

Statistical Principals of Data Analysis

Dr. Felipe Teles

Graz - Austria, 2023

Table of Contents

- 1 Linear Models
 - Fundamentals of Linear Models
 - Goodness of Fit
- 2 Analysis of Variance
 - Variance

Linear Models

A linear model in statistics is a model of the type

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_{p-1} X_{p-1} + \beta_p X_p + \varepsilon, \quad (1)$$

where Y is a response variable, $X_1, X_2, \dots, X_{p-1}, X_p$ are explanatory variables, $\beta_0, \beta_1, \beta_2, \dots, \beta_{p-1}, \beta_p$ are the coefficients of the linear model and $\varepsilon \sim N(0, \sigma^2)$ is the error of the model.

Linear Models

Or, in matrix notation:

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{21} & X_{31} & \cdots & X_{(p-1)1} & X_{p1} \\ 1 & X_{12} & X_{22} & X_{32} & \cdots & X_{(p-1)2} & X_{p2} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & X_{1n} & X_{2n} & X_{3n} & \cdots & X_{(p-1)n} & X_{pn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}, \quad (2)$$

which can be written as

$$Y = X\beta + \varepsilon, \quad (3)$$

where X is the design matrix (explanatory variables, treatments, groups, etc.).

Linear Models

Examples of linear models:

- $Y = \beta_0 + \beta_1 X_1;$
- $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_1^5;$
- $Y = \frac{\beta_1}{X_1};$
- $Y = \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_1 X_2 X_3;$

Linear Models

Examples of nonlinear models:

- $Y = \beta_0 + \beta_1 X_1 + \frac{\beta_2}{\beta_3 + X_1};$
- $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2^{\beta_3 X_3};$
- $Y = \cos(\beta_1 X_1) + \sin(\beta_2 X_1);$
- $Y = \beta_0 e^{\beta_1 X_1};$
- $\frac{dS}{dt} = -\frac{\mu_{max} S}{K_S + S} X Y_{X/S},$ where $\mu_{max}, K_S, Y_{X/S}$ are parameters;

Least Squares Method

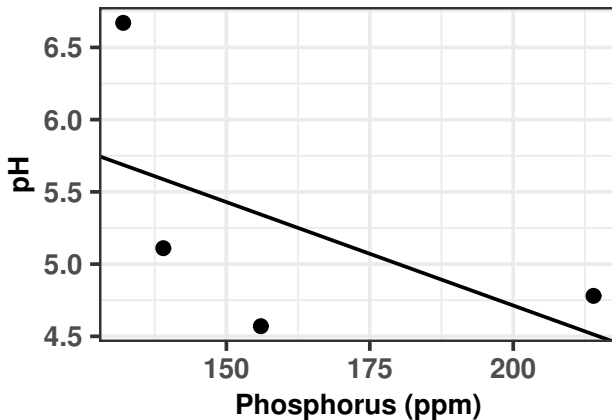


Figure: Linear Model: pH explained by nitrogen from Soil data set

Least Squares Method

Goal: minimize the residual sum of squares (RSS)

$$RSS = \|Y - X\beta\|^2 \quad (4)$$

$$= (Y - X\beta)^T(Y - X\beta) \quad (5)$$

$$= Y^T Y - Y^T X\beta - \beta^T X^T Y + \beta^T X^T X\beta \quad (6)$$

$$= Y^T Y - 2\beta^T X^T Y + \beta^T X^T X\beta. \quad (7)$$

Least Squares Method

Goal: minimize the residual sum of squares (RSS)

$$RSS = Y^T Y - 2\beta^T X^T Y + \beta^T X^T X \beta$$

$$\min_{\beta} RSS \implies \frac{\partial}{\partial \beta} RSS = 0 \quad (8)$$

$$\frac{\partial}{\partial \beta} RSS = 0 \quad (9)$$

$$-2X^T Y + 2X^T X \beta = 0 \quad (10)$$

$$X^T X \beta = X^T Y \quad (11)$$

$$(X^T X)^{-1} X^T X \beta = (X^T X)^{-1} X^T Y \quad (12)$$

$$\beta = (X^T X)^{-1} X^T Y. \quad (13)$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y. \quad (14)$$

Exercise

- 1 In *R*, install and load the packages *car*, *carData* and *ggplot2*.
- 2 Load the data set *Soils* and, by using the function *subset*, filter the data for $Gp=S1$.
- 3 Fit a linear model for pH using the explanatory variable *P*.
- 4 Check the results with *summary*.
- 5 Theoretical check - calculate $\hat{\beta}$ using Equation (14).

Exercise

```
data01=subset( Soils , Gp==" S1" )  
model1=lm( data=data01 , pH~P )  
summary(model1)
```

```
X=as.matrix( model.matrix(model1) )  
Y=as.matrix( data01$pH )  
Beta=solve( t(X)%*%X )%*%t(X)%*%Y
```

```
ggplot( data=data01 , aes( x=P , y=pH ) ) +  
geom_point( pch=19 ) + theme_bw() +  
geom_abline( intercept=coef(model1)[1] ,  
slope=coef(model1)[2] )
```

Goodness of Fit

Based on what can we say if one model is good?

- ① Coefficient of Determination R^2 ;
- ② Distribution of Errors;
- ③ Unusual observations;
- ④ Structure of the model.

Goodness of Fit - Coefficient of Determination

The coefficient of determination is the percentage of data variability explained by the model:

$$R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (15)$$

Some

- 1 $0 \leq R^2 \leq 1$;
- 2 If $R^2 = 1$, then $RSS = 0$;
- 3 R^2 is equivalent to the percentage of variability explained regarding total variability.

Goodness of Fit - Distribution of Errors

A linear model:

$$Y = X\beta + \varepsilon, \varepsilon \sim N(0, \sigma^2). \quad (16)$$

Check:

- 1 Distribution of errors around zero: $E[\varepsilon] = 0$;
- 2 Homoscedasticity: $\text{Var}[\varepsilon] = \sigma^2$;
- 3 Normality of residuals: Q-Q plot;

Goodness of Fit - Distribution of Errors

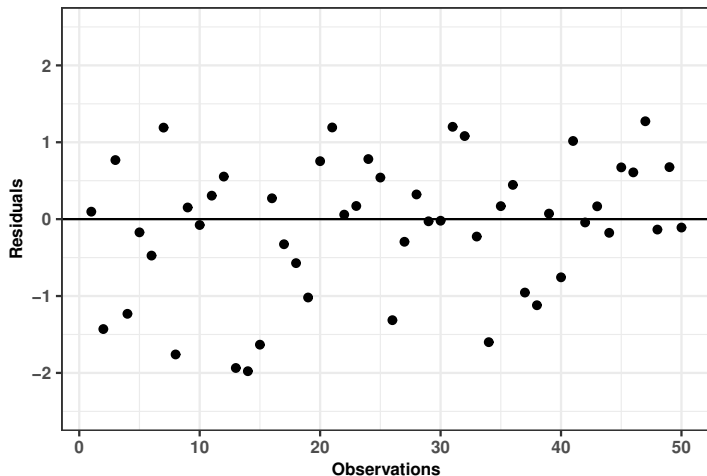


Figure: Residuals normally distributed with homoscedastic variance structure.

Goodness of Fit - Distribution of Errors

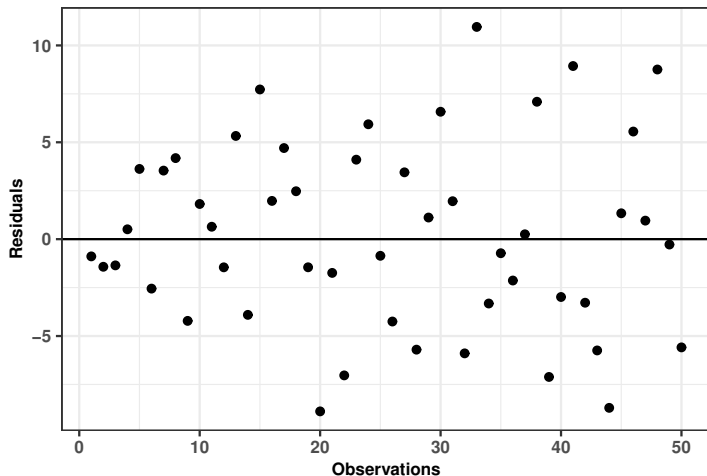


Figure: Residuals normally distributed with mild disturbance.

Goodness of Fit - Distribution of Errors

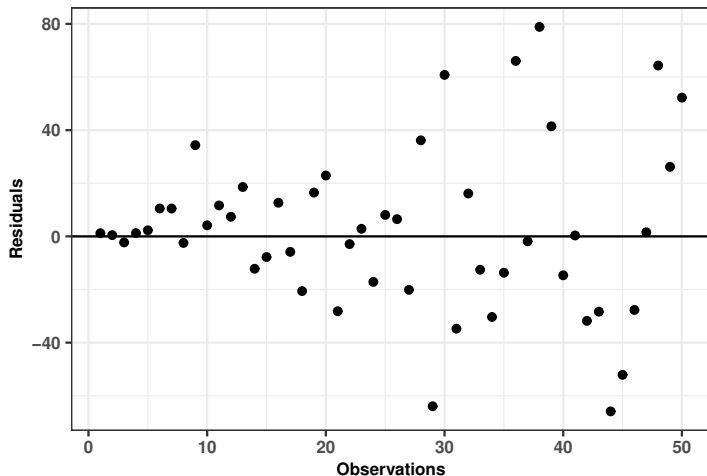


Figure: Residuals normally distributed with heteroscedastic variance structure.

Goodness of Fit - Distribution of Errors

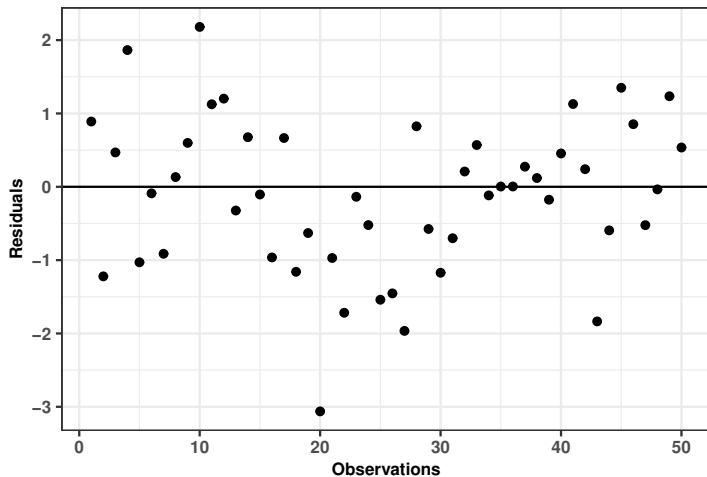


Figure: Nonlinear residuals.

Goodness of Fit - Distribution of Errors

Normal Q-Q Plot

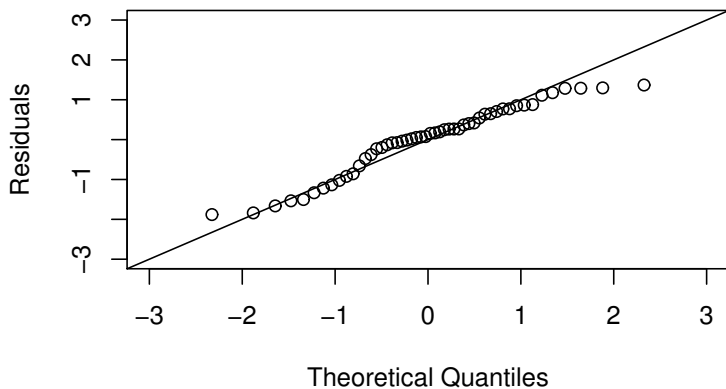


Figure: Q-Q plot.

Goodness of Fit - Distribution of Errors

Test of normality - Shapiro-Wilk test (shapiro.test)

$$H_0 : \varepsilon \sim N(\mu, \sigma^2) \quad (17)$$

$$H_1 : \varepsilon \not\sim N(\mu, \sigma^2) \quad (18)$$

Exercise

Using the data set *Soils*:

- 1 Fit the following model:

$$pH = \beta_0 + \beta_1 N + \beta_2 P + \beta_3 Ca + \beta_4 Mg + \beta_5 K + \beta_6 Na + \varepsilon.$$

- 2 Make the analysis of diagnostics - plot the residuals and identify possible leverage points and outliers.
- 3 Based on ANOVA, is there any variable that is not significant to explain pH?
- 4 What happens when we fit the following model:
 $pH = \beta_1 N + \beta_2 P + \beta_3 Ca + \beta_4 Mg + \beta_5 K + \beta_6 Na + \varepsilon$? Repeat analyses 1, 2 and 3.
- 5 Now include Block effects. What can we see?
- 6 Is there any significant block effect?
- 7 Is there any variable deviating from normality? Is there any transformation possible?
- 8 Is there some *normalization* that we could apply?

Variance

- ① We are mostly interested in analysing the variance of a data set, but why?
- ② Can we calculate the variance of a variable always with the same formula?
- ③ What is the difference between sample and population variances?
- ④ What is between group variance? And within group variance?
- ⑤ What about covariance?

Expectation

Before we start with variance, let us remember what is expectation:

$$E[X] = \sum_{X \in \Omega} X \mathbb{P}[X = X] \quad (19)$$

or

$$E[X] = \int_{\Omega} X \mathbb{P}[X = X] dx, \quad (20)$$

where X is a variable defined in $(\mathcal{F}, \Omega, \mathbb{P})$.

Expectation

For example, when $X \sim U(a, b)$, whose probability density function is $\mathbb{P}[x = X] = \frac{1}{b-a}$:

$$E[X] = \int_a^b x \mathbb{P}[x = X] dx \quad (21)$$

$$= \int_a^b x \left(\frac{1}{b-a} \right) dx \quad (22)$$

$$= \left(\frac{1}{b-a} \right) \left[\frac{x^2}{2} \right]_a^b \quad (23)$$

$$= \left(\frac{1}{b-a} \right) \frac{(b^2 - a^2)}{2} \quad (24)$$

$$= \left(\frac{1}{b-a} \right) \frac{(b-a)(b+a)}{2} \quad (25)$$

$$= \frac{b+a}{2}. \quad (26)$$

Exercise

- ① What is the expectation of X when $X \sim U(3, 7)$? Use the function *runif*(n, a, b).
- ② Plot a histogram for 10, 100, 1000, 10000 and 100000 values of X . Use the function *hist*(\cdot).
- ③ What is the mean of 4, 50 and 100 values generated from $U(3, 7)$? Is the mean the expected value?
- ④ Plot the function f , such that: $f(n) = E[X | \#\Omega = n]$. What can you observe?
- ⑤ Maurício: Exponential, Ailton: Binomial, Thalita: Gamma, Gabriel: Weibull, Fernanda: Normal.

Some distributions

Table: Expectation and Variance for some distributions.

Definition	Expected value	Variance
$N(\mu, \sigma^2)$	μ	σ^2
$U(a, b)$	$\frac{b + a}{2}$	$\frac{(b - a)^2}{12}$
$Exp(\lambda)$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$
$Poisson(\lambda)$	λ	λ
$Gamma(\alpha, \beta)$	$\frac{\alpha}{\beta}$	$\frac{\alpha}{\beta^2}$
$Weibull(\alpha, \beta)$	$\alpha \Gamma(1 + 1/\beta)$	$\alpha^2 \left(\Gamma(1 + 2/\beta) - (\Gamma(1 + 1/\beta))^2 \right)$
$Binomial(n, p)$	np	$np(1 - p)$

Exercise

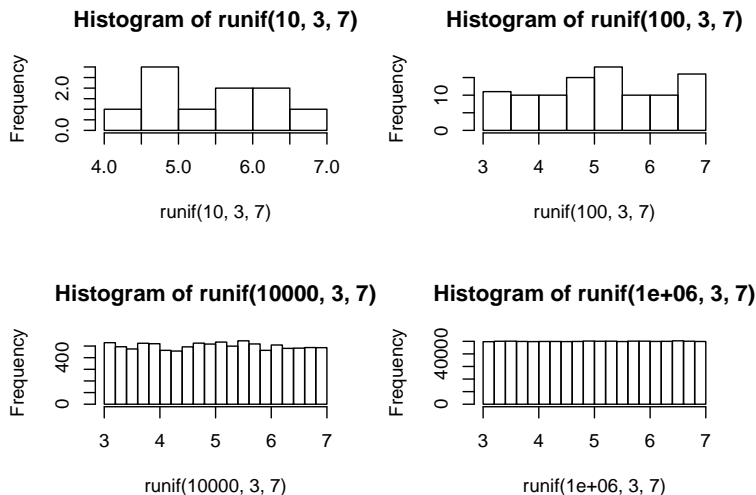


Figure: Histogram for $U(3, 7)$ for different number of generated points.

Exercise

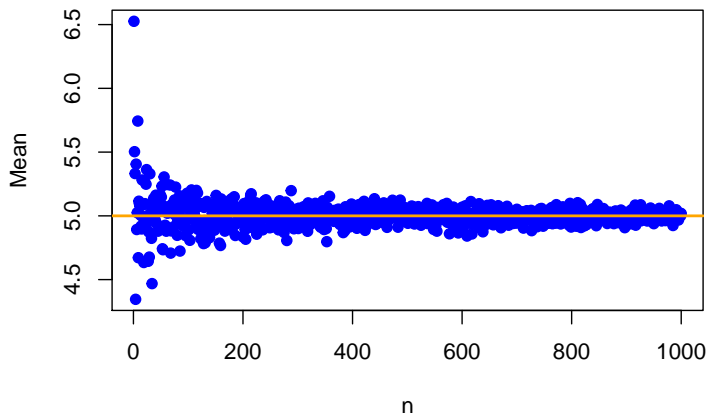
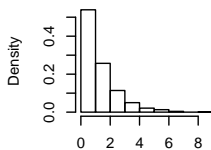
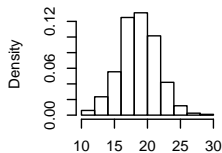


Figure: Mean of $U(3,7)$ for different values of n . In orange, the expected value.

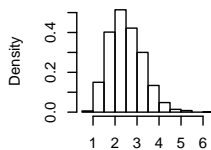
Exercise

Histogram of Exponential

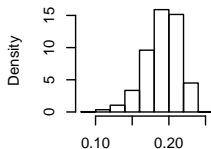
Exponential

Histogram of Binomial

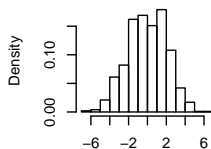
Binomial

Histogram of Gamma

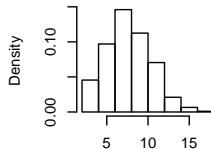
Gamma

Histogram of Weibull

Weibull

Histogram of Normal

Normal

Histogram of Poisson

Poisson

Figure: Histogram for different distributions of probability.

Exercise

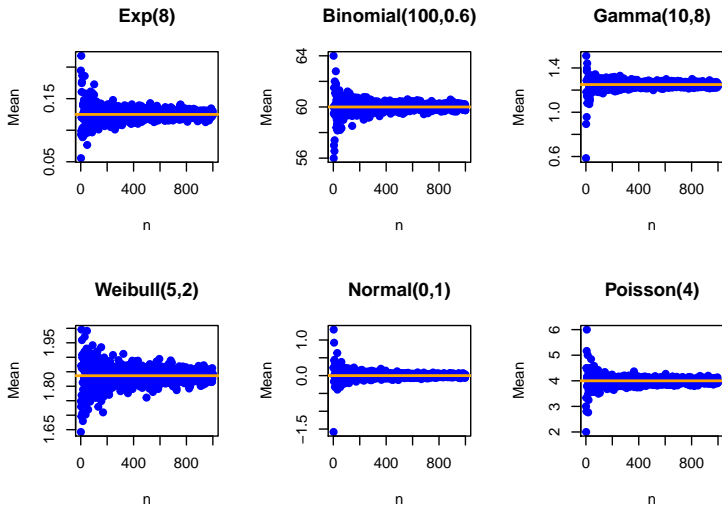


Figure: Mean of different distributions for different values of n . In orange, the expected value.