

ME - 731 Análise Multivariada

Segundo semestre de 2009

Lista de Exercícios VI

Entrega: Exercícios 2 e 4 até o dia 15/12 às 10h na sala do Professor (210 IMECC)

Obs1: Não é necessário digitar a resolução da lista, os exercícios podem ser entregues feitos à mão.

Obs2: Os exercícios seguem, a menos que mencionado o contrário, as notações definidas nas notas de aulas disponíveis no site e em classe

Exercícios

1. Na análise de correlações canônicas considere um vetor aleatório $\mathbf{X}_{(p+q)}$, $p \leq q$, com as devidas partições do vetor de médias $\boldsymbol{\mu} = \mathcal{E}(\mathbf{X})$, da matriz de covariâncias $\boldsymbol{\Sigma} = Cov(\mathbf{X})$ e da matriz de correlações $\boldsymbol{\rho} = Corre(\mathbf{X})$, definidas em classe. Considere o vetor aleatório padronizado \mathbf{Z} com as partições correspondentes.

Sejam $\mathbf{U}_{(p \times 1)}^{(\mathbf{Z})} = (U_1^{(\mathbf{Z})}, \dots, U_p^{(\mathbf{Z})})'$ as $\mathbf{V}_{(q \times 1)}^{(\mathbf{Z})} = (V_1^{(\mathbf{Z})}, \dots, V_q^{(\mathbf{Z})})'$ as variáveis canônicas obtidas a partir do vetor \mathbf{Z} . Responda os itens:

- a) Calcule $\mathcal{E}(\mathbf{U}^{(\mathbf{Z})})$, $\mathcal{E}(\mathbf{V}^{(\mathbf{Z})})$, $Cov(\mathbf{U}^{(\mathbf{Z})})$, $Cov(\mathbf{V}^{(\mathbf{Z})})$, $Corr(\mathbf{U}^{(\mathbf{Z})})$ e $Corr(\mathbf{V}^{(\mathbf{Z})})$.
 - b) Prove que: $Cov(\mathbf{z}^{(1)}, \mathbf{U}^{(\mathbf{z})}) = Corr(\mathbf{z}^{(1)}, \mathbf{U}^{(\mathbf{z})}) = \mathbf{A}^{(\mathbf{z})} \boldsymbol{\rho}_{11}$,
 $Cov(\mathbf{z}^{(2)}, \mathbf{V}^{(\mathbf{z})}) = Corr(\mathbf{z}^{(2)}, \mathbf{V}^{(\mathbf{z})}) = \mathbf{B}^{(\mathbf{z})} \boldsymbol{\rho}_{22}$.
 - c) Prove que: $Cov(\mathbf{Z}^{(1)}) = \mathbf{A}^{(\mathbf{Z})-1} \mathbf{A}'^{(\mathbf{Z})-1}$, $Cov(\mathbf{Z}^{(2)}) = \mathbf{B}^{(\mathbf{Z})-1} \mathbf{B}'^{(\mathbf{Z})-1}$.
 - d) Suponha que $\mathbf{X}_{(p+q)} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Qual a distribuição conjunta do vetor $(\mathbf{U}^{(\mathbf{Z})}, \mathbf{V}^{(\mathbf{Z})})$? As componentes de $\mathbf{U}^{(\mathbf{Z})}$ são independentes? O são as componentes de $\mathbf{V}^{(\mathbf{Z})}$? Com relação às componentes de $\mathbf{U}^{(\mathbf{Z})}$ e $\mathbf{V}^{(\mathbf{Z})}$, o que podemos afirmar no que diz respeito à independência de suas componentes?
2. Considere os arquivo de dados da iris setosa (com as 4 variáveis) e, num primeiro momento, ignore a existência dos grupos (espécie). Considere que $\mathbf{X}^{(1)} = (X_1^{(1)}, X_2^{(1)})'$ e $\mathbf{X}^{(2)} = (X_1^{(2)}, X_2^{(2)})'$, em que $X_1^{(1)}$: comprimento da sépala, $X_2^{(1)}$: largura da sépala, $X_1^{(2)}$: comprimento da pétala e $X_2^{(2)}$: largura da pétala. Queremos estudar como as

variáveis relacionadas as dimensões das pétalas influenciam aquelas relacionadas às dimensões das sépalas, através das análise de correlações canônicas, utilizando as variáveis padronizadas. Sugestão: veja os slides inerentes à análise de correlação canônica.

Responda os itens:

- a) Faça um diagrama de dispersão matricial entre as 4 variáveis. O que você pode dizer a respeito da estrutura de correlação delas?
 - b) Obtenha os dois pares de variáveis canônicas, apresentando as equações canônicas, as correlações entre cada variável canônica e as variáveis padronizadas e as correlações canônicas entre as variáveis canônicas. Interprete todos estes resultados, de modo adequado. Interprete, inclusive, cada uma das variáveis canônicas em termos das variáveis padronizadas.
 - c) Calcule o percentual da variabilidade de cada uma das variáveis canônicas explicadas pelas variáveis padronizadas. Calcule a proporção das somas das variâncias das variáveis originais explicadas pelas variáveis canônicas. Interprete todos estes resultados, de modo adequado.
 - d) Com base nos resultados dos itens anteriores, quantos pares de variáveis canônicas você escolheria para representar os dados? Justifique, adequadamente, sua resposta.
 - e) Utilizando somente o primeiro par de variáveis canônicas, ajuste um modelo de regressão normal linear da primeira variável canônica contra a segunda. Interprete este modelo da forma mais ampla possível e verifique se os coeficientes são significativos. Verifique se a suposição de normalidade é razoável para dos dados.
 - f) Faça um diagrama de dispersão entre as duas primeiras variáveis canônicas. Como as espécies (setosa, versicolor e virginica) se comportam com relação à estas variáveis? As espécies apresentam comportamentos diferentes? Como cada uma das espécies se caracteriza? Ou seja, descreva o comportamento das espécies. Sugestão: utilize símbolos e ou cores diferentes para cada grupo.
3. Considere a metodologia de análise de discriminante apresentada em classe. Considere dois vetores aleatórios independentes \mathbf{X}_1 e \mathbf{X}_2 , tais que $\mathbf{X}_i \sim N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, $i=1,2$. Considere \mathbf{x} um valor observado de algum desses dois vetores. Responda os itens:
- a) Prove que, se $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}$ então, $-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)$.

- b) Seja $y_i = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i$ e $m = \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)$. Prove que $m = \frac{1}{2} (y_1 + y_2)$
- c) Prove que, se $\boldsymbol{\Sigma}_1 \neq \boldsymbol{\Sigma}_2$ então, $-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) = -\frac{1}{2} \mathbf{x}' (\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1}) \mathbf{x} + (\boldsymbol{\mu}_1' \boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\mu}_2' \boldsymbol{\Sigma}_2^{-1}) \mathbf{x} - c$, em que $c = \frac{1}{2} \ln \left(\frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_2|} \right) + \frac{1}{2} (\boldsymbol{\mu}_1' \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2' \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2)$
- d) Repita os itens de a) - c), considerando amostras aleatórias (independentes) dos vetores aleatórios \mathbf{X}_1 e \mathbf{X}_2 , ou seja, substituindo os parâmetros em questão pelos respectivos estimadores usualmente adotados.
4. Considere os dados da Tabela 11.9 nas páginas 666 e 667 do livro do Johnson & Wichern (sexta edição) sobre características nutricionais de cereais de fabricantes americanas (note que, embora se pareçam, não são exatamente os mesmos dados de cereal considerados anteriormente). Considere todas as variáveis em questão (calorias, proteína, gordura, sódio, fibra, carboidratos, açúcar e potássio). Os grupos são os fabricantes (considere apenas dois) G e K. Responda aos itens:
- a) Faça um diagrama de dispersão matricial entre as todas as variáveis. O que você pode dizer a respeito da estrutura de correlação delas?
- b) Calcule a matriz de covariâncias ponderada \mathbf{S}_p^2 (considerando os dois grupos) bem como a respectiva matriz de correlações.
- c) Assuma distribuição normal multivariada para cada um dos 3 grupos, com iguais matrizes de covariâncias. Assuma também custos iguais de classificação errada e iguais probabilidades a priori de classificação. Considere as 10 primeiras observações de cada um dos grupos e construa a regra de classificação de Fisher, dentro desse contexto. Interprete a função discriminante de modo adequado. Estime a PTCE bem como o TAP (utilizando os outros cereais que não entraram na regra de classificação). Critique, do modo mais completo possível, a regra de classificação que você obteve.
- d) Existem indícios de que algum(ns) fabricantes estão relacionados à cereais mais “nutritivos” (alto teor de proteína, baixo teor de gordura, alto teor de fibra, baixo teor de açúcar e assim por diante)? Sugestão: faça um gráfico de dispersão, dos valores da função discriminante, para cada um dos grupos, plotando, no gráfico, as iniciais das fabricantes.