



Regressão Múltipla

Parte II: Intervalos de confiança e testes de hipóteses para os parâmetros, soma extra de quadrados do modelo e dados padronizados.

Inferência sobre os parâmetros da regressão

Os estimadores de mínimos quadrados ou de máxima verossimilhança são não tendenciosos, isto é: $E(\hat{\beta}) = \beta$

A matriz de variância-covariância dos estimadores é dada por:

$$\sigma^2(\hat{\beta}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \quad (28)$$

(p x p)

A estimativa da matriz de variância-covariância é dada por:

$$s^2(\hat{\beta}) = MS_R(\mathbf{X}'\mathbf{X})^{-1} \quad (29)$$

(p x p)

Intervalo de confiança para os parâmetros β_k

Para o modelo com erros normais, (17), temos:

$$\frac{\hat{\beta}_k - \beta_k}{s(\hat{\beta}_k)} \sim t(n - p) \quad k = 0, 1, \dots, p - 1 \quad (30)$$

Assim, o intervalo para β_k , com confiança $1 - \alpha$ é dado por:

$$\hat{\beta}_k \pm t(1 - \alpha / 2; n - p) s(\hat{\beta}_k) \quad (31)$$

Testes de hipóteses para β_k

Hipóteses:

$$\begin{aligned} H_0 : \beta_k &= 0 \\ H_a : \beta_k &\neq 0 \end{aligned} \quad (32)$$

Estatística de teste:

$$t^* = \frac{\hat{\beta}_k}{s(\hat{\beta}_k)} \quad (33)$$

Critério do teste:

Se $|t^*| \leq t(1-\alpha/2; n-p)$, aceita-se a hipótese nula, caso contrário rejeita-se a mesma.

Estimação da resposta média e predição de uma nova observação

Intervalo de confiança para $E(Y_h)$

Para valores dados de X_1, X_2, \dots, X_{p-1} , representados por: $X_{h1}, X_{h2}, \dots, X_{h,p-1}$, a resposta média é representada por $E(Y_h)$. Vamos definir o vetor:

$$\mathbf{X}_h = \begin{bmatrix} 1 \\ X_{h1} \\ \cdot \\ \cdot \\ X_{h,p-1} \end{bmatrix}$$

$p \times 1$

A resposta média estimada, correspondente ao vetor \mathbf{X}_h , é dada por :

$$\hat{Y}_h = \mathbf{X}_h' \hat{\boldsymbol{\beta}} \quad (34)$$

Estimação da resposta média e predição de uma nova observação

A variância estimada da resposta média é dada por:

$$s^2(\hat{Y}_h) = MS_R (\mathbf{X}'_h (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}_h) = \mathbf{X}'_h \mathbf{s}^2(\hat{\beta}) \mathbf{X}_h \quad (35)$$

O intervalo de confiança para a resposta média, $E(Y_h)$, é dado por:

$$\hat{Y}_h \pm t(1 - \alpha / 2; n - p) s(\hat{Y}_h) \quad (36)$$

Limites de predição para uma nova observação $Y_{h(novo)}$

Os limites de predição com confiança $1-\alpha$ para uma nova observação $Y_{h(nova)}$ correspondente ao vetor \mathbf{X}_h , os valores das variáveis explanatórias, são:

$$\hat{Y}_h \pm t(1 - \alpha / 2; n - p) s(pred) \quad (37)$$

A variância do erro de predição (é a diferença entre a nova observação e o valor estimado) é dado por:

$$s^2(pred) = MS_R (1 + \mathbf{X}_h' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}_h) \quad (38)$$

Método de soma extra de quadrados

- A idéia básica é verificar a redução na soma de quadrados do erro quando uma ou mais variáveis preditoras são adicionadas no modelo de regressão, dado que outras variáveis preditoras já estão incluídas no modelo. De outro lado, podemos pensar no acréscimo na soma de quadrados da regressão quando uma ou mais variáveis explanatórias são adicionadas no modelo.
- Utilização: verificar se certas variáveis X podem ser retiradas do modelo de regressão. (Construção de modelos).

Método de soma extra de quadrados

- Consider o modelo com $p-1$ preditores.
- Seja o vetor de coeficientes particionados como:
- Considere as hipóteses

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix} \quad \begin{array}{l} \text{com } p-r \text{ colunas} \\ \text{com } r \text{ colunas} \end{array}$$

$$H_0 : \boldsymbol{\beta}_2 = 0$$

$$H_a : \boldsymbol{\beta}_2 \neq 0$$

Método de soma extra de quadrados

- Modelo completo $Y = X_1\beta_1 + X_2\beta_2 + \varepsilon$
- X_1 $n \times (p-r)$ – representa as colunas dos valores de X associado com β_1
- X_2 $n \times r$ – representa coluna dos valores de X associado com β_2
- Para o modelo completo $Y = X_1\beta_1 + X_2\beta_2 + \varepsilon$ temos
$$\hat{\beta} = (XX^T)^{-1} X^T y$$
- Soma de Quadrados do modelo completo com $p-1$ graus de liberdade e soma de quadrados dos resíduos com $n-p$ graus de liberdade

$$SS_M = \hat{\beta}X^T y \quad \longrightarrow \quad MS_R = \frac{y^T y - \hat{\beta}X^T y}{n - p}$$

Método de soma extra de quadrados

- Para encontrar a contribuição de β_2 no modelo, ajuste o modelo assumindo que $\beta_2 = \mathbf{0}$ (hipótese nula é verdadeira).
- O modelo reduzido é $Y = X_1\beta_1 + \varepsilon$
- Encontra $\hat{\beta}_1 = (X_1 X_1^T)^{-1} X_1^T y$
- A soma de quadrados da regressão com $p-r$ graus de liberdade é
$$SS_M(\beta_1) = \hat{\beta}_1 X_1^T y$$
- A soma de quadrados devido a β_2 dado β_1 com $p-(p-r)=r$ graus de liberdade

$$SS_M(\beta_2 / \beta_1) = SS_M(\beta) - SS_M(\beta_1)$$

Soma extra de quadrados

Método de soma extra de quadrados

- Estatística do teste

$$F^* = \frac{SS_M(\beta_2 / \beta_1) / r}{MS_R}$$

Se $F^* > F(\alpha; r, n-p)$, rejeitamos a hipótese nula concluindo que pelo menos um dos parâmetros em β_2 é diferente de zero. Caso contrário, não rejeitamos a hipótese nula.

A soma extra de quadrados visa encontrar o melhor subconjunto de regressores para o problema.

Caso especial: Se as colunas de \mathbf{X}_1 são ortogonais as colunas de \mathbf{X}_2 ($\mathbf{X}_2^T \mathbf{X}_1 = \mathbf{0}$)

- Esse caso permite encontrar a soma de quadrados da regressão devido somente a β_2 ou seja $SS_M(\beta_2)$ sem levar em conta a dependência no regressor x_1 .

$$\begin{bmatrix} \mathbf{X}_1^T \mathbf{X}_1 & \mathbf{X}_1^T \mathbf{X}_2 \\ \mathbf{X}_2^T \mathbf{X}_1 & \mathbf{X}_2^T \mathbf{X}_2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1^T \mathbf{y} \\ \mathbf{X}_2^T \mathbf{y} \end{bmatrix}$$

- Nessas condições a solução das equações normais são

$$\hat{\beta}_1 = (\mathbf{X}_1 \mathbf{X}_1^T)^{-1} \mathbf{X}_1^T \mathbf{y} \quad \hat{\beta}_2 = (\mathbf{X}_2 \mathbf{X}_2^T)^{-1} \mathbf{X}_2^T \mathbf{y}$$

Caso especial: Se as colunas de \mathbf{X}_1 são ortogonais as colunas de \mathbf{X}_2 ($\mathbf{X}_2^T \mathbf{X}_1 = \mathbf{0}$)

- A soma de quadrados da regressão para o modelo completo é

$$SS_M(\boldsymbol{\beta}) = \hat{\boldsymbol{\beta}} \mathbf{X}^T \mathbf{y} = \hat{\boldsymbol{\beta}}_1 \mathbf{X}_1^T \mathbf{y} + \hat{\boldsymbol{\beta}}_2 \mathbf{X}_2^T \mathbf{y}$$

- Para cada conjunto temos

$$SS_M(\boldsymbol{\beta}_1) = \hat{\boldsymbol{\beta}}_1 \mathbf{X}_1^T \mathbf{y} \quad SS_M(\boldsymbol{\beta}_2) = \hat{\boldsymbol{\beta}}_2 \mathbf{X}_2^T \mathbf{y}$$

- Portanto,

$$SS_M(\boldsymbol{\beta}_1 / \boldsymbol{\beta}_2) = SS_M(\boldsymbol{\beta}) - SS_M(\boldsymbol{\beta}_2) \equiv SS_M(\boldsymbol{\beta}_1)$$

$$SS_M(\boldsymbol{\beta}_2 / \boldsymbol{\beta}_1) = SS_M(\boldsymbol{\beta}) - SS_M(\boldsymbol{\beta}_1) \equiv SS_M(\boldsymbol{\beta}_2)$$

Padronização dos dados

- É útil quando os regressores estão em escalas de unidades diferentes.
- Dois tipos mais comuns:

$$z_{ij} = \frac{x_{ij} - \bar{x}_j}{S_j} \quad y_i^* = \frac{y_i - \bar{y}}{S_y}$$

$i=1,2,\dots,n$ $j=1,2,\dots,p-1$

$$y_i^* = \beta_0 + \beta_1 z_{i1} + \dots + \beta_p z_{ip-1} + \varepsilon_i$$

$$w_{ij} = \frac{x_{ij} - \bar{x}_j}{\sqrt{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}} \quad y_i^* = \frac{y_i - \bar{y}}{\sqrt{SS_T}}$$

$i=1,2,\dots,n$ $j=1,2,\dots,p-1$

$$y_i^* = \beta_0 + \beta_1 w_{i1} + \dots + \beta_p w_{ip-1} + \varepsilon_i$$

Cada regressor tem média zero e $\sqrt{\sum_{i=1}^n (w_{ij} - \bar{w}_j)^2} = 1$