

## 1. ESTIMAÇÃO

Um **estimador**,  $\hat{\theta}$ , do parâmetro  $\theta$  é uma função qualquer dos elementos da amostra. **Estimativa** é o valor numérico assumido pelo estimador quando os valores observados são considerados. Os **parâmetros** são funções de valores populacionais, enquanto **estatísticas** são funções dos dados amostrais.

### Propriedades dos estimadores

#### 1) Não tendenciosidade

Um estimador  $\hat{\theta}$  do parâmetro  $\theta$  é dito um estimador não tendencioso se  $E(\hat{\theta}) = \theta$ .

**Exemplo:**  $\bar{X} = \sum_i^n X_i / n$  é um estimador não tendencioso da média populacional  $\mu$ .

**Prova**

$$\begin{aligned} E(\bar{X}) &= E\left(\sum_i^n X_i / n\right) = \frac{1}{n} E\left(\sum_i^n X_i\right) = \frac{1}{n} E(X_1 + X_2 + \dots + X_n) = \\ &= \frac{1}{n} (E(X_1) + E(X_2) + \dots + E(X_n)) = \frac{1}{n} (\mu + \mu + \dots + \mu) = \frac{1}{n} (n\mu) = \mu \end{aligned}$$

#### 2) Consistência

Um estimador  $\hat{\theta}$  é dito um estimador consistente do parâmetro  $\theta$  se:

i)  $\lim_{n \rightarrow \infty} E(\hat{\theta}) = \theta$ ;

ii)  $\lim_{n \rightarrow \infty} V(\hat{\theta}) = 0$ .

**Exemplo:**  $\bar{X} = \sum_i^n X_i / n$

i)  $\lim_{n \rightarrow \infty} E(\bar{X}) = \lim_{n \rightarrow \infty} \mu = \mu$ ;

ii)  $\lim_{n \rightarrow \infty} V(\bar{X}) = \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} = 0$ .

#### 3) Eficiência

Se  $\hat{\theta}_1$  e  $\hat{\theta}_2$  são dois estimadores não tendenciosos de  $\theta$ , então,  $\hat{\theta}_1$  é mais eficiente que  $\hat{\theta}_2$  se:

$$V(\hat{\theta}_1) < V(\hat{\theta}_2).$$

#### 4) Suficiência ou Precisão

Um estimador é suficiente se contém o máximo de informação com relação ao parâmetro por ele estimado.

$$\text{Quantidade de informação ou precisão} = \frac{1}{V(\hat{\theta})}$$

## MÉTODOS DE ESTIMAÇÃO

### 1. Estimadores de momentos (Método dos Momentos)

Considere uma variável aleatória (v.a.)  $X$  com densidade  $f(x; \theta_1, \theta_2, \dots, \theta_r)$ , a qual depende de  $r$  parâmetros. O primeiro momento populacional é definido como  $\mu_1 = E(X)$ . O  $k$ -ésimo momento populacional é  $\mu_k = E(X^k)$ .

No caso de uma distribuição normal,  $X \sim N(\mu, \sigma^2)$ , temos  $k = 2$ ,  $\theta_1 = \mu$  e  $\theta_2 = \sigma^2$ . O primeiro momento é  $\mu_1 = E(X)$ . O segundo momento é  $\mu_2 = E(X^2)$ .

Sabemos que  $Var(X) = E(X^2) - (E(X))^2$ . Portanto o segundo momento populacional é  $E(X^2) = Var(X) + (E(X))^2 = \sigma^2 + \mu^2$ .

Suponha que seja retirada uma amostra  $x_1, x_2, \dots, x_n$  dessa população. Definem-se então os momentos amostrais como sendo, o  $k$ -ésimo momento  $m_k = \frac{1}{n} \sum_{i=1}^n x_i^k$ . Tem-se, portanto, que o primeiro momento amostral corresponde à média amostral,  $m_1 = \bar{x}$ .

**Definição:** Os estimadores obtidos pelo método dos momentos são soluções da equação:  $m_k = \mu_k$ . O procedimento de obter os estimadores consiste em substituir os momentos teóricos pelos respectivos momentos amostrais.

**Exemplo.** Considere uma v.a.  $X$  com distribuição normal e parâmetros  $\mu$  e  $\sigma^2$ . Encontre os estimadores de momentos para os dois parâmetros.

### 2. Estimadores de mínimos quadrados

O método dos mínimos quadrados consiste em encontrar o estimador para o parâmetro, que minimize a soma dos quadrados dos resíduos do modelo em questão. Isso é feito derivando a equação do resíduo em relação a cada um dos parâmetros, igualando essas derivadas a zero e resolvendo as equações resultantes.

Como exemplo, considere o modelo estatístico de uma regressão linear simples:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

Os resíduos e a soma de quadrados de resíduos são obtidos do modelo de regressão:

$$\varepsilon_i = y_i - \beta_0 - \beta_1 x_i ; \quad Q = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

O método dos mínimos quadrados consiste em minimizar a soma de quadrados do erro ou resíduo do modelo ao longo de todos os  $n$  pares  $(x_i, y_i)$ . Para isto, derivamos a expressão (Q) em relação aos parâmetros  $\beta_0$  e  $\beta_1$  obtendo-se assim o Sistema de Equações Normais (SEN):

$$(SEN) \begin{cases} \frac{\partial Q}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \\ \frac{\partial Q}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i \end{cases}$$

Igualando essas derivadas a zero e substituindo  $\beta_0$  e  $\beta_1$  pelos respectivos estimadores  $\hat{\beta}_0$  e  $\hat{\beta}_1$  tem-se:

$$\begin{cases} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \end{cases} \Leftrightarrow \begin{cases} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \quad (A) \\ \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \quad (B) \end{cases} \Leftrightarrow \begin{cases} \sum_{i=1}^n y_i - n \hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n x_i = 0 \\ \sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0 \end{cases}$$

Resolvendo esse sistema para os parâmetros, tem-se os estimadores de mínimos quadrados para  $\beta_0$  e  $\beta_1$ , respectivamente:

$$\hat{\beta}_0 = b_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \text{e} \quad \hat{\beta}_1 = b_1 = \frac{SP_{xy}}{S_{xx}}$$

### 3. Estimadores de Máxima Verossimilhança

Esse princípio baseia-se em escolher o estimador que maximiza a probabilidade de se obter a amostra particular observada.

Considere uma v.a.  $X$  com densidade  $f(x; \theta)$ , sendo  $\theta$  o parâmetro desconhecido, ou o vetor de parâmetros. Retirando-se dessa população uma amostra de tamanho  $n$  ( $x = x_1, x_2, \dots, x_n$ ), define-se a função de verossimilhança dos dados por,

$$L(\theta; x_1, x_2, \dots, x_n) = L(\theta; x) = L(\theta | x) = f(x_1; \theta) \cdot f(x_2; \theta) \cdot \dots \cdot f(x_n; \theta)$$

O procedimento para obter os estimadores de máxima verossimilhança consiste em obter a função de verossimilhança e maximizar essa função  $L(\theta; x)$ , ou o seu logaritmo  $\log(L(\theta; x)) = l(\theta; x)$ , que geralmente é mais conveniente.

**Exemplo.** Considere uma v.a.  $X$  com distribuição exponencial, com parâmetro  $\beta$ . Encontre o estimador de máxima verossimilhança para o parâmetro dessa distribuição.

## 2. MODELO DE REGRESSÃO LINEAR MÚLTIPLA

A análise de regressão múltipla é o estudo de como a variável dependente  $y$  se relaciona com duas ou mais variáveis independentes. Portanto, temos uma regressão linear múltipla quando admitimos que o valor da variável dependente é função linear de duas ou mais variáveis independentes. O modelo estatístico de uma regressão múltipla com  $k$  variáveis independentes é:

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \varepsilon_i$$
$$\text{ou } y_i = \beta_0 + \sum_{j=1}^k \beta_j x_{ij} + \varepsilon_i$$
$$\text{ou } y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

em que,  $y$  é a variável dependente;  $x_1, x_2, \dots, x_k$  são as variáveis independentes;  $\beta_0, \beta_1, \beta_2, \dots, \beta_k$  são os parâmetros da regressão;  $\varepsilon$  é o termo que representa o resíduo ou erro da regressão, ou seja, é um termo de erro aleatório.

O termo  $\beta_0$  é denominado intercepto, ou coeficiente linear, e representa o valor da intersecção com o eixo  $y$ . Os termos  $\beta_1, \beta_2, \dots, \beta_k$  são chamados de coeficientes angulares. O termo de erro é responsável pela variabilidade em  $y$  que não pode ser explicada pelo efeito linear das variáveis independentes ( $k$  variáveis).

No ajuste do modelo de regressão múltipla, é conveniente expressar as operações matemáticas por meio de notação matricial. Utilizando a notação matricial o modelo fica

$$Y = X\beta + \varepsilon$$

em que,

$$Y_{n \times 1} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad X_{n \times p} = \begin{bmatrix} 1 & x_{11} & \dots & x_{k1} \\ 1 & x_{12} & \dots & x_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & \dots & x_{kn} \end{bmatrix}, \quad \beta_{p \times 1} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}, \quad e \quad \varepsilon_{n \times 1} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}.$$

### PRESSUPOSIÇÕES SOBRE O MODELO DE REGRESSÃO LINEAR MÚLTIPLA

As pressuposições sobre o modelo de regressão múltipla são:

- i) O erro  $\varepsilon$  é uma variável aleatória com média, ou valor esperado, igual a zero, ou seja,  $E(\varepsilon) = 0$ .

Implicação: Para dados valores de  $x_1, x_2, \dots, x_k$ , o valor esperado, ou média, de  $y$  é dado por:

$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k.$$

- ii) A variância do erro é constante e igual a  $\sigma^2$  para todos os valores das variáveis independentes  $x_1, x_2, \dots, x_k$  (condição de homogeneidade de variância).

- iii) Os erros são independentes, isto é, o erro de um conjunto de valores em particular das variáveis independentes não está relacionado com o erro de qualquer outro conjunto de valores ( $Cov(\varepsilon_i, \varepsilon_j) = 0$  para todo  $i \neq j$ ).

- iv) Os resíduos ou erros têm distribuição normal com média zero e variância constante ( $I\sigma^2$ ).

### EQUAÇÃO DE REGRESSÃO MÚLTIPLA ESTIMADA

Se os valores de  $\beta_0, \beta_1, \beta_2, \dots, \beta_k$  fossem conhecidos, a equação de regressão poderia ser usada para calcular o valor médio de  $y$  em relação a valores de  $x_1, x_2, \dots, x_k$ . Mas, geralmente

esses parâmetros são desconhecidos, e devem ser estimados a partir de dados amostrais. A estimação dos parâmetros pode ser obtida por meio do Método de Mínimos Quadrados. Esse método consiste em minimizar a soma de quadrados dos resíduos. O estimador de Mínimos quadrados de  $\beta$  é:

$$\hat{\beta} = (X'X)^{-1} X'Y$$

em que

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}, \quad X = \begin{bmatrix} 1 & x_{11} & \cdots & x_{k1} \\ 1 & x_{12} & \cdots & x_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & \cdots & x_{kn} \end{bmatrix}_{n \times p}, \quad X'Y = \begin{bmatrix} \sum y_i \\ \sum x_{1i}y_i \\ \vdots \\ \sum x_{ki}y_i \end{bmatrix}_{p \times 1}$$

$$X'X = \begin{bmatrix} n & \sum x_{1i} & \cdots & \sum x_{ki} \\ \sum x_{1i} & \sum x_{1i}^2 & \cdots & \sum x_{1i}x_{ki} \\ \vdots & \vdots & \ddots & \vdots \\ \sum x_{ki} & \sum x_{ki}x_{1i} & \cdots & \sum x_{ki}^2 \end{bmatrix}_{p \times p} \quad \text{e} \quad \hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_k \end{bmatrix}_{p \times 1}.$$

Nota-se que  $\hat{\beta}$  é o vetor das estimativas dos parâmetros do modelo de regressão múltipla. Portanto, a equação de regressão múltipla estimada é:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \cdots + \hat{\beta}_k x_k.$$

## INTERPRETAÇÃO DOS COEFICIENTES DA REGRESSÃO MÚLTIPLA

Na regressão linear simples, interpreta-se  $\hat{\beta}_i$  como uma estimativa da alteração em  $y$  correspondente à alteração de uma unidade na variável independente. Na regressão linear múltipla, interpreta-se  $\hat{\beta}_i$  como uma estimativa da alteração em  $y$  correspondente à alteração de uma unidade na variável  $x_i$  quando as outras variáveis independentes se mantêm constantes.

**Exemplo:** Para ilustrar os conceitos, utilizaremos um estudo realizado em uma academia de ginástica, durante 15 meses, cujo propósito foi identificar quais fatores afetavam seus gastos gerais. Foram identificados dois fatores potenciais: gastos com energia expressos em kilowatts consumidos e gastos com pessoal expressos em horas de Mão-de-Obra Direta (MOD). O resultado encontrado na pesquisa está apresentado na Tabela 1.

**Tabela 1** – Resultados da pesquisa sobre os gastos da academia.

Gastos reais da Academia ( y)	Consumo Kilowatts ( x1)	Horas de MOD (x2)
350	6	10
400	8	14
470	12	16
550	10	26
620	15	24
380	7	12
290	6	13
490	9	21
580	11	20
610	13	24
560	12	23
420	14	12
450	11	19
510	12	19
380	9	11

Pretende-se com essas informações ajustar um modelo de regressão linear múltipla. Usando os dados fornecidos, obtemos:

$$X = \begin{bmatrix} 1 & 6 & 10 \\ 1 & 8 & 14 \\ 1 & 12 & 16 \\ 1 & 10 & 26 \\ 1 & 15 & 24 \\ 1 & 7 & 12 \\ 1 & 6 & 13 \\ 1 & 9 & 21 \\ 1 & 11 & 20 \\ 1 & 13 & 24 \\ 1 & 12 & 23 \\ 1 & 14 & 12 \\ 1 & 11 & 19 \\ 1 & 12 & 19 \\ 1 & 9 & 11 \end{bmatrix} \text{ e } Y = \begin{bmatrix} 350 \\ 400 \\ 470 \\ 550 \\ 620 \\ 380 \\ 290 \\ 490 \\ 580 \\ 610 \\ 560 \\ 420 \\ 450 \\ 510 \\ 380 \end{bmatrix}$$

A partir dessas matrizes conseguimos definir as seguintes matrizes:

$$X'X = \begin{bmatrix} 15 & 155 & 264 \\ 155 & 1711 & 2847 \\ 264 & 2847 & 5050 \end{bmatrix} \text{ e } X'Y = \begin{bmatrix} 7060 \\ 75950 \\ 131000 \end{bmatrix}.$$

Usando as regras de inversão de matrizes determinou-se a matriz  $(X'X)^{-1}$ :

$$(X'X)^{-1} = \begin{bmatrix} 1,19055 & -0,06928 & -0,02318 \\ -0,06928 & 0,01347 & -0,00397 \\ -0,02318 & -0,00397 & 0,00365 \end{bmatrix}. \text{ Logo, } \hat{\beta} = (X'X)^{-1} X'Y = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} 106,7247 \\ 13,5794 \\ 12,7058 \end{bmatrix}.$$

Portanto, o modelo de regressão linear múltipla estimado foi  $\hat{y} = 106,7247 + 13,5794x_1 + 12,7058x_2$ .

### Interpretação:

O modelo de regressão linear múltipla estimado foi  $\hat{y} = 106,7247 + 13,5794x_1 + 12,7058x_2$ , então,  $\hat{\beta}_1 \cong 13,58$  e  $\hat{\beta}_2 \cong 12,71$ . Assim,  $\hat{\beta}_1 \cong 13,58$ , é a estimativa do aumento esperado no gasto mensal da academia correspondente ao aumento de 1 (um) kilowatt no consumo de energia quando os gastos com pessoal expressos em horas de mão-de-obra direta é mantido constante.

Da mesma forma,  $\hat{\beta}_2 \cong 12,71$  é a estimativa do aumento esperado no gasto mensal da academia correspondente ao aumento de uma hora de mão-de-obra direta quando o consumo de energia (kilowatts) é mantido constante.

## TESTE DE SIGNIFICÂNCIA DA REGRESSÃO MÚLTIPLA ANÁLISE DE VARIÂNCIA – ANAVA

O teste de hipótese para a significância da regressão é um teste para determinar se existe uma relação linear entre a variável de resposta  $y$  e o conjunto de regressores  $x_1, x_2, \dots, x_k$ . As hipóteses apropriadas são:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_1: \beta_j \neq 0 \text{ para no mínimo um } j.$$

Se  $H_0$  for rejeitada, o teste nos dá suficientes evidências estatísticas para concluirmos que um ou mais dos parâmetros não são iguais a zero e que a relação global entre  $y$  e o conjunto de variáveis independentes  $x_1, x_2, \dots, x_k$  é significativa. Entretanto, se  $H_0$  não puder ser rejeitada, não teremos evidências suficientes para concluir que uma relação significativa está presente.

Portanto, utiliza-se o teste F para concluir se existe ou não uma relação global significativa entre  $y$  e o conjunto de variáveis independentes  $x_1, x_2, \dots, x_k$ . Uma análise de variância geral fornece os resultados do teste F de um modelo de regressão múltipla. Na Tabela 2 é apresentado o procedimento para a ANAVA e consequentemente do teste F, sendo que  $p$  é o número de parâmetros do modelo e  $n$  é o número de observações da variável resposta (variável dependente).

**Tabela 2** – Análise de Variância de um modelo de regressão linear múltipla com  $k$  variáveis independentes.

<i>FV</i>	<i>GL</i>	<i>SQ</i>	<i>QM</i>	<i>F</i>
<b>Regressão</b>	$p - 1$	$SQ_{Reg}$	$SQ_{Reg}/(p - 1)$	$QM_{Reg}/QM_{Erro}$
<b>Erro</b>	$n - p$	$SQ_{Erro}$	$SQ_{Erro}/(n - p)$	
<b>Total</b>	$n - 1$	$SQ_{Total}$		

Também neste caso, existe uma relação entre a  $SQ_{Total}$ ,  $SQ_{Reg}$  e  $SQ_{Erro}$ :

$$SQ_{Total} = SQ_{Reg} + SQ_{Erro}$$

em que,  $SQ_{Total} = Y'Y - \frac{(\sum y_i)^2}{n}$ ;  $SQ_{Reg} = \hat{\beta}'X'Y - \frac{(\sum y_i)^2}{n}$  e  $SQ_{Erro} = Y'Y - \hat{\beta}'X'Y$ .

A regra de decisão é:

Critério do valor p: Rejeita-se  $H_0$  se o valor  $p \leq \alpha$ , em que  $\alpha$  é o nível de significância.

Critério do valor crítico: Rejeita-se  $H_0$  se  $F \geq F_\alpha$  em que  $F_\alpha$  se baseia em uma distribuição F com  $p - 1$  graus de liberdade no numerador e  $n - p$  graus de liberdade no denominador.

**Exemplo:** O procedimento de construção da ANAVA é ilustrado para o exemplo da academia. O modelo de regressão linear múltipla estimado foi  $\hat{y} = 106,7247 + 13,5794x_1 + 12,7058x_2$ , ou seja, o modelo possui 3 parâmetros ( $p = 3$ ). Como a coleta dos dados foi realizada durante 15 meses segue-se que  $n$  é igual a 15 (quinze) observações.

Soma de Quadrados Total:

$$SQ_{Total} = Y'Y - \frac{(\sum y_i)^2}{n} = 3.464.400 - \frac{(7060)^2}{15} = 141.493,3333.$$

Soma de Quadrados de Regressão:

$$SQ_{Reg} = \hat{\beta}'X'Y - \frac{(\sum y_i)^2}{n} = 3.449.287,3343 - \frac{(7060)^2}{15} = 126.380,6676.$$

Soma de Quadrados do Erro:

$$SQ_{Erro} = Y'Y - \hat{\beta}'X'Y = 3.464.400 - 3.449.287,3343 = 15.112,6657.$$

**Tabela 3** – Análise de Variância do modelo de regressão linear múltipla para o exemplo de Gastos da Academia.

<i>FV</i>	<i>GL</i>	<i>SQ</i>	<i>QM</i>	<i>F</i>
<b>Regressão</b>	3-1=2	126.380,6676	63.190,3338	50,1754
<b>Resíduo</b>	15-3=12	15.112,6657	1.259,3888	
<b>Total</b>	15-1=14	141.493,3333		

Como  $50,1754 = F \geq F_{0,05}(2, 12) = 3,89$  então rejeita-se  $H_0$  e concluímos que um ou mais dos parâmetros não são iguais a zero e que a relação global entre gastos reais da academia e o conjunto de variáveis independentes é significativa.

## TESTE DE HIPÓTESE PARA A SIGNIFICÂNCIA DA REGRESSÃO MÚLTIPLA

Se no teste  $F$  exibir uma significância global, o teste  $t$  é usado para determinar se cada uma das variáveis independentes individuais é significativa, ou seja, um teste  $t$  separado é realizado para cada uma das variáveis independentes do modelo; referimo-nos a cada um desses testes  $t$  como teste de significância individual.

$$\text{Hipóteses: } \begin{cases} H_0 : \beta_j = 0 \\ H_1 : \beta_j \neq 0 \end{cases} ; \quad \text{Estatística do teste: } t_c = \frac{\hat{\beta}_j}{\sqrt{S^2 C_{jj}}}$$

com  $\nu = n - p$  graus de liberdade e  $S_{\hat{\beta}_j} = \sqrt{S^2 C_{jj}} = \sqrt{\hat{V}(\hat{\beta}_j)}$  é a estimativa do erro padrão do parâmetro  $\beta_j$ .



**Regra de Decisão:** Rejeita-se  $H_0$  se  $t \leq -t_{(\alpha/2; n-p)}$  ou  $t \geq t_{(\alpha/2; n-p)}$ .

Tais testes são úteis na determinação do valor potencial de cada um dos regressores no modelo de regressão. Por exemplo, o modelo pode ser mais eficiente com a inclusão de variáveis adicionais ou talvez com a retirada de um ou mais regressores do modelo. A adição de uma variável ao modelo sempre aumenta a  $SQ_{Regressão}$  e sempre diminui a  $SQ_{Erro}$ . Temos que decidir se o aumento na soma quadrática da regressão é grande o suficiente para justificar o uso de uma variável adicional no modelo. Além disso, a adição de uma variável não importante no modelo pode, na verdade, aumentar a soma quadrática do erro, indicando que a adição de tal variável fez realmente o modelo apresentar um ajuste mais pobre aos dados.

A estimativa da matriz de variâncias e covariâncias para os parâmetros do modelo de regressão múltipla, usada para o teste t, é calculada por:

$$\hat{V}(\hat{\beta}) = (X'X)^{-1} \cdot S^2 = (X'X)^{-1} \cdot QMErro = \begin{bmatrix} \hat{V}(\hat{\beta}_0) & \hat{Cov}(\hat{\beta}_1, \hat{\beta}_0) & \cdots & \hat{Cov}(\hat{\beta}_j, \hat{\beta}_0) \\ \hat{Cov}(\hat{\beta}_0, \hat{\beta}_1) & \hat{V}(\hat{\beta}_1) & \cdots & \hat{Cov}(\hat{\beta}_j, \hat{\beta}_1) \\ \vdots & \vdots & \ddots & \vdots \\ \hat{Cov}(\hat{\beta}_0, \hat{\beta}_j) & \hat{Cov}(\hat{\beta}_1, \hat{\beta}_j) & \cdots & \hat{V}(\hat{\beta}_j) \end{bmatrix}$$

com,  $j=1,2,\dots,p$  e  $p$  é igual ao número de parâmetros.

**Exemplo:** Considerando os dados do exemplo de gastos da academia, um teste t será realizado para determinar a significância de cada um dos parâmetros individuais.

A matriz de variâncias e covariâncias é:

$$(X'X)^{-1} QMErro = \begin{bmatrix} 1,19055 & -0,06928 & -0,02318 \\ -0,06928 & 0,01347 & -0,00397 \\ -0,02318 & -0,00397 & 0,00365 \end{bmatrix} 1.259,3888$$

As estimativas de  $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, S_{\hat{\beta}_0}, S_{\hat{\beta}_1}$  e  $S_{\hat{\beta}_2}$  são:

$$\begin{aligned} \hat{\beta}_0 &= 106,7247 & S_{\hat{\beta}_0} &= \sqrt{1,19055 \times 1.259,3888} = 38,7216 \\ \hat{\beta}_1 &= 13,5794 & S_{\hat{\beta}_1} &= \sqrt{0,01347 \times 1.259,3888} = 4,1187 \\ \hat{\beta}_2 &= 12,7058 & S_{\hat{\beta}_2} &= \sqrt{0,00365 \times 1.259,3888} = 2,1440 \end{aligned}$$

As estatísticas de testes de hipóteses para cada um dos parâmetros  $\beta_0, \beta_1$  e  $\beta_2$  são:

	<i>Erro padrão</i>	<i>Estatística t</i>
$\beta_0$	38,7216	2,7562
$\beta_1$	4,1187	3,2970
$\beta_2$	2,1440	5,9262

A 5% de significância, podemos rejeitar  $H_0 : \beta_0 = 0, H_0 : \beta_1 = 0$  e  $H_0 : \beta_2 = 0$ , pois as estatísticas de teste são maiores que o valor tabelado  $t_{(\alpha/2; n-p)} = t_{(0,025; 12)} = 2,179$  em todos os casos, logo concluímos que os gastos da academia estão relacionados com consumo de kilowatts ( $x_1$ ) e horas de MOD ( $x_2$ ).

## COEFICIENTE DE DETERMINAÇÃO MÚLTIPLO

O coeficiente de determinação múltiplo mede a eficiência de ajuste da equação de regressão múltipla estimada em relação ao conjunto de dados. O coeficiente de determinação múltiplo pode ser interpretado como a proporção da variabilidade da variável dependente que pode ser explicada pela equação de regressão múltipla estimada.

O coeficiente de determinação múltiplo é:

$$R^2 = \frac{SQRe g}{SQTotal}.$$

## COEFICIENTE DE DETERMINAÇÃO MÚLTIPLO AJUSTADO

O coeficiente de determinação múltiplo ajustado também mede a eficiência de ajuste da equação de regressão múltipla estimada em relação ao conjunto de dados, porém, ponderando o número de variações ( $n$ ) e o número de variáveis ( $k$ ) para evitar a superestimação do  $R^2$ .

O coeficiente de determinação múltiplo ajustado é:

$$R_{ajust}^2 = 1 - (1 - R^2) \frac{(n-1)}{(n-k-1)} = 1 - \frac{QMErro}{QMTotal}.$$

**Exemplo:** Vamos calcular o  $R_{ajust}^2$  para discutir a eficiência do ajuste do modelo estimado para o exemplo da academia.

O coeficiente de determinação múltiplo é:  $R^2 = \frac{SQRe g}{SQTotal} = \frac{126.380,6676}{141.493,3333} = 0,8932$ .

O coeficiente de determinação múltiplo ajustado é:  $R_{ajust}^2 = 1 - (1 - 0,8932) \frac{(15-1)}{(15-2-1)} = 0,8754$ .

Interpretação: calcula-se que 87,54% da variabilidade nos gastos da academia ( $y$ ) pode ser explicada pela presença dos dois regressores, gastos com energia ( $x_1$ ) e gastos com pessoal ( $x_2$ ).

## INTERVALO DE CONFIANÇA PARA O COEFICIENTE DE REGRESSÃO

A partir da matriz de variâncias e covariâncias dos parâmetros do modelo de regressão múltipla pode-se definir o intervalo de confiança para  $\beta_j$  como:

$$IC(\beta_j): \hat{\beta}_j \pm t_{(\alpha/2, n-p)} \cdot \sqrt{S^2 C_{jj}}$$

100(1- $\alpha$ )%

em que,  $\sqrt{S^2 C_{jj}} = \sqrt{\hat{V}(\hat{\beta}_j)} = S_{\hat{\beta}_j}$  é o erro padrão do parâmetro  $\beta_j$ .

**Exemplo:** Considerando os dados do exemplo de gastos da academia, construiremos um intervalo de confiança de 95% para o parâmetro  $\beta_1$ .

A estimativa de  $\beta_1$  é  $\hat{\beta}_1 = 13,5794$ . A estimativa de seu erro padrão é  $S_{\hat{\beta}_1} = \sqrt{S^2 C_{11}} = \sqrt{\hat{V}(\hat{\beta}_1)} = \sqrt{1.259,3888 \times 0,01347} = 4,1187$ . Consequentemente, o intervalo de confiança de 95% para  $\beta_1$  é calculado a partir da expressão (39).

$$\begin{aligned} IC_{95\%}(\beta_1): \hat{\beta}_1 \pm t_{(0,05/2; 12)} \cdot \sqrt{S^2 C_{11}} \\ IC_{95\%}(\beta_1): 13,5794 \pm 2,179 \times \sqrt{1.259,3888 \times 0,01347} \\ IC_{95\%}(\beta_1): 13,5794 \pm 2,179 \times 4,1187 \\ IC_{95\%}(\beta_1): 13,5794 \pm 8,9746 \\ IC_{95\%}(\beta_1): [4,6048; 22,5540] \end{aligned}$$

Pode-se afirmar com 95% de confiança que o verdadeiro valor do parâmetro  $\beta_1$  está contido no intervalo [4,6048; 22,5540].

## INTERVALO DE CONFIANÇA PARA A RESPOSTA MÉDIA

Podemos obter também um intervalo de confiança para a resposta média em um determinado ponto,  $x_{01}, x_{02}, \dots, x_{0k}$ . Para estimar a resposta média nesse ponto, defina o vetor  $X_0$  como:

$$X_0 = \begin{bmatrix} 1 \\ x_{01} \\ x_{02} \\ \vdots \\ x_{0k} \end{bmatrix}$$

A resposta média nesse ponto é estimada por  $\hat{\mu}_{Y|X_0} = X_0' \hat{\beta}$ . Portanto, o intervalo de confiança de  $100(1-\alpha)\%$  para a resposta no ponto  $x_{01}, x_{02}, \dots, x_{0k}$  é:

$$IC_{100(1-\alpha)\%}(\mu_{Y|X_0}): X_0' \hat{\beta} \pm t_{(\alpha/2, n-p)} \cdot \sqrt{S^2 X_0' (X' X)^{-1} X_0}$$

**Exemplo:** Considerando os dados do exemplo de gastos da academia, determinaremos um intervalo de confiança de 95% para o gasto médio da academia quando se consome 10 kilowatts ( $x_1 = 10$ ) e 20 horas de mão-de-obra direta ( $x_2 = 20$ ).

Usando a equação de regressão estimada  $\hat{y} = 106,7247 + 13,5794x_1 + 12,7058x_2$ , sendo  $x_1 = 10$  e  $x_2 = 20$ , obtemos o seguinte valor de  $\hat{y}$ :

$$\hat{y} = 106,7247 + 13,5794 \times 10 + 12,7058 \times 20 = R\$496,64$$

**Nota:** Observe que

$$\hat{y} = X_0' \hat{\beta} = \begin{bmatrix} 1 & 10 & 20 \end{bmatrix} \begin{bmatrix} 106,7247 \\ 13,5794 \\ 12,7058 \end{bmatrix} = 106,7247 + 13,5794 \times 10 + 12,7058 \times 20 = R\$496,64.$$

$$\text{em que } X_0 = \begin{bmatrix} 1 \\ 10 \\ 20 \end{bmatrix} \text{ e } \hat{\beta} = \begin{bmatrix} 106,7247 \\ 13,5794 \\ 12,7058 \end{bmatrix}.$$

A estimação por ponto do gasto médio da academia quando se consome 10 kilowatts e 20 horas de mão-de-obra direta é de R\$ 496,64. A estimação por intervalo para o gasto médio da academia foi obtida a partir da expressão (41):

$$\begin{aligned} IC_{95\%}(\mu_{Y|X_0}): X_0' \hat{\beta} \pm t_{(0,025; 12)} \sqrt{S^2 X_0' (X'X)^{-1} X_0} \\ IC_{95\%}(\mu_{Y|X_0}): 496,64 \pm 2,179 \sqrt{1.259,3888 \times 0,09675} \\ IC_{95\%}(\mu_{Y|X_0}): 496,64 \pm 24,05 \\ IC_{95\%}(\mu_{Y|X_0}): [472,59; 520,69] \end{aligned}$$

em que,

$$X_0' (X'X)^{-1} X_0 = \begin{bmatrix} 1 & 10 & 20 \end{bmatrix} \begin{bmatrix} 1,19055 & -0,06928 & -0,02318 \\ -0,06928 & 0,01347 & -0,00397 \\ -0,02318 & -0,00397 & 0,00365 \end{bmatrix} \begin{bmatrix} 1 \\ 10 \\ 20 \end{bmatrix} = 0,09675$$

Pode-se afirmar com 95% de confiança que o gasto médio da academia quando se consome 10 kilowatts e 20 horas de mão-de-obra direta por mês está contido no intervalo [472,59; 520,69].

### 3. REGRESSÃO LOGÍSTICA

A regressão linear geralmente é adequada quando a variável resposta é quantitativa. No entanto, se a variável resposta é binária (0 ou 1, sucesso ou fracasso) existem apenas dois possíveis valores para o erro (resíduo do modelo) caso se assuma o modelo linear simples. Nesse caso, os erros não podem ser considerados com distribuição normal, e, além disso, a variância dos erros não é constante, pois essa é função da média. Assim, a forma da função de resposta ( $E(Y_i)$ ) deve ser não linear.

Uma função de resposta monotonamente crescente (ou decrescente) em forma de S, chamada de função logit, geralmente é empregada para modelar a variável resposta. O modelo de regressão é então chamado de regressão logística, e a função logit, considerando uma variável resposta  $y$  e apenas uma covariável, tem a forma a seguir.

$$E(Y) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)} = \frac{1}{1 + \exp(-(\beta_0 + \beta_1 x))}$$

A função logit tem a forma dos gráficos a seguir.

O termo  $\exp(\beta_0 + \beta_1 x)$  é chamado de razão de chances (odds ratio). A razão de chances pode ser interpretada como o número de vezes que o sucesso é mais provável do que o fracasso, considerando um valor particular de  $x$ . Assim, se a razão de chances é igual a 2 por exemplo, significa que o sucesso é 2 vezes mais provável que o fracasso, para aquele valor de  $x$ . A razão de chances varia de  $e^{\beta_1}$  quando  $x$  aumenta em uma unidade.

O log da razão de chances,  $\beta_0 + \beta_1 x$ , é uma função linear da variável regressora. Esse termo é chamado de preditor linear, e a inclinação  $\beta_1$  é a variação no log das chances que resulta do aumento de 1 unidade em  $x$ .

A estimação dos parâmetros é feita pelo método da máxima verossimilhança. No entanto, as funções são não lineares, e assim, são necessários métodos numéricos para se encontrar as estimativas de máxima verossimilhança.

Exemplo. A regressão logística é ilustrada usando os dados de temperatura de lançamento e falha de O-ring para 24 lançamentos de ônibus espaciais antes do desastre do Challenger em janeiro de 1986. O número 1 significa que houve pelo menos uma falha e o zero significa que não houve falha.

Dados:

Temperatura	53	56	57	63	66	67	67	67	68	69	70	70
Falha	1	1	1	0	0	0	0	0	0	0	0	1
Temperatura	70	70	72	73	75	75	76	76	78	79	80	81
Falha	1	1	0	0	0	1	0	0	0	0	0	0

O modelo ajustado é  $\hat{y} = \frac{1}{1 + \exp(-(10,875 - 0,17132x))}$ . O termo  $e^{\beta_1}$  é 0,84; logo, o aumento de um grau na temperatura reduz as chances de falha por 0,84.

No gráfico de probabilidade, fica evidente o aumento na probabilidade de falha, para pequenas temperaturas. A temperatura de lançamento do Challenger foi de 31°, e este valor está fora do intervalo considerado, não sendo aconselhável fazer previsões para esta temperatura. No entanto, analisando-se o gráfico, pode-se verificar forte indício de que uma temperatura tão baixa resultaria quase certamente em uma falha. Esse exemplo indica dramaticamente que todos os cientistas devem ter conhecimentos estatísticos.