Modelo de Regressão Linear Simples

O modelo de regressão linear simples considera apenas uma variável explicativa e a função de regressão é linear. O modelo é definido por:

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \tag{1}$$

em que:

- Y_i é o valor da variável resposta para a i-ésima observação,
- β_0 é o intercepto e β_1 é o coeficiente angular, ambos são parâmetros desconhecidos,
- X_i é uma constante conhecida, o valor da variável explicativa para a i-ésima observação. X_i é uma variável fixa (ou sem erro ou determinística),
- ε_i são independentes e $N(0, \sigma^2)$.



ALGUNS EXEMPLOS

- A) Um estudo foi realizado para verificar a influência das variáveis: X_1 capital investido e X_2 gasto em publicidade no lucro anual, Y, de empresas.
- B) Um experimento foi realizado para estudar a influência das variáveis: X_1 quantidade de adubo, X_2 quantidade de chuva e X_3 tipo de solo, na produção de cana de açúcar, Y.
- C) Com o objetivo de verificar o perfil dos vestibulando, uma Universidade estudou quais fatores, como, X_1 nota na redação, X_2 notas médias nas provas de Matemática, Química e Física, X_3 curso pretendido, X_4 sexo e X_5 relação candidato/vaga influênciam na nota, Y, obtida na prova do Vestibular.

Modelo de regressão linear com 2 variáveis explicativas

Quando há duas variáveis explicativas, X_1 e X_2 , o modelo de regressão é definido por:

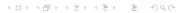
$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i,$$

é conhecido como modelo de primeira ordem, pois é linear nos parâmetros e linear nas variáveis explicativas.

- Os parâmetros desconhecidos do modelo são β_0 , β_1 e β_2 .
- ε_i é o erro do modelo com $N(0, \sigma^2)$.
- $E(Y) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2}$ é conhecida como superfície de regressão ou superfície de resposta.

Modelo de regressão linear com 2 variáveis explicativas

- Quando o escopo do modelo inclui $X_1=0$ e $X_2=0$, então β_0 representa a média da variável resposta, E(Y). Caso contrário, β_0 não tem interpretação no modelo.
- O parâmetro β_1 indica a mudança no valor esperado de Y, E(Y), com o aumento em uma unidade em X_1 , mantendo X_2 constante.
- O parâmetro β_2 indica a mudança no valor esperado de Y, E(Y), com o aumento em uma unidade em X_2 , mantendo X_1 constante.
- Quando o efeito de X_1 na resposta média não depende de X_2 e quando o efeito de X_2 na resposta média não depende de X_1 , as variáveis explicativas são consideradas ter efeito aditivo.
- Dessa forma, os coeficientes β_1 e β_2 são também conhecidos como coeficientes de regressão parcial.



Em geral, a variável resposta Y pode estar relacionada com p variáveis explicativas. O modelo de regressão é definido por:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \ldots + \beta_p X_{ip} + \varepsilon_i, \qquad (2)$$

é conhecido como modelo de primeira ordem, pois é linear nos parâmetros e linear nas variáveis explicativas. Esse modelo é conhecido como um modelo de regressão linear múltipla ou também como um modelo de regressão linear geral.

- Y_i é o valor da variável resposta para a i-ésima observação,
- β_0 , β_1 , ..., β_p são parâmetros desconhecidos,
- $X_{i1}, X_{i2}, \dots, X_{ip}$ são constantes conhecidas,
- ε_i são independentes e $N(0, \sigma^2)$, $i = 1, \ldots, n$.



Assumindo que $E(\varepsilon_i)=0$, a função da resposta para o modelo (2), é:

$$E(Y_i) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \ldots + \beta_p X_{ip}.$$
 (3)

- Essa função resposta é um hiperplano, que é um plano em mais do que duas dimensões.
- O parâmetro β_k indica a mudança no valor esperado de Y, E(Y), com o aumento em uma unidade em X_k, mantendo todas as outras variáveis constantes.
- O efeito das variáveis explicativas no modelo são considerados aditivos.
- O modelo (2) com erro normal implica que as observações Y_i são variáveis normais independentes com $E(Y_i)$ dado em (3) e variância constante, $Var(Y_i) = \sigma^2$.

O modelo de regressão linear geral engloba uma variedade de situações. Algumas serão consideradas a seguir:

VARIÁVEL EXPLICATIVA QUALITATIVA

- O modelo (2) considera não apenas variáveis explicativas quantitativas, mas também variáveis qualitativas como, por exemplo: sexo, classe social, entre outras.
- Para a inclusão das variáveis explicativas qualitativas no modelo é necessário usar variável indicadora para identificar a classe/categoria da variável qualitativa.
- Considere uma análise de regressão para prever a duração da internação em um hospital (Y) baseada na idade (X_1) e sexo (X_2) do paciente.
- A variável sexo é definida por:

$$X_2 = \left\{ egin{array}{ll} 1 & ext{se sexo feminino} \\ 0 & ext{se sexo masculino} \end{array}
ight.$$

Variável explicativa Qualitativa

O modelo de regressão linear é definido por:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i,$$

em que:

- Y_i é a duração da internação em um hospital,
- X_{i1} é a idade e X_{i2} é o sexo do *i*-ésimo paciente.
- A função de regressão é: $E(Y_i) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2}$.

Variável explicativa Qualitativa

$$E(Y_i) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2}$$

- Para pacientes do sexo masculino, $X_2=0$, a função de regressão é: $E(Y_i)=\beta_0+\beta_1X_{i1}$
- Para pacientes do sexo feminino, $X_2 = 1$, a função de regressão é: $E(Y_i) = (\beta_0 + \beta_2) + \beta_1 X_{i1}$
- Essas duas funções de regressão representam linhas retas paralelas com interceptos diferentes.

Variável explicativa Qualitativa

Um modelo de regressão linear para a variável resposta Y com idade, sexo e estado do paciente como variáveis explicativas é definido por:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \beta_4 X_{i4} + \varepsilon_i$$

em que:

- Y_i é a duração da internação em um hospital,
- X_{i1} é a idade
- $X_{i2} = \left\{ egin{array}{ll} 1 & ext{se sexo feminino} \\ 0 & ext{se sexo masculino} \end{array}
 ight.$
- a variável estado do paciente: bom, estável ou grave, é representada pelas variáveis indicadoras:
- $X_{i3} = \begin{cases} 1 \text{ se grave} \\ 0 \text{ caso contrário} \end{cases}$ $X_{i4} = \begin{cases} 1 \text{ se estável} \\ 0 \text{ caso contrário} \end{cases}$

VARIÁVEL EXPLICATIVA QUALITATIVA

 a variável estado do paciente: bom, estável ou grave, é representada pelas variáveis indicadoras:

Estado do paciente	X_{i3}	X_{i4}
grave	1	0
estável	0	1
bom	0	0

ullet Em geral, variáveis qualitativas com c categorias/classes são representadas por c-1 variáveis indicadoras.

Modelo de Regressão Polinomial

Modelos de regressão polinomial são um caso particular do modelo de regressão linear geral. Eles consideram termo quadrático ou de maior ordem da variável explicativa. Um modelo de regressão polinomial com uma variável explicativa é definido por:

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \varepsilon_i,$$

se definirmos:

- $X_{i1} = X_i$ e $X_{i2} = X_i^2$
- O modelo pode ser escrito por:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i.$$



MODELO DE REGRESSÃO COM VARIÁVEIS TRANSFORMADAS

Considere o modelo com a seguinte variável transformada:

$$\log Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \varepsilon_i.$$

• Se considerarmos $Y'_i = \log Y_i$, o modelo pode ser escrito por:

$$Y_i' = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \varepsilon_i,$$

MODELO DE REGRESSÃO COM VARIÁVEIS TRANSFORMADAS

Agora, considere o modelo com a seguinte variável transformada:

$$Y_i = \frac{1}{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i}.$$

• Se considerarmos $Y'_i = 1/Y_i$, o modelo pode ser escrito por:

$$Y_i' = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i,$$

Modelo de Regressão com Interação

Quando o efeito da variável explicativa na variável resposta não é aditivo e o efeito de uma variável explicativa depende dos níveis de outras variáveis explicativas, tem-se modelos de regressão não aditivos:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i1} X_{i2} + \varepsilon_i.$$

• Se considerarmos $X_{i3} = X_{i1}X_{i2}$, o modelo pode ser escrito por:

$$Y_{i} = \beta_{0} + \beta_{1}X_{i1} + \beta_{2}X_{i2} + \beta_{3}X_{i3} + \varepsilon_{i}.$$



Combinação de casos

- Um modelo de regressão pode combinar alguns elementos já citados anteriormente e ainda pode ser tratado como um modelo de regressão linear geral.
- Seja um modelo de segunda ordem e com interação:

$$Y_{i} = \beta_{0} + \beta_{1}X_{i1} + \beta_{2}X_{i1}^{2} + \beta_{3}X_{i2} + \beta_{4}X_{i2}^{2} + \beta_{5}X_{i1}X_{i2} + \varepsilon_{i}.$$
 (4)

- Seja: $Z_{i1} = X_{i1}$, $Z_{i2} = X_{i1}^2$, $Z_{i3} = X_{i2}$, $Z_{i4} = X_{i2}^2$ e $Z_{i5} = X_{i1}X_{i2}$.
- O modelo de regressão (4) pode ser expresso por:

$$Y_i = \beta_0 + \beta_1 Z_{i1} + \beta_2 Z_{i2} + \beta_3 Z_{i3} + \beta_4 Z_{i4} + \beta_5 Z_{i5} + \varepsilon_i$$



Observação: Por meio dos modelos apresentados, verifica-se que o *Modelo de Regressão Linear Geral* não é restrito à superfície de resposta linear. O termo *modelo linear* se refere ao fato que o modelo (2) é linear nos parâmetros.

MODELO DE REGRESSÃO LINEAR GERAL EM TERMOS MATRICIAIS

• O modelo de regressão linear geral:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \ldots + \beta_p X_{ip} + \varepsilon_i$$
 $i = 1, \ldots, n$ pode ser escrito por:

$$Y_{1} = \beta_{0} + \beta_{1}X_{11} + \beta_{2}X_{12} + \dots + \beta_{p}X_{1p} + \varepsilon_{1}$$

$$Y_{2} = \beta_{0} + \beta_{1}X_{21} + \beta_{2}X_{22} + \dots + \beta_{p}X_{2p} + \varepsilon_{2}$$

$$\vdots$$

$$Y_{n} = \beta_{0} + \beta_{1}X_{n1} + \beta_{2}X_{n2} + \dots + \beta_{p}X_{np} + \varepsilon_{n}$$
(5)

MODELO DE REGRESSÃO LINEAR GERAL EM TERMOS MATRICIAIS

 E por ser definido pelos vetores das observações, matriz de variáveis explicativas, vetor dos parâmetros e vetor de erros:

$$\mathbf{Y}_{n\times 1} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \quad \mathbf{X}_{n\times(p+1)} = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1p} \\ 1 & X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix}$$
$$\boldsymbol{\beta}_{(p+1)\times 1} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} \quad \boldsymbol{\varepsilon}_{n\times 1} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

 Dessa forma, o modelo (2) pode ser escrito em termos matriciais:

$$\mathbf{Y}_{n \times 1} = \mathbf{X}_{n \times (p+1)} \boldsymbol{\beta}_{(p+1) \times 1} + \boldsymbol{\varepsilon}_{n \times 1}$$

• As condições do erro do modelo (2), $E(\varepsilon_i) = 0$ e $Var(\varepsilon_i) = \sigma^2$, podem ser escritas em termos matriciais por:

$$E(\varepsilon)_{n\times 1}=\mathbf{0}_{n\times 1} \tag{6}$$

que implica em:

$$\begin{bmatrix} E(\varepsilon_1) \\ E(\varepsilon_2) \\ \vdots \\ E(\varepsilon_n) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

• A condição que os erros têm variância constante σ^2 e que todas as covariâncias, $Cov(\varepsilon_i, \varepsilon_j) = 0$ para todo $i \neq j$ (e que os erros são normais e independentes) é expressa em termos matriciais através da matriz de covariâncias:

$$Var(\varepsilon)_{n\times n} = \begin{bmatrix} \sigma^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma^2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \dots & \sigma^2 \end{bmatrix} = \sigma^2 \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

$$Var(\varepsilon)_{n\times n} = \sigma^2 \mathbf{I}_{n\times n} \tag{7}$$

 Assim, o modelo de regressão (2), em termos matriciais é definido por:

$$\mathbf{Y}_{n\times 1} = \mathbf{X}_{n\times(p+1)}\boldsymbol{\beta}_{(p+1)\times 1} + \boldsymbol{\varepsilon}_{n\times 1} \tag{8}$$

em que,

- ε é um vetor de variáveis aleatórias independentes e identicamente distribuídas com distribuição Normal e $E(\varepsilon) = \mathbf{0}$ e $Var(\varepsilon) = \sigma^2 \mathbf{I}$.
- Consequentemente,

$$E(\mathbf{Y}) = E(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{X}\boldsymbol{\beta} + E(\boldsymbol{\varepsilon}) = \mathbf{X}\boldsymbol{\beta}$$
$$Var(\mathbf{Y}) = Var(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{0} + Var(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$$

 O método de mínimos quadrados para estimar os parâmetros do modelo de regressão linear geral, cujo objetivo é minimizar a soma do quadrado dos erros:

$$SQ(\beta_0, \beta_1, \dots, \beta_p) = \sum_{i=1}^n \varepsilon_i^2$$

$$= \sum_{i=1}^n (Y_i - (\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}))^2.$$

Em notação matricial, tem-se que:

$$SQ(\beta) = \varepsilon' \varepsilon = (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta)$$
 (9)



• Ao expandir (9), tem-se que:

$$SQ(\beta) = (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) = \mathbf{Y}'\mathbf{Y} - \mathbf{Y}'\mathbf{X}\beta - \beta'\mathbf{X}'\mathbf{Y} + \beta'\mathbf{X}'\mathbf{X}\beta$$

• $\mathbf{Y}'\mathbf{X}\boldsymbol{\beta}$ é um escalar (1x1), que é igual a sua transposta: $(\mathbf{Y}'\mathbf{X}\boldsymbol{\beta})' = \boldsymbol{\beta}'\mathbf{X}'\mathbf{Y}$. Dessa forma, escrevemos:

$$SQ(\beta) = \mathbf{Y}'\mathbf{Y} - 2\mathbf{Y}'\mathbf{X}\beta + \beta'\mathbf{X}'\mathbf{X}\beta$$
 (10)

• Para encontrar o vetor β que minimiza (10), faz-se: $\frac{\partial SQ(\beta)}{\partial \beta}$

$$\frac{\partial SQ(\beta)}{\partial \beta} = -2\mathbf{X}'\mathbf{Y} + 2\mathbf{X}'\mathbf{X}\beta \tag{11}$$

• Igualando (11) a zero e substituindo β por $\hat{\beta}$, encontra-se as equações normais em forma matricial:

$$\begin{aligned} -2\mathbf{X}'\mathbf{Y} + 2\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} &= 0 \\ -\mathbf{X}'\mathbf{Y} + \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} &= 0 \\ \text{equações normais} &\Rightarrow \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} &= \mathbf{X}'\mathbf{Y} \end{aligned}$$

 Para obter os coeficientes estimados através das equações normais por métodos matriciais, multiplica-se ambos os lados pela inversa de X'X:

$$(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\hat{\boldsymbol{eta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

• Sabendo que $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X} = \mathbf{I}$, matriz identidade, tem-se que:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \tag{12}$$



- A matriz $(\mathbf{X}'\mathbf{X})^{-1}$ sempre existirá se as variáveis regressoras forem **linearmente independentes**.
- As variáveis regressoras são linearmente independente, se nenhuma coluna da matriz X for uma combinação linear das outras colunas.

Observação:

- Assume-se que as variáveis explicativas são fixas (não aleatórias), medidas sem erro.
- Entretanto, quando os dados são de estudos observacionais, algumas variáveis podem ser aleatórias. Porém, todos os resultados encontrados quando assume-se que as variáveis explicativas são fixas, são válidos para o caso em que essas variáveis são aleatórias.
- Quando as variáveis X_1, X_2, \dots, X_p são aleatórias é necessário que as observações de cada variável sejam independentes e que a distribuição não dependa dos coeficientes β ou de σ^2 .
- Para construir IC e testes de hipóteses, tem que assumir que a distribuição condicional de Y dado X_1, X_2, \ldots, X_p é normal com média $\beta_0 + \beta_1 X_1 + \ldots + \beta_p X_p$ e variância σ^2 .



Valores ajustados

• Seja o vetor de valores ajustados \hat{Y}_i denotado por $\hat{\mathbf{Y}}$:

$$\hat{\mathbf{Y}}_{n \times 1} = \begin{bmatrix} Y_1 \\ \hat{Y}_2 \\ \vdots \\ \hat{Y}_n \end{bmatrix}$$

• Em notação matricial tem-se que:

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} \tag{13}$$

porque:

$$\mathbf{Y}_{n\times 1} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1p} \\ 1 & X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

Valores ajustados

• A equação (13) pode ser escrita como:

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

ou, de forma equivalente:

$$\hat{\mathbf{Y}} = \mathbf{HY} \tag{14}$$

em que:

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^{'}\mathbf{X})^{-1}\mathbf{X}^{'}$$

é conhecida como a matriz chapéu. Essa matriz tem importante papel na análise de diagnóstico. A matriz **H** é simétrica e tem a propriedade de ser idempotência. Ou seja,

$$HH = H$$



Resíduo

• Seja o vetor de resíduos $e_i = Y_i - \hat{Y}_i$ denotado por **e**:

$$\mathbf{e}_{n \times 1} = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

• Em notação matricial tem-se que:

$$\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} \tag{15}$$

ou,

$$e = Y - \hat{Y} = Y - HY = (I - H)Y$$



Estimador de σ^2

• Como definido anteriormente, o estimador de σ^2 é dado por:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^{n} e_i^2}{n - (p+1)} = \frac{SQRes}{n - (p+1)} = MSRes$$

• Em termos matriciais $SQRes = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$ pode ser definido por:

$$SQRes = \mathbf{e}'\mathbf{e}$$

$$= (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

$$= (\mathbf{Y}' - \hat{\boldsymbol{\beta}}'\mathbf{X}')(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

$$= \mathbf{Y}'\mathbf{Y} - \mathbf{Y}'\mathbf{X}\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}}$$

$$= \mathbf{Y}'\mathbf{Y} - 2\hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}}$$

Estimador de σ^2

• Desde que $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$, tem-se que:

$$SQRes = \mathbf{Y}'\mathbf{Y} - 2\hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$
$$= \mathbf{Y}'\mathbf{Y} - 2\hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y}$$
$$= \mathbf{Y}'\mathbf{Y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y}$$

• Dessa forma, o estimador de σ^2 em termos matriciais é definido por:

$$\hat{\sigma}^2 = \frac{\mathbf{Y}'\mathbf{Y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y}}{n - (p+1)} = MSRes$$

Em geral, a variável resposta Y pode estar relacionada com p variáveis explicativas. O modelo de regressão é definido por:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \ldots + \beta_p X_{ip} + \varepsilon_i, \qquad (2)$$

é conhecido como modelo de primeira ordem, pois é linear nos parâmetros e linear nas variáveis explicativas. Esse modelo é conhecido como um modelo de regressão linear múltipla ou também como um modelo de regressão linear geral.

- Y_i é o valor da variável resposta para a i-ésima observação,
- β_0 , β_1 , ..., β_p são parâmetros desconhecidos,
- $X_{i1}, X_{i2}, \dots, X_{ip}$ são constantes conhecidas,
- ε_i são independentes e $N(0, \sigma^2)$, $i = 1, \ldots, n$.



 Dessa forma, o modelo (2) pode ser escrito em termos matriciais:

$$\mathbf{Y}_{n imes 1} = \mathbf{X}_{n imes (p+1)} oldsymbol{eta}_{(p+1) imes 1} + arepsilon_{n imes 1}$$

As condições do erro do modelo (2) são:

$$E(\varepsilon)_{n \times 1} = \mathbf{0}_{n \times 1}$$
 e $Var(\varepsilon)_{n \times n} = \sigma^2 \mathbf{I}_{n \times n}$

Consequentemente,

$$E(\mathbf{Y}) = E(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{X}\boldsymbol{\beta} + E(\boldsymbol{\varepsilon}) = \mathbf{X}\boldsymbol{\beta}$$
$$Var(\mathbf{Y}) = Var(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{0} + Var(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$$

ullet O estimador de mínimos quadrados de eta é dada por:

$$\hat{oldsymbol{eta}} = (\mathbf{X}^{'}\mathbf{X})^{-1}\mathbf{X}^{'}\mathbf{Y}$$





Estimador de σ^2

ullet Como definido anteriormente, o estimador de σ^2 é dado por:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n - (p+1)} = \frac{SQRes}{n - (p+1)} = MSRes$$

• Em termos matriciais:

$$\hat{\sigma}^2 = \frac{\mathbf{Y}'\mathbf{Y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y}}{n - (p+1)} = \frac{SQRes}{n - (p+1)} = MSRes$$

Soma de Quadrados

 Vamos definir a soma de quadrados total em termos matriciais. Lembrando que:

$$SQT = \sum_{i=1}^{n} (Y_i - \bar{Y})^2 = \sum_{i=1}^{n} Y_i^2 - \frac{(\sum_{i=1}^{n} Y_i)^2}{n}$$

• O termo $\frac{(\sum_{i=1}^{n} Y_i)^2}{n}$ pode ser escrito como:

$$\frac{\left(\sum_{i=1}^{n} Y_{i}\right)^{2}}{n} = \left(\frac{1}{n}\right) \mathbf{Y}' \mathbf{J} \mathbf{Y}$$

em que **J** é uma matriz de 1s. Isto é:

$$\mathbf{J} = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix}$$

Soma de Quadrados

• Dessa forma, a soma de quadrados total é definida por:

$$SQT = \mathbf{Y}'\mathbf{Y} - \left(\frac{1}{n}\right)\mathbf{Y}'\mathbf{JY}.$$

- O número de graus de liberdade associado a SQT é n-1.
- Dado as informações das SQT e SQRes, a soma de quadrados da regressão é definida por:

$$SQReg = SQT - SQRes = \hat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{Y} - \left(\frac{1}{n}\right) \mathbf{Y}' \mathbf{J} \mathbf{Y}.$$

O número de graus de liberdade associado a SQReg é p.



Análise de Variância

• Hipóteses:

$$H_0: \beta_1 = \beta_2 = \ldots = \beta_p = 0$$

 H_1 : existe pelo menos $\mathrm{um} eta_k
eq 0$

ou,

 H_0 : ausência de regressão

 H_1 : existe regressão

Análise de Variância

TABELA: Tabela de Análise de Variância

Fonte de Variação	SQ	g.l.	QM	F
Regressão	SQReg	р	MSReg	MSReg/MSRes
Resíduo	SQRes	n- $(p+1)$	MSRes	
Total	SQT	n-1		

Análise de Variância

Estatística do teste

$$F = \frac{MSReg}{MSRes},$$

tem distribuição F com p graus de liberdade no numerador e (n-(p+1)) graus de liberdade no denominador.

- Se $F \leq F_{(1-lpha,p,n-(p+1))} \Rightarrow$ Não há evidência para rejeitar H_0
- Se $F > F_{(1-lpha,p,n-(p+1))} \Rightarrow$ Há evidência para rejeitar H_0

COEFICIENTE DE DETERMINAÇÃO

 A medida R², chamada de Coeficiente de Determinação, representa a proporção da variação total explicada pela relação X e Y (regressão).

$$R^2 = \frac{SQReg}{SQT} = 1 - \frac{SQRes}{SQT}$$

• Desde de que $0 \le SQRes \le SQT$, segue que

$$0 \le R^2 \le 1$$

• Valores grandes de R^2 indicam que a variação total de \mathbf{Y} é reduzida pela introdução das variáveis explicativas X_1, X_2, \ldots, X_p .

Coeficiente de Determinação Ajustado

- Ao adicionar mais variáveis explicativas no modelo, a medida R² sempre aumentará. Porque, SQReg aumenta ao adicionar mais variáveis explicativas e a SQT é sempre a mesma para o conjunto de variável resposta.
- O Coeficiente de Determinação Ajustado, $R_{\rm ajust}^2$ é definido por:

$$R_{\mathsf{ajust}}^2 = 1 - \frac{\mathit{SQRes}/(\mathit{n} - (\mathit{p} + 1))}{\mathit{SQT}/(\mathit{n} - 1)},$$

 R². apenas aumentará quando for adicionado uma variável que reduzirá o MSRes.

Propriedades dos estimadores de Mínimos Quadrados

- O Teorema de Gauss Markov, afirma que os estimadores de mínimos quadrados de β_0 , β_1 , ..., β_p são não viesados. Em termos matriciais o Teorema é definido por:
- Lembrando que $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = \mathbf{CY}$
- Sendo que $\mathbf{C} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ é uma matriz de constante.
- Dessa forma, tem-se que:

$$E(\hat{eta}) = E(\mathbf{CY}) = \mathbf{C}E(\mathbf{Y}) = \mathbf{C}\mathbf{X}eta = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}eta$$
 $E(\hat{eta}) = eta$

Propriedades dos estimadores de Mínimos Quadrados

• A matriz de covariâncias de $\hat{\beta}$ é definida por:

$$Var(\hat{m{eta}}) = Var({f CY}) = {f C} Var({f Y}) {f C}'$$
como $Var({f Y}) = \sigma^2 {f I}$

$$\begin{aligned} \mathit{Var}(\hat{\boldsymbol{\beta}}) &= \mathbf{C}\sigma^{2}\mathbf{IC}' \\ &= \sigma^{2}\mathbf{CC}' \\ &= \sigma^{2}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']' \\ &= \sigma^{2}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\ &= \sigma^{2}(\mathbf{X}'\mathbf{X})^{-1} \end{aligned}$$

Propriedades dos estimadores de Mínimos Quadrados

• A matriz de covariâncias estimada é obtida ao substituir σ^2 por seu estimador, $\hat{\sigma^2}$. Ou seja,

$$\hat{Var(\hat{\beta})} = \hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1} = MSRes(\mathbf{X}'\mathbf{X})^{-1}$$

Distribuição amostral

- Ao considerar as hipóteses do modelo (2):
 - ε são independentes e identicamente distribuídos com distribuição Normal,
 - $E(\varepsilon) = \mathbf{0}$ e $Var(\varepsilon) = \sigma^2 \mathbf{I}$,
 - Portanto, **Y** são independentes e identicamente distribuídos com distribuição Normal com média $E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$ e variância $Var(\mathbf{Y}) = \sigma^2 \mathbf{I}$.
- Desde que, os estimadores $\hat{\beta}$ são uma combinação linear das observações, \mathbf{Y} , segue que:

$$\hat{oldsymbol{eta}} \sim \mathcal{N}(oldsymbol{eta}, \sigma^2(\mathbf{X}^{'}\mathbf{X})^{-1})$$

Intervalo de Confiança para β_k

• Para o modelo de regressão linear (2), tem-se a estatística:

$$\frac{\hat{eta}_k - eta_k}{\sqrt{\hat{\sigma}^2 C_{kk}}} \sim \text{t-Student}(n - (p+1)),$$

em que $k=0,1,\ldots,p$, $Var(\beta_k)=\sigma^2C_{kk}$ e C_{kk} é o k-ésimo elemento da diagonal da matriz $(\mathbf{X}'\mathbf{X})^{-1}$.

Intervalo de Confiança para β_k

• Um Intervalo de Confiança $(1 - \alpha)$ para o parâmetro β_k do modelo de regressão linear (2) é dado por:

$$IC(\beta_k, (1-\alpha)) = (\hat{\beta}_k - t_{(1-\alpha/2; n-(p+1))} \sqrt{\hat{\sigma}^2 C_{kk}}; \hat{\beta}_k + t_{(1-\alpha/2; n-(p+1))} \sqrt{\hat{\sigma}^2 C_{kk}})$$

Teste de hipótese para β_k

• Para testar a significância de algum coeficiente de regressão, β_k , as hipóteses são:

$$H_0: \beta_k = 0$$
$$H_1: \beta_k \neq 0$$

Supondo H₀ verdadeira, a estatística do teste é dada por:

$$t = rac{\hat{eta}_k - 0}{\sqrt{\hat{\sigma}^2 C_{kk}}} \sim ext{t-Student}(n - (p + 1)).$$

• A hipótese nula será rejeitada se |t| > t-Student(n - (p+1)).



Intervalo de Confiança para a resposta média

• Para alguma valor de X_1, \ldots, X_p denotado por X_{01}, \ldots, X_{0p} , a resposta média é $E(Y_0)$. O vetor \mathbf{X}_0 é definido por:

$$\mathbf{X}_0 = \begin{bmatrix} 1 \\ X_{01} \\ \vdots \\ X_{0p} \end{bmatrix}$$

O valor ajustado é definido por:

$$\hat{Y}_0 = \mathbf{X}_0' \hat{\boldsymbol{\beta}}$$

• O estimador \hat{Y}_0 é não viesado:

$$E(\hat{Y}_0) = \mathbf{X}_0' \boldsymbol{\beta} = E(Y_0)$$

Intervalo de Confiança para a resposta média

• A variância de \hat{Y}_0 é dada por:

$$Var(\hat{Y}_0) = \sigma^2 \mathbf{X}_0' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}_0$$

• Um intervalo de confiança $(1-\alpha)$ para $E(Y_0)$ é dado por:

$$(\hat{Y}_0 - t_{(1-\alpha/2;n-(p+1))} \sqrt{\hat{\sigma^2} \mathbf{X}_0' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}_0}; \hat{Y}_0 + t_{(1-\alpha/2;n-(p+1))} \sqrt{\hat{\sigma^2} \mathbf{X}_0' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}_0})$$