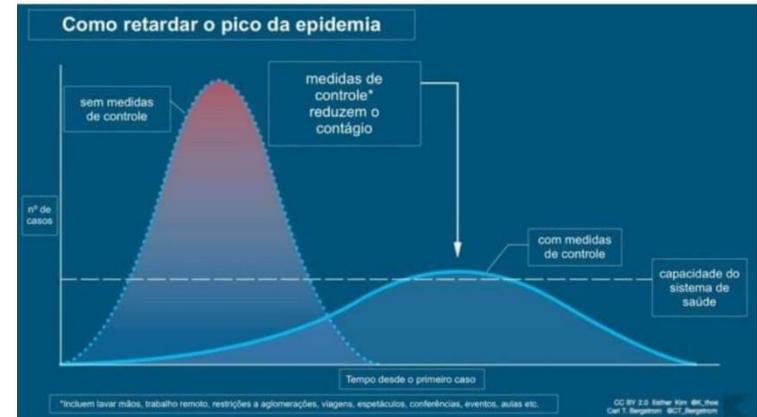


MAE 5776

ANÁLISE MULTIVARIADA



A Curva de Contágio e suas consequências no Sistema de Saúde e em nosso Dia a Dia.

Júlia M Pavan Soler
pavan@ime.usp.br

1º Semestre/2020

Já vimos ☺

MAE5776

$$Y_{n \times p} = (Y_{ij}) \in \mathfrak{R}^{n \times p}$$

Matriz de Dados: Estatísticas descritivas multivariadas

- Definidas no espaço das colunas ($\mathfrak{R}^p, \mathfrak{R}^{p \times p}$): $\bar{Y}_{p \times 1}, S_{p \times p}, R_{p \times p}, S_{p \times p}^{-1}$
- Definidas no espaço das linhas ($\mathfrak{R}^{n \times n}$): $D_{n \times n} = (d_{ij}^2)$; $d_{Eij}^2, d_{Pij}^2, d_{Mij}^2$

Regiões (elipsóides) de Concentração de Observações ($Y_i \in \mathfrak{R}^p$):

$$R(Y_i) = \left(Y_i \in \mathfrak{R}^p; d_M^2(Y_i; \mu) = (Y_i - \bar{Y})' S_u^{-1} (Y_i - \bar{Y}) \leq c^2; c^2 = \chi_p^2(\alpha) \right)$$

Matriz Aleatória: $Y_{n \times p} \sim N_{n \times p} (1_n \otimes \mu'_{p \times 1}; \Omega_{np \times np} = \Psi_{n \times n} \otimes \Sigma_{p \times p})$

- Estimadores e Distribuições Amostrais:
$$\left. \begin{array}{l} Y_{i_{p \times 1}} \stackrel{iid}{\sim} N_p(\mu; \Sigma) \\ \bar{Y}_{p \times 1} \sim N_p(\mu; \Sigma/n) \\ nS_{p \times p} \sim W_p(n-1; \Sigma) \end{array} \right\}$$

- Regiões (elipsoides) de Confiança para μ :

$$R(\mu | Y) = \left\{ \mu \in \mathfrak{R}^p; n(\bar{Y} - \mu)' S_u^{-1} (\bar{Y} - \mu) \leq c^2; c^2 = T_{(p; n-1)}^2(\alpha) = \frac{(n-1)p}{(n-p)} F_{p, (n-p)}(\alpha) \right\}$$

Caso de
Uma única
População

Inferência – Análise Multivariada

Regiões de Confiança e
Testes Multivariados

x

Intervalos de Confiança e
Testes Univariados

$$Y_{i_{p \times 1}} \stackrel{iid}{\sim} (\mu; \Sigma), \quad H_0 : \mu_{p \times 1} = 0_{p \times 1}$$

x

$$\begin{cases} H_{01} : \mu_1 = 0 \\ \dots \\ H_{0p} : \mu_p = 0 \end{cases}$$

- Há interesse na análise conjunta de múltiplas variáveis
- Realizar inferências mais “precisas” devido a incorporar a informação da covariância entre variáveis
- Realizar comparações entre os parâmetros associados às diferentes variáveis: construir contrastes entre médias das variáveis
- Construir níveis de significância coletivos × Correções para múltiplos testes

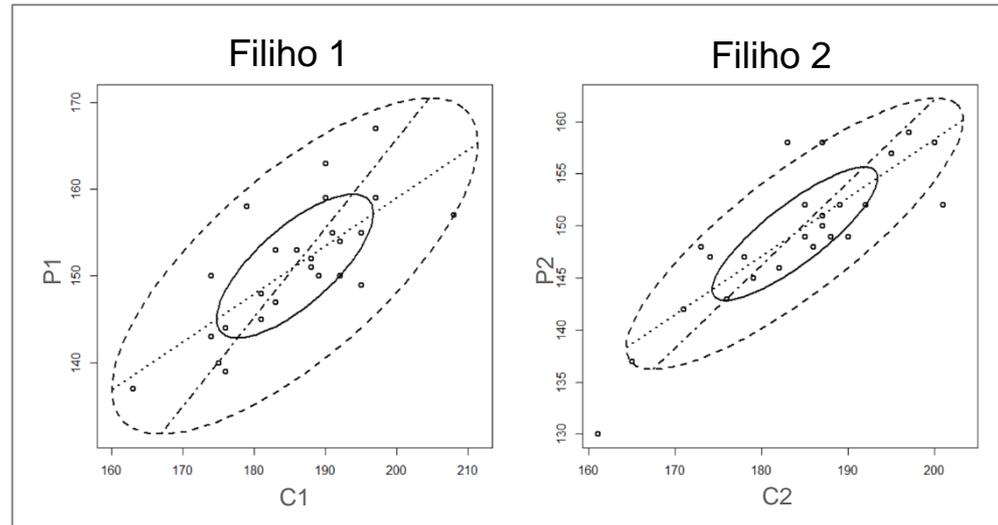
Bonferroni, FDR

Inferência sobre para $\mu \in \mathbb{R}^p$

Dados do Primeiro e Segundo filhos

Família	1° Filho		2° Filho	
	Comprimento	Perímetro	Comprimento	Perímetro
1	191	155	179	145
2	195	149	201	152
3	181	148	185	149
4	183	153	188	149
5	176	144	171	142
6	208	157	192	152
7	189	150	190	149
8	197	159	189	152
9	188	152	197	159
10	192	150	187	151
11	179	158	186	148
12	183	147	174	147
13	174	150	185	152
14	190	159	195	157
15	188	151	187	158
16	163	137	161	130
17	195	155	183	158
18	186	153	173	148
19	181	145	182	146
20	175	140	165	137
21	192	154	185	152
22	174	143	178	147
23	176	139	176	143
24	197	167	200	158
25	190	163	187	150

Boxplot Bivariado: Elipses de concentração de observações



Estatísticas Descritivas:

$$\bar{Y} = (185,72 \quad 151,12 \quad 183,84 \quad 149,24)'$$

$$S_u = \begin{pmatrix} 91,481 & 50,753 & 66,875 & 44,267 \\ & 52,186 & 49,259 & 33,651 \\ & & 96,775 & 54,278 \\ & & & 43,222 \end{pmatrix}'$$

$$Y_{25 \times 4} = (Y_1, \dots, Y_{25})'; Y_{i_{4 \times 1}} \stackrel{iid}{\sim} (\mu; \Sigma)$$

Região de Confiança e Teste para $\mu \in \mathbb{R}^p$

$$Y_{25 \times 4} = (Y_1, \dots, Y_{25})'; Y_{i_{4 \times 1}} \stackrel{iid}{\sim} N_4(\mu; \Sigma)$$

Família	1° Filho		2° Filho	
	Comprimento	Perímetro	Comprimento	Perímetro
1	191	155	179	145
2	195	149	201	152
3	181	148	185	149
4	183	153	188	149
5	176	144	171	142
6	208	157	192	152
7	189	150	190	149
8	197	159	189	152
9	188	152	197	159
10	192	150	187	151
11	179	158	186	148
12	183	147	174	147
13	174	150	185	152
14	190	159	195	157
15	188	151	187	158
16	163	137	161	130
17	195	155	183	158
18	186	153	173	148
19	181	145	182	146
20	175	140	165	137
21	192	154	185	152
22	174	143	178	147
23	176	139	176	143
24	197	167	200	158
25	190	163	187	150

Caso 1: Há interesse em inferências para o vetor de médias μ

Região de Confiança (baseada na estatística T2 de Hotelling):

$$T^2(\alpha) = n(\bar{Y} - \mu)' S_u^{-1}(\bar{Y} - \mu) \leq \frac{(n-1)p}{(n-p)} F_{p, n-p}(\alpha)$$

Nível de significância global $n = 25, p = 4, \alpha = 5\% \Rightarrow 12.98$

$$\begin{cases} H_0 : \mu = (185, 150, 180, 148)' \\ H_1 : \mu \neq (185, 150, 180, 148)' \end{cases} \quad T^2 = 6,44$$

Conclusão?

$$\begin{cases} H_0 : \mu = (185, 150, 180, 145)' \\ H_1 : \mu \neq (185, 150, 180, 145)' \end{cases} \quad T^2 = 17,98$$

Pense na situação de Intervalos de Confiança e Testes de hipóteses individuais para cada variável!!

Região de Confiança e Teste para $\mu \in \mathbb{R}^p$

$$Y_{25 \times 4} = (Y_1, \dots, Y_{25})'; Y_{i_{4 \times 1}} \stackrel{iid}{\sim} N_4(\mu; \Sigma)$$

Família	1° Filho		2° Filho	
	Comprimento	Perímetro	Comprimento	Perímetro
1	191	155	179	145
2	195	149	201	152
3	181	148	185	149
4	183	153	188	149
5	176	144	171	142
6	208	157	192	152
7	189	150	190	149
8	197	159	189	152
9	188	152	197	159
10	192	150	187	151
11	179	158	186	148
12	183	147	174	147
13	174	150	185	152
14	190	159	195	157
15	188	151	187	158
16	163	137	161	130
17	195	155	183	158
18	186	153	173	148
19	181	145	182	146
20	175	140	165	137
21	192	154	185	152
22	174	143	178	147
23	176	139	176	143
24	197	167	200	158
25	190	163	187	150

Caso Multivariado:

↓

$$T^2(\alpha) = n(\bar{Y} - \mu)' S_u^{-1}(\bar{Y} - \mu) \leq \frac{(n-1)p}{(n-p)} F_{p, n-p}(\alpha)$$

Nível de significância global

$\alpha = 5\% \Rightarrow 12.98$

$$\begin{cases} H_0 : \mu = (185, 150, 180, 145)' \\ H_1 : \mu \neq (185, 150, 180, 145)' \end{cases} \quad T^2 = 17,98$$

Conclusão?
Qual é o α global
(para os múltiplos testes)?

Casos Univariados:

$$t_j = \frac{(\bar{Y}_j - \mu_j)}{\sqrt{s_j / n}} \sim t_{n-1}(\alpha); \quad j = 1, 2, 3, 4$$

$\alpha = 5\% \Rightarrow 4, 26$

$$t_j^2 = n(\bar{Y}_j - \mu_j) s_j^{-1} (\bar{Y}_j - \mu_j) \sim F_{(1, n-1)}(\alpha); \quad j = 1, 2, 3, 4$$

$t_1^2 = 0.14; \quad t_2^2 = 0.58; \quad t_3^2 = 3.66; \quad t_4^2 = 9.98$

Inferência sobre um Vetor de Médias

Correspondência entre as Estatísticas de Teste dos casos Uni e Multivariado

$$t^2 = \frac{(\bar{Y} - \mu)^2}{s^2/n} = n \underbrace{(\bar{Y} - \mu)(s^2)^{-1}(\bar{Y} - \mu)} \sim t_{(n-1)}^2 = F_{1,(n-1)}$$


$$H_0 : \mu = \mu_0 \Rightarrow t^2 = \frac{(\bar{Y} - \mu_0)^2}{s^2/n} > t_{(n-1)}^2 = F_{1,(n-1)}(\alpha)$$

Pode ser calculada para cada variável

$$T^2 = nd' S_u^{-1} d = n \underbrace{(\bar{Y} - \mu)' S_u^{-1} (\bar{Y} - \mu)} \sim \frac{(n-1)p}{(n-p)} F_{p,(n-p)}$$

$$H_0 : \mu = \mu_0 \Rightarrow T^2 = n(\bar{Y} - \mu_0)' S^{-1} (\bar{Y} - \mu_0) > \frac{(n-1)p}{(n-p)} F_{p,(n-p)}(\alpha)$$

Teste conjunto para as p variáveis

Região de Confiança e Testes

Caso de Duas Populações - Amostras Pareadas

Amostra Pareada \Rightarrow respostas multivariadas são avaliadas na mesma unidade amostral em “duas” condições diferentes (Ex.: Antes e Depois de uma intervenção)

Duas Populações

$$Y_{1n \times p}; Y_{1i p \times 1} = (Y_{1i1}, Y_{1i2}, \dots, Y_{1ip})' \quad Y_{2n \times p}; Y_{2i p \times 1} = (Y_{2i1}, Y_{2i2}, \dots, Y_{2ip})' \quad i = 1, 2, \dots, n$$

$$Y_{1i p \times 1} \stackrel{iid}{\sim} N_p(\mu_1; \Sigma_1) \longleftrightarrow Y_{2i p \times 1} \stackrel{iid}{\sim} N_p(\mu_2; \Sigma_2)$$

Observações Pareadas



$$D_{ij} = Y_{1ij} - Y_{2ij} \quad j = 1, 2, \dots, p, \quad i = 1, 2, \dots, n$$

$$D_{i p \times 1} \stackrel{iid}{\sim} N_p(\mu_D = \mu_1 - \mu_2; \Sigma_D); \quad \bar{D}_{p \times 1} \sim N(\mu_D = \mu_1 - \mu_2; \Sigma_D / n)$$

\Rightarrow Uma Única População de Diferenças

$$T^2 = n \left(\bar{D} - \mu_D \right)' \underset{\uparrow S_D}{S_D^{-1}} \left(\bar{D} - \mu_D \right) \sim \frac{(n-1)p}{(n-p)} F_{p, n-p}$$

Elipse de Confiança para a diferença entre vetores de médias:

$$R(Y_1, Y_2) = \left\{ \mu_D \in \mathbb{R}^p; n \left(\bar{D} - \mu_D \right)' S_D^{-1} \left(\bar{D} - \mu_D \right) \leq c_\alpha^2 \right\}$$

Região de Confiança e Testes

Caso de Duas Populações Independentes

Amostras Independentes - Homocedasticidade

$$Y_{1n_1 \times p}; Y_{1i} \stackrel{iid}{\sim} N_p(\mu_1; \Sigma_1); \quad Y_{2n_2 \times p}; Y_{2i} \stackrel{iid}{\sim} N_p(\mu_2; \Sigma_2); \quad \Sigma_1 = \Sigma_2 = \Sigma$$

$$\Rightarrow \bar{D}_{p \times 1} = \bar{Y}_1 - \bar{Y}_2 \sim N_p\left(\mu_D = \mu_1 - \mu_2; \Sigma\left(\frac{1}{n_1} + \frac{1}{n_2}\right)\right)$$

$$S_c = \frac{(n_1 - 1)S_{u1} + (n_2 - 1)S_{u2}}{n_1 + n_2 - 2}$$

S_c : Matriz de covariância comum aos grupos.
Estimador de Σ

$$T^2 = (\bar{D} - \mu_D)' \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) S_c \right]^{-1} (\bar{D} - \mu_D) \sim \frac{(n_1 + n_2 - 2)p}{(n_1 + n_2 - p - 1)} F_{(p, (n_1 + n_2 - p - 1))}$$

Elipse de Confiança para a diferença entre vetores de médias:

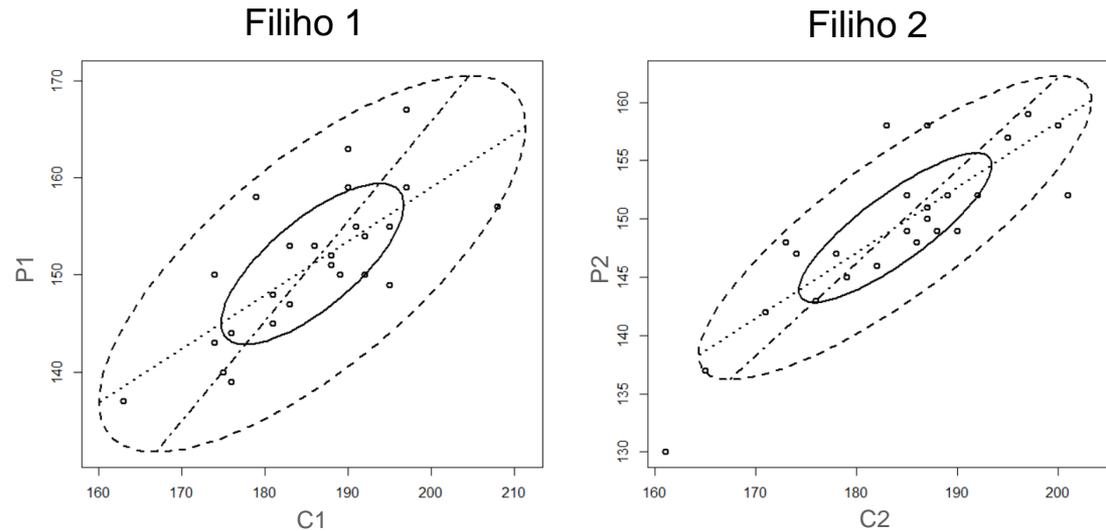
$$R(Y_1, Y_2) = \left\{ \mu_D \in \mathbb{R}^2; (\bar{D} - \mu_D)' \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) S_c \right]^{-1} (\bar{D} - \mu_D) \leq c_\alpha^2 \right\}$$

Região de Confiança e Testes - Duas Populações

Retomando os dados do Primeiro e Segundo filho em 25 famílias

Boxplot Bivariado: Elipses de concentração

	C1	P1	C2	P2
1	191	155	179	145
2	195	149	201	152
3	181	148	185	149
4	183	153	188	149
5	176	144	171	142
6	208	157	192	152
7	189	150	190	149
8	197	159	189	152
9	188	152	197	159
10	192	150	187	151
11	179	158	186	148
12	183	147	174	147
13	174	150	185	152
14	190	159	195	157
15	188	151	187	158
16	163	137	161	130
17	195	155	183	158
18	186	153	173	148
19	181	145	182	146
20	175	140	165	137
21	192	154	185	152
22	174	143	178	147
23	176	139	176	143
24	197	167	200	158
25	190	163	187	150



Vamos supor duas “possíveis estruturas para esses dados:

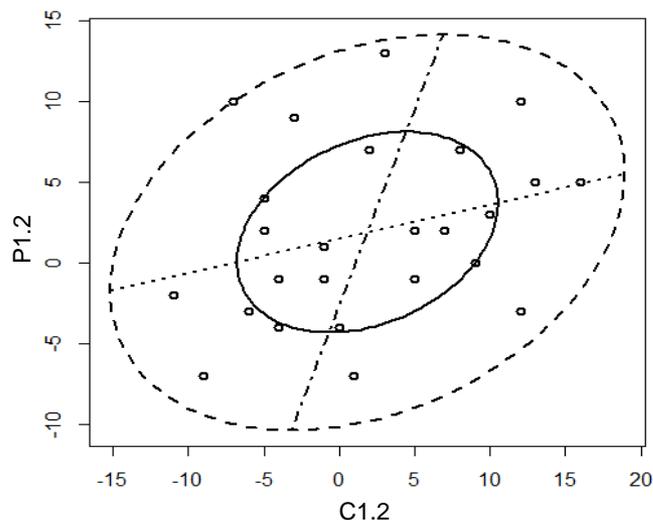
- **Observações pareadas** (estrutura original): primeiro e segundo filho avaliados na mesma família (n=25 pares de irmãos)
- **Observações independentes**: dados do primeiro e do segundo filho avaliados em famílias diferentes ($n_1 = n_2 = 25$; n=50)

Região de Confiança e Testes - Duas Populações Pareadas

Dados do Primeiro e Segundo
filho em 25 famílias

	C1	P1	C2	P2	C1-C2	P1-P2
1	191	155	179	145	12	10
2	195	149	201	152	-6	-3
3	181	148	185	149	-4	-1
4	183	153	188	149	-5	4
5	176	144	171	142	5	2
6	208	157	192	152	16	5
7	189	150	190	149	-1	1
8	197	159	189	152	8	7
9	188	152	197	159	-9	-7
10	192	150	187	151	5	-1
11	179	158	186	148	-7	10
12	183	147	174	147	9	0
13	174	150	185	152	-11	-2
14	190	159	195	157	-5	2
15	188	151	187	158	1	-7
16	163	137	161	130	2	7
17	195	155	183	158	12	-3
18	186	153	173	148	13	5
19	181	145	182	146	-1	-1
20	175	140	165	137	10	3
21	192	154	185	152	7	2
22	174	143	178	147	-4	-4
23	176	139	176	143	0	-4
24	197	167	200	158	-3	9
25	190	163	187	150	3	13

- Observações pareadas (estrutura original)
- ⇒ medida resumo: **diferença entre as variáveis**



$$\bar{D} = (1.88, 1.88) \quad S_D = \begin{pmatrix} 56.78 & 11.98 \\ 11.98 & 29.28 \end{pmatrix}$$

Região de Confiança e Testes - Duas Populações Pareadas

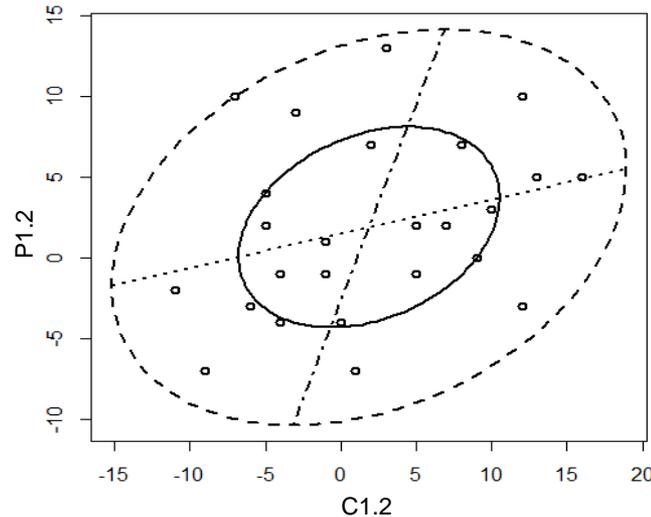
Dados do Primeiro e Segundo
filho em 25 famílias

	C1	P1	C2	P2	C1.2	P1.2
1	191	155	179	145	12	10
2	195	149	201	152	-6	-3
3	181	148	185	149	-4	-1
4	183	153	188	149	-5	4
5	176	144	171	142	5	2
6	208	157	192	152	16	5
7	189	150	190	149	-1	1
8	197	159	189	152	8	7
9	188	152	197	159	-9	-7
10	192	150	187	151	5	-1
11	179	158	186	148	-7	10
12	183	147	174	147	9	0
13	174	150	185	152	-11	-2
14	190	159	195	157	-5	2
15	188	151	187	158	1	-7
16	163	137	161	130	2	7
17	195	155	183	158	12	-3
18	186	153	173	148	13	5
19	181	145	182	146	-1	-1
20	175	140	165	137	10	3
21	192	154	185	152	7	2
22	174	143	178	147	-4	-4
23	176	139	176	143	0	-4
24	197	167	200	158	-3	9
25	190	163	187	150	3	13

C1-C2 P1-P2

- Observações pareadas (estrutura original)

⇒ medida resumo: **diferença entre as variáveis**



O centróide está distante do vetor nulo?

$$T^2(\alpha) = n(\bar{D} - \mu_D)' S_D^{-1} (\bar{D} - \mu_D) \leq \frac{(n-1)p}{(n-p)} F_{p, n-p}(\alpha)$$

$$n = 25, p = 2, \alpha = 5\% \Rightarrow 7,14$$

$$\begin{cases} H_0 : \mu_D = (0,0)' \\ H_1 : \mu_D \neq (0,0)' \end{cases} \quad T^2 = 3,61; \quad p = 0,1994$$

Região de Confiança e Testes - Duas Populações Independentes

Dados do Primeiro e Segundo filho

	C1	P1	C2	P2
1	191	155	179	145
2	195	149	201	152
3	181	148	185	149
4	183	153	188	149
5	176	144	171	142
6	208	157	192	152
7	189	150	190	149
8	197	159	189	152
9	188	152	197	159
10	192	150	187	151
11	179	158	186	148
12	183	147	174	147
13	174	150	185	152
14	190	159	195	157
15	188	151	187	158
16	163	137	161	130
17	195	155	183	158
18	186	153	173	148
19	181	145	182	146
20	175	140	165	137
21	192	154	185	152
22	174	143	178	147
23	176	139	176	143
24	197	167	200	158
25	190	163	187	150

hipoteticamente



	C	P
1	191	155
2	195	149
3	181	148
4	183	153
5	176	144
6	208	157
7	189	150
8	197	159
9	188	152
10	192	150
...		
24	197	167
25	190	163

1	179	145
2	201	152
3	185	149
4	188	149
5	171	142
6	192	152
7	190	149
8	189	152
9	197	159
10	187	151
...		
24	200	158
25	187	150

- Observações independentes: dados do primeiro e do segundo filho avaliados em famílias diferentes ($n_1 = n_2 = 25$; $n = 50$)

$$\bar{Y}_1 = (185.72, 151.12) \quad \bar{Y}_2 = (183.84, 149.24)$$

$$S_1 = \begin{pmatrix} 95.29 & 52.87 \\ 52.87 & 54.36 \end{pmatrix} \quad S_2 = \begin{pmatrix} 100.81 & 56.54 \\ 56.54 & 45.02 \end{pmatrix}$$

$$\bar{D} = (1.88, 1.88) \quad S_c = \begin{pmatrix} 98.05 & 54.70 \\ 54.70 & 49.69 \end{pmatrix}$$

Assumindo $\Sigma_1 = \Sigma_2 = \Sigma$

Teste M de Box: Homocedasticidade
Chi-Sq = 2.2674, df = 3, p-value = 0.5188

Região de Confiança e Testes - Duas Populações Independentes

Dados do Primeiro e Segundo filho

	C	P
1	191	155
2	195	149
3	181	148
4	183	153
5	176	144
6	208	157
7	189	150
8	197	159
9	188	152
10	192	150
...		
24	197	167
25	190	163

1	179	145
2	201	152
3	185	149
4	188	149
5	171	142
6	192	152
7	190	149
8	189	152
9	197	159
10	187	151
...		
24	200	158
25	187	150

$$T^2 = (\bar{D} - \mu_D)' \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) S_c \right]^{-1} (\bar{D} - \mu_D) \leq \frac{(n_1 + n_2 - 2)p}{(n_1 + n_2 - p - 1)} F_{(p, (n_1 + n_2 - p - 1))}$$

$$n_1 = n_2 = 25, p = 2, \alpha = 5\% \Rightarrow 6.53$$

$$\begin{cases} H_0 : \mu_D = (0, 0)' \\ H_1 : \mu_D \neq (0, 0)' \end{cases} \quad T^2 = 0.90; \quad p = 0,1824$$

É esperado que delineamentos com Amostras Pareadas ofereçam mais precisão e maior poder comparados aos delineamentos com Amostras Independentes!

Por que isso não foi atingido no caso dos dados dos Filhos?

Teste da Igualdade de Matrizes de Covariância

Comparação de Vetores de Médias - Amostras Independentes

$$\left. \begin{aligned}
 & Y_{1n_1 \times p}; Y_{1i} \stackrel{iid}{\sim} N_p(\mu_1; \Sigma_1) \quad Y_{2n_2 \times p}; Y_{2i} \stackrel{iid}{\sim} N_p(\mu_2; \Sigma_2) \\
 & \Rightarrow H_0 : \mu_1 - \mu_2 = 0 \quad ; \Sigma_1 = \Sigma_2 = \Sigma
 \end{aligned} \right\}$$

$$S_g = \frac{1}{n_g - 1} \sum_{i=1}^{n_g} (Y_{gi} - \bar{Y}_g)(Y_{gi} - \bar{Y}_g)'; g = 1, 2$$

$$S_c = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{n_1 + n_2 - 2};$$

É uma Hipótese condicional.
Logo a homocedasticidade deve ser verificada.

- **Teste M de Box:** $\Rightarrow H_0 : \Sigma_g = \Sigma; \mu_g \in \mathfrak{R}^p$

$$-2 \ln \lambda = n \ln S_c - \sum_{g=1}^G [n_g \ln |S_g|]$$

$$M = (1 - c) \left\{ \left[\sum_{g=1}^G (n_g - 1) \right] \ln |S_c| - \sum_{g=1}^G [(n_g - 1) \ln |S_g|] \right\} \sim \chi_{\frac{1}{2}p(p+1)(G-1)}^2$$

$$c = \left[\sum_{g=1}^G \frac{1}{(n_g - 1)} \right] \left[\frac{2p^2 + 3p - 1}{6(p+1)(G-1)} \right]$$

S (com divisor n), S_u (divisor n-1)
Para p=1 o teste de Box equivale ao teste de Bartlett.

Critério “prático” de heterocedasticidade sugerido em Johnson and Wichern 1992:

$$\sigma_{gij} = 4\sigma_{g'ij}$$

Inferência sobre Vetores de Médias de Duas Populações

Medidas de produtividade e altura de plantas de duas variedades

Variedade A		Variedade B	
X11	X12	X21	X22
5,7	2,1	4,4	1,8
8,9	1,9	7,5	1,75
6,2	1,98	5,4	1,78
5,8	1,92	4,6	1,89
6,8	2	5,9	1,9
6,2	2,01		

- Teste a igualdade do vetor de médias das duas variedades, sob homocedasticidade.
- Obtenha os intervalos de confiança simultâneos e de Bonferroni.

Inferência sobre Vetores de Médias de Duas Populações

	Variedade A		Variedade B	
	X11	X12	X21	X22
	5,7	2,1	4,4	1,8
	8,9	1,9	7,5	1,75
	6,2	1,98	5,4	1,78
	5,8	1,92	4,6	1,89
	6,8	2	5,9	1,9
	6,2	2,01		
Média	6,6	1,985	5,56	1,824
Var.	1,42	0,005	1,543	0,0045

$$\bar{D} = \begin{pmatrix} \bar{D}_1 = 6,6 - 5,56 = 1,04 \\ \bar{D}_2 = 1,985 - 1,824 = 0,161 \end{pmatrix}$$

$$S_1 = \begin{pmatrix} 1,4200 & -0,0504 \\ -0,0504 & 0,0051 \end{pmatrix}$$

$$S_2 = \begin{pmatrix} 1,543 & -0,037 \\ -0,037 & 0,0045 \end{pmatrix}$$

$$S_c = \begin{pmatrix} 1,4745 & -0,0442 \\ -0,0442 & 0,0049 \end{pmatrix}$$

$$\Rightarrow H_0 : \mu_1 - \mu_2 = 0 \quad ; \Sigma_1 = \Sigma_2 = \Sigma$$

$$T^2 = (\bar{Y}_1 - \bar{Y}_2)' \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) S_c \right]^{-1} (\bar{Y}_1 - \bar{Y}_2) = 30,584 \Rightarrow p = 0,0027$$

Conclusão?

Existe diferença significativa entre os dois grupos para alguma combinação linear das variáveis.

Inferência sobre um Vetor de Médias

Intervalos de Confiança Univariados e Simultâneos

$$I.C(\mu_k) \text{ a } 100(1-\alpha)\% = \left(\bar{Y}_k - t_{n-1}(\alpha/2) \sqrt{\frac{S_{kk}}{n}}; \bar{Y}_k + t_{n-1}(\alpha/2) \sqrt{\frac{S_{kk}}{n}} \right)$$

Sob independência $\Rightarrow P(\text{todos os } p \text{ intervalos conterem os } \mu_k \text{'s}) = (1-\alpha)^p$

 $(1-\alpha) = 0,95, p = 2, n = 25 \Rightarrow P(\text{coletivo}) = (0,95)^2 = 0,90 \quad t_{24}(\alpha/2) = 2,06$

$(1-\alpha) = 0,95, p = 4, n = 15 \Rightarrow P(\text{coletivo}) = (0,95)^4 = 0,81 \quad t_{14}(\alpha/2) = 2,14$

$$I.C.S(\mu_k) \text{ a } 100(1-\alpha)\% = \left(\bar{Y}_k - \sqrt{\frac{(n-1)p}{(n-p)} F_{p,(n-p)}(\alpha) \frac{S_{kk}}{n}}; \bar{Y}_k + \sqrt{\frac{(n-1)p}{(n-p)} F_{p,(n-p)}(\alpha) \frac{S_{kk}}{n}} \right)$$

Nível coletivo igual a $(1-\alpha) \Rightarrow$ intervalos simultâneos são mais largos que os individuais:

$(1-\alpha) = 0,95, p = 2, n = 25 \Rightarrow \sqrt{\frac{(n-1)p}{(n-p)} F_{p,n-p}(\alpha)} = 2,67$ ← mais largos

$(1-\alpha) = 0,95, p = 4, n = 15 \Rightarrow \sqrt{\frac{(n-1)p}{(n-p)} F_{p,n-p}(\alpha)} = 4,13$

Inferência sobre um Vetor de Médias

O Método de Bonferroni para Comparações Múltiplas

$$\Rightarrow P(\text{todos os } p \text{ intervalos conterem os } \mu_k \text{'s}) = P(\text{todas hipóteses } H_0 \text{ serem verdadeiras}) \\ = 1 - P(\text{pelo menos uma } H_0 \text{ ser falsa})$$

$$\geq 1 - \sum_{k=1}^p P(H_{0k} \text{ falsa}) = 1 - \underbrace{(\alpha_1 + \alpha_2 + \dots + \alpha_p)}$$

Controle da taxa de erro total independente da estrutura de covariância \Rightarrow muito conservador



$$\left(\bar{Y}_1 - t_{n-1}(\alpha/2p) \sqrt{\frac{s_{11}}{n}}; \bar{Y}_1 + t_{n-1}(\alpha/2p) \sqrt{\frac{s_{11}}{n}} \right)$$

$$\left(\bar{Y}_2 - t_{n-1}(\alpha/2p) \sqrt{\frac{s_{22}}{n}}; \bar{Y}_2 + t_{n-1}(\alpha/2p) \sqrt{\frac{s_{22}}{n}} \right)$$

$$\left(\bar{Y}_p - t_{n-1}(\alpha/2p) \sqrt{\frac{s_{pp}}{n}}; \bar{Y}_p + t_{n-1}(\alpha/2p) \sqrt{\frac{s_{pp}}{n}} \right)$$

O critério de Bonferroni para correção de múltiplos testes é conservador

Bastante utilizado para comparações de subconjuntos de médias ($m < p$)

$$\Rightarrow \alpha/2m$$

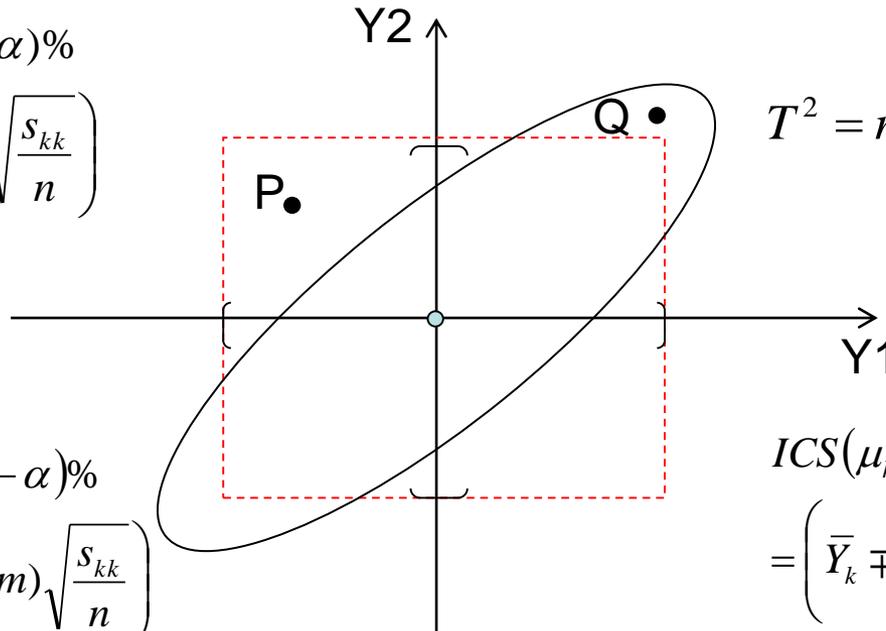
Intervalos e Regiões de Confiança

Intervalos de Confiança Univariados e Simultâneos

Everitt, 2002

$$IC(\mu_k) \text{ a } 100(1-\alpha)\%$$

$$= \left(\bar{Y}_k \mp t_{n-1}(\alpha/2) \sqrt{\frac{s_{kk}}{n}} \right)$$



$$T^2 = n (\bar{Y} - \mu)' S^{-1} (\bar{Y} - \mu)$$

$$ICB(\mu_k) \text{ a } 100(1-\alpha)\%$$

$$= \left(\bar{Y}_k \mp t_{n-1}(\alpha/2m) \sqrt{\frac{s_{kk}}{n}} \right)$$

$$ICS(\mu_k) \text{ a } 100(1-\alpha)\%$$

$$= \left(\bar{Y}_k \mp \sqrt{\frac{(n-1)p}{(n-p)} F_{p,(n-p)}(\alpha) \frac{s_{kk}}{n}} \right)$$

⇒ Compare as análises multivariadas (Região de Confiança, ICS) e univariada (IC e ICB).

⇒ Comente sobre as decisões tomadas para os pontos P e Q sob análises univariadas e multivariadas. Justifique.

Inferência sobre Vetores de Médias de Duas Populações

	Variedade A		Variedade B	
	X11	X12	X21	X22
	5,7	2,1	4,4	1,8
	8,9	1,9	7,5	1,75
	6,2	1,98	5,4	1,78
	5,8	1,92	4,6	1,89
	6,8	2	5,9	1,9
	6,2	2,01		
Média	6,6	1,985	5,56	1,824
Var.	1,42	0,005	1,543	0,0045

$$\bar{D} = \begin{pmatrix} \bar{D}_1 = 6,6 - 5,56 = 1,04 \\ \bar{D}_2 = 1,985 - 1,824 = 0,161 \end{pmatrix}$$

$$S_c = \begin{pmatrix} 1,4745 & -0,0442 \\ -0,0442 & 0,0049 \end{pmatrix}$$

$$\Rightarrow H_0 : \mu_1 - \mu_2 = 0 \quad ; \Sigma_1 = \Sigma_2 = \Sigma$$

$$T^2 = 30,584; \quad p = 0,0027$$



I.C.S. para a diferença entre os grupos em cada variável (ou combinação linear delas)

$$l'(\bar{D}) \pm \sqrt{\frac{(n_1 + n_2 - 2)p}{(n_1 + n_2 - p - 1)} F_{p, (n_1 + n_2 - p - 1)}} \sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2}\right) l' S_c l} \Rightarrow$$

$$\text{Variável X1: } 1,04 \pm \sqrt{(18/8)4,46(1/6 + 1/5)1,47} = 1,04 \pm 2,33 = (-1,29; 3,37)$$

Conclusão?

$$\text{Variável X2: } 0,161 \pm \sqrt{(18/8)4,46(1/6 + 1/5)0,0049} = 0,161 \pm 0,134 = (0,027; 0,295)$$

Inferência sobre Vetores de Médias de Duas Populações

	Variedade A		Variedade B	
	X11	X12	X21	X22
	5,7	2,1	4,4	1,8
	8,9	1,9	7,5	1,75
	6,2	1,98	5,4	1,78
	5,8	1,92	4,6	1,89
	6,8	2	5,9	1,9
	6,2	2,01		
Média	6,6	1,985	5,56	1,824
Var.	1,42	0,005	1,543	0,0045

$$\bar{D} = \begin{pmatrix} \bar{D}_1 = 6,6 - 5,56 = 1,04 \\ \bar{D}_2 = 1,985 - 1,824 = 0,161 \end{pmatrix}$$

$$S_c = \begin{pmatrix} 1,4745 & -0,0442 \\ -0,0442 & 0,0049 \end{pmatrix}$$

$$\Rightarrow H_0 : \mu_1 - \mu_2 = 0 \quad ; \Sigma_1 = \Sigma_2 = \Sigma$$

$$T^2 = 30,584; \quad p = 0,0027$$



I.C. de Bonferroni para a diferença entre os grupos em cada variável

$$ICB(\mu_{1j} - \mu_{2j}) \text{ a } 100(1 - \alpha)\% = (\bar{Y}_{1j} - \bar{Y}_{2j}) \pm t_{n_1+n_2-2}(0,05/4) \sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2}\right) S_{.jj}}$$

$$1,04 \pm 2,685 \sqrt{(1/6 + 1/5)1,47} = 1,04 \pm 1,97 = (-0,93; 3,01)$$

$$0,161 \pm 2,685 \sqrt{(1/6 + 1/5)0,0049} = 0,161 \pm 0,114 = (0,047; 0,275)$$

Conclusão?