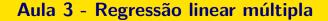
DSBD - Modelos Lineares

Slides: Cesar Taconeli, Apres.: José Padilha

11 de agosto, 2018



Introdução

 A regressão linear múltipla pode ser vista como uma extensão da regressão linear simples, em que um conjunto de variáveis independentes são utilizadas para explicar a resposta.

 Ao considerar conjuntamente o efeito de duas ou mais variáveis independentes temos condições de avaliar o efeito de uma particular variável ajustado (controlando) o efeito das demais variáveis.

 A análise de regressão linear com múltiplas variáveis é mais complexa que a regressão linear simples devido à maior dificuldade de encontrar um bom modelo, visualizar e interpretar os resultados.

• O modelo de regressão linear múltipla é definido da seguinte forma:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \epsilon_i.$$
 (1)

- Para os erros do modelo de regressão linear múltipla assumimos suposições semelhantes às consideradas para o modelo de regressão linear simples:
 - Linearidade: $E(\epsilon_i) = 0$;
 - Variância constante: $Var(\epsilon_i) = \sigma^2$;
 - Independência: ϵ_i e ϵ_j são independentes para $i \neq j$;
 - x_i é independente de ϵ_i , para todo i;
 - Normalidade: $\epsilon_i \sim N(0, \sigma^2)$.

• Como consequências da especificação do modelo, temos:

- **3** $y_i|x_i \sim N(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + ... + \beta_k x_{ik}, \sigma^2);$
- **①** Condicional aos respectivos vetores de variáveis explicativas, y_i e y_j são independentes, para todo $i \neq j$.

Observe que:

$$\frac{\partial E(y|\mathbf{x})}{\partial x_j} = \beta_j \tag{2}$$

Interpretação dos parâmetros do modelo de regressão linear múltipla

 β_j representa a alteração esperada na resposta (y) para uma unidade a mais em x_j quando todas as demais variáveis $x_k \neq x_j$ são mantidas fixas.

• Devido a essa interpretação, os parâmetros de regressão $(\beta_j's)$ são usualmente chamados **coeficientes de regressão parcial**.

• O intercepto (β_0) é a resposta esperada no ponto $x_1 = 0, x_2 = 0, \ldots, x_k = 0$, caso esse ponto pertença ao escopo do problema;

• A interpretação apresentada para os parâmetros $\beta_j's$ somente é válida na ausência de interações;

• Considere o seguinte modelo de regressão linear múltipla com termo de interação:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon.$$
 (3)

• Nesse caso, por exemplo:

$$\frac{\partial E(y|\mathbf{x})}{\partial x_1} = \beta_1 + \beta_3 x_2. \tag{4}$$

- Assim, mantendo x_2 fixa, espera-se uma variação de $\beta_1 + \beta_3 x_2$ em y para cada unidade acrescida em x_1 .
- De forma semelhante, mantendo x_1 fixa, espera-se uma variação de $\beta_2 + \beta_3 x_1$ em y para cada unidade acrescida em x_2 .
- Dessa forma, a superfície de regressão não é mais plana, pois a taxa de variação de x_1 varia conforme o valor de x_2 e vice-versa.

Notação matricial do modelo

• Considere *n* observações (y_i, \mathbf{x}_i) , em que $\mathbf{x}_i = (x_{i1}, x_{i2}, ..., x_{ik})'$:

$$y_{1} = \beta_{0} + \beta_{1}x_{11} + \beta_{2}x_{12} + \dots + \beta_{k}x_{1k} + \epsilon_{1}$$

$$y_{2} = \beta_{0} + \beta_{1}x_{21} + \beta_{2}x_{22} + \dots + \beta_{k}x_{2k} + \epsilon_{2}$$

$$\vdots$$

$$y_{n} = \beta_{0} + \beta_{1}x_{n1} + \beta_{2}x_{n2} + \dots + \beta_{k}x_{nk} + \epsilon_{n}$$

Notação matricial do modelo

 O modelo de regressão linear múltipla fica representado matricialmente por:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},\tag{5}$$

em que

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}, \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Notação matricial do modelo

• As especificações e resultados apresentados anteriormente podem ser representados na forma matricial:

$$2 Var(\epsilon) = \sigma^2 I;$$

 $\mathbf{v} | \mathbf{X} \sim Normal(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}).$

- Vamos considerar n observações (y_i, \mathbf{x}_i) , em que $\mathbf{x}_i = (x_{i1}, x_{i2}, ..., x_{ik})'$.
- A estimação de mínimos quadrados para o modelo de regressão linear múltipla baseia-se, novamente, na determinação de $\beta_0, \beta_1, ..., \beta_k$ que minimizem a soma de quadrados dos erros:

$$S = S(\beta_0, \beta_1, ..., \beta_k) = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i1} + ... + \beta_k x_{ik}))^2$$
 (6)

• Os estimadores de mínimos quadrados para $\beta_0, \beta_1, ..., \beta_k$ devem satisfazer:

$$\frac{\partial S(\beta)}{\partial \beta} = \begin{bmatrix}
\frac{\partial S(\beta)}{\partial \beta_0} \\
\frac{\partial S(\beta)}{\partial \beta_1} \\
\frac{\partial S(\beta)}{\partial \beta_2} \\
\vdots \\
\frac{\partial S(\beta)}{\partial \beta_t}
\end{bmatrix} = \begin{bmatrix}
0 \\
0 \\
0 \\
0 \\
0
\end{bmatrix}$$
(7)

 Derivando parcialmente em relação aos parâmetros de regressão obtemos:

$$\frac{\partial S}{\partial \beta_0}\Big|_{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k} = -2\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^k \beta_j x_{ij} \right) = 0$$
 (8)

е

$$\frac{\partial S}{\partial \beta_j}\Big|_{\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_k} = -2\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^k \beta_j x_{ij} \right) x_{ij} = 0, \quad j = 1, 2, ..., k. \quad (9)$$

Na forma matricial:

$$S(\beta) = \sum_{i=1}^{n} \epsilon_i^2 = \epsilon' \epsilon = (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta), \tag{10}$$

de maneira que o vetor \hat{eta} tal que:

$$\left. \frac{\partial S}{\partial \boldsymbol{\beta}} \right|_{\hat{\boldsymbol{\beta}}} = \mathbf{0} \tag{11}$$

é o estimador de mínimos quadrados de β , dado por:

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

- Observe que os estimadores de mínimos quadrados somente existem se a matriz (X'X)⁻¹ existe;
- A condição de existência de $(X'X)^{-1}$ é que as colunas de X sejam linearmente independentes, ou seja, que nenhuma coluna de X seja combinação linear das demais;
- O modelo ajustado, para um vetor $\mathbf{x} = (1, x_1, x_2, ..., x_k)$ fica denotado por:

$$\hat{y} = \mathbf{x'}\hat{\boldsymbol{\beta}} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_k x_k$$
 (13)

• O vetor de valores ajustados $\hat{\boldsymbol{y}} = (\hat{y}_1, \hat{y}_2, ..., \hat{y}_n)$ é dado por:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{H}\mathbf{y}. \tag{14}$$

- A matriz $n \times n$ $H = X(X'X)^{-1}X'$, chamada matriz chapéu (hat matrix) mapeia o vetor de valores observados no vetor de valores ajustados.
- O vetor de resíduos $\mathbf{r} = (r_1, r_2, ..., r_n)$ fica definido, em notação matricial, por:

$$\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}},\tag{15}$$

em que $r_i = y_i - \hat{y}_i$, i = 1, 2, ..., n.

- Propriedades dos estimadores:
- $E(\hat{\beta}) = \beta \ (\hat{\beta} \text{ é um estimador não viciado de } \beta);$
- **2** $Var(\hat{\beta}) = \sigma^2(X'X)^{-1};$
- $\hat{\beta}$ é o melhor (mais eficiente) estimador linear não viciado de β (teorema de Gauss Markov);
- Sob a suposição de que os erros têm distribuição normal os estimadores de mínimos quadrados equivalem aos de máxima verossimilhança.

Estimação de σ^2

• Um estimador não viciado para σ^2 , baseado na soma de quadrados de resíduos, é dado por:

$$\hat{\sigma}^2 = QM_{Res} = \frac{SQ_{Res}}{n - p} = \frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{n - p},$$
 (16)

em que p = k + 1 é o número de parâmetros do modelo.

Testes de hipóteses e intervalos de confiança para os parâmetros do modelo de RLM

• Assim como no caso de RLS, também na RLM a inferência sobre os parâmetros do modelo é um ponto importante, que permitirá:

- * Checar a significância do modelo ajustado;
- * Identificar quais variáveis explicativas são relevantes na a
- * Avaliar o erro de estimativas e predições geradas pelo mode
 - Deste ponto em diante assumiremos todas as suposições especificadas para os erros, inclusive a de normalidade.

 Na análise de variância em regressão linear múltipla, a variação total (corrigida pela média) é novamente decomposta em dois componentes: variação explicada pela regressão e variação residual, tal que:

$$SQ_{Total} = SQ_{Reg} + SQ_{Res}.$$
 (17)

Usando notação matricial, as somas de quadrados ficam definidas por:

$$SQ_{Res} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \mathbf{y'y} - \hat{\beta}' \mathbf{X'y};$$
(18)

$$SQ_{Reg} = \sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2 = \hat{\beta}' X' y - \frac{\left(\sum_{i=1}^{n} y_i\right)^2}{n};$$
(19)

$$SQ_{Total} = \sum_{i=1}^{n} (y_i - \bar{y})^2 = \mathbf{y'y} - \frac{(\sum_{i=1}^{n} y_i)^2}{n}.$$
 (20)

Tabela 1: Quadro de análise de variânciapara o modelo de RLM

Fonte de variação	Graus de liberdade	Soma de quadrados	Quadrados médios	F
Regressão	p-1	$\hat{\beta}' \mathbf{X}' \mathbf{y} - \frac{\left(\sum_{i=1}^{n} y_i\right)^2}{n}$	$QM_{Reg} = \frac{SQ_{Reg}}{p-1}$	$F = \frac{QM_{Reg}}{QM_{Res}}$
Resíduos	n - p	$y'y - \hat{eta}'X'y$	$QM_{Res} = \frac{SQ_{Res}}{n-p}$	
Total	n-1	$\mathbf{y'y} - \frac{\left(\sum_{i=1}^{n} y_i\right)^2}{n}$		

• Vale lembrar que n é o tamanho da amostra e p=k+1 o número de parâmetros do modelo.

 Podemos testar a significância do modelo ajustado com base no seguinte par de hipóteses:

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

$$H_1: \beta_j \neq 0$$
 para pelo menos um j .

- Sob a hipótese nula (não significância do modelo) a estatística F segue distribuição F—Snedecor, com p-1 e n-p graus de liberdade.
- Assim, fixado um nível de significância α , H_0 deve ser rejeitada se o valor da estatística F for maior que o quantil $1-\alpha$ da distribuição $F_{p-1,n-p}$.

• O coeficiente de determinação, como anteriormente, fica definido por:

$$R^2 = 1 - \frac{SQ_{Res}}{SQ_{Total}} = \frac{SQ_{Reg}}{SQ_{Total}},$$
 (21)

que expressa a proporção da variabilidade original dos dados explicada pelo modelo de regressão.

 Uma propriedade de R² que o torna pouco apropriado para a comparação dos ajustes de diferentes modelos é que ele nuna decresce à medida que incluímos novas variáveis ao modelo.

• Como alternativa ao \mathbb{R}^2 podemos considerar o \mathbb{R}^2 ajustado, definido por:

$$R_{Aj}^{2} = 1 - \frac{SQ_{Res}/(n-p)}{SQ_{Total}/(n-1)}.$$
 (22)

- Como $SQ_{Total}/(n-1)$ é fixo, então R_{Aj}^2 somente aumentará se houver redução do quadrado médio de resíduos.
- Diferentemente de R^2 , R_{Aj}^2 penaliza a inclusão de variáveis não importantes no modelo, permitindo comparar adequadamente modelos com diferentes complexidades (números de variáveis).

TH's e IC's para os parâmetros do modelo de regressão linear múltipla

- Primeiramente vamos considerar TH's e IC's para parâmetros individuais do modelo.
- Suponha que se deseja testar a significância de x_j no modelo. Partimos do seguinte par de hipóteses:

$$H_0: \beta_j = 0 \quad vs \quad H_1: \beta_j \neq 0.$$
 (23)

A estatística do teste é dada por:

$$t = \frac{\hat{\beta}_j}{ep(\hat{\beta}_j)},\tag{24}$$

em que $ep(\hat{\beta}_j) = \sqrt{\hat{\sigma}^2(\boldsymbol{X}'\boldsymbol{X})_{jj}^{-1}}$, sendo $(\boldsymbol{X}'\boldsymbol{X})_{jj}^{-1}$ o j – ésimo termo da diagonal de $(\boldsymbol{X}'\boldsymbol{X})^{-1}$ e $\hat{\sigma}^2 = QM_{Res}$.

TH's e IC's para os parâmetros do modelo de regressão linear múltipla

- Sob a hipótese nula a estatística t tem distribuição t-Student com n-p graus de liberdade.
- Assim, a hipótese H_0 deverá ser rejeitada, para um nível de significância α , se $|t| > |t_{n-p,\alpha/2}|$, em que $t_{n-p,\alpha/2}$ é o quantil $\alpha/2$ da distribuição t - Student com n - p graus de liberdade.
- Usando a distribuição t_{n-p} como referência, um intervalo de confiança $100(1-\alpha)\%$ para β_i fica definido por:

$$\hat{\beta}_j \mp t_{n-p,\alpha/2} \sqrt{\hat{\sigma}^2 (\boldsymbol{X}' \boldsymbol{X})_{jj}^{-1}}.$$
 (25)

• Para qualquer valor β_{i0} pertencente ao intervalo de confiança não se tem evidências, ao nível de significância α , que $\beta_i \neq \beta_{i0}$.

Intervalo de confiança para a resposta média e para uma predição

- Considere interesse em estimar a resposta média em um ponto $\mathbf{x'}_0 = (1, x_{01}, x_{02}, ..., x_{0k})$, ou seja, $E(y|\mathbf{x}_0)$.
- A estimativa pontual é dada pelo valor ajustado pelo modelo em x_0 :

$$\widehat{E(y|\mathbf{x}_0)} = \hat{y}_0 = \mathbf{x}_0' \hat{\boldsymbol{\beta}}. \tag{26}$$

 O estimador apresentado é não viciado para a real resposta média, com variância:

$$Var(\widehat{E(y|\mathbf{x}_0)}) = \sigma^2 \mathbf{x}_0'(\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0.$$
(27)

Intervalo de confiança para a resposta média e para uma predição

• Um intervalo de confiança $100(1-\alpha)\%$ para a resposta média em $\mathbf{x'}_0 = (1, x_{01}, x_{02}, ..., x_{0k})$ é dado por:

$$\widehat{E(y|\mathbf{x}_0)} \mp t_{n-p,\alpha/2} \sqrt{\hat{\sigma}^2 \mathbf{x}_0' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0}.$$
 (28)

em que $\hat{\sigma}^2 = QM_{Res}$.

- Considere agora que se deseja predizer a resposta em um ponto $\mathbf{x'}_0 = (1, x_{01}, x_{02}, ..., x_{0k}).$
- A estimativa pontual, novamente, é dada pelo valor ajustado de y em x'_0 :

$$\hat{y}_0 = \mathbf{x}_0' \hat{\boldsymbol{\beta}}. \tag{29}$$

Intervalo de confiança para a resposta média e para uma predição

• Neste caso, a variância de \hat{y}_0 fica dada por:

$$Var(\hat{y}_0) = \sigma^2 \left(1 + x_0' (X'X)^{-1} x_0 \right).$$
 (30)

• Um intervalo de confiança $100(1-\alpha)\%$ para a predição de uma nova observação em \mathbf{x}_0 fica dada por:

$$\hat{y}_0 \mp t_{n-p,\alpha/2} \sqrt{\hat{\sigma}^2 \left(1 + x_0' (X'X)^{-1} x_0\right)},$$
 (31)

em que $\hat{\sigma}^2 = QM_{Res}$.

- Em geral os estimadores dos parâmetros do modelo de RLM são correlacionados (a menos que as correspondentes variáveis sejam independentes);
- Avaliar a significância individual das variáveis e avaliar a significância conjunta das mesmas, neste caso, são coisas distintas.
- Em alguns casos temos interesse particular em analisar a significância conjunta de dois ou mais parâmetros, como no caso de modelos polinomiais e com variáveis indicadoras.
- Vamos abordar testes de hipóteses simultâneos para dois ou mais parâmetros do modelo usando o princípio da verossimilhança.

• Considere o modelo de regressão linear múltipla:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \epsilon, \tag{32}$$

 $\hat{\beta}' = (\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_k)$ o vetor de estimativas de mínimos quadrados e $\hat{\sigma}^2 = SQ_{Res}/n$ a estimativa de máxima verossimilhança para σ^2 .

• O interesse aqui é testar uma hipótese do tipo $H_0=\beta_1=\beta_2=...=\beta_q=0,\ q\leq k.$ Por simplicidade de notação, vamos considerar que a hipótese nula contemple os q primeiros parâmetros do modelo.

• O modelo induzido pela hipótese nula é dado por:

$$y = \beta_0 + 0x_1 + 0x_2 + \dots + 0x_q + \beta_{q+1}x_{q+1} + \dots + \beta_k x_k + \epsilon$$

= $\beta_0 + \beta_{q+1}x_{q+1} + \dots + \beta_k x_k + \epsilon$ (33)

• Vamos denotar por $\hat{\beta}_{\mathbf{0}}' = (\hat{\beta}_{0}', 0, 0, ..., \hat{\beta'}_{q+1}, ..., \hat{\beta'}_{k})$ o estimador de mínimos quadrados para o modelo restrito.

 A verossimilhança para o modelo completo, avaliada nas estimativas de máxima verossimilhança, é dada por:

$$L = \left(2\pi \frac{SQ_{Res}}{n}\right)^{-n/2},\tag{34}$$

em que SQ_{Res} é a soma de quadrados de resíduos do modelo.

Para o modelo restrito a verossimilhança maximizada fica dada por:

$$L_0 = \left(2\pi \frac{SQ_{Res_0}}{n}\right)^{-n/2},\tag{35}$$

em que SQ_{Res_0} é a soma de quadrados de resíduos para o modelo restrito (ajustado apenas com as k-q variáveis não restritas a zero.

• O teste da razão de verossimilhanças para testar H_0 baseia-se na estatística da razão de verossimilhanças:

$$\frac{L_0}{L} = \left(\frac{SQ_{Res_0}}{SQ_{Res}}\right)^{-n/2} = \left(\frac{SQ_{Res}}{SQ_{Res_0}}\right)^{n/2}.$$
 (36)

• Sob a hipótese H_0 , assintoticamente:

$$\Lambda = -2 \ln \left(\frac{L_0}{L} \right) \sim \chi_q^2, \tag{37}$$

em que χ_q^2 denota a distribuição qui-quadrado com q graus de liberdade.

• Para um nível de significância α , H_0 será rejeitada se Λ superar o quantil $1-\alpha$ da distribuição χ_q^2 .

- No caso de modelos lineares, no entanto, temos um teste exato como alternativa ao teste assintótico baseado na distribuição qui-quadrado;
- Sob a hipótese nula, a estatística:

$$F_0 = \frac{(SQ_{Res_0} - SQ_{Res})/q}{SQ_{Res}/(n-p)}$$
(38)

tem distribuição F-Snedecor com q e n-p graus de liberdade.

- Observe que F_0 baseia-se na variação da soma de quadrados de resíduos resultante da restrição aplicada a q parâmetros do modelo.
- A hipótese H_0 deverá ser rejeitada, ao nível de significância α , se F_0 superar o quantil $1-\alpha$ da distribuição F-Snedecor com q e n-p graus de liberdade.

• A estatística F_0 pode ser calculada por:

$$F_0 = \frac{(\hat{\beta}_q - \beta_q^{(0)})' \mathbf{V}_{11}^{-1} (\hat{\beta}_q - \beta_q^{(0)})}{qQM_{Res}},$$
 (39)

em que $\hat{\beta}_q$ denota o vetor de q entradas de $\hat{\beta}$ referente aos parâmetros restritos e V_{11} a matriz quadrada com as q entradas (linhas e colunas) de $(X'X)^{-1}$.

• Repare que nesta representação $\beta_q^{(0)}$ representa o vetor postulado para os q parâmetros restritos sob H_0 (geralmente um vetor de zeros).

• O teste da significância do modelo de regressão $(H_0: \beta_1 = \beta_2 = ... = \beta_k = 0)$ é um caso particular, em que a estatística F, apresentada no quadro da análise de variância, tem distribuição F-Snedecor com p-1 e n-p graus de liberdade sob H_0 .

Região de confiança

- Suponha interesse em estimar simultaneamente algum subconjunto de parâmetros do modelo.
- Seja β_q um subconjunto de elementos de β , contendo parâmetros sobre os quais se deseja inferir.
- Adicionalmente, seja $\hat{\beta}_{m{q}}$ o vetor de estimadores de mínimos quadrados de $\beta_{m{q}}.$
- Uma região de confiança $100(1-\alpha)\%$ para os componentes de β_q é definido pelo conjunto de todos os vetores $\beta_q^{(0)}$ tais que:

$$F_0 = \frac{(\hat{\beta}_1 - \beta_1^{(0)})' \mathbf{V}_{11}^{-1} (\hat{\beta}_1 - \beta_1^{(0)})}{qQM_{Res}} \le F_{q,n-p} (1 - \alpha)$$
(40)

em que $F_{q,n-p}(\alpha)$ é o quantil $1-\alpha$ da distribuição F-Snedecor com q e n-p graus de liberdade.

Testes de hipóteses para combinações lineares dos parâmetros

• De forma mais geral, podemos definir hipóteses lineares na forma:

$$H_0: \mathbf{L}\beta = \mathbf{c},\tag{41}$$

em que \boldsymbol{L} é uma matriz de constantes de dimensão $q \times p$, de rank linha completo, e \boldsymbol{c} um vetor de constantes de dimensão q (ambos especificados).

 Neste caso, H₀ compreende q hipóteses lineares sobre os parâmetros do modelo, do tipo:

$$\begin{split} L_{11}\beta_0 + L_{12}\beta_1 + L_{13}\beta_2 + \dots + L_{1p}\beta_k &= c_1 \\ L_{21}\beta_0 + L_{22}\beta_1 + L_{23}\beta_2 + \dots + L_{2p}\beta_k &= c_2 \\ &\vdots \\ L_{a1}\beta_0 + L_{a2}\beta_1 + L_{a3}\beta_2 + \dots + L_{ap}\beta_k &= c_a \end{split}$$

Testes de hipóteses para combinações lineares dos parâmetros

• Sob a hipótese H_0 a estatística:

$$F = \frac{(\boldsymbol{L}\beta - \boldsymbol{c})'[\boldsymbol{L}(\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{L}']^{-1}(\boldsymbol{L}\beta - \boldsymbol{c})}{qQM_{Res}}$$
(42)

tem distribuição F-Snedecor com q e n-p graus de liberdade.

• Assim, a hipótese nula será rejeitada, ao nível de significância α , se o valor calculado da estatística F exceder o quantil $1-\alpha$ da distribuição F-Snedecor com q e n-p graus de liberdade.

Intervalos de confiança para combinações lineares dos parâmetros

- Seja $\mathbf{l'} = (l_0, l_1, l_2, ..., l_p)$ um vetor de constantes e considere interesse em estimar $\theta = \mathbf{l'}\beta$.
- A estimativa pontual de $l'\beta$ é dada por $l'\hat{\beta}$.
- Um intervalo de confiança $100(1-\alpha)\%$ para $I'\beta$ tem limites:

$$\mathbf{I}'\hat{\boldsymbol{\beta}} \mp t_{n-p,\alpha/2} \sqrt{\hat{\sigma}^2 \mathbf{I}' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{I}}.$$
 (43)

Regressão com coeficientes padronizados

- Na análise de RLM, a comparação das magnitudes dos $\hat{\beta}'_j s$ em sempre é possível devido ao impacto das diferentes unidades de medidas dos $x'_j s$.
- Caso seja desejado que tais estimativas sejam comparáveis, pode-se padronizar cada uma das variáveis de forma que as variáveis resultantes tenham uma mesma escala.
- Uma allternativa de padronização consiste em 'normalizar' cada uma das variáveis, aplicando:

$$z_{ij} = \frac{x_{ij} - \bar{x}_j}{s_j}, \quad i = 1, 2, ..., n; \ j = 1, 2, ..., k,$$
 (44)

е

$$y_i^* = \frac{y_i - \bar{y}}{s_v}, \quad i = 1, 2, ..., n,$$
 (45)

Regressão com coeficientes padronizados

- Neste caso, \bar{x}_j e s_j são a média e o desvio padrão amostrais de x_j e \bar{y} e s_y a média e desvio padrão amostrais de y.
- Usando as variáveis normalizadas, o modelo de regressão linear múltipla fica definido por:

$$y_i^* = b_1 z_{i1} + b_2 z_{i2} + ... + b_k z_{ik} + \epsilon_i, \quad i = 1, 2, ..., n.$$
 (46)

 A análise segue então da maneira usual de forma que o estimador de mínimos quadrados de b fica dado por:

$$\hat{\boldsymbol{b}} = (\boldsymbol{Z'Z})^{-1}\boldsymbol{Z'y}. \tag{47}$$

Regressão com coeficientes padronizados

- Ao centrar as variáveis, o intercepto do modelo é deslocado para zero.
- As interpretações dos parâmetros do modelo devem ser feitas em termos dos valores padronizados das variáveis originais.