

Análise Multivariada

Lupércio França Bessegato
Dep. Estatística/UFJF

Roteiro

1. Introdução
2. Distribuições de Probabilidade Multivariadas
3. Representação de Dados Multivariados
4. Testes de Significância c/ Dados Multivariados
5. Análise de Componentes Principais
6. Análise Fatorial
7. Análise de Agrupamentos
8. Análise de Correlação Canônica
9. Referências

Análise Multivariada - 2022

3

Discriminação e Classificação

Roteiro

1. Introdução
2. Classificação e Discriminação
3. Análise Discriminante
4. Outras Abordagens
5. Referências

Análise Multivariada - 2022

5

Introdução

Agrupamento e Classificação

- Agrupar:
 - ✓ Processo de alocar item em grupo
 - ✓ Não há suposições sobre o número de grupos ou sobre a estrutura dos grupos
 - Técnica mais primitiva
- Classificar:
 - ✓ Predição de pertinência a grupo
 - ✓ Número de grupos é conhecido e o objetivo é alocar novas observações a um desses grupos
 - ✓ Usa status conhecido para encontrar preditores, aplicando-os a uma nova observação

Análise Multivariada - 2022

7

Conjunto de Dados

- Partição do conjunto de dados:
 - ✓ Conjunto de treinamento
 - Usado para desenvolver modelo de classificação
 - ✓ Conjunto de teste
 - Usado para determinar desempenho do modelo
 - ✓ Importante não avaliar desempenho com as mesmas observações usadas para desenvolver o modelo

8

Análise Multivariada - 2022

Passos para Classificação

1. Conjunto de dados coletado, com alocações de item em grupo já conhecidas (ou atribuídas)
 - ✓ Observação, julgamento de especialista, procedimentos de agrupamento
2. Dados são divididos em conjunto de treinamento e teste
 - ✓ Treinamento: de 50% a 80% (comum: 67%)
 - ✓ Restante atribuído ao conjunto de teste

9

Análise Multivariada - 2022

3. Construção do modelo de predição

- ✓ Predizer alocação dos dados de treinamento tão bem quanto possível

4. Avaliação do desempenho do modelo usando os dados do conjunto de teste

Análise Multivariada - 2022

10

Métodos de Classificação

- Há inúmero métodos de classificação:
 - ✓ Análise discriminante
 - ✓ Regressão logística
 - ✓ Naive Bayes Classification
 - ✓ Random Forest Classifiers
 - ✓ Método do vizinho mais próximo
 - ✓ Classification and Regression Trees – CART
 - ✓ Support Vector Machine – SVM
 - ✓ Método dos núcleos estimadores
 - ✓ Redes neurais artificiais

Análise Multivariada - 2022

11

Análise de Agrupamento e Análise Discriminante

- Análise de Agrupamentos

- ✓ Dividir os elementos da amostra (ou população) em grupos, de maneira que:
 - Elementos de um grupo são similares entre si
 - Elementos de grupos diferentes sejam heterogêneos em relação a essas características

Análise Multivariada - 2022

12

- Análise discriminante:

- ✓ Classificação de elementos de amostra (população)
 - Grupos são pré-definidos
- ✓ Procedimento:
 - Regra de classificação

Análise Multivariada - 2022

13

Análise Discriminante

- Caso especial de correlações canônicas
 - ✓ Variáveis dependentes são categóricas por natureza
- Objetivo:
 - ✓ Usar informações das variáveis independentes para a separação (discriminação) mais clara possível entre os grupos

Análise Multivariada - 2022

14

- Abordagens:
 - ✓ Fischer
 - ✓ Mahalanobis

Análise Multivariada - 2022

15

Análise Discriminante

- Caso especial de correlações canônicas
 - ✓ Variáveis dependentes são categóricas por natureza
- Objetivo:
 - ✓ Usar informações das variáveis independentes para a separação (discriminação) mais clara possível entre os grupos

Análise Multivariada - 2022

16

- Abordagens:
 - ✓ Fischer
 - ✓ Mahalanobis

Análise Multivariada - 2022

17

Aplicações Potenciais

- Perfil:

- ✓ Compreender como cada variável independente (X) influencia a variável dependente (Y: grupo)
- ✓ Descrição, em análise de regressão
- ✓ Quando os objetivos do estudo são principalmente exploratórios

Análise Multivariada - 2022

18

- ✓ Como os grupos são discriminados pelas variáveis subjacentes?

- Exame dos perfis de segmentos do mercado para entender como consumidores diferem com relação a variáveis demográficas e psicológicas
- Diferenças entre usuários de categoria de produto em relação ao tamanho da família, renda, educação, etc.

- ✓ Como potenciais consumidores de marca diferem da população em geral em relação ao seu envolvimento com a mídia?

Análise Multivariada - 2022

19

- Diferenciação:

- ✓ Capacidade de afirmar, com certo nível de confiança, se a relação entre X e Y se deve ao acaso
- ✓ Inferência, em análise de regressão
- ✓ Traçados os perfis dos grupo, pode ser importante verificar se as diferenças aparentes entre eles dão de fato significativas
- ✓ Exemplo:
 - Entender e controlar as variações associadas a certos processos de produção

Análise Multivariada - 2022

20

- Classificação:

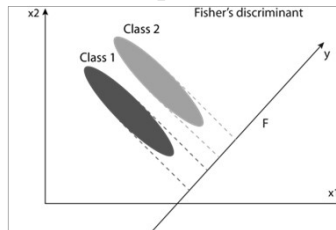
- ✓ Usar o modelo para avaliar o valor da variável dependente, com observações fora da amostra de treinamento
 - Predizer a pertinência a grupo
- ✓ Predição, em análise de regressão
- ✓ Exemplos:
 - *Credit scoring*
 - Traçar o perfil dos clientes de empréstimo e julgar se novos candidatos oferecem risco ao crédito
 - Marketing direto
 - Que perfil de clientes devem receber oferta de mala direta?

Análise Multivariada - 2022

21

Fisher – Intuição

- Baseia-se na noção de pontuação discriminante



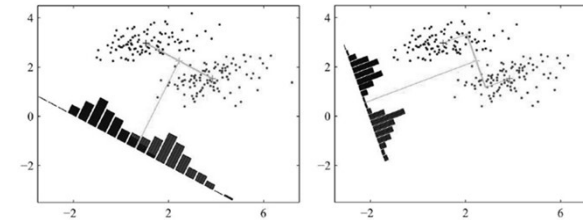
- ✓ Encontrar combinação linear das variáveis independente que produza pontuações discriminantes maximamente diferentes

Análise Multivariada - 2022

22

- Função objetivo:

✓ Quantifica a noção de “maximamente diferente”



✓ Função linear que melhor aloca as observações

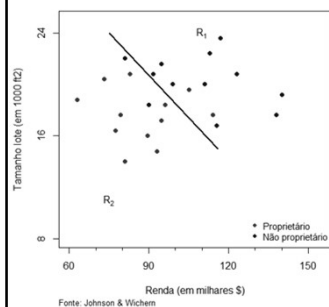
- Eixo que descreve diferença entre centróides
- Ajusta de acordo com o padrão de covariância

Análise Multivariada - 2022

23

Mahalanobis – Intuição

- Encontrar o ‘locus’ dos pontos equidistantes das médias dos 2 grupos



- ✓ 2 variáveis explicativas:
 - ‘locus’ dos pontos é uma linha
- ✓ 3 variáveis explicativas:
 - ‘locus’ dos pontos é um plano ou hiperplano
- ✓ ‘locus’ serve para discriminar os dois grupos

Análise Multivariada - 2022

24

- Medida de distância ajustada

$$D_i^2 = (\mathbf{x} - \bar{\mathbf{x}}_{(i)})' \mathbf{C}_W^{-1} (\mathbf{x} - \bar{\mathbf{x}}_{(i)}), i = 1, 2.$$

✓ Distância ao quadrado da covariância ajustada de qualquer ponto \mathbf{x} à média do grupo i

✓ Dados seguem normal multivariada:

- Distância ajustada reflete com mais precisão a probabilidade de pertinência ao grupo do que a distância euclidiana

Análise Multivariada - 2022

25

- Por definição, 'locus' dos pontos descritos por Mahalanobis é ortogonal ao eixo da função discriminante proposta por Fisher

Análise Multivariada - 2022

26

Análise Discriminante – Abordagens

- São complementares:
 - √ Fisher:
 - Reduz os dados em uma única dimensão de modo a maximizar a separação entre grupos
 - √ Mahalanobis:
 - Determina linha divisória (ou plano) que separa mais precisamente os dois grupos
 - Ortogonal à dimensão discriminante

Análise Multivariada - 2022

27

Discriminação e Classificação para Duas Populações

Objetivo

1. Separar duas classes de objetos
 2. Atribuir um novo objeto a uma das duas classes
- Classes:
 - √ π_1 e π_2
 - Objetos:
 - √ São separados ordinariamente ou classificados com base em medidas de p variáveis aleatórias associadas

Análise Multivariada - 2022

29

Conceitos

- **Classes:**
✓ π_1 e π_2
- **Objetos:**
✓ São separados ordinariamente ou classificados com base em medidas de p variáveis aleatórias associadas
$$\mathbf{X}' = [X_1, X_2, \dots, X_p].$$
- **Hipótese:**
✓ Os valores observados de \mathbf{X} diferem em alguma quantidade de uma classe a outra

Análise Multivariada - 2022

30

- **Populações das duas classes:**
✓ Podem ser descritas por suas funções de densidade
 $[f_1(\mathbf{x}) \text{ e } f_2(\mathbf{x})]$
✓ Pode-se falar em atribuir:
 - Observações a populações ou
 - Objetos a classes

Análise Multivariada - 2022

31

Exemplos

Populações π_1 e π_2	Variáveis \mathbf{X}
Risco de crédito alto/baixo	Renda, Idade, n° de cartões de crédito, tamanho família
Duas espécies de flor	Comprimento e largura de pétalas e sépalas, diâmetro do pólen, etc.
Masculino e feminino	Medidas antropológicas tais como: circunferência e volume de crânios antigos
Seleção a curso de pós-graduação	Histórico escolar, curriculum vitae, cartas de referência, experiência profissional

Análise Multivariada - 2022

32

- **Outros Exemplos:**
✓ **Agricultura:**
 - Identificar áreas de maior potencial para plantação de determinadas sementes
- ✓ **Marketing:**
 - Identificar mercados potenciais e não potenciais para determinados produtos e serviços
- ✓ **Esporte:**
 - Identificação de atletas promissores para cada modalidade
- ✓ **Estudos de criminalidade:**
 - Identificação de regiões que necessitam de política de segurança diferenciada

Análise Multivariada - 2022

33

Regras de Alocação e Classificação

- Em geral, são desenvolvidas de amostras de treinamento:
 - √ Examinadas diferenças das medidas características de objetos selecionados
 - √ O conjunto de todos os resultados amostrais possíveis são dividido em duas regiões (R_1 e R_2)
 - Se uma nova observação pertencer à região R_1 ela é alocada à população π_1 .
 - Se uma nova observação pertencer à região R_2 ela é alocada à população π_2 .

Análise Multivariada - 2022

34

Problema da Classificação

- Como saber se algumas observações pertencem a uma particular população?
 - √ Incerteza na classificação

Análise Multivariada - 2022

35

Paradoxos da Classificação

- Informação incompleta sobre desempenho futuro:
 - √ Classificação de candidato como capaz de concluir ou não um mestrado
- Informação perfeita exige destruição objeto:
 - √ Classificação de itens como bons ou defeituosos
- Informação cara ou indisponível:
 - √ Problemas médicos que podem ser identificados conclusivamente apenas com procedimentos caros

Análise Multivariada - 2022

36

Erros de Classificação

- Caso médico:
 - √ Em geral, deseja-se diagnosticar um mal a partir de sintomas externos facilmente observáveis
- Erro de classificação:
 - √ Pode não ser clara a distinção entre as características medidas das duas populações.

Análise Multivariada - 2022

37

Exemplo

- Discriminação de proprietários e não-proprietários de cortador de grama

✓ X_1 : renda

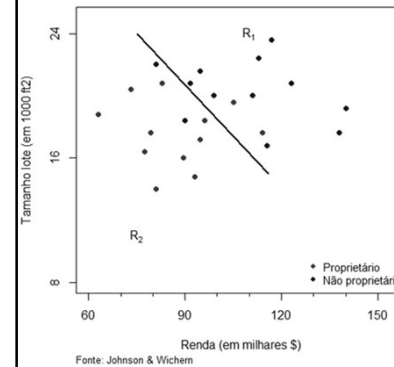
✓ X_2 : tamanho do lote



Análise Multivariada - 2022

38

- Proprietários e não proprietários: Renda e tamanho lote



- Proprietários tendem a ter rendas e lotes maiores
- Renda aparenta discriminar melhor que tamanho do lote
- Há uma certa sobreposição entre os dois grupos
 - ✓ Erros de classificação

Ideia:
Criar regra que minimize a chance de erros de classificação

Análise Multivariada - 2022

39

Crítérios para Classificação

- Bom procedimento de classificação:
 - ✓ Poucos erros de classificação
- Regra ótima deveria considerar estas probabilidades a priori
 - ✓ Pode ser que uma classe (ou população) tenha uma verossimilhança de ocorrência maior que outra
 - ✓ Uma das classes é relativamente maior que a outra
 - ✓ Ex.:
 - Há muito mais empresas solventes que insolventes

Análise Multivariada - 2022

40

- Outro aspecto a considerar:

✓ Custo associado ao erro de classificação

✓ Ex.:

- Classificar um objeto π_1 como π_2 é mais sério que classificar um objeto π_2 como π_1 .

Análise Multivariada - 2022

41

Caso de Classificação de Duas Populações Normais Multivariadas

Classificação em Duas Populações

- Supondo disponíveis:
 - ✓ Conjunto de observações independentes de duas populações π_1 e π_2
 - ✓ Distribuições de probabilidades do vetor X , associadas às populações π_1 e π_2 .
- Regra de classificação que minimize a chance de se classificar incorretamente elemento amostral:
 - ✓ Princípio da máxima verossimilhança

Análise Multivariada - 2022

44

Exemplo

- Processo de seleção de alunos:
 - ✓ Fase 1: todos fazem várias provas
 - ✓ Fase 2: apenas aprovados na fase 1
- Populações:
 - ✓ População 1:
 - Alunos que passaram na 1ª. Fase, mas reprovados na 2ª
 - ✓ População 2:
 - Alunos aprovados em ambas as fases

Análise Multivariada - 2022

45

Exemplo

- Objetivo:
 - ✓ A partir dos dados, construir uma regra de classificação que permita identificar, dentre os aprovados na 1ª. Fase, quais provavelmente serão aprovados na 2ª Fase
- Considere apenas a variável aleatória nota na prova de Matemática dos candidatos na fase 1

Análise Multivariada - 2022

46

Exemplo

- Considere apenas a variável aleatória nota na prova de Matemática dos candidatos na fase 1
- Suponha que X tenha uma distribuição normal
 - ✓ População 1: média μ_1 .
 - ✓ População 2: média μ_2
 - ✓ Ambas populações como o mesmo desvio padrão σ .

Análise Multivariada - 2022

47

Exemplo

- Razão de verossimilhança entre as 2 populações

$$\lambda(x) = \frac{f_1(x)}{f_2(x)} = \frac{\frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu_1}{\sigma}\right)^2\right\}}{\frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu_2}{\sigma}\right)^2\right\}} = \exp\left\{-\frac{1}{2}\left[\left(\frac{x-\mu_1}{\sigma}\right)^2 - \left(\frac{x-\mu_2}{\sigma}\right)^2\right]\right\}$$

- Para uma nota fixa x:

- ✓ $\lambda(x) > 1$: razoável classificar candidato em π_1 (não aprovado da Fase 2)
- ✓ $\lambda(x) < 1$: provável aprovado Fase 2 (π_2)
- ✓ $\lambda(x) = 1$: candidato poderia ser classificado em π_2 ou π_2 .
 - Obter informações adicionais sobre o candidato

Análise Multivariada - 2022

48

- Classificação de 2 populações normais, com mesma variabilidade

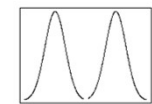


Gráfico a: nenhuma intersecção

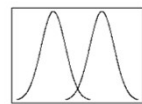


Gráfico b: pouca intersecção

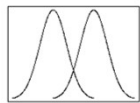


Gráfico c: intersecção moderada



Gráfico d: grande intersecção

Fonte: Mingoti

Qualidade da discriminação depende do grau de intersecção

- ✓ (a): número de classificações incorretas é 0
- ✓ (b): pequeno número de erros de classificação
- ✓ (c) e (d): número de erros de classificação tende a aumentar
- ✓ Intersecção pode chegar a valores que inviabiliza o uso da função discriminante como regra de classificação

Análise Multivariada - 2022

49

Função Discriminante

$$\lambda(x) = \exp\left\{-\frac{1}{2}\left[\left(\frac{x-\mu_1}{\sigma}\right)^2 - \left(\frac{x-\mu_2}{\sigma}\right)^2\right]\right\}$$

$$-2 \ln(\lambda(x)) = \left(\frac{x-\mu_1}{\sigma}\right)^2 - \left(\frac{x-\mu_2}{\sigma}\right)^2$$

$$= \frac{1}{\sigma^2} [(x-\mu_1)^2 - (x-\mu_2)^2]$$

- ✓ Relacionada com a diferença das distância euclidiana ponderadas ao quadrado

✓ $\lambda(x) > 1 \Rightarrow -2 \ln(\lambda(x)) < 0$: x está mais próximo de μ_1 .

✓ $\lambda(x) < 1 \Rightarrow -2 \ln(\lambda(x)) > 0$: x está mais próximo de μ_2 .

Análise Multivariada - 2022

50

- Regra de classificação:

- ✓ Se $-2\ln(\lambda(x)) < 0$, classifique o elemento amostral em π_1
- ✓ Se $-2\ln(\lambda(x)) > 0$, classifique o elemento amostral em π_2 .
- ✓ Se $-2\ln(\lambda(x)) = 0$, o elemento amostral poderá ser classificado tanto em π_1 como π_2 .

Análise Multivariada - 2022

51

Populações com Variâncias Diferentes

- Função discriminante:

$$\lambda(x) = \frac{f_1(x)}{f_2(x)} = \frac{\frac{1}{\sqrt{2\pi}\sigma_1} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu_1}{\sigma_1}\right)^2\right\}}{\frac{1}{\sqrt{2\pi}\sigma_2} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu_2}{\sigma_2}\right)^2\right\}}$$

$$= \frac{\sigma_2}{\sigma_1} \exp\left\{-\frac{1}{2}\left[\left(\frac{x-\mu_1}{\sigma_1}\right)^2 - \left(\frac{x-\mu_2}{\sigma_2}\right)^2\right]\right\}$$

✓ e, considerando $-2\ln(\lambda(x))$:

$$-2\ln(\lambda(x)) = -2\ln\left(\frac{\sigma_2}{\sigma_1}\right) + \left[\left(\frac{x-\mu_1}{\sigma_1}\right)^2 - \left(\frac{x-\mu_2}{\sigma_2}\right)^2\right].$$

- Com mesma regra de classificação

Análise Multivariada - 2022

52

Populações Multivariadas – Caso $\Sigma_1 \neq \Sigma_2$

- Populações normais multivariadas, com vetor de médias μ_i e matriz de covariâncias Σ_i , $i = 1, 2$.
- Função discriminante:

$$-2\ln(\lambda(x)) = -2\ln\left\{\frac{(2\pi)^{\frac{p}{2}}|\Sigma_1|^{-\frac{1}{2}} \left[\exp\left\{-\frac{1}{2}(\mathbf{x}-\mu_1)'\Sigma_1^{-1}(\mathbf{x}-\mu_1)\right\}\right]}{(2\pi)^{\frac{p}{2}}|\Sigma_2|^{-\frac{1}{2}} \left[\exp\left\{-\frac{1}{2}(\mathbf{x}-\mu_2)'\Sigma_2^{-1}(\mathbf{x}-\mu_2)\right\}\right]}\right\}$$

$$= [(\mathbf{x}-\mu_1)'\Sigma_1^{-1}(\mathbf{x}-\mu_1) - (\mathbf{x}-\mu_2)'\Sigma_2^{-1}(\mathbf{x}-\mu_2)]$$

$$+ [\ln|\Sigma_1| - \ln|\Sigma_2|]$$

✓ para um vetor de observações: $\mathbf{x}' = (x_1, x_2, \dots, x_p)$.

- Com mesma regra de classificação

Análise Multivariada - 2022

53

- Função discriminante quadrática

$$[(\mathbf{x}-\mu_1)'\Sigma_1^{-1}(\mathbf{x}-\mu_1) - (\mathbf{x}-\mu_2)'\Sigma_2^{-1}(\mathbf{x}-\mu_2)] + [\ln|\Sigma_1| - \ln|\Sigma_2|]$$

✓ Depende das distâncias de Mahalanobis do vetor \mathbf{x} aos vetores de médias μ_1 e μ_2

✓ Fator de correção relacionando as variâncias generalizadas das duas populações

Análise Multivariada - 2022

54

- Quando $\Sigma_1 = \Sigma_2 = \Sigma$

$$-2\ln(\lambda(\mathbf{x})) = [(\mathbf{x} - \mu_1)' \Sigma^{-1} (\mathbf{x} - \mu_1) - (\mathbf{x} - \mu_2)' \Sigma^{-1} (\mathbf{x} - \mu_2)]$$

✓ Pode ser reescrita como

$$-2\ln(\lambda(\mathbf{x})) = -2 \left[(\mu_1 - \mu_2)' \Sigma^{-1} \mathbf{x} - \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) \right].$$

Análise Multivariada - 2022

55

Função Discriminante de Fisher

$$f_d(\mathbf{x}) = \ln(\lambda(\mathbf{x})) = (\mu_1 - \mu_2)' \Sigma^{-1} \mathbf{x} - \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2).$$

- Regra de classificação:

✓ \mathbf{x} é classificado à π_1 se $f_d(\mathbf{x}) > 0$ ou seja, se

$$(\mu_1 - \mu_2)' \Sigma^{-1} \mathbf{x} > \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2).$$

✓ \mathbf{x} é classificado à π_2 se $f_d(\mathbf{x}) < 0$ ou seja, se

$$(\mu_1 - \mu_2)' \Sigma^{-1} \mathbf{x} < \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2).$$

Análise Multivariada - 2022

57

- A função discriminante de Fisher tem a forma

$$(\mu_1 - \mu_2)' \Sigma^{-1} \mathbf{x} = \mathbf{b}' \mathbf{x} = b_1 X_1 + b_2 X_2 + \dots + b_p X_p.$$

✓ Dependendo do valor numérico desta combinação, o elemento amostral é classificado em uma ou outra população

- Constante de delimitação da região de classificação

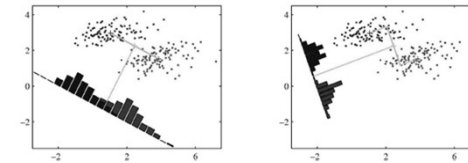
$$\frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) = \mathbf{b}' \frac{(\mu_1 + \mu_2)}{2}.$$

✓ Combinação linear dos vetores de médias das populações

Análise Multivariada - 2022

58

- Discriminação de duas populações normais:



✓ Os valores observados das combinações lineares $\mathbf{b}'\mathbf{X}$ na população π_1 são os mais separados possíveis daqueles observados da população π_2 .

Análise Multivariada - 2022

61

Padronização

- O vetor **b** é único $\mathbf{b}' = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1}$.
 ✓ Exceto para multiplicações de todos seus componentes pela mesma constante c.
- É recomendável que os componentes do vetor **b** sejam padronizados ou normalizados, como em:

$$\mathbf{b}^* = \frac{\mathbf{b}}{\sqrt{\mathbf{b}'\mathbf{b}}}.$$
 - ✓ Componentes de **b*** estarão no intervalo [-1, 1]
 - ✓ É possível comparar os loadings de f_{dn}

Análise Multivariada - 2022

63

Custos

- Regra de discriminação pelo princípio da máxima verossimilhança minimiza as probabilidades de erros de classificações
 ✓ Máxima separação entre as combinações lineares
- Não leva em consideração possíveis diferenças entre os custos associados aos erros de classificação

Análise Multivariada - 2022

66

Mistura de Duas Normais Multivariadas

- Suponha:
 - ✓ 2 populações normais p-variadas com vetor de médias $\boldsymbol{\mu}_i$, $i = 1, 2$
 - ✓ Probabilidades de mistura p_j , $j=1, 2$
 - Probabilidade de que uma observação escolhida ao acaso pertença à π_j , $j=1, 2$.
 - ✓ Função de densidade de π_j . $f(\mathbf{x}|\pi_j) = f(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$.

Análise Multivariada - 2022

68

- ✓ Dada uma observação **x**, qual a melhor maneira de distinguir de qual das duas populações ela foi amostrada?

$$f(\mathbf{x} \text{ amostrada de } \pi_j) = p_j f(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}).$$

- ✓ Densidade de **x** oriunda de uma população não especificada:

$$f(\mathbf{x}) = p_1 f(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}) + p_2 f(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}).$$

- ✓ Densidade posterior de **x** de amostras da população j para uma dada observação **x**

$$f(\pi_j|\mathbf{x}) = \frac{p_j f(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma})}{p_1 f(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}) + p_2 f(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma})}.$$

Expressão básica para estimar a população da qual a observação **x** foi amostrada

Análise Multivariada - 2022

70

- Função discriminante de Fisher:

$$f_d(\mathbf{x}) = \underbrace{\ln\left(\frac{p_1}{p_2}\right) - \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)}_{\text{constante}} + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} \mathbf{x}.$$

- Limite entre as populações π_1 e π_2 .

$$f_d(\mathbf{x}) = \text{constante} + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} \mathbf{x} = 0.$$

Análise Multivariada - 2022

72

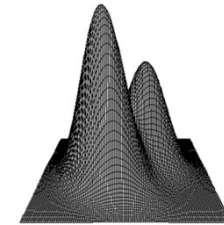
Exemplo

- Mistura de normais bivariadas com mesma estrutura de variabilidade:

$$\boldsymbol{\Sigma} = \begin{bmatrix} 1 & -0,4 \\ -0,4 & 1,2 \end{bmatrix}.$$

$$\sqrt{\pi_1}: p_1 = 0,6 \text{ e } \boldsymbol{\mu}_1 = (-1, -1)'$$

$$\sqrt{\pi_2}: p_2 = 0,4 \text{ e } \boldsymbol{\mu}_2 = (1, 1)'.$$



Análise Multivariada - 2022

73

- Função discriminante:

$$\begin{aligned} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} \mathbf{x} &= [-1 \ -1 \ -1 \ -1] \begin{bmatrix} 1 & -0,4 \\ -0,4 & 1,2 \end{bmatrix}^{-1} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= [-3,0769 \ -2,6923] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= -3,0769x_1 - 2,6923x_2. \end{aligned}$$

- Constante de delimitação:

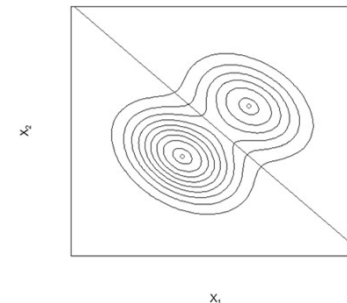
$$\begin{aligned} \text{cte} &= \ln\left(\frac{0,6}{0,4}\right) - \frac{1}{2} \left([-1 \ -1 \ -1 \ -1] \begin{bmatrix} 1 & -0,4 \\ -0,4 & 1,2 \end{bmatrix}^{-1} \begin{bmatrix} -1+1 \\ -1+1 \end{bmatrix} \right) \\ &= 0,40546. \end{aligned}$$

- Fronteira: $-3,0769x_1 - 2,6923x_2 = 0,40546$
 $x_2 = 0,1506 - 1,1428x_1 - 1.$

Análise Multivariada - 2022

74

- Fronteira entre as duas populações



$$x_2 = 0,1506 - 1,1428x_1 - 1.$$

√ Linha divisória discriminando as duas populações

Análise Multivariada - 2022

75

Estimação da Regra de Classificação

- Na prática, μ_1 , μ_2 , Σ_1 e Σ_2 não são conhecidos
- Caso 1: $\Sigma_1 = \Sigma_2 = \Sigma$

✓ Σ é estimada por S: $S_{\text{pol}} = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{n_1 + n_2 - 2}$.

✓ Função discriminante de Fisher estimada por:

$$\hat{f}_d(\mathbf{x}) = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' S_{\text{pol}}^{-1} \mathbf{x} - \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' S_{\text{pol}}^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2).$$

Análise Multivariada - 2022

76

- Caso 2: $\Sigma_1 \neq \Sigma_2$.

✓ Função discriminante quadrática estimada por:

$$-2 \ln(\hat{\lambda}(\mathbf{x})) = [(\mathbf{x} - \bar{\mathbf{x}}_1)' S_1^{-1} (\mathbf{x} - \bar{\mathbf{x}}_1) - (\mathbf{x} - \bar{\mathbf{x}}_2)' S_2^{-1} (\mathbf{x} - \bar{\mathbf{x}}_2)] + [\ln |S_1| - \ln |S_2|]$$

✓ Considera os sistemas de variabilidade das duas populações separadamente.

Análise Multivariada - 2022

77

Estrutura de Variabilidades das Populações

- Testes de hipóteses para decidir se as matrizes Σ_1 e Σ_2 são iguais ou diferentes
- Alternativa prática:
 - ✓ Ajuste aos dados dos dois modelos: linear de Fisher e quadrático
 - ✓ Escolhe-se o modelo que resultar em menores proporções de erros de classificação
 - ✓ Caso os resultados sejam semelhantes, opta-se pelo modelo linear
 - Matriz de covariâncias estimada com mais observações

Análise Multivariada - 2022

78

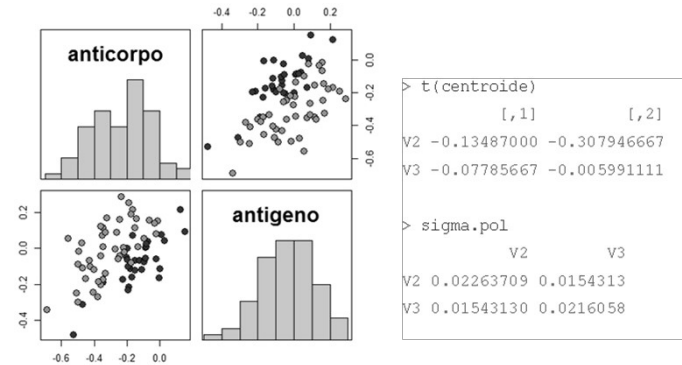
Exemplo

- Detecção de portadoras de hemofilia A
 - ✓ Grupo 1: mulheres que não têm o gene da hemofilia
 - Grupo normal ($n_1 = 30$)
 - ✓ Grupo 2: mulheres portadoras do gene da hemofilia
 - Filhas de hemofílicos, mães com mais de um filho hemofílico e outros parentes hemofílicos ($n_2 = 45$)
 - ✓ Variáveis:
 - V_1 : grupos
 - V_2 : log(atividade AHF)
 - V_3 : log(antígeno AHF)

Análise Multivariada - 2022

81

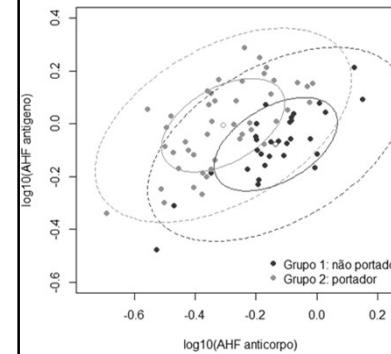
- Análise descritiva do conjunto de dados:



Análise Multivariada - 2022

82

- *Contour plot* por grupo



- ✓ Curvas contendo 50% e 95% de probabilidade para normais bivariadas centradas em \bar{x}_1 e \bar{x}_2 .
- ✓ Normal bivariada aparenta se ajustar bem aos dados

Análise Multivariada - 2022

83

- Função discriminante de Fisher

```
>> ajuste.ad = lda(dados, grupo)
> ajuste.ad
Call:
lda(dados, grupo)

Prior probabilities of groups:
  1  2 
0.4 0.6 

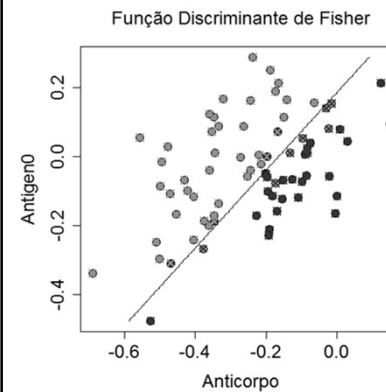
Group means:
      V2      V3
1 -0.1348700 -0.077856667
2 -0.3079467 -0.005991111

Coefficients of linear discriminants:
      LD1
V2 -9.032787
V3  8.006605
```

Análise Multivariada - 2022

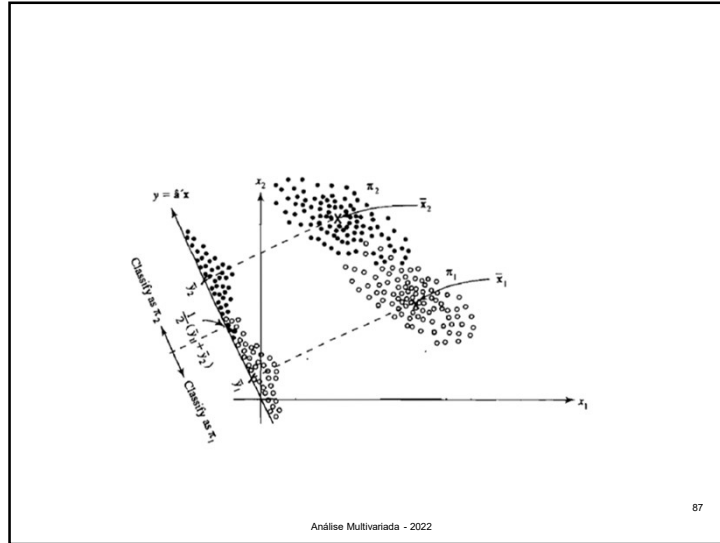
85

- Gráfico da função discriminante

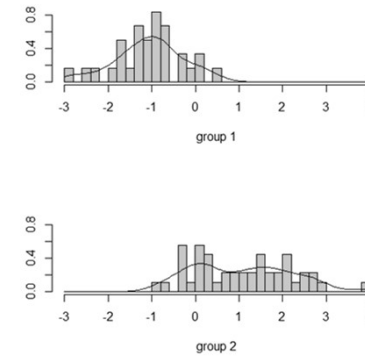


Análise Multivariada - 2022

86



- Grupos ajustados



- Predição:

✓ Mulher com $V1 = -0,210$ e $V2 = -0,044$

```
> predict(ajuste.ad, newdata = c(-0.210, -0.044))$class
[1] 2
Levels: 1 2
```

✓ Mulher pode ser portadora de hemofilia

Análise Multivariada - 2022

90

Avaliação das Funções de Classificação

Erros de Classificação

- Erros a serem avaliados:
 - ✓ Erro 1: Elemento amostral pertence a π_1 , mas é classificado em π_2 .
 - ✓ Erro 2: Elemento amostral pertence a π_2 , mas é classificado em π_1 .
- Notação:
 - ✓ $P(\text{Erro 1}) = p(2|1)$
 - ✓ $P(\text{Erro 2}) = p(1|2)$
- Quanto menores essas probabilidades, melhor será a função de discriminação

Análise Multivariada - 2022

95

Procedimentos de Estimação dos Erros de Classificação

1. Método da Ressubstituição:
 - ✓ Calculado o escore de cada elemento amostral
 - ✓ Calculada a frequência das classificações corretas e incorretas
 - ✓ Estimação da regra de classificação e dos erros de classificação com os mesmos elementos

Análise Multivariada - 2022

96

- Tabela de frequência de classificação:

Origem	Classificação		Total
	π_1	π_2	
	π_1	π_2	
	n_{11}	n_{12}	n_1
	n_{21}	n_{22}	n_2

- ✓ Estimativa das probabilidades de erros de classificação:

$$\hat{p}(2|1) = \frac{n_{12}}{n_1}$$

$$\hat{p}(1|2) = \frac{n_{21}}{n_2}$$

Análise Multivariada - 2022

97

Comentários

- Também denominado de estimação do erro aparente de classificação (APER)
- Procedimento viciado, mas consistente.
 - ✓ Vício tende a zero para n_1 e n_2 grandes
 - ✓ Tende a subestimar os verdadeiros valores de $p(1|2)$ e $p(2|1)$ para elementos que não pertencem à amostra conjunta $n = n_1 + n_2$.
 - ✓ Pode servir como etapa inicial de avaliação
 - Valores elevados indica a necessidade de reformulação da regra de discriminação

Análise Multivariada - 2022

98

2. Método de colocação de elementos à parte para classificação (Hold-out validation)

- ✓ Amostra conjunta é repartida em duas partes
 - Amostra de treinamento: construção da regra de discriminação
 - Amostra de validação: para estimação dos erros de classificação
- ✓ São selecionados aleatoriamente os elementos amostrais que constituirão cada amostra
- ✓ Estimação dos erros de classificação da maneira descrita no método de ressubstituição

Análise Multivariada - 2022

99

Comentários

- Procedimento não é enviesado
- Recomendável:
 - ✓ Separar de 25% a 50% dos elementos originais para a amostra de validação
- Desvantagem:
 - ✓ Redução do tamanho da amostra original para estimação da regra de discriminação
 - ✓ Não pode ser empregado em amostras pequenas
- Para amostras grandes, é melhor que o método 1

Análise Multivariada - 2022

100

3. Método de validação cruzada (Método de Lachenbruch)

- ✓ Retira-se um elemento amostral da amostra conjunta e constrói-se a função de discriminação
- ✓ Utiliza-se a regra de discriminação para classificar o elemento que ficou à parte
- ✓ Elemento amostral é retornado à amostra e retira-se elemento amostral diferente do anterior, repetindo-se o procedimento.
- Estimação dos erros de classificação:

$$\hat{p}(2|1) = \frac{n_{12}}{n_1}$$
$$\hat{p}(1|2) = \frac{n_{21}}{n_2}$$

Análise Multivariada - 2022

101

Comentários

- Estimativas são aproximadamente não viciadas
 - ✓ Melhores que o método da ressubstituição para populações normais e não-normais

Análise Multivariada - 2022

102

Estimação da Probabilidade Global de Acerto

- Estimação da probabilidade global de acerto da função discriminante:

$$\hat{p}(\text{acerto}) = \frac{n_{11} + n_{22}}{n_1 + n_2}.$$

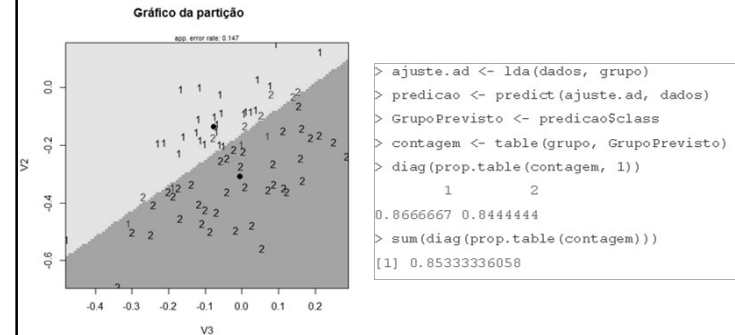
- ✓ Recomendável estimar as probabilidades de ocorrência dos erros de classificação tipo 1 e 2
 - Possível função discriminante com alta probabilidade de acerto global, mas apresentando alta probabilidade de algum dos erros parciais.

Análise Multivariada - 2022

103

Exemplo

- Conjunto de dados: Hemofilia



Análise Multivariada - 2022

107

Construção da Regra de Discriminação: Caso de Várias Populações

Regra de Discriminação para Várias Populações

- Classificar unidades amostrais em $g > 2$ populações
 - ✓ $f_i(x)$: função de densidade π_i , $i = 1, 2, \dots, g$.
- Objetivo:
 - ✓ Construir regra de classificação que minimize as probabilidades de erros de classificação

Análise Multivariada - 2022

110

Procedimento

- Para um vetor de observações \mathbf{x} :
 - ✓ Calcula-se o valor de $f_i(\mathbf{x})$, para cada i .
 - ✓ Classifica-se o elemento amostral na população k correspondente ao maior valor $f_i(\mathbf{x})$.
- No caso de população normal multivariada, corresponde a classificar na população k , tal que:

$$d_k^Q(\mathbf{x}) = \max\{d_1^Q(\mathbf{x}), d_2^Q(\mathbf{x}), \dots, d_g^Q(\mathbf{x})\}$$

- ✓ Sendo d_i^Q : escore quadrático de discriminação

$$d_i^Q = -\frac{1}{2} \ln(|\Sigma_i|) - \frac{1}{2}(\mathbf{x} - \mu_i)' \Sigma_i^{-1}(\mathbf{x} - \mu_i)$$

Análise Multivariada - 2022

111

- Na prática, os escores quadrático de discriminação são estimados por:

$$\hat{d}_i^Q = -\frac{1}{2} \ln(|S_i|) - \frac{1}{2}(\mathbf{x} - \bar{\mathbf{x}}_i)' S_i^{-1}(\mathbf{x} - \bar{\mathbf{x}}_i)$$

- Se $\Sigma_1 = \Sigma_2 = \dots = \Sigma_g = \Sigma$, usa-se S_{pol} para estimar \hat{d}_i^Q :

$$S_{\text{pol}} = \frac{1}{\sum_{i=1}^g n_i - g} \sum_{i=1}^g (n_i - 1) S_i$$

Análise Multivariada - 2022

112

Erro de Classificação

- Elemento amostral pertence a π_j , mas a regra de discriminação o classifica em π_k , $j, k = 1, 2, \dots, g$, $j \neq k$.
- Erros estimados por:

$$\hat{p}(k|j) = \frac{n_{jk}}{n_j}$$

- ✓ n_{jk} : número de elementos de π_j classificados em π_k .

Análise Multivariada - 2022

114

Exemplo

- Cultivares de Vinho: 178 tipos de vinhos

✓ Class: cultivares de vinho

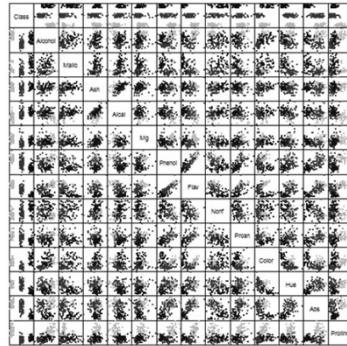
✓ Variáveis:

- | | |
|-----------|------------|
| - Alcohol | - Flav: |
| - Malic: | - Nonf |
| - Ash | - Proan |
| - Alcal | - Color |
| - Mg | - Hue |
| - Phenol | - Abs |
| | - Proline: |

Análise Multivariada - 2022

117

- Matrix scatter plot:



- Há variáveis que diferem entre os grupos

Análise Multivariada - 2022

118

- Médias das variáveis por grupo

```
> round(t(aggregate(vinho, list(vinho$Class),
  mean)[,2:14]), 2)
```

	[,1]	[,2]	[,3]
Class	1.00	2.00	3.00
Alcohol	13.74	12.28	13.15
Malic	2.01	1.93	3.33
Ash	2.46	2.24	2.44
Alcal	17.04	20.24	21.42
Mg	106.34	94.55	99.31
Phenol	2.84	2.26	1.68
Flav	2.98	2.08	0.78
Nonf	0.29	0.36	0.45
Proan	1.90	1.63	1.15
Color	5.53	3.09	7.40
Hue	1.06	1.06	0.68
Abs	3.16	2.79	1.68

- Abordagem marginal é intuitiva para identificar médias dos grupos

✓ Falha em identificar características para discriminar os indivíduos

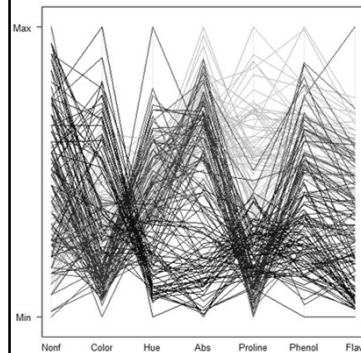
✓ Abordagem univariada não considera correlação entre variáveis

- Pode haver combinação de variáveis que fornecem nível mais alto de discriminação

Análise Multivariada - 2022

119

- Parallel coordinate plot

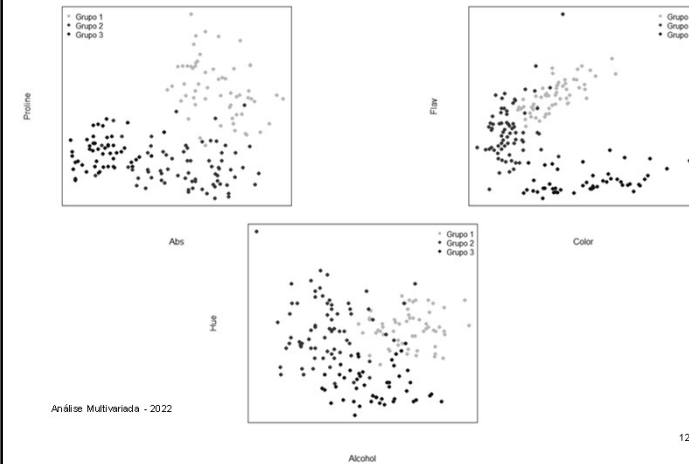


- Há alguma sobreposição entre grupos

Análise Multivariada - 2022

120

- ✓ Prováveis pares de variáveis úteis para discriminação



Análise Multivariada - 2022

121

• Análise discriminante

```
> ajuste.ld <- lda(Class~., data = vinho)
> ajuste.ld
Call:
lda(Class ~ ., data = vinho)

Prior probabilities of groups:
 1      2      3 
0.3314607 0.3988764 0.2696629
```

```
Group means:
  Alcohol  Malic  Ash  Alcal  Mg  Phenol
Flav  Nonf
1 13.74475 2.010678 2.455593 17.03729 106.3390 2.840169
2 2.9823729 0.290000
12 12.27873 1.932676 2.244789 20.23803 94.5493 2.258873
2.0808451 0.363662
13 13.15375 3.333750 2.437083 21.41667 99.3125 1.678750
0.7814583 0.447500
  Proan  Color  Hue  Abs  Proline
1 1.899322 5.528305 1.0620339 3.157797 1115.7119
2 1.630282 3.086620 1.0562817 2.785352 519.5070
1.153542 7.396250 0.6827083 1.683542 629.8958
```

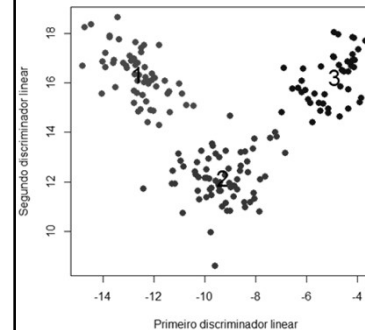
```
Coefficients of linear discriminants:
              LD1      LD2
Alcohol -0.403399781  0.8717930699
Malic    0.165254596  0.3053797325
Ash      -0.369075256  2.3458497486
Alcal    0.154797889  -0.1463807654
Mg       -0.002163496  -0.0004627565
Phenol   0.618052068  -0.0322128171
Flav     -1.661191235  -0.4919980543
Nonf     -1.495818440  -1.6309537953
Proan    0.134092628  -0.3070875776
Color    0.355055710  0.2532306865
Hue      -0.818036073  -1.5156344987
Abs      -1.157599376  0.0511839665
Proline  -0.002691206  0.0028529846

Proportion of trace:
              LD1      LD2
0.6875 0.3125
```

Análise Multivariada - 2022

122

• Discriminadores lineares dos dados



- Discriminadores:
 - ✓ Combinações lineares que melhor discriminam os grupos
- Figura indica distinção clara entre os grupos

Análise Multivariada - 2022

123

Outros Métodos de Discriminação

Métodos de Discriminação

- Método do vizinho mais próximo
- Classification and Regression Trees – CART
- Support Vector Machine – SVM
- Método dos núcleos estimadores
- Redes neurais artificiais

Análise Multivariada - 2022

141

Método do Vizinho mais Próximo

- *Nearest neighbor discriminant analysis*
 - ✓ Não é um método paramétrico
 - Não depende da suposição de normalidade multivariada
- Procedimento de discriminação:
 - ✓ Encontra-se vizinho mais próximo
 - (distância de Mahalanobis)
 - ✓ Classifica-se a observação na população do vizinho
- Variação:
 - ✓ Método dos k vizinhos mais próximos

Análise Multivariada - 2022

143

Árvores de Regressão e Classificação

- CART – Classification and Regression Trees
 - ✓ Trata simultaneamente variáveis contínuas e não contínuas
- Procedimento de classificação:
 - ✓ Baseado em informações das distribuições isoladas de cada variável
 - ✓ Resulta numa árvore com vários nós

Análise Multivariada - 2022

145

Support Vector Machine

- Grupo de algoritmos de classificação, incluem ampla variedade de modelos paramétricos e não paramétricos
 - ✓ Modelos lineares e métodos de regressão
 - ✓ Técnicas de suavização por núcleo
- Estimativa de probabilidade de classificação:
 - ✓ Validação cruzada ou por amostras de treinamento

Análise Multivariada - 2022

147

Referências

Bibliografia Recomendada

- JOHNSON, R. A.; WINCHERN, D. W. *Applied Multivariate Statistical Analysis*. Prentice Hall, 2007
- EVERITT, B.; HOTHORN, T. *An introduction to applied multivariate analysis with R*. Springer, 2011
- MINGOTI, S.A. *Análise de Dados através de Métodos de Estatística Multivariada*. Ed. UFMG, 2005.
- ZELTERMAN, D. *Applied Multivariate Statistics with R*. Springer, 2015.