



Lisbon School  
of Economics  
& Management  
Universidade de Lisboa

# Estatística II

Licenciatura em Gestão do Desporto  
2.º Ano/2.º Semestre  
2023/2024

# Aulas Teórico-Práticas N.ºs 18 a 20 (Semanas 11 e 12)

**Docente:** Elisabete Fernandes

**E-mail:** efernandes@iseg.ulisboa.pt



<https://doity.com.br/estatistica-aplicada-a-nutricao>



<https://basiccode.com.br/produto/informatica-basica/>

# Conteúdos Programáticos

## Aulas Teórico-Práticas (Semanas 1 a 5)

- **Capítulo 1:** Revisões e Distribuições de Amostragem

## Aulas Teórico-Práticas (Semanas 5 a 7)

- **Capítulo 2:** Estimação

## Aulas Teórico-Práticas (Semanas 7 a 9)

- **Capítulo 3:** Testes de Hipóteses

## Aulas Teórico-Práticas (Semanas 10 a 13)

- **Capítulo 4:** Modelo de Regressão Linear Múltipla

**Material didático:** Exercícios do Livro Murteira et al (2015), Formulário e Tabelas Estatísticas

**Bibliografia:** B. Murteira, C. Silva Ribeiro, J. Andrade e Silva, C. Pimenta e F. Pimenta; *Introdução à Estatística*, 2ª ed., Escolar Editora, 2015.

<https://cas.iseg.ulisboa.pt>



# Teste de Hipóteses de Ajustamento do Q-Q

Hipóteses, Estatística de Teste e Decisão

1

4. O tabagismo é um fator de risco para neoplasia gástrica. Pretende-se saber se se pode considerar que a ocorrência de neoplasia gástrica, em fumadores, é igualmente provável na região pilórica, no corpo gástrico e na região do cárdia. Observada uma amostra aleatória constituída por 161 indivíduos com neoplasia gástrica e fumadores de 20 cigarros/dia pelo menos durante 20 anos, obteve-se a seguinte tabela de frequências:

Região pilórica	Corpo gástrico	Região do cárdia
45	54	62

☞ Realizado o teste estatístico adequado, obteve-se o *output*:

Test Statistics	
	Neoplasia gástrica
Chi-Square <sup>a</sup>	2,696
df	2
Asymp. Sig.	,260

a. 0 cells (.0%) have expected frequencies less than 5.  
The minimum expected cell frequency is .....

- 4.1. Identifique, justificando, o teste estatístico utilizado.
- 4.2. Formule as hipóteses estatísticas associadas ao teste.
- 4.3. Calcule as frequências esperadas sob a hipótese nula, e complete o output.
- 4.4. Indique o valor observado da estatística de teste e a forma como foi obtido.
- 4.5. O que pode afirmar ao nível de significância de 5%?



## Exercício 4: Variável

Localização\_neoplasia

- Localização de neoplasia gástrica (1-Região pilórica, 2-Corpo gástrico, 3-Região do cárdia) em fumadores
- Qualitativa nominal

Freq

- Frequência absoluta

# Exercícios 4.1. e 4.2: Teste de Ajustamento do Qui-Quadrado

## Hipóteses

$H_0: p_1 = p_2 = p_3 = 1/3 = 33,3\%$  (“Ocorrência de neoplasia gástrica é igualmente provável nas 3 regiões”)

*Versus*

$H_1$ : Pelo menos uma destas probabilidades regista outro valor diferente na população (ou “Ocorrência de neoplasia gástrica não é igualmente provável nas 3 regiões”)

### Dados

$n = 161$

$p_i$  = probabilidade da ocorrência de neoplasia gástrica na  $i$ -ésima região,  $i = 1, 2, 3$

### **Objetivo do Teste de Ajustamento do Qui-Quadrado:**

- ✓ Pretende-se saber se a ocorrência de neoplasia gástrica, em fumadores, é igualmente provável em três regiões do corpo (região pilórica, corpo gástrico e região do cárdia). De outra forma, pretende-se testar se a frequência dessa doença é igual nas 3 regiões do corpo.

# Exercícios 4.3. e 4.4: Teste de Ajustamento do Qui-Quadrado

Formulário

Estatística de teste

$$\sum_i \frac{(O_i - E_i)^2}{E_i} \overset{\text{apr}}{\sim} \chi^2_{(k-1)}$$

**Teste de Ajustamento:**  $Q = \sum_{j=1}^m \frac{(N_j - fe_j)^2}{fe_j} \overset{a}{\sim} \chi^2(m-1)$

Com estimação de  $k$  parâmetros para obter as estimativas  $\hat{p}_{\circ j}$ :  $\chi^2_{(m-k-1)}$

**Pela fórmula:**

$$\text{VOE} = (-8,67)^2/53,67 + 0,33^2/53,67 + 8,33^2/53,67 = 2,695$$

**Resposta:**  $E_i$  mínimo é 53,67

Região	Freq. obser. (O <sub>i</sub> )	Freq. esper. (E <sub>i</sub> = n × p <sub>i</sub> )	Resíduos (O <sub>i</sub> - E <sub>i</sub> )
1 - Região pilórica	45	161 × 1/3 = 53,67	-8,67
2 - Corpo gástrico	54	53,67	0,33
3 - Região do cárdia	62	53,67	8,33
<b>Total</b>	161		

**Condições de Aplicabilidade dos Testes do Qui-Quadrado:**

- As frequências esperadas devem ser  $\geq 5$ .
- No caso de tal não se verificar, então pelo menos 80% das frequências esperadas  $\geq 5$  e todas  $> 1$

**A Condição de Aplicabilidade dos Testes do Qui-Quadrado é satisfeita:**

- Todas as frequências esperadas  $E_i \geq 5$ .



# Exercício 4.5: Teste de Ajustamento do Qui-Quadrado

Decisão (para  $\alpha = 0,05$ )

**Pelo valor crítico:**  $\text{VOE} = 2,695 < \chi^2_{0,95;2} = 5,99$

Região de rejeição ou crítica:  
 $2,695 \notin \text{RR} = [\chi^2_{95;2}; +\infty[ = [5,99; +\infty[$

Tabela da Distribuição do Qui-Quadrado

Quantil de probabilidade  $1-\alpha$  da distribuição Qui-Quadrado

**Regra de decisão pelo valor crítico ou região de rejeição (RR):**

$\left\{ \begin{array}{l} \text{VOE} \geq \chi^2_{1-\alpha} \\ \text{VOE} \in \text{RR} = [\chi^2_{1-\alpha}; +\infty[ \end{array} \right\} \Rightarrow \text{Rejeita-se } H_0 \text{ para } \alpha$

**Regra de decisão pelo valor-p:**

Valor-p =  $P(X^2 \geq \text{VOE}) < \alpha \Rightarrow \text{Rejeita-se } H_0 \text{ para } \alpha$

**Pelo valor-p:** valor-p =  $P(X^2 \geq 2,695) = 0,260 > 0,05$

# Cálculo do Quantil da Distribuição Qui-Quadrado de Probabilidade $1-\alpha$ e com $n-1$ g.l.'s

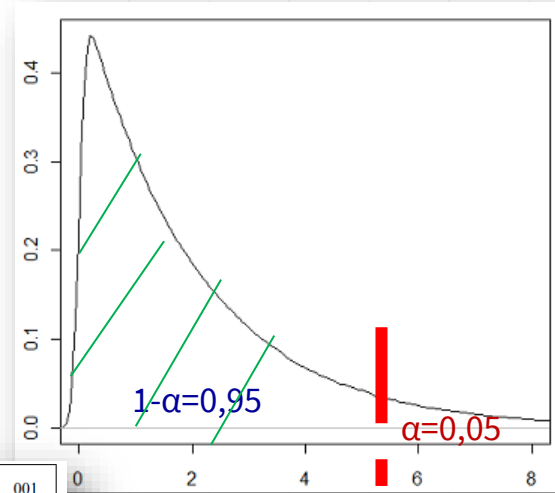
Nível de confiança ( $1-\alpha=0,95$ )

Nível de significância ( $\alpha=0,05$ )

Área total é igual a 1

O nível de significância é igual a  $\alpha = 0,05$ , então tem-se  $1-\alpha = 0,95$

$\chi^2_{0,95;2} = 5,991$  (ver tabela)



$\chi^2_{0,95;2} = 5,991$

$$\chi^2_{n,\epsilon} : P(X > \chi^2_{n,\epsilon}) = \epsilon$$

$\epsilon$	.995	.990	.975	.950	.900	.750	.500	.250	.100	.050	.025	.010	.005	.001
<b>n</b>														
1	.000	.000	.001	.004	.016	.102	.455	1.323	2.706	3.841	5.024	6.635	7.879	10.827
2	.100	.020	.051	.103	.211	.575	1.386	2.773	4.605	<b>5.991</b>	7.378	9.210	10.597	13.815
3	.072	.115	.216	.352	.584	1.213	2.366	4.108	6.251	7.879	9.348	11.345	12.838	16.266
4	.207	.297	.484	.711	1.064	1.923	3.357	5.385	7.779	9.488	11.143	13.277	14.860	18.466
5	.412	.554	.831	1.145	1.610	2.675	4.351	6.626	9.236	11.070	12.832	15.086	16.750	20.515
6	.676	.872	1.237	1.635	2.204	3.455	5.348	7.841	10.645	12.592	14.449	16.812	18.548	22.457
7	.989	1.239	1.690	2.167	2.833	4.255	6.346	9.037	12.017	14.067	16.013	18.475	20.278	24.321
8	1.344	1.647	2.180	2.733	3.490	5.071	7.344	10.219	13.362	15.507	17.535	20.090	21.955	26.124
9	1.735	2.088	2.700	3.325	4.168	5.899	8.343	11.389	14.684	16.919	19.023	21.666	23.589	27.877
10	2.156	2.558	3.247	3.940	4.865	6.737	9.342	12.549	15.987	18.307	20.483	23.209	25.188	29.588

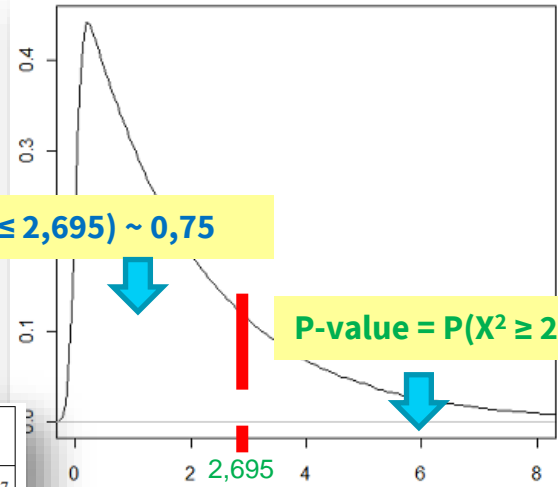
# Cálculo do Valor-p quando a Estatística de Teste tem Distribuição Qui-Quadrado

$$\text{valor-p} = P(X^2 \geq 2,695) \sim P(X^2 \geq 2,773) = 0,25$$

$$\chi_{n,\varepsilon}^2 : P(X > \chi_{n,\varepsilon}^2) = \varepsilon$$

n	ε	.995	.990	.975	.950	.900	.750	.500	.250	.100	.050	.025	.010	.005	.001
1	→	.000	.000	.001	.004	.016	.102	.455	1.323	2.706	3.841	5.024	6.635	7.879	10.827
2		.010	.020	.051	.103	.211	.575	1.386	2.773	4.605	5.991	7.378	9.210	10.597	13.815
3		.072	.115	.216	.352	.584	1.213	2.366	4.108	6.251	7.815	9.348	11.345	12.838	16.266
4		.207	.297	.484	.711	1.064	1.923	3.357	5.385	7.779	9.488	11.143	13.277	14.860	18.466
5		.412	.554	.831	1.145	1.610	2.675	4.351	6.626	9.236	11.070	12.832	15.086	16.750	20.515
6		.676	.872	1.237	1.635	2.204	3.455	5.348	7.841	10.645	12.592	14.449	16.812	18.548	22.457
7		.989	1.239	1.690	2.167	2.833	4.255	6.346	9.037	12.017	14.067	16.013	18.475	20.278	24.321
8		1.344	1.647	2.180	2.733	3.490	5.071	7.344	10.219	13.362	15.507	17.535	20.090	21.955	26.124
9		1.735	2.088	2.700	3.325	4.168	5.899	8.343	11.389	14.684	16.919	19.023	21.666	23.589	27.877
10		2.156	2.558	3.247	3.940	4.865	6.737	9.342	12.549	15.987	18.307	20.483	23.209	25.188	29.588

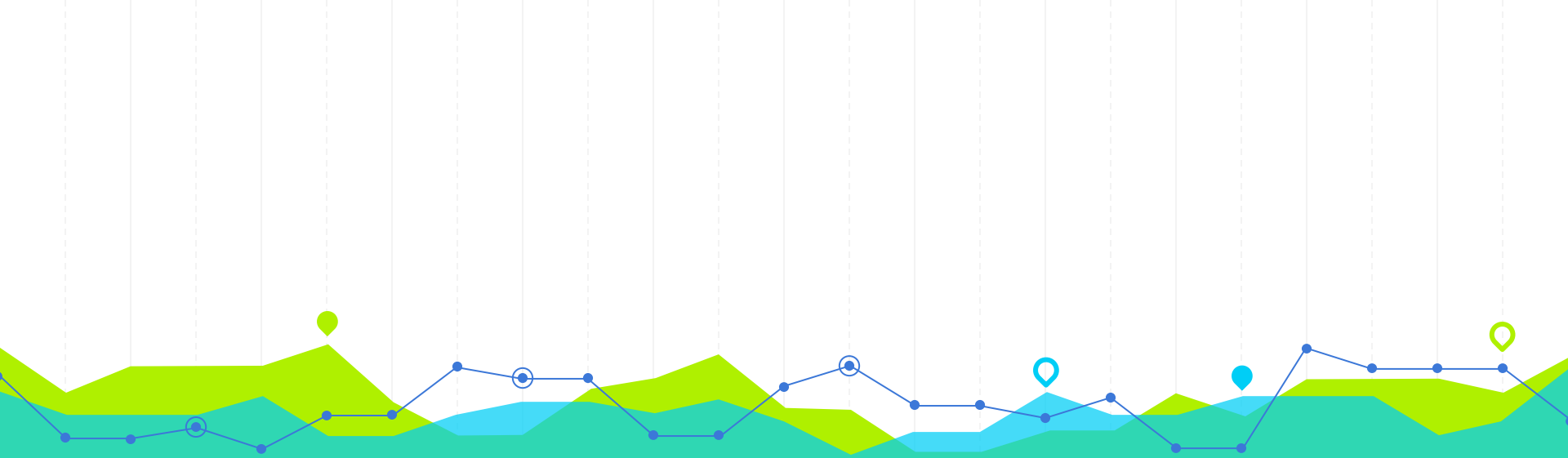
Área total é igual a 1



$P(X^2 \leq 2,695) \sim 0,75$

$P\text{-value} = P(X^2 \geq 2,695) \sim 0,25$

**Regra de decisão pelo valor-p:**  
 $\text{Valor-p} = P(X^2 \geq \text{VOE}) < \alpha \Rightarrow \text{Rejeita-se } H_0 \text{ para } \alpha$



# Teste de Hipóteses de Independência do Q-Q

Hipóteses, Estatística de Teste e Decisão

# 2

5. Um inspetor de qualidade recolheu uma amostra de 176 produtos alimentares num centro de distribuição. Sabendo que cada produto pode ser proveniente de uma de três fábricas e pode ou não estar contaminado; o inspetor avaliou todos os produtos e obteve os seguintes resultados:

	<b>Fábrica A</b>	<b>Fábrica B</b>	<b>Fábrica C</b>	<b>Total</b>
<b>Contaminado</b>	8	15	11	34
<b>Não contaminado</b>	55	67	20	142
<b>Total</b>	63	82	31	176

Pode-se afirmar que o facto de um produto estar contaminado é independente da sua fábrica de origem, considerando  $\alpha = 0,01$ ?

[Adaptado da fonte: <https://www.ime.unicamp.br/~veronica/Coordenadas1s/aula8pr.pdf>]



$$\text{Teste de Independência: } Q = \sum_{i=1}^r \sum_{j=1}^s \frac{(N_{ij} - fe_{ij})^2}{fe_{ij}} \sim \chi_{((r-1)(s-1))}^2$$

## Exercício: Teste de Independência do Qui-Quadrado

### Hipóteses

$H_0$ : As variáveis são independentes

*Versus*

$H_1$ : As variáveis não são independentes

### Estatística de teste

$$\sum_i \sum_j \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \sim \chi_{(1-1)(c-1)}^2$$

### Decisão

**Pelo valor crítico:** Valor da Estatística de Teste = 7,024 não pertence a  $RR = [\chi_{0,99;2}^2; +\infty[ = [9,21; +\infty[$

**Pelo valor-p:** valor-p = 0,030 > 0,01

Não se rejeita  $H_0$  para  $\alpha = 1\%$ . Assim, não existe evidência estatística para afirmar que as variáveis não são independentes para  $\alpha = 1\%$ .

### Dados

N = 176

VOE = 7,024

Valor-p = 0,03

$\alpha = 1\% \Rightarrow 1 - \alpha = 99\%$

df = g.l.'s = 2

Frequências esperadas

$$E_{ij} = \frac{L_i \times C_j}{N}$$

L = total linhas e C = total colunas  
N = nº total de elementos

Tabela da distribuição do Qui-Quadrado

# Exercício: Teste de Independência do Qui-Quadrado

Estatística de Teste:

$$\sum_i \sum_j \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \sim \chi^2_{(1-1)(c-1)}$$

Frequências esperadas

$$E_{ij} = \frac{L_i \times C_j}{N}$$

L = total linhas e C = total colunas  
N = n° total de elementos

Contaminado \* Fábrica Tabulação cruzada

			Fábrica			Total
			Fábrica A	Fábrica B	Fábrica C	
Contaminado	Sim	Contagem	8	15	11	34
		Expected Count	12,2	15,8	6,0	34,0
	Não	Contagem	55	67	20	142
		Expected Count	50,8	66,2	25,0	142,0
Total		Contagem	63	82	31	176
		Expected Count	63,0	82,0	31,0	176,0

Valor da Estatística de Teste

Testes de chi-quadrado

	Valor	df	Sig. Assint. (2 lados)
Chi-quadrado de Pearson	7,024	2	,030
Razão de probabilidade	6,450	2	,040
Associação Linear por Linear	6,099	1	,014
N de Casos Válidos	176		

A Condição de Aplicabilidade dos Testes do Qui-Quadrado é satisfeita:

- Todas as frequências esperadas  $E_{ij} \geq 5$ .

a. 0 células (0,0%) esperam contagem menor do que 5. A contagem mínima esperada é 5,000.

# Teste de Independência do Qui-Quadrado

Decisão (para  $\alpha = 0,01$ )

**Pelo valor crítico:**  $VOE = 7,024 > \chi^2_{0,99;2} = 9,21$

Região de rejeição ou crítica:

$7,024$  não pertence a  $RR = [\chi^2_{0,99;2}; +\infty[ = [9,21; +\infty[$

Tabela da Distribuição do Qui-Quadrado

**Regra de decisão pelo valor crítico ou região de rejeição (RR):**

$\{ VOE \geq \chi^2_{1-\alpha}$   
 $(VOE \in RR = [\chi^2_{1-\alpha}; +\infty[ \Rightarrow \text{Rejeita-se } H_0 \text{ para } \alpha$

Quantil de probabilidade  $1-\alpha$  da distribuição do Qui-Quadrado

**Pelo valor-p:** valor-p =  $0,03 > 0,01$

**Regra de decisão pelo valor-p:**

Valor-p =  $P(X^2 \geq VOE) < \alpha \Rightarrow \text{Rejeita-se } H_0 \text{ para } \alpha$

Não se rejeita-se  $H_0$  para  $\alpha = 0,01$ . Assim, não existe evidência estatística para afirmar que as variáveis não são independentes para  $\alpha = 1\%$ .



# Cálculo do Quantil da Distribuição Qui-Quadrado de Probabilidade $1-\alpha$ e com $(l-1) \times (c-1)$ g.l.'s

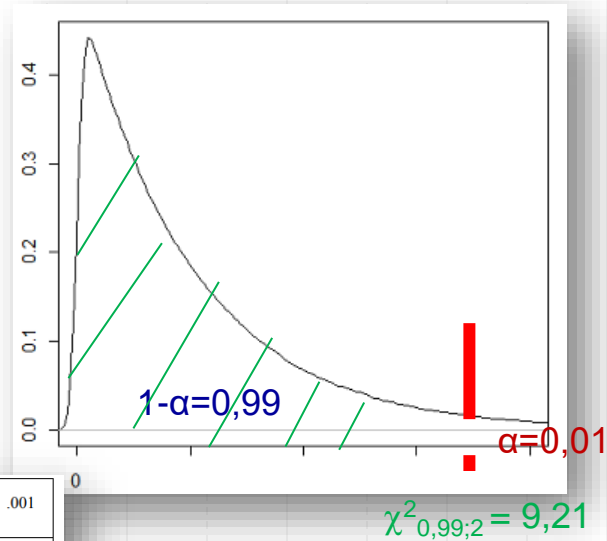
Nível de confiança ( $1-\alpha=0,99$ )

Nível de significância ( $\alpha=0,01$ )

Área total é igual a 1

O nível de significância é igual a  $\alpha = 0,01$ , então tem-se  $1-\alpha = 0,99$

$\chi^2_{0,99;2} = 9,21$  (ver tabela)



$$\chi^2_{n,\varepsilon} : P(X > \chi^2_{n,\varepsilon}) = \varepsilon$$

$\varepsilon$	.995	.990	.975	.950	.900	.750	.500	.250	.100	.050	.025	.010	.005	.001
<b>n</b>														
<b>1</b>	.000	.000	.001	.004	.016	.102	.455	1.323	2.706	3.841	5.024	6.635	7.879	10.827
<b>2</b>	.010	.020	.051	.103	.211	.575	1.386	2.773	4.605	5.991	7.378	9.210	10.597	13.815
<b>3</b>	.072	.115	.216	.352	.584	1.213	2.366	4.108	6.251	7.815	9.348	11.345	12.838	16.266
<b>4</b>	.207	.297	.484	.711	1.064	1.923	3.357	5.385	7.779	9.488	11.143	13.277	14.860	18.466
<b>5</b>	.412	.554	.831	1.145	1.610	2.675	4.351	6.626	9.236	11.070	12.832	15.086	16.750	20.515
<b>6</b>	.676	.872	1.237	1.635	2.204	3.455	5.348	7.841	10.645	12.592	14.449	16.812	18.548	22.457
<b>7</b>	.989	1.239	1.690	2.167	2.833	4.255	6.346	9.037	12.017	14.067	16.013	18.475	20.278	24.321
<b>8</b>	1.344	1.647	2.180	2.733	3.490	5.071	7.344	10.219	13.362	15.507	17.535	20.090	21.955	26.124
<b>9</b>	1.735	2.088	2.700	3.325	4.168	5.899	8.343	11.389	14.684	16.919	19.023	21.666	23.589	27.877
<b>10</b>	2.156	2.558	3.247	3.940	4.865	6.737	9.342	12.549	15.987	18.307	20.483	23.209	25.188	29.588

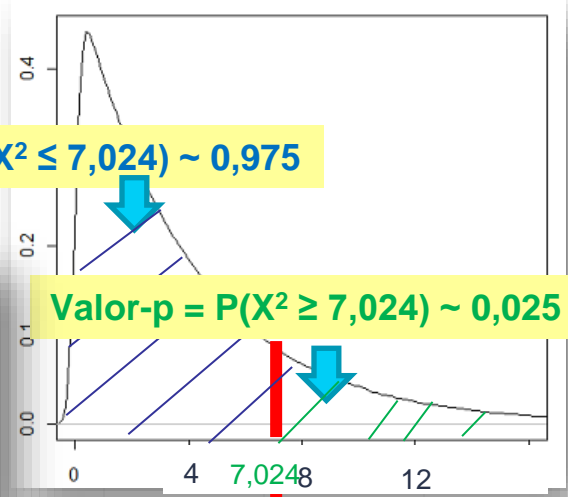
**Regra de decisão pelo valor-p:**  
 Valor-p =  $P(X^2 \geq \text{VOE}) < \alpha \Rightarrow$  Rejeita-se  $H_0$  para  $\alpha$

# Cálculo do Valor-p quando a Estatística de Teste tem Distribuição Qui-Quadrado

valor-p =  $P(X^2 \geq 7,024) \sim P(X^2 \geq 7,378) = 0,025$

$\chi^2_{n,\epsilon} : P(X > \chi^2_{n,\epsilon}) = \epsilon$

Área total é igual a 1



$P(X^2 \leq 7,024) \sim 0,975$

**Valor-p =  $P(X^2 \geq 7,024) \sim 0,025$**

$\epsilon$	.995	.990	.975	.950	.900	.750	.500	.250	.100	.050	.025	.010	.005	.001
1	.000	.000	.001	.004	.016	.102	.455	1.323	2.706	3.841	5.024	6.635	7.879	10.827
2	.010	.020	.051	.103	.211	.575	1.386	2.773	4.605	5.991	7.378	9.210	10.597	13.815
3	.072	.115	.216	.352	.584	1.213	2.366	4.108	6.251	7.815	9.348	11.345	12.838	16.266
4	.207	.297	.484	.711	1.064	1.923	3.357	5.385	7.779	9.488	11.143	13.277	14.860	18.466
5	.412	.554	.831	1.145	1.610	2.675	4.351	6.626	9.236	11.070	12.832	15.086	16.750	20.515
6	.676	.872	1.237	1.635	2.204	3.455	5.348	7.841	10.645	12.592	14.449	16.812	18.548	22.457
7	.989	1.239	1.690	2.167	2.833	4.255	6.346	9.037	12.017	14.067	16.013	18.475	20.278	24.321
8	1.344	1.647	2.180	2.733	3.490	5.071	7.344	10.219	13.362	15.507	17.535	20.090	21.955	26.124
9	1.735	2.088	2.700	3.325	4.168	5.899	8.343	11.389	14.684	16.919	19.023	21.666	23.589	27.877
10	2.156	2.558	3.247	3.940	4.865	6.737	9.342	12.549	15.987	18.307	20.483	23.209	25.188	29.588

# Exercício: Teste de Independência do Qui-Quadrado

## Condições de Aplicabilidade dos Testes do Qui-Quadrado:

- As frequências esperadas devem ser  $\geq 5$ .
- No caso de tal não se verificar, então pelo menos 80% das frequências esperadas  $\geq 5$  e todas  $> 1$  (não é válido para tabelas 2x2).

## Verificação das condições de aplicabilidade

Neste caso, todas as células têm frequências esperadas superiores a 5.

O teste do Qui-Quadrado apenas informa sobre a independência entre variáveis, mas nada diz sobre o grau de associação existente.

Para esse efeito calculam-se **medidas de associação** tais como o coeficiente Phi, o coeficiente V de Cramer e o coeficiente de contingência.

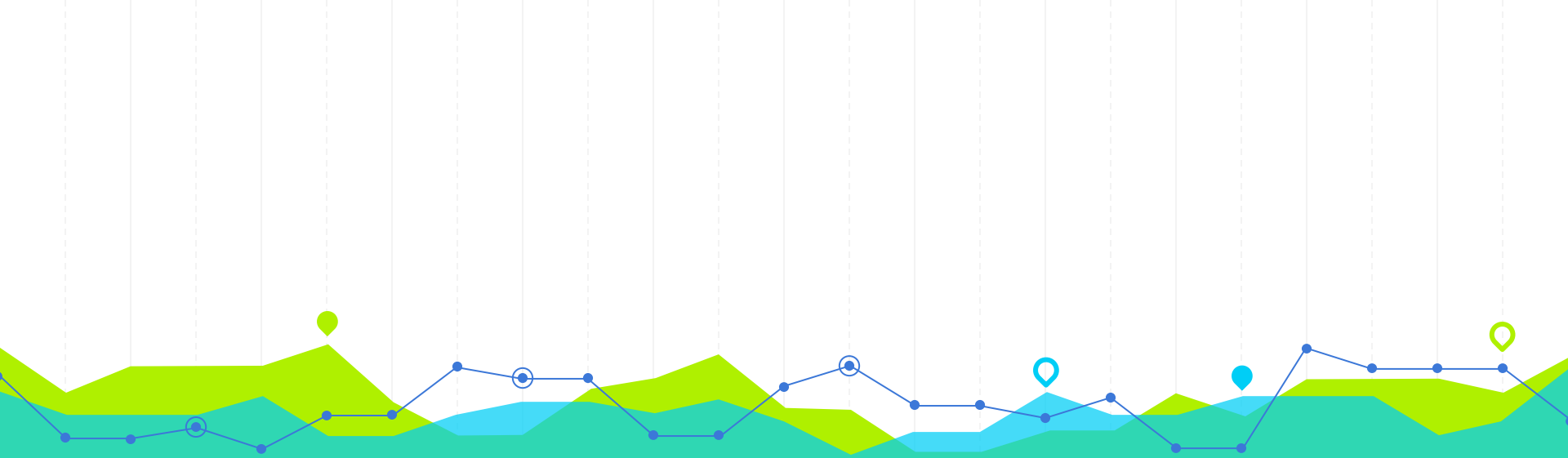
Frequências esperadas

$$E_{ij} = \frac{L_i \times C_j}{N}$$

L = total linhas e C = total colunas  
N = n° total de elementos

Contaminado \* Fábrica Tabulação cruzada

		Fábrica			Total	
		Fábrica A	Fábrica B	Fábrica C		
Contaminado	Sim	Contagem	8	15	11	34
		Expected Count	12,2	15,8	6,0	34,0
	Não	Contagem	55	67	20	142
		Expected Count	50,8	66,2	25,0	142,0
Total	Contagem	63	82	31	176	
	Expected Count	63,0	82,0	31,0	176,0	



# Teste de Hipóteses Não Paramétricos

Hipóteses, Estatística de Teste e Decisão  
Murteira et al (2015)

# 3

1. Com o objectivo de remodelar determinado centro comercial, realizou-se uma pesquisa sobre o movimento de entradas e saídas pelas suas três portas. Inquiriu-se qual a porta de entrada para uma amostra aleatória de 201 pessoas:

Entrada	1	2	3
N.º de pessoas	83	62	56

Foi afirmado que não havia preferência por qualquer uma das três entradas. Comente para uma dimensão de 0.05.



## Exercício 1

$X$  = Porta escolhida para entrar no centro comercial

$D_X = \{1, 2, 3\}$        $n = 201$  pessoas

$P_j = P(X = j)$  ( $j = 1, 2, 3$ )

$H_0: P_1 = P_2 = P_3 = \frac{1}{3}$        $H_1: \exists j: P_j \neq \frac{1}{3}$  ( $j = 1, 2, 3$ )

$P_{0j} = P(X = j | H_0) = \frac{1}{3}$  ( $j = 1, 2, 3$ )

# Exercício 1

$$Q = \sum_{j=1}^m \frac{(N_j - fe_j)^2}{fe_j} \approx \chi^2(m-1) = \chi^2(2)$$

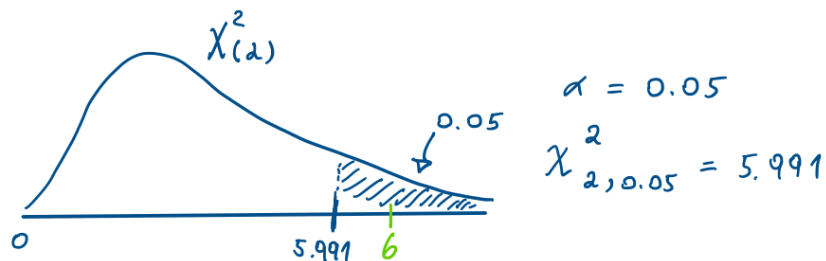
$m = m^\circ$  de categorias

Porta  $N_j$   $fe_j = m P_{0j}$

1	83	$201 \times \frac{1}{3} = 67$	$> 5$	} $m = 3$
2	62	$201 \times \frac{1}{3} = 67$	$> 5$	
3	56	$201 \times \frac{1}{3} = 67$	$> 5$	

$$Q_{obs} = \frac{(83 - 67)^2}{67} + \frac{(62 - 67)^2}{67} + \frac{(56 - 67)^2}{67} = 6$$

# Exercício 1



$$W_{\alpha} = \{Q_{obs} : Q_{obs} > \chi^2_{2,0.05} = 5.991\}$$

$$Q_{obs} = 6 \in W_{\alpha} \text{ logo rejeita-se } H_0.$$

Conclusão : Tendo como referência um teste de dimensão 0.05 a hipótese nula é rejeitada, o que significa que a evidência estatística presente na amostra recolhida não é favorável à afirmação enunciada.

Nota : `> pchisq(6, 2, lower.tail = FALSE)`  
[1] 0.04978707

$$p_{obs} = 0.0498 < \alpha = 0.05 \text{ logo rejeita-se } H_0.$$

Com  $\alpha = 0.01$  já não se rejeitaria  $H_0$ .



4. O número de erros de impressão por página de um livro é frequentemente considerado uma variável aleatória de Poisson. A contagem dos erros de impressão em 100 páginas de um livro deu o seguinte resultado:

N.º de erros	0	1	2	3
N.º de páginas	65	25	8	2

- a) Teste a hipótese, ao nível de 0.01, de que o número de erros por página é uma variável aleatória de Poisson de parâmetro  $\lambda = 0.4$ .



## Exercício 4

a)  $X \equiv$  nº de erros por página  
amostra:  $n = 100$  páginas

$H_0: X \sim \text{Poi}(0.4)$      $H_1: X \neq \text{Poi}(0.4)$     ( $\alpha = 0.01$ )

erros     $N_j$      $n p_{0j}$     sob  $H_0$ ,  $f_x(x) = \frac{e^{-0.4} 0.4^x}{x!}$     ( $x = 0, 1, 2, \dots$ )

0    65    67.03

1    25    26.81

2    8    5.36

$\geq 3$     2    0.8

$\left. \begin{array}{l} 0.8 < 5 \end{array} \right\} \rightarrow$  Temos que agrupar estas duas categorias:

$$m = 3 \quad N_3 = 10$$

$$m p_{03} = 5.36 + 0.8 = 6.16$$

## Exercício 4

$$P_{0j} = P(X = j-1 | H_0) = \frac{e^{-0.4} 0.4^{j-1}}{(j-1)!} \quad (j = 1, 2, 3)$$

$$P_{04} = P(X \geq 3 | H_0) = 1 - P_{01} - P_{02} - P_{03}$$

Assim sendo:

$$n P_{01} = 100 P(X = 0 | H_0) = 100 \times \underbrace{0.6703}_{\text{Tabela 2}} = 67.03$$

$$n P_{02} = 100 P(X = 1 | H_0) = 100 \times \underbrace{0.2681} = 26.81$$

$$n P_{03} = 100 P(X = 2 | H_0) = 100 \times \underbrace{0.0536} = 5.36$$

$$\begin{aligned} n P_{04} &= 100 P(X \geq 3 | H_0) = 100 \times (1 - P_{01} - P_{02} - P_{03}) = \\ &= 100 (1 - 0.6703 - 0.2681 - 0.0536) = 100 \times 0.008 = 0.8 \end{aligned}$$

## Exercício 4

A Hipótese que realmente se testa é a hipótese "afarentada":

$$H'_0: \begin{cases} P(X = x) = \frac{e^{-0.4} 0.4^x}{x!} \quad (x = 0, 1) \\ P(X \geq 2) = 1 - P(X=0) - P(X=1) \end{cases}$$

$$Q = \sum_{j=1}^3 \frac{(N_j - m p_{0j})^2}{m p_{0j}} \approx \chi^2_{(m-1)} = \chi^2_{(2)}$$

$$\chi^2_{2, 0.01} = 9.21$$

$$W_Q = \{Q_{obs} : Q_{obs} > 9.21\}$$

## Exercício 4

$$Q_{obs} = \frac{(65 - 67.03)^2}{67.03} + \frac{(25 - 26.81)^2}{26.81} + \frac{(10 - 6.16)^2}{6.16} =$$
$$= 2.58 \notin W_{\alpha} \text{ logo não se rejeita } H_0.$$

16. Uma companhia de seguros está interessada em saber se existe independência entre a realização de seguros de dois tipos: “ Seguro de Vida” e “Seguro de Saúde”. Para tal analisou as apólices detidas por uma amostra casual de 200 clientes tendo-se verificado que: 90 tinham simultaneamente um seguro de vida e um seguro de saúde; 35 tinham apenas seguro de vida; 42 tinham apenas seguro de saúde; os restantes não tinham qualquer destes dois seguros. Com base num teste a 0.05, que pode concluir?



## Exercício 16

$$H_0: P_{ij} = P_{i\cdot} P_{\cdot j} \quad (i, j = 1, 2)$$

$$H_1: \exists (i, j): P_{ij} \neq P_{i\cdot} P_{\cdot j}$$

$$Q = \frac{\sum_{i=1}^2 \sum_{j=1}^2 (N_{ij} - f_{ij}^e)^2}{f_{ij}^e} \approx \chi^2_{\{(2-1)(2-1)\}} = \chi^2(1)$$

# linhas      # colunas

max 1 132      00 1 a ~ ~

$$\hat{P}_{\cdot 1} = \frac{132}{200} = 0.66 \quad \frac{68}{200} = 0.34 = P_{\cdot 2}$$

## Exercício 16

$$H_0: P_{ij} = P_{i \cdot} P_{\cdot j} \quad (i, j = 1, 2)$$

$$H_1: \exists (i, j): P_{ij} \neq P_{i \cdot} P_{\cdot j}$$

$$Q = \frac{\sum_{i=1}^2 \sum_{j=1}^2 (N_{ij} - f_{e_{ij}})^2}{f_{e_{ij}}} \approx \chi^2_{\{(\overset{\# \text{ linhas}}{2-1}) (\overset{\# \text{ colunas}}{2-1})\}} = \chi^2(1)$$



## Exercício 16

Frequências esperadas sob  $H_0$  ( $f_{e_{ij}}$ ):

$$f_{e_{ij}} = n \hat{P}_{i\cdot} \hat{P}_{\cdot j} \quad (i, j = 1, 2)$$

	<u>seguro saúde</u>	
	sim	não
<u>seguro</u>	$200 \times 0.625 \times 0.66 = 82.5$	$200 \times 0.625 \times 0.34 = 42.5$
<u>vida</u>	$200 \times 0.375 \times 0.66 = 49.5$	$200 \times 0.375 \times 0.34 = 25.5$

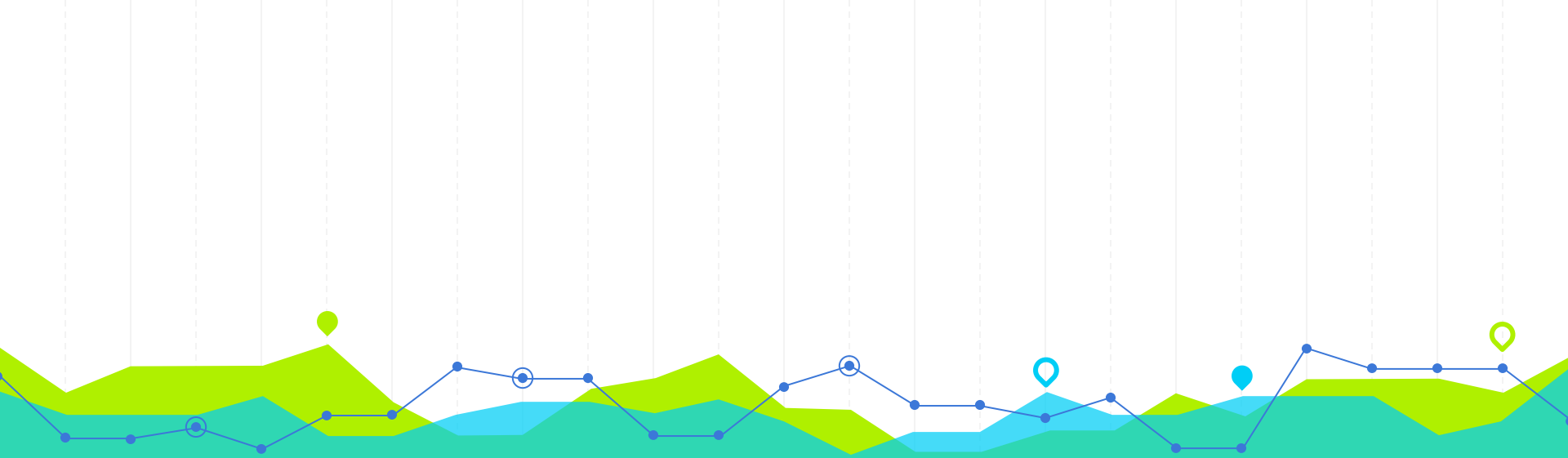
## Exercício 16

$$\chi_{1,0.05}^2 = 3.841$$

$$W_Q = \{Q_{obs} : Q_{obs} > 3.841\}$$

$$Q_{obs} = \frac{(90-82.5)^2}{82.5} + \frac{(35-42.5)^2}{42.5} + \frac{(42-49.5)^2}{49.5} + \frac{(33-25.5)^2}{25.5} =$$

$$= 5.3476 \in W_Q \rightarrow \text{Rejeitamos } H_0.$$



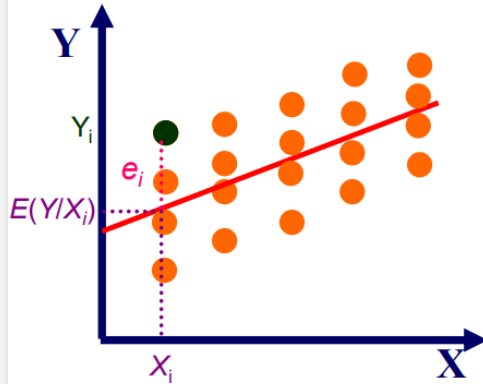
# Modelo de Regressão Linear Múltipla

Estimação dos Coeficientes da Reta de Regressão

# 4

# Modelo de Regressão Linear Simples: Revisão

Seja a relação entre  $Y$  e  $X$  na população:



$$Y_i = \alpha + \beta X_i + e_i$$

ou

$$E(Y/X_i) = \alpha + \beta X_i$$

Modelo de Regressão Linear Simples para  $Y$  na população

Onde:

$Y$  é a variável dependente ou regressando  
 $X$  é a variável independente ou regressor  
 $\alpha$  é o intercepto ou constante do modelo  
 $\beta$  é o coeficiente angular do modelo

Beta é o declive

Alfa é a ordenada na origem

**Erro de previsão:**

Seja  $X_i$  a  $i$ -ésima observação de  $X$ , teremos:

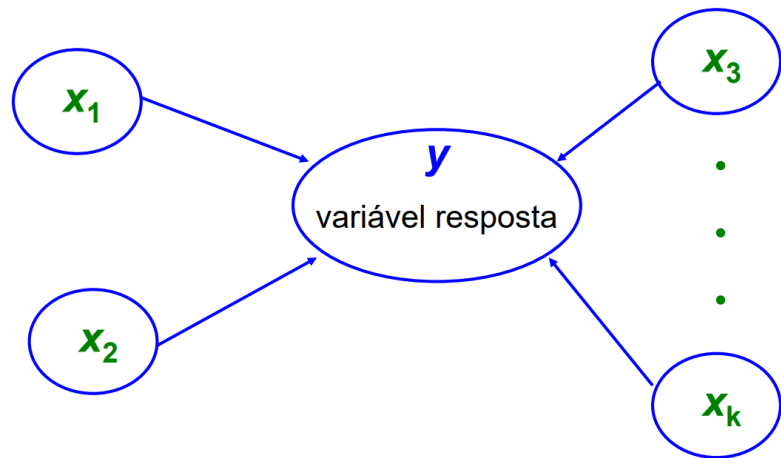
$Y_i$  é o valor observado em  $Y$  para o  $i$ -ésimo valor de  $X$

$E(Y/X_i)$  é a esperança condicional de  $Y$  e representa o valor esperado de  $Y$  para o  $i$ -ésimo valor de  $X$

$e_i$  é o erro, ou variação de  $Y_i$  não explicada pelo modelo

# Modelo de Regressão Linear Múltipla (MRLM)

Chamamos Modelo de Regressão Linear Múltipla a qualquer modelo de regressão linear com duas ou mais variáveis explicativas.



$x_1, x_2, \dots, x_k$ : variáveis explicativas (regressores)

# MRLM

Vamos admitir que  $X_1, X_2, \dots, X_k$  sejam as variáveis independentes e  $Y$  a variável dependente.

Dada uma amostra de  $n$  observações,

$$(x_{1i}, x_{2i}, \dots, x_{ki}, y_i), \quad i = 1, 2, \dots, n,$$

# MRLM

o modelo de regressão linear múltipla será dado por:

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \varepsilon_i ,$$

ou

$$E[y_i | x_{1i}, x_{2i}, \dots, x_{ki}] = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} ,$$

$i = 1, 2, \dots, n$

em que  $n > (k+1)$ .

Neste modelo,  $k$  é o  $n^\circ$  de variáveis independentes e  $k+1$  é o  $n^\circ$  de coeficientes

# Estimação dos Coeficientes do MRLM: Método dos Mínimos Quadrados (MMQ)

Para determinarmos os estimadores de mínimos quadrados de  $\beta_0, \beta_1, \dots, \beta_k$ , devemos minimizar o erro quadrático total ( $\sum \varepsilon_i^2$ ):

$$\sum_{i=1}^n (\varepsilon_i)^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_k x_{ki})^2$$

[Index of /wp-content/uploads/2014/02](http://wp-content/uploads/2014/02) (hedibert.org)



# Estimação dos Coeficientes do MRLM: MMQ

## O mínimo da função

$$S(\beta_0, \beta_1, \dots, \beta_k) = \sum_{i=1}^n (\varepsilon_i)^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_k x_{ki})^2$$

é obtido derivando-a em relação a  $\beta_0, \beta_1, \dots, \beta_k$ , e igualando o resultado a zero. Ou seja,

$$\frac{\partial}{\partial \beta_0} S(\beta_0, \beta_1, \dots, \beta_k) = 0 \quad \dots \quad \frac{\partial}{\partial \beta_k} S(\beta_0, \beta_1, \dots, \beta_k) = 0$$

## Estimação dos Coeficientes do MRLM: MMQ

$$\frac{\partial}{\partial \beta_0} S(\beta_0, \beta_1, \dots, \beta_k) = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \dots - \hat{\beta}_k x_{ki}) = 0$$

$$\frac{\partial}{\partial \beta_1} S(\beta_0, \beta_1, \dots, \beta_k) = -2 \sum_{i=1}^n [(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \dots - \hat{\beta}_k x_{ki}) x_{1i}] = 0$$

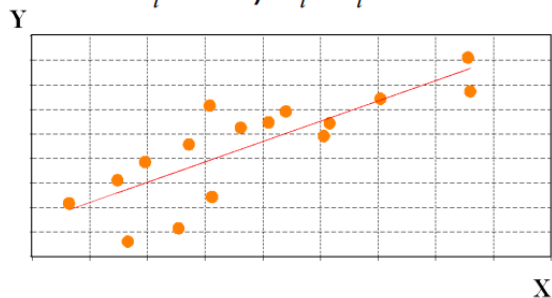
⋮

$$\frac{\partial}{\partial \beta_k} S(\beta_0, \beta_1, \dots, \beta_k) = -2 \sum_{i=1}^n [(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \dots - \hat{\beta}_k x_{ki}) x_{ki}] = 0$$

# Estimação dos Coeficientes do MRLM: MMQ

## Regressão Linear Simples:

$$Y_i = \alpha + \beta X_i + e_i$$



Onde:

$$EQT(\hat{\alpha}, \hat{\beta}) = \sum \hat{e}_i^2 = \sum [Y_i - (\hat{\alpha} + \hat{\beta}X_i)]^2$$

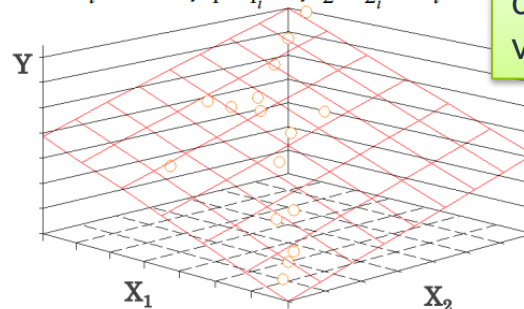
Minimizando EQT:

$$\frac{\partial EQT}{\partial \hat{\alpha}} = 0 \Rightarrow \hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

$$\frac{\partial EQT}{\partial \hat{\beta}} = 0 \Rightarrow \hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

## Regressão Linear Múltipla:

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + e_i$$



Onde:

$$EQT(\hat{\alpha}, \hat{\beta}_1, \hat{\beta}_2) = \sum \hat{e}_i^2 = \sum [Y_i - (\hat{\alpha} + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i})]^2$$

Minimizando EQT:

$$\frac{\partial EQT}{\partial \hat{\alpha}} = 0 \Rightarrow \hat{\alpha} = \bar{Y} - \hat{\beta}_1 \bar{X}_1 - \hat{\beta}_2 \bar{X}_2$$

$$\frac{\partial EQT}{\partial \hat{\beta}_1} = 0 \Rightarrow \hat{\beta}_1 = \frac{(\sum y_i x_{1i})(\sum x_{2i}^2) - (\sum y_i x_{2i})(\sum x_{1i} x_{2i})}{(\sum x_{1i}^2)(\sum x_{2i}^2) - (\sum x_{1i} x_{2i})^2}$$

$$\frac{\partial EQT}{\partial \hat{\beta}_2} = 0 \Rightarrow \hat{\beta}_2 = \frac{(\sum y_i x_{2i})(\sum x_{1i}^2) - (\sum y_i x_{1i})(\sum x_{1i} x_{2i})}{(\sum x_{1i}^2)(\sum x_{2i}^2) - (\sum x_{1i} x_{2i})^2}$$

Caso particular: MRLM com apenas duas variáveis regressoras

# MRLM: Abordagem Matricial

Devido à complexidade das fórmulas envolvidas, utilizaremos a abordagem matricial, que nos permitirá, entre outras coisas:

- i. encontrar o vetor de estimadores;
- ii. verificar as propriedades estatísticas de (i);
- iii. obter a distribuição de probabilidades de (i);

qualquer que seja o número de regressores presentes no modelo.

# MRLM: Abordagem Matricial

Assim, a equação

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \varepsilon_i, i = 1, 2, \dots, n.$$

também pode ser escrita como

$$y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{21} + \dots + \beta_k x_{k1} + \varepsilon_1$$

$$y_2 = \beta_0 + \beta_1 x_{12} + \beta_2 x_{22} + \dots + \beta_k x_{k2} + \varepsilon_2$$

$$y_3 = \beta_0 + \beta_1 x_{13} + \beta_2 x_{23} + \dots + \beta_k x_{k3} + \varepsilon_3$$

⋮

$$y_n = \beta_0 + \beta_1 x_{1n} + \beta_2 x_{2n} + \dots + \beta_k x_{kn} + \varepsilon_n$$

## MRLM: Abordagem Matricial

As igualdades anteriores podem ser alocadas facilmente em dois vetores colunas ( $n \times 1$ ), descritos a seguir:

$$\underbrace{\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}}_{(n \times 1)} = \underbrace{\begin{pmatrix} \beta_0 + \beta_1 x_{11} + \dots + \beta_k x_{k1} + \varepsilon_1 \\ \beta_0 + \beta_1 x_{12} + \dots + \beta_k x_{k2} + \varepsilon_2 \\ \vdots \\ \beta_0 + \beta_1 x_{1n} + \dots + \beta_k x_{kn} + \varepsilon_n \end{pmatrix}}_{(n \times 1)}$$

# MRLM: Abordagem Matricial

Ainda,

$$\underbrace{\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}}_{(n \times 1)} = \underbrace{\begin{pmatrix} \beta_0 + \beta_1 x_{11} + \dots + \beta_k x_{k1} \\ \beta_0 + \beta_1 x_{12} + \dots + \beta_k x_{k2} \\ \vdots \\ \beta_0 + \beta_1 x_{1n} + \dots + \beta_k x_{kn} \end{pmatrix}}_{(n \times 1)} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{(n \times 1)}$$

# MRLM: Abordagem Matricial

Finalmente,

$$\underbrace{\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}}_{(n \times 1)} = \underbrace{\begin{pmatrix} 1 & x_{11} & \cdots & x_{k1} \\ 1 & x_{12} & \cdots & x_{k2} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n} & \cdots & x_{kn} \end{pmatrix}}_{(n \times (k+1))} \cdot \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}}_{((k+1) \times 1)} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{(n \times 1)}$$



# MRLM: Abordagem Matricial

Vamos definir:

$$\underset{\sim}{\mathbf{y}} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

$$\underset{\sim}{\mathbf{X}} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{k1} \\ 1 & x_{12} & \cdots & x_{k2} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n} & \cdots & x_{kn} \end{pmatrix}$$

$$\underset{\sim}{\omega}_i = (1 \quad x_{1i} \quad \cdots \quad x_{ki})$$

$$\underset{\sim}{\boldsymbol{\beta}} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$$

$$\underset{\sim}{\boldsymbol{\varepsilon}} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

# Modelo de Regressão Linear Múltipla (MRLM)

Assim, utilizando os resultados do *slide* anterior, podemos escrever o modelo de regressão linear múltipla como:

$$\underset{\sim}{y} = \underset{\sim}{X} \underset{\sim}{\beta} + \underset{\sim}{\varepsilon},$$

que é chamado **Modelo Linear Geral**.

# Estimação do MRLM: Métodos dos Mínimos Quadrados (MMQ)

Para determinarmos os estimadores de MQO de  $\beta_0$ ,  $\beta_1, \dots, \beta_k$ , devemos minimizar

$$S = \sum_{i=1}^n (\varepsilon_i)^2 = \varepsilon_1^2 + \dots + \varepsilon_n^2 = \underset{\sim}{\boldsymbol{\varepsilon}}' \underset{\sim}{\boldsymbol{\varepsilon}}$$

ou, ainda,

$$S = \underset{\sim}{\boldsymbol{\varepsilon}}' \underset{\sim}{\boldsymbol{\varepsilon}} = \underset{\sim}{\left( \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \right)}' \underset{\sim}{\left( \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \right)}$$

# Estimação do MRLM: MMQ

Curiosidade

Abrindo a expressão anterior, vem que

$$\begin{aligned} S &= \left( \underset{\sim}{\mathbf{y}} - \underset{\sim}{\mathbf{X}} \underset{\sim}{\boldsymbol{\beta}} \right)' \left( \underset{\sim}{\mathbf{y}} - \underset{\sim}{\mathbf{X}} \underset{\sim}{\boldsymbol{\beta}} \right) = \left( \underset{\sim}{\mathbf{y}}' - \underset{\sim}{\boldsymbol{\beta}}' \underset{\sim}{\mathbf{X}}' \right) \left( \underset{\sim}{\mathbf{y}} - \underset{\sim}{\mathbf{X}} \underset{\sim}{\boldsymbol{\beta}} \right) = \\ &= \underset{\sim}{\mathbf{y}}' \underset{\sim}{\mathbf{y}} - \underset{\sim}{\mathbf{y}}' \underset{\sim}{\mathbf{X}} \underset{\sim}{\boldsymbol{\beta}} - \underset{\sim}{\boldsymbol{\beta}}' \underset{\sim}{\mathbf{X}}' \underset{\sim}{\mathbf{y}} + \underset{\sim}{\boldsymbol{\beta}}' \underset{\sim}{\mathbf{X}}' \underset{\sim}{\mathbf{X}} \underset{\sim}{\boldsymbol{\beta}} \end{aligned}$$

# Estimação do MRLM: MMQ

Curiosidade

Como

$$\underset{\sim}{y}' \underset{\sim}{X} \underset{\sim}{\beta} \quad \text{e} \quad \underset{\sim}{\beta}' \underset{\sim}{X}' \underset{\sim}{y}$$

são escalares e

$$\underset{\sim}{y}' \underset{\sim}{X} \underset{\sim}{\beta} = \left( \underset{\sim}{\beta}' \underset{\sim}{X}' \underset{\sim}{y} \right)'$$

então

$$\underset{\sim}{y}' \underset{\sim}{X} \underset{\sim}{\beta} = \underset{\sim}{\beta}' \underset{\sim}{X}' \underset{\sim}{y}$$

# Estimação do MRLM: MMQ

Curiosidade

Assim

$$S = \underbrace{\mathbf{y}'\mathbf{y}}_{\sim} - 2 \underbrace{\mathbf{y}'\mathbf{X}}_{\sim} \underbrace{\boldsymbol{\beta}}_{\sim} + \underbrace{\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}}_{\sim} \underbrace{\boldsymbol{\beta}}_{\sim}$$

Logo, nosso interesse, agora, é encontrar o resultado para

$$\frac{\partial S}{\partial \boldsymbol{\beta}}$$

# Estimação do MRLM: MMQ

Curiosidade

Lembrando que objetivamos minimizar

$$S(\underset{\sim}{\boldsymbol{\beta}}) = \underset{\sim}{\mathbf{y}}' \underset{\sim}{\mathbf{y}} - 2 \underset{\sim}{\mathbf{y}}' \underset{\sim}{\mathbf{X}} \underset{\sim}{\boldsymbol{\beta}} + \underset{\sim}{\boldsymbol{\beta}}' \underset{\sim}{\mathbf{X}}' \underset{\sim}{\mathbf{X}} \underset{\sim}{\boldsymbol{\beta}}$$

e, utilizando os resultados vistos anteriormente, temos que

$$\frac{\partial S(\underset{\sim}{\boldsymbol{\beta}})}{\partial \underset{\sim}{\boldsymbol{\beta}}} = -2 \underset{\sim}{\mathbf{X}}' \underset{\sim}{\mathbf{y}} + 2 \underset{\sim}{\mathbf{X}}' \underset{\sim}{\mathbf{X}} \underset{\sim}{\boldsymbol{\beta}}$$

# Estimação do MRLM: MMQ

Curiosidade

E, igualando o resultado anterior a zero, vem que

$$-2 \underset{\sim}{X}' \underset{\sim}{y} + 2 \underset{\sim}{X}' \underset{\sim}{X} \underset{\sim}{\hat{\beta}} = \underset{\sim}{0} \Leftrightarrow \underset{\sim}{X}' \underset{\sim}{X} \underset{\sim}{\hat{\beta}} = \underset{\sim}{X}' \underset{\sim}{y}$$

que é o sistema de equações normais na forma matricial.

Para encontrarmos o resultado de interesse, precisaremos supor que **a matriz  $X'X$  admite inversa** (ou seja, precisaremos supor que  $X'X$  é não-singular). Para tanto, assumiremos que **os regressores não apresentam relação linear perfeita.**



## Estimação do MRLM: MMQ

Assim, assumindo que  $X'X$  é não-singular, a solução do sistema de equações normais é dada por

$$\hat{\beta} = (X'X)^{-1} X'y$$

Nota:

$X'$  é a matriz transposta da matriz  $X$

que é o vetor de estimadores de mínimos quadrados do vetor de parâmetros de interesse.

# MRLS: MMQ em Notação Matricial

## Regressão Linear Simples

Dada a equação:

$$Y_i = \alpha + \beta X_i + e_i$$

Que representa o sistema:

$$Y_1 = \alpha + \beta X_1 + e_1$$

$$Y_2 = \alpha + \beta X_2 + e_2$$

...

$$Y_n = \alpha + \beta X_n + e_n$$

Para obter os estimadores de MQO:

$$EQT = \sum \hat{e}_i^2$$

e

$$\frac{\partial EQT}{\partial \hat{\alpha}} = 0 \Rightarrow \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X}$$

$$\frac{\partial EQT}{\partial \hat{\beta}} = 0 \Rightarrow \hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

A equivalente matricial:

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

Que representa o sistema:

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{pmatrix} = \alpha \begin{pmatrix} 1 \\ 1 \\ \dots \\ 1 \end{pmatrix} + \beta \begin{pmatrix} X_1 \\ X_2 \\ \dots \\ X_n \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ \dots \\ e_n \end{pmatrix} \Rightarrow \underbrace{\begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{pmatrix}}_{\mathbf{y}_{n \times 1}} = \underbrace{\begin{pmatrix} 1 & X_1 \\ 1 & X_2 \\ \dots & \dots \\ 1 & X_n \end{pmatrix}}_{\mathbf{X}_{n \times p}} \underbrace{\begin{pmatrix} \alpha \\ \beta \end{pmatrix}}_{\boldsymbol{\beta}_{p \times 1}} + \underbrace{\begin{pmatrix} e_1 \\ e_2 \\ \dots \\ e_n \end{pmatrix}}_{\mathbf{e}_{n \times 1}}$$

Para obter os estimadores de MQO:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} \Rightarrow \hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}} \Rightarrow EQT = \hat{\mathbf{e}}^T \hat{\mathbf{e}} \quad \text{onde} \quad \hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}$$

Então:

$$\frac{\partial EQT}{\partial \hat{\boldsymbol{\beta}}} = 0 \Rightarrow \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})$$

**Nota:**

$p = k+1$ , sendo

$p = n^\circ$  de parâmetros a estimar

$k = n^\circ$  de variáveis

# MRLM: MMQO em Notação Matricial

Dada a equação:

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + e_i$$

Que representa o sistema:

$$Y_1 = \alpha + \beta_1 X_{11} + \beta_2 X_{21} + \dots + \beta_k X_{k1} + e_1$$

$$Y_2 = \alpha + \beta_1 X_{12} + \beta_2 X_{22} + \dots + \beta_k X_{k2} + e_2$$

...

$$Y_n = \alpha + \beta_1 X_{1n} + \beta_2 X_{2n} + \dots + \beta_k X_{kn} + e_n$$

Para obter os estimadores de MQO:

$$EQT = \sum \hat{e}_i$$

e

$$\frac{\partial EQT}{\partial \hat{\alpha}} = 0 \Rightarrow \hat{\alpha} = \dots$$

...

$$\frac{\partial EQT}{\partial \hat{\beta}_k} = 0 \Rightarrow \hat{\beta}_k = \dots$$

A equivalente matricial:

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

Que representa o sistema:

$$\underbrace{\begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{pmatrix}}_{\mathbf{y}_{n \times 1}} = \underbrace{\begin{pmatrix} 1 & X_{11} & X_{21} & \dots & X_{k1} \\ 1 & X_{12} & X_{22} & \dots & X_{k2} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & X_{1n} & X_{2n} & \dots & X_{kn} \end{pmatrix}}_{\mathbf{X}_{n \times p}} \underbrace{\begin{pmatrix} \alpha \\ \beta_1 \\ \dots \\ \beta_k \end{pmatrix}}_{\boldsymbol{\beta}_{p \times 1}} + \underbrace{\begin{pmatrix} e_1 \\ e_2 \\ \dots \\ e_n \end{pmatrix}}_{\mathbf{e}_{n \times 1}}$$

Para obter os estimadores de MQO:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} \quad \longrightarrow \quad \hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}} \quad \longrightarrow \quad EQT = \hat{\mathbf{e}}^T \hat{\mathbf{e}}$$

Então:

$$\frac{\partial EQT}{\partial \hat{\boldsymbol{\beta}}} = 0 \Rightarrow \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})$$

**Nota:**

$p = k+1$ , sendo  
 $p = n^\circ$  de parâmetros a  
estimar  
 $k = n^\circ$  de variáveis

# Estimadores dos MMQ dos Coeficientes do MRLM

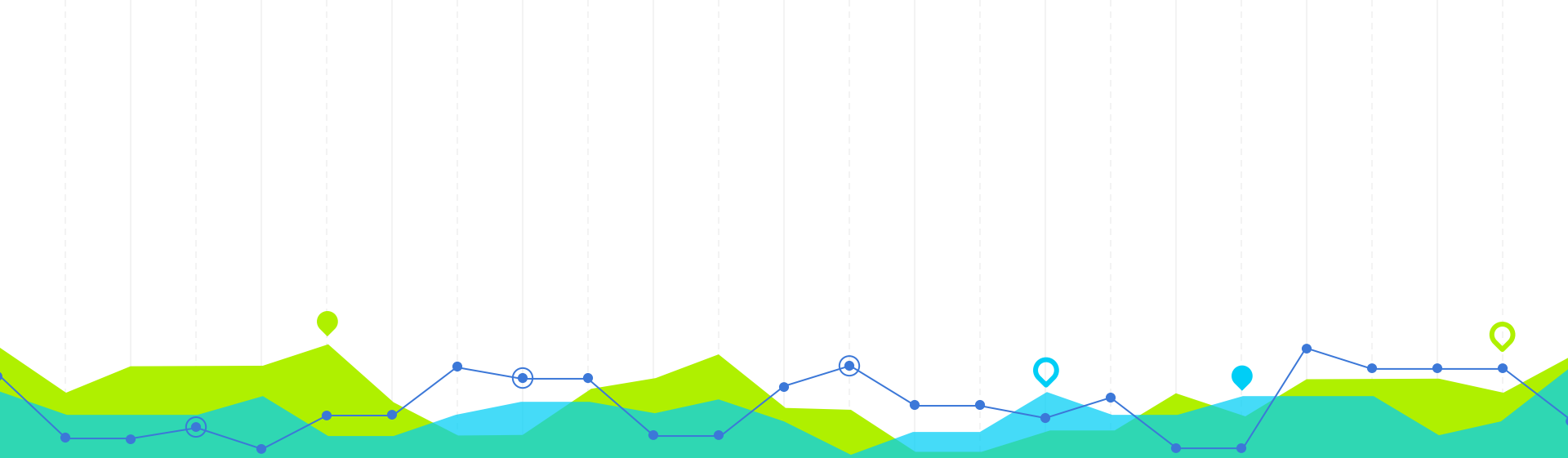
MODELO REGRESSÃO LINEAR

$$y_t = \beta_1 + \beta_2 x_{t2} + \dots + \beta_k x_{tk} + u_t, \quad t = 1, 2, \dots, n.$$

Formulário

EMQ (estimadores dos mínimos quadrados)

Caso geral	Caso particular: $y_t = \beta_1 + \beta_2 x_t + u_t$	Caso particular: Regressão Linear Simples
$b = (X^T X)^{-1} X^T Y$ $\hat{u}_t = y_t - \hat{y}_t$ $s^2 = \frac{\sum \hat{u}_t^2}{(n-k)}$ $\text{Côv}(b   X) = s^2 (X^T X)^{-1}$	$b_1 = \bar{y} - b_2 \bar{x} \quad ; \quad \hat{V}ar(b_1   X) = \frac{s^2 \sum x_t^2}{n \sum x_t^2 - (\sum x_t)^2}$ $b_2 = \frac{n \sum x_t y_t - \sum x_t \sum y_t}{n \sum x_t^2 - (\sum x_t)^2} ;$ $\hat{V}ar(b_2   X) = \frac{ns^2}{n \sum x_t^2 - (\sum x_t)^2}$ $s^2 = \frac{\sum \hat{u}_t^2}{(n-2)}$	<p><b>Nota:</b>  <math>k = n^{\circ}</math> de parâmetros a estimar  <math>k-1 = n^{\circ}</math> de variáveis</p>



# Modelo de Regressão Linear Múltipla

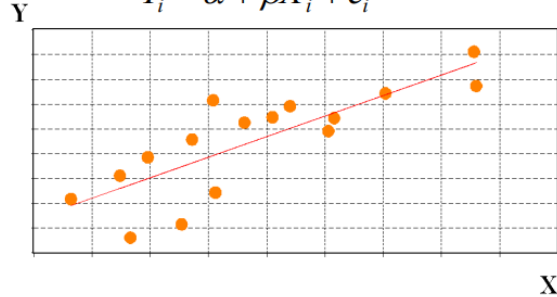
Interpretação dos Coeficientes da Reta de Regressão

# 5

# MRLM: Interpretação dos Coeficientes

## Regressão Linear Simples:

$$Y_i = \alpha + \beta X_i + e_i$$



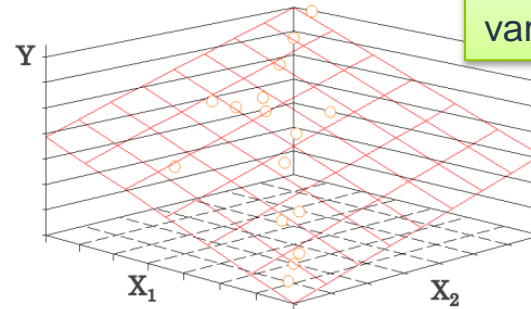
Temos que:

$E[Y / X = 0] = \alpha$  Valor esperado de  $Y$  quando  $X$  é nulo.

$\frac{dY}{dX} = \beta$  Variação marginal esperada em  $Y$  para cada variação unitária em  $X$ .

## Regressão Linear Múltipla:

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + e_i$$



Temos que:

$E[Y / X_1 = 0, X_2 = 0] = \alpha$  Valor esperado de  $Y$  quando ambos  $X_1$  e  $X_2$  são nulos.

$\frac{\partial Y}{\partial X_1} = \beta_1$  Variação marginal esperada em  $Y$  para cada variação unitária em  $X_1$ , mantendo  $X_2$  constante.

$\frac{\partial Y}{\partial X_2} = \beta_2$  Variação marginal esperada em  $Y$  para cada variação unitária em  $X_2$ , mantendo  $X_1$  constante.

Caso particular: MRLM com apenas duas variáveis regressoras

# MRLM: Interpretação dos Coeficientes...

Caso geral: MRLM com  $k$  variáveis regressoras

## Regressão Múltipla

Em um modelo de regressão múltipla, a variável dependente ( $Y$ ) será determinada por mais de uma variável independente ( $X$ ). Genericamente, um modelo de regressão linear múltipla com  $k$  variáveis independentes e  $p$  parâmetros ( $p=k+1$ ) pode ser representado por:

$$Y_i = \alpha + \beta_1 X_{1_i} + \beta_2 X_{2_i} + \dots + \beta_k X_{k_i} + e_i$$

Onde:

$\alpha$  é o valor esperado de  $Y$  quando todas as variáveis independentes forem nulas;

$\beta_1$  é a variação esperada em  $Y$  dado um incremento unitário em  $X_1$ , mantendo-se constantes todas as demais variáveis independentes;

...

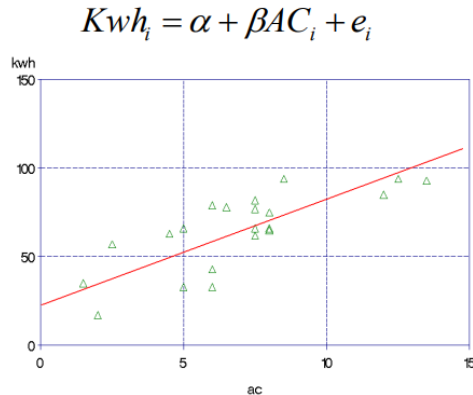
$\beta_k$  é a variação esperada em  $Y$  dado um incremento unitário em  $X_k$ , mantendo-se constantes todas as demais variáveis independentes;

$e_i$  é o erro não explicado pelo modelo;

# MRLM: Exemplo 1 - Interpretação dos Coeficientes

Seja a relação para consumo de energia ( $Kwh$ ), horas de ar condicionado ligado ( $AC$ ) e horas de secadora ligada ( $SEC$ ):

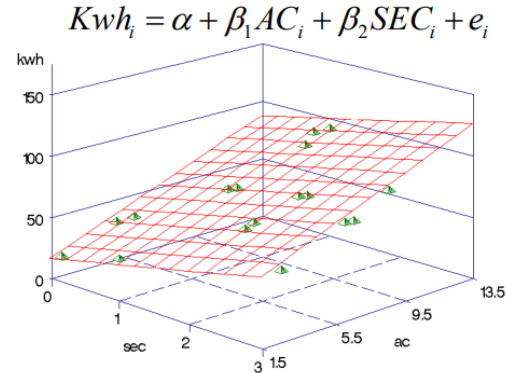
Regressão  
Linear Simples



O coeficiente  $\alpha$  indicará o consumo esperado de energia quando o ar condicionado permanecer desligado.

O coeficiente  $\beta$  indicará o consumo de energia adicional esperado para cada hora adicional com ar condicionado ligado.

Regressão  
Linear Múltipla



O coeficiente  $\alpha$  indicará o consumo esperado de energia quando ambos ar condicionado e secadora permanecerem desligados.

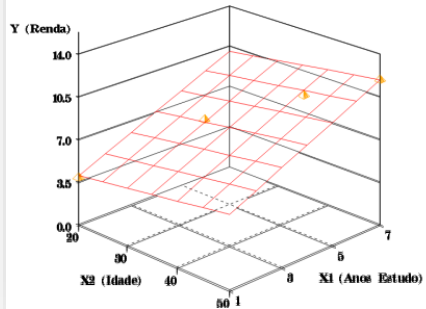
O coeficiente  $\beta_1$  indicará o aumento no consumo de energia esperado para cada hora adicional com ar condicionado ligado, mantendo-se constante o tempo de uso da secadora. Analogamente, O coeficiente  $\beta_2$  indicará efeito isolado de uma hora adicional com a secadora ligada sobre o consumo esperado de energia.



# MRLM: Exemplo 2 - Estimação dos Coeficientes e Interpretação

Seja a relação entre renda familiar em SM ( $Y$ ), anos de estudo ( $X_1$ ) e idade ( $X_2$ ) do responsável pela família:

Y (Renda)	X <sub>1</sub> (Anos Estudo)	X <sub>2</sub> (Idade)
4	1	20
8	4	30
10	6	40
12	7	50



$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + e_i \quad \Rightarrow \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

A função de regressão amostral será dada por:

$$\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{e}} \quad \Rightarrow \quad \begin{pmatrix} 4 \\ 8 \\ 10 \\ 12 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 20 \\ 1 & 4 & 30 \\ 1 & 6 & 40 \\ 1 & 7 & 50 \end{pmatrix} \begin{pmatrix} \hat{\alpha} \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} + \begin{pmatrix} \hat{e}_1 \\ \hat{e}_2 \\ \hat{e}_3 \\ \hat{e}_4 \end{pmatrix}$$

$\mathbf{y}_{4 \times 1}$        $\mathbf{X}_{4 \times 3}$        $\hat{\boldsymbol{\beta}}_{3 \times 1}$        $\hat{\mathbf{e}}_{4 \times 1}$

E as estimativas de MQO:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y}) = \begin{pmatrix} 4 & 18 & 140 \\ 18 & 102 & 730 \\ 140 & 730 & 5400 \end{pmatrix}^{-1} \begin{pmatrix} 34 \\ 180 \\ 1320 \end{pmatrix} = \begin{pmatrix} 1,9 \\ 1 \\ 0,06 \end{pmatrix}$$

O departamento de RH da empresa TEMCO objetiva estudar o comportamento dos salários dos funcionários dos mais diversos setores da empresa.

Para tanto, o gerente de RH, baseando-se numa amostra aleatória de 46 empregados, coletou informações sobre as seguintes variáveis:

- id** – número cadastral do funcionário;
- salario** – anual, em dólares;
- anosemp** – tempo (em anos) na empresa;
- expprev** – experiência anterior (em anos);
- educ** – anos de estudo após o segundo grau;
- sexo** – (feminino = 0, masculino = 1);
- dept** – departamento no qual atua (Compras = 1, Engenharia = 2, Propaganda = 3, Vendas = 4);
- super** – número de empregados sob responsabilidade do empregado.





EViews - [Group: UNTITLED Workfile: TEMCO::Temco]

File Edit Object View Proc Quick Options Window Help

View Proc Object Print Name Freeze Default Sort Transpose Edit+/- Smp+/- Title Sample

obs	ID	SALARIO	ANOSEMP	EXPPREV	EDUC	SEXO	DEPT	SUPER
1	972.0000	47536.00	15.00000	5.000000	6.000000	0.000000	3.000000	4.000000
2	539.0000	23654.00	0.000000	0.000000	0.000000	1.000000	3.000000	2.000000
3	649.0000	37548.00	19.00000	9.000000	4.000000	0.000000	3.000000	6.000000
4	824.0000	36578.00	4.000000	4.000000	8.000000	0.000000	3.000000	8.000000
5	649.0000	54679.00	20.00000	3.000000	6.000000	1.000000	3.000000	4.000000
6	624.0000	53234.00	25.00000	0.000000	6.000000	0.000000	3.000000	3.000000
7	891.0000	31425.00	7.000000	6.000000	5.000000	1.000000	3.000000	6.000000
8	974.0000	39743.00	9.000000	6.000000	5.000000	1.000000	2.000000	1.000000
9	648.0000	26452.00	1.000000	3.000000	2.000000	1.000000	2.000000	0.000000
10	321.0000	34632.00	5.000000	4.000000	4.000000	0.000000	2.000000	0.000000
11	264.0000	35631.00	6.000000	4.000000	4.000000	0.000000	2.000000	2.000000
12	291.0000	46211.00	14.00000	5.000000	6.000000	1.000000	2.000000	5.000000
13	267.0000	34231.00	6.000000	2.000000	6.000000	0.000000	2.000000	3.000000
14	548.0000	26548.00	5.000000	1.000000	0.000000	0.000000	2.000000	2.000000
15	555.0000	36512.00	6.000000	6.000000	4.000000	1.000000	2.000000	2.000000
16	366.0000	34869.00	7.000000	5.000000	4.000000	1.000000	2.000000	1.000000
17	246.0000	41255.00	9.000000	4.000000	6.000000	0.000000	2.000000	4.000000
18	215.0000	39331.00	9.000000	3.000000	6.000000	1.000000	2.000000	1.000000
19	814.0000	35487.00	8.000000	2.000000	2.000000	1.000000	2.000000	2.000000
20	212.0000	36487.00	6.000000	5.000000	2.000000	0.000000	2.000000	3.000000
21	526.0000	68425.00	25.00000	2.000000	12.00000	0.000000	2.000000	1.000000
22	778.0000	69246.00	22.00000	3.000000	10.00000	0.000000	2.000000	45.00000

**Quadro 1 - Parte de uma planilha que contém informações sobre os empregados da empresa TEMCO.**

Como parte do estudo, a gerente de RH propôs a estimação dos parâmetros do seguinte modelo de regressão múltipla:

$$\text{salario} = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{anosemp} + \varepsilon$$

- a) Em termos do problema,  $\beta_0$  apresenta algum significado prático?
- b) Qual o sinal esperado para  $\beta_1$ ? E para  $\beta_2$ ?
- c) Encontre as estimativas dos parâmetros, via mínimos quadrados ordinários, escreva a equação estimada e interprete os resultados obtidos, em termos do problema de interesse.

33



## Exercícios (a) e (b): Interpretação dos Parâmetros

Interpretação dos parâmetros do modelo proposto, em termos do problema:

$\beta_0$  – salário médio dos funcionários da empresa TEMCO, que acabaram de entrar na empresa (ou que ainda não completaram um ano) e que não apresentam nenhum ano de escolaridade após o segundo grau;

$\beta_1$  – efeito no salário médio dos funcionários da empresa TEMCO, dada a variação de um ano no tempo de escolaridade após o segundo grau, mantendo constante a variável *anosemp*; e

$\beta_2$  – efeito no salário médio dos funcionários da empresa TEMCO, dada a variação de um ano no tempo de empresa, mantendo constante a variável *educ*.

## Exercícios (a) e (b): Interpretação dos Parâmetros

Dependent Variable: SALARIO

Method: Least Squares

Date: 08/26/12 Time: 15:45

Sample: 1 46

Included observations: 46

SALARIO=C(1)+C(2)\*EDUC+C(3)\*ANOSEMP

	Coefficient	Std. Error	t-Statistic	Prob.
C(1)	23177.47	1769.732	13.09660	0.0000
C(2)	1916.489	379.2670	5.053139	0.0000
C(3)	672.3250	141.6725	4.745629	0.0000
R-squared	0.739927	Mean dependent var		39827.39
Adjusted R-squared	0.727830	S.D. dependent var		10999.24
S.E. of regression	5738.291	Akaike info criterion		20.21070
Sum squared resid	1.42E+09	Schwarz criterion		20.32996
Log likelihood	-461.8462	Hannan-Quinn criter.		20.25538
F-statistic	61.16907	Durbin-Watson stat		1.229794
Prob(F-statistic)	0.000000			

## Exercício (c): Modelo Estimado

$$\hat{\text{salário}} = 23177,47 + 1916,49 \text{educ} + 672,32 \text{anosemp}$$

**Pergunta:** qual o salário médio estimado para pessoas com 3 anos de escolaridade após o 2º grau e com 5 anos na empresa?

$$\hat{\text{salário}} = 23.177,47 + 1.916,49 * 3 + 672,33 * 5$$

$$\hat{\text{salário}} = 32288,54$$

# Obrigada!

Questões?

