

Objetivos do Curso

- Fornecer uma ferramenta estatística unificada que permita a análise de qualquer modelo que puder ser escrito na forma linear nos parâmetros. Este estudo englobaria Análise de Variância, Análise de Covariância, Correlação e Regressão
- Estudo de alguns Modelos mais específicos, que saem como casos particulares desta abordagem; Curvas de Crescimento (com pouco esforço, conhecida a formulação linear geral $\tilde{Y} = X\beta + \varepsilon$, entende-se este tópico)

Modelos com medidas repetidas (modelos mistos)

Regressão Multivariada

- Fornecer um estudo das principais distribuições utilizadas na Análise de Variância e Regressão
- Apresentar resultados matriciais úteis para estatísticos

1 - Introdução e Conceitos Básicos

1.1 - Introduções - Modelo Linear Geral

Um modelo (estatístico) linear geral é dado por

$$Y = \mu(x_1, x_2, \dots, x_K) + \epsilon$$

onde

Y e ϵ são variáveis aleatórias

x_1, x_2, \dots, x_K são variáveis não aleatórias

$\mu(x_1, x_2, \dots, x_K)$ é uma função de x_1, x_2, \dots, x_K , definida num domínio D e linear num conjunto de parâmetros desconhecidos, $\beta_0, \beta_1, \beta_2, \dots, \beta_K$ pertencentes a um espaço paramétrico Ω_p

ϵ - erro não observável, são feitas suposições sobre sua distribuição de probabilidades

$\mu \rightarrow$ parte determinística

$\epsilon \rightarrow$ parte aleatória ou estocástica do modelo

Y - variável resposta ou variável dependente

$x_1, x_2 \dots x_K$ - variáveis independentes, explicativas ou preditoras

a forma de $M(x_1, x_2 \dots x_K)$ será admitida conhecida

Ex:

$$M(x) = \beta_0 + \beta_1 x$$

$$Y = \beta_0 + \beta_1 x + \epsilon$$

→ modelo de regressão linear
simples

$$E(\epsilon) = 0$$

$$\text{Var}(\epsilon) = \sigma^2 \quad (\text{desconhecido}) \quad K=1$$

$$M(x_1, x_2, x_3) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon \quad \rightarrow \text{modelo de regressão linear
múltipla} \quad K=3$$

$$M(x_1, x_2, x_3) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

$$x_1 = x \quad x_2 = x^2 \quad x_3 = x^3$$

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \epsilon \quad \rightarrow \text{modelo de regressão
polinomial} \quad K=3$$

$$Y = \beta_0 + \beta_1 x + \beta_2 e^x + \epsilon$$

$\epsilon \sim N(0, \sigma^2)$ σ^2 desconhecido

$$x_1 = x \quad x_2 = e^x \quad k=2$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon$$

$$x_3 = x_1 x_2 \quad k=3$$

Na forma geral

$$\mu(x_1, x_2, \dots, x_k) = \beta_0 + \sum_{j=1}^k g_j(x_1, x_2, \dots, x_k) \beta_j$$

Como a distribuição de probabilidades de Y depende das variáveis explicativas, poderíamos escrever

$$Y_{(z)} = \mu(z) + \epsilon_{(z)} \quad E(\epsilon_{(z)}) = 0 \quad z \in D$$

$$\underline{z} = (z_1, z_2, \dots, z_k)$$

Obs:

- Este é um modelo populacional que representa a situação real

Obs:

Este é um modelo populacional que representa a situação em estudo. Nenhuma menção é feita à amostra ainda.

Para inferirmos sobre os parâmetros do modelo, tomamos uma amostra de n observações.

Modelo Linear Geral (amostral)

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_K x_{Ki} + \epsilon_i$$

$$E(\epsilon_i) = 0 \quad i = 1, 2, \dots, n$$

Consequência: $E(Y_i | x_{1i}, x_{2i}, \dots, x_{Ki}) = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_K x_{Ki}$

com

(1) Y_i v.a. observável

(2) x_{ji} v. não aleatórias pertencentes a um domínio D .

(3) β_j , $j = 0, 1, 2, \dots, K$, parâmetros definidos em um espaço paramétrico \mathcal{P}_β .

(4) ϵ_i v.a. não observáveis, $Cov(\epsilon_i, \epsilon_l) = \delta_{il}$

$$i, l = 1, 2, \dots, n$$

Na forma matricial, este modelo fica

$$\tilde{Y} = X \tilde{\beta} + \tilde{\epsilon}$$

onde

(1) \tilde{Y} é um vetor aleatório observável $n \times 1$

(2) X é uma matriz $n \times (k+1)$ de valores observáveis fixos
(os elementos de X não são variáveis aleatórias) $n > k+1$

(3) $\tilde{\beta}$ é um vetor $(k+1) \times 1$ de parâmetros desconhecidos,
definido num espaço paramétrico \mathcal{L}_{β}

(4) $\tilde{\epsilon}$ é um vetor aleatório não observável com

$$E(\tilde{\epsilon}) = \tilde{0} \quad \text{e} \quad \text{Var}(\tilde{\epsilon}) = \begin{bmatrix} \text{Var}(\epsilon_1) & \text{Cov}(\epsilon_1, \epsilon_2) & \dots & \text{Cov}(\epsilon_1, \epsilon_n) \\ \text{Cov}(\epsilon_2, \epsilon_1) & \text{Var}(\epsilon_2) & & \text{Cov}(\epsilon_2, \epsilon_n) \\ \vdots & & & \text{Var}(\epsilon_n) \end{bmatrix}$$

$$= \tilde{\epsilon}$$

$$\tilde{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1k} \\ 1 & x_{21} & & x_{2k} \\ \vdots & & & \vdots \\ 1 & x_{n1} & & x_{nk} \end{bmatrix}$$

$$\tilde{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$$

$$\tilde{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Este modelo admite vários casos especiais dependendo de

- i) A distribuição de probabilidades de ϵ .
- ii) A estrutura de covariância definida na matriz Σ .
- iii) O posto e a estrutura da matriz X .

(casos Particulares: Modelos de Regressão)

Modelos de Análise de Variância
Modelos Mistos

Exemplo 1

Y - Pressão Sanguínea

x - idade

$$Y = \beta_0 + \beta_1 x + \epsilon$$

$$E(\epsilon) = 0 \quad \text{Var}(\epsilon) = \sigma^2$$

para $50 \leq x \leq 60$

$$D = \{x \mid 50 \leq x \leq 60\}$$

Tomada uma amostra de 6 observações independentes
 $x_1 = 50 \quad x_2 = 52 \quad x_3 = 55 \quad x_4 = 56 \quad x_5 = 58 \quad x_6 = 60$
fixados, observa-se $y_1, y_2 \dots y_6$.

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i=1, 2 \dots 6$$

ϵ_i variáveis aleatórias independentes

$$\epsilon_i \sim N(0, \sigma^2)$$

\downarrow
sup. adicional

A partir da amostra \rightarrow inferência sobre β_0, β_1 e σ^2 .

$$\forall x \in D, \text{ admite-se que } E(Y|x) = \beta_0 + \beta_1 x.$$

$$Y|x \sim N(\beta_0 + \beta_1 x, \sigma^2)$$

$$\text{Por ex., } x = 54 \quad Y|x=54 \sim N(\beta_0 + 54\beta_1, \sigma^2)$$

$$E(Y|x=54) = \beta_0 + 54\beta_1$$

Podemos estimar a média de $Y|x=54$ sem amostrar esta distribuição

$$\hat{E}(Y|x=54) = \hat{\beta}_0 + \hat{\beta}_1 \cdot 54 \quad \hat{\beta}_0 \text{ e } \hat{\beta}_1 \text{ estimadores de } \beta_0 \text{ e } \beta_1.$$

Esta análise não é possível para $x=62$
porque $x=62 \notin D$.

No caso, $X = \begin{bmatrix} 1 & 50 \\ 1 & 52 \\ 1 & 55 \\ 1 & 56 \\ 1 & 58 \\ 1 & 60 \end{bmatrix}$ $\mathcal{Z} = 6^2 I_6$

Modelo de Regressão Linear Simples

É uma suposição forte $E(Y|x) = \beta_0 + \beta_1 x$
 $\forall x, 50 \leq x \leq 60$

$X_1, X_2 \dots X_K \rightarrow$ quantitativas; modelo de regressão

9

\rightarrow qualitativos; modelo de planejamento de
(atributos; cor, tipo de máquina,) experimentos
sexo, religião

Exemplo 2

y - Pressão sanguínea

x - sexo $n = 6$ 3 homens, 3 mulheres

$$Y_{ij} = \mu_i + \epsilon_{ij} \quad E(\epsilon_{ij}) = 0$$

$i = 1, 2$ (para cada sexo)

μ_1 - Pressão média dos Homens

$j = 1, 2, 3$ (para indivíduos)

μ_2 - Pressão média das mulheres

$$\tilde{Y} = X\beta + \tilde{\epsilon}$$

$\tilde{\epsilon} \sim N(0, \sigma^2)$

$$\begin{bmatrix} Y_{11} \\ Y_{12} \\ Y_{13} \\ Y_{21} \\ Y_{22} \\ Y_{23} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}_{6 \times 2} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}_{2 \times 1} + \begin{bmatrix} \epsilon_{11} \\ \epsilon_{12} \\ \epsilon_{13} \\ \epsilon_{21} \\ \epsilon_{22} \\ \epsilon_{23} \end{bmatrix}$$

mudar para σ^2

2 parâmetros

O objetivos: Estimar μ_1, μ_2 Testar $\mu_1 = \mu_2$

$$\mu_i = \mu + \tau_i \quad i=1,2 \quad \text{Forma alternativa de escrever o modelo}$$

μ - média geral

τ_i - efeito de sexo

$$Y_{ij} = \mu + \tau_i + \epsilon_{ij} \quad \beta$$

$$\begin{bmatrix} Y_{11} \\ Y_{12} \\ Y_{13} \\ Y_{21} \\ Y_{22} \\ Y_{23} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}_{6 \times 3} \begin{bmatrix} \mu \\ \tau_1 \\ \tau_2 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} \epsilon_{11} \\ \epsilon_{12} \\ \epsilon_{13} \\ \epsilon_{21} \\ \epsilon_{22} \\ \epsilon_{23} \end{bmatrix}$$

$$\begin{array}{l} \mu + \tau_1 \\ \mu + \tau_2 \\ \mu + \tau_3 \\ \mu - \tau_1 \\ \mu - \tau_2 \\ \mu - \tau_3 \end{array} \quad X = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}$$

$$\begin{array}{l} \mu + \tau_1 \\ \mu + \tau_2 \\ \mu + \tau_3 \\ \mu - \tau_1 \\ \mu - \tau_2 \\ \mu - \tau_3 \end{array} \quad \text{3 parâmetros}$$

Com relação à forma anterior

mais parâmetros

$\mu + \tau_i \rightarrow$ diferente forma de escrever μ_i (desnecessária neste caso pode ser útil em planejamentos mais complicados)

$X \quad 6 \times 3 \quad \text{posto} 2 = n^{\circ} \text{ de colunas linearmente independentes de } X$
 $\quad \quad \quad < 3 \rightarrow \text{posto incompleto}$

No outro caso: $X_{6 \times 2} \quad \text{posto}(X) = 2 \rightarrow \text{posto completo}$

Esta é a diferença entre os modelos dos tópicos 3 e 4

Nos dois casos, X é qualitativa: Modelo de planejamento de experimentos, a matriz X contém somente zeros e uns

Neste caso, a matriz de planejamento apresenta ainda padrões específicos, de acordo com o particular planejamento.

Exemplo 3 (Graybill, pag 165)

Y - Produção de Milho

em função de dois diferentes compostos e quatro diferentes métodos de aplicação.

$$Y_{ij} = \mu + \alpha_i + \tau_j + \epsilon_{ij} \quad i=1,2 \quad j=1,2,3,4.$$

Y_{ij} - produção de milho observada no canteiro que recebeu o i -ésimo composto aplicado pelo j -ésimo método.

μ - produção média quando não é aplicado nenhum composto e \therefore nenhum método

α_i - efeito do i -ésimo composto $i=1,2$

τ_j - efeito do j -ésimo método de aplicação $j=1,2,3,4$.

$$Y_{11} = \mu + \alpha_1 + \tau_1 + \epsilon_{11}$$

$$Y_{12} = \mu + \alpha_1 + \tau_2 + \epsilon_{12}$$

$$Y_{13} = \mu + \alpha_1 + \tau_3 + \epsilon_{13}$$

$$Y_{14} = \mu + \alpha_1 + \tau_4 + \epsilon_{14}$$

$$Y_{21} = \mu + \alpha_2 + \tau_1 + \epsilon_{21}$$

$$Y_{22} = \mu + \alpha_2 + \tau_2 + \epsilon_{22}$$

$$Y_{23} = \mu + \alpha_2 + \tau_3 + \epsilon_{23}$$

$$Y_{24} = \mu + \alpha_2 + \tau_4 + \epsilon_{24}$$

Na forma Matricial $Y = X\beta + \epsilon$

$$Y = \begin{bmatrix} Y_{11} \\ Y_{12} \\ Y_{13} \\ Y_{14} \\ Y_{21} \\ Y_{22} \\ Y_{23} \\ Y_{24} \end{bmatrix}_{8 \times 1} \quad X = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 1 \end{bmatrix}_{8 \times 7} \quad \beta = \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \tau_1 \\ \tau_2 \\ \tau_3 \\ \tau_4 \end{bmatrix}_{7 \times 1}$$

Modelo de Análise de Variância com dois fatores, X de posto incompleto posto ($X = 5$)

Definido o Modelo Linear Geral estudaremos
Inferência sobre β para modelos com matriz
de postos completo e incompleto.

Em modelos com matriz de planejamento
de postos incompleto, nem todos os elementos
de β podem ser estimados. Veremos então
nestes modelos quais funções lineares de
 β poderão ser estimadas.