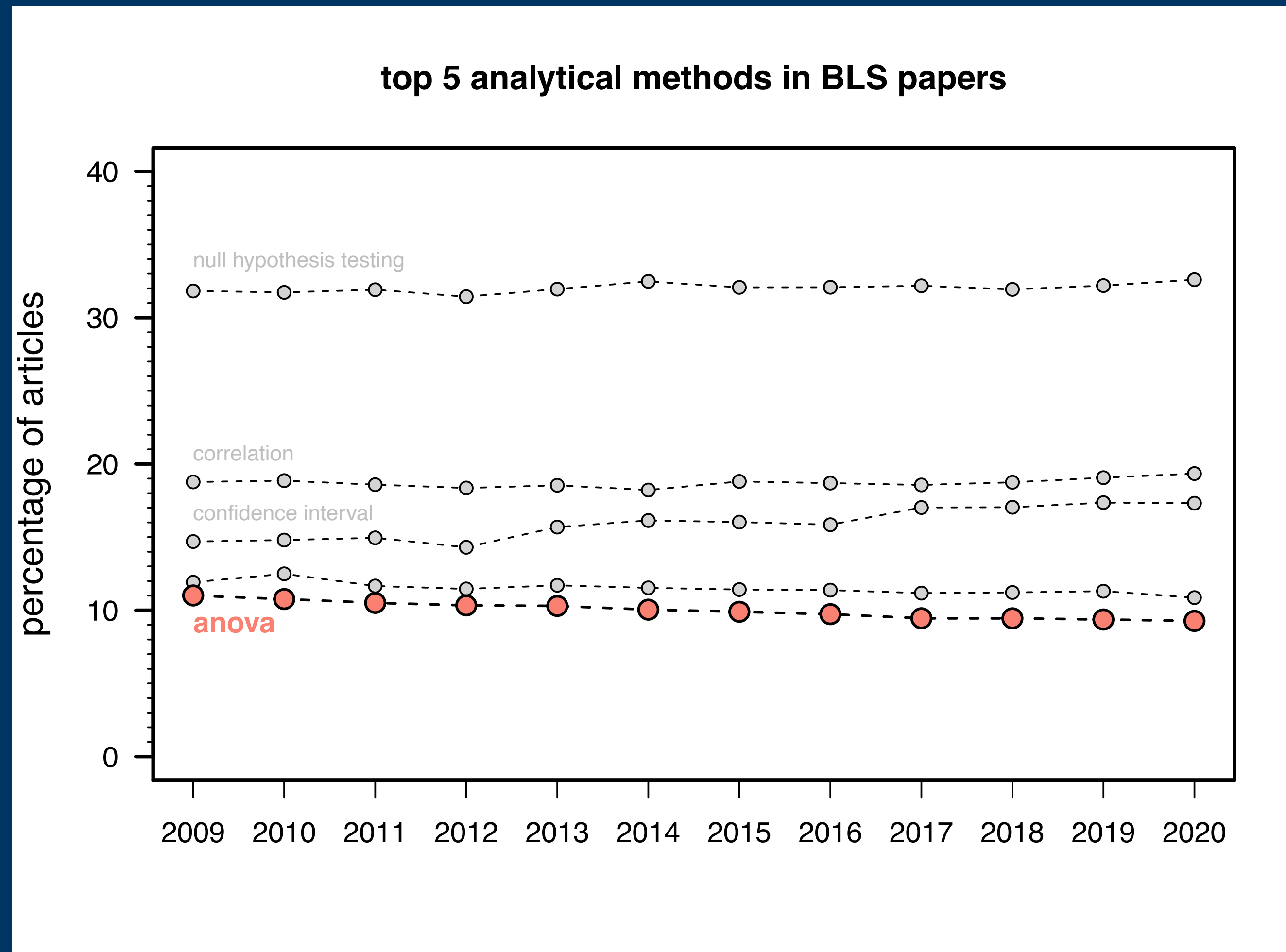
Why learning linear regression?



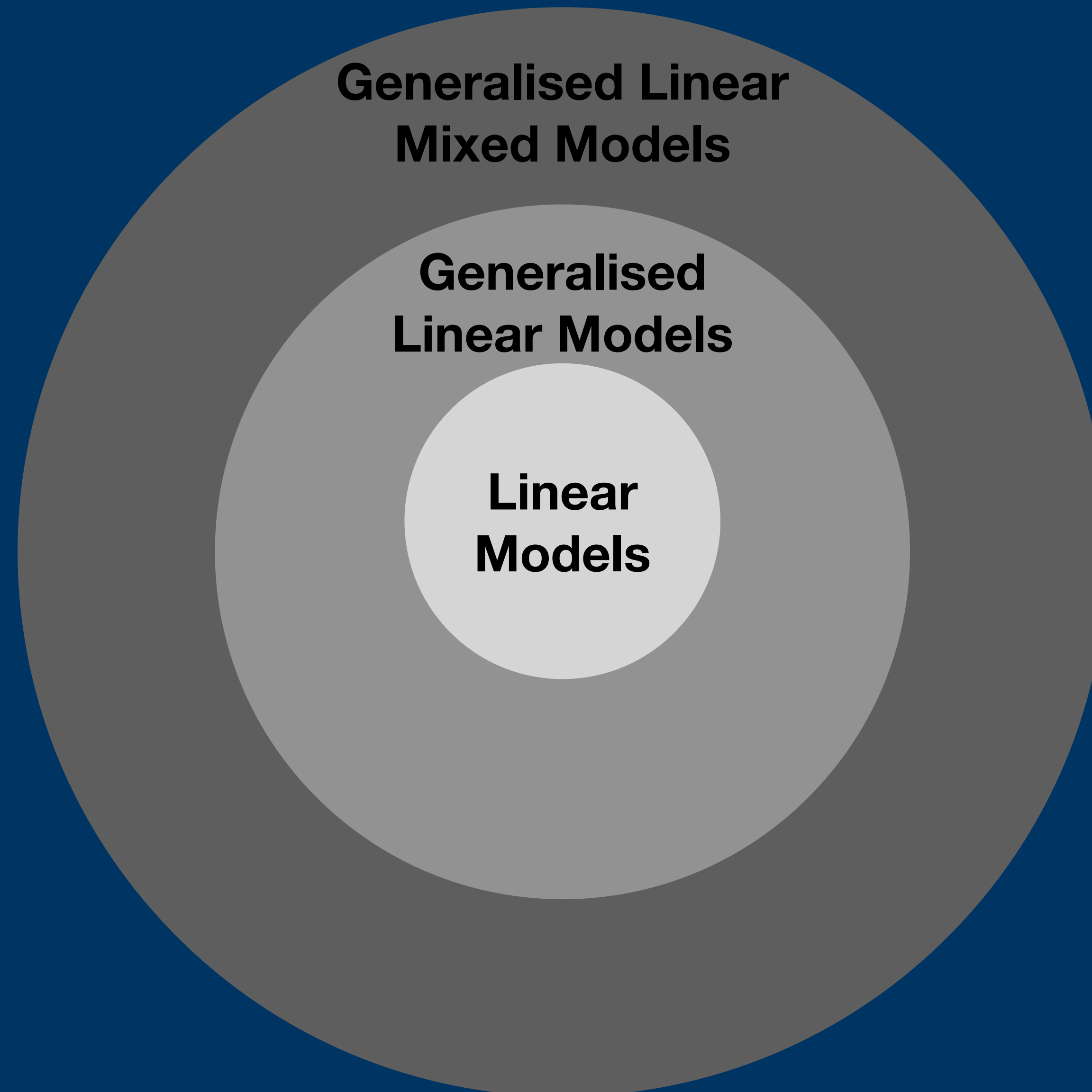
Why learning linear regression?



Why learning linear regression?



Linear regression as a foundation for more complex models



Objectives

Touch-base on simple linear regression

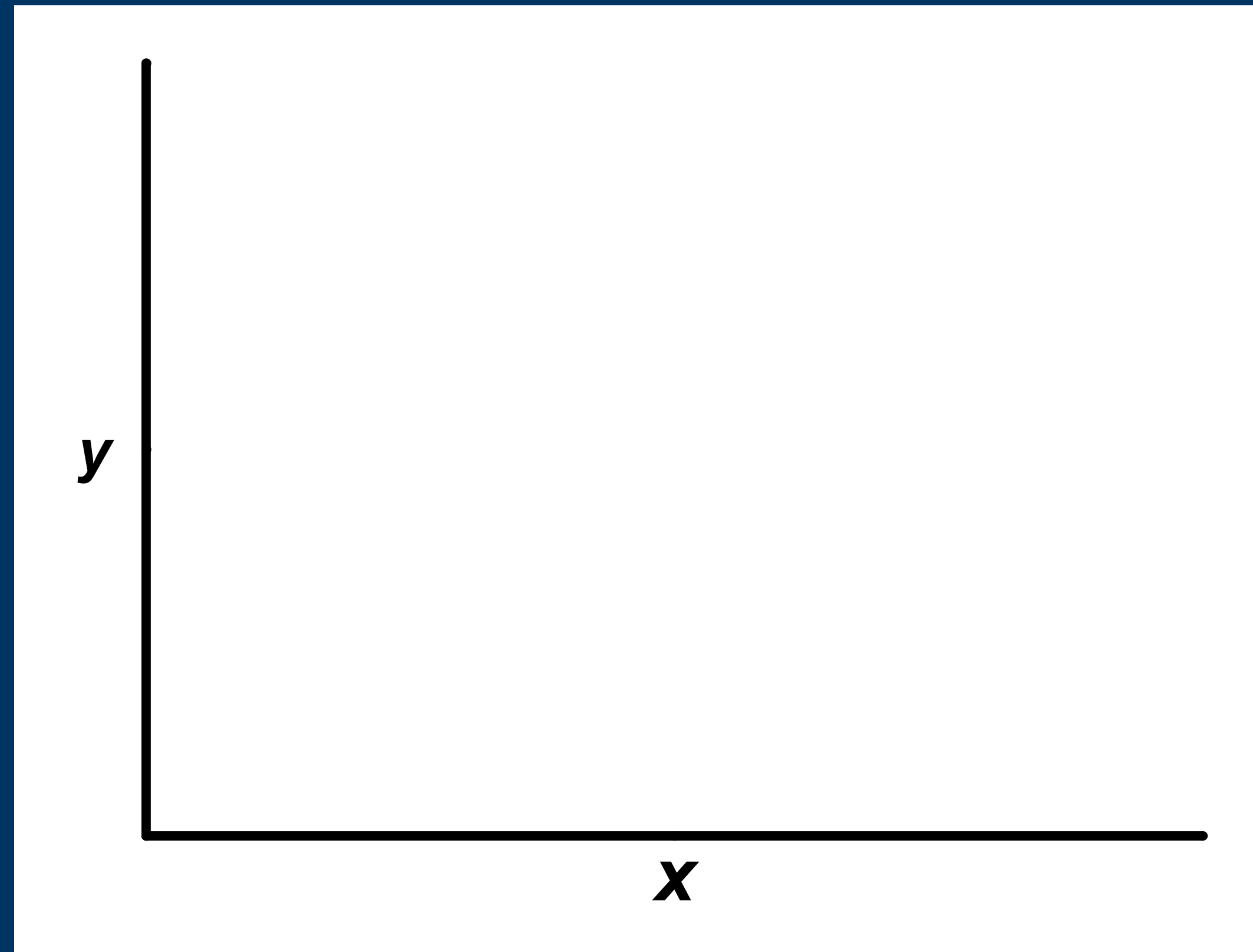
How to estimate parameters?

How to interpret the results?

How to be cautious in reporting?

Use R to conduct data analysis

Study of the relationship between x and y



Simple linear regression

$$y = a + bx$$

y = response variable

dependent variable

outcome variable

x = explicative variable

independent variable

covariate/predictor

feature

Simple linear regression

$$y = a + bx$$

a = intercept

b = slope

Simple linear regression

$$y = a + bx$$

a = intercept

b = slope

$$y = \alpha + \beta x$$

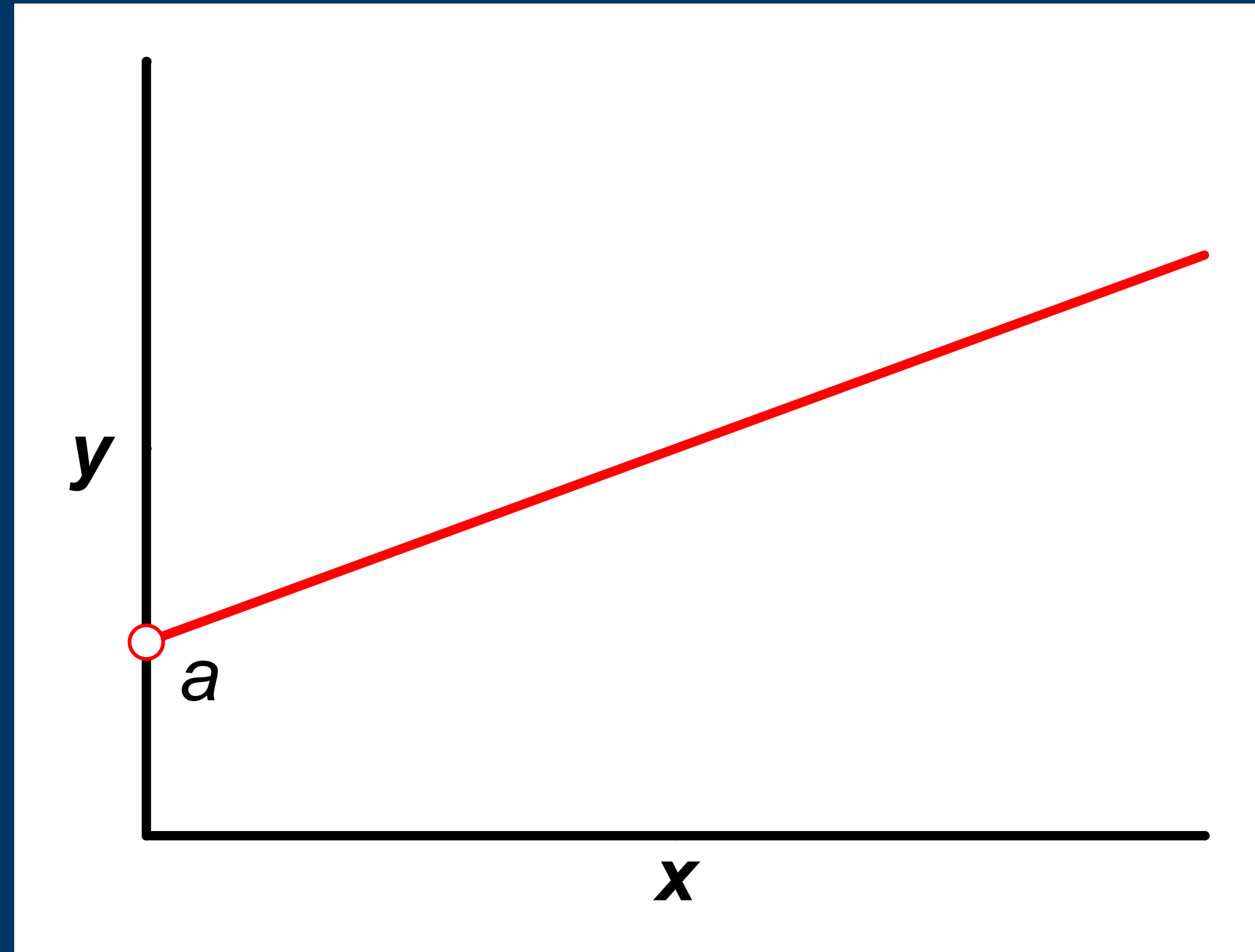
Footnote

$$y = \beta_0 + \beta_1 x$$

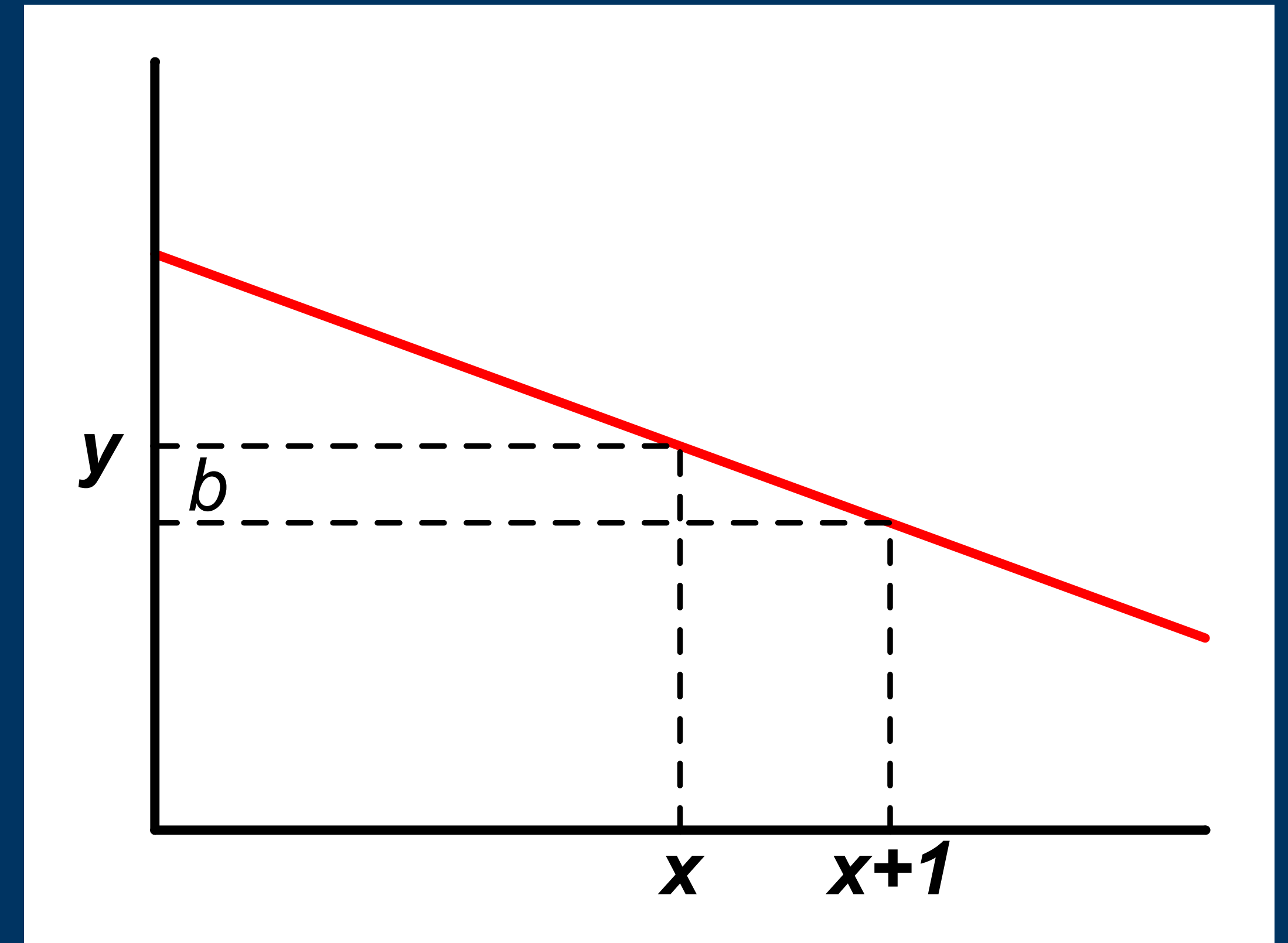
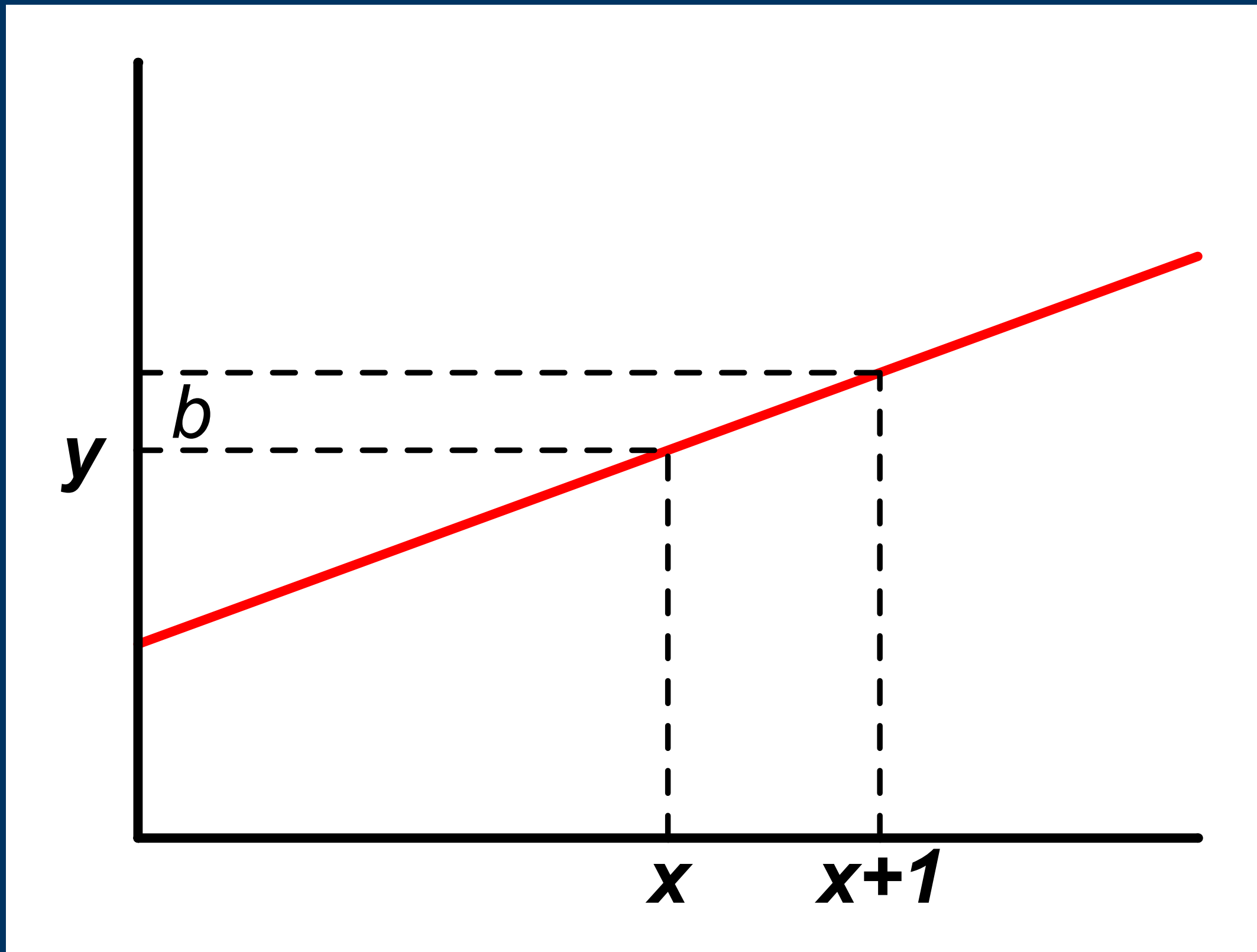
Statisticians/Mathematicians

What is the advantage of writing the equation like this?

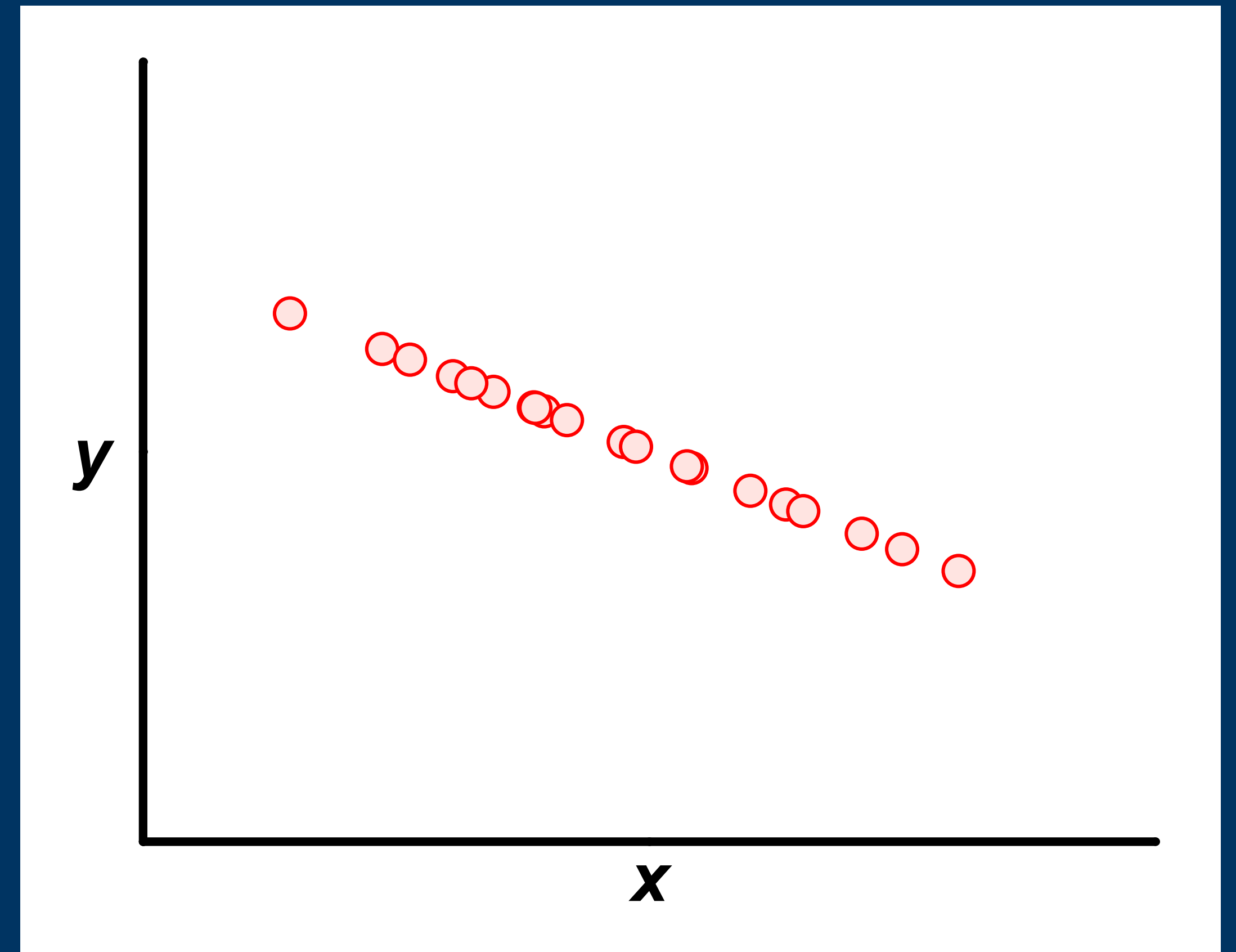
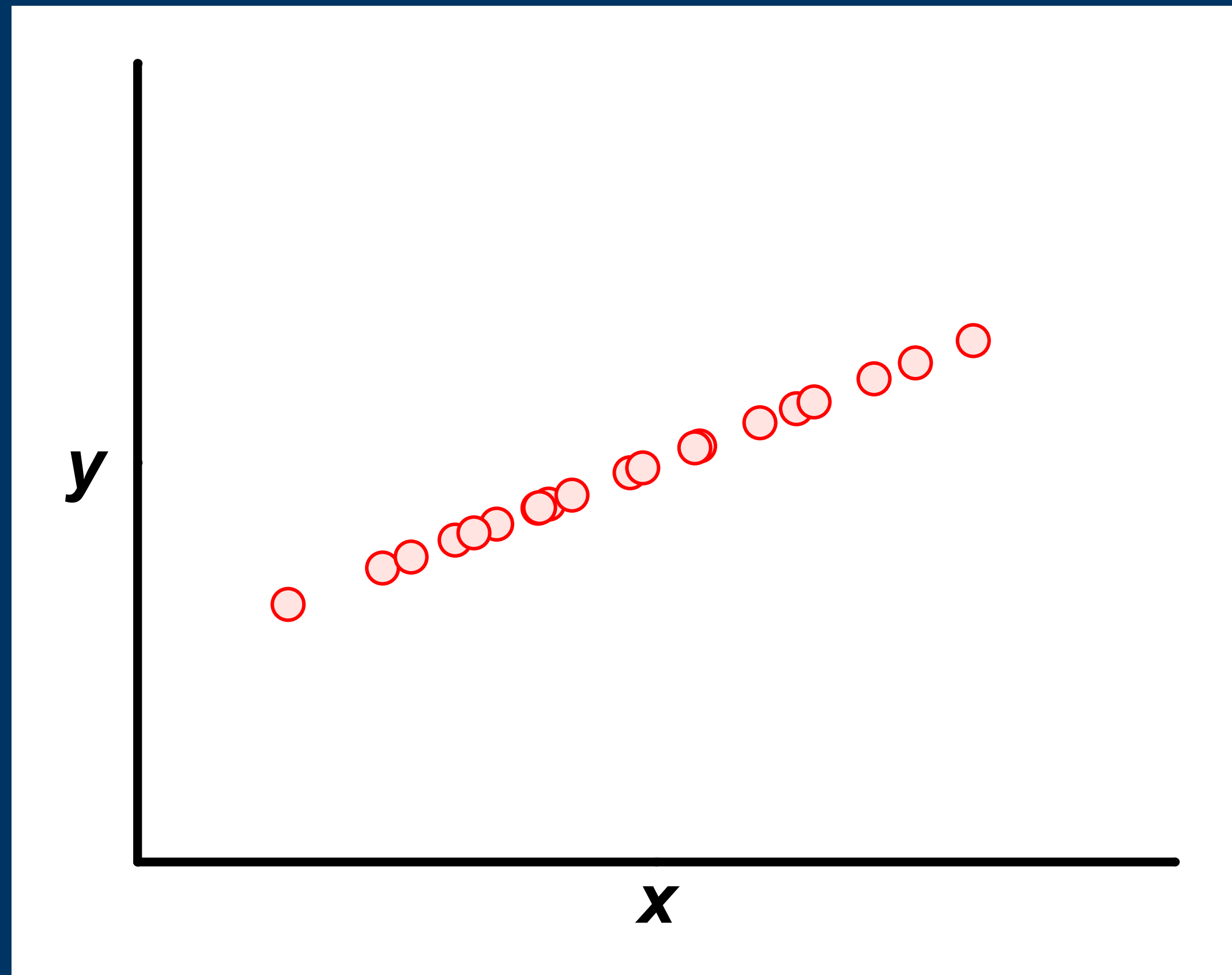
What does the intercept represent?



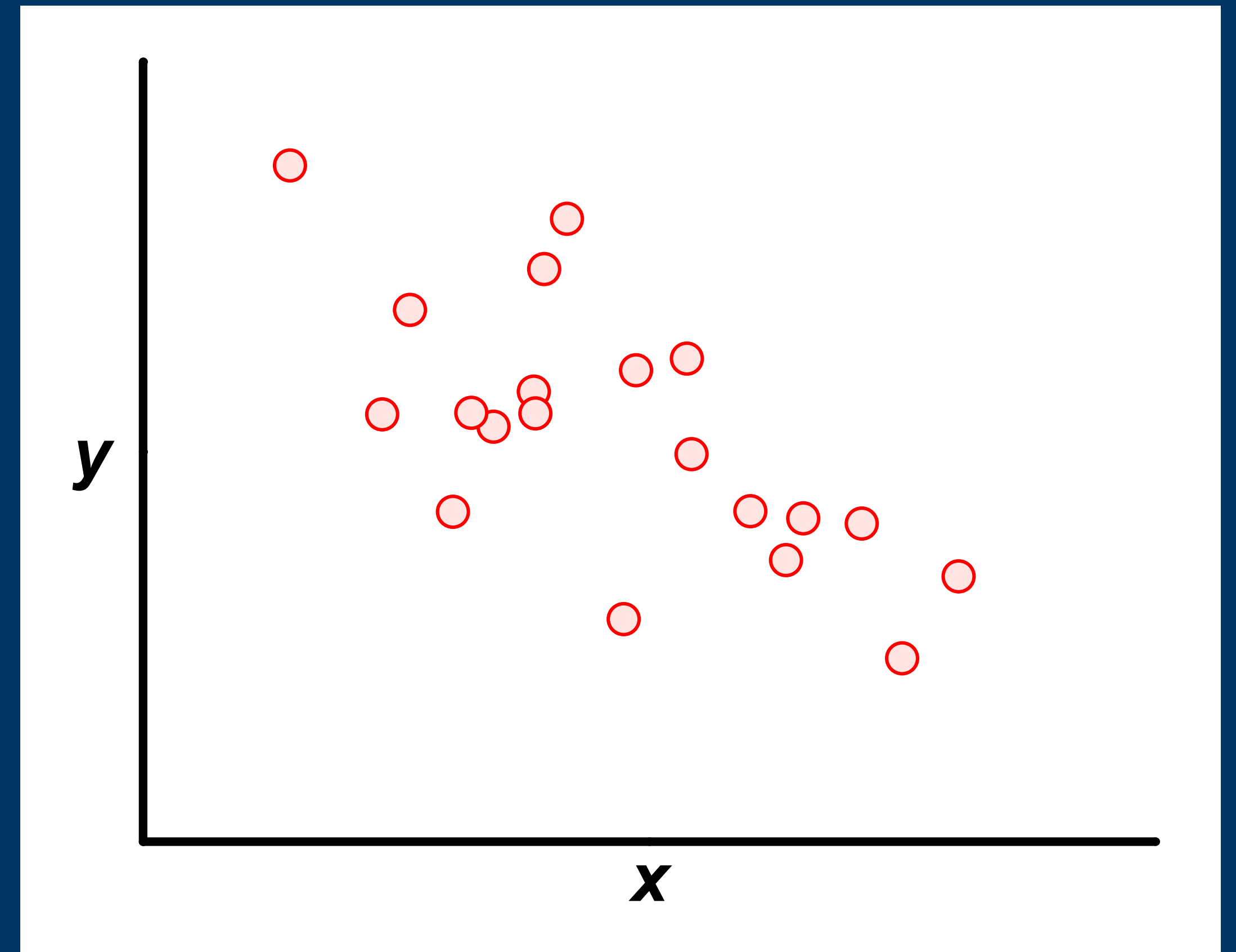
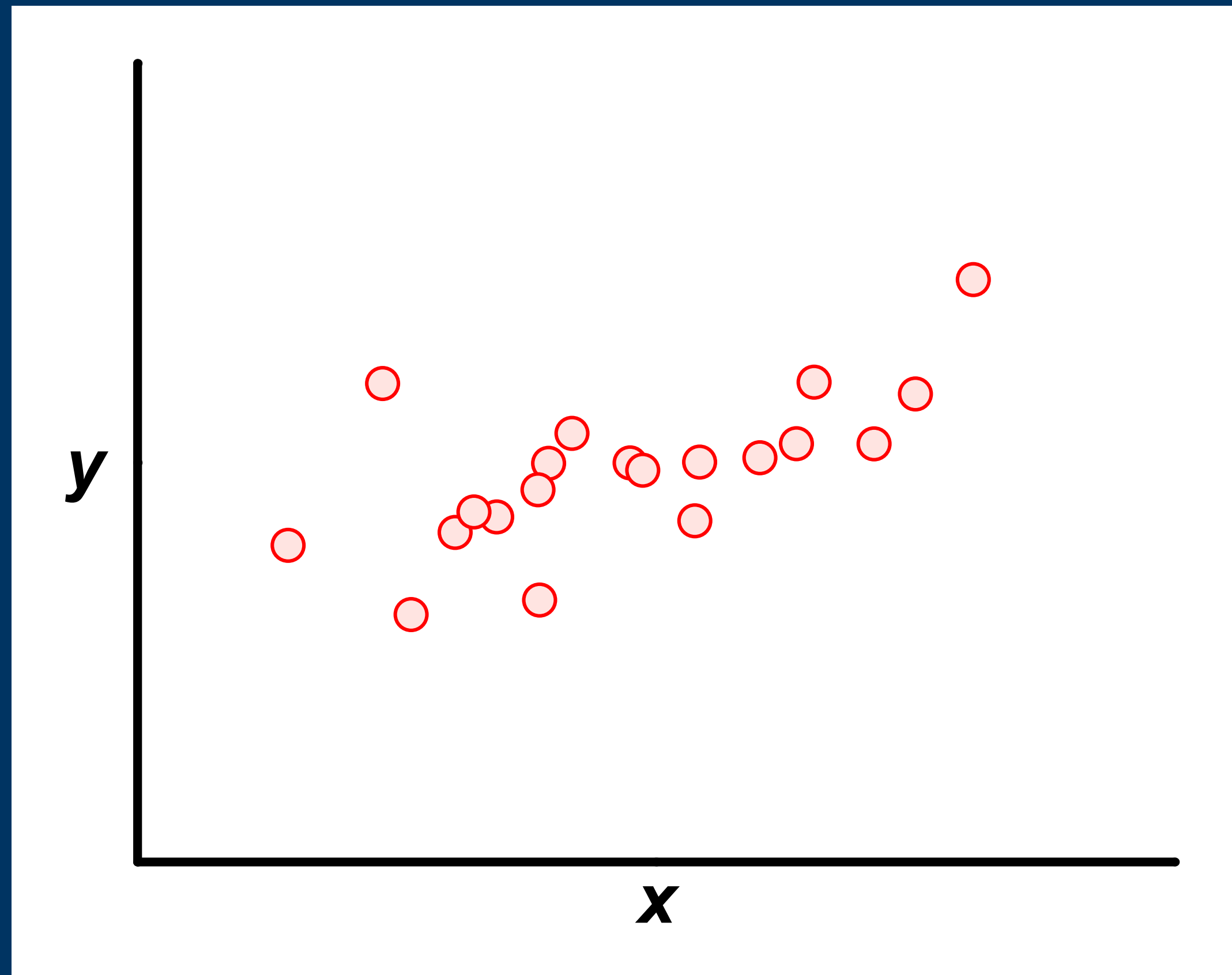
What does the slope represent?



In a perfect world



In a not-so-perfect world



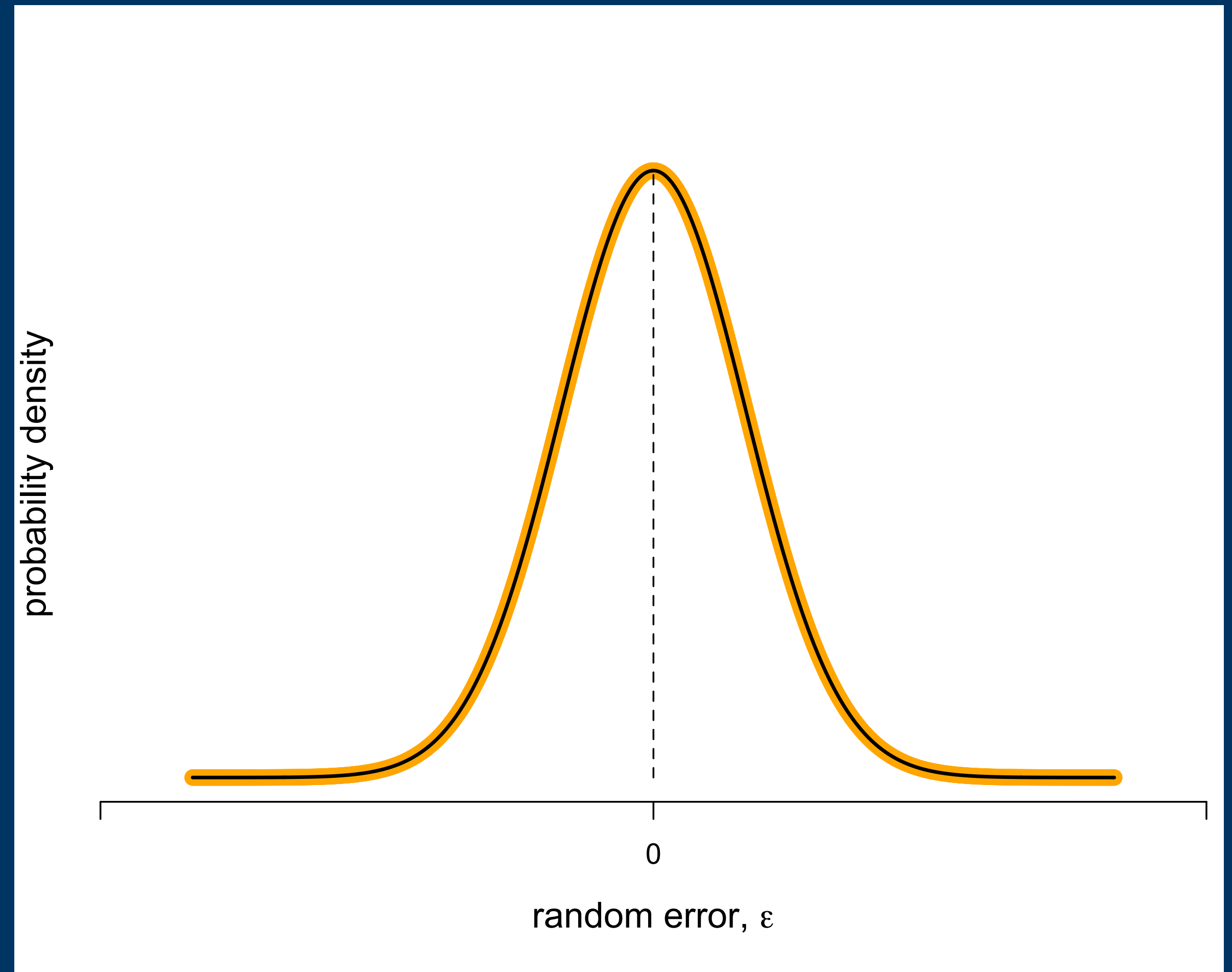
Introducing uncertainty/randomness

$$y_i = a + bx_i + \epsilon_i, \quad i = 1, \dots, n$$

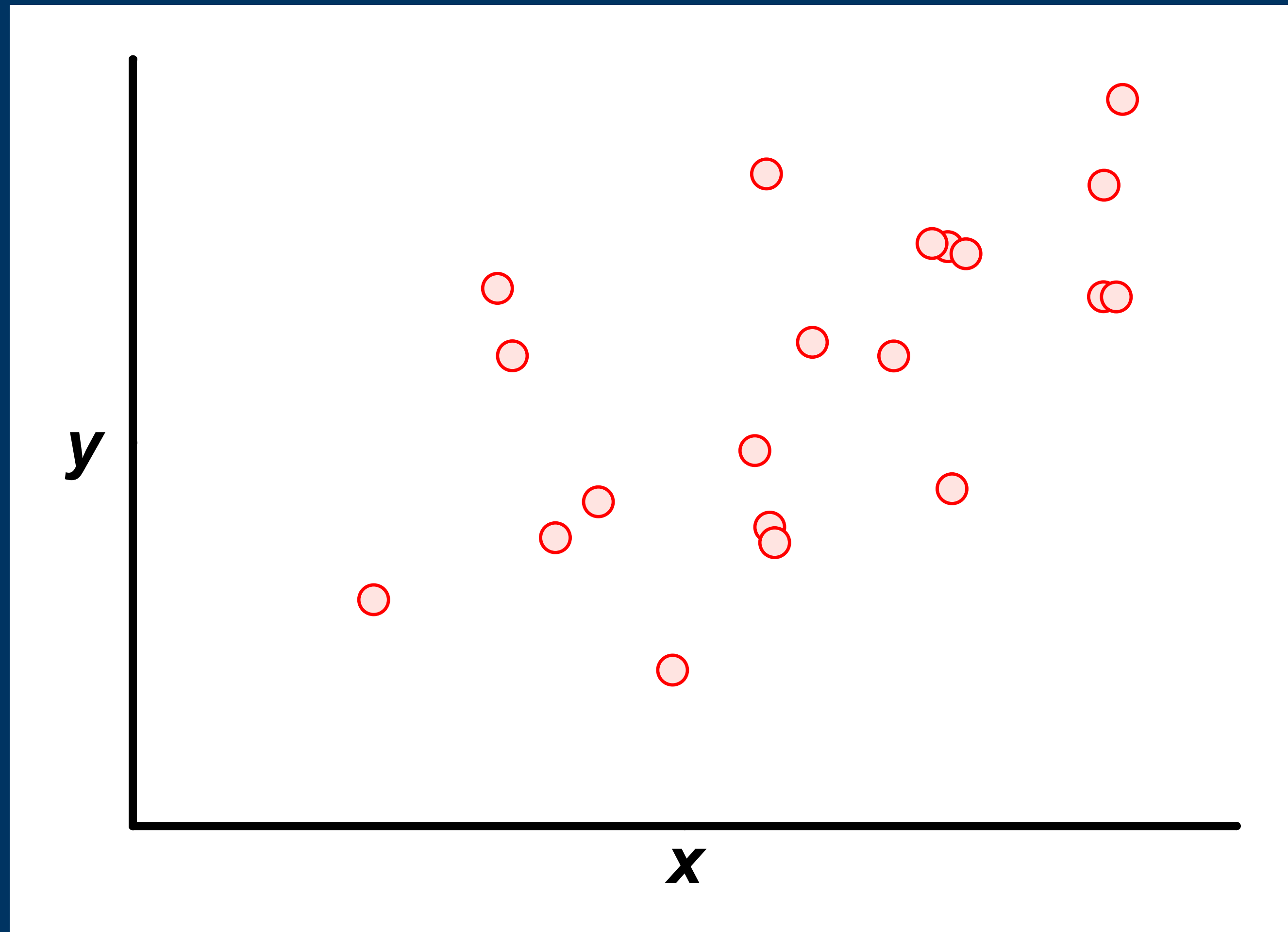
$$\epsilon_i \rightsquigarrow N(\mu = 0, \sigma)$$

Random error

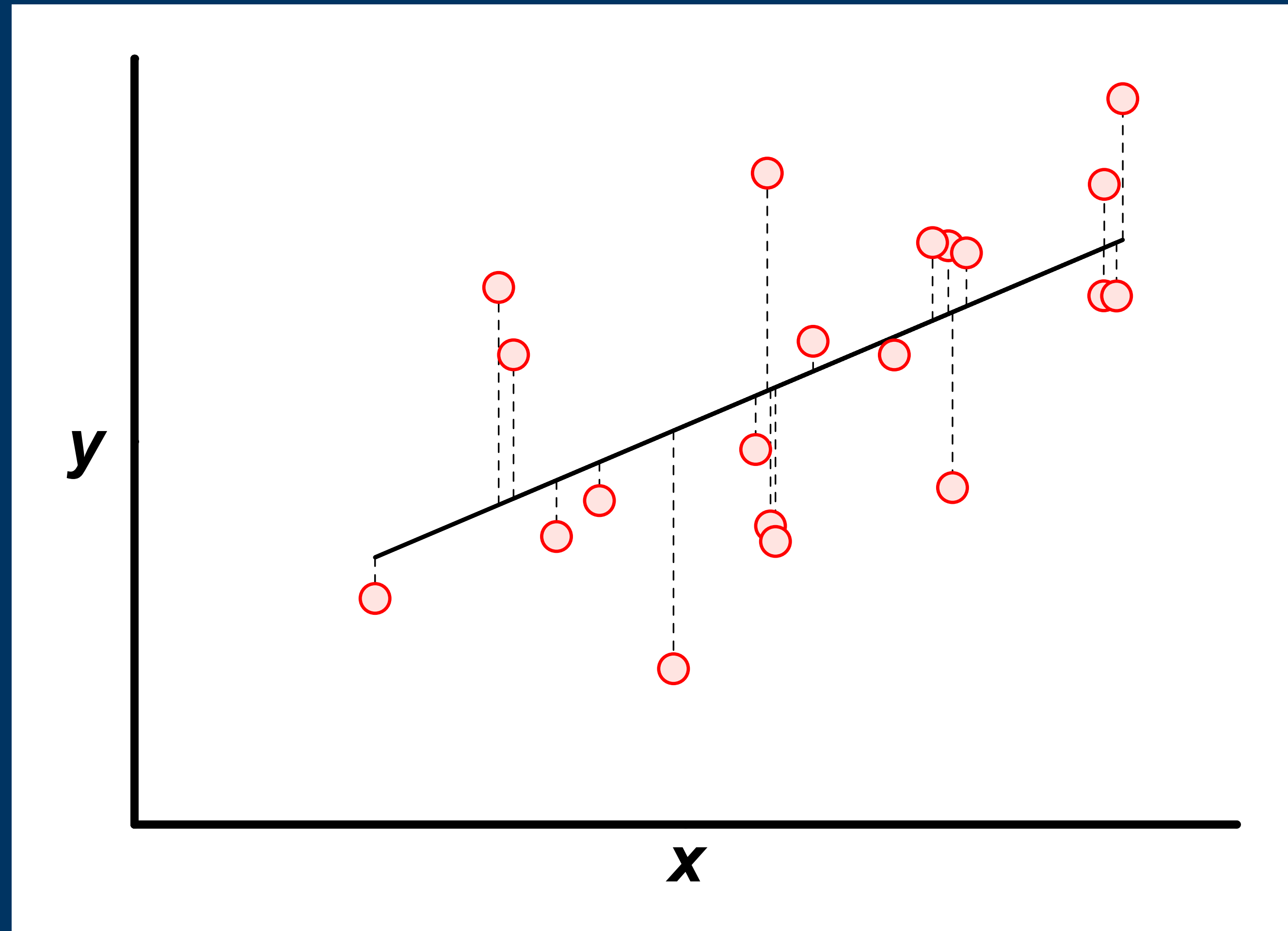
What is the source of this random error?



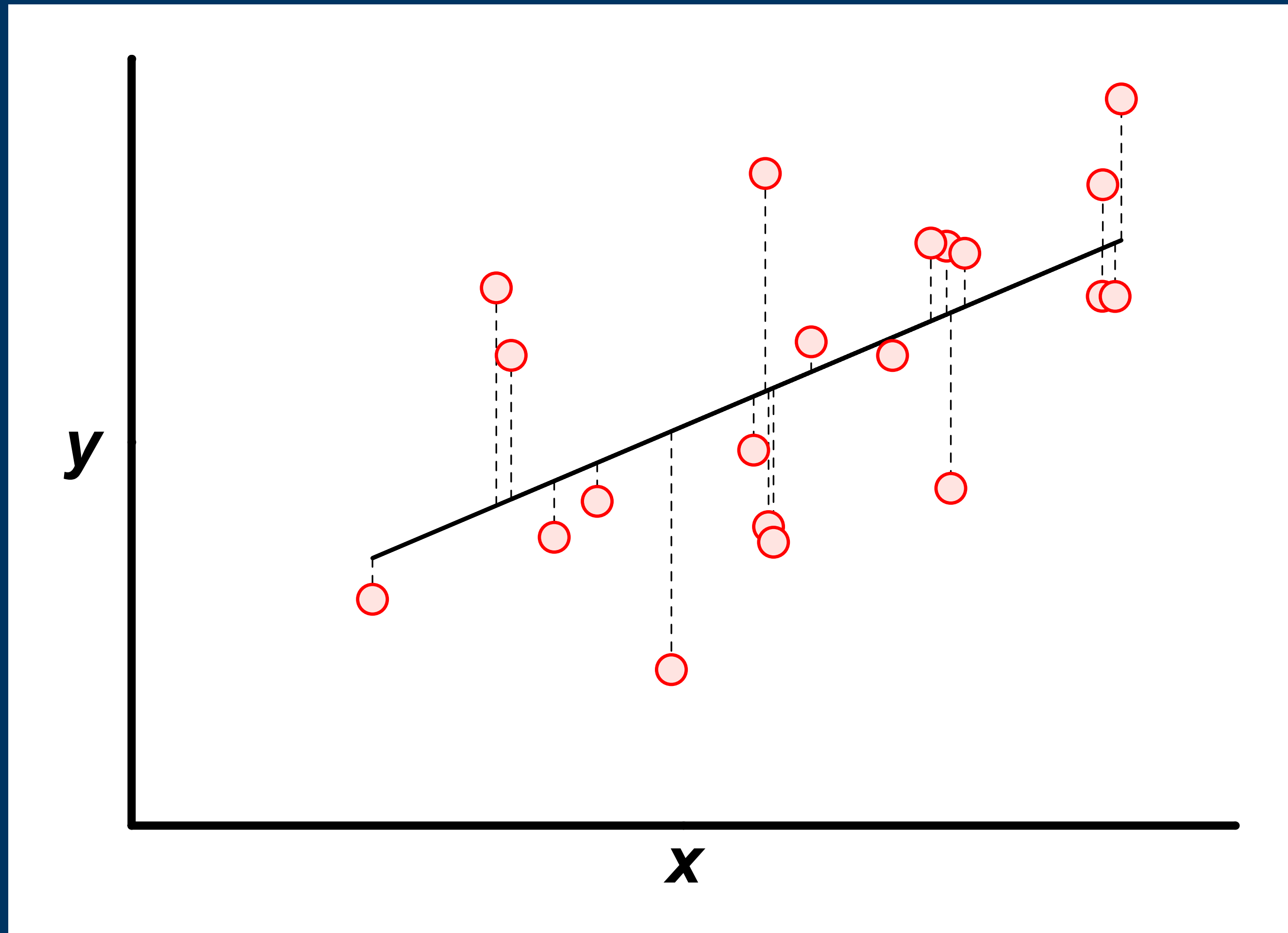
How to estimate a and b ?



Ordinary Least Squares Method



Ordinary Least Squares Method



$$\begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix} = \operatorname{argmin}_{a,b} \sum_{i=1}^n (y_i - a - bx_i)^2$$

Inference of interest (in large samples)

$$H_0 : a = 0 \text{ versus } H_1 : a \neq 0$$

$$H_0 : b = 0 \text{ versus } H_1 : b \neq 0$$

$$t = \frac{\hat{a}}{se(\hat{a})} \mid H_0 \rightsquigarrow N(\mu = 0, \sigma = 1)$$

$$t = \frac{\hat{b}}{se(\hat{b})} \mid H_0 \rightsquigarrow N(\mu = 0, \sigma = 1)$$

p-value < 0.05, reject H_0

p-value \geq 0.05, not reject H_0

0.05 is the **significance level of the test**

Warnings

Technical warning

Be aware of the model assumptions and their validity in the data

Interpretative warnings

Be aware of the dangers of extrapolating

Be aware of the dangers of inferring causality

It is R time!

