



SME0822 Análise Multivariada e Aprendizado Não-Supervisionado

Aula 6a: **Análise de Variância Multivariada (MANOVA)**

Prof. Cibeles Russo

cibele@icmc.usp.br

<http://www.icmc.usp.br/~cibele>

Baseado em Johnson, R. A., & Wichern, D. W. (2007). Applied Multivariate Statistical Analysis. Prentice Hall.

Comparação de médias em várias populações independentes

Sejam

- $\tilde{X}_{11}, \dots, \tilde{X}_{1n_1}$ a.a. de uma população com distribuição normal com $E(\tilde{X}_{1j}) = \underline{\mu}_1$ para $j = 1, \dots, n$ e $\text{Var}(\tilde{X}_{1j}) = \Sigma$,
- $\tilde{X}_{21}, \dots, \tilde{X}_{2n_2}$ a.a. de uma população com distribuição normal com $E(\tilde{X}_{2j}) = \underline{\mu}_2$ para $j = 1, \dots, n$ e $\text{Var}(\tilde{X}_{2j}) = \Sigma$
- ...
- $\tilde{X}_{g1}, \dots, \tilde{X}_{gn_g}$ a.a. de uma população com distribuição normal com $E(\tilde{X}_{gj}) = \underline{\mu}_g$ para $j = 1, \dots, n_g$ e $\text{Var}(\tilde{X}_{gj}) = \Sigma$

supondo que a todas as populações são independentes entre si.

Comparação de médias em várias populações independentes

Sejam

- $\tilde{X}_{11}, \dots, \tilde{X}_{1n_1}$ a.a. de uma população com distribuição normal com $E(\tilde{X}_{1j}) = \underline{\mu}_1$ para $j = 1, \dots, n$ e $\text{Var}(\tilde{X}_{1j}) = \Sigma$,
- $\tilde{X}_{21}, \dots, \tilde{X}_{2n_2}$ a.a. de uma população com distribuição normal com $E(\tilde{X}_{2j}) = \underline{\mu}_2$ para $j = 1, \dots, n$ e $\text{Var}(\tilde{X}_{2j}) = \Sigma$
- ...
- $\tilde{X}_{g1}, \dots, \tilde{X}_{gn_g}$ a.a. de uma população com distribuição normal com $E(\tilde{X}_{gj}) = \underline{\mu}_g$ para $j = 1, \dots, n_g$ e $\text{Var}(\tilde{X}_{gj}) = \Sigma$

supondo que a todas as populações são independentes entre si.

Comparação de médias em várias populações independentes

Sejam

- $\tilde{X}_{11}, \dots, \tilde{X}_{1n_1}$ a.a. de uma população com distribuição normal com $E(\tilde{X}_{1j}) = \underline{\mu}_1$ para $j = 1, \dots, n$ e $\text{Var}(\tilde{X}_{1j}) = \Sigma$,
- $\tilde{X}_{21}, \dots, \tilde{X}_{2n_2}$ a.a. de uma população com distribuição normal com $E(\tilde{X}_{2j}) = \underline{\mu}_2$ para $j = 1, \dots, n$ e $\text{Var}(\tilde{X}_{2j}) = \Sigma$
- ...
- $\tilde{X}_{g1}, \dots, \tilde{X}_{gn_g}$ a.a. de uma população com distribuição normal com $E(\tilde{X}_{gj}) = \underline{\mu}_g$ para $j = 1, \dots, n_g$ e $\text{Var}(\tilde{X}_{gj}) = \Sigma$

supondo que a **todas as populações são independentes entre si.**

Comparação de médias em várias populações independentes

Deseja-se avaliar as hipóteses

$$H_0 : \underline{\mu}_1 = \underline{\mu}_2 = \dots \underline{\mu}_g = \mu \text{ contra}$$

$$H_1 : \text{pelo menos um } \underline{\mu}_i \text{ diferente}$$

Para isso, vamos considerar a reparametrização

$$\underline{\mu}_k = \underline{\mu} + \underline{\tau}_k, \quad k = 1, \dots, g,$$

e então avaliar se

$$H_0 : \tau_1 = \tau_2 = \dots \tau_g = 0 \text{ contra}$$

$$H_1 : \text{pelo menos um } \tau \text{ diferente de } 0$$

Comparação de médias em várias populações independentes

Deseja-se avaliar as hipóteses

$$H_0 : \underline{\mu}_1 = \underline{\mu}_2 = \dots \underline{\mu}_g = \mu \text{ contra}$$

$$H_1 : \text{pelo menos um } \underline{\mu}_i \text{ diferente}$$

Para isso, vamos considerar a reparametrização

$$\underline{\mu}_k = \underline{\mu} + \underline{\tau}_k, \quad k = 1, \dots, g,$$

e então avaliar se

$$H_0 : \tau_1 = \tau_2 = \dots \tau_g = 0 \text{ contra}$$

$$H_1 : \text{pelo menos um } \tau \text{ diferente de } 0$$

Modelo de ANOVA multivariada - MANOVA

Seja o modelo de ANOVA multivariada (MANOVA)

$$\underline{X}_{kj} = \underline{\mu} + \underline{\tau}_k + \underline{\epsilon}_{kj},$$

para $j = 1, \dots, n_k, k = 1, \dots, g$.

Suposições:

- $\underline{\epsilon}_{kj} \stackrel{ind}{\sim} N_p(\underline{0}, \Sigma)$,
- $\underline{\mu}$ é a média geral,
- $\underline{\tau}_k$ é o efeito do k -ésimo tratamento,
- $\sum_{k=1}^g n_k \underline{\tau}_k = \underline{0}$.

Modelo de ANOVA multivariada - MANOVA

Seja o modelo de ANOVA multivariada (MANOVA)

$$\underline{X}_{kj} = \underline{\mu} + \underline{\tau}_k + \underline{\epsilon}_{kj},$$

para $j = 1, \dots, n_k, k = 1, \dots, g$.

Suposições:

- $\underline{\epsilon}_{kj} \stackrel{ind}{\sim} N_p(\underline{0}, \Sigma)$,
- $\underline{\mu}$ é a média geral,
- $\underline{\tau}_k$ é o efeito do k -ésimo tratamento,
- $\sum_{k=1}^g n_k \underline{\tau}_k = \underline{0}$.

Modelo de ANOVA multivariada - MANOVA

Considerando os vetores observados, podemos escrever

$$\underbrace{\tilde{X}_{kj}}_{\text{observação}} = \underbrace{\bar{X}}_{\text{média geral}} + \underbrace{(\bar{X}_k - \bar{X})}_{\text{efeito estimado do tratamento } k} + \underbrace{(\tilde{X}_{kj} - \bar{X}_k)}_{\text{resíduo}},$$

para $j = 1, \dots, n_k$ e $k = 1, \dots, g$.

Modelo de ANOVA multivariada - MANOVA

Como no caso univariado, queremos decompor a variabilidade dos dados em torno da média em variabilidade intra e entre tratamentos. Temos que

$$\begin{aligned}(\underline{X}_{kj} - \bar{\underline{X}})(\underline{X}_{kj} - \bar{\underline{X}})^T &= \\&= \left[(\underline{X}_{kj} - \bar{\underline{X}}_k) + (\bar{\underline{X}}_k - \bar{\underline{X}}) \right] \left[(\underline{X}_{kj} - \bar{\underline{X}}_k) + (\bar{\underline{X}}_k - \bar{\underline{X}}) \right]^T = \\&= (\underline{X}_{kj} - \bar{\underline{X}}_k)(\underline{X}_{kj} - \bar{\underline{X}}_k)^T + (\underline{X}_{kj} - \bar{\underline{X}}_k)(\bar{\underline{X}}_k - \bar{\underline{X}})^T + \\&\quad + (\bar{\underline{X}}_k - \bar{\underline{X}})(\underline{X}_{kj} - \bar{\underline{X}}_k)^T + (\bar{\underline{X}}_k - \bar{\underline{X}})(\bar{\underline{X}}_k - \bar{\underline{X}})^T\end{aligned}$$

Modelo de ANOVA multivariada - MANOVA

Como no caso univariado, queremos decompor a variabilidade dos dados em torno da média em variabilidade intra e entre tratamentos. Temos que

$$\begin{aligned}(\underline{X}_{kj} - \bar{X})(\underline{X}_{kj} - \bar{X})^T &= \\&= \left[(\underline{X}_{kj} - \bar{X}_k) + (\bar{X}_k - \bar{X}) \right] \left[(\underline{X}_{kj} - \bar{X}_k) + (\bar{X}_k - \bar{X}) \right]^T = \\&= (\underline{X}_{kj} - \bar{X}_k)(\underline{X}_{kj} - \bar{X}_k)^T + (\underline{X}_{kj} - \bar{X}_k)(\bar{X}_k - \bar{X})^T + \\&\quad + (\bar{X}_k - \bar{X})(\underline{X}_{kj} - \bar{X}_k)^T + (\bar{X}_k - \bar{X})(\bar{X}_k - \bar{X})^T\end{aligned}$$

Modelo de ANOVA multivariada - MANOVA

Como no caso univariado, queremos decompor a variabilidade dos dados em torno da média em variabilidade intra e entre tratamentos. Temos que

$$\begin{aligned}(\underline{X}_{kj} - \bar{X})(\underline{X}_{kj} - \bar{X})^T &= \\&= \left[(\underline{X}_{kj} - \bar{X}_k) + (\bar{X}_k - \bar{X}) \right] \left[(\underline{X}_{kj} - \bar{X}_k) + (\bar{X}_k - \bar{X}) \right]^T = \\&= (\underline{X}_{kj} - \bar{X}_k)(\underline{X}_{kj} - \bar{X}_k)^T + (\underline{X}_{kj} - \bar{X}_k)(\bar{X}_k - \bar{X})^T + \\&\quad + (\bar{X}_k - \bar{X})(\underline{X}_{kj} - \bar{X}_k)^T + (\bar{X}_k - \bar{X})(\bar{X}_k - \bar{X})^T\end{aligned}$$

Modelo de ANOVA multivariada - MANOVA

Vamos fazer a soma em j , mas

$$\begin{aligned} & \sum_{j=1}^{n_k} [(\underline{X}_{kj} - \bar{\underline{X}}_k)(\bar{\underline{X}}_k - \bar{\underline{X}})^T] = \\ & \sum_{j=1}^{n_k} [(\underline{X}_{kj} - \bar{\underline{X}}_k)] (\bar{\underline{X}}_k - \bar{\underline{X}})^T = \\ & (\sum_{j=1}^{n_k} \underline{X}_{kj} - n_k \bar{\underline{X}}_k)(\bar{\underline{X}}_k - \bar{\underline{X}})^T = \\ & (\sum_{j=1}^{n_k} \underline{X}_{kj} - \sum_{j=1}^{n_k} \underline{X}_{kj})(\bar{\underline{X}}_k - \bar{\underline{X}})^T = 0 \end{aligned}$$

Idem para

$$\sum_{j=1}^{n_k} (\bar{\underline{X}}_k - \bar{\underline{X}})(\underline{X}_{kj} - \bar{\underline{X}}_k)^T = 0$$

Modelo de ANOVA multivariada - MANOVA

Então, fazendo a soma em j e em k , ficamos com

$$\begin{aligned} \sum_{k=1}^g \sum_{j=1}^{n_k} (\mathbf{X}_{kj} - \bar{\mathbf{X}})(\mathbf{X}_{kj} - \bar{\mathbf{X}})^T &= \\ &= \sum_{k=1}^g n_k (\bar{\mathbf{X}}_k - \bar{\mathbf{X}})(\bar{\mathbf{X}}_k - \bar{\mathbf{X}})^T + \sum_{k=1}^g \sum_{j=1}^{n_k} (\mathbf{X}_{kj} - \bar{\mathbf{X}}_k)(\mathbf{X}_{kj} - \bar{\mathbf{X}}_k)^T \end{aligned}$$

Ou seja,

$T = B + W$, com

- T : soma de quadrados e produtos cruzados **total**
- B : soma de quadrados e produtos cruzados **entre** tratamentos
- W : soma de quadrados e produtos cruzados **intra** tratamento (**dentro**)

Modelo de ANOVA multivariada - MANOVA

Consideramos as somas de quadrados e produtos cruzados

$$T = \sum_{k=1}^g \sum_{j=1}^{n_k} (\tilde{X}_{kj} - \tilde{\bar{X}})(\tilde{X}_{kj} - \tilde{\bar{X}})^{\top} \text{ (total)}$$

$$B = \sum_{k=1}^g n_k (\tilde{\bar{X}}_k - \tilde{\bar{X}})(\tilde{\bar{X}}_k - \tilde{\bar{X}})^{\top} \text{ (entre)}$$

$$W = \sum_{k=1}^g \sum_{j=1}^{n_k} (\tilde{X}_{kj} - \tilde{\bar{X}}_k)(\tilde{X}_{kj} - \tilde{\bar{X}}_k)^{\top} \text{ (dentro)}$$

Modelo de ANOVA multivariada - MANOVA

Temos

$$T = B + W$$

$$SQT = SQTrat + SQRes$$

Como supomos que todos os grupos têm a mesma variância, podemos considerar como estimativa de Σ :

$$S_{pooled} = W = \sum_{k=1}^g \sum_{j=1}^{n_k} (\tilde{X}_{kj} - \bar{\tilde{X}}_k)(\tilde{X}_{kj} - \bar{\tilde{X}}_k)^T$$

$$W = (n_1 - 1)S_1 + (n_2 - 1)S_2 + \dots + (n_g - 1)S_g.$$

Tabela MANOVA

Tabela ANOVA multivariada

Fonte de variação	Somas de quadrados e produtos cruzados	graus de liberdade
Tratamento	B	$g - 1$
Resíduo	W	$N - g$
Total	T	$N - 1$

em que $N = \sum_{k=1}^g n_k$. Assim, para avaliar

$H_0 : \underline{\mu}_1 = \underline{\mu}_2 = \dots = \underline{\mu}_g = \underline{\mu}$ contra H_1 : pelo menos um $\underline{\mu}_i$ diferente

ou seja

$H_0 : \tau_1 = \tau_2 = \dots = \tau_g = \mathcal{Q}$ contra H_1 : pelo menos um τ diferente de \mathcal{Q} .

Tabela MANOVA

Queremos verificar se W é pequeno em relação a T .

Mas como avaliar, se são matrizes?

Como são matrizes, podemos considerar

Lambda de Wilks (variância generalizada):

$$\Lambda^* = \frac{|W|}{|W + B|}.$$

A distribuição de Λ^* foi estudada em diferentes casos e existem tabelas com os resultados (ver Tabela 6.3 de Johnson & Wichern).

Tabela MANOVA

Queremos verificar se W é pequeno em relação a T .

Mas como avaliar, se são matrizes?

Como são matrizes, podemos considerar

Lambda de Wilks (variância generalizada):

$$\Lambda^* = \frac{|W|}{|W + B|}.$$

A distribuição de Λ^* foi estudada em diferentes casos e existem tabelas com os resultados (ver Tabela 6.3 de Johnson & Wichern).

Lambda de Wilks

Exemplo: Lambda de Wilks

Se $p = 2$ e $g \geq 2$, então

$$\left(\frac{\sum_{k=1}^g n_k - g - 1}{g - 1} \right) \left(\frac{1 - \sqrt{\Lambda^*}}{\sqrt{\Lambda^*}} \right) \sim F_{2(g-1), 2 \sum n_k - g - 1}$$

Bartlett mostrou que se H_0 é verdadeira e $\sum_{k=1}^g = N$ é grande, então

$$- \left(n - 1 - \frac{p + g}{2} \right) \log(\Lambda^*) \sim \chi_{p(g-1)}^2.$$

MANOVA

Exemplo: Uma base de dados com 150 observações e 5 variáveis referentes a crânios egípcios de cinco épocas:

- época: a época à qual o crânio foi atribuído, um fator ordenado com níveis c4000BC, c3300BC, c1850BC, c200BC e cAD150, onde os anos são dados apenas aproximadamente, é claro.
- MB: largura máxima do crânio.
- bh: altura basibregmática do crânio.
- bl: comprimento basialveolar do crânio.
- nh: altura nasal do crânio.

Mais informações em:

<http://friendly.github.io/heplots/reference/Skulls.html>