



SME0822 Análise Multivariada e Aprendizado Não-Supervisionado

Aula 5a: **Testes de Hipóteses sobre a média**

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.

Testes de hipóteses para a média

Seja $\underline{X}_1, \dots, \underline{X}_n$ uma amostra aleatória de uma distribuição normal p-variada com vetor de médias $\underline{\mu}$ e matriz de variâncias e covariâncias Σ .
Sejam $\overline{\underline{X}}$ e S o vetor de médias amostrais e a matriz de variâncias e covariâncias amostrais.

Queremos avaliar se

$$\begin{aligned} H_0 : \underline{\mu} &= \underline{\mu}_0 \text{ contra} \\ H_1 : \underline{\mu} &\neq \underline{\mu}_0, \end{aligned}$$

Testes de hipóteses para a média

Relembramos o resultado anterior

Resultado

- 1 $\bar{\underset{\sim}{X}} \sim N_p \left(\underset{\sim}{\mu}, \frac{\Sigma}{n} \right).$
- 2 $(n-1)S \sim Wishart(n-1).$
- 3 $\bar{\underset{\sim}{X}}$ e S são independentes.

Além disso, sob H_0 ,

$$T^2 = \sqrt{n}(\bar{\underset{\sim}{X}} - \underset{\sim}{\mu}_0)^\top \left(\frac{(n-1)S^{-1}}{n-1} \right) \sqrt{n}(\bar{\underset{\sim}{X}} - \underset{\sim}{\mu}_0) \sim \frac{(n-1)p}{n-p} F_{p, n-p}$$

Testes de hipóteses para a média

A quantidade

$$T^2 = n(\bar{\tilde{X}} - \underline{\mu}_0)^\top S^{-1}(\bar{\tilde{X}} - \underline{\mu}_0) \sim \frac{(n-1)p}{n-p} F_{p, n-p}$$

é conhecida como a **Estatística T^2 de Hotelling**.

Testes de hipóteses para a média

Assim, rejeitamos H_0 a um nível de significância α se

$$T_{obs}^2 = n(\bar{\tilde{X}} - \underline{\mu}_0)^\top S^{-1}(\bar{\tilde{X}} - \underline{\mu}_0) > \frac{(n-1)p}{n-p} q_{F_{p,n-p},\alpha}$$

em que $q_{F_{p,n-p},\alpha}$ é o quantil α -superior de uma distribuição $F_{p,n-p}$.

Região de Confiança

Seja $\underline{X}_1, \dots, \underline{X}_n$ uma amostra aleatória de uma distribuição $N_p(\underline{\mu}, \Sigma)$.

Uma **região com** $100(1 - \alpha)\%$ **de confiança** para $\underline{\mu}$ é dada pelo elipsóide determinado pelos valores de $\underline{\mu}^*$ que satisfazem

$$n(\bar{\underline{x}} - \underline{\mu}^*)^\top S^{-1}(\bar{\underline{x}} - \underline{\mu}^*) \leq \frac{(n-1)p}{n-p} F_{p, n-p, \alpha} = c$$

em que $\bar{\underline{x}}$ e S são, respectivamente, a média e a matriz de variâncias e covariâncias amostrais observadas.

Assim, para verificar se $\underline{\mu}_0$ está dentro da região de confiança, verificamos se

$$n(\bar{\underline{x}} - \underline{\mu}_0)^\top S^{-1}(\bar{\underline{x}} - \underline{\mu}_0) \leq \frac{(n-1)p}{n-p} F_{p, n-p, \alpha}$$

Se for verdadeiro, $\underline{\mu}_0$ está dentro da região de confiança.

Região de Confiança

Seja $\underline{X}_1, \dots, \underline{X}_n$ uma amostra aleatória de uma distribuição $N_p(\underline{\mu}, \Sigma)$.

Uma **região com** $100(1 - \alpha)\%$ **de confiança** para $\underline{\mu}$ é dada pelo elipsóide determinado pelos valores de $\underline{\mu}^*$ que satisfazem

$$n(\bar{\underline{x}} - \underline{\mu}^*)^\top S^{-1}(\bar{\underline{x}} - \underline{\mu}^*) \leq \frac{(n-1)p}{n-p} F_{p, n-p, \alpha} = c$$

em que $\bar{\underline{x}}$ e S são, respectivamente, a média e a matriz de variâncias e covariâncias amostrais observadas.

Assim, para verificar se $\underline{\mu}_0$ está dentro da região de confiança, verificamos se

$$n(\bar{\underline{x}} - \underline{\mu}_0)^\top S^{-1}(\bar{\underline{x}} - \underline{\mu}_0) \leq \frac{(n-1)p}{n-p} F_{p, n-p, \alpha}$$

Se for verdadeiro, $\underline{\mu}_0$ está dentro da região de confiança.

Comparação de médias multivariadas

Considere o problema de avaliar a igualdade de médias multidimensionais de diferentes populações.

Para isso, vamos considerar três casos básicos:

- 1 Comparações pareadas ou medidas repetidas
- 2 Comparação de médias em duas populações independentes
- 3 Comparação de médias em mais que duas populações independentes

Comparação de médias multivariadas

Considere o problema de avaliar a igualdade de médias multidimensionais de diferentes populações.

Para isso, vamos considerar três casos básicos:

- 1 Comparações pareadas ou medidas repetidas
- 2 Comparação de médias em duas populações independentes
- 3 Comparação de médias em mais que duas populações independentes

Comparação de médias multivariadas

Considere o problema de avaliar a igualdade de médias multidimensionais de diferentes populações.

Para isso, vamos considerar três casos básicos:

- 1 Comparações pareadas ou medidas repetidas
- 2 Comparação de médias em duas populações independentes
- 3 Comparação de médias em mais que duas populações independentes

Comparação de médias multivariadas

Considere o problema de avaliar a igualdade de médias multidimensionais de diferentes populações.

Para isso, vamos considerar três casos básicos:

- 1 Comparações pareadas ou medidas repetidas
- 2 Comparação de médias em duas populações independentes
- 3 Comparação de médias em mais que duas populações independentes

Comparações pareadas

Sejam

- $\underline{X}_{11}, \dots, \underline{X}_{1n}$ vetores aleatórios $p \times 1$ referentes a uma população normal multivariada **antes** de um tratamento com $E(\underline{X}_{1j}) = \underline{\mu}_1$ para $j = 1, \dots, n$,
- $\underline{X}_{21}, \dots, \underline{X}_{2n}$ vetores aleatórios $p \times 1$ referentes a uma população normal multivariada **após** de um tratamento com $E(\underline{X}_{2j}) = \underline{\mu}_2$ para $j = 1, \dots, n$,

sendo que $\underline{X}_{11}, \dots, \underline{X}_{1n}$ e $\underline{X}_{21}, \dots, \underline{X}_{2n}$ são amostras aleatórias de uma mesma população em diferentes situações, em que \underline{X}_{1j} e \underline{X}_{2j} são correlacionadas (por exemplo, vetores aleatórios de medições antes e após um tratamento).

Comparações pareadas

Sejam

- $\underline{X}_{11}, \dots, \underline{X}_{1n}$ vetores aleatórios $p \times 1$ referentes a uma população normal multivariada **antes** de um tratamento com $E(\underline{X}_{1j}) = \underline{\mu}_1$ para $j = 1, \dots, n$,
- $\underline{X}_{21}, \dots, \underline{X}_{2n}$ vetores aleatórios $p \times 1$ referentes a uma população normal multivariada **após** de um tratamento com $E(\underline{X}_{2j}) = \underline{\mu}_2$ para $j = 1, \dots, n$,

sendo que $\underline{X}_{11}, \dots, \underline{X}_{1n}$ e $\underline{X}_{21}, \dots, \underline{X}_{2n}$ são amostras aleatórias de uma mesma população em diferentes situações, em que \underline{X}_{1j} e \underline{X}_{2j} são correlacionadas (por exemplo, vetores aleatórios de medições antes e após um tratamento).

Comparações pareadas

Sejam

- $\underline{X}_{11}, \dots, \underline{X}_{1n}$ vetores aleatórios $p \times 1$ referentes a uma população normal multivariada **antes** de um tratamento com $E(\underline{X}_{1j}) = \underline{\mu}_1$ para $j = 1, \dots, n$,
- $\underline{X}_{21}, \dots, \underline{X}_{2n}$ vetores aleatórios $p \times 1$ referentes a uma população normal multivariada **após** de um tratamento com $E(\underline{X}_{2j}) = \underline{\mu}_2$ para $j = 1, \dots, n$,

sendo que $\underline{X}_{11}, \dots, \underline{X}_{1n}$ e $\underline{X}_{21}, \dots, \underline{X}_{2n}$ são amostras aleatórias de uma mesma população em diferentes situações, em que \underline{X}_{1j} e \underline{X}_{2j} são correlacionadas (por exemplo, vetores aleatórios de medições antes e após um tratamento).

Comparações pareadas

Sejam $\underline{\mu}_1$ e $\underline{\mu}_2$ os vetores de médias em situações 1 e 2, respectivamente. Deseja-se testar se não há diferença entre as situações 1 e 2 para verificar, por exemplo, que o tratamento não produz nenhum efeito, ou seja, se

$$\underline{\mu}_1 = \underline{\mu}_2.$$

Para avaliar as hipóteses

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

$$H_1 : \underline{\mu}_1 \neq \underline{\mu}_2$$

vamos considerar as diferenças

$$\underline{D}_j = \underline{X}_{1j} - \underline{X}_{2j}.$$

Assim, $\underline{D}_1, \dots, \underline{D}_n$ são i.i.d e $\underline{D}_j \sim N(\underline{\mu}_D, \Sigma_D)$.

Comparações pareadas

Então, avaliamos se

$$H_0 : \underline{\mu}_D = \underline{0} \text{ contra}$$

$$H_1 : \underline{\mu}_D \neq \underline{0}$$

com a estatística T^2 de Hotelling:

$$T^2 = n(\bar{\underline{D}} - \underline{0})^\top S_D^{-1}(\bar{\underline{D}} - \underline{0}) \overset{\text{sob } H_0}{\sim} \frac{(n-1)p}{n-p} F_{p, n-p},$$

em que $\bar{\underline{D}}$ e S_D são o vetor de médias e a matriz de variâncias e covariâncias amostrais de \underline{D} .

Comparações pareadas

Um teste análogo poderia ser desenvolvido para avaliar

$$H_0 : \underline{\mu}_D = \underline{\mu}_{D0} \text{ contra}$$

$$H_1 : \underline{\mu}_D \neq \underline{\mu}_{D0}$$

com a estatística T^2 de Hotelling:

$$T^2 = n(\bar{\underline{D}} - \underline{\mu}_{D0})^\top S_D^{-1}(\bar{\underline{D}} - \underline{\mu}_{D0}) \stackrel{\text{sob } H_0}{\sim} \frac{(n-1)p}{n-p} F_{p, n-p},$$

em que $\bar{\underline{D}}$ e S_D são o vetor de médias e a matriz de variâncias e covariâncias amostrais de \underline{D} .

Comparações pareadas

A região de confiança, com nível de confiança $100(1 - \alpha)\%$ nesse caso seria

$$\{\underline{\mu}_D^*; n(\bar{\underline{d}} - \underline{\mu}_D^*)^\top S_D^{-1}(\bar{\underline{d}} - \underline{\mu}_D^*) \leq \frac{(n-1)p}{n-p} q_{F_{p, n-p}, \alpha}\}$$

em que $\bar{\underline{D}}$ e S_D são o vetor de médias e a matriz de variâncias e covariâncias amostrais de \underline{D} .

Comparação de médias de duas populações independentes

Sejam

- $\underline{X}_{11}, \dots, \underline{X}_{1n_1}$ vetores aleatórios $p \times 1$ referentes a uma população com $E(\underline{X}_{1j}) = \underline{\mu}_1$ para $j = 1, \dots, n_1$,
- $\underline{X}_{21}, \dots, \underline{X}_{2n_2}$ vetores aleatórios $p \times 1$ referentes a uma população com $E(\underline{X}_{2j}) = \underline{\mu}_2$ para $j = 1, \dots, n_2$,

supondo que a população 1 é independente da população 2.

Comparações pareadas

Temos

| | Média amostral | Matriz de covariâncias amostrais |
|-------------|---|--|
| População 1 | $\bar{X}_1 = \frac{1}{n_1} \sum_{j=1}^{n_1} X_{1j}$ | $S_1 = \frac{1}{n_1} \sum_{j=1}^{n_1} (X_{1j} - \bar{X}_1)(X_{1j} - \bar{X}_1)^\top$ |
| População 2 | $\bar{X}_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} X_{2j}$ | $S_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} (X_{2j} - \bar{X}_2)(X_{2j} - \bar{X}_2)^\top$ |

Comparação de médias de duas populações independentes

Suposições adicionais

- 1 Ambas as populações têm distribuição normal multivariada
- 2 $\Sigma_1 = \Sigma_2 = \Sigma$

Deseja-se avaliar as hipóteses

$$H_0 : \mu_1 = \mu_2 \text{ contra}$$

$$H_1 : \mu_1 \neq \mu_2$$

Comparação de médias de duas populações independentes

Primeiramente, consideramos um estimador para Σ , por exemplo

$$S = S_{pooled} = \frac{\sum_{j=1}^{n_1} (\tilde{X}_{1j} - \bar{\tilde{X}}_1)(\tilde{X}_{1j} - \bar{\tilde{X}}_1)^\top + \sum_{j=1}^{n_2} (\tilde{X}_{2j} - \bar{\tilde{X}}_2)(\tilde{X}_{2j} - \bar{\tilde{X}}_2)^\top}{n_1 + n_2 - 2}$$

ou seja

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

Comparação de médias de duas populações independentes

Então, reescrevemos as hipóteses de interesse na forma mais geral

$$H_0 : \mu_1 - \mu_2 = \delta_0 \text{ contra}$$

$$H_1 : \mu_1 - \mu_2 \neq \delta_0.$$

e rejeitamos H_0 ao nível de significância α se

$$T_{obs}^2 = (\bar{x}_1 - \bar{x}_2 - \delta_0) \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) S \right]^{-1} (\bar{x}_1 - \bar{x}_2 - \delta_0) > c^2$$

$$\text{com } c^2 = \frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - p - 1} q_{F_{p, n_1 + n_2 - p - 1, \alpha}}.$$

Comparação de médias de duas populações independentes

Note que

- $E(\bar{X}_1 - \bar{X}_2) = E(\bar{X}_1) - E(\bar{X}_2) = \mu_1 - \mu_2$
- $\text{Var}(\bar{X}_1 - \bar{X}_2) = \text{Var}(\bar{X}_1) + \text{Var}(\bar{X}_2)$ pois \bar{X}_1 é independente de \bar{X}_2

Logo, como as populações originais tem distribuição normal multivariada,

$$\bar{X}_1 - \bar{X}_2 \sim N(\mu_1 - \mu_2, \left(\frac{1}{n_1} + \frac{1}{n_2}\right) \Sigma)$$

Comparação de médias de duas populações independentes

Temos:

$$(n_1 - 1)S_1 \sim \text{Wishart}(n_1 - 1, \Sigma)$$

$$(n_2 - 1)S_2 \sim \text{Wishart}(n_2 - 1, \Sigma)$$

Como as observações da população 1 são independentes das da população 2, S_1 é independente de S_2 . Uma propriedade garante que $(n_1 - 1)S_1 + (n_2 - 1)S_2 \sim \text{Wishart}(n_1 + n_2 - 2, \Sigma)$.

Comparação de médias em duas populações independentes

Então

$$T^2 = \left[\bar{\tilde{X}}_1 - \bar{\tilde{X}}_2 - (\underline{\mu}_1 - \underline{\mu}_2) \right]^T \frac{S_{pooled}^{-1}}{n_1 + n_2 - 2} \left[\bar{\tilde{X}}_1 - \bar{\tilde{X}}_2 - (\underline{\mu}_1 - \underline{\mu}_2) \right],$$

o que é novamente o produto de uma v.a. normal multivariada pela inversa de uma v.a. Wishart dividida pelos seus g.l. e uma normal multivariada.

Comparação de médias em duas populações independentes

Se as hipóteses de interesse são

$$H_0 : \mu_1 = \mu_2 \text{ contra } H_1 : \mu_1 \neq \mu_2$$

então a estatística se simplifica, sob H_0 , em

$$T^2 = (\bar{X}_1 - \bar{X}_2)^\top \frac{S_{pooled}^{-1}}{n_1 + n_2 - 2} (\bar{X}_1 - \bar{X}_2) \sim \frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1}.$$