



SME0822 Análise Multivariada e Aprendizado Não-Supervisionado

Aula 11a: **Análise discriminante linear (LDA)**

Prof. Cibeles Russo

cibele@icmc.usp.br

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

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

Mingoti, S. A. (2007) Análise de dados através de métodos de estatística multivariada: uma abordagem aplicada. Editora UFMG.

Análise discriminante

Objetivos:

- **Classificação** de elementos de uma amostra ou população,
- **Redução de dimensionalidade** de forma que se obtenha um bom classificador com o número mínimo de dimensões possível.

Difere da análise de agrupamentos pelo fato de que são determinados previamente os grupos aos quais serão direcionados os elementos da amostra.

A análise discriminante pode ser considerada uma técnica de aprendizado supervisionado, enquanto que a análise de agrupamentos é uma técnica de aprendizado não-supervisionado

Análise discriminante

Objetivos:

- **Classificação** de elementos de uma amostra ou população,
- **Redução de dimensionalidade** de forma que se obtenha um bom classificador com o número mínimo de dimensões possível.

Difere da análise de agrupamentos pelo fato de que são determinados previamente os grupos aos quais serão direcionados os elementos da amostra.

A análise discriminante pode ser considerada uma técnica de aprendizado supervisionado, enquanto que a análise de agrupamentos é uma técnica de aprendizado não-supervisionado

Análise discriminante

A análise discriminante linear (*linear discriminant analysis* ou LDA) foi proposta originalmente por **Sir Ronald Fisher** em **1936**, a princípio para duas classes.

Em 1948, **C. R. Rao** propôs uma **generalização para múltiplas classes**.

Análise discriminante

Suponha que tenhamos n_1 e n_2 elementos amostrais procedentes das populações A e B, respectivamente, e que em cada um dos $n = n_1 + n_2$ tenham sido observadas p características.

Se observarmos um novo elemento amostral, cuja origem é incerta, como compará-lo ao perfil geral dos grupos A e B e **classificá-lo** como pertencente a um deles?

Análise discriminante

Suponha que tenhamos n_1 e n_2 elementos amostrais procedentes das populações A e B, respectivamente, e que em cada um dos $n = n_1 + n_2$ tenham sido observadas p características.

Se observarmos um novo elemento amostral, cuja origem é incerta, como compará-lo ao perfil geral dos grupos A e B e **classificá-lo** como pertencente a um deles?

Análise discriminante

Exemplo:

Em medicina, é comum querer identificar fatores de risco ou distinguir doenças que tenham alguma similaridade sintomática aos apresentados pelo paciente.

Ou seja, deseja-se identificar um perfil de portadores ou não de uma determinada doença para classificar novos pacientes como prováveis ou não prováveis portadores da patologia em questão.

Outros exemplos podem ser encontrados em áreas como educação, finanças, entre outras.

Análise discriminante

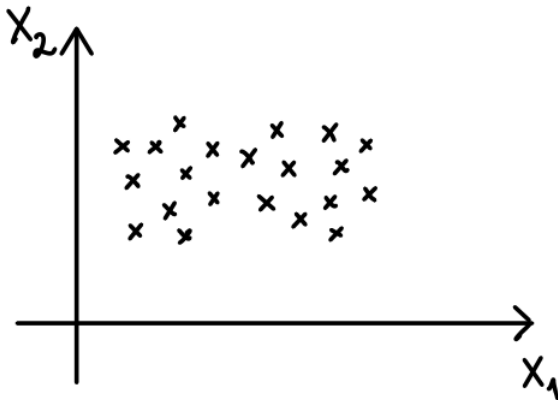
Exemplo:

Em medicina, é comum querer identificar fatores de risco ou distinguir doenças que tenham alguma similaridade sintomática aos apresentados pelo paciente.

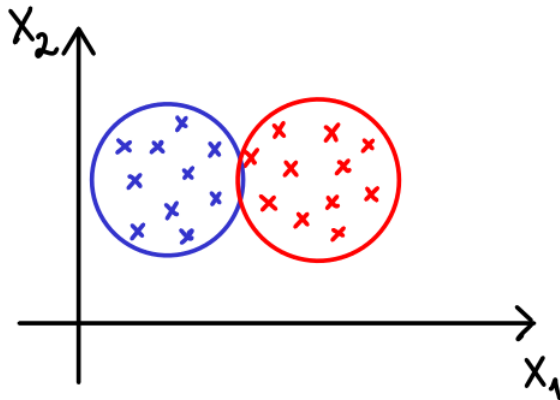
Ou seja, deseja-se identificar um perfil de portadores ou não de uma determinada doença para classificar novos pacientes como prováveis ou não prováveis portadores da patologia em questão.

Outros exemplos podem ser encontrados em áreas como educação, finanças, entre outras.

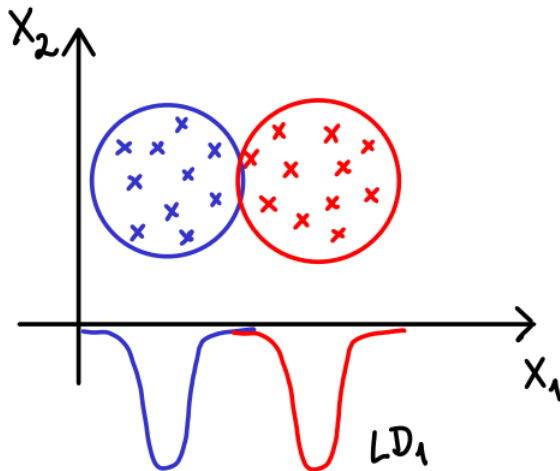
Análise discriminante - Exemplo



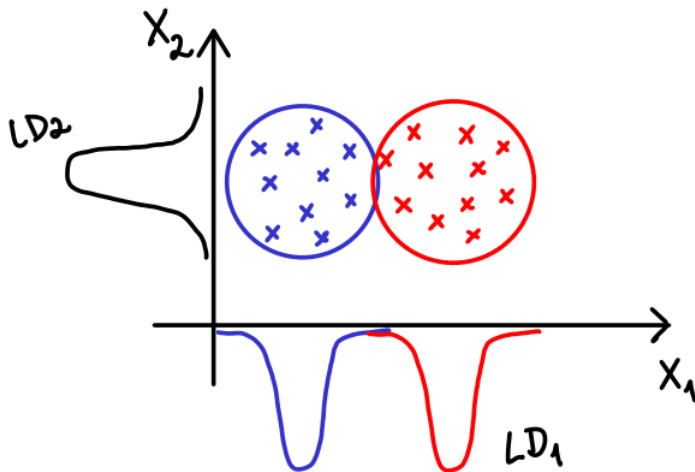
Análise discriminante - Exemplo



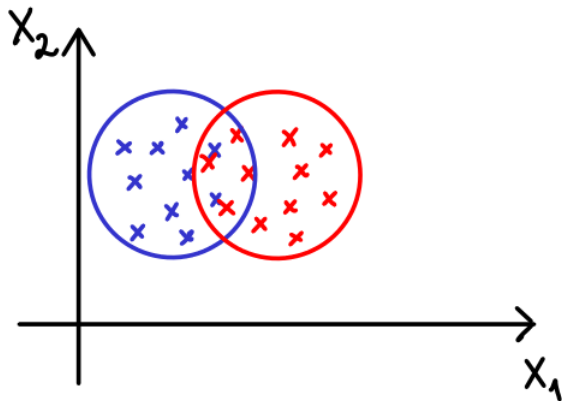
Análise discriminante - Exemplo



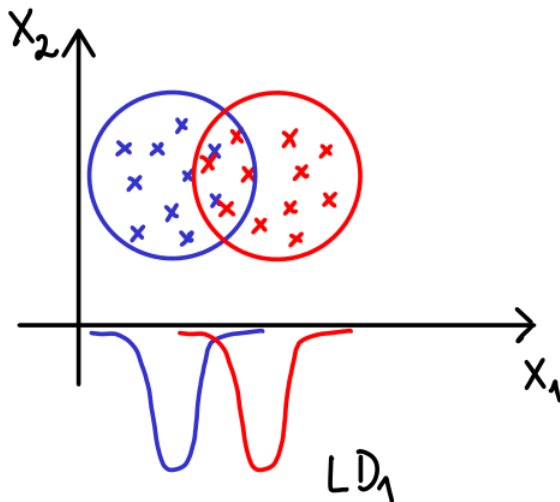
Análise discriminante - Exemplo



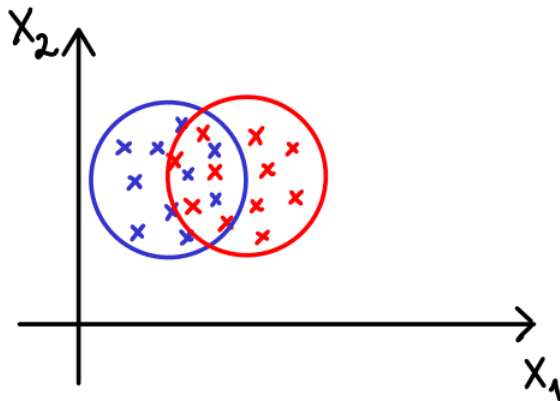
Análise discriminante - Exemplo



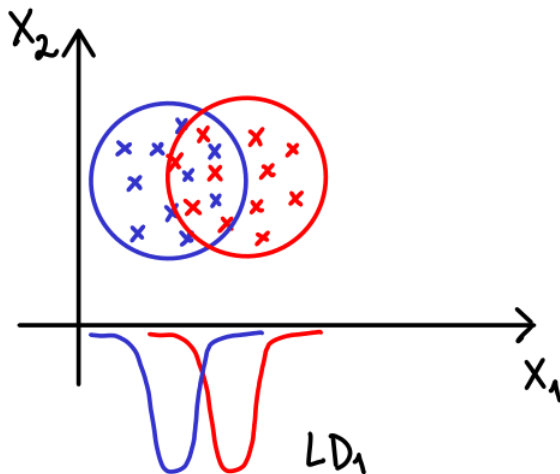
Análise discriminante - Exemplo



Análise discriminante - Exemplo



Análise discriminante - Exemplo



Análise discriminante

Um bom procedimento de classificação resulta em poucas classificações incorretas.

Se a distribuição de probabilidade das características medidas nos elementos amostrais de cada população for conhecida, é possível utilizar o princípio da máxima verossimilhança (Casella e Berger, 2002) para construir uma regra de classificação que minimize a chance de classificar um elemento amostral incorretamente.

Casella, G.; Berger, R. L., (2002) Statistical inference. Pacific Grove, CA: Duxbury.

Análise discriminante

Suponha que uma escola adote um processo seletivo de duas fases. Seja X a nota na prova de Matemática de candidatos na fase 1, e considere duas populações de alunos:

População 1: Alunos que passaram na 1^a fase mas não foram aprovados na 2^a fase.

População 2: Alunos que passaram em ambas as fases do vestibular.

Análise discriminante

A partir dos dados, deseja-se criar uma regra de classificação que permita identificar, dentre os aprovados na primeira fase, quais provavelmente serão aprovados na segunda fase. Suponha, no caso univariado, que

População 1: $X \sim N(\mu_1, \sigma^2)$

População 2: $X \sim N(\mu_2, \sigma^2)$.

Agora, para cada possível nota x de um candidato, pode-se calcular uma razão de probabilidades

$$\lambda(x) = \frac{f_1(x)}{f_2(x)},$$

que indica a razão de densidades de x na população 1 e 2.

Análise discriminante

A partir dos dados, deseja-se criar uma regra de classificação que permita identificar, dentre os aprovados na primeira fase, quais provavelmente serão aprovados na segunda fase. Suponha, no caso univariado, que

População 1: $X \sim N(\mu_1, \sigma^2)$

População 2: $X \sim N(\mu_2, \sigma^2)$.

Agora, para cada possível nota x de um candidato, pode-se calcular uma razão de probabilidades

$$\lambda(x) = \frac{f_1(x)}{f_2(x)},$$

que indica a razão de densidades de x na população 1 e 2.

Análise discriminante

Se f_1 e f_2 são densidades da distribuição normal como suposto anteriormente, temos

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

que pode ser simplificado por

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

Análise discriminante

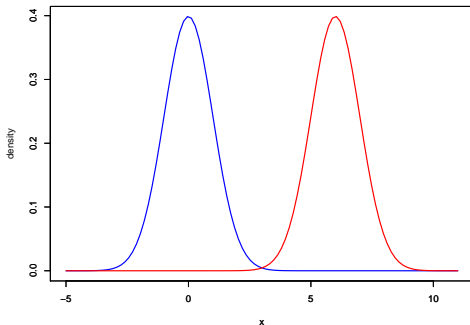
Se $\lambda(x) > 1$, então é razoável classificar o candidato como um provável não aprovado na 2^a fase.

Se $\lambda(x) < 1$, ele é um provável aprovado em ambas as fases.

Se $\lambda(x) = 1$, então as probabilidades são as mesmas de estar na população 1 e na população 2, segundo esse critério.

Análise discriminante

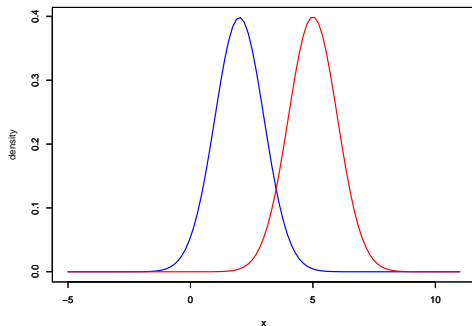
A qualidade da discriminação dependerá do grau de intersecção das duas densidades.



Casos como esse podem ocasionar poucas ou nenhuma classificação incorreta, ou seja, há **forte poder de discriminação**.

Análise discriminante

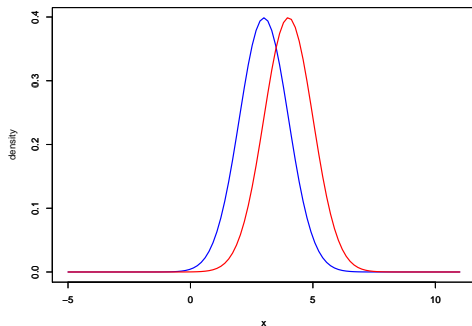
A qualidade da discriminação dependerá do grau de intersecção das duas densidades.



Casos como esse podem ocasionar poucas classificações incorretas, ou seja, há **poder razoável de discriminação**.

Análise discriminante

A qualidade da discriminação dependerá do grau de intersecção das duas densidades.



Casos como esse podem ocasionar muitas classificações incorretas, ou seja, há **poder fraco de discriminação**.

Análise discriminante

É comum considerar $-2 \log \lambda(x)$ com as seguintes correspondências:

$\lambda(x)$	$-2 \log \lambda(x)$	Situação
> 1	< 0	x mais próximo de μ_1
< 1	> 0	x mais próximo de μ_2
$= 1$	$= 0$	x igualmente próximo de μ_1 e μ_2

Análise discriminante

As funções $\lambda(x)$ e $-2 \log \lambda(x)$ são chamadas de **funções discriminantes**.

Para o caso multivariado, em que $\underline{X} \sim N(\underline{\mu}_1, \Sigma_1)$ na população 1 e $\underline{X} \sim N(\underline{\mu}_2, \Sigma_2)$ na população 2, temos

$$-2 \log \lambda(\underline{x}) = \log \left[\frac{(2\pi)^{p/2} |\Sigma_1|^{-1/2} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{\mu}_1)^\top \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1) \right\}}{(2\pi)^{p/2} |\Sigma_2|^{-1/2} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{\mu}_2)^\top \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2) \right\}} \right]$$

Análise discriminante

As funções $\lambda(x)$ e $-2 \log \lambda(x)$ são chamadas de **funções discriminantes**.

Para o caso multivariado, em que $\underline{X} \sim N(\underline{\mu}_1, \Sigma_1)$ na população 1 e $\underline{X} \sim N(\underline{\mu}_2, \Sigma_2)$ na população 2, temos

$$-2 \log \lambda(\underline{x}) = \log \left[\frac{(2\pi)^{p/2} |\Sigma_1|^{-1/2} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{\mu}_1)^\top \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1) \right\}}{(2\pi)^{p/2} |\Sigma_2|^{-1/2} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{\mu}_2)^\top \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2) \right\}} \right]$$

Análise discriminante

ou seja,

$$\begin{aligned} -2 \log \lambda(\underline{x}) &= (\underline{x} - \underline{\mu}_1)^\top \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1) - (\underline{x} - \underline{\mu}_2)^\top \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2) \\ &\quad + \log |\Sigma_1| - \log |\Sigma_2| \end{aligned}$$

e então classificamos \underline{x} na população 1 se $-2 \log \lambda(\underline{x}) < 0$ e na população 2 se $-2 \log \lambda(\underline{x}) > 0$.

Análise discriminante

Quando $\Sigma_1 = \Sigma_2 = \Sigma$, temos a **função discriminante de Fisher**:

$$fd(\underline{x}) = (\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1} \underline{x} - \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1}(\underline{\mu}_1 + \underline{\mu}_2).$$

Um elemento amostral com vetor de observações \underline{x} seria classificado na população 1 se $fd(\underline{x}) > 0$, ou seja, se

$$(\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1} \underline{x} > \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1}(\underline{\mu}_1 + \underline{\mu}_2)$$

e seria classificado na população 2 se

$$(\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1} \underline{x} < \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1}(\underline{\mu}_1 + \underline{\mu}_2).$$

Análise discriminante

Quando $\Sigma_1 = \Sigma_2 = \Sigma$, temos a **função discriminante de Fisher**:

$$fd(\underline{x}) = (\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1} \underline{x} - \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1}(\underline{\mu}_1 + \underline{\mu}_2).$$

Um elemento amostral com vetor de observações \underline{x} seria classificado na população 1 se $fd(\underline{x}) > 0$, ou seja, se

$$(\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1} \underline{x} > \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1}(\underline{\mu}_1 + \underline{\mu}_2)$$

e seria classificado na população 2 se

$$(\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1} \underline{x} < \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)^\top \Sigma^{-1}(\underline{\mu}_1 + \underline{\mu}_2).$$

Estimação da função discriminante

Se as matrizes de variâncias e covariâncias populacionais Σ_1 e Σ_2 são desconhecidas, como também as médias populacionais μ_1 e μ_2 , mas é possível calcular as médias amostrais \bar{x}_1 e \bar{x}_2 e as matrizes de variâncias e covariâncias S_1 e S_2 , para as populações 1 e 2, respectivamente, estima-se a função discriminante por

$$\begin{aligned} -2 \log \hat{\lambda}(\underline{x}) &= (\underline{x} - \bar{x}_1)^\top S_1^{-1} (\underline{x} - \bar{x}_1) - (\underline{x} - \bar{x}_2)^\top S_2^{-1} (\underline{x} - \bar{x}_2) \\ &\quad + \log |S_1| - \log |S_2| \end{aligned}$$

Funções discriminantes lineares de Fisher

Considere que p variáveis \underline{X} foram observadas em elementos amostrais de g populações distintas, e que não seja possível supor normalidade dos dados, mas que seja razoável assumir que

$$\Sigma_1 = \Sigma_2 = \dots = \Sigma_g = \Sigma.$$

Sejam

- $\underline{\mu}_i$ o vetor de médias populacionais no grupo i
- $\underline{\bar{\mu}}$ o vetor de médias das médias populacionais

Funções discriminantes lineares de Fisher

Considere a soma de produtos cruzados entre grupos

$$B_{\mu} = \sum_{i=1}^g (\underline{\mu}_i - \underline{\bar{\mu}})(\underline{\mu}_i - \underline{\bar{\mu}})^{\top}$$

e a combinação linear

$$Y = \underline{a}^{\top} \underline{X}.$$

Funções discriminantes lineares de Fisher

Sejam π_1, \dots, π_g variáveis indicadoras da população i à qual pertence uma determinada observação. Logo,

$$\mu_{i,Y} = E(Y|\pi_i) = \underline{a}^\top E(\underline{X}|\pi_i) = \underline{a}^\top \underline{\mu}_i \text{ e}$$

$$\text{Var}(Y|\pi_i) = \underline{a}^\top \text{Cov}(\underline{X})\underline{a} = \underline{a}^\top \Sigma \underline{a}$$

para todas as populações.

Funções discriminantes lineares de Fisher

Ou seja,

$\mu_{i,Y} = \underline{a}^\top \underline{\mu}_i$ muda conforme a população à qual a observação pertence.

A média geral de Y é dada por

$$\bar{\mu}_Y = E(Y) = \frac{1}{g} \sum_{i=1}^g \mu_{i,Y} = \frac{1}{g} \sum_{i=1}^g \underline{a}^\top \underline{\mu}_i = \underline{a}^\top \frac{1}{g} \sum_{i=1}^g \underline{\mu}_i$$

Funções discriminantes lineares de Fisher

Considere a razão

$$\frac{(\text{soma dos quadrados das distâncias das populações à média geral})}{(\text{variância de } Y)} =$$

$$\frac{\sum_{i=1}^g (\mu_{i,Y} - \bar{\mu}_Y)^2}{\sigma_Y^2} = \frac{\sum_{i=1}^g (\underline{a}^\top \underline{\mu}_i - \underline{a}^\top \bar{\mu}_i)^2}{\underline{a}^\top \Sigma \underline{a}} = \frac{\underline{a}^\top B_\mu \underline{a}}{\underline{a}^\top \Sigma \underline{a}}$$

A razão acima mede a variabilidade entre grupos sobre a variabilidade intra grupos. O objetivo é buscar \underline{a} que maximize essa razão para obter a maior discriminação possível.

Funções discriminantes lineares de Fisher

Em geral, como μ_i e Σ são desconhecidos, utilizamos as estimativas \bar{x}_i e S_i 's, para $i = 1, \dots, g$.

Assim

$$B = \sum_{i=1}^g (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})^\top$$

e

$$W = \sum_{i=1}^g (n_i - 1)S_i.$$

Funções discriminantes lineares de Fisher

Note que $W = (n_1 + n_2 + \dots + n_g - g)S_{pooled}$, então a mesma constante $\hat{\underline{a}}$ que maximiza $\frac{\hat{\underline{a}}^\top B \hat{\underline{a}}}{\hat{\underline{a}}^\top S_{pooled} \hat{\underline{a}}}$ maximiza $\frac{\hat{\underline{a}}^\top B \hat{\underline{a}}}{\hat{\underline{a}}^\top W \hat{\underline{a}}}$.

É possível mostrar que o autovetor $\hat{\underline{e}}_1$ referente ao maior autovalor λ_1 de $W^{-1}B$ levam ao máximo de $\frac{\hat{\underline{a}}^\top B \hat{\underline{a}}}{\hat{\underline{a}}^\top W \hat{\underline{a}}}$.

Funções discriminantes lineares de Fisher

Discriminantes lineares de Fisher

Sejam $\hat{\lambda}_1 > \dots > \hat{\lambda}_s > 0$ os autovalores de $W^{-1}B$ e $\hat{e}_1, \dots, \hat{e}_s$ os autovetores ortonormais correspondentes. Então o vetor de coeficientes \hat{a} que maximiza a razão $\frac{\hat{a}^\top B \hat{a}}{\hat{a}^\top W \hat{a}}$ é $\hat{a}_1 = \hat{e}_1$.

A combinação linear $\hat{a}_1^\top X$ é chamada de **primeiro discriminante linear**.
A combinação linear $\hat{a}_2^\top X$ com $\hat{a}_2 = \hat{e}_2$ é chamada de **segundo discriminante linear** e assim por diante.

Como usar os discriminantes para classificar observações?

Seja o k -ésimo discriminante linear amostral,

$$\hat{Y}_k = \hat{a}_k^\top \underline{X}.$$

Uma possibilidade para classificar uma observação \underline{x} na l -ésima população é utilizar o k -ésimo discriminante linear amostral fazendo a verificação de que

$$(\hat{y}_k - \bar{y}_l)^2 \leq (\hat{y}_k - \bar{y}_i)^2 \text{ para todo } i \neq l \text{ ou}$$

$$(\hat{a}_k^\top \underline{x} - \hat{a}_k^\top \bar{\underline{x}}_l)^2 \leq (\hat{a}_k^\top \underline{x} - \hat{a}_k^\top \bar{\underline{x}}_i)^2 \text{ para todo } i \neq l.$$

Como usar os discriminantes para classificar observações?

Seja o k -ésimo discriminante linear amostral,

$$\hat{Y}_k = \hat{a}_k^\top \underline{X}.$$

Uma possibilidade para classificar uma observação \underline{x} na l -ésima população é utilizar o k -ésimo discriminante linear amostral fazendo a verificação de que

$$(\hat{y}_k - \bar{y}_l)^2 \leq (\hat{y}_k - \bar{y}_i)^2 \text{ para todo } i \neq l \text{ ou}$$

$$(\hat{a}_k^\top \underline{x} - \hat{a}_k^\top \bar{\underline{x}}_l)^2 \leq (\hat{a}_k^\top \underline{x} - \hat{a}_k^\top \bar{\underline{x}}_i)^2 \text{ para todo } i \neq l.$$

Como usar os discriminantes para classificar observações?

Seja o k -ésimo discriminante linear amostral,

$$\hat{Y}_k = \hat{a}_k^\top \underline{X}.$$

Uma possibilidade para classificar uma observação \underline{x} na l -ésima população é utilizar o k -ésimo discriminante linear amostral fazendo a verificação de que

$$(\hat{y}_k - \bar{y}_l)^2 \leq (\hat{y}_k - \bar{y}_i)^2 \text{ para todo } i \neq l \text{ ou}$$

$$(\hat{a}_k^\top \underline{x} - \hat{a}_k^\top \bar{\underline{x}}_l)^2 \leq (\hat{a}_k^\top \underline{x} - \hat{a}_k^\top \bar{\underline{x}}_i)^2 \text{ para todo } i \neq l.$$