MAE 5776

ANÁLISE MULTIVARIADA

Júlia M Pavan Soler

pavan@ime.usp.br

g vimos

MAE5776

Matriz de Dados:
$$Y_{n \times p} = (Y_{ij}) \in \Re^{n \times p}; \quad Y_{ip \times 1} \stackrel{\text{iid}}{\sim} (\mu_{p \times 1}; \Sigma_{p \times p}), i = 1, ..., n$$

Estatísticas descritivas multivariadas:
$$\overline{Y}_{p\times 1}, S_{p\times p}, R_{p\times p}, S_{p\times p}^{-1}$$
 $D_{n\times n} = (d_{ij}^2), d_{Pij}^2, d_{Mij}^2$

Regiões (elipsóides) de Concentração de Observações:

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

Matriz de Dados Aleatórios:
$$Y_{n \times p} \in \Re^{n \times p}$$
; $Y_i \stackrel{\text{iid}}{\sim} N_p \left(\mu_{p \times 1}; \Sigma_{p \times p}\right), i = 1, ..., n$

Caso de Uma única População: H₀:μ=μ₀
 Teste T² de Hotelling

$$T^{2} = n\left(\overline{Y} - \mu\right)' S_{u}^{-1}\left(\overline{Y} - \mu\right) \le c^{2}; \quad c^{2} = \frac{(n-1)p}{(n-p)} F_{p,(n-p)}(\alpha)$$

- Caso de **Duas Populações**: H₀:μ₁=μ₂ ⇔ H₀:μρ=0 Teste T² de Hotelling

$$T^{2} = n(\bar{D} - \mu_{D})' S_{D}^{-1}(\bar{D} - \mu_{D}) \le c^{2}; \quad c^{2} = \frac{(n-1)p}{(n-p)} F_{p,(n-p)}(\alpha)$$
 Dados dependentes

$$T^{2} = \left(\overline{D} - \mu_{D}\right)' \left(S_{c} \left(\frac{1}{n_{1}} + \frac{1}{n_{2}}\right)\right)^{-1} \left(\overline{D} - \mu_{D}\right) \leq c^{2}; \quad c^{2} = \frac{\left(n_{1} + n_{2} - 2\right)p}{\left(n_{1} + n_{2} - p - 1\right)} F_{(p;n_{1} + n_{2} - p - 1)}(\alpha)$$

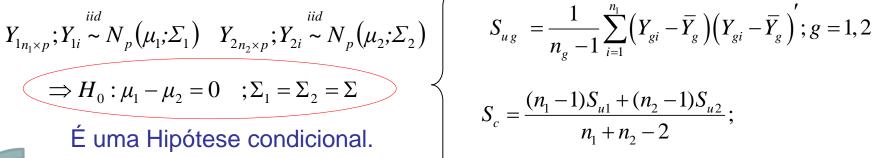
Dados independentes e Homecedásticos, $\Sigma_1 = \Sigma_2$

Teste da Igualdade de Matrizes de Covariância

Comparação de Vetores de Médias - Amostras Independentes e Homocedasticidade

$$Y_{1n_{1}\times p}; Y_{1i} \stackrel{iid}{\sim} N_{p}(\mu_{1}; \Sigma_{1}) \quad Y_{2n_{2}\times p}; Y_{2i} \stackrel{iid}{\sim} N_{p}(\mu_{2}; \Sigma_{2})$$

$$\Rightarrow H_{0}: \mu_{1} - \mu_{2} = 0 \quad ; \Sigma_{1} = \Sigma_{2} = \Sigma$$





Logo a homocedasticidade deve ser verificada.

• Teste M de Box: $\Rightarrow H_0: \Sigma_g = \Sigma; \ \mu_g \in \Re^p$



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

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

É a estatística da Razão de Verossimilhanças (sob N_p) Para p=1 o teste de Box equivale ao teste de Bartlet.

$$\sim \chi^2_{\frac{1}{2}p(p+1)(G-1)}$$

Critério "prático" de heterocedasticidade sugerido em Johnson and Wichern 1992:

$$\sigma_{gjj} = 4\sigma_{g'jj}$$

Inferência sobre Vetores de Médias de Duas Populações

Caso Multivariado - Amostras Independentes — Heterocedasticidade

$$\overline{D}_{p \times 1} = \overline{Y}_1 - \overline{Y}_2 \sim N_p \left(\mu_D = \delta = \mu_1 - \mu_2; \Sigma_{\overline{D}} = \left(\frac{\Sigma_1}{n_1} + \frac{\Sigma_1}{n_2} \right) \right)$$

$$\Rightarrow H_0: \mu_D = \delta_0; \quad \Sigma_g \in \Re^{p \times p}, g = 1, 2$$

 $\Rightarrow H_0: \mu_D = \delta_0; \quad \Sigma_g \in \Re^{p \times p}, g = 1,2$ Hipótese condicional, sob heterocedasticidade (resultados assintóticos) (resultados assintóticos).

$$T^{2} = \left(\overline{D} - \delta_{0}\right)' \left(\frac{S_{1}}{n_{1}} + \frac{S_{2}}{n_{2}}\right)^{-1} \left(\overline{D} - \delta_{0}\right) \stackrel{n_{1} - p \to \infty}{\sim} \chi_{p}^{2}$$



Pesquise:

- Método de Welch para aproximações dos graus de liberdade da estatística de teste sob heterocedasticidade (Kaiser e Bowden, 1983)
- Algoritmo de Behrens-Fisher: Testar a igualdade dos vetores de médias "SEM" suposições sobre as matrizes de covariâncias

Inferência: Vetores de Médias de Duas Populações Correções para Múltiplos Testes

$$Y_{1i} \sim N_p(\mu_1; \Sigma_1) \qquad Y_{2i} \sim N_p(\mu_2; \Sigma_2)$$

Dados Dependentes*1:

$$n_1 = n_2 \implies D_i, \quad i = 1, ..., n$$

$$\bar{D}$$
 $S_{\bar{D}} = S_D \frac{1}{n}$

Dados Independentes*2: $n_1 \Rightarrow \overline{Y}_1 \quad S_1 \quad n_2 \Rightarrow \overline{Y}_2 \quad S_2$

$$n_1 \Rightarrow \overline{Y_1} S_1$$

$$n_2 \Rightarrow \overline{Y}_2 S_2$$

$$\overline{D} \quad S_{\overline{D}} = S_c \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$$

Intervalos de Confiança Simultâneos (para combinações lineares das p variáveis)

$$\Rightarrow ICS \left(\mu_{Dj}\right) a \left(1-\alpha\right) \times 100\% = \left(\overline{D}_{j} \mp \sqrt{\frac{v_{1}}{v_{2}}} F_{v_{1},v_{2}}(\alpha) \sqrt{S_{\overline{D}jj}}\right) *1: v_{1} = (n-1) p, \quad v_{2} = (n-p) \\ *2: v_{1} = (n_{1}+n_{2}-2) p, \quad v_{2} = (n_{1}+n_{2}-p-1)$$

Intervalos de Confiança de Bonferroni (correção para múltiplos testes)

$$\Rightarrow ICB(\mu_{Dj}) a (1-\alpha) \times 100\% = (\overline{Y}_{1j} - \overline{Y}_{2j}) \pm t_{v}(\alpha/2p) \sqrt{S_{\overline{D}jj}}$$
*1: $v = (n-1)$
*2: $v = (n-1)$

Inferência sobre Vetores de Médias Comparações Múltiplas e o Problema de Múliplos Testes

$$\begin{aligned} Y_{1i} &\sim N_p \left(\mu_1; \Sigma_1 \right) &\rightarrow n_1 \\ & & \Rightarrow \begin{cases} H_0: \mu_{1:p \times 1} = \mu_{2:p \times 1} \\ H_1: \mu_{1:p \times 1} \neq \mu_{2:p \times 1} \end{cases} \\ Y_{2i} &\sim N_p \left(\mu_2; \Sigma_2 \right) &\rightarrow n_2 \end{aligned}$$

Estatística T²: Teste simultâneo para as p variáveis a um nível de significância α. Ao rejeitar H₀, há diferença entre os grupos para pelo menos uma variável ou alguma combinação linear delas!

Exemplo: Dados Iris
Comparar Virginica x Setosa
Estatística T² é significante
Neste caso p=4 Variáveis.
Para qual variável os grupos
são diferentes?

$$H_{01}: \mu_{1g=1} = \mu_{1g=2} \qquad \Rightarrow \qquad \alpha_{1}$$

$$H_{02}: \mu_{2g=1} = \mu_{2g=2} \qquad \Rightarrow \qquad \alpha_{2}$$

$$H_{03}: \mu_{3g=1} = \mu_{3g=2} \qquad \Rightarrow \qquad \alpha_{3}$$

$$H_{04}: \mu_{4g=1} = \mu_{4g=2} \qquad \Rightarrow \qquad \alpha_{4}$$

Para testar as 4 hipóteses, qual é o nível de significância global?

Em Estudos envolvendo Comparações Múltiplas é necessário obter um nível de significância corrigido para a Família de Comparações de interesse.

Inferência sobre Vetores de Médias Comparações Múltiplas e o Problema de Múliplos Testes

$$H_{0k} \times H_{1k} \quad k = 1, ..., K$$

 $P(\text{pelo menos uma Rej H}_0) = 1 - P(\text{nenhuma Rej H}_0)$

$$=1-\prod_{l=1}^{K}P(p_{l}>\alpha)=1-(1-\alpha)^{K} \stackrel{\mathsf{K}\to\infty}{\approx} 1$$

```
> B<-10000
> K<-500
> set.seed(1299)
minpval <- replicate(B,
min(runif(K,0,1))<0.01)
table(minpval)
table(minpval)[2]/B</pre>
```

```
> table(minpval)
minpval
FALSE TRUE
    78 9922
> table(minpval)[2]/B
    TRUE
0.9922
```

Quando o número de testes (K), realizados simultaneamente, cresce, é alta a probabilidade de obter resultados significantes ao acaso (Falsos Positivos)

Critérios de Correção do nível de significância individual para garantir um nível de significância coletivo igual a α (em geral, 5%): **Bonferroni, Holm, FDR**

Inferência sobre Vetores de Médias Comparações Múltiplas e o Problema de Múliplos Testes

Rejeitar H _{0j} se:	Método
$p_{(k)} < \alpha / K$	Correção de <u>Bonferroni</u> para múltiplos testes
$p_{(j)} < \alpha/(K-j+1)$	Correção de Holm (Controle "forte" da taxa de erro para os múltiplos testes)
$p_{(j)} < j\alpha/K$	Correção FDR (Taxa de Falsa Descoberta): Benjamini-Hochberg Controle menos conservador da taxa de erro para os múltiplos testes)

K: número total de testes α: nível de significância global fixado

 $p_{(j)}$: nível descritivo (p-valor) ordenado, $p_{(1)} \le p_{(2)} \le ... \le p_{(K)}$

Oficina no R Comparação de Vetores de Médias de Duas Populações

Exemplo 1: Gerar dados de duas populações N₂

Exemplo 2: Analisar os dados de duas espécies de Moscas (Johnson and Wichern, 1992)