

# Statistical Methods for Economists

## Lecture (7 & 8)a

One-Way Analysis of Variance  
(ANOVA)



**SILESIAN  
UNIVERSITY**

SCHOOL OF BUSINESS  
ADMINISTRATION IN KARVINA

**David Bartl**

Statistical Methods for Economists  
INM/BASTE

# Outline of the lecture

---



- One-Way ANOVA: Introduction and Motivation
  - One-Way ANOVA: Summary, Assumptions, and the Goal of the analysis
  - One-Way ANOVA as a model of Multiple Linear Regression
  - One-Way ANOVA: the  $F$ -test
-



The methods of the **Analysis of Variance (ANOVA) are special cases of the methods of the **Multiple Linear Regression**.**

- In essence, we measure one numerical statistical variable, denoted by “ $y$ ”, that is a quantitative data item of some statistical data units.
  - In addition, we consider one or more (qualitative / quantitative) factors.
  - Depending upon the number of the factors, we distinguish:
    - one-factor = one-way ANOVA
    - two-factor = two-way ANOVA
-

# ANOVA: Introduction

---



- We assume that each of the factors can attain only finitely many distinct values.
  - Hence, there are only finitely many possible combinations (a Cartesian product) of all possible values of the factors.
  - Each combination of the values constitutes one group of the statistical units.
  - The purpose is to study whether the expected value  $E[y]$  of the statistical variable “ $y$ ” depends upon the values of the parameters, or not; that is, whether the expected value in each group is the same, or not.
  - ANOVA is based upon the theory of Multiple Linear Regression.
  - We shall study one-factor ANOVA in this lecture.
-

# One-factor ANOVA: Motivation: Example 1

---



Consider a gross sample of  $n$  patients who have been cured for some disease.

The patients were divided into  $k$  groups:

- There are  $n_1$  patients in the 1<sup>st</sup> group.
- There are  $n_2$  patients in the 2<sup>nd</sup> group.
- ...
- There are  $n_k$  patients in the  $k^{\text{th}}$  group.

---

→ There are  $n_1 + n_2 + \dots + n_k = n$  patients in total.

---

# One-factor ANOVA: Motivation: Example 1

---



Each of the  $n_1$  patients in the 1<sup>st</sup> group has been cured by method no. 1.

Each of the  $n_2$  patients in the 2<sup>nd</sup> group has been cured by method no. 2.

...

Each of the  $n_k$  patients in the  $k^{\text{th}}$  group has been cured by method no.  $k$ .

We then measure the success of the medical treatment.

That is, let “ $y$ ” be a (quantitative) variable meaning the health of the patient.

---

# One-factor ANOVA: Motivation: Example 1

---



Measuring the health of each of the  $n_1$  patients in the 1<sup>st</sup> group,  
we have the values  $y_{11}, y_{12}, y_{13}, \dots, y_{1n_1}$

Measuring the health of each of the  $n_2$  patients in the 2<sup>nd</sup> group,  
we have the values  $y_{21}, y_{22}, \dots, y_{2n_2}$

...

...

Measuring the health of each of the  $n_k$  patients in the  $k^{\text{th}}$  group,  
we have the values  $y_{k1}, y_{k2}, y_{k3}, y_{k4}, \dots, y_{kn_k}$

---

# One-factor ANOVA: Motivation: Example 1

---



Having got the results of the measurements of the health

group no. 1:  $y_{11}, y_{12}, y_{13}, \dots, y_{1n_1}$   
group no. 2:  $y_{21}, y_{22}, \dots, y_{2n_2}$   
group no. 3:  $y_{31}, y_{32}, y_{33}, y_{34}, y_{35}, \dots, y_{3n_3}$   
...  
group no.  $k$ :  $y_{k1}, y_{k2}, \dots, y_{kn_k}$

**we test the null hypothesis that:**

- there are no differences among the treatments, that is
  - the expected values of the health in the groups are approximately the same
-



## One-factor ANOVA: Motivation: Example 2

---



We test  $k$  distinct cars. We test the 1<sup>st</sup> car  $n_1$  times, we test the 2<sup>nd</sup> car  $n_2$  times, etc., and we test the  $k^{\text{th}}$  car  $n_k$  times for mileage (fuel consumption per 100 km).

Then  $y_{11}, y_{12}, \dots, y_{1n_1}, y_{21}, y_{22}, \dots, y_{2n_2}, \dots, y_{k1}, y_{k2}, \dots, y_{kn_k}$  are the results of the measurements, i.e. the fuel consumptions per 100 km.

We test the null hypothesis that:

the average fuel consumption of each car is the same.

---

# One-Way ANOVA



- Summary
- Assumptions
- The classical assumptions
- The purpose of the Analysis of Variance

# One-way ANOVA: Summary

---



We have got the sample of the  $k$  groups  
of the  $n = n_1 + n_2 + \dots + n_k$  observations

$y_{11}, y_{12}, y_{13}, \dots, y_{1n_1}$

$y_{21}, y_{22}, \dots, y_{2n_2}$

$y_{31}, y_{32}, y_{33}, y_{34}, y_{35}, \dots, y_{3n_3}$

...

$y_{k1}, y_{k2}, \dots, y_{kn_k}$

where  $y_{ij} \in \mathbb{R}$  for  $j = 1, 2, \dots, n_i$  for  $i = 1, 2, \dots, k$ .

The sample could have been obtained in either of the following two ways:

---

# One-way ANOVA: Summary

---



## First:

- A sample of  $n$  statistical units was selected from a larger population.
- Each of the statistical units was placed into its respective group  $i \in \{1, 2, \dots, k\}$  and measured, so we have obtained the values  $y_{ij}$  for  $j = 1, 2, \dots, n_i$  for  $i = 1, 2, \dots, k$ . (Where  $n_i$  is the number of the units finally found in the  $i$ -th group for  $i = 1, 2, \dots, k$ . In the end, it holds  $n = n_1 + n_2 + \dots + n_k$ .)
- We assume  $y_{ij} \approx \mu_i$  and we have  $y_{ij} = \mu_i + \varepsilon_{ij}$ , where  $\varepsilon_{ij}$  is a random deviation (error).
- The random deviation is caused by the intrinsic properties of the statistical unit

# One-way ANOVA: Summary

---



## **Second:**

- We prepare the groups  $i = 1, 2, \dots, k$  at the beginning.
- When measuring the value  $y_{ij}$ , we select a unit from the group  $i$  first and measure its value  $y_{ij}$  then for  $j = 1, 2, \dots, n_i$  for  $i = 1, 2, \dots, k$ .
- The random deviation  $\varepsilon_i$  here is caused either by the intrinsic properties of the system (further unknown / “random” / unconsidered factors),  
or by random errors in the measurement itself. (!)

# One-way ANOVA: Summary

---



## Remarks:

- In practice, the data may be obtained in either way (first or second).
  - In either case (first or second), the group which the unit belongs to is assumed to be known exactly, i.e. without any doubts.
  - Assuming  $y_{ij} \approx \mu_i$ , even the dependent values  $y_{ij}$  may be measured exactly, i.e. without any measurement error, the random deviation  $\varepsilon_{ij} = y_{ij} - \mu_i$  being caused by the intrinsic properties (other unknown / “random” / unconsidered factors).
  - For the purpose of the mathematical analysis, we assume the second case only.
-

# One-way ANOVA: Assumptions

---



- There are  $k$  groups of sizes  $n_1, n_2, \dots, n_k$  given, known, and fixed before the measurements.
- Moreover, we are given  $k$  numbers  $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}$ , which are the expected values of the variable “ $y$ ” in the respective groups.
- We have  $n = n_1 + n_2 + \dots + n_k$  random variables

$$Y_{11}, Y_{12}, \dots, Y_{1n_1}, Y_{21}, Y_{22}, \dots, Y_{2n_2}, \dots \dots \dots, Y_{k1}, Y_{k2}, \dots, Y_{kn_k}$$

which are assumed to be independent (or uncorrelated).

- We also have  $n$  random variables

$$\varepsilon_{11}, \varepsilon_{12}, \dots, \varepsilon_{1n_1}, \varepsilon_{21}, \varepsilon_{22}, \dots, \varepsilon_{2n_2}, \dots \dots \dots, \varepsilon_{k1}, \varepsilon_{k2}, \dots, \varepsilon_{kn_k}$$

# One-way ANOVA: Assumptions



- We stack the group expected values  $\mu_1, \mu_2, \dots, \mu_k$  into a vector:

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_1 \\ \mu_2 \\ \vdots \\ \mu_2 \\ \vdots \\ \vdots \\ \vdots \\ \mu_k \\ \vdots \\ \mu_k \end{pmatrix} \begin{matrix} \left. \vphantom{\begin{matrix} \mu_1 \\ \vdots \\ \mu_1 \\ \mu_2 \\ \vdots \\ \mu_2 \\ \vdots \\ \vdots \\ \vdots \\ \mu_k \\ \vdots \\ \mu_k \end{matrix}} \right\} n_1\text{-times} \\ \left. \vphantom{\begin{matrix} \mu_1 \\ \vdots \\ \mu_1 \\ \mu_2 \\ \vdots \\ \mu_2 \\ \vdots \\ \vdots \\ \vdots \\ \mu_k \\ \vdots \\ \mu_k \end{matrix}} \right\} n_2\text{-times} \\ \left. \vphantom{\begin{matrix} \mu_1 \\ \vdots \\ \mu_1 \\ \mu_2 \\ \vdots \\ \mu_2 \\ \vdots \\ \vdots \\ \vdots \\ \mu_k \\ \vdots \\ \mu_k \end{matrix}} \right\} \dots\text{-times} \\ \left. \vphantom{\begin{matrix} \mu_1 \\ \vdots \\ \mu_1 \\ \mu_2 \\ \vdots \\ \mu_2 \\ \vdots \\ \vdots \\ \vdots \\ \mu_k \\ \vdots \\ \mu_k \end{matrix}} \right\} n_k\text{-times} \end{matrix}$$



# One-way ANOVA: Assumptions



- We stack the random variables  $Y_{11}, \dots, Y_{1n_1}, Y_{21}, \dots, Y_{2n_2}, \dots, Y_{k1}, \dots, Y_{kn_k}$  into a random vector:

$$\mathbf{Y} = \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ Y_{21} \\ \vdots \\ Y_{2n_2} \\ \vdots \\ \vdots \\ \vdots \\ Y_{k1} \\ \vdots \\ Y_{kn_k} \end{pmatrix} \begin{array}{l} \left. \vphantom{\begin{matrix} Y_{11} \\ \vdots \\ Y_{1n_1} \end{matrix}} \right\} n_1\text{-times} \\ \left. \vphantom{\begin{matrix} Y_{21} \\ \vdots \\ Y_{2n_2} \end{matrix}} \right\} n_2\text{-times} \\ \left. \vphantom{\begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix}} \right\} \dots\text{-times} \\ \left. \vphantom{\begin{matrix} Y_{k1} \\ \vdots \\ Y_{kn_k} \end{matrix}} \right\} n_k\text{-times} \end{array}$$

# One-way ANOVA: Assumptions



- We stack the random variables  $\varepsilon_{11}, \dots, \varepsilon_{1n_1}, \varepsilon_{21}, \dots, \varepsilon_{2n_2}, \dots, \dots, \varepsilon_{k1}, \dots, \varepsilon_{kn_k}$  into a random vector:

$$\boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1n_1} \\ \varepsilon_{21} \\ \vdots \\ \varepsilon_{2n_2} \\ \vdots \\ \vdots \\ \vdots \\ \varepsilon_{k1} \\ \vdots \\ \varepsilon_{kn_k} \end{pmatrix} \begin{array}{l} \left. \vphantom{\begin{pmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1n_1} \end{pmatrix}} \right\} n_1\text{-times} \\ \left. \vphantom{\begin{pmatrix} \varepsilon_{21} \\ \vdots \\ \varepsilon_{2n_2} \end{pmatrix}} \right\} n_2\text{-times} \\ \left. \vphantom{\begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix}} \right\} \dots\text{-times} \\ \left. \vphantom{\begin{pmatrix} \varepsilon_{k1} \\ \vdots \\ \varepsilon_{kn_k} \end{pmatrix}} \right\} n_k\text{-times} \end{array}$$

# One-way ANOVA: The Classical Assumptions

---



- We have the underlying probability space  $(\Omega, \mathcal{F}, P)$ .
- Let  $\omega \in \Omega$  be the outcome of the random experiment.
- We have

$$\mathbf{y} = \mathbf{Y}(\omega) = \boldsymbol{\mu} + \boldsymbol{\varepsilon}(\omega)$$

## In other words:

- The measured values  $y_{11}, \dots, y_{1n_1}, y_{21}, \dots, y_{2n_2}, \dots, y_{k1}, \dots, y_{kn_k}$  are the numerical outcomes  $Y_{11}(\omega), \dots, Y_{kn_k}(\omega)$  of the random experiment.
- The numerical outcomes  $Y_{11}(\omega), \dots, Y_{kn_k}(\omega)$  are obtained so that the numerical outcomes  $\varepsilon_{11}(\omega), \dots, \varepsilon_{kn_k}(\omega)$  of the random experiment

# One-way ANOVA: The Classical Assumptions

---



**The classical assumptions of One-Way ANOVA are:**

$$Y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \text{and} \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$$

where  $I$  denotes the  $n \times n$  identity matrix

and  $\mathbf{0}$  denotes the  $n \times 1$  zero vector.

It follows that:

$$E[Y_{ij}] = \mu_i \quad \text{and} \quad E[\varepsilon_{ij}] = 0$$

# One-way ANOVA: The Classical Assumptions



The classical assumptions  $Y \sim \mathcal{N}(\mu, \sigma^2 I)$  and  $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$  also mean that

$$\text{Var}(Y) = \text{Var}(\varepsilon) = \sigma^2 I = \begin{pmatrix} \sigma^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma^2 & 0 & \dots & 0 \\ 0 & 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma^2 \end{pmatrix}$$

That is,

- $\text{Var}(Y_{ij}) = \text{Var}(\varepsilon_{ij}) = \sigma^2$  for  $j = 1, 2, \dots, n_i$

for  $i = 1, 2, \dots, k$  for some  $\sigma^2 \in \mathbb{R}$

- and the random variables  $Y_{11}, \dots, Y_{kn_k}$  or  $\varepsilon_1, \dots, \varepsilon_{kn_k}$

**homoscedasticity**,  
i.e. the variance  
is the same

# One-way ANOVA: The purpose of the analysis

---



By the Classical Assumptions, we have:

$$E[Y_{11}] = E[Y_{12}] = \dots = E[Y_{1n_1}] = \mu_1$$

$$E[Y_{21}] = E[Y_{22}] = \dots = E[Y_{2n_2}] = \mu_2$$

...

$$E[Y_{k1}] = E[Y_{k2}] = \dots = E[Y_{kn_k}] = \mu_k$$

The purpose is to test the null hypothesis

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

# One-way ANOVA: Solution

---



Given the sample

$$y_{11}, y_{12}, y_{13}, \dots, y_{1n_1}$$

$$y_{21}, y_{22}, \dots, y_{2n_2}$$

$$y_{31}, y_{32}, y_{33}, y_{34}, y_{35}, \dots, y_{3n_3}$$

...

$$y_{k1}, y_{k2}, \dots, y_{kn_k}$$

of the random variables

$$Y_{11}, Y_{12}, Y_{13}, \dots, Y_{1n_1}$$

$$Y_{21}, Y_{22}, \dots, Y_{2n_2}$$

$$Y_{31}, Y_{32}, Y_{33}, Y_{34}, Y_{35}, \dots, Y_{3n_3}$$

...

$$Y_{k1}, Y_{k2}, \dots, Y_{kn_k}$$

that is

---

# One-way ANOVA: Solution

---



...and

assuming (among others) that the expected values are the same in each group

$$E[Y_{ij}] = \mu_i \quad \text{for } j = 1, 2, \dots, n_i \quad \text{for } i = 1, 2, \dots, k$$

it is our purpose to test the null hypothesis that

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

## SOLUTION:

- We shall express this situation as a model of Multiple Linear Regression.
  - We shall then use the theory of Multiple Linear Regression ("Theorem 8") to make up an  $F$ -test for the null hypothesis  $H_0$ .
-



# One-Way ANOVA as a model of Multiple Linear Regression



# One-way ANOVA: Model of Multiple Linear Regression



Given the random variables  $Y_{11}, Y_{12}, \dots, Y_{1n_1}, Y_{21}, Y_{22}, \dots, Y_{2n_2}, \dots, Y_{k1}, Y_{k2}, \dots, Y_{kn_k}$ , we can express the original assumption (on the left) in terms of Multiple Linear Regression (on the right):

$$\begin{array}{llll} \mathbf{E}[Y_{1j}] = \mu_1 & \rightarrow & \mathbf{E}[Y_{1j}] = \mu_1 & \text{for } j = 1, 2, \dots, n_1 \\ \mathbf{E}[Y_{2j}] = \mu_2 & \rightarrow & \mathbf{E}[Y_{2j}] = \mu_1 + \beta_2 & \text{for } j = 1, 2, \dots, n_2 \\ \mathbf{E}[Y_{3j}] = \mu_3 & \rightarrow & \mathbf{E}[Y_{3j}] = \mu_1 + \beta_3 & \text{for } j = 1, 2, \dots, n_3 \\ \vdots & & \vdots & \\ \mathbf{E}[Y_{kj}] = \mu_k & \rightarrow & \mathbf{E}[Y_{kj}] = \mu_1 + \beta_2 + \beta_3 + \dots + \beta_k & \text{for } j = 1, 2, \dots, n_k \end{array}$$

# One-way ANOVA: Model of Multiple Linear Regression

---



Having

$$E[Y_{1j}] = \mu_1 \quad \text{for } j = 1, 2, \dots, n_1$$

$$E[Y_{ij}] = \mu_i = \mu_1 + \beta_j \quad \text{for } j = 1, 2, \dots, n_i \quad \text{for } i = 2, \dots, k$$

the original hypothesis

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

is equivalent with

$$H_0: \beta_2 = \dots = \beta_k = 0$$

which we can test (recall "Theorem 8").

---

# One-way ANOVA: Model of Multiple Linear Regression



Putting yet  $\beta_1 = \mu_1$  and by using the notation

$$\mathbf{Y} = \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ Y_{21} \\ \vdots \\ Y_{2n_2} \\ \vdots \\ \vdots \\ \vdots \\ Y_{k1} \\ \vdots \\ Y_{kn_k} \end{pmatrix} \begin{array}{l} \left. \vphantom{\begin{matrix} Y_{11} \\ \vdots \\ Y_{1n_1} \end{matrix}} \right\} n_1\text{-times} \\ \left. \vphantom{\begin{matrix} Y_{21} \\ \vdots \\ Y_{2n_2} \end{matrix}} \right\} n_2\text{-times} \\ \left. \vphantom{\begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix}} \right\} \dots\text{-times} \\ \left. \vphantom{\begin{matrix} Y_{k1} \\ \vdots \\ Y_{kn_k} \end{matrix}} \right\} n_k\text{-times} \end{array} \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{pmatrix} \quad \text{and} \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1n_1} \\ \varepsilon_{21} \\ \vdots \\ \varepsilon_{2n_2} \\ \vdots \\ \vdots \\ \vdots \\ \varepsilon_{k1} \\ \vdots \\ \varepsilon_{kn_k} \end{pmatrix} \begin{array}{l} \left. \vphantom{\begin{matrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1n_1} \end{matrix}} \right\} n_1\text{-times} \\ \left. \vphantom{\begin{matrix} \varepsilon_{21} \\ \vdots \\ \varepsilon_{2n_2} \end{matrix}} \right\} n_2\text{-times} \\ \left. \vphantom{\begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix}} \right\} \dots\text{-times} \\ \left. \vphantom{\begin{matrix} \varepsilon_{k1} \\ \vdots \\ \varepsilon_{kn_k} \end{matrix}} \right\} n_k\text{-times} \end{array}$$

# One-way ANOVA: Model of Multiple Linear Regression



... and by considering the matrix

$$\mathbf{X} = \begin{pmatrix} \underbrace{\beta_1}_{\underbrace{\quad}} & \underbrace{\beta_2}_{\underbrace{\quad}} & \underbrace{\dots}_{\underbrace{\quad}} & \underbrace{\beta_k}_{\underbrace{\quad}} \\ \vdots & & & \\ 1 & & & \\ 1 & 1 & & \\ \vdots & \vdots & & \\ 1 & 1 & & \\ \vdots & & \ddots & \\ \vdots & & & \\ 1 & & & 1 \\ \vdots & & & \vdots \\ 1 & & & 1 \end{pmatrix} \begin{array}{l} \left. \vphantom{\begin{matrix} \vdots \\ 1 \\ 1 \\ \vdots \\ 1 \\ \vdots \\ \vdots \\ 1 \\ \vdots \\ 1 \end{matrix}} \right\} n_1\text{-times} \\ \left. \vphantom{\begin{matrix} 1 \\ 1 \\ \vdots \\ 1 \\ \vdots \\ \vdots \\ 1 \\ \vdots \\ 1 \end{matrix}} \right\} n_2\text{-times} \\ \left. \vphantom{\begin{matrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{matrix}} \right\} \dots\text{-times} \\ \left. \vphantom{\begin{matrix} 1 \\ \vdots \\ 1 \end{matrix}} \right\} n_k\text{-times} \end{array}$$

Notice that

$$\text{rank}(\mathbf{X}) = k$$

# One-way ANOVA: Model of Multiple Linear Regression

---



... we can write the One-Way ANOVA in terms of Multiple Linear Regression as follows:

$$Y = X\beta + \varepsilon$$

where  $\beta_1 = \mu_1$  plays the rôle of the intercept term

and  $\beta_2 = \mu_2 - \mu_1, \dots, \beta_k = \mu_k - \mu_1$  are the other regression coefficients.

Recall that we wish to test the null hypothesis

$$H_0: \beta_2 = \dots = \beta_k = 0$$

---

# The Coefficient of Determination ( $R^2$ ): Th. 8: Corollary



Let:  $\bar{y} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}$

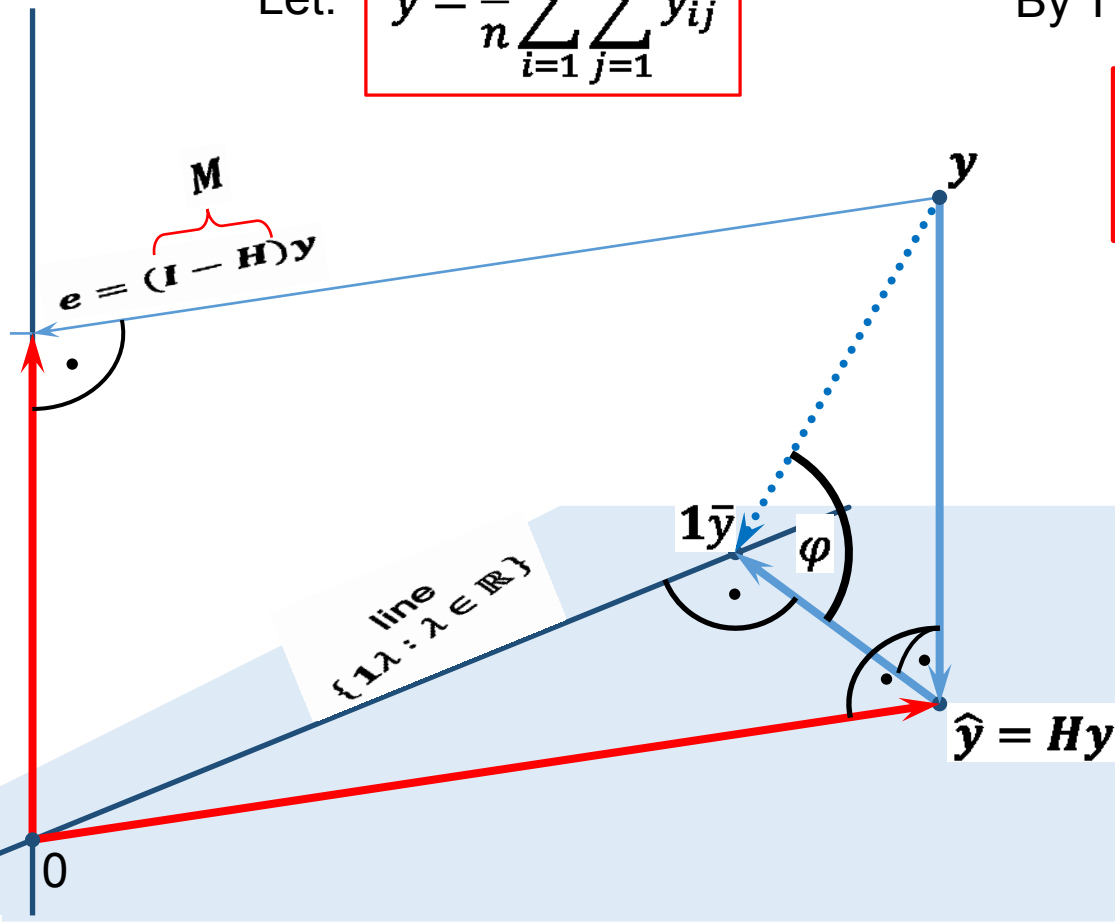
By Theorem 8:

$$(\cotan \varphi)^2 / \frac{k-1}{n-k} \sim F_{k-1, n-k}$$

$$\cotan^2 \varphi = \frac{(\hat{y} - \mathbf{1}\bar{y})^T (\hat{y} - \mathbf{1}\bar{y})}{\text{RSS}} = \frac{R^2}{1 - R^2}$$

$\{X\beta : \beta \in \mathbb{R}^k\}^\perp$   
 (the orthogonal complement =  
 = the space of the residuals)  
 subspace  
 of dimension

$$n - k$$



the line is a subspace  
 of dimension  $1$

the dimension of its complement within  
 the subspace of dimension  $k$  is  $-1$

$\{X\beta : \beta \in \mathbb{R}^k\}$   
 (the linear hull of the columns of  $X$ )  
 subspace of dimension  $k$



# The Coefficient of Determination ( $R^2$ ): Th. 8: Corollary

Let: 
$$\bar{y} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}$$

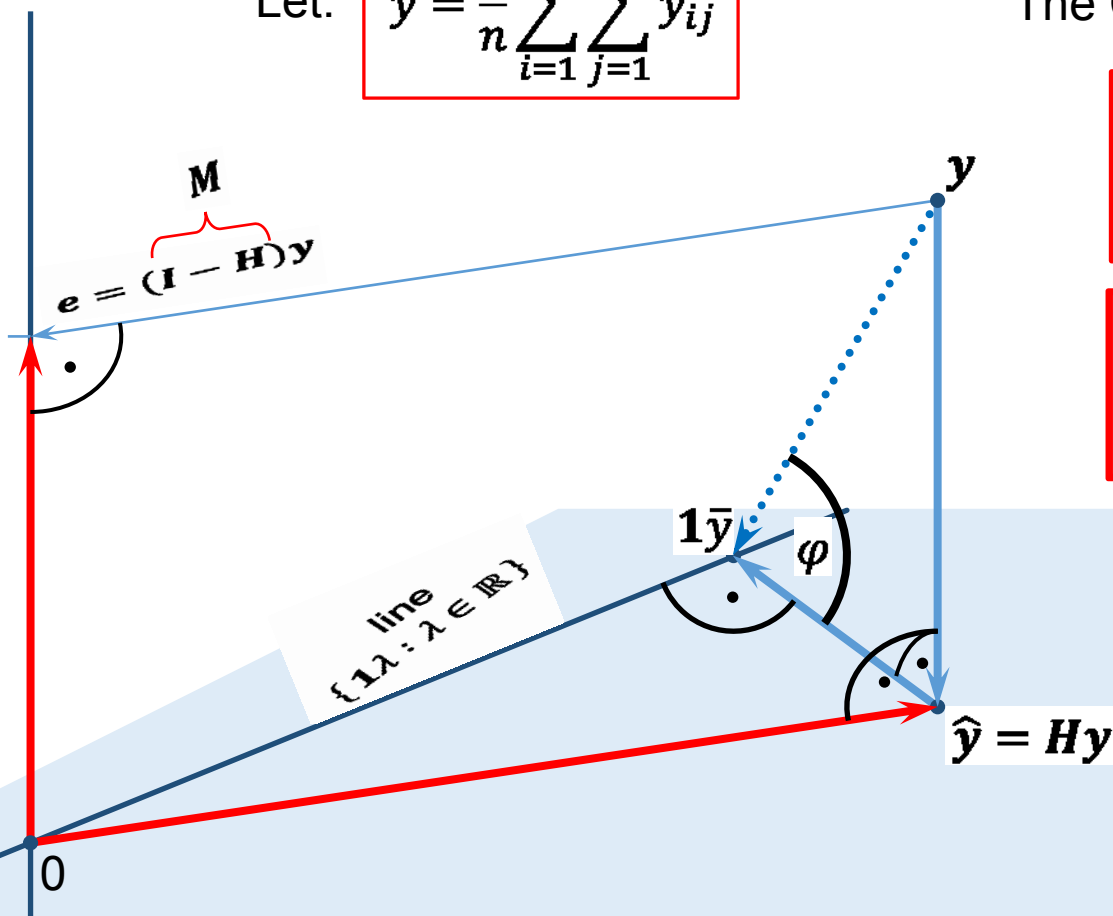
The Coefficient of Determination:

$$R^2 = \cos^2 \varphi = \frac{(\hat{y} - \mathbf{1}\bar{y})^T (\hat{y} - \mathbf{1}\bar{y})}{(\mathbf{y} - \mathbf{1}\bar{y})^T (\mathbf{y} - \mathbf{1}\bar{y})}$$

$$\frac{R^2}{1 - R^2} = \frac{\cos^2 \varphi}{\sin^2 \varphi} = \cotan^2 \varphi$$

$\{X\beta : \beta \in \mathbb{R}^k\}^\perp$   
(the orthogonal complement =  
the space of the residuals)  
subspace  
of dimension

$$n - k$$



$\{X\beta : \beta \in \mathbb{R}^k\}$   
(the linear hull of the columns of  $X$ )  
subspace of dimension

$$k$$

the line is a subspace  
of dimension  $1$

the dimension of its complement within  
the subspace of dimension  $k$  is  $-1$





# The Coefficient of Determination ( $R^2$ ): $TSS=RSS+RegSS$

Let: 
$$\bar{y} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}$$

$$TSS = (\mathbf{y} - \mathbf{1}\bar{y})^T (\mathbf{y} - \mathbf{1}\bar{y})$$

$$RegSS = (\hat{\mathbf{y}} - \mathbf{1}\bar{y})^T (\hat{\mathbf{y}} - \mathbf{1}\bar{y})$$

$$RSS = (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{e}^T \mathbf{e}$$

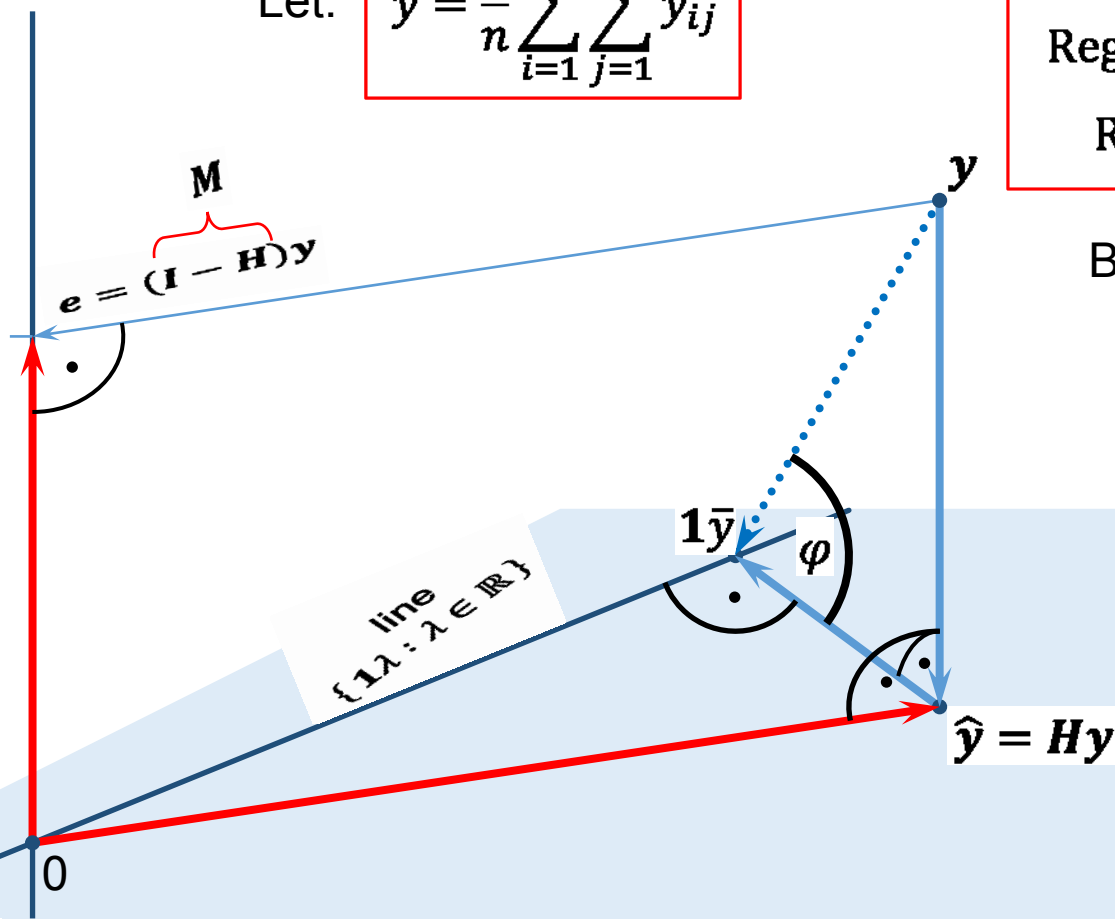
By the Pythagoras Theorem:

$$TSS = RSS + RegSS$$

$\{\mathbf{X}\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^k\}^\perp$

(the orthogonal complement = the space of the residuals) subspace of dimension

$$n - k$$



line  $\{\mathbf{1}\lambda : \lambda \in \mathbb{R}\}$

$\{\mathbf{X}\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^k\}$

(the linear hull of the columns of  $\mathbf{X}$ ) subspace of dimension

$$k$$

the line is a subspace of dimension  $1$

the dimension of its complement within the subspace of dimension  $k$  is  $-1$

# The Coefficient of Determination ( $R^2$ )



Assuming  $\mathbf{1} \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^k\}$ , define the

**Coefficient of Determination:**

$$R^2 = \frac{\text{RegSS}}{\text{TSS}} = \frac{\text{TSS} - \text{RSS}}{\text{TSS}}$$

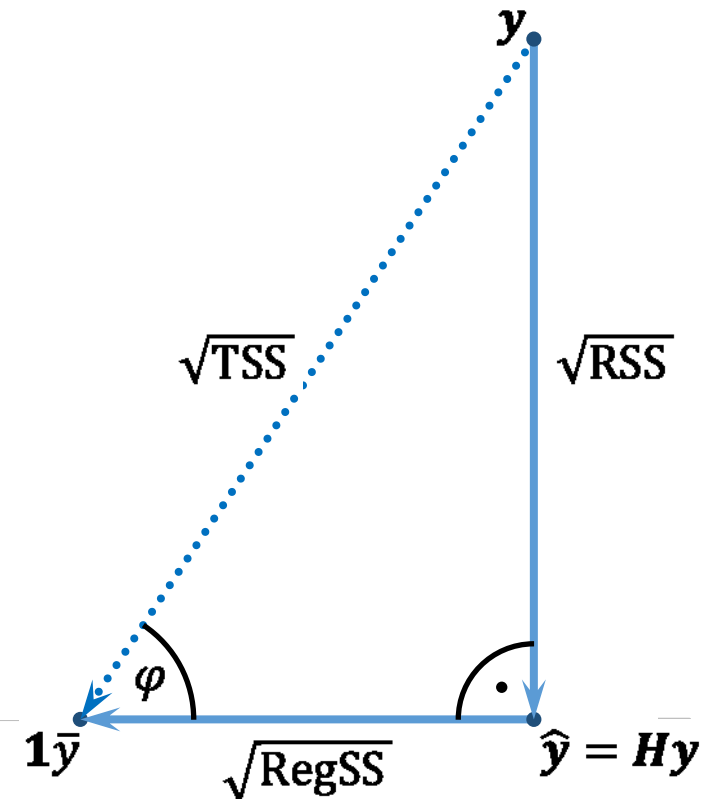
$$R^2 = \cos^2 \varphi = \frac{\text{RegSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

$$\cotan^2 \varphi = \frac{\cos^2 \varphi}{\sin^2 \varphi} = \frac{R^2}{1 - R^2} = \frac{\text{RegSS}}{\text{RSS}}$$

$$\text{TSS} = (\mathbf{y} - \mathbf{1}\bar{y})^T (\mathbf{y} - \mathbf{1}\bar{y})$$

$$\text{RegSS} = (\hat{\mathbf{y}} - \mathbf{1}\bar{y})^T (\hat{\mathbf{y}} - \mathbf{1}\bar{y})$$

$$\text{RSS} = (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{e}^T \mathbf{e}$$



# The Coefficient of Determination ( $R^2$ ): Th. 8: Corollary



**Theorem 8: Corollary:** Assume for simplicity that  $\text{rank}(X) = k$  and assume that  $\mathbf{1} \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^k\}$ . Under the hypothesis that

$$\beta_2 = \dots = \beta_k = 0$$

it holds

$$\begin{aligned} (\cotan \varphi)^2 / \frac{k-1}{n-k} &= \frac{R^2}{1-R^2} / \frac{k-1}{n-k} \sim F_{k-1, n-k} \\ &= \frac{\text{RegSS}}{\text{RSS}} / \frac{k-1}{n-k} \sim F_{k-1, n-k} \end{aligned}$$

# One-way ANOVA: Model of Multiple Linear Regression



We have

$$\text{RegSS} = (\hat{\mathbf{y}} - \mathbf{1}\bar{y})^T (\hat{\mathbf{y}} - \mathbf{1}\bar{y}) = \sum_{i=1}^k \sum_{j=1}^{n_i} (\hat{y}_{ij} - \bar{y})^2 \quad \text{where} \quad \bar{y} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}$$

where

$$\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$$

where

$$\mathbf{H} = \mathbf{X}\mathbf{C}\mathbf{X}^T$$

where

$$\mathbf{C} = (\mathbf{X}^T\mathbf{X})^{-1}$$

# One-way ANOVA: Model of Multiple Linear Regression



Calculate:

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} n & n_2 & n_3 & n_4 & \dots & n_k \\ n_2 & n_2 & & & & \\ n_3 & & n_3 & & & \\ n_4 & & & n_4 & & \\ \vdots & & & & \ddots & \\ n_k & & & & & n_k \end{pmatrix}$$

where

$$n = n_1 + n_2 + \dots + n_k = \sum_{i=1}^k n_i$$

# One-way ANOVA: Model of Multiple Linear Regression



Calculate  $C = (X^T X)^{-1}$ :

$$(X^T X)^{-1} = \begin{pmatrix} +1/n_1 & -1/n_1 & -1/n_1 & -1/n_1 & \dots & -1/n_1 \\ -1/n_1 & 1/n_2 + 1/n_1 & +1/n_1 & +1/n_1 & \dots & +1/n_1 \\ -1/n_1 & +1/n_1 & 1/n_3 + 1/n_1 & +1/n_1 & \dots & +1/n_1 \\ -1/n_1 & +1/n_1 & +1/n_1 & 1/n_4 + 1/n_1 & \dots & +1/n_1 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ -1/n_1 & +1/n_1 & +1/n_1 & +1/n_1 & \dots & 1/n_k + 1/n_1 \end{pmatrix}$$

# One-way ANOVA: Model of Multiple Linear Regression



Calculate  $H = XCX^T = X(X^T X)^{-1} X^T$ :

$$XCX^T = H = \begin{pmatrix} \overbrace{\begin{matrix} 1/n_1 & \dots & 1/n_1 \\ \vdots & \ddots & \vdots \\ 1/n_1 & \dots & 1/n_1 \end{matrix}}^{n_1} & & & \\ & \overbrace{\begin{matrix} 1/n_2 & \dots & 1/n_2 \\ \vdots & \ddots & \vdots \\ 1/n_2 & \dots & 1/n_2 \end{matrix}}^{n_2} & & \\ & & \dots & \\ & & & \overbrace{\begin{matrix} 1/n_k & \dots & 1/n_k \\ \vdots & \ddots & \vdots \\ 1/n_k & \dots & 1/n_k \end{matrix}}^{n_k} \end{pmatrix} \begin{matrix} \left. \begin{matrix} \dots \\ \dots \\ \dots \end{matrix} \right\} n_1\text{-times} \\ \left. \begin{matrix} \dots \\ \dots \\ \dots \end{matrix} \right\} n_2\text{-times} \\ \left. \begin{matrix} \dots \\ \dots \\ \dots \end{matrix} \right\} \dots\text{-times} \\ \left. \begin{matrix} \dots \\ \dots \\ \dots \end{matrix} \right\} n_k\text{-times} \end{matrix}$$

# One-way ANOVA: Model of Multiple Linear Regression

---



Calculate:

$$\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$$

That is:

$$\hat{y}_{ij} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij} \quad \text{for } i = 1, 2, \dots, k$$

Denote:

$$\bar{y}_i = \hat{y}_{ij} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij} \quad \text{for } i = 1, 2, \dots, k$$

and call this ↗ quantity **the sample mean of the  $i$ -th group** for  $i = 1, 2, \dots, k$ .

---



# One-way ANOVA: Model of Multiple Linear Regression



Calculate:

$$\text{RegSS} = (\hat{\mathbf{y}} - \mathbf{1}\bar{y})^T (\hat{\mathbf{y}} - \mathbf{1}\bar{y}) = \sum_{i=1}^k \sum_{j=1}^{n_i} (\hat{y}_{ij} - \bar{y})^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{y}_i - \bar{y})^2$$

Denote:

$$\text{SS}_B = \text{RegSS} = \sum_{i=1}^k n_i \times (\bar{y}_i - \bar{y})^2$$

and call this ↗ quantity **the sum of squares “between”**

where  $\bar{y} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}$  and  $\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$

# One-way ANOVA: Model of Multiple Linear Regression



Calculate:

$$\text{RSS} = (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \hat{y}_{ij})^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$$

Denote:

$$SS_W = \text{RSS} = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$$

and call this ↗ quantity **the sum of squares “within”**

where

$$\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$$

## One-way ANOVA: Remarks: Sample variances

---



The quantity

$$MS_B = \frac{SS_B}{DF_B} = \frac{\text{RegSS}}{k - 1} = \frac{\sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{y}_i - \bar{y})^2}{k - 1}$$

is the sample variance between the groups.

The quantity

$$MS_W = \frac{SS_W}{DF_W} = \frac{\text{RSS}}{n - k} = \frac{\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2}{n - k}$$

is the sample variance within the groups.

---

# One-Way ANOVA: The $F$ -test



- The  $F$ -test

# One-Way ANOVA: the $F$ -test

---



To test the null hypothesis

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

against the alternative hypothesis

$$H_1: \mu_i \neq \mu_j \quad \text{for some } i, j \in \{1, 2, \dots, k\}$$

calculate the statistic

$$F = \frac{SS_B}{SS_W} \bigg/ \frac{DF_B}{DF_W} = \frac{\text{RegSS}}{\text{RSS}} \bigg/ \frac{k-1}{n-k} = \frac{\sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{y}_i - \bar{y})^2}{\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2} \bigg/ \frac{k-1}{n-k}$$

# One-Way ANOVA: the $F$ -test

---



Recall that, if  $H_0$  holds true, then

$$F \sim F_{k-1, n-k}$$

**Choose the level of significance**, a small number  $\alpha > 0$ , such as  $\alpha = 5\%$ .

**and calculate the critical value**

$$c = F_{k-1, n-k}(1 - \alpha)$$

If  $F \in [c, +\infty)$ , **the critical region**, then **reject** the null hypothesis.

If  $F \in [0, c)$ , then **do not reject** (or **fail to reject**) the null hypothesis.

---