

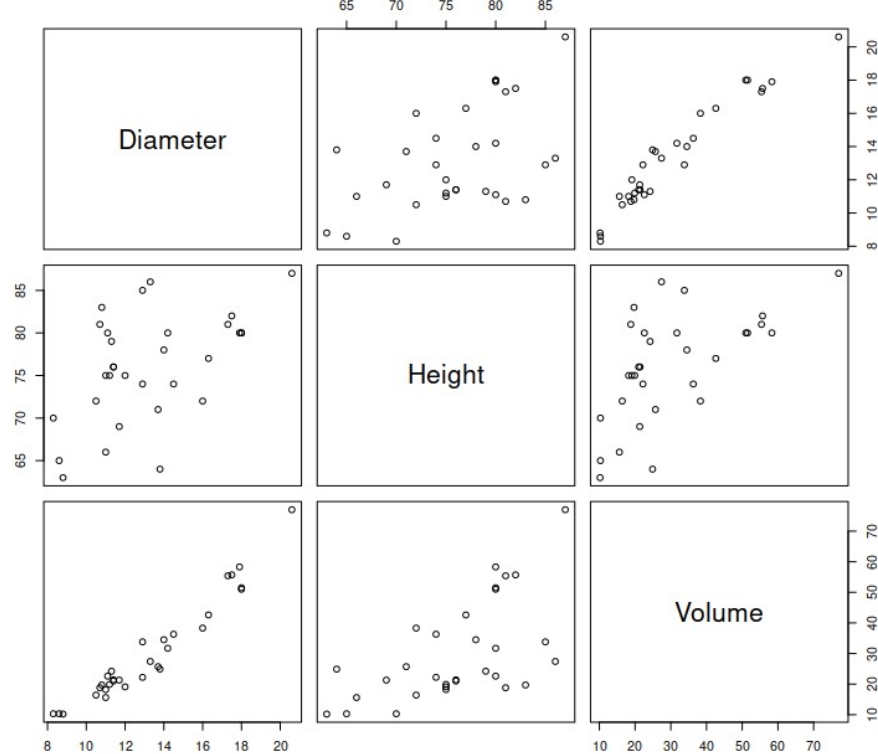
MULTIPLE LINEAR REGRESSION

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

With 3 or more variables we can no longer visualize the relationship with a scatterplot. We have to use a "matrix of scatterplots" which shows all the PAIRWISE combinations.



The goal is to predict the volume given the other 2 measures. If we think of the shape of a tree, we could think of approximating it to a cylinder.



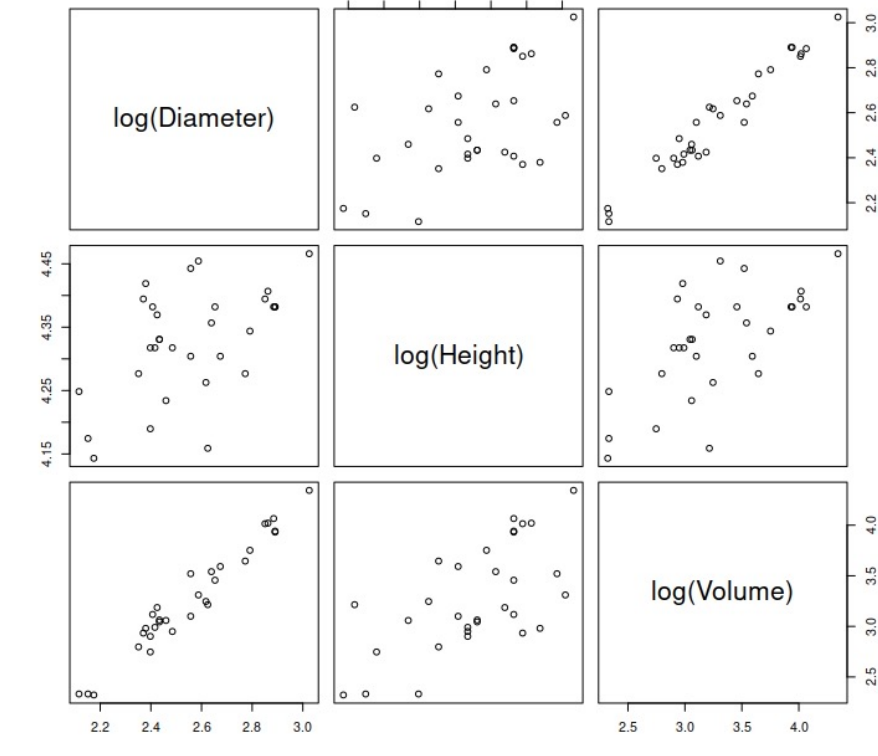
volume = $\pi \cdot \text{radius}^2 \cdot \text{height}$
 = $\pi \cdot (d/2)^2 \cdot \text{height}$

Hence we could specify a model where $\text{volume}_i \approx \pi \cdot \left(\frac{\text{diameter}_i}{2}\right)^2 \cdot \text{height}_i$ ← NOT LINEAR

However, the relationship can be linearized
 eg. $\log(\text{volume}_i) \approx \log \pi + \log 4 + 2 \log(\text{diameter}_i) + \log(\text{height}_i)$

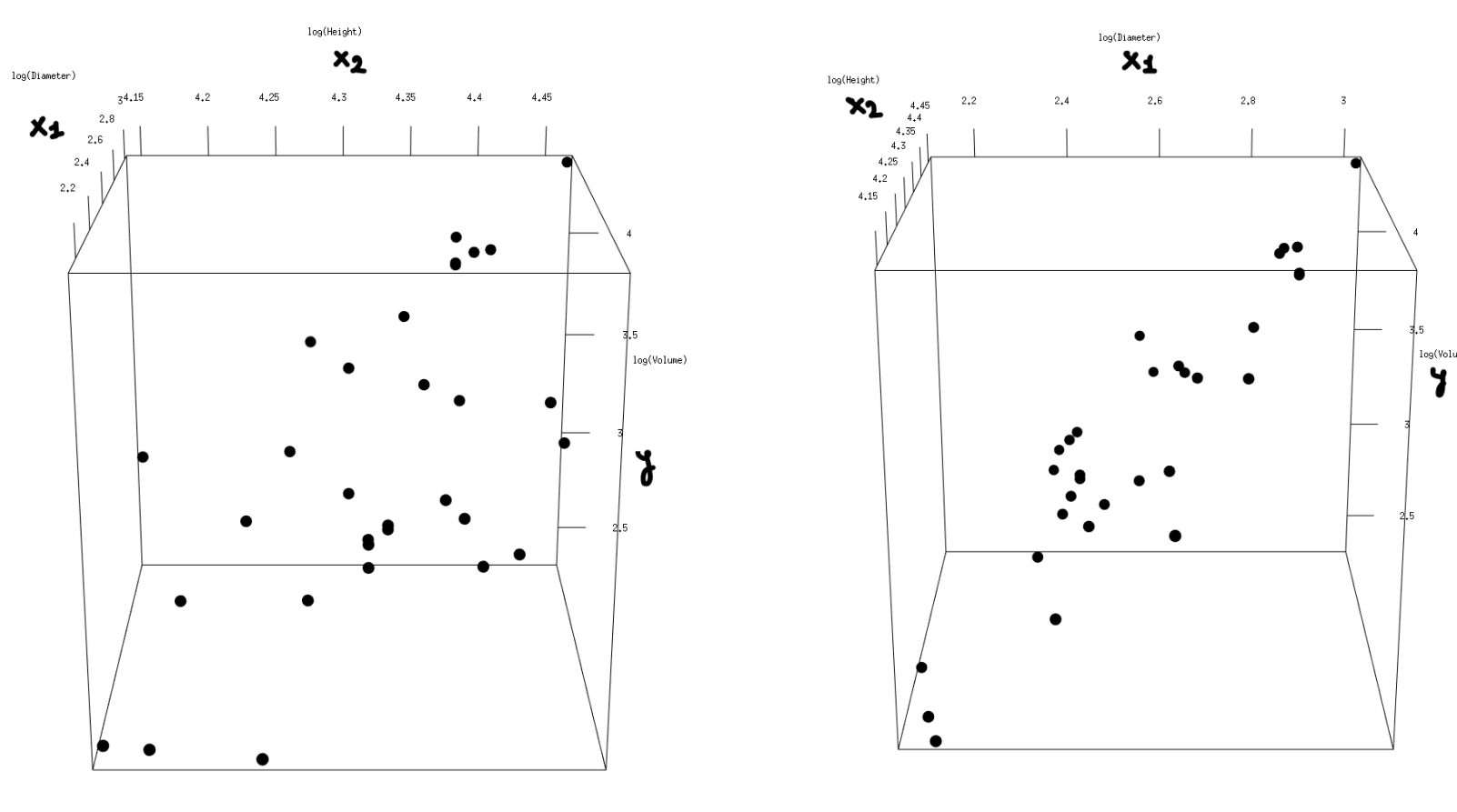
We can consider the transformed variables

$Y = \log(\text{volume})$
 $x_1 = \log(\text{diameter})$
 $x_2 = \log(\text{height})$



with 2 or more covariates we can not see the JOINT effect they have on y, but only the individual (marginal) effect of 1 covariate at a time

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



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

MODEL SPECIFICATION

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

$y_i = \mu_i + \epsilon_i$
 = $\beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$ $i = 1, \dots, n$
 i.e. if we include the intercept

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

- 1 - normality, homoscedasticity, corr = 0 → $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ $i = 1, \dots, n$
- 2 - linearity: $\mu_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip}$
- 3 - 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:

$$\begin{cases} y_1 = \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_p x_{1p} + \epsilon_1 \\ y_2 = \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_p x_{2p} + \epsilon_2 \\ \vdots \\ y_n = \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_p x_{np} + \epsilon_n \end{cases} \Rightarrow \underline{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \underline{X}_j = \begin{bmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{bmatrix} \quad \underline{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

for $j = 1, \dots, p$

$\Rightarrow \underline{Y} = \beta_1 \underline{X}_1 + \dots + \beta_p \underline{X}_p + \underline{\epsilon}$
 $\Rightarrow \underline{Y} = \sum_{j=1}^p \beta_j \underline{X}_j + \underline{\epsilon}$
 $\Rightarrow \underline{Y} = \underline{X} \underline{\beta} + \underline{\epsilon}$

with

$\underline{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1j} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2j} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nj} & \dots & x_{np} \end{bmatrix}^{n \times p} = \begin{bmatrix} \underline{X}_1 & \underline{X}_2 & \dots & \underline{X}_j & \dots & \underline{X}_p \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n x_{i1} \\ \sum_{i=1}^n x_{i2} \\ \vdots \\ \sum_{i=1}^n x_{ij} \\ \vdots \\ \sum_{i=1}^n x_{ip} \end{bmatrix}$

- \underline{X}_j is the j-th covariate observed on the n units (n-dim vector)
- \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. (n x 1)
- \underline{X} is a matrix of known constants (n x p)
- $\underline{\beta}$ is a vector of unknown constants (p x 1)
- $\underline{\epsilon}$ is a vector of r.v. (n x 1)

let's analyze the 3 hypotheses:

1) ABSENCE OF MULTICOLLINEARITY

What is the meaning of this hypothesis on $\underline{X}_1, \dots, \underline{X}_p$ (i.e. on the matrix X)?
 Intuitively, it means that each covariate \underline{X}_j should have an individual contribution for predicting Y → 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 (eg. $x_1 = \text{total years of education}$; $x_2 = \text{years of pre-university education}$; $x_3 = \text{years of post-university education}$; $\Rightarrow x_1 = x_2 + x_3$)

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 $\Rightarrow \text{rank}(X) = p$ (p is the number of columns of X, including the intercept $x_0 = 1$)

2) LINEARITY $\mu = \sum_{j=1}^p \beta_j x_j = \underline{X} \underline{\beta}$

3) DISTRIBUTION: normality, homoscedasticity, uncorrelation

$\underline{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_i \\ \vdots \\ \epsilon_n \end{bmatrix}$ vector of the errors

• expectation: $E[\underline{\epsilon}] = 0$ n-dimensional vector of zeros

• variance

$\text{var}(\underline{\epsilon}) = E[(\underline{\epsilon} - E[\underline{\epsilon}])(\underline{\epsilon} - E[\underline{\epsilon}])^T] = E[\underline{\epsilon} \underline{\epsilon}^T]$
 = $E \begin{bmatrix} E[\epsilon_1^2] & E[\epsilon_1 \epsilon_2] & \dots & E[\epsilon_1 \epsilon_n] \\ E[\epsilon_2 \epsilon_1] & E[\epsilon_2^2] & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ E[\epsilon_n \epsilon_1] & \dots & \dots & E[\epsilon_n^2] \end{bmatrix}$ since $E[\epsilon_i \epsilon_k] = 0$ for $i \neq k$
 $E[\epsilon_i^2] = \sigma^2$ for $i = 1, \dots, n$
 = $\begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix}$
 = $\sigma^2 I_n$ (n x n) matrix, diagonal elements = σ^2
 off-diagonal elements = 0

Hence $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ $i = 1, \dots, n \Rightarrow \underline{\epsilon} \sim N_n(0, \sigma^2 I_n)$

consequence for the response variable

$E[\underline{Y}] = E[\underline{X} \underline{\beta} + \underline{\epsilon}] = \underline{X} \underline{\beta}$
 $\text{var}(\underline{Y}) = \text{var}(\underline{X} \underline{\beta} + \underline{\epsilon}) = \text{var}(\underline{\epsilon}) = \sigma^2 I_n$

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

INTERPRETATION OF THE COEFFICIENTS β_1, \dots, β_p

We have seen that in the simple linear model $y_i = \beta_1 + \beta_2 x_i + \epsilon_i$, β_2 is the expected change in y_i (i.e., the change in $\mu = E[y_i]$) when we increase x_i by one unit. (or, equivalently, the expected difference in y when we consider two individuals i and k which differ in x by 1 unit: $\beta_2 = E[y_k] - E[y_i]$, when $x_i = x_0$ and $x_k = x_0 + 1$)

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

$y_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$

β_j now represents the expected change in y_i (i.e., the change in μ_i), when we increase x_{ij} by one unit, while keeping all other covariates fixed.

Let's consider the mean of y of two units i and k , $E[y_i] = \mu_i$ and $E[y_k] = \mu_k$. Assume that the values of the j-th covariate on these individuals are $x_{ij} = x_0$ and $x_{kj} = x_0 + 1$ while the other covariates are all equal: $x_{i2} = x_{k2}, x_{i3} = x_{k3}, \dots, x_{i,j-1} = x_{k,j-1}, x_{i,j+1} = x_{k,j+1}, \dots, x_{ip} = x_{kp}$

We get
 $\mu_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_j x_{ij} + \dots + \beta_p x_{ip}$ mean of individual i
 = $\beta_1 + \beta_2 x_{i2} + \dots + \beta_j x_0 + \dots + \beta_p x_{ip}$
 $\mu_k = \beta_1 + \beta_2 x_{k2} + \dots + \beta_j x_{kj} + \dots + \beta_p x_{kp}$
 = $\beta_1 + \beta_2 x_{k2} + \dots + \beta_j (x_0 + 1) + \dots + \beta_p x_{kp}$ mean of individual k
 = $\beta_1 + \beta_2 x_{k2} + \dots + \beta_j x_0 + \beta_j + \dots + \beta_p x_{kp}$

If we study the difference in their means

$\Rightarrow \mu_k - \mu_i = \beta_j$