

PARTE I – PLANEJAMENTO DE EXPERIMENTOS

ANOVA 1 FATOR

*Prof. Anna Carla Araujo
COPPE/UFRJ*



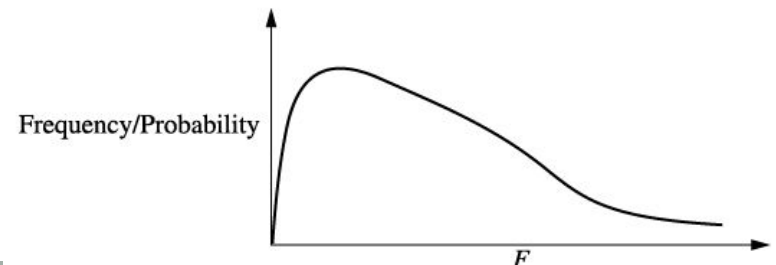
Limitações do teste com amostras não pareadas

- O número de grupos está limitado a dois
- A interação entre os grupos não pode ser testada
- Não há como relacionar as variáveis pareadas

Distribuição F

- + Avalia a razão entre a variância de duas variáveis independentes com variâncias S_1 e S_2 e número de amostras n_1 e n_2 .

$$F = \frac{[(n_1 - 1)S_1^2 / \sigma_1^2] / (n_1 - 1)}{[(n_2 - 1)S_2^2 / \sigma_2^2] / (n_2 - 1)}$$



Introdução a Análise de Variância

Table 2.2 Sales Data of Salesmen

	Salesman j					
	1	2	3	j
1	Y_{11}	Y_{12}	Y_{13}	Y_{1j}
2	Y_{21}	Y_{22}	Y_{23}	Y_{2j}
3	Y_{31}	Y_{32}	Y_{33}	Y_{3j}
Replication i
i	Y_{i1}	Y_{i2}	Y_{i3}	Y_{ij}
...
n	Y_{n1}	Y_{n2}	Y_{n3}	Y_{nj}

- Y_{ij} is the i th sales data (i th replication) of the j th salesman.
- Salesman is the factor which has effect on the response variable Y_{ij} . Let it be factor A .
- a is the number of salesmen of the company. This is also known as the number of levels/treatments of factor A .
- n is the number of sales data under each salesman. This is also known as the number of replications under each level of factor A .

Modelo:

$$Y_{ij} = \bar{y} + A_j + e_{ij}$$

where

\bar{y} is the overall mean of the sales revenue.

A_j is the effect of the j th treatment of factor A (salesman) on the response, and

e_{ij} is the random error associated with the i th replication of the j th treatment of factor A .

The decomposition of total variability into its component parts is called analysis of variance (ANOVA). The partitioning of the total variability is presented below.

The total corrected sum of squares is given by

$$SS_{\text{Total}} = \sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{Y}_{..})^2$$

The above equation is re-written by adding and subtracting the term $\bar{Y}_{.j}$ as per the following presentation:

$$\sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{Y}_{..})^2 = \sum_{i=1}^n \sum_{j=1}^a [(\bar{Y}_{.j} - \bar{Y}_{..}) + (Y_{ij} - \bar{Y}_{.j})]^2$$

$$\sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{Y}_{..})^2 = \sum_{i=1}^n \sum_{j=1}^a (\bar{Y}_{.j} - \bar{Y}_{..})^2 + \sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{Y}_{.j})^2 + 2 \left(\sum_{i=1}^n \sum_{j=1}^a (\bar{Y}_{.j} - \bar{Y}_{..}) \right) \left(\sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{Y}_{.j}) \right)$$

The following part of the last term on the right hand side of the equation is equal to 0 because of the following:

$$\sum_{i=1}^n (Y_{ij} - \bar{Y}_{.j}) = (Y_{.j} - n\bar{Y}_{.j}) = \left[Y_{.j} - n \left(\frac{Y_{.j}}{n} \right) \right] = 0$$

Therefore, the total sum of squares corrected is reduced to the following:

$$\sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{Y}_{..})^2 = \sum_{i=1}^n \sum_{j=1}^a (\bar{Y}_{.j} - \bar{Y}_{..})^2 + \sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{Y}_{.j})^2$$

$$\sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{Y}_{..})^2 = n \sum_{j=1}^a (\bar{Y}_{.j} - \bar{Y}_{..})^2 + \sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{Y}_{.j})^2$$

□ Total sum of squares = Sum of squares of treatment + Sum of squares of error

$$SS_{\text{Total}} = SS_A + SS_{\text{Error}}$$

□ Total sum of squares = Sum of squares of treatment + Sum of squares of error

$$SS_{\text{Total}} = SS_A + SS_{\text{Error}}$$

Mean sum of squares of treatment =

$$\frac{\text{Sum of squares of treatment}}{\text{Degrees of freedom of treatment}} = \frac{SS_A}{(a - 1)}$$

Mean sum of squares of error = $\frac{\text{Sum of squares of error}}{\text{Degrees of freedom of error}}$

$$= \frac{SS_{\text{Error}}}{a(n - 1)} = \frac{SS_{\text{Error}}}{(N - a)}$$

where N is the total number of observations in the experiment (an).

The hypotheses of the model are given as:

Null hypothesis, $H_0: \bar{Y}_1 = \bar{Y}_2 = \bar{Y}_3 = \bar{Y}_4 = \dots = \bar{y}_a$

Alternate hypothesis, H_1 : Treatment means are not equal for at least one pair of the treatment means. The generalized results are summarized in Table 2.3.

Table 2.3 Generalized Results of ANOVA with Single Factor

Source of variation	Degrees of freedom	Sum of squares	Mean sum of squares	F-ratio
Between treatments	$a - 1$	$SS_{\text{treatment}}$	$SS_{\text{treatment}}/(a - 1)$	$MSS_{\text{treatment}}/MSS_{\text{error}}$
Within treatments	$a(n - 1)$	SS_{error}	$SS_{\text{error}}/[a(n - 1)]$	
Total	$N - 1$	SS_{total}		

Teste da média para amostras pareadas

Hipótese H₀: médias são iguais, logo a diferença é zero.

Hipótese H₁: médias são diferentes, logo a diferença não é zero.

Variância Total:

$$\bar{y}_{ij} = \sum_{i=1}^n \sum_{j=1}^a y_{ij} / N$$

$$SS_{total} = \sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{y}_{ij})^2$$

Variância entre as m amostras:

Média da amostra j
j: [1,a]

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

$$SS_{treat} = \sum_{i=1}^n n_j (\bar{y}_i - \bar{y}_{ij})^2$$

Variância dentro das amostras:

$$SS_{erro} = \sum_{i=1}^n \sum_{j=1}^a (Y_{ij} - \bar{y}_i)^2$$

Tabela de Análise de Variância

Fonte de Variação	Soma dos Quadrados	Graus de Liberdade	Quadrado s Médios	Variavel de teste F
Entre amostras	S_{entre}	$a-1$	$\frac{SS_{\text{treat}}}{a - 1}$	$\frac{MS_{\text{treat}}}{MS_{\text{erro}}}$
Dentro das amostras	S_{erro}	$N-a$	$\frac{SS_{\text{erro}}}{N - a}$	
Total	S_{total}	$N-1$	$\frac{S_{\text{total}}}{(N - 1)}$	

Teste de Hipótese da ANOVA

Hipótese H_0 : médias são iguais, logo a diferença é zero.

Hipótese H_1 : médias são diferentes, logo a diferença não é zero.

Uma vez fixado o nível de significância α do teste H_0 deve ser rejeitado se $F(\text{observações}) > F(1-\alpha)$

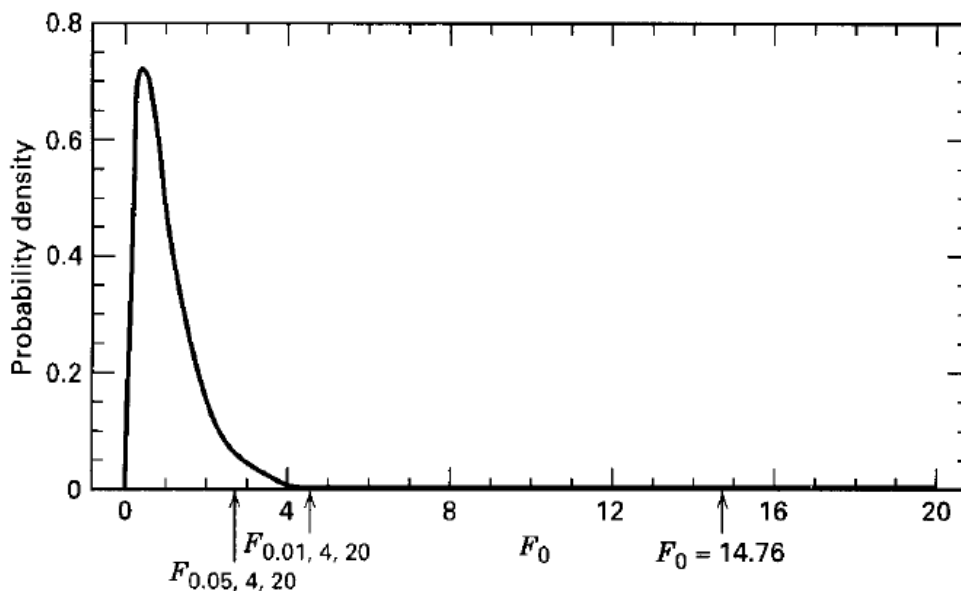
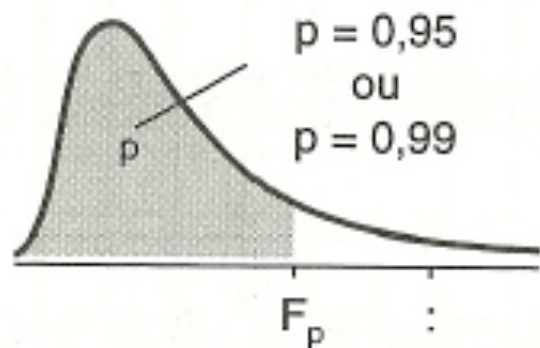


Figure 3-3 The reference distribution ($F_{4,20}$) for the test statistic F_0 in Example 3-1.

Exemplo

Fornecer os quantis $F_{0,95}$ (em cima) e $F_{0,99}$ (embaixo) em função do nº de g.l. numerador v_1 (coluna) e do nº de g.l. denominador v_2 (linha).
 F tem distribuição F com v_1 g.l. no numerador e v_2 g.l. no denominador $P(F \leq F_{0,95}) = 0,95$ e $P(F \leq F_{0,99}) = 0,99$

$v_2 \backslash v_1$	1	2
1	161,45	199,50
2	4052,18	4999,50
3	18,51	19,00
4	98,50	99,00
5	10,13	9,55
6	34,12	30,82
7	7,71	6,94
8	21,20	18,00
9	6,61	5,79
10	16,26	13,27
20	5,99	5,14
40	13,75	10,92
	5,59	4,74
	12,25	9,55
	5,32	4,46
	11,26	8,65
	5,12	4,26
	10,56	8,02
	4,96	4,10
	10,04	7,56
	4,35	3,49
	8,10	5,85
	3,88	3,23

Fonte de Variação	Soma dos Quadrados	Graus de Liberdade	Quadrados Médios	F
Entre amostras	48.05	1	48.05	1.9581
Dentro das amostras	441.7	18	24.54	
Total	489.75	19	25.78	

Uma vez fixado o nível de significância α do teste, H_0 (médias iguais) deve ser rejeitado se $F(\text{observações}) > F_{\text{teste}}(1-\alpha)$

$F_{\text{obs}} = 1.96 < F_{\text{teste}} = 4.5$

Ho é aceito, as médias são iguais!!!!

Exemplo 2 (Montgomery)

Weight Percentage of Cotton	Observed Tensile Strength (lb/in ²)					Totals $y_{i.}$	Averages $\bar{y}_{i.}$
	1	2	3	4	5		
15	7	7	15	11	9	49	9.8
20	12	17	12	18	18	77	15.4
25	14	18	18	19	19	88	17.6
30	19	25	22	19	23	108	21.6
35	7	10	11	15	11	54	10.8
						$y_{..} = 376$	$\bar{y}_{..} = 15.04$

We will use the analysis of variance to test $H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5$ against the alternative H_1 : some means are different. The sums of squares required are computed as follows:

Weight Percentage of Cotton	Observed Tensile Strength (lb/in ²)					Totals y_i	Averages \bar{y}_i
	1	2	3	4	5		
15	7	7	15	11	9	49	9.8
20	12	17	12	18	18	77	15.4
25	14	18	18	19	19	88	17.6
30	19	25	22	19	23	108	21.6
35	7	10	11	15	11	54	10.8
						$y_{..} = 376$	$\bar{y}_{..} = 15.04$

Table 3-4 Analysis of Variance for the Tensile Strength Data

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	F_0	P -Value
Cotton weight percentage	475.76	4	118.94	$F_0 = 14.76$	<0.01
Error	161.20	20	8.06		
Total	636.96	24			

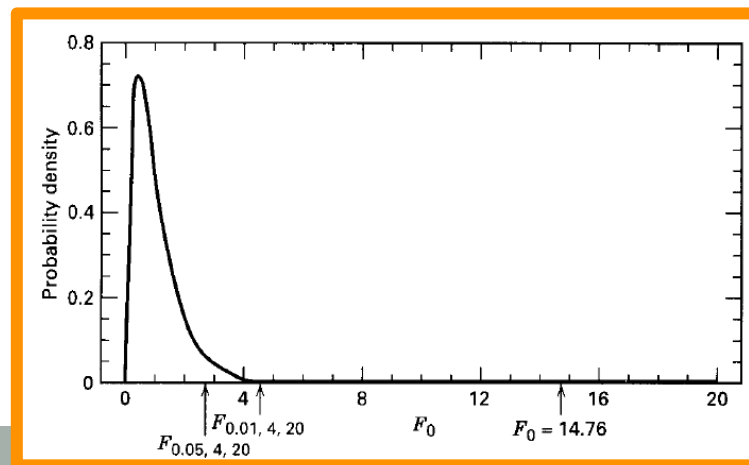
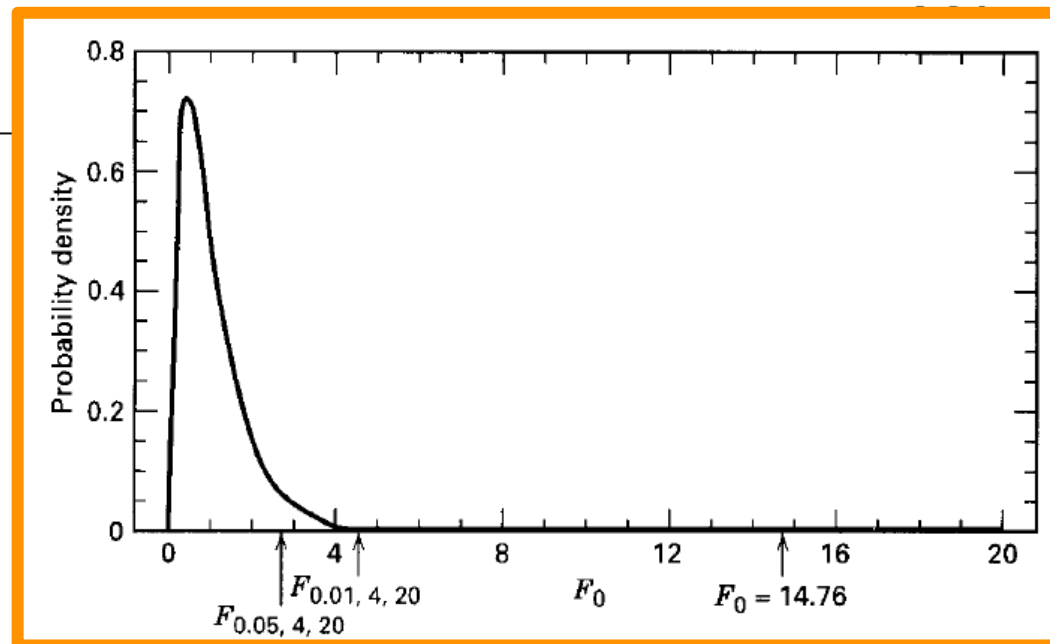


Table 3-4 Analysis of Variance for the Tensile Strength Data

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	F_0	P -Value
Cotton weight percentage	475.76	4	118.94	$F_0 = 14.76$	<0.01
Error					
Total					



The analysis of variance is summarized in Table 3-4. Note that the between-treatment mean square (118.94) is many times larger than the within-treatment or error mean square (8.06). This indicates that it is unlikely that the treatment means are equal. More formally, we can compute the F ratio $F_0 = 118.94/8.06 = 14.76$ and compare this to an appropriate upper-tail percentage point of the $F_{4,20}$ distribution. Suppose that the experimenter has selected $\alpha = 0.05$. From Appendix Table IV we find that $F_{0.05,4,20} = 2.87$. Because $F_0 = 14.76 > 2.87$, we reject H_0 and conclude that the treatment means differ; that is,

Relação entre ANOVA e ajuste de Modelo

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i, \quad i = 1, \dots, N,$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- Cap 1 – Wu e Hamada

The unknown parameters in the model are the regression coefficients $\boldsymbol{\beta}$ and the error variance σ^2 . Thus, the purpose for collecting the data is to estimate and make inferences about these parameters. For estimating $\boldsymbol{\beta}$, the least squares criterion is used; i.e., the least squares estimators (LSEs), denoted by $\hat{\boldsymbol{\beta}}$, minimize the following quantity:

$$\sum_{i=1}^N (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}))^2 \quad (1.5)$$

which in matrix notation is

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}). \quad (1.6)$$

The solution to this equation is the **least squares estimate** which is

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (1.10)$$

In fitting the model, one wants to know if any of the variables (regressors, predictors, covariates) has explanatory power. None of them has explanatory power if the null hypothesis

$$H_0: \beta_1 = \dots = \beta_p = 0 \quad (1.11)$$

In fitting the model, one wants to know if any of the variables (regressors, predictors, covariates) has explanatory power. None of them has explanatory power if the null hypothesis

$$H_0 : \beta_1 = \dots = \beta_p = 0 \quad (1.11)$$

holds. In order to test this null hypothesis, one needs to assess how much of the total variation in the response data can be explained by the model relative to the remaining variation after fitting the model, which is contained in the residuals.

Table 1.1 ANOVA Table for General Linear Model

Source	Degrees of Freedom	Sum of Squares	Mean Squares
regression	p	$\hat{\beta}^T X^T X \hat{\beta} - N\bar{y}^2$	$(\hat{\beta}^T X^T X \hat{\beta} - N\bar{y}^2)/p$
residual	$N - p - 1$	$(y - X\hat{\beta})^T (y - X\hat{\beta})$	$(y - X\hat{\beta})^T (y - X\hat{\beta}) / (N - p - 1)$
total (corrected)	$N - 1$	$y^T y - N\bar{y}^2$	

If the null hypothesis (1.11) holds, the F statistic

$$\frac{(\hat{\beta}^T \mathbf{X}^T \mathbf{X} \hat{\beta} - N\bar{y}^2)/p}{(\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta})/(N-p-1)} \quad (1.15)$$

(the regression mean square divided by the residual mean square) has an F distribution with parameters p and $N-p-1$, which are the degrees of freedom of its numerator and denominator, respectively. The p value is calculated by evaluating

$$Prob(F_{p, N-p-1} > F_{obs}), \quad (1.16)$$

where $Prob(\cdot)$ denotes the probability of an event, $F_{p, N-p-1}$ has an F distribution with parameters p and $N-p-1$, and F_{obs} is the observed value of the F statistic. The F critical values can be found in Appendix D. The p value in (1.16) can be obtained from certain pocket calculators or by interpolating the values given in Appendix D. An example of an F distribution is given in Figure 1.4 along with its critical values.

Note that the **p value** gives the probability under the null hypothesis that the F statistic value for an experiment conducted in comparable conditions will exceed the observed value F_{obs} . The smaller the p value, the stronger is the evidence that the null hypothesis does not hold. Therefore it provides a quantitative measure of the significance of effects in the experiment under study. The same interpretation can be applied when other test statistics and null hypotheses are considered.

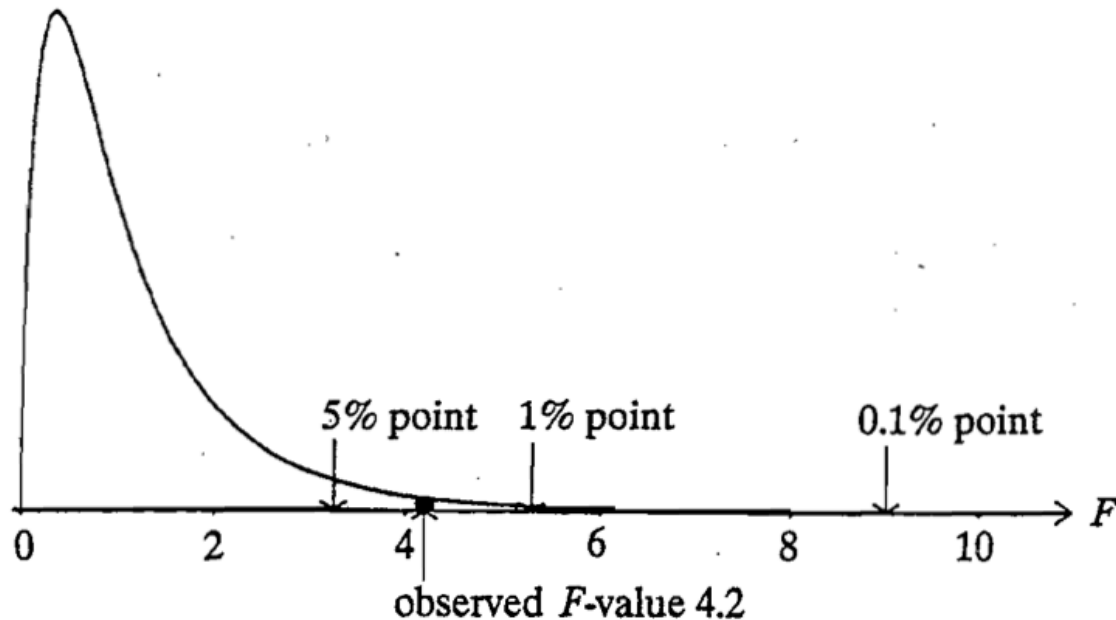


Figure 1.4. Observed F Value of 4.20 in Relation to an F Distribution With 3 and 16