

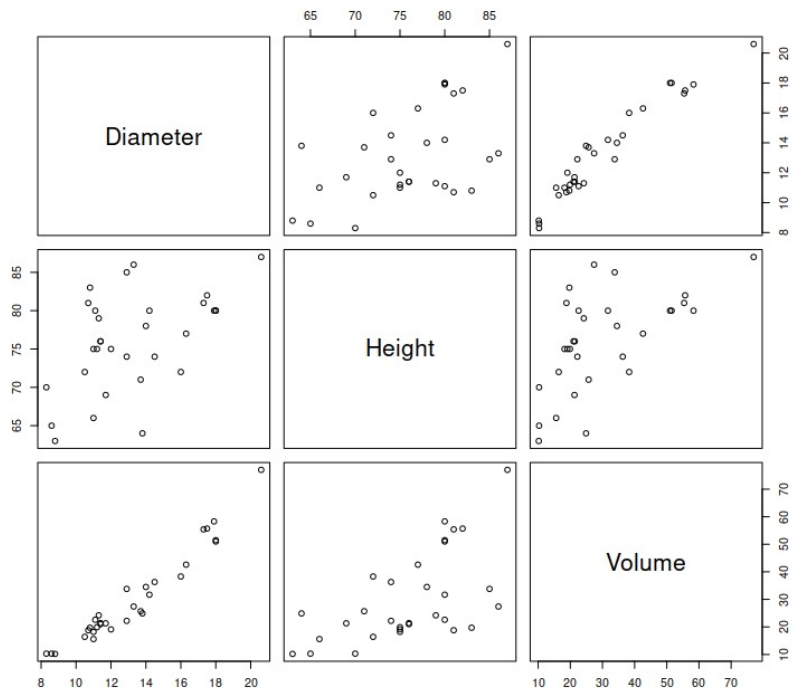
MULTIPLE LINEAR REGRESSION

There are now $p > 1$ covariates x_1, \dots, x_p .

Example: "trees" R dataset contains data on 31 cherry trees. In particular, we have

- diameter (inches)
- height (feet)
- volume

The goal is to predict the volume given the other 2 measures



Actually, if we think of the shape of a tree, we could think of approximating it to a cylinder

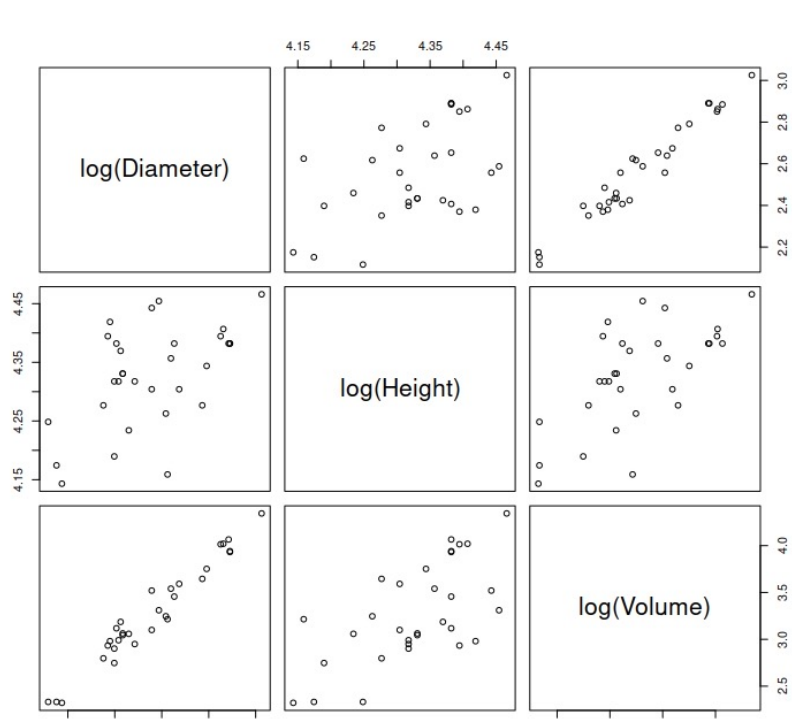


$$\text{volume} = \pi \cdot \text{radius}^2 \cdot \text{height} = \pi \cdot (d/2)^2 \cdot \text{height} \quad (\text{not linear!})$$

but

$$\log(\text{volume}) = \log \pi + 2 \log d - 2 \log 4 + \log \text{height}$$

We can consider the transformed variables



$$Y = \log(\text{volume})$$

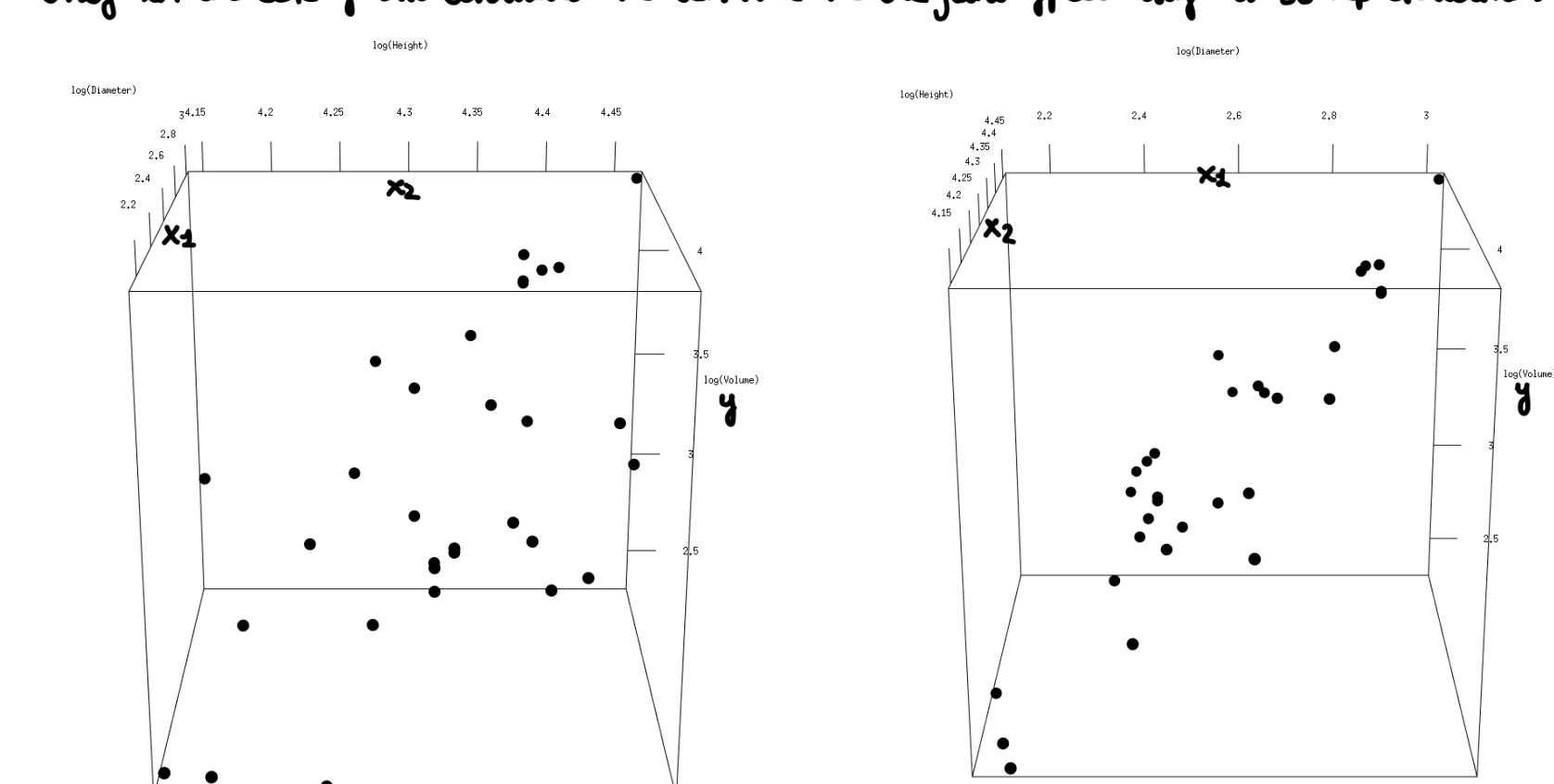
$$x_1 = \log(\text{diameter})$$

$$x_2 = \log(\text{height})$$

with 2 or more covariates we can no longer see the JOINT effect they have on y , but only the INDIVIDUAL effect of 1 predictor if we use a scatterplot.

The goal of the multiple lm is to study the JOINT EFFECT of the covariates on y .

Only in the case of two covariates we can still see the joint effect using a 3D representation



MODEL SPECIFICATION

We now observe $(y_i, x_{i1}, x_{i2}, \dots, x_{ip})$ for $i = 1, \dots, n$.

$$y_i = \mu_i + \epsilon_i$$

$$= \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i \quad i=1, \dots, n$$

± ϵ_i if we include the intercept

The assumptions don't change (they are just adjusted for the general case)

- normality, homoscedasticity, $\text{corr} = 0 \rightarrow \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \quad i=1, \dots, n$
- linearity: $\mu_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip}$
- absence of multicollinearity of the x_j : the covariates must be linearly independent (in the simple lm we had an analogous assumption: $\text{var}(x) \neq 0$)

notation:

$$\begin{cases} Y_1 = \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_p x_{1p} + \epsilon_1 \\ Y_2 = \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_p x_{2p} + \epsilon_2 \\ \vdots \\ Y_n = \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_p x_{np} + \epsilon_n \end{cases} \Rightarrow \underline{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \quad \underline{X}_j = \begin{bmatrix} x_{1j} \\ \vdots \\ x_{ij} \\ \vdots \\ x_{nj} \end{bmatrix} \quad \underline{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_i \\ \vdots \\ \epsilon_n \end{bmatrix}$$

for $j = 1, \dots, p$

$$\Rightarrow \underline{Y} = \beta_1 \underline{x}_1 + \dots + \beta_p \underline{x}_p + \underline{\epsilon}$$

$$\Rightarrow \underline{Y} = \sum_{j=1}^p \beta_j \underline{x}_j + \underline{\epsilon}$$

$$\Rightarrow \underline{Y} = \underline{X} \underline{\beta} + \underline{\epsilon}$$

$\begin{matrix} n \times 1 & n \times p & p \times 1 & n \times 1 \\ & \underline{\beta} & & \underline{\epsilon} \end{matrix}$

with

$$\underline{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1j} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots & & \vdots \\ x_{i1} & x_{i2} & \dots & x_{ij} & \dots & x_{ip} \\ \vdots & \vdots & & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nj} & \dots & x_{np} \end{bmatrix} \stackrel{n \times p}{\underline{X}} = [\underline{x}_1 \quad \underline{x}_2 \quad \dots \quad \underline{x}_j \quad \dots \quad \underline{x}_p] = \begin{bmatrix} \underline{x}_{11} \\ \vdots \\ \underline{x}_{i1} \\ \vdots \\ \underline{x}_{n1} \end{bmatrix}$$

$\rightarrow \underline{x}_j$ is the j -th covariate (n-dim vector) observed on the n units

$\rightarrow \underline{x}_i^T$ is the vector of the values of the p covariates on the i -th unit (p-dim vector)

and $\underline{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$

- \underline{Y} is a vector of r.v.
- \underline{X} is a matrix of constants (known)
- $\underline{\beta}$ is a vector of constants (unknown)
- $\underline{\epsilon}$ is a vector of r.v.

let's analyze the 3 hypotheses:

③ ABSENCE OF MULTICOLLINEARITY

What is the meaning of this hypothesis on $\underline{x}_1, \dots, \underline{x}_p$ (i.e. on the matrix \underline{X})?

Intuitively, it means that each covariate \underline{x}_j should have an individual contribution for predicting \underline{Y} \Rightarrow the information contained in \underline{x}_j can not be derived from the other variables.

- Examples of collinearity:
- the same variable is expressed using two measurement units (cm/m)
 - one variable is a linear combination of the others (e.g. $x_1 = \text{total years of education}$; $x_2 = \text{years of pre-university education}$; $x_3 = \text{years of post-university education}$; $\Rightarrow x_1 = x_2 + x_3$)

what happens when this hypothesis is not satisfied?

assume $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_p$ are linearly dependent: this means that there are p scalars a_1, \dots, a_p not all zero, such that $a_1 \underline{x}_1 + a_2 \underline{x}_2 + \dots + a_p \underline{x}_p = 0$

This means that I can write the j -th variable as $\underline{x}_j = -\frac{a_1}{a_j} \underline{x}_1 - \dots - \frac{a_{j-1}}{a_j} \underline{x}_{j-1} - \frac{a_{j+1}}{a_j} \underline{x}_{j+1} - \dots - \frac{a_p}{a_j} \underline{x}_p$

$$\Rightarrow \underline{Y} = \beta_1 \underline{x}_1 + \beta_2 \underline{x}_2 + \dots + \beta_{j-1} \underline{x}_{j-1} + \beta_j \underline{x}_j + \beta_{j+1} \underline{x}_{j+1} + \dots + \beta_p \underline{x}_p + \underline{\epsilon}$$

$$= \beta_1 \underline{x}_1 + \beta_2 \underline{x}_2 + \dots + \beta_{j-1} \underline{x}_{j-1} + \beta_j \left(-\frac{a_1}{a_j} \underline{x}_1 - \dots - \frac{a_{j-1}}{a_j} \underline{x}_{j-1} - \frac{a_{j+1}}{a_j} \underline{x}_{j+1} - \dots - \frac{a_p}{a_j} \underline{x}_p \right) + \dots + \beta_p \underline{x}_p + \underline{\epsilon}$$

$$= \underbrace{\left(\beta_1 - \beta_j \frac{a_1}{a_j} \right)}_{\beta_1^*} \underline{x}_1 + \dots + \underbrace{\left(\beta_{j-1} - \beta_j \frac{a_{j-1}}{a_j} \right)}_{\beta_{j-1}^*} \underline{x}_{j-1} + \underbrace{\left(\beta_{j+1} - \beta_j \frac{a_{j+1}}{a_j} \right)}_{\beta_{j+1}^*} \underline{x}_{j+1} + \dots + \underbrace{\left(\beta_p - \beta_j \frac{a_p}{a_j} \right)}_{\beta_p^*} \underline{x}_p + \underline{\epsilon}$$

We have expressed the same model using only $p-1$ variables.

Hence we need to require that the covariates are linearly independent $\Rightarrow \text{rank}(\underline{X}) = p$ (rank = # columns including the intercept! $\underline{x}_1 = 1$)

• linearity $\underline{\mu} = \sum_{j=1}^p \beta_j \underline{x}_j = \underline{X} \underline{\beta}$

• DISTRIBUTION: normality, homoscedasticity, incoherence

$$\underline{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_i \\ \vdots \\ \epsilon_n \end{bmatrix} \quad \mathbb{E}[\underline{\epsilon}] = \underline{0} \quad n\text{-dimensional vector of zeros}$$

$$\Rightarrow \mathbb{E}[\underline{Y}] = \mathbb{E}[\underline{X} \underline{\beta} + \underline{\epsilon}] = \underline{X} \underline{\beta}$$

$$\text{var}(\underline{\epsilon}) = \mathbb{E}[(\underline{\epsilon} - \mathbb{E}[\underline{\epsilon}])(\underline{\epsilon} - \mathbb{E}[\underline{\epsilon}])^T] \quad * \underline{\epsilon} \underline{\epsilon}^T = \begin{bmatrix} \sigma^2 & \epsilon_1 \epsilon_2 & \dots & \epsilon_1 \epsilon_n \\ \epsilon_1 \epsilon_2 & \sigma^2 & & \vdots \\ \vdots & & \ddots & \\ \epsilon_n \epsilon_1 & \dots & & \sigma^2 \end{bmatrix}$$

$$= \mathbb{E}[\underline{\epsilon} \underline{\epsilon}^T] * \quad = \sigma^2 \underline{I}_n$$

$$= \begin{bmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix} \quad \text{since } \mathbb{E}[\epsilon_i \epsilon_k] = 0$$

$$\mathbb{E}[\epsilon_i^2] = \sigma^2$$

$$\Rightarrow \text{var}(\underline{Y}) = \sigma^2 \underline{I}_n$$

Finally, the normality of $\underline{\epsilon}$ implies the normality of $\underline{Y} \Rightarrow \underline{Y} \sim N_n(\underline{X} \underline{\beta}, \sigma^2 \underline{I}_n)$

• INTERPRETATION of the coefficients β_1, \dots, β_p

we have seen that in the simple linear model

$$Y = \beta_1 + \beta_2 X + \epsilon$$

β_2 is the change in μ when we change x of one unit.

How do we interpret β_j , $j = 1, \dots, p$, in the case of multiple linear regression?

$$Y = \beta_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

β_j is now the change in μ when we change x_j of one unit, while keeping all other covariates fixed.

Let's consider the mean μ at two points x_j and $(x_j + 1)$

$$\mu^{(1)} = \beta_1 + \beta_2 x_2 + \dots + \beta_j x_j + \dots + \beta_p x_p$$

$$\mu^{(2)} = \beta_1 + \beta_2 x_2 + \dots + \beta_j (x_j + 1) + \dots + \beta_p x_p$$

$$\Rightarrow \mu^{(2)} - \mu^{(1)} = \beta_j$$