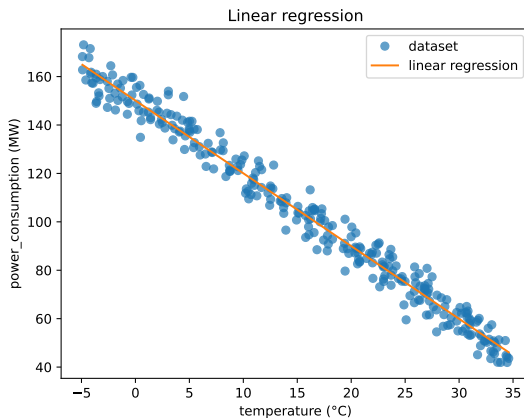


Solution to the linear regression in dimension 1



Linear regression

Formalization :

- ▶ input space (temperature) : $\mathcal{X} = \mathbb{R}$
- ▶ output space (power consumption) : $\mathcal{Y} = \mathbb{R}$
- ▶ dataset : $D_n = \{(x_1, y_1), \dots, (x_n, y_n), i \in [1, n]\}$.

When doing linear regression, our estimator is of the form :

$$h(x) = \theta x + b \tag{1}$$

with $\theta \in \mathbb{R}$, $b \in \mathbb{R}$.

Empirical risk minimization

$$R_n(\theta, b) = \sum_{i=1}^n (\theta x_i + b - y_i)^2 \quad (2)$$

We want to find θ and b such that $R_n(\theta, b)$ has the **smallest possible value**.

Derivatives

$$\begin{aligned}\frac{\partial R_n}{\partial \theta}(\theta, b) &= \sum_{i=1}^n 2(\theta x_i + b - y_i)x_i \\ &= 2\left[\theta \sum_{i=1}^n x_i^2 + b \sum_{i=1}^n x_i - \sum_{i=1}^n x_i y_i\right]\end{aligned}\tag{3}$$

$$\begin{aligned}\frac{\partial R_n}{\partial b}(\theta, b) &= \sum_{i=1}^n 2(\theta x_i + b - y_i) \\ &= 2\left[\theta \sum_{i=1}^n x_i + nb - \sum_{i=1}^n y_i\right]\end{aligned}\tag{4}$$

Hence we have a system of 2 equations with 2 unknowns (dropping the θ^* notation)

$$\theta \sum_{i=1}^n x_i^2 + b \sum_{i=1}^n x_i - \sum_{i=1}^n x_i y_i = 0 \quad (5)$$

$$\theta \sum_{i=1}^n x_i + nb - \sum_{i=1}^n y_i = 0 \quad (6)$$

Which means

$$b = \frac{1}{n} \left(\sum_{i=1}^n y_i - \theta \sum_{i=1}^n x_i \right) \quad (7)$$

$$\theta \sum_{i=1}^n x_i^2 + \frac{1}{n} \left(\sum_{i=1}^n y_i - \theta \sum_{i=1}^n x_i \right) \sum_{i=1}^n x_i - \sum_{i=1}^n x_i y_i = 0 \quad (8)$$

Finally :

$$\theta \left(\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right) + \frac{1}{n} \sum_{i=1}^n x_i \sum_{i=1}^n y_i - \sum_{i=1}^n x_i y_i = 0 \quad (9)$$

or

$$\theta^* = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{\sum_{i=1}^n x_i^2 - \frac{1}{n} \left[\sum_{i=1}^n x_i \right]^2} \quad (10)$$