



**LEANDRO DA SILVA PEREIRA**

**GEOMETRIA DOS MÉTODOS DE  
REGRESSÃO LARS, LASSO E ELASTIC NET  
COM UMA APLICAÇÃO EM SELEÇÃO  
GENÔMICA**

**LAVRAS - MG**

**2017**

**LEANDRO DA SILVA PEREIRA**

**GEOMETRIA DOS MÉTODOS DE REGRESSÃO LARS, LASSO E  
ELASTIC NET COM UMA APLICAÇÃO EM SELEÇÃO GENÔMICA**

Tese apresentada à Universidade Federal de Lavras, como parte das exigências do Programa de Pós-graduação em Estatística e Experimentação Agropecuária, área de concentração em Estatística e Experimentação Agropecuária, para a obtenção do título de Doutor.

Orientador  
Dr. Lucas Monteiro Chaves

**LAVRAS - MG  
2017**

**Ficha catalográfica elaborada pelo Sistema de Geração de Ficha Catalográfica da Biblioteca  
Universitária da UFLA, com dados informados pelo(a) próprio(a) autor(a).**

Pereira, Leandro da Silva.

Geometria dos métodos de regressão LARS, LASSO e Elastic  
Net com uma aplicação em seleção genômica. / Leandro da Silva  
Pereira. - 2017.

167 p. : il.

Orientador(a): Lucas Monteiro Chaves.

Tese (doutorado) - Universidade Federal de Lavras, 2017.  
Bibliografia.

1. Seleção de covariáveis. 2. Geometria de modelos lineares. 3.  
Regressão Elastic Net. I. Chaves, Lucas Monteiro. . II. Título.

**LEANDRO DA SILVA PEREIRA**

**GEOMETRIA DOS MÉTODOS DE REGRESSÃO LARS, LASSO E  
ELASTIC NET COM UMA APLICAÇÃO EM SELEÇÃO GENÔMICA**

Tese apresentada à Universidade Federal de Lavras, como parte das exigências do Programa de Pós-graduação em Estatística e Experimentação Agropecuária, área de concentração em Estatística e Experimentação Agropecuária, para a obtenção do título de Doutor.

APROVADA em 27 de Janeiro de 2017.

|  |      |
|--|------|
| Dr. Lucas Monteiro Chaves (Orientador) | UFLA |
| Dra. Thelma Sáfadi                     | UFLA |
| Dr. Daniel Furtado Ferreira            | UFLA |
| Dr. José Airton Rodrigues Nunes        | UFLA |
| Dr. Fabyano Fonseca e Silva            | UFV  |

**LAVRAS - MG  
2017**

“Por vezes sentimos que aquilo que fazemos não é senão uma gota de  
água no mar. Mas o mar seria menor se lhe faltasse uma gota.”

Santa Madre Teresa de Calcutá

## AGRADECIMENTOS

A Deus, pelo sol que ilumina os meus dias, pelos corpos celestes que clareiam as minhas noites mais escuras, pelas chuvas que lavam a minha alma, pelas primaveras que colorem a minha vida com as suas flores, pelos pássaros que alegam o meu espírito com os seus cantos, pelos alimentos que fortificam o meu corpo, pela verdade que aos poucos me liberta, pelo trabalho que me faz sentir útil, pela vida que me proporciona novas experiências e oportunidades de crescimento, e por tudo que ainda não tenho consciência de que ELE faz por mim.

Aos meus pais Messias e Valdelice que, pelo exemplo de humildade e dignidade, me mostram o sentido da vida.

À minha eterna companheira Mayumi, sinônimo de vida e esperança, por ser meu forte nas horas difíceis e meu encanto nas horas de paz.

À minha irmã Celiani e sua família, pelo apoio direto e indireto.

A todos os meus familiares e amigos, pelo apoio e carinho.

Ao meu orientador Lucas Monteiro Chaves, que pelos conhecimentos e esclarecimentos intelectuais e morais confiados a mim, me ensinou o verdadeiro sentido de se ensinar.

Ao Prof. Devanil Jaques de Souza, que de forma síncrona com as ideias do professor Lucas, contribui para minha formação moral e intelectual.

À Universidade Federal de Lavras (UFLA) e ao Departamento de Estatística (DES), pela oportunidade de cursar o doutorado.

Aos professores do programa de pós-graduação em estatística e experimentação agropecuária da UFLA, pelas contribuições na minha formação durante as suas disciplinas, em especial aos professores Renato Ribeiro de Lima, Júlio de Sousa Bueno Filho, Thelma Sáfadi e Daniel Furtado Ferreira.

Às funcionárias do DEX: Kelly, Josi, Maria e Lu, pela amizade e cari-

nho, em especial a Nádia, pelo companheirismo e compreensão nas minhas várias ligações e solicitações burocráticas.

À FAPEMIG, pela concessão da bolsa de estudos no início do doutorado, tornando financeiramente possível a realização do mestrado. Aos colegas da Engenharia de Controle e Automação pelos momentos de lazer e pelos conselhos “etilizados”.

A todos os colegas de turma pela companhia nos estudos, nos momentos de lazer e nos momentos de tensão de véspera das provas.

A UTFPR pelo apoio com a licença para finalização do doutorado.

A todos que de alguma forma contribuíram, mesmo inconsciente, para a realização deste trabalho.

## RESUMO

Os métodos de estimação e seleção de variáveis em modelos lineares, LASSO (Least Absolute Shrinkage and Selection Operator), LARS (Least Angle Regression) e Elastic Net, são abordados utilizando ênfase em aspectos geométricos. O texto se propõe a ser uma leitura auxiliar aos artigos clássicos de Tibshirani, Hastie, Efron, Zou e Johnstone, apresentando de forma mais detalhada alguns dos resultados citados nos referidos artigos. Tal ponto de vista não ocorre na literatura básica relativa a estes métodos, e neste sentido o trabalho representa uma contribuição original ao assunto. Simulações utilizando a linguagem R (pacote glmnet) são desenvolvidas para se estudar o comportamento dos estimadores. Para um conjunto de dados de suínos *Sus scrofa*, são analisados 237 marcadores genéticos (SNPs) pelo método Elastic Net, para a resposta relativa ao pH da carne e peso de carcaça.

Palavras-chave: Seleção de covariáveis. Geometria de modelos lineares. Marcadores genéticos. Algoritmo LARS.



## ABSTRACT

The methods of estimation and variable selection in linear models, LASSO (Least Absolute Shrinkage and Selection Operator), LARS (Least Angle Regression) and Elastic Net, are addressed using emphasis in terms of a geometric approach. The present work proposes to be an auxiliary reading to the classic articles of Tibshirani, Hastie, Efron, Zou and Johnstone, presenting in more detail some of the results cited in such articles. Such a point of view does not occur in the basic literature regarding these methods, and in this sense the work represents an original contribution to the subject. Simulations using R code (glmnet package) are developed to study the behavior of the estimators. For a data set of *Sus scrofa* pork, 237 genetic markers (SNPs) are analyzed using the Elastic Net method, using as response the pork pH and carcass weight.

Keywords: Covariates selection. Linear models geometry. Genetic markers. LARS Algorithm.

Dr. Lucas Monteiro Chaves  
Orientador

## LISTA DE FIGURAS

|           |   |    |
|-----------|---|----|
| Figura 1  | Representação geométrica de um vetor n-dimensional. . . . .   | 16 |
| Figura 2  | Reta ajustada por quadrados mínimos. . . . .  | 17 |
| Figura 3  | Geometria da regressão linear simples. . . . .  | 19 |
| Figura 4  | Geometria da regressão linear múltipla. . . . .   | 20 |
| Figura 5  | Interpretação geométrica das matrizes $A$ e $A'$ . . . . .  | 23 |
| Figura 6  | O triângulo fundamental. . . . .  | 25 |
| Figura 7  | Projeções em termos de somas direta. . . . .  | 27 |
| Figura 8  | Demonstração geométrica da Proposição 3 . . . . .   | 28 |
| Figura 9  | Reparametrização em regressão linear. . . . .   | 30 |
| Figura 10 | Demonstração da unicidade da projeção em convexos. . . . .  | 36 |
| Figura 11 | Projeção do vetor $\mathbf{y}$ no convexo $K$ . . . . .   | 37 |
| Figura 12 | Relação entre $d(\mathbf{y}, \mathbf{k})$ e $d(\mathbf{y}, P_{\text{Im}(X)}(\mathbf{y}))$ . . . . . | 37 |
| Figura 13 | Isometria entre o espaço de dados e o espaço de parâmetros. . . . .                                 | 38 |
| Figura 14 | O Estimador de quadrados mínimos. . . . .   | 40 |
| Figura 15 | Transformações ortogonais preservam elipsoides. . . . .   | 42 |
| Figura 16 | Transformações de elipsoides em esferas. . . . .  | 43 |
| Figura 17 | Transformação de esferas em elipsoides. . . . .   | 43 |
| Figura 18 | Vetores equidistantes à $X\hat{\beta}_{\text{ols}}$ . . . . .                                       | 45 |
| Figura 19 | Relação entre a hipersfera e o elipsoide. . . . .   | 46 |
| Figura 20 | Varição do raio da hipersfera. . . . .  | 46 |
| Figura 21 | O estimador ridge. . . . .  | 47 |
| Figura 22 | Exemplo de <i>ridgetrace</i> em $\mathbb{R}^6$ . . . . .  | 48 |
| Figura 23 | Penalização na definição do estimador Ridge. . . . .  | 49 |
| Figura 24 | A geometria do estimador Ridge. . . . .   | 50 |
| Figura 25 | Projeção definida pela métrica de Mahalanobis. . . . .  | 53 |
| Figura 26 | Restrição $K_p$ LASSO. . . . .  | 55 |
| Figura 27 | Obtenção do estimador LASSO. . . . .  | 56 |
| Figura 28 | O método de estimação LASSO. . . . .  | 57 |
| Figura 29 | Estimador LASSO no caso bidimensional. . . . .  | 59 |
| Figura 30 | Estimador LASSO para o caso ortogonal. . . . .  | 60 |
| Figura 31 | Estimador LASSO (caso geral). . . . .   | 62 |
| Figura 32 | Matriz de pesos do estimador Ridge. . . . .   | 65 |
| Figura 33 | Aproximação do estimador LASSO. . . . .   | 70 |
| Figura 34 | Gráfico das funções coordenadas. . . . .  | 71 |
| Figura 35 | $\text{Ker}(X)$ e o convexo $K_p$ . . . . .   | 73 |
| Figura 36 | O convexo $C_p$ em $\mathbb{R}^2$ e $\mathbb{R}^3$ . . . . .  | 74 |
| Figura 37 | Penalização do método GARROTE. . . . .  | 76 |

|           |  |     |
|-----------|--|-----|
| Figura 38 | $K_p$ na estimação <i>elastic net</i> . . . . .                                      | 78  |
| Figura 39 | Geometria do estimador <i>elastic net</i> . . . . .                                  | 79  |
| Figura 40 | O método de regressão <i>elastic net</i> . . . . .                                   | 80  |
| Figura 41 | Vetor de dados e das covariáveis. . . . .  | 89  |
| Figura 42 | Vetores normalizados das covariáveis. . . . .  | 90  |
| Figura 43 | Geometria da correlação amostral. . . . .  | 90  |
| Figura 44 | Normalização das covariáveis. . . . .  | 90  |
| Figura 45 | Stepwise e variáveis altamente correlacionadas. . . . .                              | 92  |
| Figura 46 | Algoritmo <i>stagewise</i> . . . . .   | 93  |
| Figura 47 | Geometria do LARS para a dimensão 2. . . . .   | 94  |
| Figura 48 | Geometria do LARS para a dimensão 3. . . . .   | 95  |
| Figura 49 | Vetores coluna para X ortonormais. . . . .   | 99  |
| Figura 50 | Relação LARS e OLS. . . . .  | 104 |
| Figura 51 | Geometria do conjunto $S_A$ . . . . .  | 107 |
| Figura 52 | Geometria do erro de predição . . . . .  | 112 |
| Figura 53 | Elementos representativos da estrutura $C_p$ . . . . .                               | 126 |
| Figura 54 | Comportamento do df em relação ao parâmetro Ridge. . . . .                           | 134 |
| Figura 55 | Matriz de projeção. . . . .  | 135 |
| Figura 56 | Estimação de parâmetros <i>Elastic Net</i> . . . . .                                 | 138 |
| Figura 57 | Exemplo de um plot com legenda. . . . .  | 139 |
| Figura 58 | Comparação entre LASSO Trace e Elastic Net Trace. . . . .                            | 142 |
| Figura 59 | Elastic Net trace para a situação 1. . . . .   | 143 |
| Figura 60 | Elastic Net trace para a situação 2. . . . .   | 144 |
| Figura 61 | Elastic Net trace para os dados de câncer de próstata. . . . .                       | 144 |
| Figura 62 | Elastic Net trace para covariáveis categóricos. . . . .                              | 145 |
| Figura 63 | Estimativa do vetor de parâmetros dos dados categóricos. . . . .                     | 146 |
| Figura 64 | Elastic Net trace das 237 covariáveis. . . . .                                       | 147 |
| Figura 65 | Elastic Net trace para as 237 covariáveis com a resposta PCARC. . . . .              | 149 |
| Figura 66 | As 12 covariáveis mais relevantes para a resposta PCARC. . . . .                     | 150 |
| Figura 67 | Posição das 12 covariáveis mais relevantes em relação aos grupos de ligação. . . . . | 150 |
| Figura 68 | As 13 covariáveis mais relevantes para a resposta PH45. . . . .                      | 151 |
| Figura 69 | As 9 covariáveis mais relevantes para a resposta PH45 no cromossomo SSC4. . . . .    | 152 |
| Figura 70 | Elastic Net trace para os marcadores de 90 a 99. . . . .                             | 153 |

## **LISTA DE TABELAS**

## SUMÁRIO

|       |  |     |
|-------|--|-----|
| 1     | INTRODUÇÃO . . . . .   | 14  |
| 2     | Uma abordagem geométrica à regressão linear múltipla . . . . | 17  |
| 2.1   | O triângulo fundamental . . . . .                            | 24  |
| 2.2   | Reparametrização em sistemas lineares . . . . .              | 29  |
| 2.3   | Reparametrização na presença de covariáveis categóricas . .  | 32  |
| 3     | Teoria geral da estimação . . . . .                          | 35  |
| 3.1   | Estimador de quadrados mínimos . . . . .                     | 38  |
| 4     | A regressão Ridge . . . . .                                  | 46  |
| 4.1   | O estimador Ridge . . . . .                                  | 48  |
| 4.2   | LASSO . . . . .  | 53  |
| 4.3   | Abordagem Bayesiana aos estimadores Ridge e LASSO . . .      | 62  |
| 4.4   | Estimação da variância do estimador LASSO . . . . .          | 63  |
| 4.5   | Erro de predição e estimação do parâmetro $t$ . . . . .      | 66  |
| 4.6   | O caso $p > n$ . . . . .                                     | 72  |
| 4.7   | GARROTE . . . . .  | 73  |
| 5     | O método de regressão <i>Elastic Net</i> . . . . .           | 77  |
| 5.1   | O estimador <i>elastic net</i> . . . . .                     | 77  |
| 6     | Regressão de Ângulos Mínimos . . . . .                       | 88  |
| 6.1   | Forward Stagewise Regression . . . . .                       | 89  |
| 6.2   | O Algoritmo LARS . . . . .                                   | 93  |
| 6.3   | O algoritmo LARS aplicado à delineamentos ortogonais . . .   | 99  |
| 6.4   | Relação entre LARS e OLS . . . . .                           | 102 |
| 6.5   | Relação entre o algoritmo LARS e o método LASSO . . . . .    | 103 |
| 7     | Estimação do Erro de Predição . . . . .                      | 109 |
| 7.1   | Erro de predição para modelos lineares . . . . .             | 109 |
| 7.2   | Estimação do erro de predição para o caso geral . . . . .    | 111 |
| 7.3   | Seleção de modelos . . . . .                                 | 123 |
| 7.3.1 | A Estatística $C_p$ de Mallows . . . . .                     | 123 |
| 7.4   | O conceito de grau de liberdade . . . . .                    | 126 |
| 7.5   | Grau de liberdade na estimação por encolhimento . . . . .    | 129 |
| 7.6   | Graus de liberdade e estimativas $C_p$ . . . . .             | 131 |
| 7.7   | Grau de liberdade do estimador Ridge . . . . .               | 133 |
| 7.8   | O grau de liberdade do estimador LASSO . . . . .             | 134 |
| 8     | Aspectos computacionais e aplicações . . . . .               | 136 |
| 8.1   | Simulações no R . . . . .                                    | 140 |
| 8.2   | Aplicação em dados genéticos . . . . .                       | 146 |
| 9     | Resultados e discussão . . . . .                             | 149 |

|           |                              |            |
|-----------|------------------------------|------------|
| <b>10</b> | <b>CONCLUSÃO . . . . .</b>   | <b>155</b> |
|           | <b>REFERÊNCIAS . . . . .</b> | <b>156</b> |
|           | <b>APÊNDICE . . . . .</b>    | <b>159</b> |

## 1 INTRODUÇÃO

Com o advento de um número cada vez maior de dados é colocado o problema de se obter novos e mais eficientes métodos estatísticos. Em particular, para a teoria da regressão linear o número elevado de covariáveis demanda, além de um robusto método de estimação dos parâmetros, também um método de seleção de covariáveis. O método de estimação utilizando quadrados mínimos apresenta problemas quando se tem, por exemplo, quasi-colinearidade, isto é, presença de covariáveis altamente correlacionadas. Os dois métodos clássicos mais utilizados para se sanar tal deficiência são a estimação Ridge e o método denominado Subset Selection. No entanto, ambos os métodos apresentam problemas. O método Subset Selection é excessivamente variável, pois é um processo discreto. A regressão Ridge é muito mais estável, mas raramente gera estimativas nulas para os parâmetros. Em razão destes fatos, foi proposto Tibshirani (1996) um novo método denominado LASSO (Least Absolute Shrinkage and Selection Operator). A principal característica deste método é que no processo de encolhimento, a estimativa de muitas covariáveis se anulam, e portanto o método é também um método automático de seleção de covariáveis. Definida em termos de quadrados mínimos penalizados, apresenta um desenvolvimento analítico bastante complexo. Um dos objetivos deste trabalho é apresentar a teoria do método de forma mais geométrica possível. Além disso, várias detalhadas demonstrações, não apresentadas na literatura, são desenvolvidas no sentido de tornar o texto uma referência para os artigos básicos da teoria.

O método LASSO também apresenta problemas, por exemplo, no caso em que se tem mais covariáveis do que observações, o método seleciona no máximo o número de covariáveis igual ao número de observações. Tem-se também que se duas covariáveis são altamente correlacionadas, o método LASSO tem tendência

a selecionar apenas uma delas. No sentido de se obter um método que possua as propriedades do LASSO, mas que consiga minimizar suas deficiências, um novo método denominado Elastic Net é proposto (ZOU; HASTIE, 2005). É também um método baseado na penalização da soma de quadrados. Esta penalização é obtida como uma combinação ponderada entre as penalizações do LASSO e do método Ridge. Este novo método, além das propriedades de seleção de covariáveis, tem a vantagem de fazer seleção por grupo, isto é, possui a tendência de que covariáveis altamente correlacionadas serem simultaneamente selecionadas. Ambos os métodos Lasso e Elastic Net não possuem forma explícita para se determinar suas estimativas. Estas são obtidas através de procedimentos algorítmicos. Nesse sentido foi feito também um estudo bastante aprofundado do algoritmo LARS (Least Angle Regressio) proposto por Efron et al. (2004). Este algoritmo de seleção de modelos é uma evolução dos algoritmos clássicos Forward Stepwise e Forward Stagewise. Possui importância em si mesmo e com pequenas alterações, calcula tanto a estimativa LASSO quanto a estimativa Elastic Net, sendo a base fundamental do pacote *glmnet* na linguagem R.

Sobre a notação utilizada, matrizes serão denotadas por letras maiúsculas ( $A, B, X$ ), vetores coluna por letras minúsculas em negrito ( $\mathbf{a}, \mathbf{b}, \mathbf{v}, \boldsymbol{\mu}$ ), escalares por letras minúsculas em itálico ( $a, b, \gamma$ ) e a transposta de uma matriz representada por um apóstrofo junto a esta, isto é,  $X^t = X'$ . Será utilizada também a notação  $\mathbf{Y} = (Y_1, \dots, Y_n)'$  para representação de vetores aleatórios e  $\mathbf{y} = (y_1, \dots, y_n)'$  para valores observados do vetor aleatório  $\mathbf{Y}$ . Para ser coerente com a notação de alguns artigos, o vetor  $\mathbf{y}$  poderá significar também um vetor aleatório. A imagem de uma transformação linear  $X$  será representada por  $\text{Im}(X)$ , a dimensão do subespaço gerado por uma transformação linear  $X$  é denotado por  $\text{Dim}(X)$ , o núcleo (ou *Kernel*) de uma transformação é denotado por  $\text{Ker}(X)$  e o traço de uma matriz



quadrada  $X$  é representado por  $\text{tr}(X)$ .

Geometricamente, um vetor linha ou coluna com  $p$  elementos pode ser associado com um ponto num espaço  $p$ -dimensional. Os elementos no vetor são as coordenadas do ponto. Nesse sentido, um conjunto de  $n$  observações pode ser geometricamente interpretado como um vetor  $\mathbf{y} = (y_1, \dots, y_n)'$ , de dimensão  $n \times 1$ , em que cada coordenada deste corresponde a uma observação. A Figura 1 ilustra este tipo de representação.

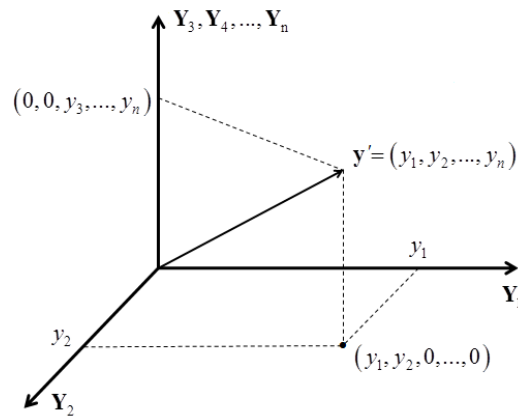


Figura 1 Representação geométrica de um vetor  $n$ -dimensional.

Com esta linguagem geométrica, o presente trabalho pretende apresentar em detalhes a teoria dos estimadores LASSO e Elastic Net. Algumas simulações são desenvolvidas para se estudar o comportamento destes estimadores na presença de vários níveis de correlação entre as covariáveis. Uma aplicação em seleção genômica é desenvolvida e analisada.

## 2 Uma abordagem geométrica à regressão linear múltipla

A regressão linear múltipla admite uma abordagem geométrica que permite um tratamento similar para três teorias importantes: a regressão *ridge*, o LASSO e o *Elasticnet*. Tal abordagem não é usual, sendo que cada uma das teorias anteriormente citadas aparentam demandar resultados teóricos diferentes.

Primeiramente será analisado o caso de regressão linear simples.

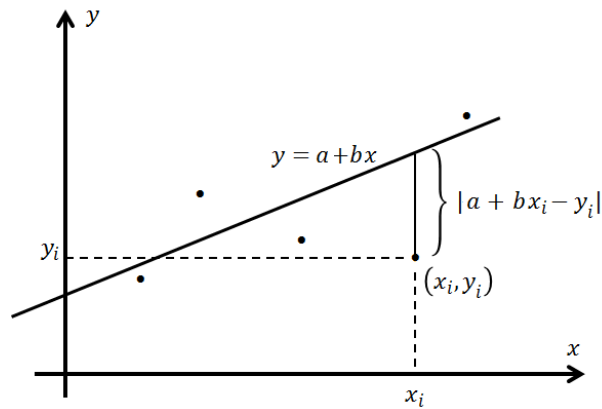


Figura 2 Reta ajustada por quadrados mínimos.

Sejam  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  os valores observados. Deseja-se então ajustar uma reta  $\hat{y} = bx + a$  de tal forma que a equação (2.1) seja minimizada.

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (a + b x_i))^2 \quad (2.1)$$

A Figura 2 ilustra tal contexto.

O método usual para se obter  $a$  e  $b$  é derivar (2.1) em relação a estes, igualar as equações derivadas a 0 e resolver o sistema. Uma abordagem geométrica é obtida da seguinte forma:

Se  $\mathbf{y}_{n \times 1} = (y_1, \dots, y_n)'$ ,  $\boldsymbol{\beta}_{2 \times 1} = (a, b)'$ ,  $\boldsymbol{\varepsilon}_{n \times 1} = (\varepsilon_1, \dots, \varepsilon_n)'$  e

$$\mathbf{X}_{n \times 2} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \text{ então}$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \Rightarrow \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

isto é,  $\boxed{y_i = a + bx_i + \varepsilon_i}$ ,  $i = 1, \dots, n$ . Portanto, reescrevendo (2.1) tem-se

$$\sum_{i=1}^n (y_i - (a + bx_i))^2 = \|\mathbf{y}_{n \times 1} - \mathbf{X}_{n \times 2}\boldsymbol{\beta}_{2 \times 1}\|^2.$$

Em termos de transformações lineares em que a matriz  $\mathbf{X}_{n \times 2}$  atua como uma transformação do espaço de parâmetros para o espaço dos dados, o conceito pode ser representado pela Figura 3.

Assim,

$$\sum_{i=1}^n (y_i - (a + bx_i))^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2.$$

Logo, minimizar  $\sum_{i=1}^n (y_i - (a + bx_i))^2$  equivale a minimizar  $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$  que, claramente, ocorre quando  $\mathbf{X}\boldsymbol{\beta}$  é a projeção ortogonal do vetor de dados  $\mathbf{y}$  no subespaço  $\text{Im}(\mathbf{X})$ . Este vetor projeção será denominado  $P_{\text{Im}(\mathbf{X})}\mathbf{y}$ , também denotado por  $\hat{\mathbf{y}}$ , vetor de dados ajustados. Note que, como o posto (coluna) de  $\mathbf{X}$  é 2, a dimensão da imagem de  $\mathbf{X}$  é  $\text{Dim}(\text{Im}(\mathbf{X})) = 2$ .

Na regressão linear múltipla, se tem várias covariáveis, isto é,  $(x_1, \dots, x_p)$ , que foram observadas  $n$  vezes. Deseja-se então ajustar o modelo

$$\boxed{y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i} \quad i = 1, \dots, n$$

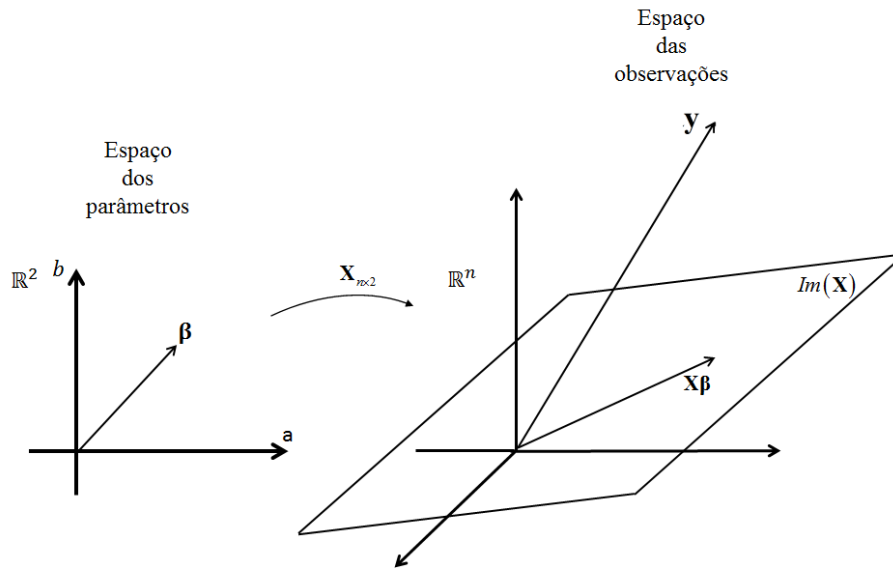


Figura 3 Geometria da regressão linear simples.

o que matricialmente equivale a

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}_{n \times 1} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix}_{n \times p+1} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}_{p+1 \times 1} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}_{n \times 1} \quad (2.2)$$

e se quer minimizar  $\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ . Observe então que, geometricamente, o modelo para o vetor de médias do vetor  $\mathbf{y}$  é simplesmente a definição do subespaço vetorial  $\text{Im}(\mathbf{X})$ . O valor mínimo é obtido pela projeção ortogonal do vetor de dados  $\mathbf{y}$  no subespaço  $\text{Im}(\mathbf{X})$ , que possui nesse caso dimensão  $p + 1$ . A Figura 4 ilustra esta situação.

Geometricamente a representação é similar a Figura 3, diferindo apenas

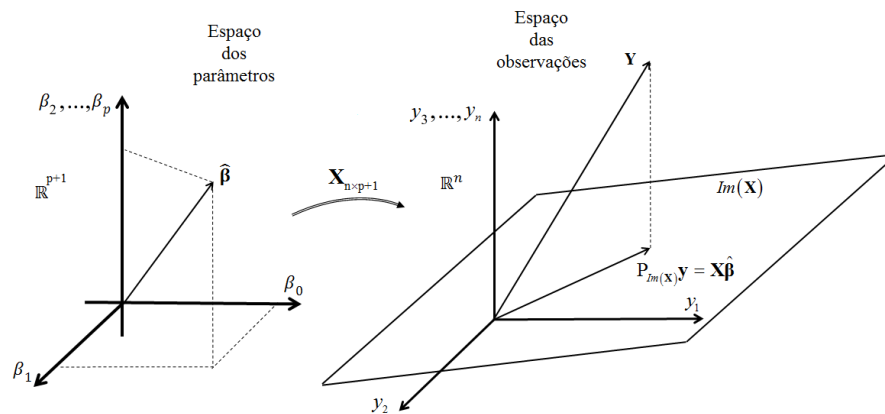


Figura 4 Geometria da regressão linear múltipla.

na dimensão do espaço dos parâmetros que agora passa a ser o  $\mathbb{R}^p$  (ver Figura 4). Considerando o posto de  $X$  igual a  $p$  (posto completo), tem-se que  $\text{Dim}(\text{Im}(X)) = p$ . Para se obter uma expressão para  $P_{\text{Im}(X)}\mathbf{y}$  é necessário a teoria das matrizes de projeção.

**Definição 1.** Uma matriz quadrada é dita matriz de projeção ou projetor se é idempotente, isto é,  $A^2 = A$ .

Pela definição 1, tem-se que um projetor  $A$  restrito a  $\text{Im}(A)$  é a identidade, ou seja,

$$\begin{aligned} A(A\mathbf{z}) &= A^2\mathbf{z} \\ &= A\mathbf{z}. \end{aligned}$$

Observe que, se  $I$  é a matriz identidade,  $I - A$  também é um projetor pois

$$\begin{aligned} (I - A)^2 &= I - 2A + A^2 \\ &= I - 2A + A \\ &= I - A. \end{aligned}$$

Como

$$\begin{aligned}
 A((I - A)\mathbf{z}) &= A(\mathbf{z} - A\mathbf{z}) \\
 &= A\mathbf{z} - A^2\mathbf{z} \\
 &= A\mathbf{z} - A\mathbf{z} \\
 &= \mathbf{0},
 \end{aligned}$$

segue que  $\text{Im}(I - A) \subset \text{Ker}(A)$ . Além do mais, se  $\mathbf{z} \in \text{Ker}(A)$  então

$$\begin{aligned}
 (I - A)\mathbf{z} &= \mathbf{z} - A\mathbf{z} \\
 &= \mathbf{z} - \mathbf{0}.
 \end{aligned}$$

Logo,  $\mathbf{z} \in \text{Im}(I - A)$  e, portanto,  $\text{Im}(I - A) = \text{Ker}(A)$ . Assim, segue então que  $\text{Im}(A) = \text{Ker}(I - A)$ .

**Definição 2.** Uma matriz de projeção  $A$  é dita um projetor ortogonal se, para um dado vetor  $\mathbf{w}$ ,  $A\mathbf{w} - \mathbf{w}$  é perpendicular ao subespaço  $\text{Im}(A)$ .

**Proposição 1.** Uma matriz de projeção é simétrica se, e somente se, é um projetor ortogonal.

*Demonstração.* Suponha  $A$  simétrica e  $\mathbf{v}$  e  $\mathbf{w}$  vetores quaisquer. Utilizando as propriedades de produto interno e sendo  $A^2 = A$ , tem-se que

$$\begin{aligned}
 \langle \mathbf{v} - A\mathbf{v}, A\mathbf{w} \rangle &= \langle \mathbf{v}, A\mathbf{w} \rangle - \langle A\mathbf{v}, A\mathbf{w} \rangle \\
 &= \langle \mathbf{v}, A\mathbf{w} \rangle - \langle \mathbf{v}, A'A\mathbf{w} \rangle \\
 &= \langle \mathbf{v}, A\mathbf{w} \rangle - \langle \mathbf{v}, AA\mathbf{w} \rangle \\
 &= \langle \mathbf{v}, A\mathbf{w} \rangle - \langle \mathbf{v}, A^2\mathbf{w} \rangle \\
 &= \langle \mathbf{v}, A\mathbf{w} \rangle - \langle \mathbf{v}, A\mathbf{w} \rangle
 \end{aligned}$$

$$= \mathbf{0}, \forall \mathbf{v}, \mathbf{w}.$$

Portanto,  $\mathbf{v} - A\mathbf{v}$  é perpendicular a  $\text{Im}(A)$ . Seja agora  $A\mathbf{v} - \mathbf{v}$  perpendicular a  $\text{Im}(A)$ . Tem-se então que, para  $\mathbf{w}$  qualquer,

$$\begin{aligned} 0 &= \langle A\mathbf{v} - \mathbf{v}, A\mathbf{w} \rangle \\ &= \langle A\mathbf{v}, A\mathbf{w} \rangle - \langle \mathbf{v}, A\mathbf{w} \rangle, \end{aligned}$$

ou seja,

$$\langle A\mathbf{v}, A\mathbf{w} \rangle = \langle \mathbf{v}, A\mathbf{w} \rangle. \quad (2.3)$$

Também, para  $(A\mathbf{w}) - \mathbf{w}$  perpendicular a  $A\mathbf{v}$ ,

$$\begin{aligned} 0 &= \langle A\mathbf{w} - \mathbf{w}, A\mathbf{v} \rangle \\ &= \langle A\mathbf{w}, A\mathbf{v} \rangle - \langle \mathbf{w}, A\mathbf{v} \rangle, \end{aligned}$$

ou seja,  $\langle A\mathbf{w}, A\mathbf{v} \rangle = \langle \mathbf{w}, A\mathbf{v} \rangle$  e portanto

$$\langle A\mathbf{v}, A\mathbf{w} \rangle = \langle A\mathbf{v}, \mathbf{w} \rangle. \quad (2.4)$$

Consequentemente, das equações (2.3) e (2.4)

$$\langle A\mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{v}, A\mathbf{w} \rangle.$$

□

Uma interpretação geométrica relativa à matriz transposta  $A'$  (operador adjunto) pode ser visualizada na Figura 5.

Para se obter uma expressão para a projeção  $P_{\text{Im}(X)}\mathbf{y}$ , será apresentada na Proposição 2 uma dedução baseada apenas em argumentos geométricos.

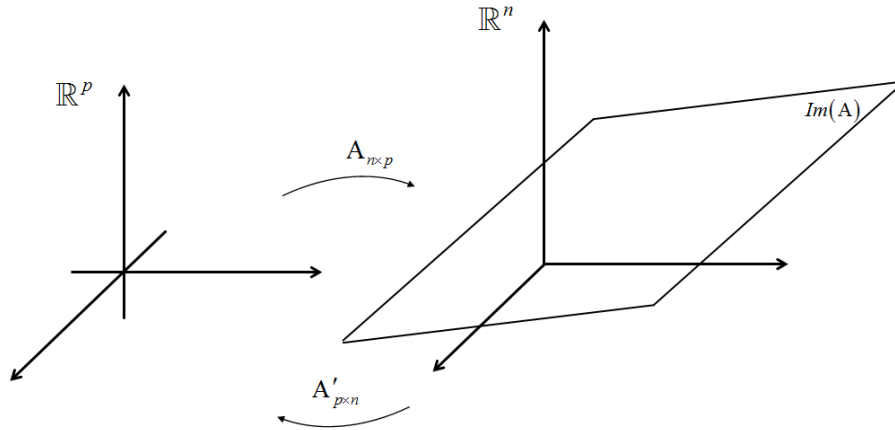


Figura 5 Interpretação geométrica das matrizes  $A$  e  $A'$

**Proposição 2.** A projeção ortogonal do vetor de dados  $\mathbf{y}$  no subespaço  $\text{Im}(X)$  é

$$P_{\text{Im}(X)}\mathbf{y} = X(X'X)^{-1}X'\mathbf{y}.$$

*Demonstração.* Seja  $\mathbf{y} = (y_1, \dots, y_n)$ . Em razão dos erros, tem-se que a probabilidade  $P[\mathbf{y} \in \text{Im}(X)] = 0$ . Como queremos minimizar  $L(\boldsymbol{\beta}) = \|\mathbf{y} - X\boldsymbol{\beta}\|^2$  e o vetor  $X\boldsymbol{\beta}$  pertence à imagem de  $X$ , segue que  $L(\boldsymbol{\beta})$  é mínimo quando  $X\boldsymbol{\beta}$  é a projeção ortogonal de  $\mathbf{y}$  na  $\text{Im}(X)$ , com notação  $P_{\text{Im}(X)}(\mathbf{y})$ , e, conseqüentemente, existe um vetor  $\hat{\boldsymbol{\beta}}$  de parâmetros tal que  $X\hat{\boldsymbol{\beta}} = P_{\text{Im}(X)}(\mathbf{y})$  já que a matriz  $X$  é de posto completo. Considere ainda o vetor de parâmetros  $\tilde{\boldsymbol{\beta}} = X'P_{\text{Im}(X)}(\mathbf{y})$ . Observe que, como  $\mathbf{y} - P_{\text{Im}(X)}(\mathbf{y})$  é ortogonal a qualquer vetor na imagem de  $X$ ,  $\mathbf{y} - P_{\text{Im}(X)}(\mathbf{y})$  está no espaço nulo de  $X$ , isto é,  $\text{Ker}(X)$ . Segue então que  $X'(\mathbf{y} - P_{\text{Im}(X)}(\mathbf{y})) = \mathbf{0}$ . Portanto, pode-se escrever que

$$\begin{aligned} \tilde{\boldsymbol{\beta}} &= X'P_{\text{Im}(X)}(\mathbf{y}) \\ &= X'P_{\text{Im}(X)}(\mathbf{y}) + X'(\mathbf{y} - P_{\text{Im}(X)}(\mathbf{y})) \\ &= X'\mathbf{y}. \end{aligned} \tag{2.5}$$



Como  $X\hat{\beta} = P_{\text{Im}(X)}(\mathbf{y})$ , pode-se também escrever que

$$\beta = X'X\hat{\beta}. \quad (2.6)$$

Assim, de (2.7) e (2.6) segue que

$$\begin{aligned} \beta &= X'X\hat{\beta} \\ &= X'\mathbf{y} \\ \hat{\beta} &= (X'X)^{-1}X'\mathbf{y}, \end{aligned}$$

conforme representado na Figura 4. □

Segue da Proposição 2 que a estimativa do vetor de parâmetros  $\beta$ , que nos dá o ajuste de quadrados mínimos, é

$$\hat{\beta} = (X'X)^{-1}X'\mathbf{y},$$

pois

$$\begin{aligned} X\hat{\beta} &= X(X'X)^{-1}X'\mathbf{y} \\ &= P_{\text{Im}(X)}\mathbf{y}. \end{aligned}$$

O modelo ajustado fica então da forma

$$\hat{\mathbf{y}} = \hat{\beta}_0 + \hat{\beta}_1\mathbf{x}_1 + \hat{\beta}_2\mathbf{x}_2 + \cdots + \hat{\beta}_p\mathbf{x}_p.$$

## 2.1 O triângulo fundamental

A geometria dos modelos de regressão linear múltipla é definida por um triângulo retângulo conforme destacado na Figura 6, em que SST, SSR e SSE são, respectivamente, abreviações do inglês para *Total Sum of Square*, *Regression Sum of Square* e *Explained Sum of Square*, conforme Rencher (2008). Portanto, segue

que

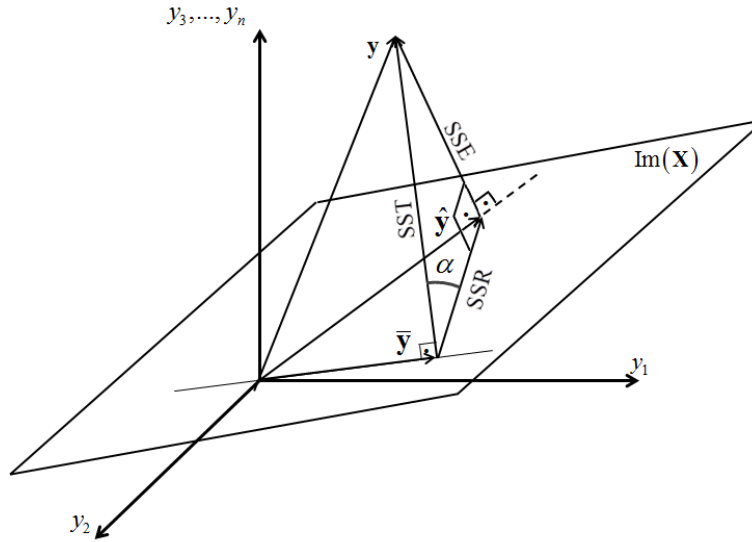


Figura 6 O triângulo fundamental.

$$(\mathbf{y} - \hat{\mathbf{y}}) + (\hat{\mathbf{y}} - \bar{\mathbf{y}}) = \mathbf{y} - \bar{\mathbf{y}}$$

$$(\mathbf{y} - \hat{\mathbf{y}})'(\hat{\mathbf{y}} - \bar{\mathbf{y}}) = 0$$

$$\|\mathbf{y} - \hat{\mathbf{y}}\|^2 + \|\hat{\mathbf{y}} - \bar{\mathbf{y}}\|^2 = \|\mathbf{y} - \bar{\mathbf{y}}\|^2,$$

e então

$$SST = \|\mathbf{y} - \bar{\mathbf{y}}\|^2 = \sum_{i=1}^n (y_i - \bar{y})^2 = (\text{hipotenusa})^2$$

$$SSR = \|\hat{\mathbf{y}} - \bar{\mathbf{y}}\|^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = (\text{cateto adjacente})^2$$

$$SSE = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = (\text{cateto oposto})^2.$$

Logo,

$$SST = SSR + SSE.$$

As somas de quadrados podem ser descritas em termos das projeções em seus respectivos subespaços, conforme se segue:

$$\begin{aligned}
\text{SSR} &= \|\hat{\mathbf{y}} - \bar{\mathbf{y}}\|^2 \\
&= (\hat{\mathbf{y}} - \bar{\mathbf{y}})' (\hat{\mathbf{y}} - \bar{\mathbf{y}}) \\
&= \left( \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - \frac{1}{n}\mathbf{J}\mathbf{y} \right)' \left( \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - \frac{1}{n}\mathbf{J}\mathbf{y} \right) \\
&= \mathbf{y}' \left( \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \frac{1}{n}\mathbf{J} \right)' \left( \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \frac{1}{n}\mathbf{J} \right) \mathbf{y} \\
&= \mathbf{y}' \left( \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \frac{1}{n}\mathbf{J} \right)^2 \mathbf{y} \\
&= \mathbf{y}' \left( \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\frac{1}{n}\mathbf{J} - \frac{1}{n}\mathbf{J}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' + \left(\frac{1}{n}\mathbf{J}\right)^2 \right) \mathbf{y}.
\end{aligned}$$

Assume-se aqui que  $\vec{\mathbf{1}} \in \text{Im}(\mathbf{X})$ . Esta é a hipótese do modelo linear com intercepto, a qual é mais utilizada nos livros texto. Neste caso, como  $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  é a projeção em  $\text{Im}(\mathbf{X})$  e  $\frac{1}{n}\mathbf{J}$  a projeção no subespaço gerado por  $\vec{\mathbf{1}}$ , então  $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\frac{1}{n}\mathbf{J}\mathbf{y} = \frac{1}{n}\mathbf{J}\mathbf{y}$  para todo  $\mathbf{y}$ , i.e.,  $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\frac{1}{n}\mathbf{J} = \frac{1}{n}\mathbf{J}$ . Em contra partida,  $\frac{1}{n}\mathbf{J}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \frac{1}{n}\mathbf{J}\hat{\mathbf{y}} = \bar{\mathbf{y}}$ . Tem-se então que

**Proposição 3.** Se  $\vec{\mathbf{1}} = (1, 1, \dots, 1)' \in \text{Im}(\mathbf{X})$  então a média do vetor ajustado é igual à média do vetor observado, i.e.,  $\bar{\hat{\mathbf{y}}} = \bar{\mathbf{y}}$ .

*Demonstração.* Considere a seguinte decomposição em soma direta ortogonal (Figura 7). Se  $V_{\vec{\mathbf{1}}}$  é o subespaço gerado pelo vetor  $\vec{\mathbf{1}}$ ,  $V_{\vec{\mathbf{1}}} = \{(a, a, \dots, a), a \in \mathbb{R}\}$ , então

$$\mathbb{R}^n = V_{\vec{\mathbf{1}}} \oplus \left\{ V_{\vec{\mathbf{1}}}^\perp \cap \text{Im}(\mathbf{X}) \right\} \oplus (\text{Im}(\mathbf{X}))^\perp.$$

Assim, tem-se que  $\mathbf{v} = \mathbf{v}_1 + \mathbf{v}_2 + \mathbf{v}_3$  com  $\mathbf{v}_1$ ,  $\mathbf{v}_2$  e  $\mathbf{v}_3$  pertencendo respectivamente a  $V_{\vec{\mathbf{1}}}$ ,  $(V_{\vec{\mathbf{1}}}^\perp \cap \text{Im}(\mathbf{X}))$  e  $\text{Im}(\mathbf{X})^\perp$ . Também,  $P_{V_{\vec{\mathbf{1}}}}(\mathbf{v}) = \mathbf{v}_1$  e

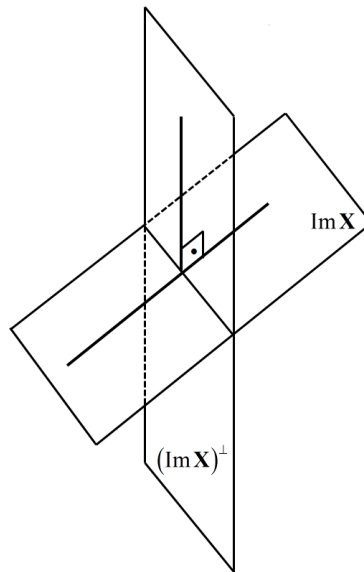


Figura 7 Projeções em termos de somas direta.

$P_{\text{Im}(X)}(\mathbf{v}) = \mathbf{v}_1 + \mathbf{v}_2$ . Portanto,

$$P_{V_1}(P_{\text{Im}(X)}(\mathbf{v})) = P_{V_1}(\mathbf{v}).$$

Se  $\mathbf{y}$  é o vetor observado,  $P_{\text{Im}(X)}(\mathbf{y}) = \hat{\mathbf{y}}$  e  $P_{V_1}(\mathbf{y}) = \bar{\mathbf{y}}$ . Desta forma,

$$\bar{\mathbf{y}} = P_{V_1}(\mathbf{y}) = P_{V_1}(P_{\text{Im}(X)}(\mathbf{y})) = P_{V_1}(\hat{\mathbf{y}}) = \widehat{\bar{\mathbf{y}}}.$$

□

Uma outra demonstração mais geométrica pode ser baseada nos dois triângulos representados na Figura 8, conforme se segue.

Suponha que a projeção ortogonal de  $\hat{\mathbf{y}}$  na direção de  $\vec{\mathbf{1}}$  não seja  $\bar{\mathbf{y}}$ . Tem-se então os dois triângulos retângulos  $c^2 = b^2 + a^2$  e  $e^2 = d^2 + a^2$ . Como  $c < e$  e  $d < b$  uma vez que as projeções são ortogonais, tem-se  $c^2 < e^2$  e  $d^2 < b^2$ . Então,  $b^2 + a^2 < e^2 = a^2 + d^2 \Rightarrow b^2 < d^2$ , o que é uma contradição.

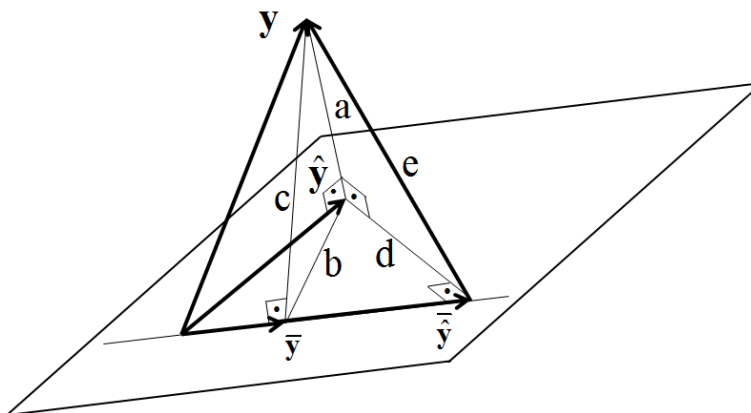


Figura 8 Demonstração geométrica da Proposição 3

Segue agora um exemplo em que  $\bar{y} \neq \hat{\bar{y}}$ .

Para  $X = \begin{pmatrix} 1 & 1 \\ 1 & -2 \\ -2 & 1 \end{pmatrix}$ ,  $y = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$ , se tem  $\vec{1} \notin \text{Im}(X)$ .  $\hat{y} =$

$$\begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}, \text{ e então } \hat{\bar{y}} = 0 \neq 2 = \bar{y}.$$

Pela Proposição 3,  $\frac{1}{n}JX(X'X)^{-1}X'y = \frac{1}{n}J\hat{y} = \hat{\bar{y}} = \bar{y} = \frac{1}{n}Jy$ .

Portanto, a SSR resulta em

$$\text{SSR} = y' \left( X(X'X)^{-1}X' - \frac{1}{n}J \right) y.$$

Para aqueles que preferem uma abordagem mais algébrica, um bom desafio é provar algebricamente que  $\frac{1}{n}JX(X'X)^{-1}X' = X(X'X)^{-1}X'\frac{1}{n}J = \frac{1}{n}J$ . A SSE é calculada então por

$$\begin{aligned} \text{SSE} &= \|y - \hat{y}\|^2 \\ &= (y - \hat{y})'(y - \hat{y}) \end{aligned}$$

$$\begin{aligned}
&= (\mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y})' (\mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}) \\
&= \left( (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \mathbf{y} \right)' (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \mathbf{y} \\
&= \mathbf{y}' (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')' (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \mathbf{y} \\
&= \mathbf{y}' (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')^2 \mathbf{y} \\
&= \mathbf{y}' (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \mathbf{y}.
\end{aligned}$$

Novamente, a identidade fundamental pode ser demonstrada:

$$\begin{aligned}
\text{SSE} + \text{SSR} &= \mathbf{y}' (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \mathbf{y} + \mathbf{y}' \left( \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \frac{1}{n}\mathbf{J} \right) \mathbf{y} \\
&= \mathbf{y}' \left( \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \frac{1}{n}\mathbf{J} \right) \mathbf{y} \\
&= \mathbf{y}' \left( \mathbf{I} - \frac{1}{n}\mathbf{J} \right) \mathbf{y} \\
&= \mathbf{y}' (\mathbf{y} - \bar{\mathbf{y}}) \\
&= (\mathbf{y}' - \bar{\mathbf{y}})' (\mathbf{y} - \bar{\mathbf{y}}) \\
&= \text{SST}.
\end{aligned}$$

## 2.2 Reparametrização em sistemas lineares

Em problemas de regressão linear, é usual que por exemplo se altere a unidade de medida de covariáveis. Neste caso é necessário uma reparametrização do modelo. Formalmente, a reparametrização pode ser descrita como se segue. Considerando o modelo  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  em que  $E[\boldsymbol{\varepsilon}] = 0$  e  $\text{cov}[\boldsymbol{\varepsilon}, \boldsymbol{\varepsilon}] = \sigma^2\mathbf{I}$ , o estimador de quadrados mínimos é dado por  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ . Uma reparametrização linear é aquela em que os novos parâmetros são funções lineares dos parâmetros originais. Suponha que os componentes do vetor  $\boldsymbol{\beta}$  sejam combinações lineares dos componentes de um vetor  $\boldsymbol{\gamma}_{p \times 1}$ , expressas por  $\boldsymbol{\beta} = \mathbf{A}_{p \times p}\boldsymbol{\gamma}$ . Desta forma,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{X}\mathbf{A}\boldsymbol{\gamma} + \boldsymbol{\varepsilon} = (\mathbf{X}\mathbf{A})\boldsymbol{\gamma} + \boldsymbol{\varepsilon},$$

conforme representação na Figura 9. Supondo as matrizes  $A$  e  $X$  injetivas então

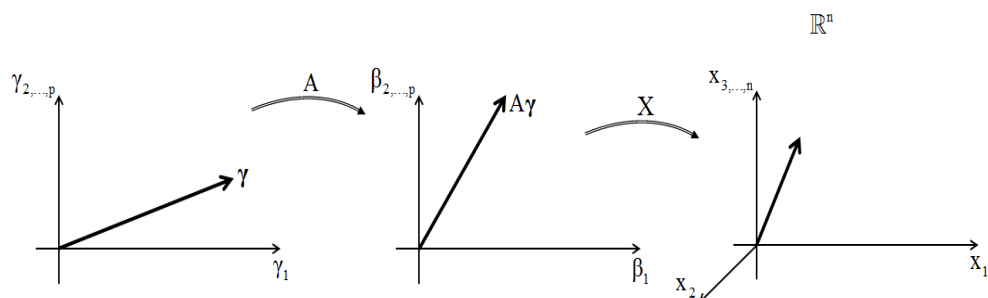


Figura 9 Reparametrização em regressão linear.

$\text{Im}(XA) = \text{Im}(X)$ . Em relação aos parâmetros dados pelas coordenadas do vetor  $\gamma$ , a nova matriz da regressão é  $XA$  e o estimador de quadrados mínimos para  $\gamma$  é dado por

$$\begin{aligned}\hat{\gamma} &= [(XA)'XA]^{-1}(XA)'\mathbf{y} \\ &= [A'X'XA]^{-1}A'X'\mathbf{y} \\ &= A^{-1}(X'X)^{-1}(A')^{-1}A'X'\mathbf{y} \\ &= A^{-1}(X'X)^{-1}X'\mathbf{y} \\ &= A^{-1}\hat{\beta}.\end{aligned}$$

Denominando as novas covariáveis pelo vetor  $\mathbf{z}$  e  $XA = Z_{n \times p}$ , então

$$\begin{bmatrix} z_{11} & z_{12} & \cdots & z_{1p} \\ z_{21} & z_{22} & \cdots & z_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ z_{n1} & z_{n2} & \cdots & z_{np} \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ a_{21} & a_{22} & \cdots & a_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{np} \end{bmatrix}.$$

Observe que os elementos da  $j$ -ésima coluna de  $Z$  são dados por

$$z_{ij} = \begin{bmatrix} x_{i1} & x_{i2} & \cdots & x_{ip} \end{bmatrix} \begin{bmatrix} a_{1j} \\ a_{2j} \\ \vdots \\ a_{pj} \end{bmatrix} = x_{i1}a_{1j} + x_{i2}a_{2j} + \cdots + x_{ip}a_{pj}$$

e toda uma coluna de  $Z$  é então

$$\begin{aligned} \begin{bmatrix} z_{1j} \\ z_{2j} \\ \vdots \\ z_{nj} \end{bmatrix} &= \begin{bmatrix} x_{11}a_{1j} + x_{12}a_{2j} + \cdots + x_{1p}a_{pj} \\ x_{21}a_{1j} + x_{22}a_{2j} + \cdots + x_{2p}a_{pj} \\ \vdots \\ x_{n1}a_{1j} + x_{n2}a_{2j} + \cdots + x_{np}a_{pj} \end{bmatrix} \\ &= \begin{bmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{n1} \end{bmatrix} a_{1j} + \begin{bmatrix} x_{12} \\ x_{22} \\ \vdots \\ x_{n2} \end{bmatrix} a_{2j} + \cdots + \begin{bmatrix} x_{1p} \\ x_{2p} \\ \vdots \\ x_{np} \end{bmatrix} a_{pj}. \end{aligned}$$

Em outras palavras, a  $j$ -ésima coluna de  $Z$  é uma combinação linear das colunas de  $X$ . As constantes da combinação linear são os elementos da  $j$ -ésima coluna de  $A$ . Logo, uma linha de  $Z$  é da forma

$$\mathbf{z}'_i = \begin{bmatrix} z_{i1} \\ z_{i2} \\ \vdots \\ z_{ip} \end{bmatrix} = \begin{bmatrix} a_{11}x_{i1} + a_{21}x_{i2} + \cdots + a_{p1}x_{ip} \\ a_{12}x_{i1} + a_{22}x_{i2} + \cdots + a_{p2}x_{ip} \\ \vdots \\ a_{1p}x_{i1} + a_{2p}x_{i2} + \cdots + a_{pp}x_{ip} \end{bmatrix} = \mathbf{A}'\mathbf{x}_i,$$

isto é, no novo modelo linear  $\mathbf{y} = (\mathbf{XA})\boldsymbol{\gamma} + \boldsymbol{\varepsilon} = \mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\varepsilon}$ , as colunas de  $Z$  são novas covariáveis, dadas por combinações lineares das covariáveis  $X$  originais. Segue então um resultado bastante importante: o valor predito não se altera quando se faz



essa mudança linear nas covariáveis, ou seja, nas covariáveis originais, a equação de predição é  $\hat{y} = \hat{\beta}'\mathbf{x}$ , enquanto que nas novas covariáveis a equação de predição é dada por

$$\begin{aligned}\hat{y} &= \hat{\gamma}'\mathbf{z} \\ &= (\mathbf{A}^{-1}\hat{\beta})'\mathbf{A}'\mathbf{x} \\ &= \hat{\beta}'(\mathbf{A}^{-1})'\mathbf{A}'\mathbf{x} \\ &= \hat{\beta}'(\mathbf{A}')^{-1}\mathbf{A}'\mathbf{x} \\ &= \hat{\beta}'\mathbf{x}.\end{aligned}$$

### 2.3 Reparametrização na presença de covariáveis categóricas

Quando se tem covariáveis qualitativas, é necessário atribuir valores arbitrários às categorias definidas. Por exemplo, para a covariável estado civil, pode-se atribuir 0 para solteiro e 1 para casado, ou -1 para solteiro e 1 para casado. Como essa arbitrariedade afetará a equação de predição? Suponha por simplicidade que a última coluna da matriz do delineamento seja definida por uma covariável binária assumindo os valores a e b. Fazendo uma reparametrização apenas desta covariável, para uma outra covariável que assumira os valores c e d, respectivamente aos valores a e b, tem-se que se obter uma matriz de forma que

$$\begin{pmatrix} 1 & x_{12} & \cdots & x_{1,p-2} & a \\ 1 & x_{22} & \cdots & x_{2,p-2} & b \\ \vdots & \vdots & \ddots & & \vdots \\ 1 & x_{n2} & \cdots & x_{n,p-2} & a \end{pmatrix}_{n \times p} = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 & x \\ 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & 0 & \cdots & 0 & y \end{pmatrix}_{p \times p}$$

$$= \begin{pmatrix} 1 & x_{12} & \cdots & x_{1p-2} & c \\ 1 & x_{22} & \cdots & x_{2p-2} & d \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{n2} & \cdots & x_{np-2} & c \end{pmatrix}_{n \times p} = \begin{pmatrix} 1 & x_{12} & \cdots & x_{1p-2} & x+ay \\ 1 & x_{22} & \cdots & x_{2p-2} & x+by \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{n2} & \cdots & x_{np-2} & x+ay \end{pmatrix}_{n \times p}.$$

Portanto, se tem o sistema  $\begin{cases} x + ya = c \\ x + yb = d \end{cases}$ , cuja solução é  $y = \frac{c-d}{a-b}$  e  $x = \frac{ad-cb}{a-b}$ .

O novo vetor de estimativa dos parâmetros, relativo às novas covariáveis, é então  $\hat{\gamma} = K^{-1}\hat{\beta}$ . De forma a se determinar a inversa  $K^{-1}$  será utilizado a fórmula Sherman-Morrison-Woodbury (RENCHER, 2008, p.23): Dada a matriz particionada  $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$ , sua inversa é dada por

$$A^{-1} = \begin{bmatrix} (A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1} & -A_{11}^{-1}A_{12}(A_{22} - A_{21}A_{11}^{-1}A_{12})^{-1} \\ -A_{22}^{-1}A_{21}(A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1} & (A_{22} - A_{21}A_{11}^{-1}A_{12})^{-1} \end{bmatrix}.$$

Para a matriz  $K$ ,  $A_{11} = I$ ,  $A_{21}$  é o vetor nulo,  $A'_{12} = \begin{bmatrix} x & 0 & \cdots & 0 \end{bmatrix}$  e  $A_{22} = y$ , e portanto

$$K^{-1} = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 & \frac{x}{y} \\ 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & 0 & \cdots & 0 & \frac{1}{y} \end{pmatrix}_{p \times p} = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 & \frac{cb-ad}{c-d} \\ 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & 0 & \cdots & 0 & \frac{a-b}{c-d} \end{pmatrix}.$$

Segue então que o vetor estimado  $\hat{\gamma}$  para o sistema com a nova covariável difere de  $\hat{\beta}$  apenas no intercepto e na componente relativa à covariável original. Neste trabalho a aplicação a dados reais utiliza uma covariável categórica codificada da

forma homozigoto recessivo (0), heterozigoto (1) e homozigoto dominante (2). Neste caso, pode-se recodificar a covariável por a e b quaisquer, porém o terceiro valor deve obedecer a relação linear expressa no sistema anterior.

### 3 Teoria geral da estimação

O processo de estimação será abordado em um contexto bastante geral e com forte ênfase na geometria, conforme em Kato (2009). Seja  $\mathbf{y}$  um vetor aleatório  $n$ -dimensional com distribuição  $n$ -variada qualquer. Considere também  $\boldsymbol{\mu}$  um vetor  $n$ -dimensional de parâmetros desta distribuição. No caso mais relevante, este último é o vetor de médias  $\boldsymbol{\mu} = E(\mathbf{y})$ . A estimação de  $\boldsymbol{\mu}$  será obtida pelo processo descrito a seguir. Seja  $K$  um subconjunto de  $\mathbb{R}^n$  onde, por hipótese, será suposto estar o vetor  $\boldsymbol{\mu}$ . Portanto, nesse sentido o subconjunto  $K$  é o modelo estatístico. Uma vez observado o vetor  $\mathbf{y}$ , a estimativa de  $\boldsymbol{\mu}$  será obtida simplesmente projetando  $\mathbf{y}$  no subconjunto  $K$ . Tal projeção será denominada  $\hat{\boldsymbol{\mu}} = \mu(\mathbf{y})$ , e portanto  $\mu : \mathbb{R}^n \mapsto K \subseteq \mathbb{R}^n$ . Este esquema, em sua generalidade, certamente apresentará problemas na obtenção das propriedades da função  $\mu$ . Serão supostas então algumas hipóteses matemáticas. Será exigido que o subconjunto  $K$  seja fechado e que a projeção seja obtida pela minimização da distância

$$\mu(\mathbf{y}) = \arg \min_{k \in K} \{d(\mathbf{y}, k), k \in K\},$$

em que a distância  $d$ , em geral, é a distância euclideana. Um dos problemas que a função  $\mu$  pode apresentar é não estar bem definida no sentido de que pode haver mais de um vetor  $k \in K$  que minimiza a distância, ou simplesmente não existir um ponto que minimize a distância. Faz-se necessário então considerar uma maior restrição do subconjunto  $K$  para que a função  $\mu$  seja bem definida.

Para que se possa garantir a existência de elementos no subconjunto  $K$  que minimizem a distância, supõe-se  $K$  um subconjunto fechado e convexo. Lembrando que  $K$  é convexo se, dados  $k_1, k_2 \in K$ , então o vetor  $t k_1 + (1 - t)k_2$ , com  $0 \leq t \leq 1$ , também pertence à  $K$ . Portanto, a menos que se diga o contrário, será

suposto  $K$  como sendo fechado e convexo.

**Proposição 4.** *Se  $K$  é convexo e fechado, então a projeção  $\mu : \mathbb{R}^n \mapsto K$ , em que  $\mu(\mathbf{y}) = \arg \min_K \{d(\mathbf{y}, \mathbf{k}), \mathbf{k} \in K\}$ , está bem definida.*

*Demonstração.* Sejam  $\mathbf{k}_1$  e  $\mathbf{k}_2$  pontos de  $K$  que minimizam a distância ao vetor  $\mathbf{y}$ . Logo,  $d(\mathbf{y}, \mathbf{k}_1) = d(\mathbf{y}, \mathbf{k}_2)$ . Observando que o triângulo ABC da Figura 10 é

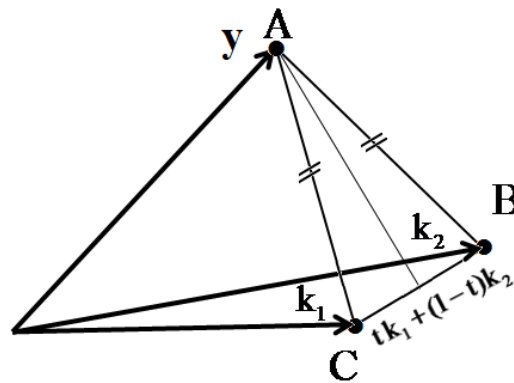


Figura 10 Demonstração da unicidade da projeção em convexos.

isósceles e portanto  $d(t\mathbf{k}_1 + (1-t)\mathbf{k}_2, \mathbf{y}) < d(\mathbf{k}_2, \mathbf{y})$ , tem-se um absurdo, pois pela convexidade de  $K$ ,  $t\mathbf{k}_1 + (1-t)\mathbf{k}_2 \in K$ .

A existência do vetor em  $K$  que minimiza a distância segue do fato de  $K$  ser fechado e não será aqui demonstrada.  $\square$

Uma observação interessante é que em algumas situações o problema de se minimizar a distância do vetor de dados  $\mathbf{y}$  ao modelo  $K$  pode ser resolvido como um problema de minimização no espaço paramétrico. Suponha o caso de regressão  $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  com  $X : \mathbb{R}^p \mapsto \mathbb{R}^n$  injetiva e que  $K \subseteq \text{Im}(X)$  (Figura 11). Para  $\mathbf{k} \in K$  e  $P_{\text{Im}(X)}(\mathbf{y})$  considere o triângulo da Figura 12. Minimizar a norma da hipotenusa  $\mathbf{y} - \mathbf{k}$  é equivalente à minimizar a norma do cateto  $(P_{\text{Im}(X)}(\mathbf{y})) - \mathbf{k}$ , pois a norma do outro cateto não depende do convexo  $K$  mas apenas da distância da projeção

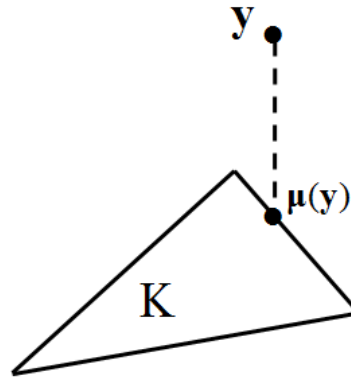


Figura 11 Projeção do vetor  $y$  no convexo  $K$ .

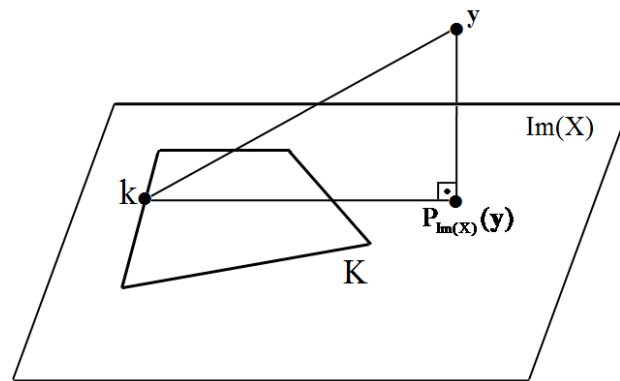


Figura 12 Relação entre  $d(y, k)$  e  $d(y, P_{\text{Im}(X)}(y))$ .

ortogonal e não do ponto  $k$ . Como por hipótese  $X : \mathbb{R}^p \mapsto \text{Im}(X) \subset \mathbb{R}^n$  é injetiva e  $K \subset \text{Im}(X)$ , tomando a pré-imagem  $K_p = X^{-1}(K)$  e observando que se o produto interno de  $\mathbb{R}^p$  é definido por  $\langle \beta_1, \beta_2 \rangle_p = \beta_1' X' X \beta_2$ , então a transformação  $X$  é uma isometria sobre a imagem, isto é,  $\langle X\beta_1, X\beta_2 \rangle = \beta_1' X' X \beta_2 = \langle \beta_1, \beta_2 \rangle_p$ . Desta forma,  $X$  preserva distâncias. Se  $x \in \text{Im}(X)$  e  $Xz = x$ , então a distância de  $z$  à  $K_p$  é igual à distância de  $x$  à  $K$ . Segue então que o problema de se encontrar  $\mu(y) \in K$  mais próximo de  $y$  é equivalente a se encontrar  $\hat{\beta}$  em  $K_p$  o mais perto

possível de  $\hat{\beta}_{\text{ols}}$  e  $X\hat{\beta} = \mu(\mathbf{y})$  (Figura 13).

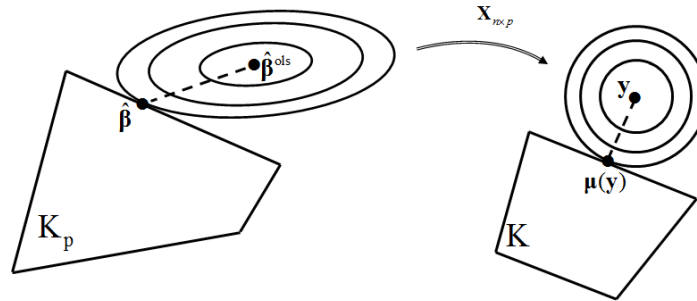


Figura 13 Isometria entre o espaço de dados e o espaço de parâmetros.

Desta forma,  $\hat{\beta}$  é a estimativa do vetor de parâmetros obtida pelo processo. As propriedades deste estimador estão claramente vinculadas às propriedades geométricas do conjunto convexo  $K_p$ .

### 3.1 Estimador de quadrados mínimos

Considere o modelo de regressão linear  $\mathbf{y} = X\beta + \varepsilon$  em que  $X$  é uma matriz definida pelos valores das covariáveis e  $\beta$  um vetor  $p$ -dimensional de parâmetros. A matriz  $X$  pode ser vista como uma transformação linear do espaço de parâmetros  $\mathbb{R}^p$  no espaço de dados  $\mathbb{R}^n$ ,  $X : \mathbb{R}^p \mapsto \mathbb{R}^n$ . Como  $\mu = E[\mathbf{y}] = E[X\beta + \varepsilon] = X\beta$ , o vetor de médias está contido no subespaço imagem de  $X$ ,  $\text{Im}(X)$ , e portanto o modelo  $K$  é justamente este subespaço. Observe que este é o caso mais simples possível uma vez que  $K$  é fechado e convexo. Mais simples ainda, se tem que a projeção  $\mu : \mathbb{R}^n \mapsto \text{Im}(X) \subseteq \mathbb{R}^n$  é a projeção linear ortogonal. Desta forma, a expressão de  $\mu(\mathbf{y})$  pode ser obtida de forma explícita, para o caso de  $X$  ser injetiva, como a fórmula matricial  $\mu(\mathbf{y}) = X(X'X)^{-1}X'\mathbf{y}$  e  $\hat{\beta}_{\text{ols}} = (X'X)^{-1}X'\mathbf{y}$  (ols - *ordinary least square*). Os estimadores  $\mu(\mathbf{y})$  e  $\hat{\beta}_{\text{ols}}$  possuem propriedades

ótimas. Este último é não viesado uma vez que

$$E[\hat{\beta}_{\text{ols}}] = E[(X'X)^{-1}X'y] = (X'X)^{-1}X'E[y] = (X'X)^{-1}X'X\beta = \beta.$$

A matriz de covariâncias de  $\hat{\beta}_{\text{ols}}$ , pela propriedade  $\text{cov}(A\mathbf{y}) = A\text{cov}(\mathbf{y})A'$ , é dada por

$$\begin{aligned} \text{cov}(\hat{\beta}_{\text{ols}}) &= \text{cov}((X'X)^{-1}X'y) \\ &= [(X'X)^{-1}X'] \text{cov}(\mathbf{y}) [(X'X)^{-1}X']' \\ &= (X'X)^{-1}X'y\mathbf{y}' [X(X'X)^{-1}] \\ &= (X'X)^{-1}X'\sigma^2IX(X'X)^{-1} \\ &= \sigma^2(X'X)^{-1}X'IX(X'X)^{-1} \\ &= \sigma^2(X'X)^{-1}. \end{aligned}$$

Nesse contexto, a variância total é

$$\text{var}_{\text{tot}}(\hat{\beta}_{\text{ols}}) = \sigma^2 \text{tr}((X'X)^{-1}) = \sigma^2 \sum_{i=1}^p \frac{1}{\lambda_i},$$

em que  $\lambda_i$  são os autovalores de  $(X'X)$ . Pelo teorema de Gauss-Markov,  $\hat{\beta}_{\text{ols}}$  possui variância total mínima dentre todos os estimadores lineares não viesados. A geometria de estimação de quadrados mínimos está descrita na Figura 14.

O estimador de quadrados mínimos obtido acima, a partir de uma consideração puramente geométrica, é obtido a partir da menor distância de um vetor  $\mathbf{y}$  a um subespaço, que é a projeção ortogonal. É possível também se deduzir analiticamente a expressão do estimador, minimizando a distância de  $\mathbf{y}$  à  $\text{Im}(X)$  por

$$\boldsymbol{\mu}(\mathbf{y}) = \arg \min \|\mathbf{y} - X\boldsymbol{\beta}\|^2$$



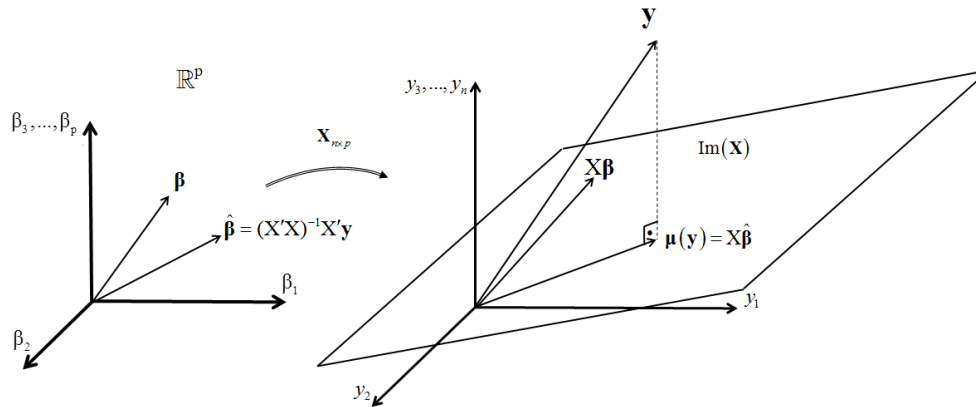


Figura 14 O Estimador de quadrados mínimos.

$$\begin{aligned}
 &= \arg \min \langle \mathbf{y} - \mathbf{X}\boldsymbol{\beta}, \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \rangle \\
 &= \arg \min \{ \langle \mathbf{y}, \mathbf{y} \rangle - 2 \langle \mathbf{y}, \mathbf{X}\boldsymbol{\beta} \rangle + \langle \mathbf{X}\boldsymbol{\beta}, \mathbf{X}\boldsymbol{\beta} \rangle \}
 \end{aligned}$$

que em termos matriciais pode ser expresso por

$$\boldsymbol{\mu}(\mathbf{y}) = \arg \min \{ \mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta} \}.$$

Derivando esta equação e igualando a zero obtém-se as equações normais  $\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}$ . No caso em que a matriz  $\mathbf{X}$  é de posto completo,  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ . É adequado, como tem-se o termo  $\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}$ , considerar no espaço paramétrico um produto interno modificado, definido por

$$\langle \boldsymbol{\beta}, \boldsymbol{\beta} \rangle_p = \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}.$$

Este produto interno introduz em  $\mathbb{R}^p$  uma geometria (métrica) modificada. Por exemplo, os vetores  $\boldsymbol{\beta}$  com  $\|\boldsymbol{\beta}\|_p = c$ , isto é, a esfera de raio  $c$  centrada na origem, é de fato um elipsoide pois

$$\|\boldsymbol{\beta}\| = c \Rightarrow \|\boldsymbol{\beta}\|^2 = c^2 \Rightarrow \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = c^2.$$

Para uma melhor descrição desse elipsoide, como  $X'X$  é simétrica e suposta de posto completo, seja  $\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$  uma base ortonormal formada de autovetores de  $X'X$ , ou seja,  $X'X\mathbf{v}_i = \lambda_i\mathbf{v}_i$ . Defina a transformação ortogonal  $A(\mathbf{e}_i) = \mathbf{v}_i$ , em que  $\{\mathbf{e}_1, \dots, \mathbf{e}_p\}$  é a base canônica de  $\mathbb{R}^p$ . Fica então definida a reparametrização  $\tilde{\boldsymbol{\beta}} = A\boldsymbol{\beta}$ , o qual implica que

$$\begin{aligned} c &= \langle \tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\beta}} \rangle_p \\ &= \tilde{\boldsymbol{\beta}}' X' X \tilde{\boldsymbol{\beta}} \\ &= (A\boldsymbol{\beta})' (X'X) A\boldsymbol{\beta} \\ &= \boldsymbol{\beta}' A' (X'X) A\boldsymbol{\beta}. \end{aligned}$$

Como

$$\begin{aligned} A' (X'X) A(\mathbf{e}_i) &= A' (X'X)(\mathbf{v}_i) \\ &= A' (\lambda_i \mathbf{v}_i) \\ &= \lambda_i A' (\mathbf{v}_i) \\ &= \lambda_i A^{-1} (\mathbf{v}_i) \\ &= \lambda_i \mathbf{e}_i \end{aligned}$$

e a matriz  $A' (X'X) A$  é da forma diagonal  $(\lambda_1, \dots, \lambda_k)$ , então com estes novos parâmetros se tem

$$\begin{aligned} \langle \tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\beta}} \rangle_p &= \langle A\boldsymbol{\beta}, A\boldsymbol{\beta} \rangle_p \\ &= \boldsymbol{\beta}' (A)' (X'X) A\boldsymbol{\beta} \\ &= \boldsymbol{\beta}' A' (X'X) A\boldsymbol{\beta} \\ &= \sum_{i=1}^k \lambda_i \beta_i^2 = c. \end{aligned}$$

Esta equação define um elipsoide com eixos nos eixos coordenados e tamanhos  $\lambda_i$ . Uma vez que  $A$  é uma transformação ortogonal e portanto uma isometria em relação à métrica usual, esta preserva a forma (Figura 15) e segue que  $c = \langle \beta, \beta \rangle_p$  é um elipsoide com eixos definidos pelos vetores  $v_i$  e tamanho  $\lambda_i$ .

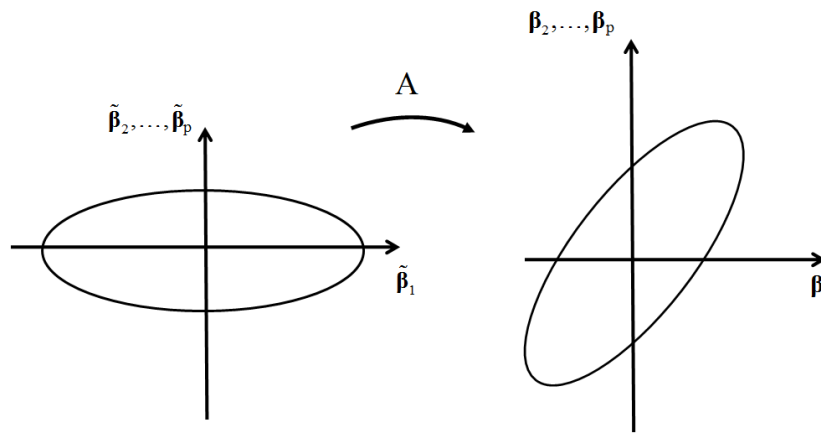


Figura 15 Transformações ortogonais preservam elipsóides.

Uma questão natural é como é a imagem deste elipsoide no subespaço  $\text{Im}(X)$ . Se  $c = \langle \beta, \beta \rangle_p$ , então sua norma em  $\mathbb{R}^n$  é

$$\|X\beta\|^2 = \langle X\beta, X\beta \rangle = \beta'X'X\beta = c,$$

em que  $\langle \cdot, \cdot \rangle$  é o produto usual no espaço  $\mathbb{R}^n$ . Portanto, a imagem do elipsoide é uma esfera  $p$ -dimensional no subespaço  $\text{Im}(X)$  (Figura 16). Prova-se, através do Teorema de Decomposição Espectral para matrizes, que também vale o resultado inverso. A imagem de uma esfera em  $\mathbb{R}^p$  é um elipsoide em  $\text{Im}(X)$  (Figura 17). Esta construção geométrica será fundamental para a abordagem que será feita para os demais estimadores tratados mais adiante neste trabalho.

Apesar do estimador de quadrados mínimos ser o estimador linear não viesado de variância mínima (BLUE), a estimação de quadrados mínimos apre-

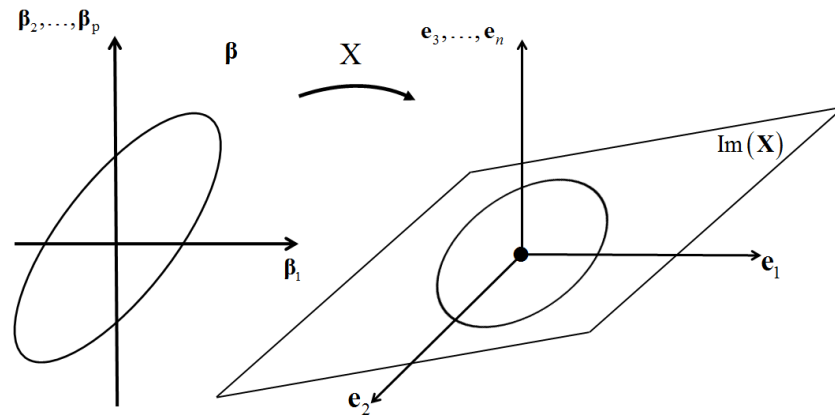


Figura 16 Transformações de elipsoides em esferas.

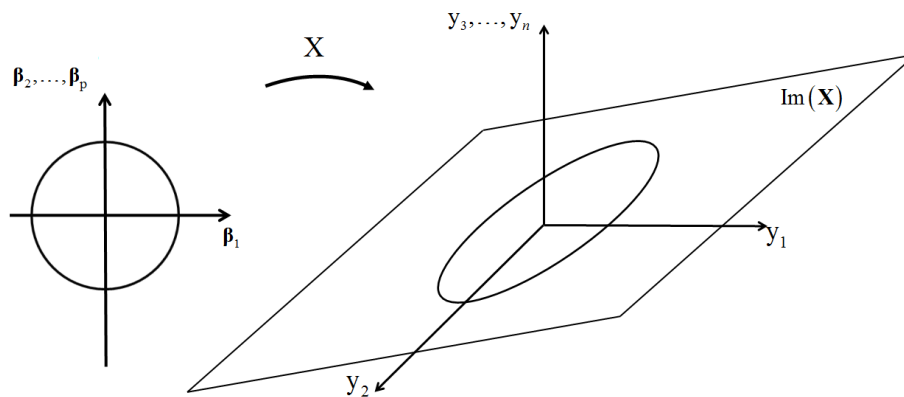


Figura 17 Transformação de esferas em elipsoides.

senta problemas em pelo menos duas situações. Uma delas é quando ocorre quasi-colinearidade, que pode ocorrer quando se tem, por exemplo, pelo menos duas covariáveis altamente correlacionados. Nesse caso, com alta probabilidade, pode-se ter que as duas colunas relativas a estas covariáveis na matriz de delineamento  $X$  sejam muito próximas. Tal fato acarreta que pelo menos um dos autovalores da matriz  $X'X$  é próximo de zero. Como a variância total do estimador de quadrados mínimos  $\hat{\beta}_{\text{ols}}$ , dada pelo traço da matriz de covariâncias,  $\text{cov}(\hat{\beta}_{\text{ols}}) = \sigma^2(X'X)^{-1}$ ,

depende dos inversos dos autovalores, esta poderá ser excessiva, inviabilizando seu uso. A outra situação ocorre quando  $p > n$ , isto é, o número de covariáveis é maior do que o número de dados, situação comum em estudos genéticos. Neste caso, a matriz  $X'X$  é singular, o que leva à necessidade da utilização de inversas generalizadas.

Qual seria uma alternativa ao método de estimação de quadrados mínimos? Ora, tal método é baseado na distância mínima. Outra possibilidade é permitir a utilização de uma distância que não seja a mínima. Esse procedimento recebe o nome genérico de método de quadrados mínimos penalizado. No sentido de explicitar as várias penalizações usualmente utilizadas, será então abordada em mais detalhes a geometria da regressão linear.

Uma vez observado o vetor de dados  $\mathbf{y}$ , considere no subespaço  $\text{Im}(X)$  todos os vetores  $X\boldsymbol{\beta}$  que estão a uma distância  $d$  do vetor  $\mathbf{y}$ . Para visualizar os vetores com essa propriedade, recorre-se ao triângulo fundamental.

$$\begin{aligned} d^2 &= \|\mathbf{y} - X\boldsymbol{\beta}\|^2 \\ &= \left\| \mathbf{y} - X\hat{\boldsymbol{\beta}}_{\text{ols}} + X\hat{\boldsymbol{\beta}}_{\text{ols}} - X\boldsymbol{\beta} \right\|^2 \\ &= \left\| \mathbf{y} - X\hat{\boldsymbol{\beta}}_{\text{ols}} \right\|^2 + \left\| X\hat{\boldsymbol{\beta}}_{\text{ols}} - X\boldsymbol{\beta} \right\|^2, \end{aligned}$$

em que a segunda igualdade se justifica pelo fato do vetor  $\mathbf{y} - X\hat{\boldsymbol{\beta}}_{\text{ols}}$  ser perpendicular ao subespaço  $\text{Im}(X)$ . Como  $\left\| \mathbf{y} - X\hat{\boldsymbol{\beta}}_{\text{ols}} \right\|$  é um valor fixo, tem-se que os vetores  $X\boldsymbol{\beta}$  satisfazem a equação  $\left\| X\hat{\boldsymbol{\beta}}_{\text{ols}} - X\boldsymbol{\beta} \right\| = r$ , com  $r$  constante, isto é, forma uma hipersfera de raio  $r$  no subespaço  $\text{Im}(X)$ , como descrito na Figura 18.

Os vetores nesta hipersfera, do ponto de vista de distância aos dados, são indistintos. Portanto, o processo de estimação, isto é, escolha de um vetor de parâmetros  $\boldsymbol{\beta}$  particular deve seguir algum tipo de restrição. Novamente, a chave

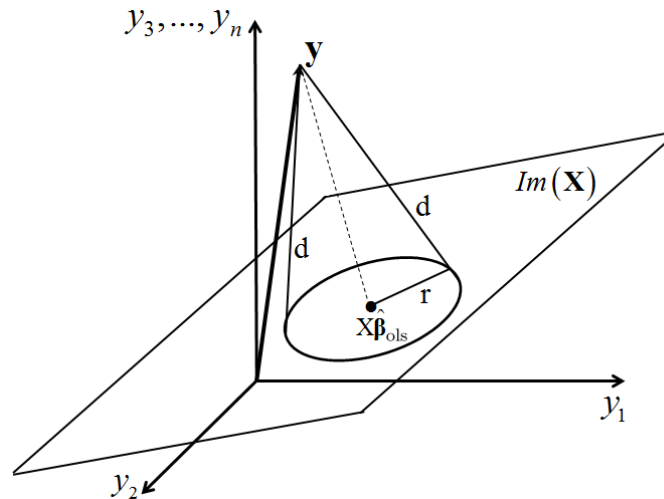


Figura 18 Vetores equidistantes à  $X\hat{\beta}_{ols}$

para se entender os processos de estimação das teorias acima citadas é a geometria. Que propriedades possuem os vetores  $\beta$  que são levados pela transformação linear  $X$  nessa hipersfera?

Observando que  $\|v\|^2 = v'v$ , tem-se

$$\begin{aligned}
 r^2 &= \|X\beta - X\hat{\beta}\|^2 \\
 &= (X\beta - X\hat{\beta})' (X\beta - X\hat{\beta}) \\
 &= (X(\beta - \hat{\beta}))' (X(\beta - \hat{\beta})) \\
 &= (\beta - \hat{\beta})' X'X (\beta - \hat{\beta})
 \end{aligned}$$

Portanto, os vetores  $\beta$  que são levados na hipersfera formam um elipsoide centrado em  $\hat{\beta}_{ols}$ , conforme representado na Figura 19.

Variando-se o raio  $r$ , varia-se de forma equivalente o elipsoide sem no entanto alterar as direções de seus eixos principais e sua excentricidade (Figura 20). Nesse sentido, o raio  $r$  pode ser considerado como um parâmetro de ajuste

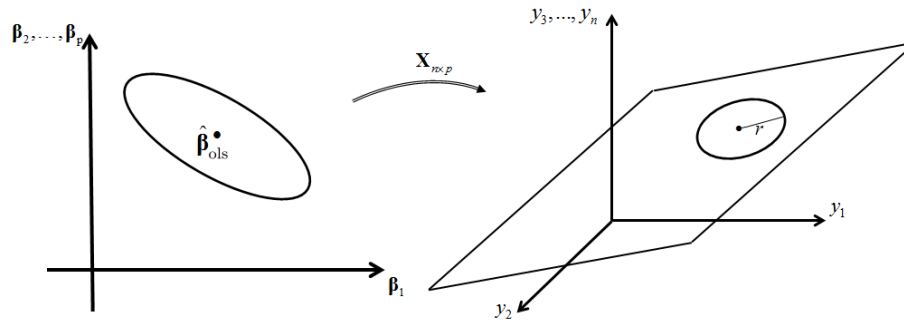


Figura 19 Relação entre a hipersfera e o elipsoide.

para o processo de estimação.

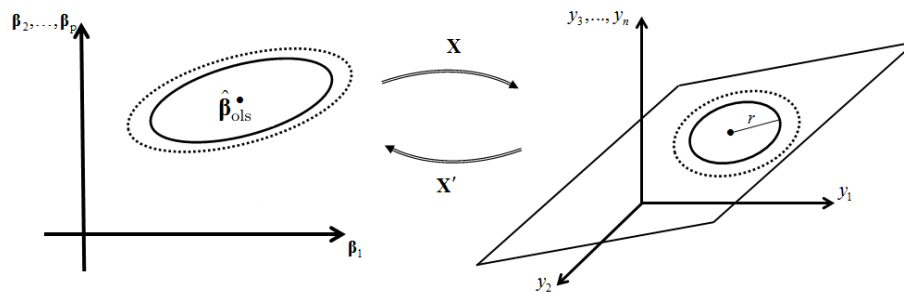


Figura 20 Variação do raio da hipersfera.

#### 4 A regressão Ridge

Um procedimento usual na estatística é a abordagem conservadora de, entre duas estimativas igualmente plausíveis, se optar pela de menor tamanho ou de menor norma. Utilizando essa filosofia, fixado o valor de  $r$ , deve-se escolher na elipse centrada em  $\hat{\beta}_{ols}$  o vetor  $\beta$  de menor norma, que será denominado  $\hat{\beta}_R(r)$  (Figura 21). Esse procedimento de estimação é denominado estimação Ridge (HÖRERL; KENNARD 1970),

Dessa forma, o estimador Ridge  $\hat{\beta}_R(r)$  é um estimador de encolhimento no sentido que  $\|\hat{\beta}_R(r)\| < \|\hat{\beta}_{OLS}\|$  e também, evidentemente, um estimador

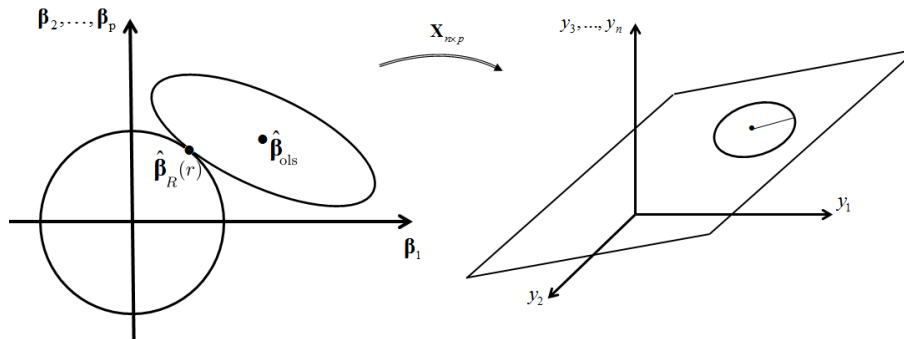


Figura 21 O estimador ridge.

viesado.

Uma das grandes vantagens do método de estimação Ridge é que, variando-se o parâmetro  $r$  obtém-se uma curva  $\hat{\beta}_R(r)$  no espaço de parâmetros. Tomando-se as coordenadas  $(\hat{\beta}_R)_i(r)$  de  $\hat{\beta}_R(r)$ , é possível uma descrição bidimensional do comportamento multidimensional do estimador, traçando-se simultaneamente os gráficos de  $(\hat{\beta}_R)_i(r)$  em função do parâmetro  $r$ . Tais gráficos são denominados *ridgetrace*, e permitem uma análise gráfica bastante útil do processo de estimação. A Figura 22 apresenta um destes *ridgetraces*.

A questão de se determinar um valor ótimo para o parâmetro  $r$ , isto é, o quanto se dever ter de encolhimento para que seja mínimo o erro quadrático médio  $E_{\beta} \left[ \left\| \hat{\beta}_R(r) - \beta \right\|^2 \right]$ , é uma das questões mais estudadas na teoria, envolvendo problemas analíticos complexos. Tal valor depende do vetor de parâmetro populacional  $\beta$ , que deve ser estimado obtendo-se estimadores do valor ótimo para  $r$  bastante complexos. Uma forma gráfica de se estimar uma região para o valor ótimo de  $r$  é o intervalo em que as curvas de *ridgetrace* se tornam aproximadamente paralelas ao eixo das abscissas.



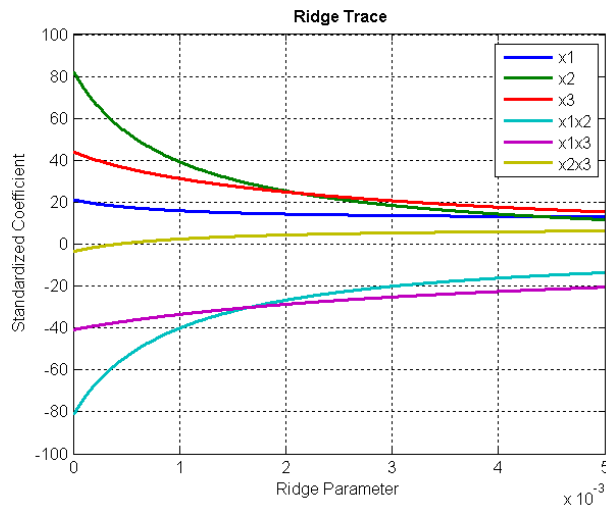


Figura 22 Exemplo de *ridgetrace* em  $\mathbb{R}^6$ .

#### 4.1 O estimador Ridge

Nesta subsecção será explicitada a obtenção do estimador  $\hat{\beta}_R(r)$ .

O estimador de mínimos quadrados, pelo teorema de Gauss-Markov, é o melhor estimador linear não viesado (BLUE). Uma pergunta natural que se pode fazer é, existe algum estimador viesado com erro quadrático menor? Neste sentido, a regressão de encolhimento trata de estimadores que, ao se permitir viés, fornecem estimativas melhores no sentido de se diminuir o erro quadrático médio. Para tratar o problema de quasi-colinearidade, Hoerl; Kennard (1970) definiram um estimador denominado Estimador Ridge utilizando uma penalização no método de quadrados mínimos. Considere o problema de regressão linear  $y = X\beta + \varepsilon$ . Novamente, o modelo  $K$  é o subespaço  $\text{Im}(X)$ . O estimador de quadrados mínimos é obtido ao se projetar ortogonalmente o vetor  $y$  em  $\text{Im}(X)$ . Introduce-se então uma penalização afirmando que as estimativas possíveis estão a uma distância  $r$  da projeção ortogonal, isto é, estão em uma esfera de raio  $r$  contida em  $\text{Im}(X)$ , esfera

esta centrada na estimativa de quadrados mínimos (Figura 23). Esta penalização gera um problema, uma vez que em razão da projeção ser ortogonal, todos os pontos desta esfera estão à mesma distância de  $\mathbf{y}$ . Como escolher então uma estimativa? A ideia é a seguinte: os vetores  $\mathbf{z}$  que pertencem a esta esfera satisfazem a condição  $\langle \mathbf{z} - P_{\text{Im}(\mathbf{X})}(\mathbf{y}), \mathbf{z} - P_{\text{Im}(\mathbf{X})}(\mathbf{y}) \rangle = r^2$ . Considere então os vetores  $\beta$  no espaço paramétrico tais que  $X\beta = \mathbf{z}$ . Logo,

$$\begin{aligned} r^2 &= \langle \mathbf{z} - P_{\text{Im}(\mathbf{X})}\mathbf{y}, \mathbf{z} - P_{\text{Im}(\mathbf{X})}\mathbf{y} \rangle \\ &= \langle X\beta - X\hat{\beta}_{ols}, X\beta - X\hat{\beta}_{ols} \rangle \\ &= \langle X(\beta - \hat{\beta}_{ols}), X(\beta - \hat{\beta}_{ols}) \rangle \\ &= [X(\beta - \hat{\beta}_{ols})]' X(\beta - \hat{\beta}_{ols}) \\ &= (\beta - \hat{\beta}_{ols})' X' X (\beta - \hat{\beta}_{ols}), \end{aligned}$$

de onde segue que a pré-imagem da esfera é um elipsoide centrado em  $\hat{\beta}_{ols}$  no espaço paramétrico  $\mathbb{R}^p$ .

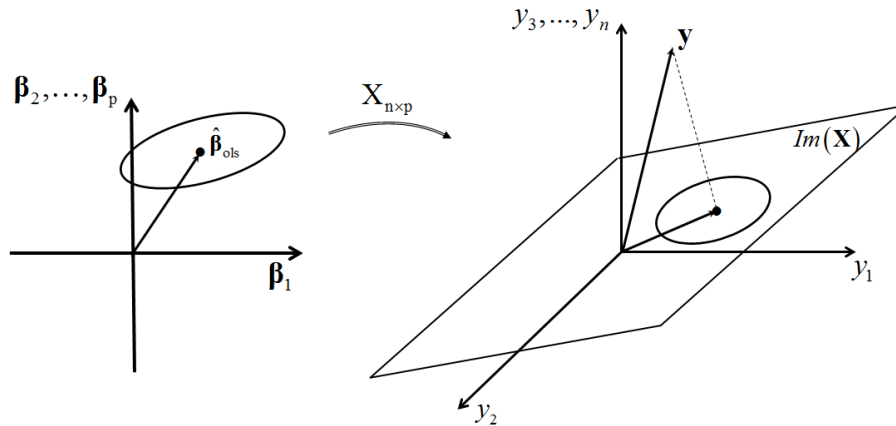


Figura 23 Penalização na definição do estimador Ridge.

Adota-se então uma atitude conservadora. Todos os  $\beta$  neste elipsoide são

estimativas viáveis. Opta-se então por aquela de menor norma, isto é, toma-se o  $\beta$  obtido como tangente entre o elipsoide e uma esfera centrada na origem. A dedução analítica da expressão deste estimador é como se segue, e é representada pela Figura 24.

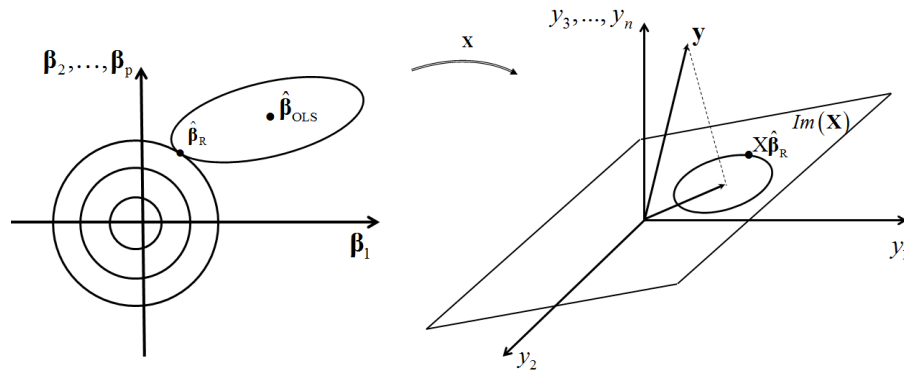


Figura 24 A geometria do estimador Ridge.

Para cada valor fixo de  $r$ , a estimativa  $\hat{\beta}_R(r)$  é, como descrito geometricamente, obtida como um problema de minimização. Analiticamente este problema pode ser descrito de duas formas equivalentes. Primeiramente no espaço de parâmetros.

Se quer minimizar a função  $\min_{\beta} \|\beta\|^2$  sujeito à restrição

$$(\beta - \hat{\beta})' X'X (\beta - \hat{\beta}) = r^2.$$

Utilizando o método dos multiplicadores de Lagrange, a Lagrangiana é

$$\begin{aligned} L(\beta, \lambda) &= \|\beta\|^2 + \lambda \left( (\beta - \hat{\beta})' X'X (\beta - \hat{\beta}) - r^2 \right) \\ &= \beta' \beta + \lambda \left( (\beta - \hat{\beta})' X'X (\beta - \hat{\beta}) - r^2 \right). \end{aligned}$$

Derivando em relação aos parâmetros e igualando à 0 tem-se

$$\frac{\partial L}{\partial \beta} = 2\beta + \lambda [2X'X (\beta - \hat{\beta})] = 0$$

e

$$\frac{\partial L}{\partial \lambda} = (\beta - \hat{\beta})' X'X (\beta - \hat{\beta}) - r^2 = 0.$$

Logo,

$$\beta + \lambda X'X\beta = \lambda X'X\hat{\beta} \Rightarrow$$

$$(I + \lambda X'X) \beta = \lambda X'X\hat{\beta}.$$

Portanto a solução é dada explicitamente por

$$\hat{\beta}_R(r) = (I + \lambda X'X)^{-1} \lambda X'X\hat{\beta} = \left( \frac{1}{\lambda} I + X'X \right)^{-1} X'X\hat{\beta}.$$

Uma vez que  $\hat{\beta} = (X'X)^{-1} X'y$ , então

$$\hat{\beta}_R(r) = \left( \frac{1}{\lambda} I + X'X \right)^{-1} X'y = (kI + X'X)^{-1} X'y,$$

sendo  $k = \frac{1}{\lambda}$ .O valor de  $k$  em função de  $r$  é obtido substituindo-se  $\hat{\beta}_R(r)$  na restrição

$$\left( \hat{\beta}_R(r) - \hat{\beta} \right)' X'X \left( \hat{\beta}_R(r) - \hat{\beta} \right) = r^2$$

$$\left( (kI + X'X)^{-1} X'y - (X'X)^{-1} X'y \right)' X'X \left( (kI + X'X)^{-1} X'y - (X'X)^{-1} X'y \right) = r^2$$

Como é fácil atribuir um valor para  $k$  e obter o valor correspondente de  $r$ , geralmente o estimador Ridge é expresso em função de  $k$  no lugar de  $r$ .

$$\hat{\beta}_R(k) = (kI + X'X)^{-1} X'y$$

Tal substituição é adequada, porém  $k$  não tem um significado geométrico como o tem  $r$ . O problema variacional que define o estimador Ridge admite uma

outra forma equivalente, no espaço de dados. Considere  $\beta$  sobre uma esfera de raio  $r$ , isto é,  $\beta'\beta = r^2$ . A imagem desta esfera pela transformação  $X$  é uma elipse.

A ideia agora é obter  $\beta$  tal que  $X\beta$  esteja o mais próximo possível de  $\mathbf{y}$ , isto é,  $\min_{\beta} \|\mathbf{y} - X\beta\|^2$  restrito à elipse imagem da esfera  $\|\beta\|^2 = r^2$ .

A Lagrangeana desse problema é

$$\begin{aligned} L(\beta, \lambda) &= \|\mathbf{y} - X\beta\|^2 + \lambda(\beta'\beta - r^2) \\ &= (\mathbf{y} - X\beta)'(\mathbf{y} - X\beta) + \lambda(\beta'\beta - r^2) \\ &= \mathbf{y}'\mathbf{y} - \mathbf{y}'X\beta - (X\beta)'\mathbf{y} + (X\beta)'X\beta + \lambda(\beta'\beta - r^2) \\ &= \mathbf{y}'\mathbf{y} - 2\beta'X'\mathbf{y} + \beta'X'X\beta + \lambda(\beta'\beta - r^2) \end{aligned}$$

e então,

$$\begin{aligned} \frac{\partial L}{\partial \beta} &= -2\mathbf{y}'X + 2X'X\beta + \lambda(2\beta) = 0 \\ &-\mathbf{y}'X + (X'X + \lambda I)\beta = 0 \\ \hat{\beta}_R(r) &= (X'X + \lambda I)^{-1}X'\mathbf{y}. \end{aligned}$$

No espaço de dados, tem-se a interpretação

$$\|\mathbf{y} - X\beta\|^2 = \|\mathbf{y} - X\hat{\beta}_{\text{ols}}\|^2 + \|X\beta - X\hat{\beta}_{\text{ols}}\|^2 = \text{cte} + \|X\hat{\beta}_{\text{ols}} - X\beta\|^2.$$

O que se quer então é a minimização

$$\min_{\beta} \|X\beta - X\hat{\beta}\|^2 = \min_{\beta} (\beta - \hat{\beta})'X'X(\beta - \hat{\beta}),$$

sujeito à restrição  $\beta'\beta = r^2$ , isto é, para os vetores  $\beta$  sobre a esfera obter aquele na elipse centrada em  $\hat{\beta}_{\text{ols}}$  de menor tamanho pela métrica de Mahalanobis (Figura 25).

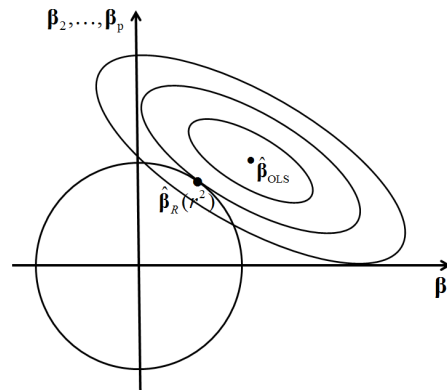


Figura 25 Projeção definida pela métrica de Mahalanobis.

## 4.2 LASSO

Em regressão linear além do problema de estimação do vetor de parâmetros também se coloca o problema de seleção de covariáveis. Na presença de muitas covariáveis é razoável que o pesquisador queira selecionar apenas algumas delas que mais afetam a variável resposta, isto é, um modelo mais parcimonioso. Uma possibilidade é se obter  $\hat{\beta}_{OLS}$  e eliminar as covariáveis relativas às coordenadas de  $\hat{\beta}_{OLS}$  de valores relativamente muito menores que as outras. Se o estimador de quadrados mínimos apresentar alta variabilidade (quasi multicolinearidade), este procedimento evidentemente apresenta problemas e é conhecido que não é estável, ou seja, se um novo vetor de respostas  $\mathbf{y}$  é observado, com alta probabilidade, covariáveis diferentes serão selecionadas.

Uma alternativa seria a de se obter um processo de estimação que além de apresentar pouca variabilidade, geralmente obtida por algum processo de encolhimento, gere com alta probabilidade estimativas  $\hat{\beta}$  em que várias de suas componentes sejam nulas. Desta forma teria-se um processo automático de seleção de covariáveis. Esta é de fato a proposta de um processo de estimação denominado

LASSO, acrônimo de *Least Absolute Shrinkage and Selection Operator*, proposto por Tibshirani (1996). Tal processo é discutido em detalhes neste texto.

Como um dos objetivos deste trabalho é facilitar a leitura do artigo original, será utilizada a mesma notação do artigo. Em particular,  $\hat{\beta}_{\text{ols}}$  será denotado por  $\hat{\beta}^\circ$ .

Sejam então os dados  $(\mathbf{x}_i, y_i)$ , com  $i = 1, \dots, n$  para  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})$  sendo  $\mathbf{x}_i$  o vetor linha da matriz de delineamento, dado pelos valores das variáveis preditoras (covariáveis) e  $y_i$  as respostas observadas. Se  $\beta = (\beta_1, \dots, \beta_p) \in \mathbb{R}^p$ , o estimador LASSO é definido por

$$(\hat{\alpha}, \hat{\beta}) = \arg \min \left\{ \sum_{i=1}^n \left( y_i - \alpha - \sum_j \beta_j x_{ij} \right)^2 \right\},$$

sujeito à restrição  $\sum_{j=1}^p |\beta_j| \leq t$ .

É fácil ver que  $\hat{\alpha} = \bar{y}$ . É adequado por uma transformação de dados colocar o problema na forma normalizada supondo  $\frac{1}{n} \sum_{i=1}^n x_{ij} = 0$ ,  $\frac{1}{n} \sum_{i=1}^n x_{ij}^2 = 1$  e  $\frac{1}{n} \sum_{i=1}^n y_i = \bar{y} = 0$ . Desta forma, o estimador pode ser escrito como

$$\hat{\beta}^{\text{lasso}} = \arg \min \left( \sum_{i=1}^n \left( y_i - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right),$$

ou também pode ser interpretado como  $\arg \min \|\mathbf{y} - \mathbf{X}\beta\|^2$  sujeito a restrição  $\sum_{j=1}^p |\beta_j| \leq t$ .

Esta última forma fica representada na forma Lagrangeana por

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta \in \mathbb{R}^p} \left( \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1 \right).$$

Em termos de modelo geométrico, o subconjunto fechado e convexo  $K_p$ ,

definido por  $\sum_{j=1}^p |\beta_j| \leq t$ , é um hipercubo cujas diagonais estão sobre os eixos coordenados, e o centro sobre a origem. O valor  $t \geq 0$  pode ser considerado como um parâmetro de ajuste. Para o caso bidimensional e tridimensional,  $K_p$  é como se segue na Figura 26.

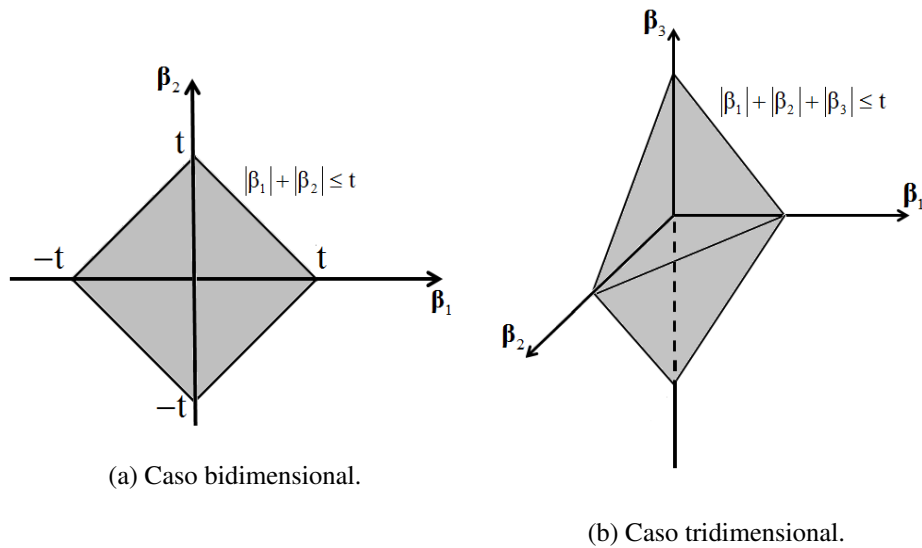


Figura 26 Restrição  $K_p$  LASSO.

No espaço de dados  $\mathbb{R}^n$ , o convexo  $K = X(K_p)$  é semelhante a  $K_p$  substituindo os eixos coordenados por eixos definidos pelos vetores das covariáveis  $x_i$ . O estimador LASSO  $\hat{\beta}^{\text{lasso}}$  é então obtido da forma anteriormente descrita. Obtém-se sobre o convexo  $K$  o ponto mais próximo do vetor de dados  $\mathbf{y}$ , o qual será denominado  $P_K(\mathbf{y})$ . Desta forma,  $X(\hat{\beta}^{\text{lasso}}) = P_K(\mathbf{y})$ . O ponto  $P_K(\mathbf{y})$  também pode ser obtido, conforme descrito anteriormente em termos do triângulo fundamental, como o ponto de  $K$  mais próximo de  $X(\hat{\beta}^0) = X(X'X)^{-1}X'\mathbf{y}$ .

A obtenção de  $\hat{\beta}^{\text{lasso}}$  também pode ser descrita no espaço paramétrico, uma vez que minimizar a distância do vetor  $\mathbf{y}$  ao convexo  $K$  é, como também descrito anteriormente, equivalente a se minimizar a distância da estimativa de



mínimos quadrados  $\hat{\beta}^o$  ao convexo  $K_p$ . Entretanto como esta distância é definida pela métrica  $\langle \mathbf{v}_1, \mathbf{v}_2 \rangle_p = \mathbf{v}'_1 X' X \mathbf{v}_2$ ,  $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^p$ , a obtenção do ponto de  $K_p$  mais próximo de  $\hat{\beta}^o$  pode ser descrita como o primeiro ponto de  $K_p$  que tangencia uma elipse centrada em  $\hat{\beta}^o$ , conforme Figura 27.

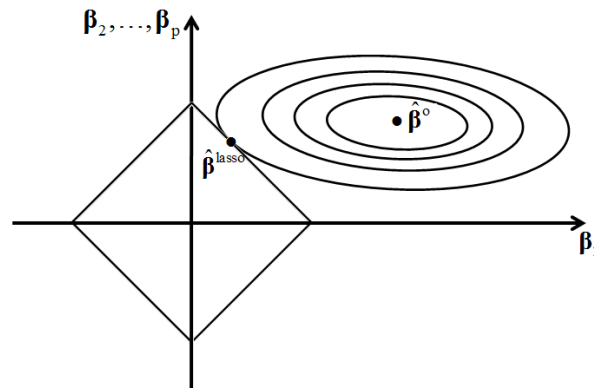


Figura 27 Obtenção do estimador LASSO.

O conjunto convexo  $K_p$  impõe duas restrições que possuem uma justificativa estatística. A primeira delas é manter os valores das estimativas  $(\hat{\beta}^{\text{lasso}})_i$  limitadas. Esta é uma posição conservadora, para evitar estimativas excessivas. A segunda justificativa é bem mais sofisticada. O convexo  $K_p$  tem como borda uma hipersuperfície formada de planos tais que estas se intersectam em superfícies formadas de planos de dimensões menores, que podemos pensar como “arestas”. Como as dimensões das arestas variam, ocorre que os pontos nas arestas possuem algumas coordenadas nulas. É intuitivo pensar que estas arestas têm uma maior probabilidade de conterem o ponto mais próximo de  $\hat{\beta}^o$ . Portanto, o estimador  $\hat{\beta}^{\text{lasso}}$  tem uma maior probabilidade de possuir algumas coordenadas nulas. Neste sentido o estimador funciona também como um processo de seleção de covariáveis e é esta propriedade que mais impulsionou seu estudo e utilização.

Para exemplificar o método, segue um caso particular em que se tem orto-

gonalidade, isto é,  $X'X = I$ .

Passando então a discutir o caso ortogonal, como exemplo mais simples, serão explicitados os cálculos para a situação de 2 covariáveis. Observe que, como por hipótese  $X'X = I$ , o produto interno  $\langle \cdot, \cdot \rangle_p$  é o produto usual e portanto no lugar de elipses centradas em  $\hat{\beta}^o$  tem-se círculos. O estimador LASSO é obtido pela projeção do estimador de quadrados mínimos  $\hat{\beta}^o = (\hat{\beta}_1^o, \hat{\beta}_2^o)$  no retângulo  $|\beta_1| + |\beta_2| \leq t$ . Vamos supor  $\hat{\beta}_1^o > 0$ ,  $\hat{\beta}_2^o > 0$  e  $\hat{\beta}_1^o + \hat{\beta}_2^o > t$ . Uma vez que a projeção ortogonal é a usual, tem-se que os vetores na faixa hachurada, representados na Figura 28, se projetam sobre a aresta, enquanto os vetores fora da faixa hachurada projetam-se no vértice  $(0, t)$  ou  $(t, 0)$ .

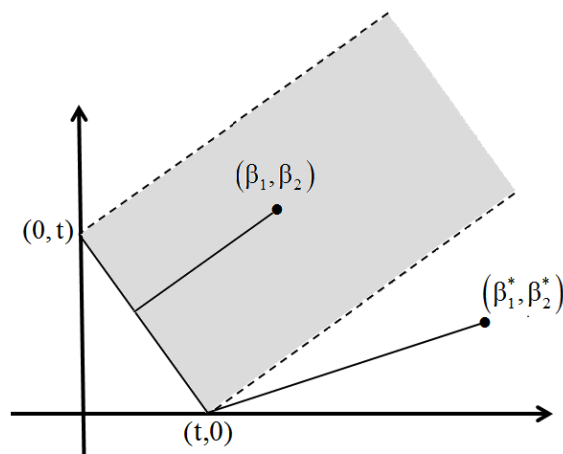


Figura 28 O método de estimação LASSO.

Suponha então que  $\hat{\beta}^o$  esteja na faixa hachurada. A reta parametrizada que passa por  $\hat{\beta}^o$  e é ortogonal à aresta do quadrado é  $\vec{r}(s) = (\hat{\beta}_1^o + s, \hat{\beta}_2^o + s)$ , com  $s \leq 0$  e  $\vec{r}(0) = (\hat{\beta}_1^o, \hat{\beta}_2^o)$ . Esta reta toca a aresta quando

$$\hat{\beta}_1^o + s + \hat{\beta}_2^o + s = t$$

$$\begin{aligned} &\Rightarrow 2s = t - (\hat{\beta}_1^{\circ} + \hat{\beta}_2^{\circ}) \\ &\Rightarrow s = \frac{t}{2} - \frac{1}{2} (\hat{\beta}_1^{\circ} + \hat{\beta}_2^{\circ}). \end{aligned}$$

A intersecção da reta com a aresta é então

$$\left( \hat{\beta}_1^{\circ} + \frac{t}{2} - \frac{1}{2} (\hat{\beta}_1^{\circ} + \hat{\beta}_2^{\circ}), \hat{\beta}_2^{\circ} + \frac{t}{2} - \frac{1}{2} (\hat{\beta}_1^{\circ} + \hat{\beta}_2^{\circ}) \right),$$

isto é,

$$\left( \frac{t}{2} + \frac{\hat{\beta}_1^{\circ} - \hat{\beta}_2^{\circ}}{2}, \frac{t}{2} - \frac{\hat{\beta}_1^{\circ} - \hat{\beta}_2^{\circ}}{2} \right).$$

Portanto, o estimador LASSO é dado por

$$\hat{\beta}^{\text{lasso}} = \left( \frac{t}{2} + \frac{\hat{\beta}_1^{\circ} - \hat{\beta}_2^{\circ}}{2}, \frac{t}{2} - \frac{\hat{\beta}_1^{\circ} - \hat{\beta}_2^{\circ}}{2} \right).$$

Os vetores que se projetam em  $(0, t)$  satisfazem a propriedade de que a reta  $s = \hat{\beta}^{\circ} + s \vec{1}$  intersecta a reta definida pela aresta  $\beta_1 + \beta_2 = t$  no segundo ou quarto quadrante. Para os vetores  $\hat{\beta}$  do primeiro quadrante que não estão na faixa hachurada da Figura 28 ocorre que para algum valor  $\tilde{s}$  uma das coordenadas se anula, pois a reta intersecta os eixos coordenados (Figura 29). Esta condição implica que o ponto de  $K_p$  mais próximo de  $\hat{\beta}^{\circ}$  é o ponto  $(t, 0)$  ou  $(0, t)$ , isto é,  $\hat{\beta}^{\circ}$  se projeta em  $(t, 0)$  ou  $(0, t)$ .

Para  $\hat{\beta}^{\circ}$  como na Figura 29, tem-se que para algum  $\tilde{s}$  a reta  $\vec{r}(s) = (\hat{\beta}_1^{\circ} - s, \hat{\beta}_2^{\circ} - s)$  é tal que

$$\hat{\beta}_2^{\circ} - s = 0 \Rightarrow \hat{\beta}_2^{\circ} = \tilde{s}.$$

Como neste caso  $\hat{\beta}^{\circ}$  está na região definida pelas retas  $(t + s, s)$  e o eixo positivo, então

$$\hat{\beta}_2^{\circ} < \hat{\beta}_1^{\circ} - s \Rightarrow s < \hat{\beta}_1^{\circ} - \hat{\beta}_2^{\circ} \Rightarrow \frac{s}{2} < \frac{\hat{\beta}_1^{\circ} - \hat{\beta}_2^{\circ}}{2}.$$

Desta forma, o estimador LASSO pode ser dado por uma única relação

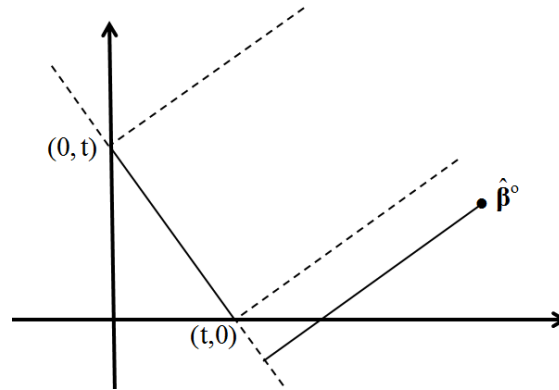


Figura 29 Estimador LASSO no caso bidimensional.

(TIBSHIRANI, 1996, p.272),

$$\hat{\beta}_1^o = \left( \frac{t}{2} + \frac{\hat{\beta}_1^o - \hat{\beta}_2^o}{2} \right)^+ \text{ e } \hat{\beta}_2^o = \left( \frac{t}{2} - \frac{\hat{\beta}_1^o - \hat{\beta}_2^o}{2} \right)^+,$$

em que  $(x)^+$  significa que  $(x)^+ = x$  se  $x \geq 0$  e  $(x)^+ = 0$  se  $x \leq 0$ .

Nas situações em que  $\hat{\beta}^o$  está em outros quadrantes, a análise segue de forma análoga.

Para o caso ortogonal em uma dimensão qualquer, o estimador de quadrados mínimos é da forma

$$\hat{\beta}^o = (X'X)^{-1}X'y = X'y = \left( \sum_{j=1}^n x_{j1}y_j, \dots, \sum_{j=1}^n x_{jk}y_j \right).$$

Se  $\hat{\beta}^o \in K_p$  então  $\hat{\beta}^{\text{lasso}} = \hat{\beta}^o$ . Entretanto, se  $\hat{\beta}^o \notin K_p$  a projeção de  $\hat{\beta}^o$  ocorrerá no bordo de  $K_p$ ,  $(\partial K_p)$ , definido por  $\sum_{i=1}^p |\beta_i| = t$ . O vetor normal a  $\partial K_p$  nas faces regulares (de dimensão  $p - 1$ ) é o vetor  $\vec{1} = (\pm 1, \dots, \pm 1)$ , em que o número de sinais negativos depende da face. Para as faces singulares (intersecção de hiperfaces), o vetor normal é formado por números  $\pm 1$  e 0. Para fixar a ideia, suponha que  $(\hat{\beta}^o)_i \neq 0, \forall i = 1, \dots, p$ . A reta que passa por  $\hat{\beta}^o$  e é ortogonal a  $\partial K_p$  é dada

por  $\vec{r}(s) = \hat{\beta}^\circ + s \vec{1}$  (ver Figura 30).

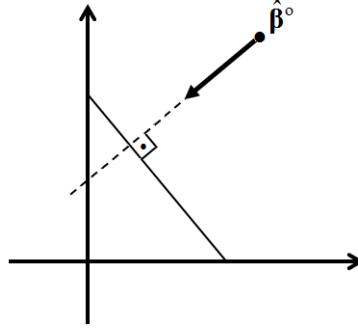


Figura 30 Estimador LASSO para o caso ortogonal.

O estimador ocorre para o primeiro valor  $s$  tal que

$$\vec{r}(s) \in \partial K_p \Rightarrow \sum_{i=1}^p \left| (\hat{\beta}^\circ)_i - s \right| = t.$$

Se  $(\hat{\beta}^\circ)_i > 0$  então  $\left| (\hat{\beta}^\circ)_i - s \right| = (\hat{\beta}^\circ)_i - s$ , para um  $s$  pequeno. Por outro lado, se  $(\hat{\beta}^\circ)_i < 0$  tem-se que  $\left| (\hat{\beta}^\circ)_i - s \right| = -(\hat{\beta}^\circ)_i + s$ , e então

$$\left| (\hat{\beta}^\circ)_i - s \right| = \left[ \text{sinal}(\hat{\beta}^\circ)_i \right] \left| (\hat{\beta}^\circ)_i - s \right|.$$

A reta  $\hat{\beta}^\circ - s \vec{1}$  atinge o bordo  $\partial K_p$  se

$$\begin{aligned} \sum_{i=1}^p \left| (\hat{\beta}^\circ)_i - s \right| = t &\Rightarrow \sum_{i=1}^p \text{sinal}(\hat{\beta}^\circ)_i \left( \left| (\hat{\beta}^\circ)_i \right| - s \right) = t \\ &\Rightarrow \sum_{i=1}^p \text{sinal}(\hat{\beta}^\circ)_i \left( \left| (\hat{\beta}^\circ)_i \right| \right) - ps = t \\ &\Rightarrow s = \frac{t - \sum_{i=1}^p \text{sinal}(\hat{\beta}^\circ)_i \left| (\hat{\beta}^\circ)_i \right|}{p}. \end{aligned}$$

Desta forma, o estimador LASSO é dado por

$$\left(\hat{\beta}^{\text{lasso}}\right)_i = \text{sinal}\left(\hat{\beta}^{\circ}\right)_i \left(\left|\left(\hat{\beta}^{\circ}\right)_i - \gamma\right|\right)^+, \quad (4.1)$$

em que (TIBSHIRANI, 1996, p. 269)

$$\gamma = \frac{\sum_{i=1}^p \text{sinal}\left(\hat{\beta}^{\circ}\right)_i \left|\left(\hat{\beta}^{\circ}\right)_i\right|}{p}.$$

Passemos então a analisar o caso não ortogonal. Para  $p = 2$ , o caso não ortogonal é descrito da mesma forma que a situação ortogonal. O que ocorre agora é que a projeção na aresta do quadrado não é mais uma projeção ortogonal em relação ao produto interno usual  $\langle, \rangle$ , mas sim em relação ao produto interno  $\langle, \rangle_p$ . Para se obter um vetor ortogonal à face, observando que o vetor  $(-1, 1)$  pertence à aresta, é necessário se resolver a equação

$$\langle (-1, 1), (1, a) \rangle_p = (-1, 1)(X'X) \begin{pmatrix} 1 \\ a \end{pmatrix} = 0.$$

A região para a qual  $\beta_1^{\text{lasso}} \neq 0$  e  $\beta_2^{\text{lasso}} \neq 0$  não é mais definida por retas, sendo então definida por curvas mais complexas.

Obtido o valor de  $a$ , define-se a reta  $\vec{r}(s) = \hat{\beta}^{\circ} + s(1, a)$  e obtém-se o valor de  $s$  para o qual a reta  $\vec{r}(s)$  intersecta a aresta e o raciocínio segue como no caso ortogonal. Representação desta ideia é conforme Figura 31.

De forma análoga ao estimador Ridge para o qual se pode traçar o ridge-trace para as componentes do vetor de estimativas, considerando como parâmetro de ajuste o valor  $t$  da definição do estimador LASSO, é possível traçar curvas para as componentes da estimativa LASSO. Um aspecto interessante para estas curvas é que se obtém a sequência ordenada em que as covariáveis deixam de ser zero. Desta forma, tem-se uma descrição gráfica para a seleção das covariáveis

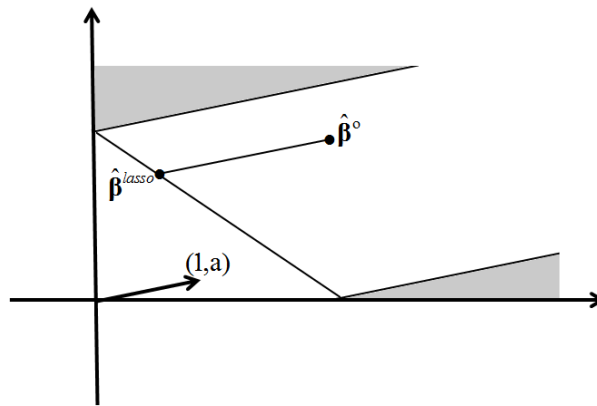


Figura 31 Estimador LASSO (caso geral).

(TIBSHIRANI, 1996, p. 271-273).

### 4.3 Abordagem Bayesiana aos estimadores Ridge e LASSO

O estimador ridge e o estimador LASSO podem ser obtidos facilmente a partir do paradigma bayesiano. Considere  $\mathbf{y}_{n \times 1}$  um vetor aleatório com densidade normal multivariada dada por

$$f(\mathbf{y}) = \frac{1}{(2\pi)^{\frac{n}{2}}} e^{-\frac{1}{2}(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})'(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})},$$

em que  $\mathbf{X}$  é de dimensão  $n \times p$  e  $\boldsymbol{\beta}$  é de dimensão  $p \times 1$ .

Supondo como priori para o vetor de parâmetros  $\boldsymbol{\beta}$  uma normal multivariada da forma

$$\pi(\boldsymbol{\beta}) = \frac{(\lambda)^{\frac{p}{2}}}{(2\pi)^{\frac{p}{2}}} e^{-\frac{\lambda}{2}\|\boldsymbol{\beta}\|^2},$$

a posteriori é proporcional a

$$\pi(\boldsymbol{\beta}|\mathbf{y}) \propto e^{-\frac{1}{2}[(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})'(\mathbf{y}-\mathbf{X}\boldsymbol{\beta}) - \lambda\|\boldsymbol{\beta}\|^2]}.$$

Uma estimativa pode então ser obtida tomando-se a moda da posteriori que ocorre

quando o expoente é minimizado, isto é,

$$\begin{aligned}\hat{\beta}^{(\text{Bayes})} &= \min_{\beta} \left\{ \frac{1}{2} \left[ (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) + \lambda \|\beta\|^2 \right] \right\} \\ &= \min_{\beta} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|^2 \right\} \\ &= \hat{\beta}^{\text{ridge}}.\end{aligned}$$

Para o LASSO, considere uma priori para o vetor  $\beta$  dada pela exponencial dupla

$$\begin{aligned}\pi(\beta) &= \prod_{i=1}^p \frac{1}{2\tau} e^{-\frac{|\beta_i|}{\tau}} \\ &= \frac{1}{(2\tau)^p} e^{-\frac{1}{\tau} \sum_{i=1}^p |\beta_i|} \\ &= \frac{1}{(2\tau)^p} e^{-\frac{1}{\tau} \|\beta\|_1}.\end{aligned}$$

Portanto, a posteriori é proporcional a

$$\pi(\beta | \mathbf{y}) \propto e^{-\frac{1}{2}(\mathbf{y}-\mathbf{X}\beta)'(\mathbf{y}-\mathbf{X}\beta) - \frac{1}{\tau} \sum_{i=1}^p |\beta_i|}.$$

Tomando-se a moda da posteriori tem-se

$$\begin{aligned}\hat{\beta}^{(\text{Bayes})} &= \min_{\beta} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta) + \frac{1}{\tau} \|\beta\|_1 \right\} \\ &= \min_{\beta} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|^2 + 2\lambda \|\beta\|_1 \right\} \\ &= \hat{\beta}^{\text{lasso}}.\end{aligned}$$

#### 4.4 Estimação da variância do estimador LASSO

Como o estimador LASSO é não linear e não diferenciável, é difícil obter estimativas acuradas para sua variância. De fato dois métodos foram propostos no artigo original e serão aqui discutidos em detalhes, sendo o primeiro deles utili-



zando Bootstrap.

Considere o modelo linear usual  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  com  $E[\boldsymbol{\varepsilon}] = 0$  e  $\text{var}[\boldsymbol{\varepsilon}] = \sigma^2\mathbf{I}$ . Um vetor de dados  $\mathbf{y}_{n \times 1}$  é observado e uma estimativa  $\boldsymbol{\beta}_{p \times 1}^{\text{lasso}}$  é obtida. Como as componentes  $y_i$  do vetor  $\mathbf{y}$  possuem esperanças diferentes, não se pode reamostrar as componentes de  $\mathbf{y}$  para se obter uma amostra bootstrap  $\mathbf{y}^*$ . A estratégia para se obter este novo vetor  $\mathbf{y}^*$  é como se segue. Com a estimativa  $\hat{\boldsymbol{\beta}}^{\text{lasso}}$  obtida a partir do vetor de dados originais  $\mathbf{y}$ , obtém-se o vetor ajustado  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}^{\text{lasso}}$ . Note que  $\hat{\mathbf{y}}$  não é um estimador não viesado de  $\boldsymbol{\mu} = E[\mathbf{y}]$ . Toma-se então os resíduos  $\{\varepsilon_i^* = \hat{y}_i - y_i, i = 1, \dots, n\}$ . Conforme observado, não se tem  $E[\varepsilon_i^*] = 0$  e certamente  $E[(\hat{y}_i - y_i)^2] \neq \sigma^2$ . No entanto temos aqui uma analogia com o ajuste de quadrados mínimos pois se  $\hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\beta}^{\text{ols}}$  então  $E[\hat{\mathbf{y}}] = \boldsymbol{\mu}$ ,  $E[(\hat{y}_i - y_i)] = 0$  e  $\frac{1}{n-p-1} \sum_{i=1}^n (\hat{y}_i - y_i)^2$  é um estimador não viesado de  $\sigma^2$ . Neste sentido é razoável considerar  $\hat{y}_i - y_i$  como um estimador de  $\varepsilon_i$ , apesar de estes resíduos não serem independentes, isto é,  $\hat{y}_i - y_i$  e  $\hat{y}_j - y_j$  não são independentes, ao passo que  $\varepsilon_i$  e  $\varepsilon_j$  o são. Baseado nestas considerações, para  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}^{\text{lasso}}$  também será considerado  $\varepsilon_i^* = \hat{y}_i - y_i$  um estimador para o erro  $\varepsilon_i$ . Pode-se agora se obter amostras bootstrap da forma:  $\{\varepsilon_1^*, \varepsilon_2^*, \dots, \varepsilon_n^*\}$  é reamostrado com repetição obtendo-se vetores  $\boldsymbol{\varepsilon}_{n \times 1}^*$ . Com estes vetores constrói-se um vetor de pseudodados  $\mathbf{y}^* = \mathbf{y} + \boldsymbol{\varepsilon}^*$ , em que  $\mathbf{y}$  é o vetor de dados originais. Com o vetor  $\mathbf{y}^*$  obtém-se então uma nova estimativa  $\hat{\boldsymbol{\beta}}^{*(\text{lasso})}$ . Realizando este processo várias vezes, tem-se amostras bootstrap de  $\hat{\boldsymbol{\beta}}^{\text{lasso}}$  e se pode então calcular uma estimativa para sua variância.

Uma outra abordagem é se obter uma forma fechada para uma aproximação de um estimador da variância do estimador LASSO. A ideia é utilizar uma aproximação em termos do estimador Ridge. Como  $|\beta_j| = \frac{(\beta_j)^2}{|\beta_j|}$  tem-se que

$$\|\boldsymbol{\beta}\|_1 = \sum_{i=1}^p |\beta_i| = \sum_{i=1}^p \frac{\beta_i^2}{|\beta_i|}$$

e então

$$\begin{aligned}\hat{\beta}^{\text{lasso}} &= \arg \min_{\beta} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|^2 + t \sum_{i=1}^p |\beta_i| \right\} \\ &= \arg \min_{\beta} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|^2 + t \sum_{i=1}^p \frac{\beta_i^2}{|\beta_i|} \right\}.\end{aligned}$$

Obtido o estimador LASSO  $\hat{\beta}^{\text{lasso}}$ , considere o estimador definido por

$$\hat{\beta}^* = \arg \min_{\beta} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|^2 + t \sum_{i=1}^p \frac{(\beta_i)^2}{|(\beta^{\text{lasso}})_i|} \right\}.$$

Ora,  $\hat{\beta}^*$  é simplesmente o estimador Ridge para o caso em que cada coordenada sofre um encolhimento diferente. Este caso mais geral de estimador Ridge admite uma fórmula fechada semelhante ao estimador Ridge usual dada por

$$\hat{\beta}^* = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{W}^-)^{-1}\mathbf{X}'\mathbf{y},$$

em que  $\mathbf{W}^-$  é da forma da Figura 32. Se  $(\beta^{\text{lasso}})_i = 0$ , então a entrada correspon-

$$\mathbf{W}^- = \begin{bmatrix} \frac{1}{|\beta_1^{\text{lasso}}|} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \frac{1}{|\beta_p^{\text{lasso}}|} \end{bmatrix}$$

Figura 32 Matriz de pesos do estimador Ridge.

dente na diagonal é zero, isto é,  $\mathbf{W}^-$  é a inversa generalizada da matriz diagonal. A quantidade  $\lambda$  é escolhida de forma que  $\sum_{i=1}^p |\hat{\beta}_i^*| = t$ , ou seja, que este estimador ridge satisfaça a mesma restrição do estimador LASSO. Tem-se então que  $\text{cov}(\hat{\beta}^*) = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{W}^-)^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \lambda\mathbf{W}^-)^{-1}\sigma^2$ , e um estimador é então

$$\widehat{\text{cov}}(\hat{\beta}^*) = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{W}^-)^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \lambda\mathbf{W}^-)^{-1}\hat{\sigma}^2,$$

em que  $\hat{\sigma}^2$  é um estimador para a variância do modelo, obtido por exemplo pelo método dos quadrados mínimos.

#### 4.5 Erro de predição e estimação do parâmetro $t$

No artigo original (TIBSHIRANI, 1996) são apresentados, de forma resumida, três métodos para a estimação do parâmetro  $t$ . No sentido de completude do texto os três métodos serão explicitados, utilizando a notação do artigo original.

Considerando o modelo linear  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  com  $E[\boldsymbol{\varepsilon}] = 0$  e  $\text{var}[\boldsymbol{\varepsilon}] = \sigma^2\mathbf{I}$ , o erro médio (*Mean Error*) de um estimador  $\hat{\boldsymbol{\beta}}$  é

$$\text{ME} = E_{\boldsymbol{\beta}} \left[ \left\| \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}\boldsymbol{\beta} \right\|^2 \right].$$

Já o erro de predição (*Prediction Error*) será aqui considerado como

$$\text{PE} = E_{\boldsymbol{\beta}} \left[ \left\| \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \right\|^2 \right],$$

com  $\mathbf{X}\hat{\boldsymbol{\beta}}$  fixo, isto é, o modelo ajustado fica fixo, e toma-se a esperança em relação a novas observações  $\mathbf{y}$ . Assim,

$$\begin{aligned} E_{\boldsymbol{\beta}} \left[ \left\| \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \right\|^2 \right] &= E_{\boldsymbol{\beta}} \left[ \left\| \mathbf{y} - \mathbf{X}\boldsymbol{\beta} + \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\hat{\boldsymbol{\beta}} \right\|^2 \right] \\ &= E_{\boldsymbol{\beta}} \left[ \left\| \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \right\|^2 \right] + E_{\boldsymbol{\beta}} \left[ 2 \langle \mathbf{y} - \mathbf{X}\boldsymbol{\beta}, \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\hat{\boldsymbol{\beta}} \rangle \right] \\ &\quad + E_{\boldsymbol{\beta}} \left[ \left\| \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\hat{\boldsymbol{\beta}} \right\|^2 \right] \\ &= E_{\boldsymbol{\beta}} \left[ \left\| \boldsymbol{\varepsilon} \right\|^2 \right] + 2E_{\boldsymbol{\beta}} \left[ \langle \mathbf{y} - \mathbf{X}\boldsymbol{\beta}, \mathbf{X}\boldsymbol{\beta} \rangle \right] \\ &\quad - 2E \left[ \langle \mathbf{y} - \mathbf{X}\boldsymbol{\beta}, \mathbf{X}\hat{\boldsymbol{\beta}} \rangle \right] + E_{\boldsymbol{\beta}} \left[ \left\| \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\hat{\boldsymbol{\beta}} \right\|^2 \right] \\ &= n\sigma^2 + 2 \langle E[\mathbf{y} - \mathbf{X}\boldsymbol{\beta}], \mathbf{X}\boldsymbol{\beta} \rangle - 2E \left[ \langle \mathbf{y} - \mathbf{X}\boldsymbol{\beta}, \mathbf{X}\hat{\boldsymbol{\beta}} \rangle \right] \\ &\quad + \text{ME} \end{aligned}$$

$$\begin{aligned}
&= n\sigma^2 + 0 - 2E \left[ \langle \mathbf{y} - \mathbf{X}\boldsymbol{\beta}, \mathbf{X}\hat{\boldsymbol{\beta}} \rangle \right] + \text{ME} \\
&= n\sigma^2 - 2 \left[ \langle E[\mathbf{y} - \mathbf{X}\boldsymbol{\beta}], \mathbf{X}\hat{\boldsymbol{\beta}} \rangle \right] + \text{ME} \\
&= n\sigma^2 + \text{ME}.
\end{aligned}$$

O erro de predição do estimador LASSO pode ser estimado por exemplo com o método *fivefold cross-validation*. Neste método os dados são divididos aleatoriamente em cinco partes iguais. Em quatro delas simultaneamente o modelo é ajustado e a outra restante é utilizada como conjunto de teste do modelo ajustado. O processo é então repetido cinco vezes e em seguida toma-se a média do erro de predição. No caso do LASSO o vetor  $\mathbf{y}$  é particionado em cinco partes iguais. Uma delas é escolhida. Retira-se então da matriz  $\mathbf{X}$  as linhas correspondentes a esta parte escolhida. Com esta matriz reduzida e o vetor de dados reduzido é obtida uma estimativa LASSO  $\hat{\boldsymbol{\beta}}_1$ . O processo é repetido cinco vezes obtendo assim  $\hat{\boldsymbol{\beta}}_1, \dots, \hat{\boldsymbol{\beta}}_5$ . Observa-se que apesar destes vetores terem sido obtidos pela matriz  $\mathbf{X}$  reduzida de um quinto das linhas, os vetores  $\hat{\boldsymbol{\beta}}_i$  obtidos são vetores em  $\mathbb{R}^p$ . Para o erro de predição a matriz original  $\mathbf{X}$  é aplicada nos vetores  $\hat{\boldsymbol{\beta}}_i$  obtendo-se vetores ajustados  $\hat{\mathbf{y}}_i$  em  $\mathbb{R}^n$ . Calcula-se então, para as componentes ao um quinto das coordenadas do vetor de dados originais o quadrado da diferença do valor ajustado e dos valores observados. O processo é repetido cinco vezes e toma-se a média destes. A escolha do valor do parâmetro  $t$  é obtida fazendo a reparametrização  $s = \frac{t}{\sum_{i=1}^p |\beta_i^{\text{ols}}|}$ . Como o valor de  $t$  na estimação LASSO só é interessante para  $0 < t \leq \sum_{i=1}^p |\beta_i^{\text{ols}}|$ , o novo parâmetro  $s$  varia de zero a um. Faz-se então uma partição de  $[0, 1]$  da forma  $s = 0, \frac{1}{\ell}, \frac{2}{\ell}, \dots, \frac{\ell-1}{\ell}, 1$  para  $\ell$  suficientemente grande e calcula-se a estimativa do erro de predição para os estimadores LASSO para cada um dos valores de  $s$ . Toma-se então por estimativa de  $t$  o valor de  $s$  correspondente que gera o menor erro de predição. Observa-se que como o erro de predição PE

e o erro quadrático médio ME diferem pelo valor da variância  $\sigma^2$ , a computação pode ser feita em relação a ME. Tem-se então uma fórmula computacionalmente mais simples, uma vez que para modelos lineares se tem

$$\text{ME} = E \left[ \left\| X\hat{\beta}^{\text{lasso}} - X\beta \right\|^2 \right] = E \left[ \left( \hat{\beta}^{\text{lasso}} - \beta \right)' X'X \left( \hat{\beta}^{\text{lasso}} - \beta \right) \right],$$

e portanto pode ser estimado por  $\frac{1}{5} \sum_{i=1}^5 \left( \hat{\beta}_i^{\text{lasso}} - \beta \right)' X'X \left( \hat{\beta}_i^{\text{lasso}} - \beta \right)$ .

O segundo método é baseado na aproximação do estimador  $\beta^{\text{lasso}}$  pelo estimador Ridge  $\beta^*(\lambda) = (X'X + \lambda W^-)^{-1} X'y$ , para  $W = \text{diag} \left( \left| \beta_j^{\text{lasso}} \right| \right)$ , conforme descrito em 4.4. O erro de predição do estimador LASSO  $\hat{\beta}^{\text{lasso}}$  será aproximado pelo erro de predição do estimador  $\beta^*(\lambda)$  que, por ser um estimador linear, é calculado por um processo de *cross-validation* tipo *leave one out*. Seja  $\beta^{*(k)}(\lambda)$  o estimador obtido omitindo-se a k-ésima coordenada do vetor de dados  $y$  e a linha correspondente na matriz  $(X'X + \lambda W^-)^{-1} X'$ . Obtém-se então a estimativa do erro de predição

$$P(\lambda) = \frac{1}{n} \left\| X\beta^{*(k)}(\lambda) - y \right\|^2 = \frac{1}{n} \sum_{k=1}^n \left( \left( X\beta^{*(k)}(\lambda) \right)_{(k)} - y_k \right)^2,$$

denominada na literatura de PRESS, acrônimo de *The Predicted Residual Sum of Squares*. De forma surpreendente Golub; Heath; Wahba (1979) demonstram que  $P(\lambda)$  admite a seguinte expressão matricial simples

$$P(\lambda) = \frac{1}{n} \left\| D(\lambda) \left( I - X(X'X + \lambda W^-)^{-1} X'y \right) \right\|^2,$$

em que  $D(\lambda) = \text{diag} \left( \frac{1}{1 - a_{ii}} \right)$  e  $a_{ii} = \left( X(X'X + \lambda W^-)^{-1} X' \right)_{ii}$ .

A utilização desta fórmula permite obter para o caso do estimador Ridge o valor ótimo para o parâmetro  $\lambda$  simplesmente obtendo o valor de  $\lambda$  que minimiza  $P(\lambda)$ , isto é, que minimiza o erro de predição. Golub, op. cit., apresentam uma generalização deste método, que consiste em uma versão invariante por rotação.

Aplica-se ao vetor de dados  $\mathbf{y}$  uma matriz especial obtendo-se uma nova estimativa para o erro de predição dada por

$$V(\lambda) = \frac{1}{n} \sum_{k=1}^n \left( \left( \mathbf{X} \boldsymbol{\beta}^{*(k)}(\lambda) \right)_k - y_k \right)^2 \omega_k(\lambda),$$

em que

$$\omega_k(\lambda) = \frac{1 - a_{kk}(\lambda)}{1 - \frac{1}{n} \text{tr}(\mathbf{A}(\lambda))}.$$

A estimativa para o parâmetro  $t$  é obtida como um minimizador de  $V(\lambda)$  ( $t$  e  $\lambda$  estão univocamente relacionados). Este procedimento recebeu o nome de GCV (Generalized Cross Validation). No artigo original do LASSO, é proposta uma estatística semelhante da forma

$$\text{GCV}(t) = \frac{1}{n} \frac{\|\mathbf{X} \boldsymbol{\beta}^{\text{lasso}}(t) - \mathbf{y}\|^2}{\left\{ 1 - \frac{1}{n} \text{tr}(\mathbf{X}(\mathbf{X}'\mathbf{X} + \lambda \mathbf{W}^-)^{-1} \mathbf{X}') \right\}^2}.$$

O terceiro método é baseado no estimador não viesado de Stein (SURE) (ver Seção 7.2). Suponha  $\mathbf{Z} \sim N_n(\boldsymbol{\mu}, \mathbf{I})$ . Se  $\hat{\boldsymbol{\mu}}$  é um estimador de  $\boldsymbol{\mu}$ , fazendo  $\hat{\boldsymbol{\mu}} = \mathbf{z} + \mathbf{g}(\mathbf{z})$ , então

$$E_{\boldsymbol{\mu}} \left[ \|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|^2 \right] = n + E_{\boldsymbol{\mu}} \left[ \|\mathbf{g}(\mathbf{z})\|^2 + 2 \sum_{i=1}^n \frac{\partial g_i}{\partial z_i} \right],$$

e pelo método dos momentos tem-se o estimador não viesado para o risco

$E_{\boldsymbol{\mu}} \left[ \|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|^2 \right]$  do estimador  $\hat{\boldsymbol{\mu}}$ , dado por  $n + \|\mathbf{g}(\mathbf{z})\|^2 + 2 \sum_{i=1}^n \frac{\partial g_i}{\partial z_i}$ . Supondo ortogonalidade, isto é,  $\mathbf{X}'\mathbf{X} = \mathbf{I}$ , tem-se uma fórmula explícita para o estimador LASSO,

$$\hat{\beta}_j^{\text{lasso}} = \text{sinal}(\hat{\beta}_j^{\circ}) \left( \left| \hat{\beta}_j^{\circ} \right| - \gamma \right)^+,$$

em que  $\gamma$  é definido por uma relação com  $t = \sum_{i=1}^p \left| \hat{\beta}_i^{\text{lasso}} \right|$ . Estimando o erro padrão de  $\hat{\beta}_j^{\circ}$  por  $\hat{\tau} = \frac{\hat{\sigma}}{\sqrt{n}}$  em que  $\hat{\sigma}^2 = \frac{1}{n-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ , então  $\frac{\hat{\beta}_j^{\circ}}{\hat{\tau}}$ ,  $j = 1, \dots, n$ ,

são aproximadamente independentes e normais padrão. Em razão da natureza geométrica do processo de estimação LASSO, pode-se afirmar que o estimador  $\hat{\beta}_j^{\text{lasso}} = \text{sinal} \left( \hat{\beta}_j^{\circ} \right) \left( \left| \hat{\beta}_j^{\circ} \right| - \gamma \right)^+$  é próximo ao estimador definido por

$$\hat{\mu}_i = \text{sinal} \left( \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right) \left( \left| \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right| - \gamma \right)^+,$$

conforme Figura 33. Fazendo  $z_i = \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}}$  tem-se

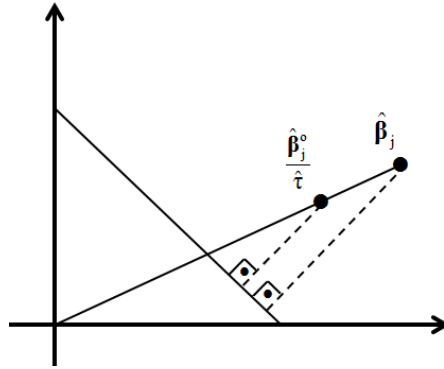


Figura 33 Aproximação do estimador LASSO.

$$\hat{\mu}_i = \text{sinal} \left( \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right) \left( \left| \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right| - \gamma \right)^+ = \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} - \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} + \text{sinal} \left( \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right) \left( \left| \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right| - \gamma \right)^+.$$

Definindo a função  $g(\mathbf{z}) = (g_1(z_1), \dots, g_p(z_p))$  por

$$g_i(z_i) = g_i \left( \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right) = -\frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} + \text{sinal} \left( \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right) \left( \left| \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right| - \gamma \right)^+,$$

$$\text{se } \hat{\beta}_j^{\circ} < 0 \text{ e } \left| \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right| < \gamma \text{ então } g_i \left( \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right) = -\frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} - 0.$$

$$\text{Se } \hat{\beta}_j^{\circ} < 0 \text{ e } \left| \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right| > \gamma \text{ então } g_i \left( \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right) = -\frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} - \left| \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right| + \gamma = \gamma.$$

$$\text{Se } \hat{\beta}_j^{\circ} > 0 \text{ e } \left| \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right| < \gamma \text{ então } g_i \left( \frac{\hat{\beta}_j^{\circ}}{\hat{\tau}} \right) = -\frac{\hat{\beta}_j^{\circ}}{\hat{\tau}}.$$

Finalmente, para  $\hat{\beta}_j^o > 0$  e  $\left|\frac{\hat{\beta}_j^o}{\hat{\tau}}\right| > \gamma$  a função  $g$  é  $g_i\left(\frac{\hat{\beta}_j^o}{\hat{\tau}}\right) = -\gamma$ . Note que  $\frac{\partial g_i}{\partial \frac{\hat{\beta}_j^o}{\hat{\tau}}} = 1$  sempre que  $\left|\frac{\hat{\beta}_j^o}{\hat{\tau}}\right| < \gamma$ .

Em suma, a função  $\hat{\boldsymbol{\mu}} = \mathbf{z} + \mathbf{g}(\mathbf{z})$  para  $\mathbf{g}(\mathbf{z}) = -\mathbf{z} + \text{sinal}(\mathbf{z})(|\mathbf{z}| - \gamma)^+$  e  $\mathbf{z} = \frac{\hat{\boldsymbol{\beta}}^o}{\hat{\tau}}$  possui coordenadas com gráfico dado pela Figura 34. Portanto, o estimador não viesado de Stein para o erro quadrático médio,  $\widehat{E}_{\boldsymbol{\mu}} \left[ \|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|^2 \right] = p + \|\mathbf{g}(\mathbf{z})\|^2 + 2 \sum_{i=1}^n \frac{\partial g_i}{\partial z_i}$ , tem sua expressão dada por

$$R \left\{ \hat{\boldsymbol{\beta}}(\gamma) \right\} \approx \hat{\tau}^2 \left\{ p - 2\# \left( j; \left| \frac{\hat{\beta}_j^o}{\hat{\tau}} \right| < \gamma \right) + \sum_{j=1}^p \max \left( \left| \frac{\hat{\beta}_j^o}{\hat{\tau}} \right|, \gamma \right)^2 \right\}.$$

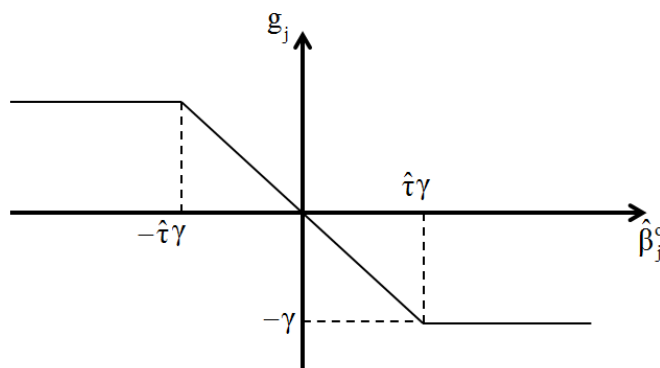


Figura 34 Gráfico das funções coordenadas.

A equação que relaciona os parâmetros  $t$  e  $\gamma$  pode ser explicitada a partir do fato de que  $\hat{\beta}_j^o = \text{sinal}(\hat{\beta}_j^o) \left( \left| \hat{\beta}_j^o \right| - \gamma \right)^+$  e, uma vez que  $\sum_{j=1}^p \left| \hat{\beta}_j^o \right| = t$ , tem-se então

$$\begin{aligned} \sum_{j=1}^p \left[ \left| \text{sinal}(\hat{\beta}_j^o) \right| \left( \left( \left| \hat{\beta}_j^o \right| - \gamma \right)^+ \right) \right] &= t, \\ \Rightarrow \hat{t} &= \sum_{j=1}^p \left( \hat{\beta}_j^o - \hat{\gamma} \right)^+. \end{aligned}$$

O artigo apresenta uma justificativa heurística para a razão de que o caso ortogo-



nal também pode ser aplicado como uma aproximação para o risco do estimador LASSO no caso geral.

#### 4.6 O caso $p > n$

É usual em problemas de regressão, principalmente em aplicações genéticas, que o número de covariáveis  $p$  seja maior do que o número de observações  $n$ . É observado sem demonstração em Zou; Hastie (2005) que o método LASSO seleciona no máximo  $n$  covariáveis, sendo esta uma limitação do método que foi superada com o método Elasticnet. Na literatura não foi encontrada uma demonstração explícita deste fato. Segue portanto uma demonstração baseada em argumentos geométricos que justifica tal fato.

No caso  $p > n$ , a matriz  $X_{n \times p}$  não é injetiva e possui kernel, isto é, o subespaço  $\text{Ker}(X) = \{\beta, X\beta = \mathbf{0}\}$ . O estimador LASSO é definido então por

$$\hat{\beta}^{lasso} = \arg \min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - X\beta\|^2 + \lambda \|\beta\|_1.$$

$\hat{\beta}^{lasso}$  é obtido geometricamente projetando-se o estimador de quadrados mínimos, pela métrica de Mahalanobis, no convexo definido por  $\|\beta\|_1 \leq t$ . Neste caso, para se obter o estimador de quadrados mínimos faz-se necessário a utilização de uma inversa generalizada de  $X'X$ , definindo a estimativa dada por  $\hat{\beta}^{ols} = (X'X)^- X'y$ . Neste caso, tem-se que se  $\beta \in \text{Ker}(X)$  então

$$\begin{aligned} \left\| \mathbf{y} - X \left( \hat{\beta}^{ols} + \beta \right) \right\|^2 + \lambda \left\| \hat{\beta}^{ols} + \beta \right\|_1 &= \left\| \mathbf{y} - X \left( \hat{\beta}^{ols} \right) \right\|^2 + \lambda \left\| \hat{\beta}^{ols} + \beta \right\|_1 \\ &\leq \left\| \mathbf{y} - X \left( \hat{\beta}^{ols} \right) \right\|^2 + \lambda \left\| \hat{\beta}^{ols} \right\|_1 + \lambda \|\beta\|_1. \end{aligned}$$

Tem-se então que se deve minimizar  $\|\beta\|_1$  para  $\beta \in \text{Ker}(X)$ . Para a norma  $\|\cdot\|_1$  isto ocorre quando o subespaço paralelo ao  $\text{Ker}(X)$ , e que contém o  $\hat{\beta}^{ols}$ , intersec-

ciona os planos coordenados, conforme Figura 35 . Supondo sem perda de generalidade que a intersecção seja transversal, então se tem no máximo uma intersecção com o plano coordenado de dimensão  $n$ , e portanto no máximo  $n$  covariáveis são selecionadas.

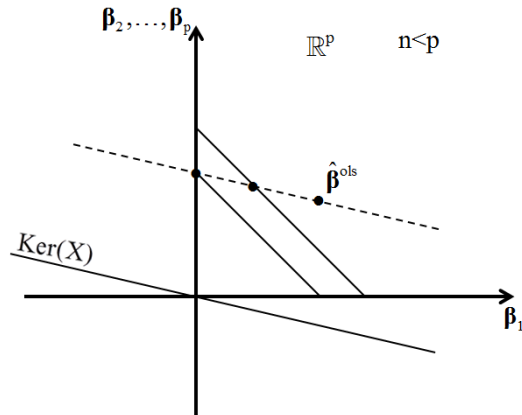


Figura 35  $\text{Ker}(X)$  e o convexo  $K_p$

#### 4.7 GARROTE

O método foi proposto originalmente no artigo clássico *Better Subset Regression Using the Nonnegative Garrote* (BREIMAN, 1995), como uma forma seleção de covariáveis em regressão linear. Este método foi que motivou a definição do estimador LASSO. Neste trabalho, o método GARROTE será discutido de forma simplificada. No problema de regressão linear  $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , com  $\hat{\boldsymbol{\beta}}^{\text{ols}}$  sendo o estimador de quadrados mínimos, considere o conjunto

$$C_p = \left\{ (c_1, \dots, c_p); c_j \geq 0 \text{ e } \sum_{j=1}^p c_j \leq t \right\}.$$

Para  $p = 2$  e  $p = 3$  estes conjuntos são, respectivamente, um triângulo e um tetraedro, conforme Figura 36. A partir do convexo  $C_p$ , é construído um convexo  $K_p$  no

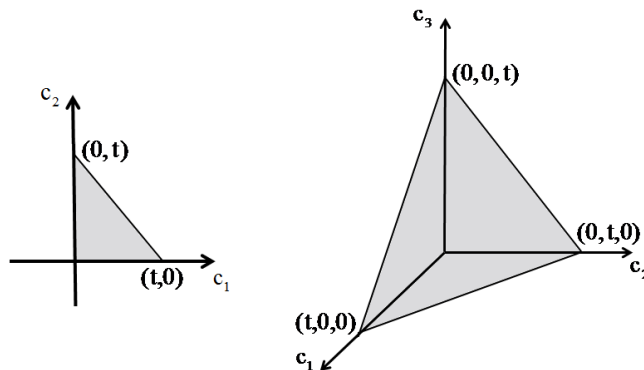


Figura 36 O convexo  $C_p$  em  $\mathbb{R}^2$  e  $\mathbb{R}^3$

espaço paramétrico da forma como se segue. Seja  $\hat{\beta}^o = (\hat{\beta}_1, \dots, \hat{\beta}_p)$  a estimativa de quadrados mínimos. Considere a transformação linear  $A : \mathbb{R}^p \rightarrow \mathbb{R}^p$  em que  $A$  é a matriz diagonal  $A = \text{diag}(\hat{\beta}_i^o)$ . Logo, a imagem de  $C_p$  por  $A$ ,  $K_p = A(C_p)$ , é

$$K_p = \left\{ \left( c_1 \hat{\beta}_1^o, \dots, c_p \hat{\beta}_p^o \right), c_j \geq 0, \sum_{j=1}^p c_j \leq t \right\}.$$

A forma de  $K_p$  é similar à de  $C_p$ , possuindo entretanto inclinações diferentes. Para  $p = 2$ , tem-se a representação na Figura 37. É possível observar que se  $t = p$  então  $(1, \dots, 1) \in C_p$  e conseqüentemente a estimativa de quadrados mínimos  $(\hat{\beta}_1, \dots, \hat{\beta}_p) \in K_p$ . Se  $t < p$  então  $\hat{\beta}^o \notin K_p$ . O método de estimação GARROTE consiste em se projetar ortogonalmente, em relação à métrica  $\langle \cdot, \cdot \rangle_p$ ,  $\hat{\beta}^o$  em  $K_p$ . No espaço de dados, a projeção de  $y$  é ortogonal em relação ao produto usual em  $K = X(K_p)$ .

Analiticamente, o GARROTE pode ser definido como

$$\hat{\beta}^{\text{garrote}} = \arg \min \sum_{i=1}^n \left( \mathbf{y}_i - \sum_{j=1}^p x_{ij} c_j \hat{\beta}_j^{\circ} \right)^2$$

sujeito à restrição  $\sum_{j=1}^p c_j \leq t$ ,  $c_j \geq 0$  (TIBSHIRANI, 1996; BREIMAN, 1995). A projeção no processo GARROTE possui tendência a ocorrer nas faces de menor dimensão e portanto tende a obter o estimador GARROTE com vários componentes nulos, o que pode ser visto como um método de seleção de variáveis. Como se pode observar, o método LASSO é semelhante ao GARROTE, diferindo no fato de que o convexo  $K_p$  no GARROTE depende da estimativa de quadrados mínimos.

Um caso particular é quando se tem  $t \geq p$ . Nesta situação,  $\sum_{j=1}^p c_j = p \leq t$  é satisfeito por  $c_1 = c_2 = \dots = c_p = 1$ , implicando que  $\hat{\beta}^{\circ} \in K_p$  e então o estimador GARROTE coincide com o estimador de quadrados mínimos  $\hat{\beta}^{\circ}$ .

Como exemplo seja o caso ortogonal, em que  $X'X = I$  e  $p = 2$  (Figura 37). O vetor normal à reta aresta que passa por  $(0, t\hat{\beta}_2^{\circ})$  e  $(t\hat{\beta}_1^{\circ}, 0)$  é

$$\langle (t\hat{\beta}_1^{\circ}, -t\hat{\beta}_2^{\circ}), (1, a) \rangle = 0$$

$$t\hat{\beta}_1^{\circ} - a t\hat{\beta}_2^{\circ} = 0 \Rightarrow a = \frac{\hat{\beta}_1^{\circ}}{\hat{\beta}_2^{\circ}}.$$

A reta que passa por  $\hat{\beta}^{\circ}$  e tem este vetor como vetor diretor é

$$r(s) = \hat{\beta}^{\circ} + s \left( 1, \frac{\hat{\beta}_1^{\circ}}{\hat{\beta}_2^{\circ}} \right).$$

A intersecção desta reta com a reta aresta ocorre para

$$\frac{\hat{\beta}_2^{\circ} + s \frac{\hat{\beta}_1^{\circ}}{\hat{\beta}_2^{\circ}}}{t\hat{\beta}_1^{\circ} - (\hat{\beta}_1^{\circ} + s)} = \frac{\hat{\beta}_2^{\circ}}{\hat{\beta}_1^{\circ}}$$

$$\hat{\beta}_1^o + s \left( \frac{\hat{\beta}_1^o}{\hat{\beta}_2^o} \right)^2 = t\hat{\beta}_1^o - \hat{\beta}_1^o - s$$

$$s \left( \left( \frac{\hat{\beta}_1^o}{\hat{\beta}_2^o} \right)^2 + 1 \right) = t\hat{\beta}_1^o - 2\hat{\beta}_1^o$$

$$s = \frac{\hat{\beta}_1^o - 2\hat{\beta}_1^o}{1 + \left( \frac{\hat{\beta}_1^o}{\hat{\beta}_2^o} \right)^2}.$$

Obtendo os estimadores:

$$\left( \left( 1 + \frac{t-2}{1 + \left( \frac{\hat{\beta}_1^o}{\hat{\beta}_2^o} \right)^2} \right)^+ \hat{\beta}_1^o, \left( 1 + \frac{t-2}{1 + \left( \frac{\hat{\beta}_2^o}{\hat{\beta}_1^o} \right)^2} \right)^+ \hat{\beta}_2^o \right).$$

A fórmula anterior difere da fórmula apresentada em Tibshirani (1996, p. 269).

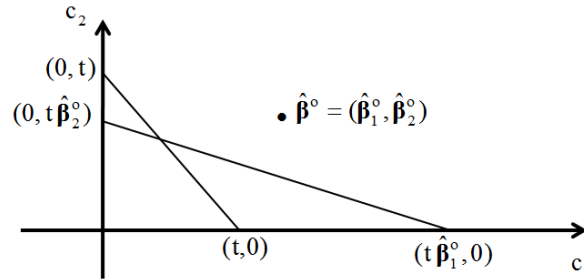


Figura 37 Penalização do método GARROTE.

O grande problema relativo ao método GARROTE é que como este utiliza a estimativa de quadrados mínimos na construção do convexo  $K_p$ , se o estimador de quadrados mínimos apresenta grande variabilidade, o mesmo ocorrerá com o estimador GARROTE.

## 5 O método de regressão *Elastic Net*

Este método de regressão pretende obter o melhor entre os dois métodos de regressão, Ridge e LASSO, porém com novas propriedades que estes métodos carecem. A ideia é bastante simples, consistindo em minimizar a soma de quadrados do resíduo restrito a uma combinação linear das restrições do método Ridge e do LASSO.

Para uma regressão múltipla  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , considerando o modelo centrado e as covariáveis na forma de correlação, isto é,  $\sum y_i = 0$ ,  $\sum_{i=1}^n x_{ij} = 0$ ,  $\sum_{i=1}^n x_{ij}^2 = 1$  com  $j = 1, \dots, p$ , a penalização da soma de quadrados será dada por

$$\lambda_1 \|\boldsymbol{\beta}\|^2 + \lambda_2 \|\boldsymbol{\beta}\|_1,$$

em que  $\|\boldsymbol{\beta}\|^2 = \sum_{j=1}^p \beta_j^2$  e  $\|\boldsymbol{\beta}\|_1 = \sum_{j=1}^p |\beta_j|$ . Portanto, a função a ser minimizada é dada por  $L(\boldsymbol{\beta}, \lambda_1, \lambda_2) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda_1 \|\boldsymbol{\beta}\|^2 + \lambda_2 \|\boldsymbol{\beta}\|_1$ . O estimador obtido será denominado estimador *elastic net* ingênuo, representado por

$$\hat{\boldsymbol{\beta}}_{eni} = \arg \min_{\boldsymbol{\beta}} L(\boldsymbol{\beta}, \lambda_1, \lambda_2)$$

### 5.1 O estimador *elastic net*

Em termos de soma de quadrados restrita (penalizada) minimizar a função  $L(\boldsymbol{\beta}, \lambda_1, \lambda_2)$  é equivalente a minimizar  $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|$  sujeito à restrição

$$\begin{aligned} \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|^2 \leq t &\Rightarrow \frac{\lambda_1}{\lambda_1 + \lambda_2} \|\boldsymbol{\beta}\|_1 + \frac{\lambda_2}{\lambda_1 + \lambda_2} \|\boldsymbol{\beta}\|^2 \leq \frac{t}{\lambda_1 + \lambda_2} \\ &\Rightarrow \left(1 - \frac{\lambda_2}{\lambda_1 + \lambda_2}\right) \|\boldsymbol{\beta}\|_1 + \frac{\lambda_2}{\lambda_1 + \lambda_2} \|\boldsymbol{\beta}\|^2 \leq \frac{t}{\lambda_1 + \lambda_2} \\ &\Rightarrow (1 - \alpha) \|\boldsymbol{\beta}\|_1 + \alpha \|\boldsymbol{\beta}\|^2 \leq t'. \end{aligned}$$

$(1 - \alpha) \|\beta\|_1 + \alpha \|\beta\|^2$  é denominada de penalidade *elastic net* e é uma combinação convexa entre as penalidades que definem a estimação Ridge e LASSO. Pode-se observar que para  $\alpha = 1$  se tem a estimação Ridge e se  $\alpha = 0$  a estimação resultante é a do LASSO.

O conjunto  $\|\beta\|_1 \leq t$  é convexo, mas não estritamente convexo. Já o conjunto  $\|\beta\|^2 \leq t$  é estritamente convexo, e portanto para  $\alpha \neq 0$ , o conjunto definido por  $(1 - \alpha) \|\beta\|_1 + \alpha \|\beta\|^2 \leq t$  é estritamente convexo. Logo o modelo  $K_p$  é o convexo

$$K_p = \left\{ \beta, (1 - \alpha) \|\beta\|_1 + \alpha \|\beta\|^2 \leq t \right\}.$$

Para analisar a forma deste conjunto convexo basta analisar o seu bordo, isto é, analisar a igualdade  $(1 - \alpha) \|\beta\|_1 + \alpha \|\beta\|^2 = t$ . Como exemplo, para  $t = 2$ ,  $\alpha = 0,5$  e  $p = 2$  se tem  $(0,5) \|\beta\|_1 + (0,5) \|\beta\|^2 = 2$ . Assim, se  $\|\beta\|_1 = 2$  e  $\|\beta\|^2 = 2$ , tem-se a representação conforme Figura 38.

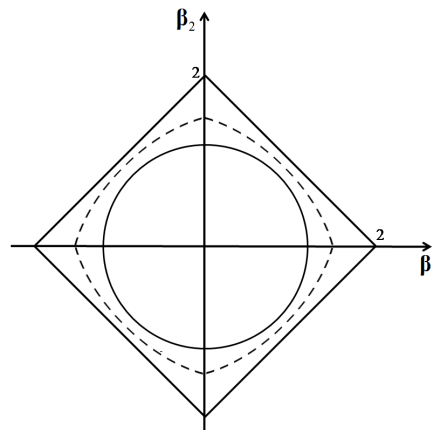


Figura 38  $K_p$  na estimação *elastic net*.

Novamente o processo de estimação *elastic net* consiste em projetar com distância mínima o estimador de quadrados mínimos  $\hat{\beta}^o$  em  $K_p$ . Geometricamente a projeção é obtida como o ponto de tangência entre a família de elipses centrada

em  $\hat{\beta}^o$  e a curva  $(1 - \alpha) \|\beta\|_1 + \alpha \|\beta\|^2 = t$ , conforme representação na Figura 39.

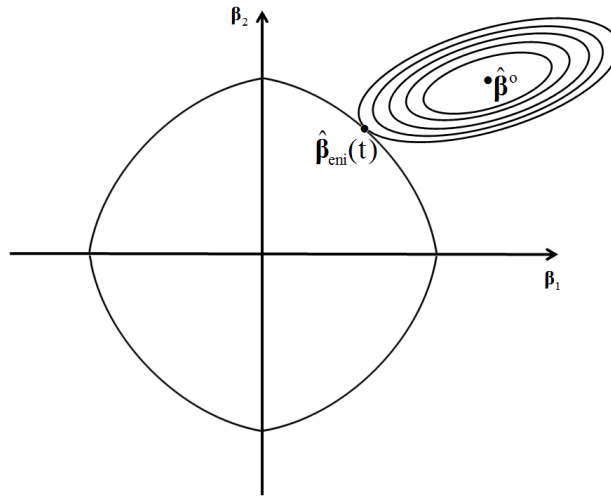


Figura 39 Geometria do estimador *elastic net*.

Observando novamente a Figura 38 tem-se que nos pontos  $(\pm t, 0)$  e  $(0, \pm t)$  a curva não é diferenciável, isto é, se tem “pontas”. Esta é a razão pela qual o método *elasticnet* mantém as propriedades do LASSO de seleção de covariáveis.

No espaço de dados o método pode ser visualizado conforme descrito a seguir. A aplicação  $X$  leva o conjunto convexo  $K_p$ ,  $(1 - \alpha) \|\beta\|_1 + \alpha \|\beta\|^2 = t$ , em um conjunto convexo  $K \subset \text{Im}(X)$ . O processo é simplesmente obter o vetor  $X\beta$  em  $K$  que esteja o mais próximo de  $y$  (Figura 40).

O processo do *elastic net* ingênuo pode ser visto como um processo de duas etapas: uma regressão Ridge e uma regressão tipo LASSO. Este fato será fundamental na construção de algoritmos eficientes para o cálculo de estimativas *elasticnet*. A ideia é que a regressão Ridge pode ser obtida com o emprego de estimadores mistos (GRUBER, 1998; COSTA, 2015), isto é, utilizando um modelo



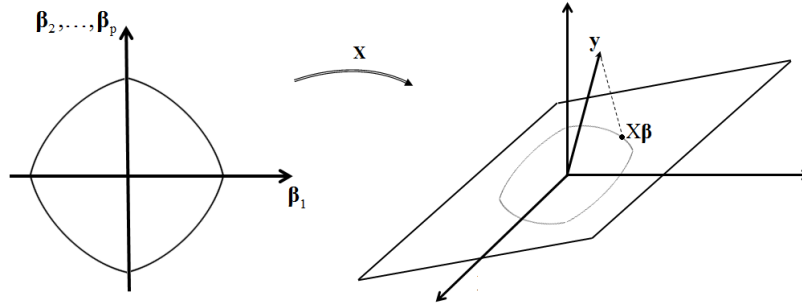


Figura 40 O método de regressão *elastic net*.

linear aumentado a partir do modelo original. Este procedimento está descrito a seguir.

Primeiramente faz-se uma reparametrização, em que os novos parâmetros são da forma  $\beta^* = \sqrt{1 + \lambda_2} \beta$ . A partir do conjunto de dados  $(\mathbf{y}_{n \times 1}, \mathbf{X}_{n \times p})$  define-se um novo conjunto de dados (dados aumentados) que podem por exemplo ser provenientes de um outro experimento anteriormente realizado. Com estes dados aumentados a regressão fica da forma  $(\mathbf{y}_{(n+p) \times 1}^*, \mathbf{X}_{(n+p) \times p}^*)$ , em que  $\mathbf{y}_{(n+p) \times 1}^* = \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix}$  e  $\mathbf{X}_{(n+p) \times p}^* = \frac{1}{\sqrt{1 + \lambda_2}} \begin{pmatrix} \mathbf{X}_{n \times p} \\ \sqrt{\lambda_2} \mathbf{I}_{p \times p} \end{pmatrix}$ . Assim, tem-se então uma nova regressão em relação aos novos parâmetros  $\beta^*$  dada por  $\mathbf{y}^* = \mathbf{X}^* \beta^* + \varepsilon$ . O estimador de quadrados mínimos desta regressão é

$$\begin{aligned} \hat{\beta}^{*o} &= (\mathbf{X}^{*'} \mathbf{X}^*)^{-1} \mathbf{X}^{*'} \mathbf{y}^* \\ &= \frac{1 + \lambda_2}{\sqrt{1 + \lambda_2}} \left( (\mathbf{X}', \sqrt{\lambda_2} \mathbf{I}) \begin{pmatrix} \mathbf{X} \\ \sqrt{\lambda_2} \mathbf{I} \end{pmatrix} \right)^{-1} (\mathbf{X}', \sqrt{\lambda_2} \mathbf{I}) \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} \\ &= \sqrt{1 + \lambda_2} (\mathbf{X}' \mathbf{X} + \lambda_2 \mathbf{I})^{-1} \mathbf{X}' \mathbf{y}. \end{aligned}$$

Logo, como  $\hat{\beta}^{*o} = \sqrt{1 + \lambda_2} \hat{\beta}$ , tem-se que  $\hat{\beta} = (\mathbf{X}' \mathbf{X} + \lambda_2 \mathbf{I})^{-1} \mathbf{X}' \mathbf{y} = \hat{\beta}_R(\lambda_2)$ .

Portanto o estimador de quadrados mínimos  $\hat{\beta}^{o*}$  satisfaz  $\hat{\beta}^{o*} = \sqrt{1 + \lambda_2} \hat{\beta}^r(\lambda_2)$ , isto é, é exatamente o valor obtido pela reparametrização da estimativa Ridge da regressão original  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ .

A estimativa LASSO para a regressão  $\mathbf{y}^* = \mathbf{X}^*\boldsymbol{\beta}^* + \boldsymbol{\varepsilon}$  é obtida minimizando a Lagrangeana

$$L(\gamma, \boldsymbol{\beta}^*) = \|\mathbf{y} - \mathbf{X}^*\boldsymbol{\beta}^*\|^2 + \gamma \|\boldsymbol{\beta}^*\|_1,$$

isto é,  $\boldsymbol{\beta}^{\text{lasso}*} = \arg \min_{\boldsymbol{\beta}^*} L(\gamma, \boldsymbol{\beta}^*)$ , em que  $\gamma = \frac{\lambda_1}{\sqrt{1 + \lambda_2}}$ .

Como

$$\begin{aligned} \|\mathbf{y}^* - \mathbf{X}^*\boldsymbol{\beta}^*\|^2 + \gamma \|\boldsymbol{\beta}^*\|_1 &= \left\| \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix} - \frac{1}{\sqrt{1 + \lambda_2}} \begin{pmatrix} \mathbf{X} \\ \sqrt{\lambda_2} \mathbf{I} \end{pmatrix} \sqrt{1 + \lambda_2} \boldsymbol{\beta} \right\|^2 \\ &\quad + \gamma \left\| \sqrt{1 + \lambda_2} \boldsymbol{\beta} \right\|_1 \\ &= \left\| \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix} - \frac{1}{\sqrt{1 + \lambda_2}} \begin{pmatrix} \mathbf{X}\boldsymbol{\beta} \\ \sqrt{\lambda_2} \boldsymbol{\beta} \end{pmatrix} \right\|^2 + \sqrt{1 + \lambda_2} \gamma \|\boldsymbol{\beta}\|_1 \\ &= \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda_2 \|\boldsymbol{\beta}\|^2 + \sqrt{1 + \lambda_2} \gamma \|\boldsymbol{\beta}\|_1 \\ &= \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda_2 \|\boldsymbol{\beta}\|^2 + \lambda_1 \|\boldsymbol{\beta}\|_1, \end{aligned}$$

tem-se que o estimador LASSO  $\hat{\boldsymbol{\beta}}^{\text{lasso}*}$  é exatamente o estimador  $\hat{\boldsymbol{\beta}}_{eni}$  pois são definidos pela minimização da mesma função. Em termos dos parâmetros originais,  $\hat{\boldsymbol{\beta}}_{eni} = \frac{1}{\sqrt{1 + \lambda_2}} \hat{\boldsymbol{\beta}}^{\text{lasso}*}$ . Desta forma, o cálculo do estimador *elastic net* ingênuo pode ser obtido como o estimador LASSO de um sistema aumentado.

Para exemplificar, vamos calcular a estimativa *elastic net* ingênuo para o caso ortogonal  $\mathbf{X}'\mathbf{X} = \mathbf{I}$ . Conforme demonstrado (fórmula 4.1), o estimador LASSO para uma regressão ortogonal é dado por

$$\left(\hat{\boldsymbol{\beta}}^{\text{lasso}}\right)_j = \text{sign} \left(\hat{\boldsymbol{\beta}}^{\text{ols}}\right)_j \left( \left| \left(\hat{\boldsymbol{\beta}}^{\text{ols}}\right)_j \right| - \gamma \right)^+.$$

A regressão com os dados aumentados satisfaz

$$\begin{aligned}
 X^* &= \frac{1}{\sqrt{1+\lambda_2}} \begin{pmatrix} X \\ \sqrt{\lambda_2} I \end{pmatrix} \Rightarrow \\
 (X^*)'(X^*) &= \frac{1}{1+\lambda_2} (X', \sqrt{\lambda_2} I) \begin{pmatrix} X \\ \sqrt{\lambda_2} I \end{pmatrix} \\
 &= \frac{1}{1+\lambda_2} (X'X + \lambda_2 I) \\
 &= \frac{1}{1+\lambda_2} (I + \lambda_2 I) \\
 &= I.
 \end{aligned}$$

Portanto o sistema aumentado também é ortogonal. O LASSO do sistema aumentado é

$$\left( \hat{\beta}^{\text{lasso}*} \right)_j = \left( \left| \left( \hat{\beta}^{\text{ols}*} \right)_j \right| - \gamma \right)^+ \text{ sinal} \left( \hat{\beta}^{\text{ols}*} \right)_j.$$

Como demonstrado,

$$\begin{aligned}
 \hat{\beta}^{\text{ols}*} &= \sqrt{1+\lambda_2} \hat{\beta}_R(\lambda_2) \\
 &= \sqrt{1+\lambda_2} (X'X + \lambda_2 I)^{-1} X'y \\
 &= \sqrt{1+\lambda_2} ((1+\lambda_2)I)^{-1} X'y \\
 &= \frac{1}{\sqrt{1+\lambda_2}} X'y.
 \end{aligned}$$

Mas  $X'y$  é o estimador de quadrados mínimos da regressão original e portanto

$\hat{\beta}^{\text{ols}*} = \frac{1}{\sqrt{1+\lambda_2}} \hat{\beta}^{\text{ols}}$ . Como o estimador *elastic net* ingênuo satisfaz  $\hat{\beta}_{eni} = \frac{1}{\sqrt{1+\lambda_2}} \hat{\beta}^{\text{lasso}*}$ , então

$$\begin{aligned}
 \left( \hat{\beta}_{eni} \right)_j &= \frac{1}{\sqrt{1+\lambda_2}} \left( \hat{\beta}^{\text{lasso}*} \right)_j \\
 &= \frac{1}{\sqrt{1+\lambda_2}} \left( \left| \left( \hat{\beta}^{\text{ols}*} \right)_j \right| - \gamma \right)^+ \text{ sinal} \left( \hat{\beta}^{\text{ols}*} \right)_j
 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\sqrt{1+\lambda_2}} \left( \left| \frac{1}{\sqrt{1+\lambda_2}} (\hat{\beta}^{\text{ols}})_j \right| - \gamma \right)^+ \text{sinal}(\hat{\beta}^{\text{ols}})_j \\
&= \frac{\left( \left| (\hat{\beta}^{\text{ols}})_j \right| - \sqrt{(1+\lambda_2)\gamma} \right)^+}{1+\lambda_2} \text{sinal}(\hat{\beta}^{\text{ols}})_j \\
&= \frac{\left( \left| (\hat{\beta}^{\text{ols}})_j \right| - \frac{\lambda_1}{2} \right)^+}{1+\lambda_2} \text{sinal}(\hat{\beta}^{\text{ols}})_j,
\end{aligned}$$

Zou (2005, p. 305).

A grande vantagem do método *elasticnet* em relação ao método LASSO é que se tem a propriedade de efeito de grupo, isto é, variáveis altamente correlacionadas têm alta probabilidade de gerar estimativas relacionadas a estas covariáveis com valores próximos. Esta propriedade decorre do fato que a penalidade empregada no método *elasticnet* define um modelo  $K_p$  estritamente convexo, ao contrário do método LASSO que define um modelo  $K_p$  apenas convexo.

A diferença entre uma penalização convexa e uma estritamente convexa é explicitada na Proposição 5. Seja  $J(\cdot)$  uma função positiva e o problema de minimização

$$\hat{\beta} = \arg \min_{\beta} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda J(\beta) \right\}.$$

Tem-se

**Proposição 5.** Se  $\mathbf{x}_i = \mathbf{x}_j$ ,  $i, j \in \{1, \dots, p\}$ ,

1. Se  $J(\cdot)$  é estritamente convexa, então  $\hat{\beta}_i = \hat{\beta}_j$ ,  $\forall \lambda > 0$ .

2. Se  $J(\beta) = \|\beta\|_1$ , portanto convexa mas não estritamente convexa, então

$$\hat{\beta}_i \hat{\beta}_j \geq 0 \text{ e}$$

$$\hat{\beta}_k^* = \begin{cases} \hat{\beta}_k, & k \neq i, j \\ \left( \hat{\beta}_i + \hat{\beta}_j \right) s, & k = i \\ \left( \hat{\beta}_i + \hat{\beta}_j \right) (1 - s), & k = j \end{cases}$$

é uma outra solução do problema de minimização para  $s \in [0, 1]$ .

*Demonstração.* 1) Se  $J$  é estritamente convexa, então  $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda J(\boldsymbol{\beta})$  também o é e, portanto, admite um único ponto de mínimo  $\hat{\boldsymbol{\beta}}$ . Seja agora o caso em que  $\mathbf{x}_i = \mathbf{x}_j$ . Definindo um novo vetor  $\hat{\boldsymbol{\beta}}^*$  da forma

$$\hat{\beta}_k^* = \begin{cases} \hat{\beta}_k & \text{se } k \neq i, j \\ \frac{1}{2}(\hat{\beta}_i + \hat{\beta}_j) & \text{se } k = i, j \end{cases}.$$

Como as  $i$ -ésima e  $j$ -ésima colunas de  $\mathbf{X}$  são iguais, tem-se que  $\mathbf{X}\hat{\boldsymbol{\beta}}^* = \mathbf{X}\hat{\boldsymbol{\beta}}$ , de onde segue que  $\hat{\boldsymbol{\beta}}^*$  também é um mínimo e pela unicidade  $\hat{\boldsymbol{\beta}}^* = \hat{\boldsymbol{\beta}}$ , o que implica que  $\hat{\beta}_i = \hat{\beta}_j$ .

2) Suponha  $\hat{\beta}_i \hat{\beta}_j < 0$  e considere novamente o vetor  $\hat{\boldsymbol{\beta}}^*$ . Assim,

$$\begin{aligned} (\hat{\beta}_j^*)^2 &= (\hat{\beta}_i^*)^2 = \hat{\beta}_i^* \hat{\beta}_j^* = \frac{1}{2}(\hat{\beta}_i + \hat{\beta}_j) \frac{1}{2}(\hat{\beta}_i + \hat{\beta}_j) \\ &= \frac{1}{4}(\hat{\beta}_i + \hat{\beta}_j)^2 = \frac{1}{4}(\hat{\beta}_i^2 + 2\hat{\beta}_i \hat{\beta}_j + \hat{\beta}_j^2) \leq \frac{1}{4}(\hat{\beta}_i^2 + \hat{\beta}_j^2). \end{aligned}$$

Portanto,

$$\|\hat{\boldsymbol{\beta}}^*\|_1 = \sum_{s=1}^p |\hat{\beta}_s^*| = \sum_{s \neq i, j} |\hat{\beta}_s^*| + |\hat{\beta}_i^*| + |\hat{\beta}_j^*| = \sum_{s \neq i, j} |\hat{\beta}_s| + 2|\hat{\beta}_i^*|.$$

Como  $\hat{\beta}_i^* = \frac{1}{2}(\hat{\beta}_i + \hat{\beta}_j)$ , se  $\hat{\beta}_i \hat{\beta}_j < 0$ ,  $\hat{\beta}_i$  e  $\hat{\beta}_j$  possuem sinais opostos e  $|\hat{\beta}_i^*| = \frac{1}{2}|\hat{\beta}_i + \hat{\beta}_j| < \frac{1}{2}(|\hat{\beta}_i| + |\hat{\beta}_j|)$ . Logo,

$$\|\hat{\boldsymbol{\beta}}^*\|_1 < \sum_{s \neq i, j} |\hat{\beta}_s| + 2 \left( \frac{1}{2}(|\hat{\beta}_i| + |\hat{\beta}_j|) \right) = \|\hat{\boldsymbol{\beta}}\|_1.$$

Tal fato não pode ocorrer pois  $\hat{\boldsymbol{\beta}}^*$  e  $\hat{\boldsymbol{\beta}}$  minimizam  $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\| + \lambda \|\boldsymbol{\beta}\|_1$  e portanto  $\hat{\beta}_i \hat{\beta}_j \geq 0$ .

Observando que se define

$$\hat{\beta}_k^* = \begin{cases} \hat{\beta}_k & , k \neq i, j \\ \left( \hat{\beta}_i + \hat{\beta}_j \right) s & , k = i \quad , s \in [0, 1], \\ \left( \hat{\beta}_i + \hat{\beta}_j \right) (1 - s) & , k = j \end{cases}$$

então

$$\begin{aligned} \|\hat{\beta}^*\|_1 &= \sum_{s \neq i, j}^P |\hat{\beta}_s| + |\hat{\beta}_i + \hat{\beta}_j| s + |\hat{\beta}_i + \hat{\beta}_j| (1 - s) \\ &= \sum_{s \neq i, j}^P |\hat{\beta}_s| + |\hat{\beta}_i + \hat{\beta}_j| = \sum_{s \neq i, j}^P |\hat{\beta}_s| + |\hat{\beta}_i| + |\hat{\beta}_j| = \|\hat{\beta}\|_1, \end{aligned}$$

pois  $\hat{\beta}_i$  e  $\hat{\beta}_j$  possuem o mesmo sinal. Portanto  $\hat{\beta}_s^*$  também é solução do problema de minimização.  $\square$

Como o caso  $\mathbf{x}_i = \mathbf{x}_j$  não é de interesse prático, um resultado de maior aplicabilidade vai garantir uma cota para a diferença entre as estimativas  $\left(\hat{\beta}_{eni}\right)_i$  e  $\left(\hat{\beta}_{eni}\right)_j$  em termos da correlação amostral  $\rho = \mathbf{x}'_i \mathbf{x}_j$ .

**Teorema 1.** Se  $\hat{\beta}_{eni}(\lambda_1, \lambda_2)$  é o vetor de estimativas elasticnet ingênua e se

$$\left(\hat{\beta}_{eni}\right)_i(\lambda_1, \lambda_2) \left(\hat{\beta}_{eni}\right)_j(\lambda_1, \lambda_2) > 0,$$

então

$$\frac{1}{\|\mathbf{y}\|_1} \left| \left(\hat{\beta}_{eni}\right)_i(\lambda_1, \lambda_2) - \left(\hat{\beta}_{eni}\right)_j(\lambda_1, \lambda_2) \right| \leq \frac{1}{\lambda_2} \sqrt{2(1 - \rho)}.$$

*Demonstração.* Observando que a diferenciabilidade da função  $L(\lambda_1, \lambda_2, \beta)$  é garantida apenas fora das intersecções com os eixos coordenados, então

$$\left. \frac{\partial}{\partial \beta_k} L(\lambda_1, \lambda_2, \beta) \right|_{\beta = \hat{\beta}(\lambda_1, \lambda_2)} = 0$$

faz sentido se  $\hat{\beta}_k(\lambda_1, \lambda_2) \neq 0$ . Com essa suposição, a derivada da Lagrangeana

em relação a  $\beta_i$  é

$$-2X' \left\{ \mathbf{y} - X\hat{\beta}(\lambda_1, \lambda_2) \right\} + \lambda_1 \text{sinal} \hat{\beta}_i(\lambda_1, \lambda_2) + 2\lambda_2 \hat{\beta}_i(\lambda_1, \lambda_2) = 0.$$

Subtraindo esta equação para  $i$  e  $j$ , se tem

$$\begin{aligned} (\mathbf{x}_j - \mathbf{x}_i)' \left\{ \mathbf{y} - X\hat{\beta}(\lambda_1, \lambda_2) \right\} + \lambda_2 \left( \hat{\beta}_i(\lambda_1, \lambda_2) - \hat{\beta}_j(\lambda_1, \lambda_2) \right) &= 0 \\ \Rightarrow \hat{\beta}_i(\lambda_1, \lambda_2) - \hat{\beta}_j(\lambda_1, \lambda_2) &= \frac{1}{\lambda_2} (\mathbf{x}'_i - \mathbf{x}'_j) \left( \mathbf{y} - X\hat{\beta} \right). \end{aligned} \quad (5.1)$$

Como

$$\begin{aligned} \|\mathbf{x}_i - \mathbf{x}_j\|^2 &= \langle \mathbf{x}_i - \mathbf{x}_j, \mathbf{x}_i - \mathbf{x}_j \rangle \\ &= \|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2 - 2\mathbf{x}'_i \mathbf{x}_j \\ &= 1 + 1 - 2\rho \\ &= 2(1 - \rho). \end{aligned}$$

e também

$$\begin{aligned} L(\lambda_1, \lambda_2, \hat{\beta}) &\leq L(\lambda_1, \lambda_2, \beta = 0) = \|\mathbf{y}\|^2 \\ \Rightarrow \left\| \mathbf{y} - X\hat{\beta} \right\|^2 + \lambda_1 \|\hat{\beta}\|_1 + \lambda_2 \|\hat{\beta}\|^2 &\leq \|\mathbf{y}\|^2, \end{aligned}$$

substituindo na equação 5.1 e utilizando-se do fato de que  $\hat{\beta}_i(\lambda_1, \lambda_2)$  e  $\hat{\beta}_j(\lambda_1, \lambda_2)$  têm o mesmo sinal, tem-se  $\left| \hat{\beta}_i(\lambda_1, \lambda_2) - \hat{\beta}_j(\lambda_1, \lambda_2) \right| \leq \frac{1}{\lambda_2} 2(1 - \rho) \|\mathbf{y}\|^2$ , e portanto

$$D_{\lambda_1, \lambda_2}(i, j) = \frac{\left| \hat{\beta}_i(\lambda_1, \lambda_2) - \hat{\beta}_j(\lambda_1, \lambda_2) \right|}{\|\mathbf{y}\|_1} \leq \frac{1}{\lambda_2} 2(1 - \rho) \frac{\|\mathbf{y}\|^2}{\|\mathbf{y}\|_1} \leq \frac{2(1 - \rho)}{\lambda_2}$$

□

Segue do Teorema 1 que se  $\mathbf{x}_i$  e  $\mathbf{x}_j$  são altamente correlacionados, então a diferença das estimativas dos parâmetros relativos a estas covariáveis é próxima de

zero. Esta propriedade é que dá ao estimador sua propriedade de seleção por grupo, isto é, grupos de variáveis altamente correlacionadas tendem a ter estimativas correspondentes próximas. Esta propriedade é muito interessante em aplicações. O LASSO, em geral, não possui tal propriedade.

Como observado anteriormente, o estimador *elastic net* ingênuo é obtido por um procedimento em dois estágios: uma estimativa RIDGE e em seguida um encolhimento tipo LASSO. Portanto, se tem um duplo encolhimento e estudos de simulação comprovam que tal fato afeta a performance do estimador. Para que se possa corrigir tal situação, o estimador *elastic net* é definido como um reescalonamento do estimador *elastic net* ingênuo  $\hat{\beta}_{en} = (1 + \lambda_2)\hat{\beta}_{eni}$ . O teorema seguinte fornece a forma explícita para o estimador.

**Teorema 2.** *O estimador elastic net satisfaz*

$$\hat{\beta}_{en} = \underset{\beta}{\operatorname{argmin}} \beta' \frac{(X'X + \lambda_2 I)}{1 + \lambda_2} \beta - 2\mathbf{y}'X\beta + \lambda_1 \|\beta\|_1.$$

Observe que a estimativa LASSO é dada por

$$\begin{aligned} \hat{\beta}^{\text{lasso}} &= \underset{\beta}{\operatorname{argmin}} \|\mathbf{y} - X\beta\|^2 + \lambda_1 \|\beta\|_1 \\ &= \underset{\beta}{\operatorname{argmin}} \{ \mathbf{y}'\mathbf{y} - \mathbf{y}'X\beta - (X\beta)'\mathbf{y} + (X\beta)'(X\beta) + \lambda_1 \|\beta\|_1 \} \\ &= \underset{\beta}{\operatorname{argmin}} \{ \mathbf{y}'\mathbf{y} - \beta'X'X\beta - 2\mathbf{y}'X\beta + \lambda_1 \|\beta\|_1 \} \\ &= \underset{\beta}{\operatorname{argmin}} \{ \beta'X'X\beta - 2\mathbf{y}'X\beta + \lambda_1 \|\beta\|_1 \}. \end{aligned}$$

Desta forma fica claro que o Teorema 2 nos garante que o estimador *elastic net* é uma versão estabilizada do estimador LASSO.



*Demonstração.*  $(\mathbf{y}^*, X^*)$  é o sistema aumentado descrito anteriormente e

$$\hat{\boldsymbol{\beta}}_{\text{en}} = \sqrt{(1 + \lambda_2)} \hat{\boldsymbol{\beta}}^* = \sqrt{(1 + \lambda_2)} \underset{\boldsymbol{\beta}^*}{\operatorname{argmin}} \|\mathbf{y}^* - X\boldsymbol{\beta}^*\|^2 + \frac{\lambda_1}{\sqrt{(1 + \lambda_2)}} \|\boldsymbol{\beta}^*\|_1,$$

em que  $\boldsymbol{\beta}^*$  é uma reparametrização dada por  $\boldsymbol{\beta}^* = \sqrt{(1 + \lambda_2)}\boldsymbol{\beta}$ . Portanto,

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{\text{en}} &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left\| \mathbf{y} - X^* \frac{\boldsymbol{\beta}}{\sqrt{1 + \lambda_2}} \right\|^2 + \frac{\lambda_1}{\sqrt{1 + \lambda_2}} \left\| \frac{\boldsymbol{\beta}}{\sqrt{1 + \lambda_2}} \right\|_1 \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left( \mathbf{y}'\mathbf{y} - 2\mathbf{y}'X^* \frac{\boldsymbol{\beta}}{\sqrt{1 + \lambda_2}} + \boldsymbol{\beta}'(X^*)'X^*\boldsymbol{\beta} + \frac{\lambda_1}{1 + \lambda_2} \|\boldsymbol{\beta}\|_1 \right) \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left( \boldsymbol{\beta}'(X^*)'X^*\boldsymbol{\beta} - 2\mathbf{y}'X^* \frac{\boldsymbol{\beta}}{\sqrt{1 + \lambda_2}} + \frac{\lambda_1}{1 + \lambda_2} \|\boldsymbol{\beta}\|_1 \right) \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left( \boldsymbol{\beta}' \frac{(X'X + \lambda_2 I)}{(1 + \lambda_2)^2} \boldsymbol{\beta} - \frac{2\mathbf{y}'X\boldsymbol{\beta}}{1 + \lambda_2} + \frac{\lambda_1}{1 + \lambda_2} \|\boldsymbol{\beta}\|_1 \right) \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left( \boldsymbol{\beta}' \frac{(X'X + \lambda_2 I)}{1 + \lambda_2} \boldsymbol{\beta} - 2\mathbf{y}'X\boldsymbol{\beta} + \lambda_1 \|\boldsymbol{\beta}\|_1 \right). \end{aligned}$$

□

## 6 Regressão de Ângulos Mínimos

Métodos de seleção de modelos, isto é, seleção de covariáveis, são clássicos na teoria dos modelos lineares. Pode-se citar entre eles: *Forward selection*, *Backward selection*, *All subset selection*, e mais recentemente os métodos LASSO e *Forward stagewise linear regression*. Este último deu origem a um método computacionalmente muito mais simples denominado *Least Angle Regression*. Estes métodos podem ser vistos como um aprimoramento de um método mais elementar denominado *Forward stepwise regression*.

### 6.1 Forward Stagewise Regression

O método *Forward stepwise* consiste em primeiramente se escolher dentre as covariáveis, isto é, entre os vetores  $\mathbf{x}_i$  (colunas da matriz  $X$ ) aquele que tem a maior correlação com o vetor de dados  $\mathbf{y}$ . A Figura 41 representa o vetor  $\mathbf{y}$  e os vetores das covariáveis  $\mathbf{x}_1, \dots, \mathbf{x}_p$ .

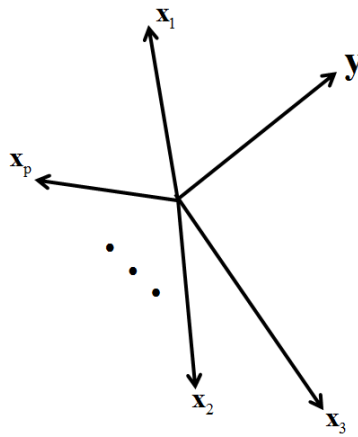


Figura 41 Vetor de dados e das covariáveis.

De forma a simplificar, é possível considerar a matriz  $X$  normalizada, isto é,  $\sum_{i=1}^n x_{ij} = 0$ ,  $\sum_{i=1}^n x_{ij}^2 = 1$  e também  $\sum_{i=1}^n y_i = 0$ . Com isso, os vetores das covariáveis estão representados em uma hipersfera de raio unitário (ver Figura 42).

A correlação amostral é dada por  $\mathbf{x}'_i \mathbf{y} = \|\mathbf{x}_i\| \|\mathbf{y}\| \cos \theta = \|\mathbf{y}\| \cos \theta = P_{\mathbf{x}_i} \mathbf{y}$  (ver representação na Figura 43).

Pode-se supor que o vetor  $\mathbf{y} - P_{\mathbf{x}_i} \mathbf{y}$  é a parte dos dados  $\mathbf{y}$  não explicada pela covariável  $\mathbf{x}_i$ . A ideia então é explicar o vetor  $\mathbf{y} - P_{\mathbf{x}_i} \mathbf{y}$  em termos das outras covariáveis  $\mathbf{x}_j$ ,  $j \neq i$ . Para tanto basta considerar as componentes dos vetores  $\mathbf{x}_j$  que são perpendiculares à  $\mathbf{x}_i$  (ver Figura 44).

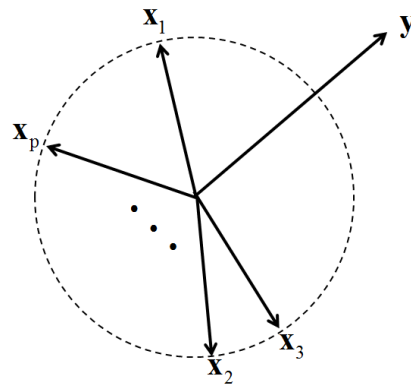


Figura 42 Vetores normalizados das covariáveis.

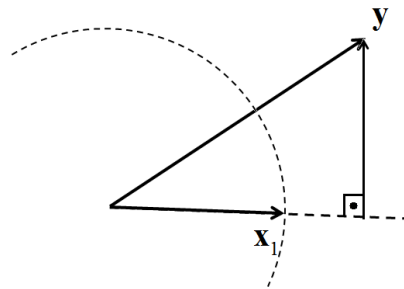


Figura 43 Geometria da correlação amostral.

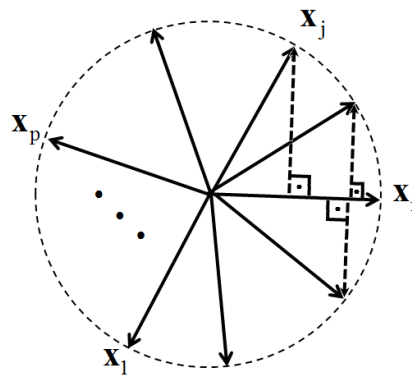


Figura 44 Normalização das covariáveis.

O processo agora é repetido considerando como dados o ortogonal do vetor  $\mathbf{y}$  em relação ao vetor  $\mathbf{x}_i$  e como vetores das covariáveis os ortogonais de  $\mathbf{x}_j$  em relação à  $\mathbf{x}_i$ . Novamente se escolhe a covariável que apresenta maior correlação amostral com  $\mathbf{y} - P_{\mathbf{x}_i}\mathbf{y}$ . Desta forma se tem uma sequência encaixada de modelos. O vetor de regressão é obtido da forma: Supondo por simplicidade que o modelo tenha escolhido na sequência as variáveis  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$ , se teria então

$$\begin{aligned}\mathbf{y} &= P_{\mathbf{x}_1}\mathbf{y} + (\mathbf{y} - P_{\mathbf{x}_1}\mathbf{y}) \\ &= P_{\mathbf{x}_1}\mathbf{y} + P_{\mathbf{x}_2}(\mathbf{y} - P_{\mathbf{x}_1}\mathbf{y}) + (\mathbf{y} - P_{\mathbf{x}_1}\mathbf{y} - (P_{\mathbf{x}_2}(\mathbf{y} - P_{\mathbf{x}_1}\mathbf{y}))).\end{aligned}$$

O processo é repetido e o vetor de regressão é obtido desprezando-se o último resíduo.

Duas covariáveis altamente correlacionadas devem fazer parte do modelo ou pelo princípio da parcimônia é razoável que apenas uma delas esteja presente no modelo? Posteriormente, por algum critério, um destes modelos é então selecionado. Este processo tende a ser agressivo (EFRON et al., 2004) no sentido de que logo na segunda iteração uma covariável  $\mathbf{x}_j$ , altamente correlacionada com  $\mathbf{x}_i$ , não seja escolhida e fique excluída do modelo obtido. Geometricamente tal fato se justifica pois supondo  $\mathbf{x}_j$  altamente correlacionada com  $\mathbf{x}_i$  então os dois vetores estão próximos. Como  $\mathbf{x}_i$  e  $\mathbf{x}_j$  estão próximos, então  $\mathbf{y} - P_{\mathbf{x}_i}\mathbf{y}$  e  $\mathbf{y} - P_{\mathbf{x}_j}\mathbf{y}$  também estarão próximos, mas  $\mathbf{x}_j - P_{\mathbf{x}_i}\mathbf{x}_j$  é um vetor tangente à hipersfera no ponto  $\mathbf{x}_i$ , e portanto não necessariamente muito correlacionado com  $\mathbf{y} - P_{\mathbf{x}_i}\mathbf{y}$  (Figura 45).

Em razão desta deficiência, foi proposto um outro algoritmo de seleção, mais cauteloso, denominado *Forward Stagewise*. A descrição do algoritmo é como se segue. Seja  $\hat{\mu}$  a estimativa de  $E[\mathbf{y}]$  em um passo do algoritmo. No

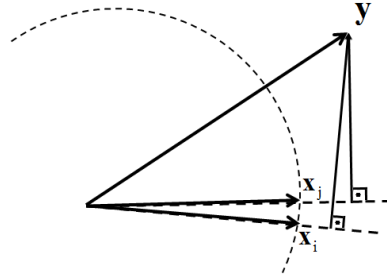


Figura 45 Stepwise e variáveis altamente correlacionadas.

primeiro passo se toma  $\hat{\boldsymbol{\mu}} = 0$ . Considere o vetor de correlações amostrais (de fato proporcional às correlações)  $\hat{\mathbf{c}} = c(\hat{\boldsymbol{\mu}}) = (\hat{c}_1, \dots, \hat{c}_p) = \mathbf{X}'(\mathbf{y} - \hat{\boldsymbol{\mu}}) = (\mathbf{x}'_1(\mathbf{y} - \hat{\boldsymbol{\mu}}), \dots, \mathbf{x}'_p(\mathbf{y} - \hat{\boldsymbol{\mu}}))$  e seja  $\hat{j} = \arg \max |\hat{c}_j|$ . O estimador  $\hat{\boldsymbol{\mu}}$  é então atualizado por

$$\hat{\boldsymbol{\mu}} \leftarrow \hat{\boldsymbol{\mu}} + \varepsilon \text{sinal}(\hat{c}_j) \mathbf{x}_j,$$

em que  $\varepsilon$  é um parâmetro de controle do processo. Observe que não se está fazendo distinção entre positivamente e negativamente correlacionado. Do ponto de vista das aplicações, se quer selecionar covariáveis que estejam altamente correlacionadas, positiva ou negativamente com a variável resposta.

O índice  $\hat{j}$  escolhido é aquele determinado pelo menor ou maior ângulo entre as covariáveis  $\mathbf{x}_i$  relativo à  $\mathbf{y} - \hat{\boldsymbol{\mu}}$ .

No caso em que se toma  $\varepsilon = |\hat{c}_j|$  ocorre que o processo se reduz ao *stepwise*. De fato, no primeiro passo do algoritmo  $\hat{\boldsymbol{\mu}} = 0$  e

$$\hat{\boldsymbol{\mu}} \leftarrow 0 + \varepsilon \text{sinal}(\hat{c}_j) \mathbf{x}_j = |\hat{c}_j| \text{sinal}(\hat{c}_j) \mathbf{x}_j = \hat{c}_j \mathbf{x}_j = \mathbf{x}'_j(\mathbf{y} - 0) \mathbf{x}_j = P_{\mathbf{x}_j} \mathbf{y},$$

que é justamente a projeção de  $\mathbf{y}$  no vetor  $\mathbf{x}_j$  a menos do sinal. No passo seguinte,  $\hat{\mathbf{c}} = \mathbf{x}'(\mathbf{y} - P_{\mathbf{x}_j} \mathbf{y})$  e portanto está se calculando a correlação das covariáveis com um vetor ortogonal à  $\mathbf{x}_j$ , e o processo se reduz ao *Forward Stepwise*. Portanto, o

valor de  $\varepsilon$  deve ser pequeno para que o processo seja capaz de detectar covariáveis muito correlacionadas. Assim, antes de se escolher uma nova covariável, vários passos ocorrem para a mesma covariável. Ver Figura 46.

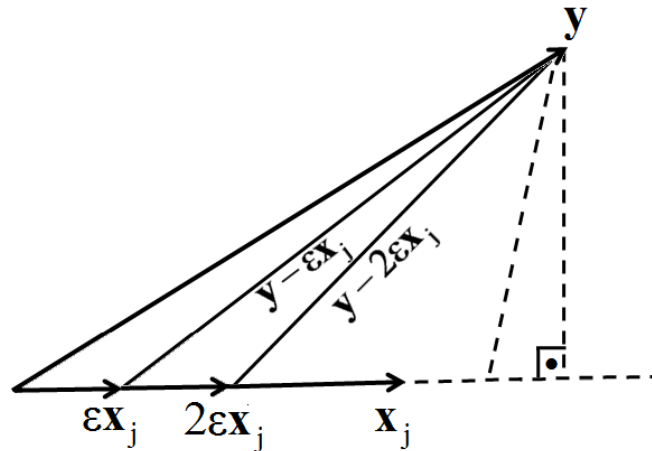


Figura 46 Algoritmo *stagewise*.

## 6.2 O Algoritmo LARS

O algoritmo LARS segue o mesmo procedimento do método *Forward Stagewise*, só que com uma ideia matemática simples, que consiste na escolha adequada para o valor de  $\varepsilon$  que define o passo do algoritmo. O número de passos é drasticamente reduzido igualando-se ao número de covariáveis. Primeiramente vamos descrever o algoritmo para o caso simples de 2 covariáveis  $X = (x_1, x_2)$ .

O vetor de correlações amostrais é

$$\mathbf{c}(\hat{\boldsymbol{\mu}}) = (c_1(\hat{\boldsymbol{\mu}}), c_2(\hat{\boldsymbol{\mu}})) = \mathbf{X}'(\mathbf{y} - \hat{\boldsymbol{\mu}}) = (\mathbf{x}'_1(\mathbf{y} - \hat{\boldsymbol{\mu}}), \mathbf{x}'_2(\mathbf{y} - \hat{\boldsymbol{\mu}})).$$

Inicialmente  $\boldsymbol{\mu}_0 = 0$ . Para fixar a construção suponha que o ângulo entre  $\mathbf{y}$  e  $\mathbf{x}_1$  seja menor do que com  $\mathbf{x}_2$  (Figura 47), isto é,  $c_1(\hat{\boldsymbol{\mu}}_0) > c_2(\hat{\boldsymbol{\mu}}_0)$ . Logo o algoritmo aumenta  $\boldsymbol{\mu}_0$  na direção de  $\mathbf{x}_1$  para  $\hat{\boldsymbol{\mu}}_1 = \hat{\boldsymbol{\mu}}_0 + \hat{\gamma}_1 \mathbf{x}_1$ . A ideia então é escolher

o valor de  $\hat{\gamma}_1$  de forma adequada, tal que  $\mathbf{y} - \hat{\boldsymbol{\mu}}_1$  seja igualmente correlacionado com  $\mathbf{x}_1$  e  $\mathbf{x}_2$ . Isto ocorre se  $\mathbf{y} - \hat{\boldsymbol{\mu}}_1$  é equiangular com  $\mathbf{x}_1$  e  $\mathbf{x}_2$ . Observe que pela ortogonalidade é possível substituir  $\mathbf{y}$  por  $\bar{\mathbf{y}}_2 = P_{\text{Im}(X)}\mathbf{y}$ . Portanto se quer que  $\bar{\mathbf{y}}_2 - \hat{\boldsymbol{\mu}}_1$  seja equiangular com  $\mathbf{x}_1$  e  $\mathbf{x}_2$ , isto é,

$$\mathbf{x}'_1(\bar{\mathbf{y}}_2 - \hat{\gamma}_1\mathbf{x}_1) = \mathbf{x}'_2(\bar{\mathbf{y}}_2 - \hat{\gamma}_1\mathbf{x}_1) \Rightarrow \hat{\gamma}_1 = \frac{(\mathbf{x}'_2 - \mathbf{x}'_1)\bar{\mathbf{y}}_2}{\mathbf{x}'_2\mathbf{x}_1 - \mathbf{x}'_1\mathbf{x}_1}.$$

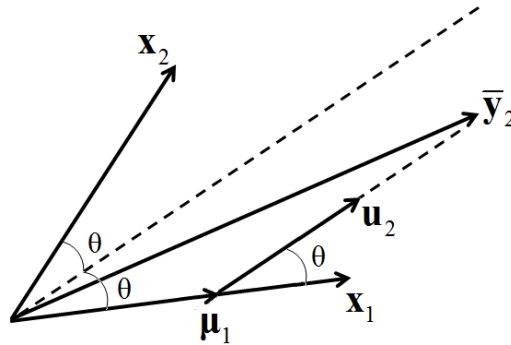


Figura 47 Geometria do LARS para a dimensão 2.

Se  $\mathbf{u}_2$  é o vetor unitário na direção equiangular, o próximo passo do algoritmo atualiza  $\boldsymbol{\mu}$  da forma  $\hat{\boldsymbol{\mu}}_2 = \hat{\boldsymbol{\mu}}_1 + \hat{\gamma}_2\mathbf{u}_2$ . No caso bidimensional,  $\hat{\gamma}_2$  é escolhido tal que  $\hat{\boldsymbol{\mu}}_2 = \bar{\mathbf{y}}_2$ .

Se a dimensão fosse maior do que 2, isto é, com uma nova covariável  $\mathbf{x}_3$  e vetor equiangular  $\mathbf{u}_3$ ,  $\hat{\gamma}_2$  seria menor e se teria uma outra alteração de direção. O valor de  $\hat{\gamma}_2$  é obtido pelas igualdades

$$\mathbf{x}'_1(\bar{\mathbf{y}}_3 - \hat{\gamma}_2\hat{\boldsymbol{\mu}}_2) = \mathbf{x}'_2(\bar{\mathbf{y}}_3 - \hat{\gamma}_2\hat{\boldsymbol{\mu}}_2) = \mathbf{x}'_3(\bar{\mathbf{y}}_3 - \hat{\gamma}_2\hat{\boldsymbol{\mu}}_2),$$

e o valor de  $\hat{\gamma}_3$  é dado por  $\hat{\boldsymbol{\mu}}_3 = \bar{\mathbf{y}}_3 = \hat{\boldsymbol{\mu}}_2 + \hat{\gamma}_3\mathbf{u}_3$ , conforme descrito na Figura 48.

O procedimento algébrico para a obtenção de um vetor  $\mathbf{u}$  em  $\mathbb{R}^n$  equiangular com os vetores colunas  $\mathbf{x}_i$  da matriz  $X$  será descrito com detalhes. Ob-

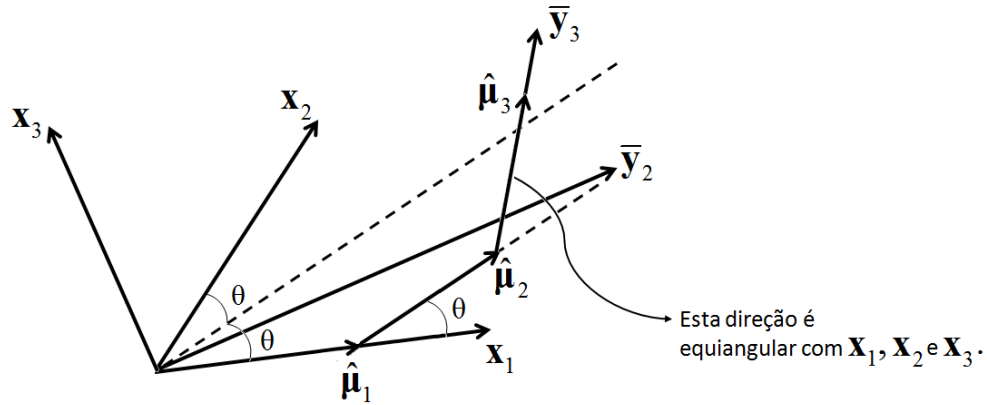


Figura 48 Geometria do LARS para a dimensão 3.

serve que se  $\|x_i\| = 1$  e  $\sum_{i=1}^n x_i = 0$ , o vetor  $X'v = (x'_1v, \dots, x'_pv)'$  é um vetor em que as coordenadas são proporcionais às correlações amostrais. O vetor  $v$  será equiangular se  $X'v = a = (a, \dots, a)' = a\mathbf{1}$ . Vamos tomar o vetor  $(X'X)^{-1}\mathbf{1}$ . A norma deste vetor, com o produto interno natural do espaço paramétrico  $\langle v_1, v_2 \rangle_p = v'_1 (X'X) v_2$  é

$$\left\| (X'X)^{-1}\mathbf{1} \right\|_p^2 = \mathbf{1}'(X'X)^{-1} (X'X) (X'X)^{-1}\mathbf{1} = \mathbf{1}'(X'X)^{-1}\mathbf{1}.$$

Como se quer vetores unitários, seja  $w = \frac{(X'X)^{-1}\mathbf{1}}{(\mathbf{1}'(X'X)^{-1}\mathbf{1})^{1/2}}$ , e  $u = Xw$ . O vetor  $u$  é unitário pois

$$\langle u, u \rangle = (Xw)'(Xw) = w'(X'X)w = 1.$$

O vetor  $u$  é um vetor equiangular em relação a todas as colunas de  $X$  pois

$$X'u = X'Xw = \frac{(X'X)(X'X)^{-1}\mathbf{1}}{\mathbf{1}'(X'X)^{-1}\mathbf{1}} = \frac{1}{\mathbf{1}'(X'X)^{-1}\mathbf{1}}\mathbf{1}.$$

Portanto, em termos de complexidade computacional, a obtenção de vetores equiangulares é da mesma complexidade da multiplicação matricial.



De forma geral, o algoritmo é descrito algebricamente pela construção conforme se segue. Seja  $\mathcal{A} \subset \{1, 2, \dots, p\}$  um subconjunto de índices, os quais indexam as covariáveis. Se  $|\mathcal{A}|$  é o número de elementos de  $\mathcal{A}$ , seja  $X_{\mathcal{A}}$  a matriz  $n \times (|\mathcal{A}|)$  obtida de  $X$  tomando-se as colunas correspondentes às covariáveis indexadas por  $\mathcal{A}$ . Um vetor  $\mathbf{u}_{\mathcal{A}}$  será equiangular em relação às colunas de  $X_{\mathcal{A}}$  se vale a relação  $X'_{\mathcal{A}}\mathbf{u}_{\mathcal{A}} = a\mathbf{1}_{\mathcal{A}}$ , em que  $\mathbf{1}_{\mathcal{A}} = \mathbf{1}_{(|\mathcal{A}|) \times 1} = (1, \dots, 1)'$ . Sejam  $G_{\mathcal{A}} = X'_{\mathcal{A}}X_{\mathcal{A}}$ ,  $(\mathbf{1}'_{\mathcal{A}}G_{\mathcal{A}}^{-1}\mathbf{1}_{\mathcal{A}})^{-\frac{1}{2}} = A_{\mathcal{A}}$ ,  $\mathbf{w}_{\mathcal{A}} = A_{\mathcal{A}}G_{\mathcal{A}}^{-1}\mathbf{1}_{\mathcal{A}}$  e  $\mathbf{u}_{\mathcal{A}} = X_{\mathcal{A}}\mathbf{w}_{\mathcal{A}}$ . Tem-se então que

$$\begin{aligned} X'_{\mathcal{A}}\mathbf{u}_{\mathcal{A}} &= X'_{\mathcal{A}}X_{\mathcal{A}}\mathbf{w}_{\mathcal{A}} \\ &= X'_{\mathcal{A}}X_{\mathcal{A}}(A_{\mathcal{A}}G_{\mathcal{A}}^{-1}\mathbf{1}_{\mathcal{A}}) \\ &= A_{\mathcal{A}}X'_{\mathcal{A}}X_{\mathcal{A}}\left((X'_{\mathcal{A}}X_{\mathcal{A}})^{-1}\mathbf{1}_{\mathcal{A}}\right) \\ &= A_{\mathcal{A}}\mathbf{1}_{\mathcal{A}}, \end{aligned}$$

e portanto  $\mathbf{u}_{\mathcal{A}}$  é equiangular com os vetores coluna de  $X$  definidos pelo subconjunto  $\mathcal{A}$ . Além disso,

$$\begin{aligned} \|\mathbf{u}_{\mathcal{A}}\|^2 &= \langle \mathbf{u}_{\mathcal{A}}, \mathbf{u}_{\mathcal{A}} \rangle \\ &= \langle X_{\mathcal{A}}\mathbf{w}_{\mathcal{A}}, X_{\mathcal{A}}\mathbf{w}_{\mathcal{A}} \rangle \\ &= \mathbf{w}'_{\mathcal{A}}X'_{\mathcal{A}}X_{\mathcal{A}}\mathbf{w}_{\mathcal{A}} \\ &= (A_{\mathcal{A}}G_{\mathcal{A}}^{-1}\mathbf{1}_{\mathcal{A}})'X'_{\mathcal{A}}X_{\mathcal{A}}(A_{\mathcal{A}}G_{\mathcal{A}}^{-1}\mathbf{1}_{\mathcal{A}}) \\ &= A_{\mathcal{A}}^2\mathbf{1}'_{\mathcal{A}}(X'_{\mathcal{A}}X_{\mathcal{A}})^{-1}(X'_{\mathcal{A}}X_{\mathcal{A}})(X'_{\mathcal{A}}X_{\mathcal{A}})^{-1}\mathbf{1}_{\mathcal{A}} \\ &= A_{\mathcal{A}}^2\mathbf{1}'_{\mathcal{A}}(X'_{\mathcal{A}}X_{\mathcal{A}})\mathbf{1}_{\mathcal{A}} \\ &= \left[ \left( \mathbf{1}'_{\mathcal{A}}(X'_{\mathcal{A}}X_{\mathcal{A}})^{-1}\mathbf{1}_{\mathcal{A}} \right)^{-\frac{1}{2}} \right]^2 \mathbf{1}'_{\mathcal{A}}(X'_{\mathcal{A}}X_{\mathcal{A}})\mathbf{1}_{\mathcal{A}} \\ &= \mathbf{1}. \end{aligned}$$

Tem-se portanto um procedimento algébrico para a obtenção de um vetor que seja equiangular à um subconjunto qualquer de colunas da matriz  $X$ .

Pode-se agora construir o algoritmo LARS. No primeiro passo considera-se  $\hat{\boldsymbol{\mu}}_0 = \mathbf{0}$ . O vetor de correlações amostrais é dado por

$$\mathbf{c}(\hat{\boldsymbol{\mu}}_0) = (c_1(\hat{\boldsymbol{\mu}}_0), \dots, c_n(\hat{\boldsymbol{\mu}}_0)) = \mathbf{X}'(\mathbf{y} - \hat{\boldsymbol{\mu}}_0) = \mathbf{X}'\mathbf{y}.$$

A ideia é se tomar, assim como no *Step Forward*, a covariável que define a coluna de maior correlação amostral com o vetor de dados  $\mathbf{y}$ . Seja  $\hat{C} = \max_{j \in \{1, \dots, p\}} \{|c_j(\hat{\boldsymbol{\mu}}_0)|\}$ . Se  $\tilde{j}$  é o índice correspondente à covariável de correlação amostral máxima em módulo, então toma-se inicialmente  $\mathcal{A} = \{\tilde{j}\}$ ,  $\mathbf{X}_{\mathcal{A}} = \mathbf{X}_{\tilde{j}}$ ,  $\mathbf{G}_{\mathcal{A}} = \mathbf{X}'_{\tilde{j}}\mathbf{X}_{\tilde{j}} = 1$ ,  $\mathbf{A}_{\mathcal{A}} = \left(\mathbf{1}'(\mathbf{X}'_{\tilde{j}}\mathbf{X}_{\tilde{j}})^{-1}\mathbf{1}\right)^{-\frac{1}{2}} = (\mathbf{1}'\mathbf{1})^{-\frac{1}{2}} = p^{-\frac{1}{2}}$ ,  $\mathbf{w}_{\mathcal{A}} = \mathbf{A}_{\mathcal{A}}(\mathbf{G}_{\mathcal{A}})^{-1}\mathbf{1}_{\mathcal{A}} = \frac{1}{\sqrt{p}}\mathbf{1}_{\mathcal{A}}$ ,  $\mathbf{u}_{\mathcal{A}} = \mathbf{X}_{\mathcal{A}}\mathbf{w}_{\mathcal{A}} = \mathbf{x}_{\tilde{j}}\frac{1}{\sqrt{p}}\mathbf{1}_{\mathcal{A}} = \frac{1}{\sqrt{p}}\mathbf{x}_{\tilde{j}}$  e  $\hat{\boldsymbol{\mu}}_1 = \hat{\boldsymbol{\mu}}_0 + \hat{\gamma}\mathbf{u}_{\mathcal{A}} = \hat{\gamma}\frac{1}{\sqrt{p}}\mathbf{x}_{\tilde{j}}$ .

O valor de  $\hat{\gamma}$  é escolhido da seguinte forma: considere a curva  $\mu(\gamma) = \gamma\mathbf{u}_{\mathcal{A}} = \gamma\frac{1}{\sqrt{p}}\mathbf{x}_{\tilde{j}}$ . Tem-se então as curvas de correlação  $c_j(\gamma) = \mathbf{x}'_j\left(\mathbf{y} - \gamma\frac{1}{\sqrt{p}}\mathbf{x}_{\tilde{j}}\right)$ , para  $j = 1, \dots, n$ . Para o índice  $\tilde{j}$ , se a correlação é positiva, o valor  $c_{\tilde{j}}(\gamma)$  é decrescente, pois

$$c_{\tilde{j}}(\gamma) = \mathbf{x}'_{\tilde{j}}\left(\mathbf{y} - \gamma\frac{1}{\sqrt{p}}\mathbf{x}_{\tilde{j}}\right) = \mathbf{x}'_{\tilde{j}}\mathbf{y} - \mathbf{x}'_{\tilde{j}}\gamma\frac{1}{\sqrt{p}}\mathbf{x}_{\tilde{j}} = \hat{C} - \frac{\gamma}{\sqrt{p}}.$$

Toma-se então como  $\hat{\gamma}$  o menor valor de  $\gamma$  para o qual existe outra covariável de índice  $\tilde{i}$  tal que  $c_{\tilde{i}}(\gamma) = c_{\tilde{j}}(\gamma)$ . Tal procedimento será melhor explicado no passo genérico a seguir.

Suponha agora que se esteja em um passo do algoritmo onde já se definiu um conjunto  $\mathcal{A}$  de índices. Tem-se então computados  $\hat{\boldsymbol{\mu}}_{\mathcal{A}}$ ,  $\mathbf{X}_{\mathcal{A}}$ ,  $\mathbf{A}_{\mathcal{A}}$ ,  $\mathbf{G}_{\mathcal{A}}$  e o vetor equiangular  $\mathbf{u}_{\mathcal{A}}$ . Vale destacar que quando se expressa  $\hat{\boldsymbol{\mu}}_{\mathcal{A}}$  no estágio em que o conjunto de índices ativo é  $\mathcal{A}$ , o vetor  $\hat{\boldsymbol{\mu}}_{\mathcal{A}}$  foi obtido no passo anterior e portanto está em um subespaço de dimensão  $|\mathcal{A}| - 1$ . Assim, o vetor  $\mathbf{u}_{\mathcal{A}}$ , que é o vetor equiangular com todos os vetores  $s_j\mathbf{x}_j$  para  $j \in \mathcal{A}$ , não está no mesmo subespaço que  $\hat{\boldsymbol{\mu}}_{\mathcal{A}}$ . Neste estágio do algoritmo, o vetor de correlações amostrais é

$c(\hat{\boldsymbol{\mu}}_{\mathcal{A}}) = X'(\mathbf{y} - \hat{\boldsymbol{\mu}}_{\mathcal{A}})$  e para os índices em  $\mathcal{A}$  possui o mesmo valor máximo em módulo, que foi denominado  $\hat{C}$ , isto é, com  $\hat{C} \geq |c_j(\hat{\boldsymbol{\mu}}_{\mathcal{A}})|, \forall j = 1, \dots, n$ .

No próximo passo, um novo índice é acrescentado ao conjunto de índice  $\mathcal{A}$ , definindo o conjunto  $\mathcal{A}+$  e o novo vetor de estimativas  $\hat{\boldsymbol{\mu}}_{\mathcal{A}+} = \hat{\boldsymbol{\mu}}_{\mathcal{A}} + \hat{\gamma}\mathbf{u}_{\mathcal{A}+}$ . O valor  $\hat{\gamma}$  é obtido pela relação

$$\hat{\gamma} = \min^+ \left\{ \frac{\hat{C} - \hat{c}_j}{A_{\mathcal{A}} - a_j}, \frac{\hat{C} + \hat{c}_j}{A_{\mathcal{A}+} + a_j} \right\}, j \notin \mathcal{A},$$

em que  $\min^+$  é o mínimo tomado apenas para valores positivos. Esta escolha do valor  $\hat{\gamma}$  se justifica da forma como se segue. Se  $\boldsymbol{\mu}(\gamma) = \hat{\boldsymbol{\mu}}_{\mathcal{A}} + \gamma\mathbf{u}_{\mathcal{A}}$ , então o vetor de correlações é  $c(\gamma) = X'(\mathbf{y} - \boldsymbol{\mu}(\gamma))$ , o que implica em

$$\begin{aligned} c_j(\gamma) &= \mathbf{x}'_j(\mathbf{y} - \boldsymbol{\mu}(\gamma)) \\ &= \mathbf{x}'_j(\mathbf{y} - \hat{\boldsymbol{\mu}}_{\mathcal{A}} - \gamma\mathbf{u}_{\mathcal{A}}) \\ &= \mathbf{x}'_j(\mathbf{y} - \hat{\boldsymbol{\mu}}_{\mathcal{A}}) - \gamma\mathbf{x}'_j\mathbf{u}_{\mathcal{A}}. \end{aligned}$$

No caso em que  $j \in \mathcal{A}$ ,  $\mathbf{x}'_j(\mathbf{y} - \hat{\boldsymbol{\mu}}_{\mathcal{A}}) = \hat{C}$  e  $\mathbf{x}'_j\mathbf{u}_{\mathcal{A}} = a_j = A_{\mathcal{A}}$ , ou seja,

$$c_j(\gamma) = \hat{C} - \gamma a_j = \hat{C} - \gamma A_{\mathcal{A}}.$$

Portanto as correlações correspondentes aos índices em  $\mathcal{A}$  decaem de forma linear.

O valor  $\hat{\gamma}$  é obtido quando este valor  $c - \gamma A_{\mathcal{A}}$  se iguala pela primeira vez a  $c_j(\gamma) = \hat{c}_j - \gamma a_j$ , para algum  $j$ , com  $j \notin \mathcal{A}$ , e portanto,

$$\hat{C} - \hat{\gamma}A_{\mathcal{A}} = \hat{c}_j - \hat{\gamma}a_j$$

$$\hat{\gamma}(A_{\mathcal{A}} - a_j) = \hat{C} - \hat{c}_j$$

$$\hat{\gamma} = \frac{\hat{C} - \hat{c}_j}{A_{\mathcal{A}} - a_j}.$$

É necessário também considerar o caso  $-c_j(\gamma)$  relativo à  $-\mathbf{x}_j$ . Assim,  $a_j = \mathbf{x}'_j\boldsymbol{\mu}_{\mathcal{A}}$

troca de sinal e

$$\hat{C} - \hat{\gamma}A_{\mathcal{A}} = -\hat{c}_j + \hat{\gamma}a_j$$

$$\hat{\gamma} = \frac{\hat{C} + \hat{c}_j}{A_{\mathcal{A}} + a_j}.$$

Portanto,  $\hat{\gamma}$  é o menor valor tal que um novo índice  $\hat{j}$  é incluído em  $\mathcal{A}$ ,  $\mathcal{A}_+ = \mathcal{A} \cup \{\hat{j}\}$  de tal forma que a nova correlação máxima em valor absoluto é  $\hat{C}_+ = \hat{C} - \hat{\gamma}A_{\mathcal{A}}$ .

### 6.3 O algoritmo LARS aplicado à delineamentos ortogonais

Para o caso em que  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , com  $\mathbf{X}'\mathbf{X} = \mathbf{I}$ , o algoritmo LARS fica bastante simplificado. Neste caso, as colunas da matriz  $\mathbf{X}$  formam em  $\mathbb{R}^n$  um conjunto com  $p$  vetores ortonormais (Figura 49).

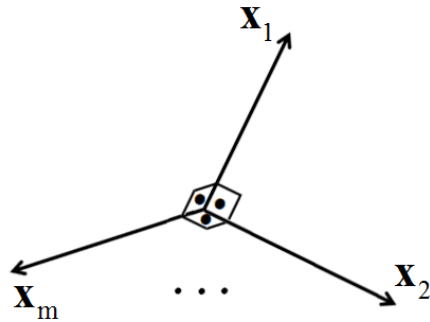


Figura 49 Vetores coluna para  $\mathbf{X}$  ortonormais.

Por meio de uma transformação linear ortogonal  $\mathbf{A}$  de  $\mathbb{R}^n$ , isto é, uma transformação dos dados, tem-se  $\mathbb{R}^p \xrightarrow{\mathbf{X}} \mathbb{R}^n \xrightarrow{\mathbf{A}} \mathbb{R}^n$ ,  $\mathbf{A} \cdot \mathbf{X} : \mathbb{R}^p \mapsto \mathbb{R}^n$  tal que  $\mathbf{A} \cdot \mathbf{X}(\mathbf{e}_i) = \mathbf{e}_i$  e portanto a imagem de  $\mathbf{A} \cdot \mathbf{X}$  é a inclusão canônica de  $\mathbb{R}^p$  em  $\mathbb{R}^n$ . Com esta observação pode-se supor, sem perda de generalidade, que as colunas da matriz  $\mathbf{X}$  são da forma  $\mathbf{x}_i = \mathbf{e}_i$ ,  $i = 1, \dots, p$ . O vetor de dados  $\mathbf{y}' = (y_1, \dots, y_n)$  se projeta em  $\mathbb{R}^p$  como  $(P_{\mathbb{R}^p} \mathbf{y})' = (y_1, \dots, y_p, 0, \dots, 0)$ , e portanto

os valores  $y_{p+1}, y_{p+2}, \dots, y_n$  não afetarão o processo de estimação e podem ser considerados nulos. Em outras palavras, todo o processo de estimação ocorre em  $\mathbb{R}^p$  e a matriz  $X$  é a identidade. Desta forma tem-se para um conjunto de índices  $\mathcal{A}$  em um passo do algoritmo que  $G_{\mathcal{A}} = X'_{\mathcal{A}}X_{\mathcal{A}} = I_{\mathcal{A}}$ . Também se tem que

$$A_{\mathcal{A}} = (\mathbf{1}'_{\mathcal{A}}G_{\mathcal{A}}\mathbf{1}_{\mathcal{A}})^{-1/2} = (\mathbf{1}'_{\mathcal{A}}I_{\mathcal{A}}\mathbf{1}_{\mathcal{A}})^{-1/2} = (\mathbf{1}'_{\mathcal{A}}\mathbf{1}_{\mathcal{A}})^{-1/2} = \frac{1}{\sqrt{|\mathcal{A}|}},$$

$$\mathbf{w}_{\mathcal{A}} = A_{\mathcal{A}}G_{\mathcal{A}}^{-1}\mathbf{1}_{\mathcal{A}} = \frac{1}{\sqrt{|\mathcal{A}|}}I_{\mathcal{A}}^{-1}\mathbf{1}_{\mathcal{A}} = \frac{1}{\sqrt{|\mathcal{A}|}}I_{\mathcal{A}}\mathbf{1}_{\mathcal{A}} = \frac{1}{\sqrt{|\mathcal{A}|}}\mathbf{1}_{\mathcal{A}},$$

e

$$\mathbf{u}_{\mathcal{A}} = X_{\mathcal{A}}\mathbf{w}_{\mathcal{A}} = \frac{1}{\sqrt{|\mathcal{A}|}}X_{\mathcal{A}}\mathbf{1}_{\mathcal{A}} = \frac{1}{\sqrt{|\mathcal{A}|}}\mathbf{1}_{\mathcal{A}}.$$

O vetor de correlações é o próprio vetor  $\mathbf{y}$ , uma vez que  $X'\mathbf{y} = I'\mathbf{y} = \mathbf{y}$ . Como temos que tomar no primeiro passo o índice da covariável de maior correlação com  $\mathbf{y}$ , basta então tomar o índice relativo a maior coordenada do vetor  $\mathbf{y}$ , a menos do sinal. É conveniente utilizar a notação de estatística de ordem. Considere então as estatísticas de ordem tomadas em relação aos valores absolutos das coordenadas. Para simplificar, será suposto que empates não ocorrem. Assim,  $|\mathbf{y}|_{(1)} > |\mathbf{y}|_{(2)} > \dots > |\mathbf{y}|_{(p)} \geq 0$ . Portanto no primeiro passo do algoritmo é escolhida a covariável de índice  $j(1)$  relativo à coordenada  $|\mathbf{y}|_{(1)}$ , isto é,  $\mathcal{A} = \{j(1)\}$ . Logo,

$$\hat{\boldsymbol{\mu}}_1 = \hat{\boldsymbol{\mu}}_0 + \hat{\gamma}_1\mathbf{u}_{\mathcal{A}} = \mathbf{0} + \hat{\gamma}_1\frac{1}{\sqrt{|\mathcal{A}|}}\mathbf{1}_{\mathcal{A}} = \hat{\gamma}_1\mathbf{e}_{j(1)}.$$

O valor de  $\hat{\gamma}_1$  é obtido tal que este seja o menor valor de tal forma que a correlação entre  $\mathbf{x}_{j(1)}$  e  $\mathbf{y} - \hat{\gamma}_1\mathbf{x}_{j(1)}$  seja igual à correlação de outra covariável  $\mathbf{x}_s$ , isto é,

$$\mathbf{x}'_{j(1)}(\mathbf{y} - \hat{\gamma}_1\mathbf{x}_{j(1)}) = |\mathbf{y}|_{(1)} - \hat{\gamma}_1 = \mathbf{x}'_s(\mathbf{y} - \hat{\gamma}_1\mathbf{x}_{j(1)}) = \mathbf{x}'_s\mathbf{y} = \mathbf{y}_s.$$

Logo,  $|\mathbf{y}|_{(1)} - \hat{\gamma}_1 = \mathbf{y}_s \Rightarrow |\mathbf{y}|_{(1)} - \mathbf{y}_s = \hat{\gamma}_1$ . O menor valor ocorre para o maior valor possível de  $\mathbf{y}_s$  que é  $|\mathbf{y}|_{(2)}$ . Portanto  $\hat{\gamma}_1 = |\mathbf{y}|_{(1)} - |\mathbf{y}|_{(2)}$  e o índice

$j(2)$  relativo a  $|\mathbf{y}|_{(2)}$  é adicionado ao conjunto de índices ativos obtendo-se  $\mathcal{A} = \{j(1), j(2)\}$ . Para se evitar o emprego de indução um passo a mais no algoritmo será desenvolvido deixando claro como é a forma geral.

O vetor unitário equiangular em relação a  $\mathbf{e}_{j(1)}$  e  $\mathbf{e}_{j(2)}$  é

$$\mathbf{u}_{\mathcal{A}} = \frac{1}{\sqrt{2}} \mathbf{1}_{\mathcal{A}} = \frac{1}{\sqrt{2}} (0, \dots, 1, 0, \dots, 1, 0, \dots, 0),$$

com 1 nas  $j(1)$  e  $j(2)$  coordenadas. A estimativa LARS é atualizada para  $\hat{\boldsymbol{\mu}}_2 = \hat{\boldsymbol{\mu}}_1 + \hat{\gamma}_2 \frac{\mathbf{1}_{\mathcal{A}}}{\sqrt{2}}$ . Novamente o valor de  $\hat{\gamma}_2$  é o menor valor para o qual uma outra covariável  $\mathbf{x}_s$  tem a mesma correlação de  $\mathbf{x}_{j(1)}$  e  $\mathbf{x}_{j(2)}$  com  $\mathbf{y} - \hat{\boldsymbol{\mu}}_2$ . Como  $\mathbf{x}_{j(1)} = \mathbf{e}_{j(1)}$ ,  $\mathbf{x}_{j(2)} = \mathbf{e}_{j(2)}$  e  $\mathbf{x}_s = \mathbf{e}_s$ , tem-se

$$\begin{aligned} \mathbf{e}'_{j(1)}(\mathbf{y} - \hat{\boldsymbol{\mu}}_2) &= \mathbf{e}'_{j(2)}(\mathbf{y} - \hat{\boldsymbol{\mu}}_2) = \mathbf{e}'_s(\mathbf{y} - \hat{\boldsymbol{\mu}}_2) \\ \Rightarrow \mathbf{e}'_{j(1)}\mathbf{y} - \mathbf{e}'_{j(1)}\hat{\boldsymbol{\mu}}_2 &= \mathbf{e}'_{j(2)}\mathbf{y} - \mathbf{e}'_{j(2)}\hat{\boldsymbol{\mu}}_2 = \mathbf{e}'_s\mathbf{y} = \mathbf{y}_s \\ \Rightarrow |\mathbf{y}|_{(1)} - \mathbf{e}'_{j(1)} \left( \hat{\gamma}_1 \mathbf{e}_{j(1)} + \hat{\gamma}_2 \frac{\mathbf{1}_{\mathcal{A}}}{\sqrt{2}} \right) &= |\mathbf{y}|_{(2)} - \mathbf{e}'_{j(2)} \left( \hat{\gamma}_1 \mathbf{e}_{j(1)} + \hat{\gamma}_2 \frac{\mathbf{1}_{\mathcal{A}}}{\sqrt{2}} \right) = \mathbf{y}_s \\ \Rightarrow |\mathbf{y}|_{(1)} - \hat{\gamma}_1 - \frac{\hat{\gamma}_2}{\sqrt{2}} &= |\mathbf{y}|_{(2)} - \frac{\hat{\gamma}_2}{\sqrt{2}} = \mathbf{y}_s \\ \Rightarrow |\mathbf{y}|_{(1)} - \left( |\mathbf{y}|_{(1)} - |\mathbf{y}|_{(2)} \right) - \frac{\hat{\gamma}_2}{\sqrt{2}} &= |\mathbf{y}|_{(2)} - \frac{\hat{\gamma}_2}{\sqrt{2}} = \mathbf{y}_s \\ \Rightarrow |\mathbf{y}|_{(2)} - \frac{\hat{\gamma}_2}{\sqrt{2}} &= \mathbf{y}_s \\ \Rightarrow |\mathbf{y}|_{(2)} - \mathbf{y}_s &= \frac{\hat{\gamma}_2}{\sqrt{2}}. \end{aligned}$$

Portanto o menor valor de  $\hat{\gamma}_2$  ocorre quando  $\mathbf{y}_s$  é a coordenada relativa a  $|\mathbf{y}|_{(3)}$  e se tem  $\hat{\gamma}_2 = \sqrt{2} \left( |\mathbf{y}|_{(2)} - |\mathbf{y}|_{(3)} \right)$ . Desta forma, fica claro como evolui o algoritmo.

Observe que

$$\hat{\boldsymbol{\mu}}_2 = \hat{\boldsymbol{\mu}}_1 + \hat{\gamma}_2 \mathbf{u}_{\mathcal{A}} = \hat{\gamma}_1 \mathbf{e}_{j(1)} + \hat{\gamma}_2 \mathbf{1}_{\mathcal{A}}$$

$$\begin{aligned}
&= \left( |\mathbf{y}|_{(1)} - |\mathbf{y}|_{(2)} \right) \mathbf{e}_{j(1)} + \frac{1}{\sqrt{2}} (0, \dots, \hat{\gamma}_2, 0, \dots, \hat{\gamma}_2, 0, \dots, 0) \\
&= \left( |\mathbf{y}|_{(1)} - |\mathbf{y}|_{(2)} \right) \mathbf{e}_{j(1)} + \hat{\gamma}_2 \frac{\mathbf{e}_{j(1)}}{\sqrt{2}} + \hat{\gamma}_2 \frac{\mathbf{e}_{j(2)}}{\sqrt{2}} \\
&= \left( |\mathbf{y}|_{(1)} - |\mathbf{y}|_{(2)} \right) \mathbf{e}_{j(1)} + \left( |\mathbf{y}|_{(2)} - |\mathbf{y}|_{(3)} \right) \mathbf{e}_{j(1)} + \left( |\mathbf{y}|_{(3)} - |\mathbf{y}|_{(2)} \right) \mathbf{e}_{j(2)} \\
&= \left( |\mathbf{y}|_{(1)} - |\mathbf{y}|_{(3)} \right) \mathbf{e}_{j(1)} + \left( |\mathbf{y}|_{(2)} - |\mathbf{y}|_{(3)} \right) \mathbf{e}_{j(2)}.
\end{aligned}$$

Como  $\mathbf{y}_{j(1)} = |\mathbf{y}|_{(1)} \mathbf{e}_{j(1)}$  e  $\mathbf{y}_{j(2)} = |\mathbf{y}|_{(2)} \mathbf{e}_{j(2)}$ , então

$$\hat{\boldsymbol{\mu}}_2 = \left( \mathbf{y}_{j(1)} - |\mathbf{y}|_{(3)} \right) \mathbf{e}_{j(1)} + \left( \mathbf{y}_{j(2)} - |\mathbf{y}|_{(3)} \right) \mathbf{e}_{j(2)}.$$

Portanto, fica demonstrado o lema:

**Lema 1.** *Para um delineamento ortogonal com  $\mathbf{x}_j = \mathbf{e}_j$ ,  $j = 1, \dots, n$ , a  $k$ -ésima estimativa LARS ( $0 \leq k \leq n$ ) é dada por*

$$(\hat{\boldsymbol{\mu}}_k(\mathbf{y}))_i = \begin{cases} y_i - |\mathbf{y}|_{(k+1)} & \text{se } y_i > |\mathbf{y}|_{(k+1)} \\ 0 & \text{se } |\mathbf{y}| \leq |\mathbf{y}|_{(k+1)} \\ y_i + |\mathbf{y}|_{(k+1)} & \text{se } y_i < -|\mathbf{y}|_{(k+1)}. \end{cases}$$

#### 6.4 Relação entre LARS e OLS

É possível se dar uma descrição do algoritmo LARS em termos de estimadores de quadrados mínimos. Suponha que o algoritmo tenha completado o  $k - 1$  passo e portanto com a estimativa  $\hat{\boldsymbol{\mu}}_{k-1}$  e iniciando o  $k$ -ésimo passo. O conjunto de índices ativo é  $\mathcal{A}_k$  (foi acrescentado o índice  $k$  para numerar o passo do algoritmo). Calcula-se então  $\mathbf{X}_k$ ,  $\mathbf{G}_k$ ,  $\mathbf{A}_k$  e  $\boldsymbol{\mu}_k$  (o sub-índice  $\mathcal{A}$  foi substituído por  $k$ ). Seja  $\bar{\mathbf{y}}_k$  a projeção de  $\mathbf{y}$  no subespaço gerado pelas colunas de  $\mathbf{X}$  indexadas por  $\mathcal{A}_k$ , ( $\mathcal{L}(\mathbf{X}_k)$ ), isto é,  $\bar{\mathbf{y}}_k = \mathbf{X}_k(\mathbf{X}'_k \mathbf{X}_k)^{-1} \mathbf{X}'_k \mathbf{y}$  e portanto  $\bar{\mathbf{y}}_k$  define o estimador de quadrados mínimos quando se considera apenas as covariáveis definidas por  $\mathcal{A}_k$ . Observe que  $\hat{\boldsymbol{\mu}}_{k-1} \in \mathcal{L}(\mathbf{X}_{k-1})$ . Pode-se então escrever, observando que

$X_k G_k^{-1} X_k' \hat{\boldsymbol{\mu}}_{k-1} = \hat{\boldsymbol{\mu}}_{k-1}$  pois  $X_k G_k^{-1} X_k' = X_k (X_k' X_k)^{-1} X_k'$  é um projetor no espaço  $\mathcal{L}(X)$  e  $\hat{\boldsymbol{\mu}}_{k-1} \in \mathcal{L}(X_{k-1})$ , que

$$\begin{aligned}
 \bar{\mathbf{y}}_k &= \hat{\boldsymbol{\mu}}_{k-1} + \bar{\mathbf{y}}_k - \hat{\boldsymbol{\mu}}_{k-1} \\
 &= \hat{\boldsymbol{\mu}}_{k-1} + X_k G_k^{-1} X_k' (\mathbf{y} - \hat{\boldsymbol{\mu}}_{k-1}) \\
 &= \hat{\boldsymbol{\mu}}_{k-1} + X_k G_k^{-1} \hat{\mathbf{C}} \mathbf{1}_{\mathcal{A}} \\
 &= \hat{\boldsymbol{\mu}}_{k-1} + \frac{1}{A_{\mathcal{A}}} X_k A_{\mathcal{A}} G_k^{-1} \hat{\mathbf{C}} \mathbf{1}_{\mathcal{A}} \\
 &= \hat{\boldsymbol{\mu}}_{k-1} + \frac{1}{A_{\mathcal{A}}} \hat{\mathbf{C}} X_k A_{\mathcal{A}} G_k^{-1} \mathbf{1}_{\mathcal{A}} \\
 &= \hat{\boldsymbol{\mu}}_{k-1} + \frac{\hat{\mathbf{C}}}{A_{\mathcal{A}}} X_k \mathbf{w}_{\mathcal{A}} \\
 &= \hat{\boldsymbol{\mu}}_{k-1} + \frac{\hat{\mathbf{C}}}{A_{\mathcal{A}}} \mathbf{u}_k,
 \end{aligned}$$

e portanto,

$$\|\bar{\mathbf{y}}_k - \hat{\boldsymbol{\mu}}_{k-1}\| = \left\| \frac{\hat{\mathbf{C}}}{A_{\mathcal{A}}} \mathbf{u}_k \right\| = \frac{\hat{\mathbf{C}}}{A_{\mathcal{A}}} = \bar{\gamma}.$$

Como  $\hat{\boldsymbol{\mu}}_k = \hat{\boldsymbol{\mu}}_{k-1} + \hat{\gamma} \mathbf{u}_k$  e  $\bar{\mathbf{y}}_k = \hat{\boldsymbol{\mu}}_{k-1} + \bar{\gamma} \mathbf{u}_k$ , os vetores  $\hat{\boldsymbol{\mu}}_k - \hat{\boldsymbol{\mu}}_{k-1}$  e  $\bar{\mathbf{y}}_k - \hat{\boldsymbol{\mu}}_{k-1}$  são colineares. Prova-se que  $\hat{\gamma} < \bar{\gamma}$  (Figura 50). Ao se atingir  $\mathcal{A}_n = \{1, \dots, n\}$ , pela definição do algoritmo  $\hat{\boldsymbol{\mu}}_n = \bar{\mathbf{y}}_n = \hat{\boldsymbol{\mu}}_{\text{ols}}$ .

## 6.5 Relação entre o algoritmo LARS e o método LASSO

Uma propriedade surpreendente do algoritmo LARS é que com pequenas modificações o algoritmo é capaz de computar as estimativas obtidas pelo método LASSO. Observa-se que em relação ao parâmetro  $t$  do método LASSO, as estimativas  $\hat{\boldsymbol{\beta}}^{\text{lasso}}(t)$  definem um caminho que se inicia em  $\boldsymbol{\beta} = 0$  e termina no estimador de quadrados mínimos  $\hat{\boldsymbol{\beta}}^{\text{ols}}$ . Tem-se que este caminho no espaço de parâmetros define também um caminho no espaço de dados  $\hat{\boldsymbol{\mu}}(t) = X \hat{\boldsymbol{\beta}}^{\text{lasso}}(t)$ . De forma semelhante, o algoritmo LARS define um caminho no espaço de dados formado por



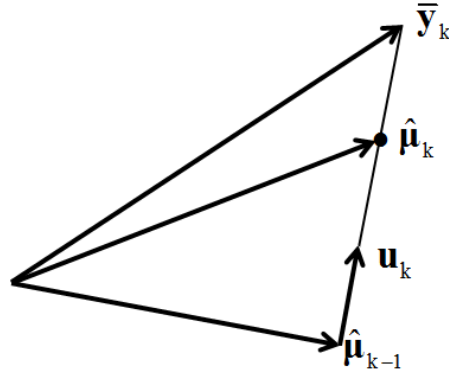


Figura 50 Relação LARS e OLS.

segmentos de retas, no  $k$ -ésimo passo dado por  $\hat{\mu}_k(\gamma) = \hat{\mu}_{k-1} + \gamma(\bar{y}_k - \hat{\mu}_{k-1})$ . Com as devidas modificações, as duas curvas serão as mesmas. As modificações são introduzidas a partir de uma longa sequência de lemas. Estes são bastante técnicos e neste trabalho serão abordados apenas alguns deles. Como referência completa, tem-se o artigo original (EFRON et al., 2004). Os três primeiros lemas são propriedades gerais do LARS.

**Lema 2.** *Suponha que ao fim do  $k - 1$  passo do algoritmo LARS a covariável definida digamos pela coluna  $\mathbf{x}_k$  seja adicionada ao conjunto de índices  $\mathcal{A}_{k-1}$  (sem perda de generalidade as covariáveis serão enumeradas por  $1, 2, 3, \dots$  na ordem em que estas são incluídas nos passos consecutivos do algoritmo). O vetor  $\mathbf{w}_k = A_k(G_k)^{-1}\mathbf{1}_k$  que define o vetor equiangular  $\mathbf{u}_k = X_k\mathbf{w}_k$  possui sua  $k$ -ésima componente  $w_{kk}$  com o mesmo sinal da correlação  $c_{kk} = \mathbf{x}'_k(\mathbf{y} - \hat{\mu}_{k-1})$ . Além disso, o vetor  $\hat{\beta}_k$ , com  $\hat{\mu}_k = X\hat{\beta}_k$ , tem sua  $k$ -ésima componente  $\hat{\beta}_{kk}$  também com o mesmo sinal.*

*Demonstração.* Como  $\mathbf{X}'_k (\mathbf{y} - \hat{\boldsymbol{\mu}}_{k-1}) = \hat{\mathbf{C}}_k \mathbf{1}_k$ , então

$$\begin{aligned} \mathbf{w}_k &= A_k (G_k)^{-1} \mathbf{1}_k \\ &= A_k (G_k)^{-1} (\hat{\mathbf{C}}_k)^{-1} \mathbf{X}'_k (\mathbf{y} - \hat{\boldsymbol{\mu}}_{k-1}) \\ &= A_k (\hat{\mathbf{C}}_k)^{-1} (\mathbf{X}'_k \mathbf{X}_k)^{-1} \mathbf{X}'_k (\mathbf{y} - \hat{\boldsymbol{\mu}}_{k-1}) \\ &= A_k (\hat{\mathbf{C}}_k)^{-1} \mathbf{w}_k^*. \end{aligned}$$

Tem-se que  $\mathbf{X}'_k (\mathbf{y} - \bar{\mathbf{y}}_{k-1}) = \mathbf{X}'_k (\bar{\mathbf{y}}_k - \bar{\mathbf{y}}_{k-1}) = (\underbrace{0}_{k-1}, \underbrace{\delta}_1)'$  pois  $\bar{\mathbf{y}}_k$  e  $\bar{\mathbf{y}}_{k-1}$  possuem as primeiras  $k-1$  coordenadas iguais, isto é,  $\mathbf{X}'_{k-1} (\bar{\mathbf{y}}_k - \bar{\mathbf{y}}_{k-1}) = 0$ . Considere as curvas  $c_i(\gamma) = \mathbf{x}'_i (\mathbf{y} - \gamma \mathbf{u}_{k-1})$ ,  $i = 1, \dots, k$ . Pelo fato de o vetor  $\mathbf{u}_{k-1}$  ser equiangular com as covariáveis  $\mathbf{x}_i$ ,  $i = 1, \dots, k-1$ , tem-se

$$c_k(\gamma) \begin{cases} < c_j(\gamma) & \gamma < \hat{\gamma}_{k-1} \\ = c_j(\gamma) = \hat{c}_k & \gamma = \hat{\gamma}_{k-1} \\ > c_j(\gamma) & \hat{\gamma}_{k-1} < \gamma < \bar{\gamma}_{k-1} \end{cases} \quad j < k.$$

Portanto,

$$\begin{aligned} \mathbf{w}_k^* &= (\mathbf{X}'_k \mathbf{X}_k)^{-1} \mathbf{X}'_k (\bar{\mathbf{y}}_k - \bar{\mathbf{y}}_{k-1} + \bar{\mathbf{y}}_{k-1} - \hat{\boldsymbol{\mu}}_{k-1}) \\ &= (\mathbf{X}'_k \mathbf{X}_k)^{-1} \begin{pmatrix} 0 \\ \delta \end{pmatrix} + (\mathbf{X}'_k \mathbf{X}_k)^{-1} (\bar{\mathbf{y}}_{k-1} - \hat{\boldsymbol{\mu}}_{k-1}) \\ &= (\mathbf{X}'_k \mathbf{X}_k)^{-1} \begin{pmatrix} 0 \\ \delta \end{pmatrix} + 0. \end{aligned}$$

$(\mathbf{X}'_k \mathbf{X}_k)$  é positiva definida e portanto  $(0, \delta)(\mathbf{X}'_k \mathbf{X}_k)^{-1} \begin{pmatrix} 0 \\ \delta \end{pmatrix} = a_{kk} \delta^2 > 0$ , em

que  $a_{kk}$  é a entrada correspondente de  $(X'_k X_k)^{-1} \Rightarrow a_{kk} > 0$ . Assim,

$$\mathbf{w}_k^* = (X'_k X_k)^{-1} \begin{pmatrix} 0 \\ \delta \end{pmatrix} = \begin{pmatrix} * \\ * \\ \vdots \\ a_{kk} \delta \end{pmatrix}$$

e  $a_{kk} \delta > 0$ . Portanto, a  $k$ -ésima componente de  $\mathbf{w}_k = A_k \hat{c}_k^{-1} \mathbf{w}_k^*$  tem o mesmo sinal de  $\hat{c}_k^{-1} = c_{kk}$ .  $\square$

O segundo lema interpreta a quantidade  $A_{\mathcal{A}} = \left( \mathbf{1}'_{\mathcal{A}} (G_{\mathcal{A}})^{-1} \mathbf{1}_{\mathcal{A}} \right)^{-\frac{1}{2}}$ . Seja  $\mathcal{S}_{\mathcal{A}}$  o simplexo estendido

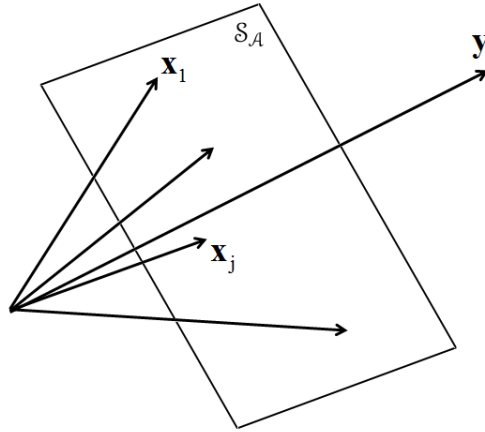
$$\mathcal{S}_{\mathcal{A}} = \left\{ \mathbf{v} = \sum_{j \in \mathcal{A}} s_j \mathbf{x}_j P_j ; \sum_{j \in \mathcal{A}} P_j = 1 \right\},$$

em que estendido significa que  $P_j$  pode ser negativo. O simplexo estendido  $\mathcal{S}_{\mathcal{A}}$  está contido em um hiperplano que passa pelos pontos definidos pelos vetores  $\{s_j \mathbf{x}_j, j \in \mathcal{A}\}$ . Recordando que  $s_j = \pm 1$  definidos de forma que  $s_j \mathbf{x}_j$  tenha um ângulo menor do que  $90^\circ$  com o vetor de dados  $\mathbf{y}$ , tal situação é representada na Figura 51.

**Lema 3.** *O ponto de  $\mathcal{S}_{\mathcal{A}}$  mais próximo à origem é definido pelo vetor  $\mathbf{v}_{\mathcal{A}} = A_{\mathcal{A}} \mathbf{u}_{\mathcal{A}} = A_{\mathcal{A}} X_{\mathcal{A}} \mathbf{w}_{\mathcal{A}}$  em que  $\mathbf{w}_{\mathcal{A}} = A_{\mathcal{A}} (G_{\mathcal{A}})^{-1} \mathbf{1}_{\mathcal{A}}$  e  $\|\mathbf{v}_{\mathcal{A}}\| = A_{\mathcal{A}}$ . Se  $\mathcal{A} \subseteq \mathcal{B}$  então  $A_{\mathcal{A}} \geq A_{\mathcal{B}}$ , sendo que o maior valor possível é  $A_{\mathcal{A}} = 1$  para um  $\mathcal{A}$  com apenas um ponto.*

*Demonstração.* Considere o vetor  $\mathbf{v} = X_{\mathcal{A}} \mathbf{P}$ , em que  $\mathbf{P} = (p_1, \dots, p_{|\mathcal{A}|})$  e portanto

$$\|\mathbf{v}\|^2 = (X_{\mathcal{A}} \mathbf{P})' (X_{\mathcal{A}} \mathbf{P}) = \mathbf{P}' X'_{\mathcal{A}} X_{\mathcal{A}} \mathbf{P} = \mathbf{P}' G_{\mathcal{A}} \mathbf{P}.$$

Figura 51 Geometria do conjunto  $S_A$ 

Tem-se o problema de minimizar  $P'G_AP$  sujeito à restrição  $\sum_{j \in A} P_j = 1$  que pode ser descrito como  $\mathbf{1}'_A P = 1$ . Pelo método dos multiplicadores de Lagrange, tem-se a lagrangiana  $P'G_AP - \lambda(\mathbf{1}'_A P - 1)$ . Derivando em relação à  $P$  tem-se que o mínimo ocorre para  $P_A = \lambda(G_A)^{-1}\mathbf{1}_A$ . Como a soma das componentes de  $P_A$  é 1, tem-se então

$$1 = \mathbf{1}'_A P_A = \mathbf{1}'_A \lambda(G_A)^{-1}\mathbf{1}_A = \lambda \mathbf{1}'_A (G_A)^{-1}\mathbf{1}_A = \lambda A^{-2} \Rightarrow \lambda = A_A^2,$$

o que implica em

$$P_A = A_A^2 (G_A)^{-1}\mathbf{1}_A = A_A A_A (G_A)^{-1}\mathbf{1}_A = A_A \mathbf{w}_A.$$

Logo, como  $\mathbf{v}_A = X_A P_A = \sum_{j \in A} s_j \mathbf{x}'_j P_j$ , tem-se  $\mathbf{v}_A \in S_A$ , e

$$\|\mathbf{v}_A\|^2 = P'_A (G_A)^{-1} P_A = A_A \mathbf{1}_A (G_A)^{-1}\mathbf{1}_A = A_A^2.$$

Se  $A \subseteq B \Rightarrow S_A \subseteq S_B$ , o ponto mais próximo à origem de  $S_B$  deve ser menor do que o ponto mais próximo de  $S_A$  implicando então que  $A_A \geq A_B$ .  $\square$

Como o algoritmo LARS evolui por segmentos de retas, as alterações são realizadas em cada segmento linear.

Considere o caminho  $\beta(\gamma) = \beta + \gamma\mathbf{d}$ , em que  $\beta$  e  $\mathbf{d}$  são vetores em  $\mathbb{R}^p$  e  $S(\gamma) = \|\mathbf{y} - \mathbf{X}\beta(\gamma)\|^2$ . Segue então o lema

**Lema 4.** Se  $\mathbf{c} = \mathbf{X}'(\mathbf{y} - \mathbf{X}\beta)$  é o vetor de correlações no passo em que  $\mu = \mathbf{X}\beta$ , então  $S(\gamma) - S(0) = -2(\mathbf{c}'\mathbf{d})\gamma + \mathbf{d}'(\mathbf{X}'\mathbf{X})\mathbf{d}\gamma^2$ .

*Demonstração.*

$$\begin{aligned}
S(\gamma) - S(0) &= (\mathbf{y} - \mathbf{X}\beta(\gamma))'(\mathbf{y} - \mathbf{X}\beta(\gamma)) - (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \\
&= (\mathbf{y} - \mathbf{X}(\beta + \gamma\mathbf{d}))'(\mathbf{y} - \mathbf{X}(\beta + \gamma\mathbf{d})) - (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \\
&= (\mathbf{y} - \mathbf{X}\beta - \mathbf{X}\gamma\mathbf{d})'(\mathbf{y} - \mathbf{X}\beta - \mathbf{X}\gamma\mathbf{d}) - (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \\
&= ((\mathbf{y} - \mathbf{X}\beta)' - (\mathbf{X}\gamma\mathbf{d})')((\mathbf{y} - \mathbf{X}\beta) - (\mathbf{X}\gamma\mathbf{d})) \\
&\quad - (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \\
&= (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) - (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{X}\gamma\mathbf{d}) \\
&\quad - (\mathbf{X}\gamma\mathbf{d})'(\mathbf{y} - \mathbf{X}\beta) + (\mathbf{X}\gamma\mathbf{d})'(\mathbf{X}\gamma\mathbf{d}) - (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \\
&= -(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{X}\gamma\mathbf{d}) - (\mathbf{X}\gamma\mathbf{d})'(\mathbf{y} - \mathbf{X}\beta) + (\mathbf{X}\gamma\mathbf{d})'(\mathbf{X}\gamma\mathbf{d}) \\
&= -(\mathbf{y} - \mathbf{X}\beta)(\mathbf{X}\gamma\mathbf{d})' - (\mathbf{X}\gamma\mathbf{d})'(\mathbf{y} - \mathbf{X}\beta) + (\mathbf{X}\gamma\mathbf{d})'(\mathbf{X}\gamma\mathbf{d}) \\
&= -2(\mathbf{X}\gamma\mathbf{d})'(\mathbf{y} - \mathbf{X}\beta) + (\mathbf{X}\gamma\mathbf{d})'(\mathbf{X}\gamma\mathbf{d}) \\
&= -2(\mathbf{X}\mathbf{d})'(\mathbf{y} - \mathbf{X}\beta)\gamma + (\mathbf{X}\gamma\mathbf{d})'(\mathbf{X}\gamma\mathbf{d}) \\
&= -2\mathbf{d}'\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta)\gamma + (\mathbf{X}\mathbf{d})'(\mathbf{X}\mathbf{d})\gamma^2 \\
&= -2\mathbf{d}'\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta)\gamma + \mathbf{d}'\mathbf{X}'\mathbf{X}\mathbf{d}\gamma^2 \\
&= -2\mathbf{d}'\mathbf{c}\gamma + \mathbf{d}'\mathbf{X}'\mathbf{X}\mathbf{d}\gamma^2.
\end{aligned}$$

□

Com as propriedades do algoritmo LARS expressas nos três lemas anteri-

ores e com uma sequência longa de lemas técnicos, é possível demonstrar que pelo algoritmo LARS se pode obter todas as estimativas LASSO, conforme em Efron et al. (2004, p. 11).

## 7 Estimação do Erro de Predição

Quando se ajusta a um conjunto de dados  $\mathbf{y}$  (vetor de dimensão  $n \times 1$ ) um modelo  $\mu(\mathbf{y})$  uma pergunta natural é o de se medir o quanto este modelo é adequado para prever um novo conjunto de dados. Deseja-se então, ao menos, se ter um estimador para o erro de predição. A medida usual da soma dos quadrados dos resíduos do modelo  $\|\mathbf{y} - \mu(\mathbf{y})\|^2$  é certamente otimista, uma vez que o modelo foi justamente ajustado com os dados observados  $\mathbf{y}$ . Tal questão pode ser rigorosamente colocada no sentido em que o que se quer é: Os dados  $\mathbf{y}$  são aleatórios, o que implica que  $\mu(\mathbf{y})$  também o é. Novos dados  $\mathbf{y}^0$  também são aleatórios. O quadrado do erro de predição  $\|\mathbf{y}^0 - \mu(\mathbf{y})\|^2$  também é aleatório. Um critério para se abordar este erro é tomar a esperança do quadrado do erro de predição. Observa-se que este erro depende de duas variáveis aleatórias independentes e igualmente distribuídas, e é dado pela esperança dupla  $E_0 \left[ E \left[ \|\mathbf{y}^0 - \mu(\mathbf{y})\|^2 \right] \right]$ .

### 7.1 Erro de predição para modelos lineares

O caso mais simples ocorre quando o modelo é linear e o estimador é não viesado,  $\mu(\mathbf{y}) = M\mathbf{y}$ , com  $\boldsymbol{\mu} = E[M\mathbf{y}] = M\boldsymbol{\mu}$ . Assim,

$$\begin{aligned} c &= E \left[ E_0 \left[ \langle \mathbf{y}^0 - M\mathbf{y}, \mathbf{y}^0 - M\mathbf{y} \rangle \right] \right] \\ &= E \left[ E_0 \left[ \langle \mathbf{y}^0, \mathbf{y}^0 \rangle - 2 \langle \mathbf{y}^0, M\mathbf{y} \rangle + \langle M\mathbf{y}, M\mathbf{y} \rangle \right] \right] \end{aligned}$$

$$\begin{aligned}
&= E \left[ E_0 \left[ \langle \mathbf{y}^0, \mathbf{y}^0 \rangle \right] - 2 \langle E_0 \left[ \mathbf{y}^0 \right], \mathbf{My} \rangle + \langle \mathbf{My}, \mathbf{My} \rangle \right] \\
&= E \left[ \|\boldsymbol{\mu}\|^2 + n\sigma^2 - 2 \langle \boldsymbol{\mu}, \mathbf{My} \rangle + \langle \mathbf{My}, \mathbf{My} \rangle \right] \\
&= \|\boldsymbol{\mu}\|^2 + n\sigma^2 - 2 \langle \boldsymbol{\mu}, E[\mathbf{My}] \rangle + E[\langle \mathbf{My}, \mathbf{My} \rangle] \\
&= \|\boldsymbol{\mu}\|^2 + n\sigma^2 - 2 \langle \boldsymbol{\mu}, \mathbf{M}\boldsymbol{\mu} \rangle + E[\langle \mathbf{My}, \mathbf{My} \rangle] \\
&= \|\boldsymbol{\mu}\|^2 + n\sigma^2 - 2 \langle \boldsymbol{\mu}, \mathbf{M}\boldsymbol{\mu} \rangle + \text{tr}(\sigma^2 \mathbf{M}'\mathbf{M}) + \|\mathbf{M}\boldsymbol{\mu}\|^2.
\end{aligned}$$

Mas como

$$\begin{aligned}
E \left[ \|\mathbf{y} - \mathbf{My}\|^2 \right] &= E[\langle \mathbf{y} - \mathbf{My}, \mathbf{y} - \mathbf{My} \rangle] \\
&= E[\langle \mathbf{y}, \mathbf{y} \rangle] - 2E[\langle \mathbf{y}, \mathbf{My} \rangle] + E[\langle \mathbf{My}, \mathbf{My} \rangle] \\
&= \|\boldsymbol{\mu}\|^2 + n\sigma^2 - 2E[\langle \mathbf{y}, \mathbf{My} \rangle] + E[\mathbf{y}'\mathbf{M}'\mathbf{My}] \\
&= \|\boldsymbol{\mu}\|^2 + n\sigma^2 - 2 \{ \sigma^2 \text{tr}(\mathbf{M}) + \boldsymbol{\mu}'\mathbf{M}\boldsymbol{\mu} \} \\
&\quad + \sigma^2 \text{tr}(\mathbf{M}'\mathbf{M}) + \|\mathbf{M}\boldsymbol{\mu}\|^2,
\end{aligned}$$

então  $E \left[ E_0 \left[ \|\mathbf{y}^0 - \mathbf{My}\|^2 \right] \right] = E \left[ \|\mathbf{y} - \mathbf{My}\|^2 \right] + 2\sigma^2 \text{tr}(\mathbf{M})$ , de onde segue que um estimador não viesado para a esperança do quadrado do erro do modelo  $\hat{\mu} = \mathbf{My}$  é dado por  $\|\mathbf{y} - \mathbf{My}\|^2 + 2\sigma^2 \text{tr}(\mathbf{M})$ . Uma vez que se tem a soma de resíduos observada do modelo e uma penalização dada por  $2\sigma^2 \text{tr}(\mathbf{M})^2$ , tal penalização admite uma interpretação em termos de covariância, uma vez que

$$\begin{aligned}
\text{cov}(\mathbf{y}, \mathbf{My}) &= E[(\mathbf{y} - \boldsymbol{\mu})'(\mathbf{My} - \boldsymbol{\mu})] \\
&= E[(\mathbf{y} - \boldsymbol{\mu})'(\mathbf{My})] \\
&= E[(\mathbf{y})'(\mathbf{My})] - E[\boldsymbol{\mu}'\mathbf{My}] \\
&= \boldsymbol{\mu}'\mathbf{M}\boldsymbol{\mu} + \sigma^2 \text{tr}(\mathbf{M}) - \boldsymbol{\mu}'\mathbf{M}\boldsymbol{\mu} \\
&= \sigma^2 \text{tr}(\mathbf{M}).
\end{aligned}$$

Desta forma um estimador não viesado para o erro de predição é dado por  $\|\mathbf{y} - M\mathbf{y}\|^2 + 2\text{cov}(\mathbf{y}, M\mathbf{y})$ . Claramente este estimador não é operacional pois depende de um parâmetro populacional.

Em geral, os modelos lineares são dados por matrizes de projeção em um subespaço de dimensão  $p$ , que é justamente a dimensão do modelo. Neste caso,  $\text{tr}(M) = p$  e o estimador fica da forma  $\|\mathbf{y} - M\mathbf{y}\|^2 + 2p\sigma^2$ . O estimador da esperança do quadrado do erro de predição evidentemente pode ser usado como uma ferramenta de seleção de modelos. De fato, o estimador não viesado obtido acima está intimamente relacionado à famosa fórmula do  $C_p$  de Mallows para seleção de modelos. Segue então uma demonstração de tal fórmula utilizando argumentos geométricos.

## 7.2 Estimação do erro de predição para o caso geral

No clássico artigo de Efron (2004), é apresentado um estimador não viesado para o erro de predição de modelos de regressão.

Considere  $\mathbf{y}$  um vetor aleatório cujas coordenadas sejam as observações de um dado experimento. Se  $\boldsymbol{\mu} = E[\mathbf{y}]$ , seja  $\boldsymbol{\mu}(\mathbf{y})$  um estimador para a média  $\boldsymbol{\mu}$  do vetor  $\mathbf{y}$ . Uma questão fundamental é mensurar a qualidade do estimador  $\boldsymbol{\mu}(\mathbf{y})$  em medir futuras observações  $\mathbf{y}^0$  geradas pelo mesmo mecanismo aleatório que gerou  $\mathbf{y}$ . Neste sentido, pode-se definir o erro de predição como o valor esperado do seguinte procedimento:

Observa-se os dados  $\mathbf{y} \Rightarrow$  Obtém-se a estimativa da média  $\boldsymbol{\mu}(\mathbf{y}) \Rightarrow$  Novos dados  $\mathbf{y}^0$  do mesmo fenômeno são observados  $\Rightarrow$  Calcula-se o quadrado do desvio  $\|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2$ .

Portanto, o erro de predição é o parâmetro populacional (uma vez que é dado por uma esperança) definido pela esperança dupla  $E \left[ E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] \right]$ ,



em que  $E[\cdot]$  é a esperança em relação ao vetor aleatório  $\mathbf{y}$  e  $E_0[\cdot]$  representa a esperança em relação ao novo vetor aleatório  $\mathbf{y}^0$ . Um estimador não viesado para o erro de predição será obtido a partir da seguinte construção. Considere os triângulos conforme Figura 52.

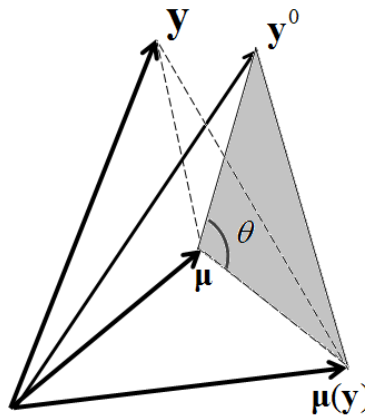


Figura 52 Geometria do erro de predição

Na Figura 52, é importante destacar novamente que o vetor  $\mathbf{y}$  é o vetor cujas coordenadas são os valores observados, o vetor  $\mathbf{y}^0$  é um vetor cujas coordenadas representam os valores de um novo experimento realizado e o vetor  $\boldsymbol{\mu}$  é o vetor média populacional. Finalmente,  $\boldsymbol{\mu}(\mathbf{y})$  é o vetor ajustado ao modelo, ou seja, é a projeção de  $\mathbf{y}$  no modelo.

Aplicando-se a lei dos cossenos ao triângulo hachurado da Figura 52 tem-se:

$$\begin{aligned} \|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2 &= \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 + \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \\ &\quad - 2 \|\mathbf{y}^0 - \boldsymbol{\mu}\| \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\| \cos(\theta) \\ &= \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 + \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 - 2 \langle \mathbf{y}^0 - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle \end{aligned}$$

Tomando-se a esperança da equação em relação a  $\mathbf{y}^0$

$$\begin{aligned}
 E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] &= E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 + \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \right. \\
 &\quad \left. - 2 \langle \mathbf{y}^0 - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle \right] \\
 &= E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 \right] + E_0 \left[ \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \right] \\
 &\quad - 2E_0 \left[ \langle \mathbf{y}^0 - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle \right] \quad (7.1)
 \end{aligned}$$

como  $\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}$  não depende de  $\mathbf{y}^0$  e uma vez que a esperança tomada é em relação à  $\mathbf{y}^0$ , pode-se escrever

$$\begin{aligned}
 -2E_0 \left[ \langle \mathbf{y}^0 - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle \right] &= -2 \langle E_0 [\mathbf{y}^0 - \boldsymbol{\mu}], E_0 [\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}] \rangle \\
 &= -2 \langle E_0 [\mathbf{y}^0 - \boldsymbol{\mu}], \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}_0 \rangle
 \end{aligned}$$

e sendo  $E_0 [\mathbf{y}^0 - \boldsymbol{\mu}] = 0$ , então

$$-2E_0 \left[ \langle \mathbf{y}^0 - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle \right] = -2 \langle 0, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}_0 \rangle = 0.$$

Portanto, a Equação 7.1 se escreve

$$E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] = E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 \right] + E_0 \left[ \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \right].$$

É interessante observar que, em termos de esperança vale o teorema de Pitágoras, isto é, em média o triângulo hachurado é um triângulo retângulo. Novamente,  $E_0 \left[ \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \right] = \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2$ , pois  $\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}$  não depende de  $\mathbf{y}^0$ . Logo,  $E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] = E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 \right] + \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2$ , e portanto o erro de predição é dado por

$$\begin{aligned}
 E \left[ E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] \right] &= E \left[ E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 \right] \right] + E \left[ \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \right]. \quad (7.2)
 \end{aligned}$$

$E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 \right]$  é uma grandeza populacional que não depende de  $\mathbf{y}$ , e portanto  $E \left[ E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 \right] \right] = E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 \right]$ . Logo, a Equação 7.2 fica

$$E \left[ E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] \right] = E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 \right] + E \left[ \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \right].$$

Como  $E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}\|^2 \right] = E \left[ \|\mathbf{y} - \boldsymbol{\mu}\|^2 \right]$  uma vez que  $\mathbf{y}^0$  e  $\mathbf{y}$  podem ser considerados a mesma variável aleatória, a Equação 7.2 se escreve

$$E \left[ E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] \right] = E \left[ \|\mathbf{y} - \boldsymbol{\mu}\|^2 \right] + E \left[ \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \right] \quad (7.3)$$

Ao se aplicar a lei dos cossenos no triângulo tracejado da Figura 52 se tem

$$E \left[ \|\mathbf{y} - \boldsymbol{\mu}\|^2 \right] + E \left[ \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \right] = E \left[ \|\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] + 2E \left[ \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle \right],$$

e portanto, substituindo na Equação 7.3

$$E \left[ E_0 \left[ \|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] \right] = E \left[ \|\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] + 2E \left[ \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle \right].$$

Note que

$$\begin{aligned} E \left[ \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle \right] &= E \left[ \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - E[\boldsymbol{\mu}(\mathbf{y})] + E[\boldsymbol{\mu}(\mathbf{y})] - \boldsymbol{\mu} \rangle \right] \\ &= E \left[ \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - E[\boldsymbol{\mu}(\mathbf{y})] \rangle \right] \\ &\quad + E \left[ \langle \mathbf{y} - \boldsymbol{\mu}, E[\boldsymbol{\mu}(\mathbf{y})] - \boldsymbol{\mu} \rangle \right] \\ &= E \left[ \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - E[\boldsymbol{\mu}(\mathbf{y})] \rangle \right] \\ &\quad + \left[ \langle E[\mathbf{y} - \boldsymbol{\mu}], E[\boldsymbol{\mu}(\mathbf{y})] - \boldsymbol{\mu} \rangle \right] \\ &= E \left[ \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - E[\boldsymbol{\mu}(\mathbf{y})] \rangle \right] \end{aligned}$$

$$\begin{aligned}
& + [\langle 0, E[\boldsymbol{\mu}(\mathbf{y})] - \boldsymbol{\mu} \rangle] \\
& = E[\langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - E[\boldsymbol{\mu}(\mathbf{y})] \rangle] \\
& = E\left[\sum_{i=1}^n (y_i - \mu_i) (\boldsymbol{\mu}(\mathbf{y})_i - E[\boldsymbol{\mu}(\mathbf{y})]_i)\right] \\
& = \sum_{i=1}^n E[(y_i - \mu_i) (\boldsymbol{\mu}(\mathbf{y})_i - E[\boldsymbol{\mu}(\mathbf{y})]_i)] \\
& = \sum_{i=1}^n E[(y_i - \mu_i) (\boldsymbol{\mu}(\mathbf{y})_i - \mu_i)] \\
& = \sum_{i=1}^n \text{cov}(y_i, \boldsymbol{\mu}(\mathbf{y})_i) \\
& = \text{cov}(\mathbf{y}, \boldsymbol{\mu}(\mathbf{y})),
\end{aligned}$$

e então

$$E\left[E_0\left[\|\mathbf{y}^0 - \boldsymbol{\mu}(\mathbf{y})\|^2\right]\right] = E\left[\|\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})\|^2\right] + 2\text{cov}(\mathbf{y}, \boldsymbol{\mu}(\mathbf{y})) \quad (7.4)$$

O erro de predição é então a esperança do erro do ajuste do modelo acrescido de uma penalidade relativa à covariância entre os dados e o estimador.

Como exemplo, suponha que se quer calcular o erro de predição de um novo valor na seguinte situação: Seja  $\mathbf{y} = (y_1, \dots, y_n)$  uma amostra iid e  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ . O erro de predição para uma nova observação  $\mathbf{y}^0$  é dado pela esperança dupla

$$\begin{aligned}
E_0\left[E\left[\|\mathbf{y}^0 - \bar{\mathbf{y}}_n\|^2\right]\right] & = E_0\left[E\left[\|\mathbf{y}^0 - E_0[\mathbf{y}^0] + E_0[\mathbf{y}^0] - \bar{\mathbf{y}}_n\|^2\right]\right] \\
& = E_0 E\left[\|\mathbf{y}^0 - E_0[\mathbf{y}^0]\|^2\right] + E_0 E\left[\|E_0[\mathbf{y}^0] - \bar{\mathbf{y}}_n\|^2\right] \\
& \quad + 2(\mathbf{y}^0 - E_0[\mathbf{y}^0]) (E_0[\mathbf{y}^0] - \bar{\mathbf{y}})
\end{aligned}$$

$$\begin{aligned}
&= E_0 E \left[ \|\mathbf{y}^0 - E_0[\mathbf{y}^0]\|^2 \right] + E_0 E \left[ \|E_0[\mathbf{y}^0] - \bar{\mathbf{y}}_n\|^2 \right] \\
&\quad + 2(\mathbf{y}^0 - E_0[\mathbf{y}^0]) \cdot (0) \\
&= E_0 E \left[ \|\mathbf{y}^0 - E_0[\mathbf{y}^0]\|^2 \right] + E E_0 \left[ \|E_0[\mathbf{y}^0] - \bar{\mathbf{y}}_n\|^2 \right] \\
&= E_0 E \left[ \|\mathbf{y}^0 - E_0[\mathbf{y}^0]\|^2 \right] + E \left[ \|E_0[\mathbf{y}^0] - \bar{\mathbf{y}}_n\|^2 \right] \\
&= E_0 \left[ \|\mathbf{y}^0 - E_0[\mathbf{y}^0]\|^2 \right] + E \left[ \|E_0[\mathbf{y}^0] - \bar{\mathbf{y}}_n\|^2 \right] \\
&= \sigma^2 + \frac{\sigma^2}{n} \\
&= \frac{n+1}{n} \sigma^2. \tag{7.5}
\end{aligned}$$

Um outro exemplo é o erro de predição do estimador de quadrados mínimos em regressão linear. Chamando  $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  de  $\mathbf{H}$ , tem-se

$$\begin{aligned}
E \left[ E_0 \left[ \|\mathbf{y}_0 - \mathbf{H}\mathbf{y}\|^2 \right] \right] &= E \left[ E_0 \left[ \langle \mathbf{y}_0 - \mathbf{H}\mathbf{y}, \mathbf{y}_0 - \mathbf{H}\mathbf{y} \rangle \right] \right] \\
&= E \left[ E_0 \left[ \langle \mathbf{y}_0, \mathbf{y}_0 \rangle - 2 \langle \mathbf{y}_0, \mathbf{H}\mathbf{y} \rangle \right. \right. \\
&\quad \left. \left. + \langle \mathbf{H}\mathbf{y}, \mathbf{H}\mathbf{y} \rangle \right] \right] \\
&= E \left[ E_0 \left[ \langle \mathbf{y}_0, \mathbf{y}_0 \rangle \right] - 2E_0 \left[ \langle \mathbf{y}_0, \mathbf{H}\mathbf{y} \rangle \right] \right. \\
&\quad \left. + E_0 \left[ \langle \mathbf{H}\mathbf{y}, \mathbf{H}\mathbf{y} \rangle \right] \right] \\
&= E \left[ E_0 \left[ \langle \mathbf{y}_0, \mathbf{y}_0 \rangle \right] \right] - 2E \left[ \langle E_0[\mathbf{y}_0], \mathbf{H}\mathbf{y} \rangle \right] \\
&\quad + E \left[ \langle \mathbf{H}\mathbf{y}, \mathbf{H}\mathbf{y} \rangle \right] \\
&= E \left[ E_0 \left[ \langle \mathbf{y}_0, \mathbf{y}_0 \rangle \right] \right] - 2 \left[ \langle E_0[\mathbf{y}_0], \mathbf{H}E[\mathbf{y}] \rangle \right] \\
&\quad + E \left[ \langle \mathbf{H}\mathbf{y}, \mathbf{H}\mathbf{y} \rangle \right] \\
&= E \left[ E_0 \left[ \langle \mathbf{y}_0, \mathbf{y}_0 \rangle \right] \right] - 2 \left[ \langle \boldsymbol{\mu}, \mathbf{H}\boldsymbol{\mu} \rangle \right] \\
&\quad + E \left[ \langle \mathbf{H}\mathbf{y}, \mathbf{H}\mathbf{y} \rangle \right] \\
&= E_0 \left[ \mathbf{y}'_0 \mathbf{I} \mathbf{y}_0 \right] - 2 \langle \boldsymbol{\mu}, \boldsymbol{\mu} \rangle + E \left[ \mathbf{y}' \mathbf{H} \mathbf{y} \right] \\
&= \text{tr}(\mathbf{I}\sigma^2) + \boldsymbol{\mu}' \mathbf{I} \boldsymbol{\mu} - 2\boldsymbol{\mu}' \mathbf{I} \boldsymbol{\mu} + \text{tr}(\mathbf{H}\sigma^2) + \boldsymbol{\mu}' \mathbf{H} \boldsymbol{\mu}
\end{aligned}$$

$$\begin{aligned}
&= n\sigma^2 - \boldsymbol{\mu}'\boldsymbol{\mu} + \sigma^2 \text{tr}(\mathbf{H}) + \boldsymbol{\mu}'\boldsymbol{\mu} \\
&= n\sigma^2 + \sigma^2 p \\
&= \sigma^2(n + p).
\end{aligned}$$

Para se obter um estimador não viesado para o erro de predição, toma-se como estimador de  $E[\|\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})\|^2]$ , utilizando o método dos momentos, a soma de quadrados dos resíduos  $\|\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})\|^2$ . É necessário agora se obter um estimador não viesado para  $\text{cov}(\mathbf{y}, \boldsymbol{\mu}(\mathbf{y}))$ . Para tanto, utiliza-se o famoso Lema de Stein (1981 apud GAJO, 2016, p. 62).

**Lema 1. (Stein)** *Seja  $y \sim N(\theta, 1)$ . Então*

$$\text{cov}(y, h(y)) = E[h(y)(y - \theta)] = E[h'(y)].$$

*Demonstração.* Observe que

$$\begin{aligned}
\text{cov}(y, h(y)) &= E[(y - \theta)(h(y) - E[h(y)])] \\
&= E[(y - \theta)h(y)] - E[(y - \theta)E[h(y)]] \\
&= E[(y - \theta)h(y)] - E[(y - \theta)]E[h(y)] \\
&= E[(y - \theta)h(y)].
\end{aligned}$$

Utilizando integração por partes

$$\begin{aligned}
E[h(y)(y - \theta)] &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} h(y)(y - \theta) \exp\left[-\frac{1}{2}(y - \theta)^2\right] dy \\
&= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} h(y) \left[ \frac{d}{dy} \left( -\exp\left[-\frac{1}{2}(y - \theta)^2\right] \right) \right] dy \\
&= -\frac{1}{\sqrt{2\pi}} h(y) \exp\left[-\frac{1}{2}(y - \theta)^2\right] \Big|_{-\infty}^{\infty}
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} h'(y) \exp \left[ -\frac{1}{2}(y - \theta)^2 \right] dy \\
& = \mathbb{E} [h'(y)].
\end{aligned}$$

□

O lema pode ser facilmente estendido para variável aleatória normal e variância qualquer. Se  $y \sim N(\theta, \sigma^2)$ , seja  $x = \frac{y - \theta}{\sigma} \sim N(0, 1)$ . Como  $y = \sigma x + \theta$ , então

$$h(y) = h(\sigma x + \theta) = g(x)$$

e,

$$\begin{aligned}
\mathbb{E} \left[ h(y) \left( \frac{y - \theta}{\sigma} \right) \right] &= \mathbb{E} [g(x)x] \\
&= \mathbb{E} [g'(x)] \quad (\text{pelo lema}) \\
&= \mathbb{E} [\sigma h'(\sigma x + \theta)] \\
&= \mathbb{E} [\sigma h'(y)].
\end{aligned}$$

Assim,  $\frac{1}{\sigma^2} \mathbb{E} [h(y)(y - \theta)] = \mathbb{E} [h'(y)]$ .

O Lema de Stein admite uma generalização para variáveis normais multivariadas  $\mathbf{y} \sim N_n(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ . No entanto, para esta generalização é necessário supor que a função  $h: \mathbb{R}^n \mapsto \mathbb{R}$  admita uma propriedade técnica denominada *quasi diferenciável*. As funções que têm esta propriedade possuem derivadas parciais em quase todo ponto. Denominando  $\nabla h = \left( \frac{\partial h}{\partial y_1}, \dots, \frac{\partial h}{\partial y_n} \right)$  tem-se o Lema de Stein Multivariado (STEIN, 1981 apud GAJO, 2016, p. 66).

**Lema 2.** *Seja  $\mathbf{y} \sim N_n(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$  e  $h: \mathbb{R}^n \mapsto \mathbb{R}$  quasi diferenciável. Então*

$$\frac{1}{\sigma^2} \mathbb{E} [h(\mathbf{y})(\mathbf{y} - \boldsymbol{\mu})] = \mathbb{E} [\nabla h(\mathbf{y})].$$

Como observação final, o lema de Stein pode ser ainda colocado em termos mais gerais. Considere novamente  $\mathbf{y} \sim N_n(\mu, \sigma^2 \mathbf{I})$  e  $f : \mathbb{R}^n \mapsto \mathbb{R}^n$ ,  $f = (f_1, \dots, f_n)$  um campo vetorial. Pelo Lema 2, então

$$\frac{1}{\sigma^2} \mathbb{E}[(\mathbf{y} - \mu) f_i(\mathbf{y})] = \mathbb{E}[\nabla f_i(\mathbf{y})].$$

Portanto,  $\frac{1}{\sigma^2} \text{cov}(y_i, f_i(\mathbf{y})) = \frac{1}{\sigma^2} \mathbb{E}[(y_i - \mu_i) f_i(\mathbf{y})] = \mathbb{E}\left[\frac{\partial}{\partial y_i} f_i(\mathbf{y})\right]$ . Somando em  $i$  tem-se

$$\begin{aligned} \frac{1}{\sigma^2} \text{cov}(\mathbf{y}, f(\mathbf{y})) &= \frac{1}{\sigma^2} \sum_{i=1}^n \text{cov}(y_i, f_i(\mathbf{y})) \\ &= \mathbb{E}\left[\sum_{i=1}^n \frac{\partial}{\partial y_i} f_i(\mathbf{y})\right] \\ &= \mathbb{E}[\text{div } f(\mathbf{y})]. \end{aligned}$$

Como consequência do Lema de Stein, tem-se um estimador não viesado para  $\text{cov}(\mathbf{y}, \mu(\mathbf{y}))$  dado por  $\text{div } \mu(\mathbf{y}) = \sum_{i=1}^n \frac{\partial}{\partial y_i} \mu_i(\mathbf{y})$ . Utilizando este estimador obtém-se a partir da Equação 7.4 um estimador não viesado para o erro de predição

$$\|\mathbf{y} - \mu(\mathbf{y})\|^2 + 2\sigma^2 \text{div } \mu(\mathbf{y}) = \|\mathbf{y} - \mu(\mathbf{y})\|^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial}{\partial y_i} \mu_i(\mathbf{y}).$$

**Definição 3.** *O Grau de Liberdade, ou degree of freedom (df), de um modelo  $\mu(\mathbf{y})$  é a grandeza populacional*

$$\text{df} = \frac{\text{cov}(\mathbf{y}, \mu(\mathbf{y}))}{\sigma^2} = \sum_{i=1}^n \frac{\text{cov}(y_i, \mu_i(\mathbf{y}))}{\sigma^2}.$$

Portanto, um estimador não viesado para o grau de liberdade é  $\widehat{\text{df}} = \sum_{i=1}^n \frac{\partial}{\partial y_i} \mu_i(\mathbf{y})$ . Uma outra importante aplicação do Lema de Stein é a obtenção de um estimador não viesado para o risco de estimadores. O risco de um estimador é a esperança da função perda. No caso da função perda quadrática, o risco é nada



mais do que o erro quadrático médio. Aplicando a lei dos cossenos ao triângulo tracejado da Figura 52,

$$\|\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})\|^2 = \|\mathbf{y} - \boldsymbol{\mu}\|^2 + \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 - 2 \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle$$

$$\|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 = \|\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})\|^2 - \|\mathbf{y} - \boldsymbol{\mu}\|^2 + 2 \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle$$

Logo,

$$\begin{aligned} \mathfrak{R} &= \mathbb{E} \left[ \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \right] \\ &= \mathbb{E} \left[ \|\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})\|^2 \right] - \mathbb{E} \left[ \|\mathbf{y} - \boldsymbol{\mu}\|^2 \right] + 2\mathbb{E} \langle \mathbf{y} - \boldsymbol{\mu}, \boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu} \rangle \\ &= \|\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})\|^2 - n\sigma^2 + 2\text{cov}(\mathbf{y}, \boldsymbol{\mu}(\mathbf{y})). \end{aligned}$$

E portanto, um estimador não viesado para o risco, isto é, para o erro de predição, é

$$\hat{\mathfrak{R}} = -n\sigma^2 + \|\boldsymbol{\mu}(\mathbf{y}) - \mathbf{y}\|^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial}{\partial y_i} \mu_i(\mathbf{y}).$$

Este estimador é denominado na literatura de SURE (*Stein's Unbiased Risk Estimate*).

**Exemplo 1:** Considere  $\mathbf{y} = (y_1, \dots, y_n)$  uma amostra iid com  $\text{var}(y_i) = \sigma^2$  e  $\boldsymbol{\mu}(\mathbf{y}) = (\bar{y}, \dots, \bar{y})$  com  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ . Logo,

$$\begin{aligned} \hat{\mathfrak{R}} &= -n\sigma^2 + \|\boldsymbol{\mu}(\mathbf{y}) - \mathbf{y}\|^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial}{\partial y_i} \bar{y} \\ &= -n\sigma^2 + \sum_{i=1}^n (y_i - \bar{y})^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial}{\partial y_i} \left( \frac{1}{n} \sum_{i=1}^n y_i \right) \\ &= -n\sigma^2 + \sum_{i=1}^n (y_i - \bar{y})^2 + 2\sigma^2 \sum_{i=1}^n \frac{1}{n} \\ &= -n\sigma^2 + \sum_{i=1}^n (y_i - \bar{y})^2 + 2\sigma^2 \end{aligned}$$

Observe que  $E \left[ \sum_{i=1}^n (y_i - \bar{y})^2 \right] = (n-1)\sigma^2$ . Logo  $E \left[ \hat{\mathfrak{R}} \right] = \sigma^2$ .

**Exemplo 2:** Considere o problema de regressão linear múltipla  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  em que  $\mathbf{X}$  é uma matriz  $n \times p$  de posto completo,  $\mathbf{y} = (y_1, \dots, y_n)$  e  $\text{var}(y_i) = \sigma^2$ . Considerando o subespaço  $\text{Im}(\mathbf{X})$  como o modelo, o estimador  $\boldsymbol{\mu}(\mathbf{y})$  é dado pela projeção ortogonal de  $\mathbf{y}$  em  $\text{Im}(\mathbf{X})$ . Esta projeção é dada explicitamente pela matriz de projeção (*Hat Matrix*)  $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  e  $\boldsymbol{\mu}(\mathbf{y}) = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ . Logo,

$$\begin{aligned}
 \hat{\mathfrak{R}} &= -n\sigma^2 + \|\boldsymbol{\mu}(\mathbf{y}) - \mathbf{y}\|^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial}{\partial y_i} \boldsymbol{\mu}(\mathbf{y}) \\
 &= -n\sigma^2 + \left\| \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - \mathbf{y} \right\|^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial}{\partial y_i} \left( \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{y}) \right)_i \\
 &= -n\sigma^2 + \left\| \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - \mathbf{y} \right\|^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial}{\partial y_i} \sum_{j=1}^n \left[ \left( \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \right)_{ij} y_j \right] \\
 &= -n\sigma^2 + \left\| \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - \mathbf{y} \right\|^2 + 2\sigma^2 \sum_{i=1}^n \left( \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \right)_{ii} \\
 &= -n\sigma^2 + \left\langle \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - \mathbf{y}, \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - \mathbf{y} \right\rangle + 2\sigma^2 \text{tr}(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \\
 &= -n\sigma^2 + \langle \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}, \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \rangle \\
 &\quad - 2 \langle \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}, \mathbf{y} \rangle + \langle \mathbf{y}, \mathbf{y} \rangle + 2\sigma^2 \text{tr}(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \\
 &= -n\sigma^2 + \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} + \mathbf{y}'\mathbf{y} + 2p\sigma^2 \\
 &= -n\sigma^2 + \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} + \mathbf{y}'\mathbf{y} + 2p\sigma^2 \\
 &= -n\sigma^2 - \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} + \mathbf{y}'\mathbf{y} + 2p\sigma^2 \\
 &= (2p-n)\sigma^2 + \mathbf{y}' \left( \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \right) \mathbf{y}.
 \end{aligned}$$

De fato,

$$\begin{aligned}
 \mathfrak{R} &= E \left[ \|\boldsymbol{\mu}(\mathbf{y}) - \boldsymbol{\mu}\|^2 \right] \\
 &= E \left[ \left\langle \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - \boldsymbol{\mu}, \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - \boldsymbol{\mu} \right\rangle \right]
 \end{aligned}$$

$$\begin{aligned}
&= E \left[ \langle X(X'X)^{-1}X'y, X(X'X)^{-1}X'y \rangle - 2 \langle X(X'X)^{-1}X'y, \mu \rangle + \langle \mu, \mu \rangle \right] \\
&= E \left[ y'X(X'X)^{-1}X'X(X'X)^{-1}X' - 2y'X(X'X)^{-1}X'\mu + \mu'\mu \right] \\
&= E \left[ y'X(X'X)^{-1}X'y - 2y'X(X'X)^{-1}X'\mu + \mu'\mu \right] \\
&= E \left[ y'X(X'X)^{-1}X'y \right] - 2E \left[ y'X(X'X)^{-1}X'\mu \right] + \mu'\mu \\
&= E \left[ y'X(X'X)^{-1}X'y \right] - 2E \left[ y' \right] X(X'X)^{-1}X'\mu + \mu'\mu \\
&= E \left[ y'X(X'X)^{-1}X'y \right] - 2\mu'X(X'X)^{-1}X'\mu + \mu'\mu \\
&= E \left[ y'X(X'X)^{-1}X'y \right] - 2\mu'\mu + \mu'\mu \\
&= \text{tr}(X(X'X)^{-1}X'\sigma^2I) + \mu'X(X'X)^{-1}X'\mu - \mu'\mu \quad (\text{Rencher 5.2.a}) \\
&= p\sigma^2 + \mu'\mu - \mu'\mu = p\sigma^2.
\end{aligned}$$

**Exemplo 3:** O resultado fundamental de que o estimador de James-Stein para a média de uma normal multivariada domina estritamente a média amostral pode ser facilmente demonstrado utilizando o SURE. Para  $\mathbf{y} \sim N_n(\boldsymbol{\mu}, I)$  o estimador de James-Stein tem a expressão  $\mu_{JS}(\mathbf{y}) = \left(1 - \frac{n-2}{\|\mathbf{y}\|^2}\right) \mathbf{y}$ . O SURE para este estimador é

$$\begin{aligned}
\hat{\mathfrak{R}} &= -n\sigma^2 + \|\boldsymbol{\mu}(\mathbf{y}) - \mathbf{y}\|^2 + 2\sigma^2 \sum_{i=1}^n \frac{\partial}{\partial y_i} \boldsymbol{\mu}(\mathbf{y}) \\
&= -n + \left\| \left(1 - \frac{n-2}{\|\mathbf{y}\|^2}\right) \mathbf{y} - \mathbf{y} \right\|^2 + 2 \sum_{i=1}^n \frac{\partial}{\partial y_i} \left[ \left(1 - \frac{n-2}{\|\mathbf{y}\|^2}\right) \mathbf{y} \right] \\
&= -n + \left\| \left(\frac{n-2}{\|\mathbf{y}\|^2}\right) \mathbf{y} \right\|^2 + 2 \sum_{i=1}^n \frac{\partial}{\partial y_i} \left[ \left(1 - \frac{n-2}{\|\mathbf{y}\|^2}\right) \mathbf{y} \right] \\
&= -n + \left(\frac{n-2}{\|\mathbf{y}\|^2}\right)^2 \|\mathbf{y}\|^2 + 2 \sum_{i=1}^n \left[ 1 - \frac{n-2}{\|\mathbf{y}\|^2} + \frac{n-2}{\|\mathbf{y}\|^4} 2(y_i)^2 \right] \\
&= -n + \frac{(n-2)^2}{\|\mathbf{y}\|^2} + 2 \left( n - n \frac{n-2}{\|\mathbf{y}\|^2} + \frac{n-2}{\|\mathbf{y}\|^4} 2 \sum_{i=1}^n (y_i)^2 \right)
\end{aligned}$$

$$\begin{aligned}
&= -n + \frac{(n-2)^2}{\|\mathbf{y}\|^2} + 2 \left( n - n \frac{n-2}{\|\mathbf{y}\|^2} + \frac{n-2}{\|\mathbf{y}\|^4} 2\|\mathbf{y}\|^2 \right) \\
&= -n + \frac{(n-2)^2}{\|\mathbf{y}\|^2} + 2 \left( n - \frac{(n-2)}{\|\mathbf{y}\|^2} (n-2) \right) \\
&= -n + \frac{(n-2)^2}{\|\mathbf{y}\|^2} + 2 \left( n - \frac{(n-2)^2}{\|\mathbf{y}\|^2} \right) \\
&= n - \frac{(n-2)^2}{\|\mathbf{y}\|^2}.
\end{aligned}$$

Portanto, o risco do estimador de James-Stein é

$$E[\hat{\mathcal{R}}] = E \left[ n - \frac{(n-2)^2}{\|\mathbf{y}\|^2} \right] = n - (n-2)^2 E \left[ \frac{1}{\|\mathbf{y}\|^2} \right] < n = E[\|\mathbf{y} - \mu\|^2],$$

ou seja, o risco do estimador de James-Stein é estritamente menor do que o risco do estimador natural, isto é, a própria observação  $\mathbf{y}$ .

### 7.3 Seleção de modelos

#### 7.3.1 A Estatística $C_p$ de Mallows

O uso da estatística  $C_p$  de Mallows é uma técnica para seleção de modelos em regressão. Define um critério de ajuste para comparação de modelos com diferente número de parâmetros. Para modelo linear  $M$ ,  $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , com  $E[\mathbf{y}] = X\boldsymbol{\beta} = \boldsymbol{\theta}$ , se supõe que o parâmetro esperança de  $\mathbf{y}$  pertença ao subespaço linear  $\text{Im}(X)$ . Se o modelo não for correto, então  $E[\mathbf{y}] = \boldsymbol{\theta} \notin \text{Im}(X)$  e pode-se decompor ortogonalmente  $\boldsymbol{\theta} = \boldsymbol{\theta}_X + \boldsymbol{\theta}_{X^\perp}$  em relação ao subespaço  $\text{Im}(X)$ . A estimação pelo modelo errado é dada pela projeção ortogonal em  $\text{Im}(X)$ ,  $P_{\text{Im}(X)}\mathbf{y} = \boldsymbol{\theta}_X + \boldsymbol{\varepsilon}_X$ . Desta forma, a matriz de dispersão deste estimador é

$$E[(\hat{\mathbf{y}} - \boldsymbol{\theta})(\hat{\mathbf{y}} - \boldsymbol{\theta})'] = E[(\boldsymbol{\theta}_X + \boldsymbol{\varepsilon}_X - (\boldsymbol{\theta}_X + \boldsymbol{\theta}_{X^\perp}))$$

$$\begin{aligned}
& (\boldsymbol{\theta}_X + \boldsymbol{\varepsilon}_X - (\boldsymbol{\theta}_X + \boldsymbol{\theta}_{X^\perp}))' \\
&= \text{E} [(\boldsymbol{\varepsilon}_X - \boldsymbol{\theta}_{X^\perp}) (\boldsymbol{\varepsilon}_X - \boldsymbol{\theta}_{X^\perp})'] \\
&= \text{E} [\boldsymbol{\varepsilon}_X \boldsymbol{\varepsilon}'_X + \boldsymbol{\theta}_{X^\perp} \boldsymbol{\theta}'_{X^\perp}] \\
&= \text{E} [\boldsymbol{\varepsilon}_X \boldsymbol{\varepsilon}'_X] + \boldsymbol{\theta}_{X^\perp} \boldsymbol{\theta}'_{X^\perp} \\
&= \sigma^2 \mathbf{I}_{p \times p} + \boldsymbol{\theta}_{X^\perp} \boldsymbol{\theta}'_{X^\perp}.
\end{aligned}$$

O erro quadrático médio do estimador  $\hat{\mathbf{y}}$  é dado por

$$\begin{aligned}
\text{EQM}(\hat{\mathbf{y}}) &= \text{E} [(\hat{\mathbf{y}} - \boldsymbol{\theta})' (\hat{\mathbf{y}} - \boldsymbol{\theta})] \\
&= \text{tr} (\text{E} [(\hat{\mathbf{y}} - \boldsymbol{\theta})' (\hat{\mathbf{y}} - \boldsymbol{\theta})]) \\
&= \text{tr} (\text{E} [(\hat{\mathbf{y}} - \boldsymbol{\theta}) (\hat{\mathbf{y}} - \boldsymbol{\theta})']) \\
&= \text{tr} [\sigma^2 \mathbf{I} + \boldsymbol{\theta}_{X^\perp} \boldsymbol{\theta}'_{X^\perp}] \\
&= p\sigma^2 + \|\boldsymbol{\theta}_{X^\perp}\|^2.
\end{aligned}$$

Ao se optar entre vários modelos, deve-se escolher aquele que tem menor  $\text{EQM}(\hat{\mathbf{y}})$ . Entretanto, uma vez que  $\text{EQM}(\hat{\mathbf{y}})$  depende de parâmetros populacionais, é necessário buscar um estimador não viesado para esta grandeza. Primeiramente, calcula-se a esperança da estatística  $\hat{\sigma}^2 = \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}{n-p}$ , por

$$\begin{aligned}
\text{E} [\hat{\sigma}^2] &= \text{E} \left[ \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}{n-p} \right] \\
&= \text{E} \left[ \frac{\mathbf{y}' (\mathbf{I} - \mathbf{P}_X) \mathbf{y}}{n-p} \right] \\
&= \frac{1}{n-p} \text{E} [(\boldsymbol{\theta} + \boldsymbol{\varepsilon})' (\mathbf{I} - \mathbf{P}_X) (\boldsymbol{\theta} + \boldsymbol{\varepsilon})] \\
&= \frac{1}{n-p} \text{E} [\text{tr} \{(\mathbf{I} - \mathbf{P}_X) (\boldsymbol{\theta} + \boldsymbol{\varepsilon}) (\boldsymbol{\theta} + \boldsymbol{\varepsilon})'\}] \\
&= \frac{1}{n-p} \text{E} [\text{tr} \{(\boldsymbol{\theta}_{X^\perp} + \boldsymbol{\varepsilon}_{X^\perp}) (\boldsymbol{\theta} + \boldsymbol{\varepsilon})'\}]
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n-p} \mathbf{E} [(\boldsymbol{\theta}_{X^\perp} + \boldsymbol{\varepsilon}_{X^\perp})(\boldsymbol{\theta}_X + \boldsymbol{\varepsilon}_X)'] \\
&= \frac{1}{n-p} (\boldsymbol{\theta}_{X^\perp} \boldsymbol{\theta}'_{X^\perp} + \mathbf{E}[\boldsymbol{\varepsilon}_{X^\perp} \boldsymbol{\varepsilon}'_{X^\perp}]) \\
&= \frac{1}{n-p} \|\boldsymbol{\theta}_{X^\perp}\|^2 + \sigma^2.
\end{aligned}$$

Logo, para  $(n-p)\hat{\sigma}^2 - (n-2p)\sigma^2$ , tem-se

$$\begin{aligned}
\mathbf{E} [(n-p)\hat{\sigma}^2 - (n-2p)\sigma^2] &= \|\boldsymbol{\theta}_X\|^2 + (n-p)\sigma^2 - (n-2p)\sigma^2 \\
&= \|\boldsymbol{\theta}_X\|^2 + p\sigma^2 \\
&= \text{EQM}(\hat{\mathbf{y}}).
\end{aligned}$$

Portanto,  $\hat{\sigma}^2$  é um estimador viesado de  $\sigma^2$ . Se é obtido um estimador não viesado  $\hat{\sigma}^*$  para  $\sigma^2$ , por exemplo, considerando um modelo de dimensão muito maior do que o modelo  $M$  considerado, se tem então um estimador não viesado para  $\text{EQM}(\hat{\mathbf{y}})$  dado por

$$\begin{aligned}
\text{EQM}(\hat{\mathbf{y}}) &= \mathbf{E} [(n-p)(\hat{\sigma}^2 - \hat{\sigma}^{2*}) + p\hat{\sigma}^2] \\
&= \mathbf{E} [\|\mathbf{y} - \hat{\mathbf{y}}\|^2 - (n-p)\hat{\sigma}^{2*} + p\hat{\sigma}^2].
\end{aligned}$$

Desta forma,  $\|\mathbf{y} - \hat{\mathbf{y}}\|^2 - (n-p)\hat{\sigma}^{2*} + p\hat{\sigma}^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 + 2p\hat{\sigma}^{2*} - n\hat{\sigma}^{2*}$  é um estimador não viesado de  $\text{EQM}(\hat{\mathbf{y}})$ . A estatística  $C_p$  de Mallows é então definida como

$$C_p = \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2 + 2p\hat{\sigma}^{2*} - n\hat{\sigma}^{2*}}{\hat{\sigma}^{2*}} = \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}{\hat{\sigma}^{2*}} + 2p - n = \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}{\hat{\sigma}^{2*}} - (n-2p).$$

Assim, a estatística  $C_p$  é um estimador não viesado de  $\frac{1}{n-p} \|\boldsymbol{\theta}_{X^\perp}\|^2 + \sigma^2$ , isto é, de uma medida ponderada pela dimensão de quanto o modelo se afasta do modelo real. No caso de se escolher entre modelos, deve-se optar por aquele

com menor  $C_p$ , pois neste caso é razoável admitir que este possui o menor viés ponderado pela dimensão. A constituição deste estimador pode ser explicada pela geometria da construção (Figura 53).

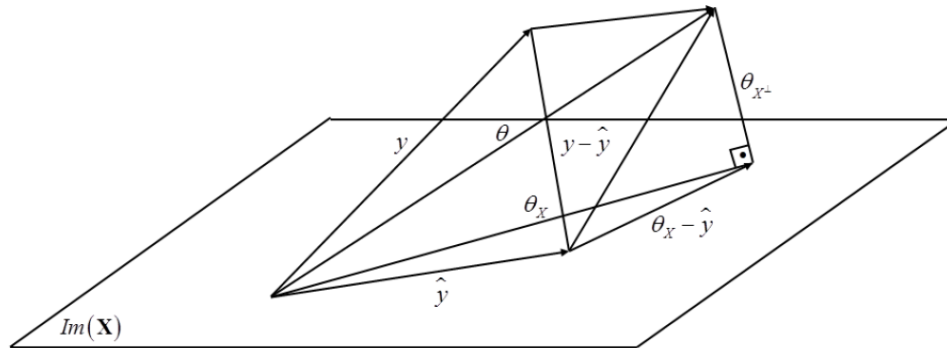


Figura 53 Elementos representativos da estrutura  $C_p$ .

#### 7.4 O conceito de grau de liberdade

Grau de liberdade é um dos termos mais utilizados na literatura estatística, sendo onipresente nos artigos científicos da área. No entanto, a experiência mostra que o conceito é utilizado sem que de fato se tenha consciência de seu verdadeiro significado. Uma razão para tal fato, talvez, seja que o grau de liberdade pode e é utilizado em vários sentidos. Vejamos alguns deles. Ao se referir a uma qui-quadrado com  $n$  graus de liberdade, ou a uma distribuição  $t$  com  $n$  graus de liberdade,  $n$  significa simplesmente um parâmetro para estas distribuições, ou seja, ambas as distribuições são famílias uniparamétricas indexadas pelos números naturais, os quais são denominados Grau de Liberdade. A origem do termo decorre do fato de que uma variável aleatória obtida como soma de quadrados de  $n$  normais padrão independentes tem distribuição qui-quadrado com  $n$  graus de liberdade. Nesse sentido, pode-se supor que, se é observado um fenômeno que depende, por exemplo, de  $n$  fatores e estes se combinem de forma que nenhum de-

les restrinja o comportamento dos demais, então é razoável dizer que este sistema tem  $n$  graus de liberdade. Como definir então grau de liberdade para métodos de estimação mais gerais?

No caso de regressão linear  $\mathbf{y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}$ , com  $\mathbf{X}$  de posto completo e  $p < n$ , o grau de liberdade é a dimensão da imagem da transformação  $\mathbf{X}$ , isto é, a dimensão do modelo. O grau de liberdade, em uma definição geral, deve refletir isto, a complexidade do modelo.

O modelo supõe que todos os valores das componentes do vetor de parâmetros  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$  podem assumir qualquer valor real. Isto significa que a única suposição sobre o vetor de médias  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n) = \text{E}[\mathbf{y}]$  é que este pertença ao subespaço imagem de  $\mathbf{X}$  (denotado por  $\text{Im}(\mathbf{X})$ ). A dimensão deste subespaço é  $p$ . O processo de estimação derivado do método dos quadrados mínimos é simplesmente a projeção ortogonal do vetor observado  $\mathbf{y}$  no subespaço  $\text{Im}(\mathbf{X})$ . Observe que todas as amostras da forma  $\mathbf{y} + \mathbf{w}$ , com  $\mathbf{w}$  pertencendo ao subespaço ortogonal à  $\text{Im}(\mathbf{X})$ , resultam na mesma estimativa para o vetor de médias, uma vez que o vetor  $\mathbf{w}$  é projetado na origem. Desta forma, estas amostras não contêm novas informações para a estimação do vetor de médias. Como, informalmente, o número de tais amostras, isto é, o número de vetores  $\mathbf{w}$ , pode ser medido pela dimensão do espaço ortogonal, que é  $n - p$ , então a quantidade de informação da amostra em relação à estimação da média fica restrita à  $\text{Im}(\mathbf{X}) = p$ . Logo, o grau de liberdade pode ser explicado também em termos da dimensão do modelo, isto é, a dimensão do espaço em que se supões estar o vetor de médias  $\boldsymbol{\mu}$ . No aspecto algébrico, um subespaço de dimensão  $k$  é definido por  $k$  equações lineares. Escolhendo-se uma base formada por vetores nesse subespaço e por vetores no ortogonal, o subespaço é simplesmente definido zerando-se as variáveis correspondentes ao subespaço ortogonal, ficando livres as demais variáveis. Esta



construção leva à definição mais utilizada de grau de liberdade, que é a diferença entre o número total de variáveis e o número de vínculos, ou seja, o número de equações. Portanto, a interpretação geométrica é clara no sentido de que o grau de liberdade é igual à dimensão de  $\text{Im}(X) = p$ , que é o número de variáveis livres para determinação do vetor de médias  $\mu$ , sendo que as variáveis relativas ao subespaço ortogonal de  $\text{Im}(X)$  são dependentes desta.

Como visto na seção relativa à estimação do erro de predição de um modelo, vê-se que a esperança do erro de predição é o erro quadrático médio mais uma penalização dada por  $(\mathbf{y}, \mu(\mathbf{y}))$ . Como também foi visto para o caso linear  $\mu(\mathbf{y}) = M\mathbf{y}$ ,  $\frac{\text{cov}(\mathbf{y}, \mu(\mathbf{y}))}{\sigma^2}$  é exatamente  $\text{tr}(M)$  e, portanto, o grau de liberdade usual.

Certamente a penalização  $\text{cov}(\mathbf{y}, \mu(\mathbf{y}))$  reflete a complexidade do modelo. A ideia de complexidade e a de graus de liberdade em um sistema físico estão intuitivamente relacionadas. Pensando em um sistema físico a complexidade, e portanto o grau de liberdade, não deve depender da amplitude de variação do sistema. Um sistema planetário não é necessariamente mais complexo do que um sistema orbital de um átomo. Portanto, é razoável que  $\text{cov}(\mathbf{y}, \mu(\mathbf{y}))$  seja dividido por  $\sigma^2$  na penalização do erro de predição. Com as considerações acima descritas é razoável então a Definição 3 (EFRON et al., 2004) na Subseção 7.2. A partir desta definição,  $df$  é evidentemente um parâmetro populacional e precisa ser estimado. Neste sentido tem-se o resultado fundamental dado pelo Lema 1 (Stein). Sob suposição de normalidade  $N_n(\mu, \sigma^2 I)$ , tem-se que

$$\begin{aligned} \frac{\text{cov}(\mathbf{y}, \mu(\mathbf{y}))}{\sigma^2} &= \frac{1}{\sigma^2} E[(\mathbf{y} - \boldsymbol{\mu})(\mu(\mathbf{y}) - E[\mu(\mathbf{y})])] \\ &= \frac{1}{\sigma^2} E[(\mathbf{y} - \boldsymbol{\mu})\mu(\mathbf{y})] \\ &= E[\nabla \mu(\mathbf{y})], \end{aligned}$$

e também  $\nabla g(\mathbf{x}) = \text{div } g(\mathbf{x}) = \sum_{i=1}^n \frac{\partial g_i(\mathbf{x})}{\partial x_i}$  é o divergente da função  $g$ .

Pelo lema de Stein tem-se então como estimador não viesado de  $df$  de um modelo  $M(\mathbf{y})$ ,  $\nabla M(\mathbf{y}) = \text{div } \mu(\mathbf{y})$ . Como mostrado no Exemplo 2 (Subseção 7.2), se  $\mu(\mathbf{y}) = M\mathbf{y}$ , então  $\widehat{df} = \text{div } \mu(\mathbf{y}) = \text{tr}(M) = p$  e neste caso simples o estimador não possui variabilidade e é igual ao  $df$ .

No lema de Stein, tem-se uma hipótese matemática bastante técnica que o modelo  $M : \mathbb{R}^n \mapsto \mathbb{R}$ , em que  $M_i(\mathbf{y}) = (\mu(\mathbf{y}))_i$ , tem que ser quasidiferenciável. Tal fato torna o cálculo do  $\text{div } \mu(\mathbf{y})$ , para os casos em questão (LASSO, *Elastic Net*, Ridge) bastante complexos. Neste trabalho apenas alguns aspectos serão abordados.

### 7.5 Grau de liberdade na estimação por encolhimento

Esta subseção está baseada no artigo de Kato (2009). No contexto geral do processo de estimação,  $K_p$  é um convexo fechado no espaço paramétrico e o estimador é

$$\hat{\beta}_{K_p}(\mathbf{y}) = \min_{\beta \in K_p} \left\| \beta - \hat{\beta}^{\circ}(\mathbf{y}) \right\|_p,$$

com  $\hat{\beta}^{\circ} = \hat{\beta}_{\text{ols}}$ ,  $X\hat{\beta}_{K_p}(\mathbf{y}) = \mu_K(\mathbf{y})$  e  $\mu : \mathbb{R}^n \mapsto K$  é a projeção de distância mínima no convexo  $K = X K_p$ .

O primeiro resultado, conforme em Kato (2009, p. 4), é bastante técnico e garante que  $\mu_K$  e  $\beta_{K_p}$  são funções absolutamente contínuas e, portanto, satisfazem a condição do lema de Stein. Pelo Lema de Stein, um estimador não viesado para o grau de liberdade é o divergente de  $\hat{\mu}_K$  (o divergente é relativo à variável  $\mathbf{y}$ ). Segue agora uma relação entre os divergentes das funções  $\mu_K$  e  $\beta_{K_p}$ , em que o divergente de  $\hat{\beta}_{K_p}$  é em relação à variável  $\beta$ .

Seja  $P_{K_p} : \mathbb{R}^p \mapsto K_p$  a projeção em  $K_p$  de distância mínima em relação à métrica  $\langle \cdot, \cdot \rangle_p$  e  $\hat{\beta}^{\circ} : \mathbb{R}^n \mapsto \mathbb{R}^p$ , com  $\hat{\beta}^{\circ}(\mathbf{y}) = (X'X)^{-1}X'\mathbf{y}$  o estimador de quadrados mínimos. Logo,  $\hat{\beta}_{K_p}(\mathbf{y}) = P_{K_p}(\hat{\beta}^{\circ}(\mathbf{y}))$  e  $\mu_K(\mathbf{y}) = X(P_{K_p}(\hat{\beta}^{\circ}(\mathbf{y})))$

e pode-se então utilizar a regra da cadeia.

Se  $f : \mathbb{R}^n \mapsto \mathbb{R}^n$  então  $df(\mathbf{x})$  é uma transformação linear  $df(\mathbf{x}) : \mathbb{R}^n \mapsto \mathbb{R}^n$ . A definição de divergente independente de coordenadas é  $\text{div } f(\mathbf{x}) = \text{tr}(df(\mathbf{x}))$ . Se  $f$  é uma transformação linear  $f(\mathbf{y}) = M\mathbf{y}$  então  $df(\mathbf{y}) = M$  e  $\text{div } f(\mathbf{y}) = \text{tr}(M)$ . No caso de composição,  $\beta : \mathbb{R}^n \xrightarrow{f} \mathbb{R}^p \xrightarrow{g} \mathbb{R}^n$ , a derivada é dada pela composição de transformações lineares  $d(g \circ f)(\mathbf{x}) = dg(f(\mathbf{x})) \cdot df(\mathbf{x})$ . Tem-se então

$$\begin{aligned}\boldsymbol{\mu}_K(\mathbf{y}) &= X\hat{\boldsymbol{\beta}}_{K_p}(\mathbf{y}) \\ &= XP_{K_p}\hat{\boldsymbol{\beta}}^\circ(\mathbf{y}) \\ &= XP_{K_p}(X'X)^{-1}X'\mathbf{y}.\end{aligned}$$

Logo,  $d\boldsymbol{\mu}_K(\mathbf{y}) = X \cdot dP_{K_p}(\hat{\boldsymbol{\beta}}^\circ(\mathbf{y})) (X'X)^{-1}X'$ . Segue então que

$$\begin{aligned}\text{div } \boldsymbol{\mu}_K(\mathbf{y}) &= \text{tr}\left(X dP_{K_p}(\hat{\boldsymbol{\beta}}^\circ(\mathbf{y})) (X'X)^{-1}X'\right) \\ &= \text{tr}\left(dP_{K_p}(\hat{\boldsymbol{\beta}}^\circ(\mathbf{y})) (X'X)^{-1}X'X\right) \\ &= \text{tr}\left(dP_{K_p}(\hat{\boldsymbol{\beta}}^\circ(\mathbf{y}))\right).\end{aligned}$$

Portanto o divergente de  $\boldsymbol{\mu}_K$  no ponto  $\mathbf{y}$  é igual ao divergente da função  $P_{K_p} : \mathbb{R}^p \mapsto K_p \subset \mathbb{R}^p$ , em relação à variável  $\boldsymbol{\beta}$ , aplicado no ponto  $\hat{\boldsymbol{\beta}}^\circ(\mathbf{y})$ .

Como exemplo, suponha  $K_p$  uma bola de raio  $r$  centrada na origem. Para um ponto  $\boldsymbol{\beta}$  fora da bola,  $P_{K_p}(\boldsymbol{\beta}) = r \frac{\boldsymbol{\beta}}{\|\boldsymbol{\beta}\|}$ . Logo, ao se derivar em relação a  $\beta_i$  tem-se

$$\frac{\partial}{\partial \beta_i} \left( \frac{r\beta_i}{\sqrt{\sum_{j=1}^n \beta_j^2}} \right) = r \frac{\sqrt{\sum_{j=1}^n \beta_j^2} - \frac{1}{2} \left( \sum_{j=1}^n \beta_j^2 \right)^{-\frac{1}{2}} (2\beta_i) \beta_i}{\sum_{j=1}^n \beta_j^2}.$$

Como o  $\text{div} (P_{K_p}(\boldsymbol{\beta}))$  é o somatório, em  $i$ , de tal derivada, então

$$\begin{aligned}
 \text{div} (P_{K_p}(\boldsymbol{\beta})) &= \sum_{i=1}^n \left( r \frac{\sqrt{\sum_{j=1}^n \beta_j^2} - \frac{1}{2} \left( \sum_{j=1}^n \beta_j^2 \right)^{-\frac{1}{2}} (2\beta_i) \beta_i}{\sum_{j=1}^n \beta_j^2} \right) \\
 &= \frac{r \left[ n \sqrt{\sum_{j=1}^n \beta_j^2} - \left( \sum_{j=1}^n \beta_j^2 \right)^{-\frac{1}{2}} \left( \sum_{i=1}^n \beta_i^2 \right) \right]}{\sum_{j=1}^n \beta_j^2} \\
 &= \frac{r \left[ n \sqrt{\sum_{j=1}^n \beta_j^2} - \left( \sum_{j=1}^n \beta_j^2 \right)^{\frac{1}{2}} \right]}{\sum_{j=1}^n \beta_j^2} \\
 &= \frac{r \left[ (n-1) \sqrt{\sum_{j=1}^n \beta_j^2} \right]}{\sum_{j=1}^n \beta_j^2} \\
 &= \frac{r(n-1)}{\sqrt{\sum_{j=1}^n \beta_j^2}} \\
 &= \frac{r(n-1)}{\|\boldsymbol{\beta}\|}.
 \end{aligned}$$

## 7.6 Graus de liberdade e estimativas Cp

Como os métodos de estimação LARS, LASSO e Elasticnet são também métodos de seleção de covariáveis é necessário algum critério de seleção de modelos encaixados para a conclusão do processo de estimação. Uma das possibilidades é a de se utilizar estatísticas tipo Cp de Mallows.

Conforme demonstrado anteriormente, se  $\hat{\boldsymbol{\mu}} = g(\mathbf{y})$  então

$$E \left[ \frac{\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|^2}{\sigma^2} \right] = E \left[ \frac{\|\mathbf{y} - \hat{\boldsymbol{\mu}}\|^2}{\sigma^2} \right] - n + 2 \frac{\sum_{i=1}^m \text{cov}(\hat{\mu}_i, y_i)}{\sigma^2},$$

em que o grau de liberdade do estimador é dado por  $\text{df}_{\boldsymbol{\mu}, \sigma^2} = \sum_{i=1}^m \text{cov}(\hat{\mu}_i, y_i)$  e a estatística Cp fica da forma  $C_p(\hat{\boldsymbol{\mu}}) = \frac{\|\mathbf{y} - \hat{\boldsymbol{\mu}}\|^2}{\sigma^2} n + 2 \text{df}_{\boldsymbol{\mu}, \sigma^2}$ .

Como  $\boldsymbol{\mu}$ ,  $\sigma^2$  e  $\text{df}_{\boldsymbol{\mu}, \sigma^2}$  são parâmetros populacionais, para que Cp possa ser utilizada é necessário estimar estes parâmetros. É sugerido então (Eron et al., 2004) o procedimento que se segue. Obtém-se os estimadores de quadrados mínimos de  $\bar{\boldsymbol{\mu}}$  e  $\bar{\sigma}^2$  considerando o modelo completo. Para estimar  $\text{df}_{\boldsymbol{\mu}, \sigma^2}$  é utilizado o método *bootstrap*. Amostras *bootstrap*  $\mathbf{y}^*$  são geradas supondo  $\mathbf{y}^* \sim N(\bar{\boldsymbol{\mu}}, \bar{\sigma}^2)$  e estimativas  $\hat{\boldsymbol{\mu}}^* = g(\mathbf{y}^*)$ . Repetindo o processo de forma independente por, digamos, B vezes obtém-se estimativas para as covariâncias

$$\widehat{\text{cov}}_i = \sum_{b=1}^B \frac{\hat{\mu}_i^*(b) [y_i^*(b) - y_i^*(\cdot)]}{B-1},$$

em que  $\mathbf{y}^*(\cdot) = \frac{\sum_{b=1}^B \mathbf{y}^*(b)}{B}$ , obtendo-se então  $\widehat{\text{df}} = \sum_{i=1}^n \frac{\widehat{\text{cov}}_i}{\bar{\sigma}^2}$ .

A suposição de normalidade parece não ser essencial nos resultados, conforme indica a simulação computacional em Efron et al. (2004). Os resultados obtidos sugerem que o grau de liberdade para o estimador LARS no k-ésimo passo é  $\text{df}(\hat{\boldsymbol{\mu}}_k) \approx k$ , e portanto igual ao grau de liberdade do estimador OLS com k covariáveis. Neste caso, a estimativa para a seleção dos modelos fica da forma

$$C_p = \frac{\|\mathbf{y} - \hat{\boldsymbol{\mu}}_k\|^2}{\bar{\sigma}^2} - n + 2k.$$

O resultado computacional de fato revela uma propriedade verdadeira.

**Teorema 3.** Se as covariáveis  $\mathbf{x}_1, \dots, \mathbf{x}_m$  são ortogonais, então o estimador  $\hat{\boldsymbol{\mu}}_k(\mathbf{y})$  dado pelo  $k$ -ésimo passo do algoritmo LARS possui  $df(\hat{\boldsymbol{\mu}}_k) = k$ .

*Demonstração.* Do Lema 3 do artigo, cuja demonstração será omitida, tem-se que  $\hat{\boldsymbol{\mu}}_k$  é quasi diferenciável e portanto se pode aplicar o Lema de Stein. Assim,

$$\nabla \cdot \hat{\boldsymbol{\mu}}_k = \sum_i \frac{\partial \hat{\mu}_{k,i}}{\partial \mathbf{y}_i}(\mathbf{y}) = \sum_i \mathbf{I} \left[ |\mathbf{y}_i| > |\mathbf{y}_i|_{(k+1)} \right] = k,$$

uma vez que existem  $k$  coordenadas em que  $|\mathbf{y}_i| > |\mathbf{y}_i|_{(k+1)}$ .  $\square$

### 7.7 Grau de liberdade do estimador Ridge

Uma das formas de se determinar o grau de liberdade do estimador Ridge é como se segue.

Sendo  $\boldsymbol{\mu}(\mathbf{y}) = \mathbf{X}(\mathbf{X}'\mathbf{X} + k\mathbf{I})^{-1}\mathbf{X}'\mathbf{y}$ , pode-se escrever

$$\begin{aligned} df &= \text{cov}(\mathbf{y}, \boldsymbol{\mu}(\mathbf{y})) \\ &= \sum_{i=1}^n \text{cov}(\mathbf{y}_i, \boldsymbol{\mu}(\mathbf{y})_i) \\ &= \sum_{i=1}^n \text{cov} \left( \mathbf{y}_i, \left( \mathbf{X}(\mathbf{X}'\mathbf{X} + k\mathbf{I})^{-1}\mathbf{X}'\mathbf{y} \right)_i \right) \\ &= \sum_{i=1}^n \text{cov} \left( \mathbf{y}_i, \left( \mathbf{X}(\mathbf{X}'\mathbf{X} + k\mathbf{I})^{-1}\mathbf{X}' \right)_i \mathbf{y} \right) \\ &= \sum_{i=1}^n \left( \mathbf{X}(\mathbf{X}'\mathbf{X} + k\mathbf{I})^{-1}\mathbf{X}' \right)_{ii} \text{cov}(\mathbf{y}_i, \mathbf{y}_i) \\ &= \text{tr} \left( \mathbf{X}(\mathbf{X}'\mathbf{X} + k\mathbf{I})^{-1}\mathbf{X}' \right) \sigma^2 \\ &= \sigma^2 \text{tr} \left[ (\mathbf{X}'\mathbf{X} + k\mathbf{I})^{-1}\mathbf{X}'\mathbf{X} \right] \\ &= \sigma^2 \text{tr} \left[ (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{X} + k\mathbf{I}) \right]^{-1} \\ &= \sigma^2 \text{tr} \left[ \mathbf{I} + k(\mathbf{X}'\mathbf{X})^{-1} \right]^{-1}. \end{aligned}$$

Como o traço é invariante por conjugação, isto é,  $\text{tr}(ABA^{-1}) = \text{tr}(B)$ , pode-se considerar  $X'X$  na forma diagonal  $X'X = \text{diag}(\lambda_i)$ . Logo,

$$\begin{aligned} df &= \sigma^2 \text{tr}[\mathbf{I} + k \text{diag}(\lambda_i^{-1})]^{-1} \\ &= \sigma^2 \text{tr} \left[ \text{diag} \left( 1 + \frac{k}{\lambda_i} \right) \right]^{-1} \\ &= \sigma^2 \text{tr} \left[ \text{diag} \left( \frac{\lambda_i + k}{\lambda_i} \right) \right]^{-1} \\ &= \sigma^2 \sum_{i=1}^p \frac{\lambda_i}{\lambda_i + k}. \end{aligned}$$

Observa-se que se  $k = 0$ , tem-se o estimador de quadrados mínimos e então  $df = \sigma^2 \sum_{i=1}^p \frac{\lambda_i}{\lambda_i + 0} = p\sigma^2$ . Entretanto, para  $k \rightarrow \infty$ ,  $df \rightarrow 0$ . Desta forma, se tem o comportamento conforme descrito na Figura 54.

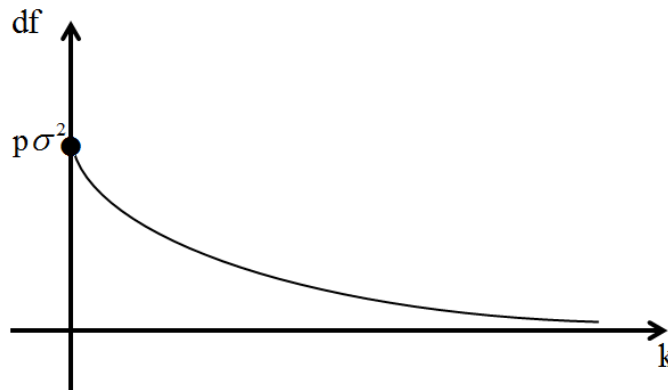


Figura 54 Comportamento do  $df$  em relação ao parâmetro Ridge.

## 7.8 O grau de liberdade do estimador LASSO

Diversos autores (KATO, 2009; ZOU, HASTIE, 2005; TIBSHIRANI, TAYLOR, 2012) demonstram que um estimador não viesado para o grau de liberdade do estimador LASSO é o número de variáveis selecionadas. Este resultado,





## 8 Aspectos computacionais e aplicações

Conforme destacado anteriormente, os métodos LASSO e Elastic Net não possuem forma fechada para se determinar seus respectivos estimadores. Neste sentido, faz-se necessário o emprego de métodos computacionais iterativos para se determinar estimativas de  $\beta$ . Dentre os vários programas que permitem o cálculo de tais estimativas, optou-se por programas implementados no software R em razão de este ser uma linguagem livre (*freeware*) e por ser amplamente empregado no meio estatístico devido a sua forma simples de programação.

Existem dois principais pacotes no R para aplicação dos métodos LASSO, Ridge e Elastic Net: o `elasticnet` e o `glmnet`. Criados respectivamente em 2012 e 2016, ambos são excelentes ferramentas para a computação de estimativas via LASSO/Elastic Net. Entretanto optou-se por utilizar o segundo pacote em detrimento do primeiro, uma vez que nas próprias palavras de um dos autores do pacote (e também criador do método Elastic Net), Trevor Hastie, o pacote contém procedimentos extremamente eficientes para o ajuste do LASSO ou Elastic Net. Tal eficiência se dá em função da metodologia empregada no algoritmo de minimização da penalização, o qual é baseado em uma metodologia conhecida por *cyclical coordinate descent*, que por sua vez tem o LARS como fundamento. O `glmnet` é um pacote relativamente simples de ser utilizado. Ainda assim, sua grandiosidade se apresenta principalmente por sua eficiência computacional e por sua robustez, esta última no sentido de que o pacote é ideal para aplicação em *big data*, ou seja, grande quantidade de dados. Segundo Efron et al. (2003), o pacote suportou com alta eficiência conjuntos de dados da ordem de 800 mil covariáveis. Desta forma, os métodos LASSO e Elastic Net, junto com o pacote `glmnet`, se apresentam como uma ferramenta poderosa para análise de dados genéticos, dada a grande quantidade de dados presentes em amostras genéticas (mesmo com um número pequeno

de observações, geralmente se tem uma informação genética muito grande em tais amostras, isto é,  $p > n$ ). Esta situação é exatamente a mesma descrita anteriormente neste texto em que se tem  $p \gg n$ . Segue então uma descrição sobre o pacote a ser utilizado.

O pacote `glmnet` foi originalmente criado com o intuito de prover uma ferramenta de ajuste de modelo generalizado via máxima verossimilhança penalizada. Mais especificamente fornece, em linguagem R, uma forma de se determinar a estimativa do vetor de parâmetros pelos métodos Ridge, LASSO e Elastic Net. Seus autores, a se saber Friedman, F.; Hastie, T.; Simon, N.; Tibshirani, R., são em sua maioria os próprios autores dos artigos originais dos métodos LASSO e Elastic Net. O comando `glmnet` (`arguments`) cria um objeto do tipo `glmnet` que possui um vasto conjunto de informações referentes ao ajuste realizado. Dois principais (e essenciais) argumentos a serem definidos para a aplicação do pacote são, respectivamente, a matriz das covariáveis (por exemplo  $X$ ) e um vetor resposta (por exemplo  $y$ ), sendo que o número de linhas da matriz  $X$  deve coincidir com o número de linhas do vetor  $y$ . Portanto, uma simples aplicação do pacote seria da forma

```
reg=glmnet (X, y)
```

É importante frisar que com o comando anterior o ajuste, bem como vários parâmetros inerentes a este, é armazenado no objeto `reg`. Vários outros parâmetros possíveis de se determinar como argumento do pacote são bastante específicos. Desta forma seguem destacados apenas os mais relevantes a nossa pesquisa.

`alpha = c`, com  $c \in [0, 1]$ , determina se o ajuste a ser feito é do tipo Ridge, LASSO ou Elastic Net. Se  $c = 0$  o pacote executa um ajuste Ridge. Se  $c = 1$ , o ajuste realizado é pelo método LASSO. Já para  $0 < c < 1$ , é ajustado um modelo do tipo Elastic Net, sendo que quanto mais próximo de 0 o parâmetro  $c$  estiver,

mais próximo de uma restrição Ridge será o formato da penalização, enquanto que um valor de  $c$  mais próximo de 1 resulta em uma penalização com formato mais próximo a restrição de norma  $\| \cdot \|_1$ . Portanto, um comando do tipo

```
reg=glmnet(X, y, alpha=1)
```

irá realizar um ajuste LASSO com o vetor de respostas  $y$  e com a matriz  $X$ , cujas colunas são as covariáveis analisadas. Já um comando do tipo

```
reg=glmnet(X, y, alpha=0.5)
```

irá ajustar um modelo no formato Elastic Net. A função `plot(arguments)` permite que se visualize os traços gráficos das componentes do vetor estimado para cada valor da norma L1, lembrando que tal norma está intrinsecamente ligada ao parâmetro  $s$  da regressão, que por sua vez representa a distância da origem até o extremo do convexo de restrição no espaço paramétrico. Como exemplo, a Figura 56 representa um gráfico feito com um conjunto de dados de 20 covariáveis. Ao

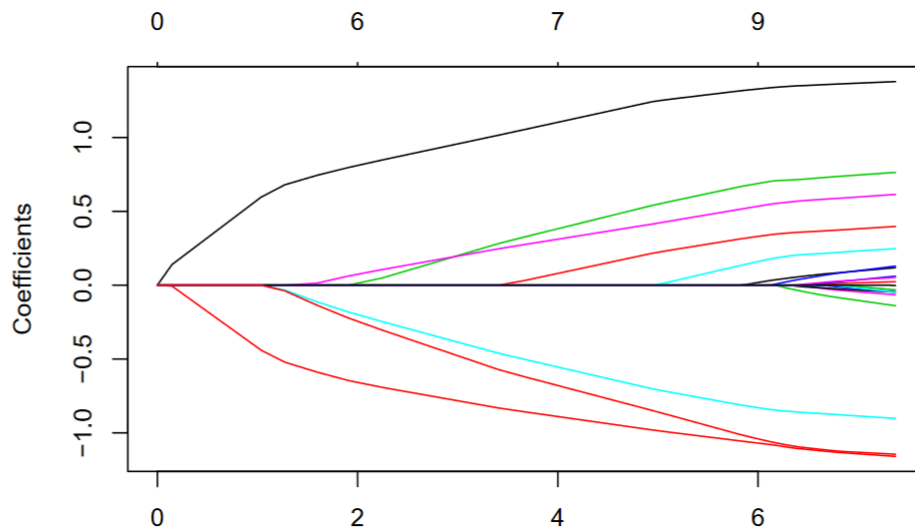


Figura 56 Estimação de parâmetros *Elastic Net*.

se incluir na função plot o argumento `label = TRUE`, do lado esquerdo do gráfico plotado irão aparecer legendas representando cada uma das covariáveis. Na Figura 57 é possível ver no lado direito da imagem os números que indicam a qual covariável cada uma das linhas traçadas representa. Uma vez realizado o ajuste, é

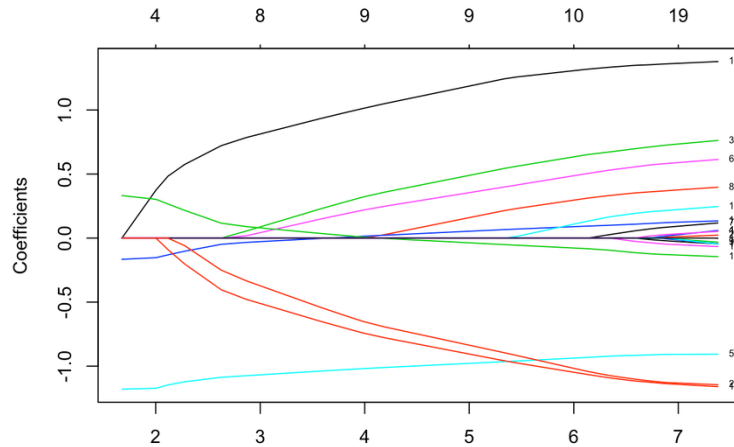


Figura 57 Exemplo de um plot com legenda.

possível se determinar uma estimativa específica dado um valor de  $s$ . Assim, ao se digitar

```
coef(reg, s=0.1)
```

temos como resultado uma estimativa pontual do vetor  $\beta$ , a qual foi extraída do ajuste reg. Uma pergunta natural é então como se determinar qual dentre as várias estimativas dadas por um ajuste é a mais indicada para a escolha do modelo ajustado? O pacote fornece tal estimativa baseado no conceito de *Generalized Cross Validation* (GCV), indicando qual das estimativas resulta em um modelo com Erro Quadrático Preditivo menor, no sentido generalizado. Para tanto, é necessário a criação de um novo objeto, do tipo `cv.glmnet`. Tal objeto armazena informações sobre o processo de Cross-Validation executado pelo pacote. Assim, para se de-

terminar o vetor estimativa com menor erro preditivo, segundo o GCV, deve-se executar os comandos da forma

```
cvreg = cv.glmnet(X, y, alpha=0.5)
coef(cvreg, s = "lambda.min")
```

Este último comando fornece os coeficientes do vetor estimativa  $\hat{\beta}$  com menor GCV, que é então utilizado como vetor de coeficientes para o modelo selecionado.

Como aplicação em dados reais, foram escolhidos dados genéticos referentes a mensuração de pH carne de suínos *Sus Scrofa* 45 min e 24 h após o abate. Foram então analisadas as covariáveis genéticas em termos de regressão Elastic Net, de forma a se agrupar covariáveis correlacionadas e se determinar, através do método, aquelas mais relevantes à variável resposta pH. Um estudo semelhante, porém utilizando *Partial Least Squares* (PLS), sobre os mesmos dados pode ser encontrado em Silveira et al. (2017). Mais detalhes referentes aos dados serão apresentados nas subseções seguintes.

## 8.1 Simulações no R

Foram realizadas simulações para se testar a aplicabilidade do pacote. Como exemplificação, uma primeira simulação foi feita aos moldes de um exemplo idealizado no artigo original (ZOU; HASTIE 2005, p. 313). Neste exemplo, o objetivo principal é mostrar a diferença básica entre os métodos LASSO e Elastic Net. Foram então criadas duas variáveis aleatórias independentes  $Z_1$  e  $Z_2$ , ambas com uma distribuição  $U(0, 20)$ . A resposta  $y$  foi gerada por uma  $N(Z_1 + 0.1Z_2, 1)$ . É importante destacar que, na resposta  $y$ , os pesos das variáveis aleatórias são bem distintos, já que a variável  $Z_2$  tem peso dez vezes menor do que a variável  $Z_1$ , ou seja,  $Z_2$  contribui bem menos para a resposta. Suponha que foram observadas seis covariáveis “latentes”  $x_1, \dots, x_6$  identicamente distribuídas da forma que, para

erros  $\varepsilon_i$  independentes e identicamente distribuídos por uma  $N(0, 1/16)$  se tem a relação

$$\begin{aligned} \mathbf{x}_1 &= \mathbf{Z}_1 + \varepsilon_1, & \mathbf{x}_2 &= -\mathbf{Z}_1 + \varepsilon_2, & \mathbf{x}_3 &= \mathbf{Z}_1 + \varepsilon_3 \\ \mathbf{x}_4 &= \mathbf{Z}_2 + \varepsilon_4, & \mathbf{x}_5 &= -\mathbf{Z}_2 + \varepsilon_5, & \mathbf{x}_6 &= \mathbf{Z}_2 + \varepsilon_6. \end{aligned}$$

Claramente, as covariáveis  $\mathbf{x}_1$ ,  $\mathbf{x}_2$  e  $\mathbf{x}_3$  formam um grupo cujo fator subjacente é a variável  $\mathbf{Z}_1$ , enquanto  $\mathbf{x}_4$ ,  $\mathbf{x}_5$  e  $\mathbf{x}_6$  possuem  $\mathbf{Z}_2$  como fator subjacente. A Figura 58 (ZOU; HASTIE 2005, p. 311) apresenta os traços do LASSO (a) e do Elastic Net (b). É possível ver que, apesar de o LASSO ter o efeito de seleção de covariáveis (pois as linhas do traço surgem de forma sequencial, ao contrário do Ridge em que estas entram simultaneamente), cada linha (relativa a uma covariável) entra separada das demais. Já no traço do Elastic Net, é possível perceber claramente o efeito de agrupamento. Além de os traços das covariáveis  $\mathbf{x}_1$ ,  $\mathbf{x}_2$  e  $\mathbf{x}_3$  surgirem no mesmo momento, o método ainda as apresenta antes das três últimas covariáveis, ou seja, para um valor de  $s$  menor. Isso se deve ao fato de que as três primeiras covariáveis estão relacionadas à variável  $\mathbf{Z}_1$ , que possui uma influência na resposta 10 vezes maior do que a variável  $\mathbf{Z}_2$ . Além do mais, pode-se perceber no Elastic Net trace que a linha referente à variável  $\mathbf{x}_2$  aparece de forma decrescente, bem como a linha referente à variável  $\mathbf{x}_5$ . Isto ocorre devido a natureza destas duas variáveis, a qual é influenciada de forma negativa pelas covariáveis  $\mathbf{Z}_1$  e  $\mathbf{Z}_2$ , respectivamente. Nota-se que o LASSO, para esse exemplo específico, não foi capaz de perceber esta natureza decrescente referente a variável  $\mathbf{x}_5$ .

Em resumo o exemplo anterior, apesar de hipotético, ilustra de forma clara as características do método Elastic Net em relação ao LASSO. O efeito de agrupamento, dependendo do tipo de aplicação, pode se apresentar como uma ferramenta crucial e inovadora na análise de dados em que se tem uma gama muito grande de

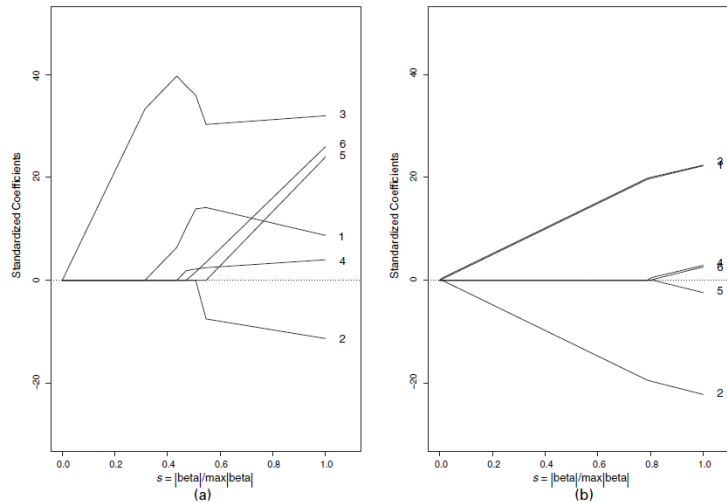


Figura 58 Comparação entre LASSO Trace e Elastic Net Trace.

covariáveis.

O código em R da simulação anterior se encontra no Apêndice na parte referente ao **Código 1**.

Uma segunda simulação (**Código 2** no Apêndice) foi feita no sentido de verificar os efeitos do Elastic Net não em variáveis latentes, mas sim diretamente nas covariáveis resposta. Foram então criadas 5 covariáveis, de forma que a variável  $x_1$  tenha distribuição  $N(0, 1)$ . A covariável  $x_2$  também é normal, porém com média  $x_1$  e desvio padrão  $sd = 0.1$ , ou seja, é altamente correlacionado com a variável  $x_1$ .  $x_3$ , independente de  $x_1$  e  $x_2$ , também segue uma normal  $N(0, 1)$  e  $x_4$ , também normal, possui média  $x_3$  e desvio padrão de  $sd = 0.008$ . Apenas  $x_5$  foi criada de forma a depender de duas das covariáveis anteriores, seguindo então uma normal  $N(x_4 + x_3, 0.005)$ . O vetor resposta  $y$  segue uma distribuição normal com desvio padrão 1, porém duas situações foram analisadas: uma em que  $y$  possui média  $0.01x_1 + x_2 + 2x_3 - x_4 + 2x_5$  e outra em que se altera o peso da variável  $x_4$  na média de  $y$  de -1 para -10. As Figuras 59 e 60 apresentam, respectivamente,

o traço do Elastic Net para as duas situações. Comparando ambas as situações por

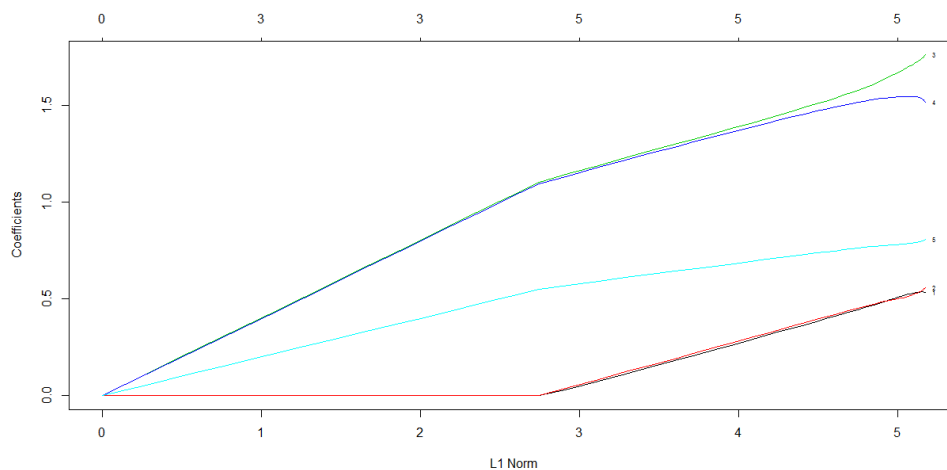


Figura 59 Elastic Net trace para a situação 1.

seus respectivos gráficos, ao se aplicar um peso relativamente grande na covariável  $x_4$ , as outras covariáveis correlacionadas a esta ( $x_3$  e  $x_5$ ), mesmo tendo coeficientes positivos, apresentam comportamento similar ao da covariável de coeficiente negativo ( $x_4$ ). Tal fato demonstra como o efeito de agrupamento no Elastic Net é considerável.

No **Código 3** (Apêndice), foi reproduzido a plotagem de um exemplo clássico, que primeiro surge em Tibshirani (1996), referente a dados de 8 covariáveis relacionadas ao câncer de próstata (*prostate cancer*) e uma variável resposta, **lpsa**, acrônimo de *log prostate specific antigen*. Os dados, presentes no pacote `lasso2`, estão armazenados em um *dataframe* de nome `Prostate`. O gráfico da variável resposta em função das 8 covariáveis é apresentado na Figura 61. Conforme se pode observar, as covariáveis com maior peso na resposta são `lcavol`, `svi` e `lweight`. Posteriormente, Zou; Hastie (2005) utilizaram este mesmo conjunto de dados no artigo do Elastic Net para se fazer um comparativo entre os dois métodos.



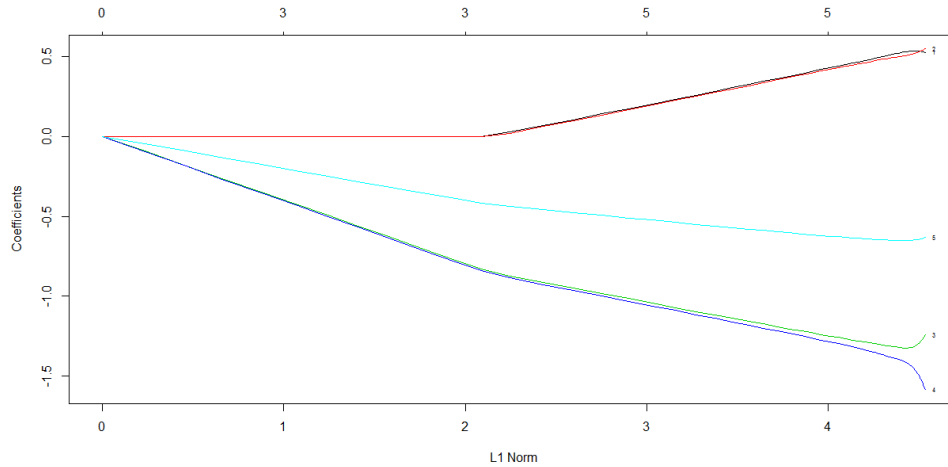


Figura 60 Elastic Net trace para a situação 2.

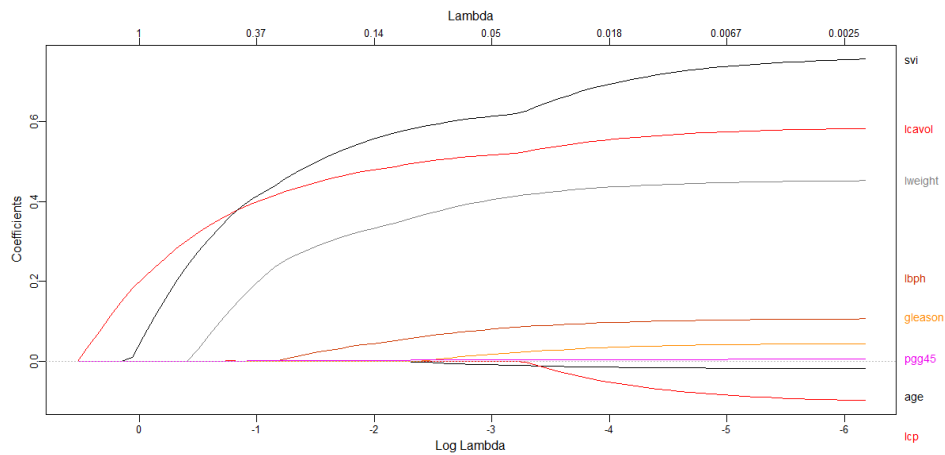


Figura 61 Elastic Net trace para os dados de câncer de próstata.

Uma última simulação (**Código 4** no Apêndice) foi feita com o intuito de se analisar o efeito de covariáveis categóricas no processo dos métodos LASSO e Elastic Net. Tal simulação se faz de interesse uma vez que, em geral, ao se coletar dados genéticos referentes a marcadores moleculares (SNPs), os resultados são

categóricos (heterozigoto, homozigoto dominante e homozigoto recessivo). Um grupo de 8 covariáveis categóricas (com níveis 0, 1 e 2) foi criado por meio de cópulas, de forma que as três primeiras covariáveis formem um grupo correlacionado, as três próximas covariáveis (4, 5 e 6) formem um segundo grupo correlacionado (independente do primeiro) e as duas últimas covariáveis formam um terceiro grupo correlacionado entre si. Em seguida, uma variável resposta foi produzida da forma

$$y = 130x_1 + 120.1x_2 + 150x_3 - 15x_4 + 0.2x_5 - 0.15x_6 + 60x_7 - 59x_8 + \varepsilon,$$

em que  $\varepsilon \sim N(0, 0.1)$ . O gráfico do Elastic Net trace é conforme Figura 62.

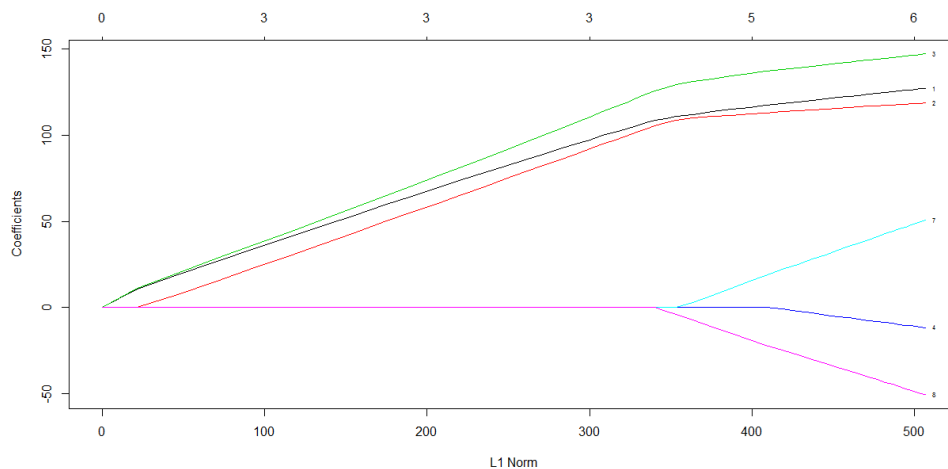
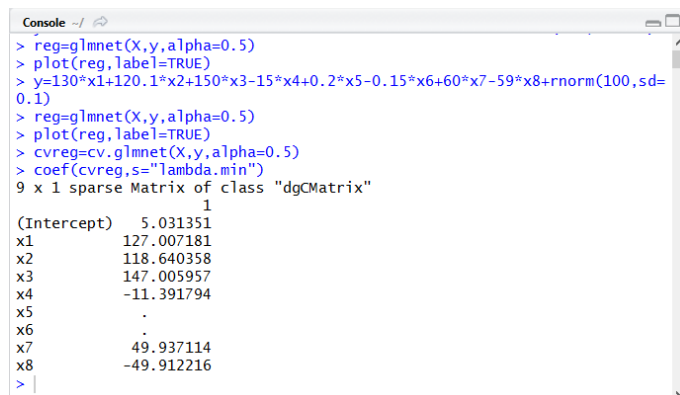


Figura 62 Elastic Net trace para covariáveis categóricas.

Alguns pontos sobre esta última simulação merecem destaque. Primeiramente, e mais importante, é possível dizer que o método Elastic Net (em particular o pacote glmnet) opera normalmente com covariáveis categóricas. Portanto, a aplicação do método em dados categóricos é perfeitamente possível. Novamente, se pode perceber que pela Figura 62 as covariáveis mais relevantes para o vetor

resposta  $y$  surgem primeiro no Elastic Net trace, ou seja, as covariáveis 1, 2 e 3. Além do mais, o efeito de agrupamento se manteve inalterado em covariáveis categóricas. Outro ponto que merece destaque é o fato de, no segundo grupo de covariáveis correlacionadas (4, 5 e 6), apenas a covariável 4 aparece no gráfico. As outras duas covariáveis (5 e 6) não aparecem uma vez que o valor de seus coeficientes (ou pesos) são praticamente desprezíveis em relação aos valores dos demais coeficientes. Fato que corrobora esta última afirmação é o valor da estimativa que resulta no menor erro quadrático, por meios de GCV, já que as estimativas dos coeficientes 5 e 6 são zeradas por este critério. A Figura 63 apresenta os valores estimados por GCV para a simulação em R.



```

Console -/ ↻
> reg=glmnet(X,y,alpha=0.5)
> plot(reg,label=TRUE)
> y=130*x1+120.1*x2+150*x3-15*x4+0.2*x5-0.15*x6+60*x7-59*x8+rnorm(100,sd=
0.1)
> reg=glmnet(X,y,alpha=0.5)
> plot(reg,label=TRUE)
> cvreg=cv.glmnet(X,y,alpha=0.5)
> coef(cvreg,s="lambda.min")
9 x 1 sparse Matrix of class "dgCMatrix"
      1
(Intercept)  5.031351
x1           127.007181
x2           118.640358
x3           147.005957
x4           -11.391794
x5            .
x6            .
x7            49.937114
x8           -49.912216
>

```

Figura 63 Estimativa do vetor de parâmetros dos dados categóricos.

## 8.2 Aplicação em dados genéticos

Como aplicação em dados reais, optou-se por adotar dados genéticos referentes a mensuração de pH da carne e peso de carcaça de suínos *Sus scrofa*. Os dados são constituídos por 345 observações e 300 covariáveis, as quais estão alocadas nas colunas do arquivo dos dados. Entretanto, os 237 marcadores genéticos (SNPs) estão concentrados entre as colunas 64 e 300 (ALGA0000087 ~ ASGA0080951)

do arquivo de dados. Conforme descrito em Silveira et. al. (2017), os 237 marcadores estão distribuídos em seis cromossomos do *Sus scrofa* da seguinte maneira: SSC1(56), SSC4(54), SSC7(59), SSC8(30), SSC17(25) e SSCX(13). Apesar de todas as propriedades dos métodos LASSO e Elastic Net, a análise gráfica simultânea das 237 covariáveis se torna um tanto quanto impraticável, conforme apresentado na Figura 64, em que foi escolhida como variável resposta o pH mensurado a 45 minutos após o abate.

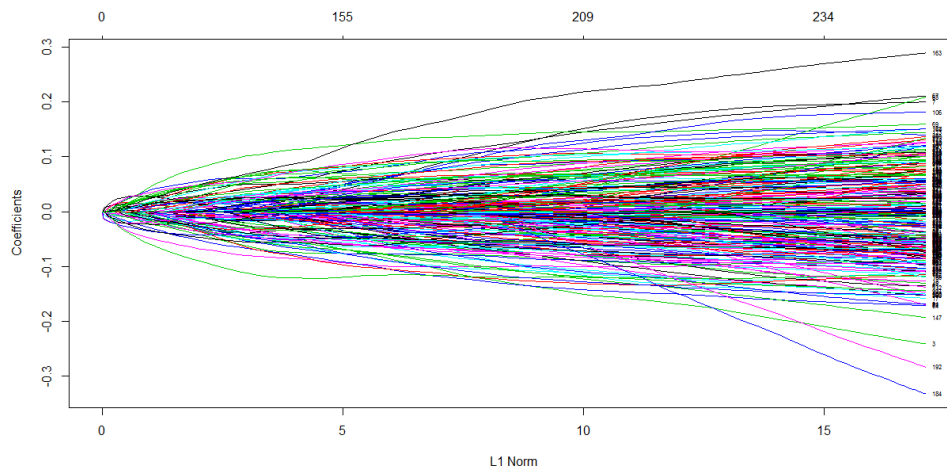


Figura 64 Elastic Net trace das 237 covariáveis.

Desta forma, o código em R foi implementado de modo que os índices  $j$  e  $k$  selecionem, respectivamente, a primeira e a última coluna da matriz  $X$  utilizada na análise, respeitando a relação  $1 \leq j < k \leq 237$ . As linhas de comando do código constam no Apêndice (**Código 5**). As linhas de comando

```
X<-X[-c(which(is.na(y))), ]
y=y[-c(which(is.na(y)))]
```

foram incluídas em função dos valores das respostas, já que algumas observações

possuem o valor NA, provavelmente para representarem parcelas perdidas. Assim, os comandos anteriores excluem tanto do vetor resposta  $y$  quanto da matriz  $X$  as linhas cujas entradas relativas ao vetor  $y$  contenham o valor NA. Como referência, a partir de agora os marcadores serão denominados de acordo com sua posição a partir da coluna 64, isto é, o marcador 3 é aquele referente a coluna 66 dos dados brutos.

## 9 Resultados e discussão

Buscou-se determinar dentre todos os 237 marcadores os 12 (aproximadamente 5%) que mais influenciam na resposta peso de carcaça PCARC. Novamente, a Figura 65 representa a dificuldade de se realizar tal análise graficamente. Para se contornar tal problema foi utilizada a ferramenta que informa o valor de

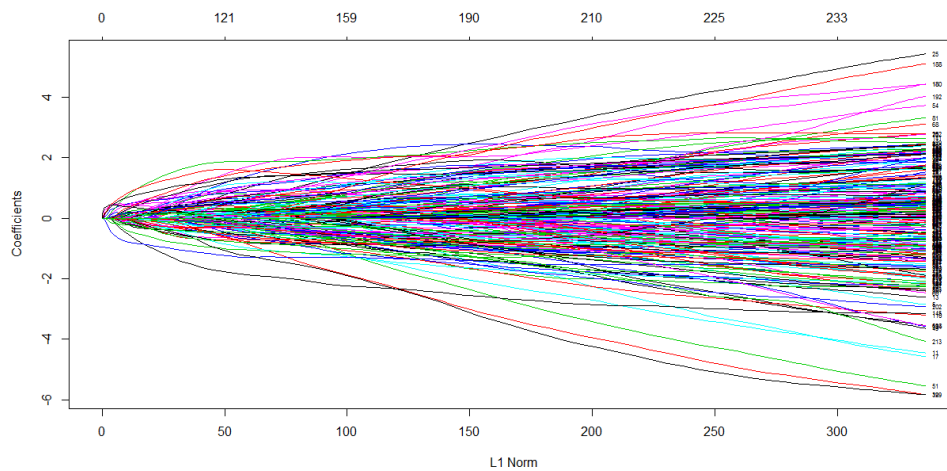


Figura 65 Elastic Net trace para as 237 covariáveis com a resposta PCARC.

uma estimativa específica para um dado valor de  $s$  e então, atribuindo valores altos para  $s$  (o que implica em valores baixos da norma  $L1$ ) e fazendo-o decrescer gradativamente, as covariáveis mais importantes surgem também de forma gradativa.

Por esta metodologia, as 12 covariáveis selecionadas foram as de número 36, 82, 129, 144, 148, 162, 190, 206, 208, 211, 219 e 229. Plotando estas covariáveis em um Elastic Net trace tem-se a Figura 66. O código para a execução desta última análise sofreu algumas alterações (em relação ao **Código 5**), as mais relevantes sendo nas linhas que montam a matriz  $X$  cujas colunas são formadas pelas 12 covariáveis selecionadas, e também foram excluídos os comandos que retiram

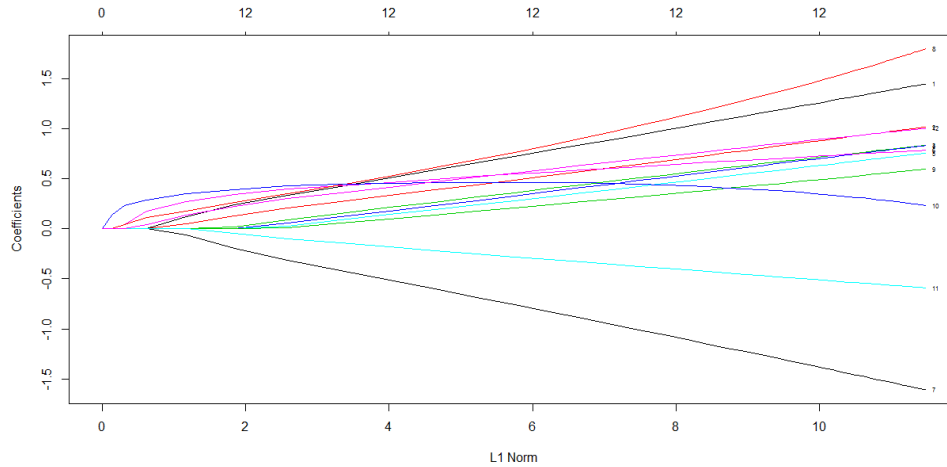


Figura 66 As 12 covariáveis mais relevantes para a resposta PCARC.

linhas da matriz  $X$  cujos respectivos elementos em  $y$  sejam NA, uma vez que utilizando PCARC como resposta, não se tem nenhuma parcela perdida. O código sumarizado está no Apêndice, **Código 6**. A Figura 67 representa de forma esquemática a localização, através das setas vermelhas, de cada uma das 12 covariáveis selecionadas na sequência total dos 220 marcadores, bem como suas posição relativa a cada um dos 6 cromossomos. Uma análise da Figura 67 mostra que a grande

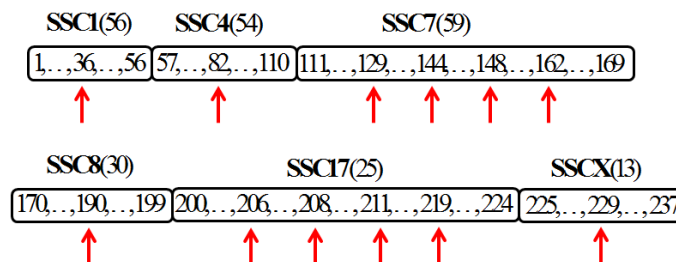


Figura 67 Posição das 12 covariáveis mais relevantes em relação aos grupos de ligação.

maioria dos marcadores (dentre os 12 mais influentes para a resposta PCARC)

estão presentes nos cromossomos SSC7 e SSC17.

Foi realizada também uma análise similar utilizando como resposta o fenótipo pH45. Desta forma, 13 covariáveis (5%) foram selecionadas como sendo mais relevantes dentre as 237, a se saber, 11, 13, 49, 54, 64, 86, 126, 133, 154, 157, 186, 199 e 234. As linhas de comando do R constam no Apêndice (**Código 7**). As curvas Elastic Net para tal situação estão representadas na Figura 68.

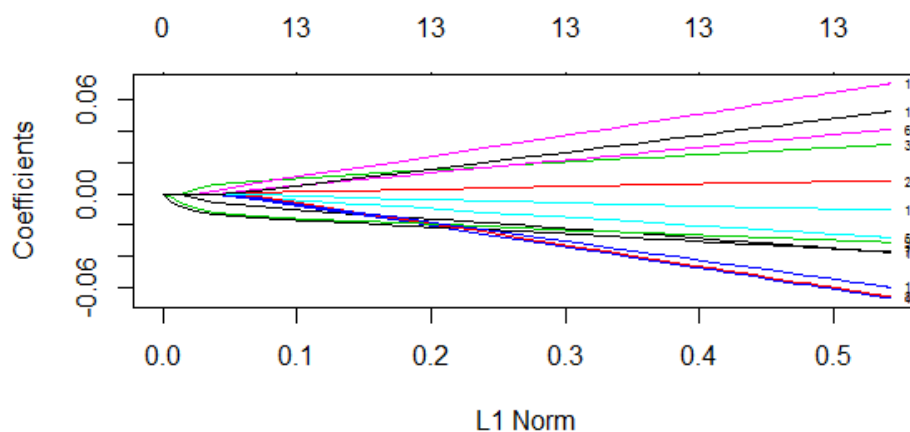


Figura 68 As 13 covariáveis mais relevantes para a resposta PH45.

No sentido de se observar melhor as covariáveis em um cromossomo específico, optou-se por buscar as 9 covariáveis mais relevantes no grupo de ligação SSC4, composto por 54 marcadores genéticos (colunas 57 a 110). O traço Elastic Net para esta situação é conforme Figura 69. Os comandos em R constam no Apêndice (**Código 8**). O cromossomo SSC4 foi escolhido em detrimento dos demais para que possa ser feita uma análise comparativa com os resultados de Silveira et al. (2017). Conforme se pode observar, os 9 marcadores selecionados no grupo de ligação SSC4 foram 64, 67, 69, 72, 76, 81, 86, 87 e 101.

Estes resultados são discrepantes com aqueles obtidos por Silveira, op.



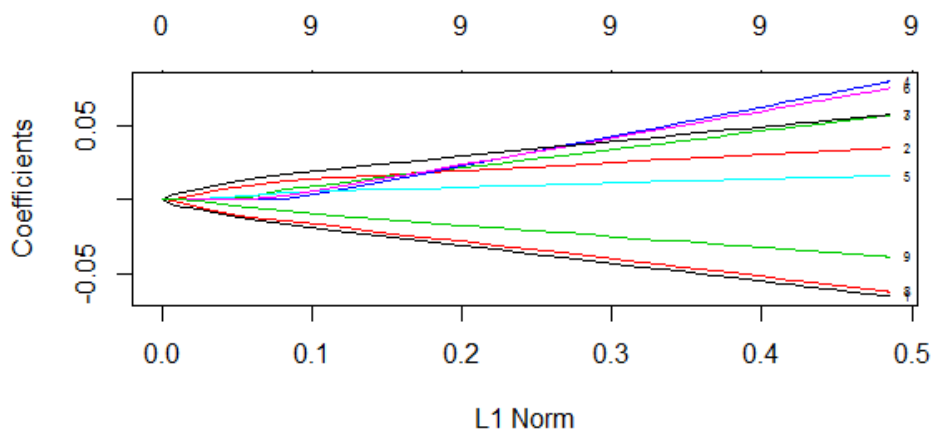


Figura 69 As 9 covariáveis mais relevantes para a resposta PH45 no cromossomo SSC4.

cit., uma vez que o método PLS identificou como mais relevantes os marcadores 62, 94, 98, 91 e 93. Mesmo a análise específica no cromossomo SSC4 observa-se que os marcadores selecionados pelo método PLS diferem daqueles selecionados pelo método Elastic Net.

Para uma análise mais específica, o método Elastic Net foi aplicado aos marcadores entre 90 e 99 (localizados no cromossomo 4). A Figura 70 apresenta o gráfico Elastic Net para estes 10 marcadores, no o caso em que a variável resposta é o pH 45 minutos após o abate. No eixo das abscissas, foi colocado entre parêntesis o número dos marcadores, de forma a facilitar a visualização e interpretação dos resultados. Conforme se pode observar, os marcadores 98 e 94 possuem maior influência na resposta quando comparado neste grupo de 10 covariáveis, inclusive apresentando um efeito de agrupamento acentuado. Este resultado segue em conformidade com o artigo citado. Em seguida, as covariáveis 96, 93 e 92 formam um segundo grupo correlacionado e apresentam significativa relevância

para o resultado, sendo que o marcador 93 foi identificado como de alta relevância pelo referido artigo. Fica claro também pelo gráfico que existe um agrupamento entre os marcadores 99 e 97, porém o peso destes na resposta é bem menor quando comparado com os demais, uma vez que estes aparecem mais a direita no eixo das abscissas. Um resultado evidenciado pelo método Elastic Net e não captado pelo

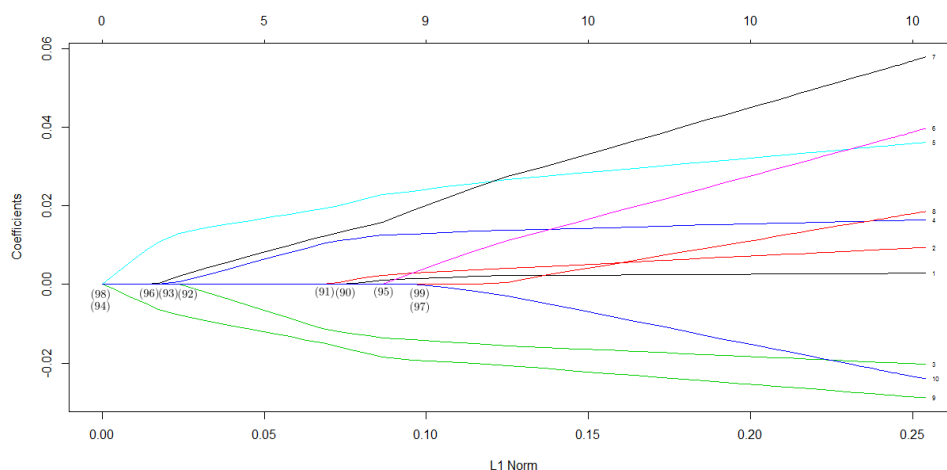


Figura 70 Elastic Net trace para os marcadores de 90 a 99.

método PLS, para estas 10 covariáveis analisadas, é que os marcadores 96 e 92 apresentam efeito similar ao marcador 93, em função do efeito de agrupamento destas três covariáveis. Entretanto, o valor da estimativa da covariável 92 na resposta apresenta valor negativo, enquanto que a estimativa das covariáveis 96 e 93 possuem valores positivos.

O método do Elastic Net, de forma análoga aos exemplos anteriores, pode ser aplicado em qualquer uma das várias respostas contidas no arquivo de dados (PCD, RCARC, PROLOM, PCOPA, etc.). Entretanto, o presente trabalho tomou apenas as respostas pH45 e PCARC para análise, uma vez que o objetivo principal de tais análises é apresentar a aplicabilidade do método Elastic Net, bem

como suas vantagens e características. Com os códigos que constam no Apêndice, uma vasta gama de análises pode ser feita com o conjunto de dados utilizado, bastando apenas pequenas alterações nas linhas de comando apresentadas. O importante a se destacar é que o método Elastic Net se apresenta como uma ferramenta bastante aplicável na análise de dados genéticos, principalmente no que tange a seleção de covariáveis por grupos correlacionados.

## 10 CONCLUSÃO

- Uma abordagem geométrica à teoria dos estimadores LASSO, LARS e Elastic Net se revelou eficiente no sentido de ser intuitiva e explicitar aspectos essenciais da teoria.
- O texto representa uma contribuição para a leitura dos artigos clássicos da teoria LASSO, LARS e Elastic Net.
- As rotinas desenvolvidas utilizando o pacote `glmnet` apresentaram bons resultados e podem ser utilizadas para futuras análises em dados genéticos.
- Os resultados obtidos na análise dos dados genéticos de *Sus scrofa* apresentam uma relevante similaridade em relação aos resultados de Silveira et al. (2017).
- Novos resultados são apresentados na detecção de SNPs relevantes na variável resposta Peso de Carcaça.

## REFERÊNCIAS

BOLDRINI, J. L. et al. **Álgebra Linear**. 2. ed. São Paulo: Harbra, 1986. 411 p.

BREIMAN, L. Better subset regression using the nonnegative garrote. **Technometrics**, 1995, vol. 37, n. 4. p. 373 to 384

COSTA, L. A. Novo estimador de cumeieira de Rao com aplicação em seleção genômica. **Tese (Mestrado em Estatística e Experimentação Agropecuária)** - Universidade Federal de Lavras, Lavras-MG, 2015., p. 126.

EFRON, B; HASTIE, T; JOHNSTONE, I; TIBSHIRANI, R. Least angle regression. **The Annals of Statistics**, 2004, vol. 32, n. 2. p. 407 to 499

EFRON, B. The estimation of prediction error: covariance penalties and cross-validation. **Journal of the American Statistical Association**, 2004, vol. 99, n. 467. p. 619 to 632

GAJO, C. A. Propriedades e aspectos geométricos de estimadores tipo James-Stein e do estimador de Hartigan. **Tese (Doutorado em Estatística e Experimentação Agropecuária)** - Universidade Federal de Lavras, Lavras-MG, 2016, p. 194.

GENTLE, J. E. **Matrix algebra: theory, computations and applications in statistics**. New York: springer, 2007. 528 p.

GOLUB G. H., HEATH, M., WAHBA, G. Generalized cross-validation as a method for choosing a good ridge parameter. **Thecnometrics**, 1979, vol. 21, n. 21. p. 215 to 223

GRUBER, M. H. J. **Improving efficiency by shrinkage: the James-Stein and ridge regression estimator**. 2nd ed. New York: M. Dekker, 1998. 632 p.

HASTIE, T; TIBSHIRANI, R; FRIEDMAN, J. **The elements of statistical learning. Data mining, Inference and Prediction**. New York: Springer, 2nd ed.

2008. 739 p.

HOERL, A.E; KENNARD, R. W. Ridge regression: Applications to nonorthogonal problems. **Thecnometrics**, 1970, vol. 12, n. 1. p. 69 to 82

HOERL, A.E; KENNARD, R. W. Ridge regression: Biased estimation for nonorthogonal problems. **Thecnometrics**, 1970, vol. 12, n. 1. p. 51 to 67

JAMES, G; WITTEN, D; HASTIE, T; TIBSHIRANI, R. **An introduction to statistical learning with application in R**. New York: Springer, 2013. 426 p.

KATO, K. On the degrees of freedom in shrinkage estimation. **Journal of Multivariate Analysis**, v. 100, p. 1338-1352, 2009.

PEREIRA, L. S. Abordagem geométrica à teoria dos modelos de Gauss-Markov. **Dissertação (Mestrado em Estatística e Experimentação Agropecuária) - Universidade Federal de Lavras, Lavras-MG, 2013, p. 130.**

PEREIRA, L. S.; CHAVES, L. M, et al.; Geometry of Basic Properties on Linear Regression and The Mallow's Cp Statistics. **Revista Brasileira de Biometria**, v. 33, n. 3, p. 357-377, 2015.

RENCHER, A. C; SCHAALJE, G.B. **Linear Models in statistics**. New Jersey: John Wiley & Sons, 2008. 672 p.

SAVILLE, D. J; WOOD, G. L. A Method for Teaching Statistics Using N-Dimensional Geometry. **The American Statistician**, v. 40, p. 205-214, 1986.

SAVILLE, D. J; WOOD, G. L. **Statistical Methods: The Geometric Approach**. New York: Springer-Verlag, 1991. 560 p.

SILVEIRA, F. G; COSTA, L. A.; PEREIRA, L. S, et al.; Uma demonstração geométrica para uma identidade de Fisher para o modelo de dois fatores. **Revista Brasileira de Biometria**, v. 30, n. 2, p. 199-222, 2012.

SILVEIRA, F. G; DUARTE, D. A. S.; CHAVES, L. M; SILVA, F. F; FILHO, I. C; DUARTE, M. S; LOPES, P. S; GUIMARÃES, S. E. F; The optimal number of partial least squares components in genomic selection for pork pH **Ciência Rural, Santa Maria**, v. 47, No.1, p. 1-5, 2017.

TIBSHIRANI, R. J; TAYLOR, J. Degrees of freedom in lasso problems. **The Annals of Statistics**, v. 40, No.2, p. 1198-1232, 2012.

TIBSHIRANI, R. Regression Shrinkage and Selection via the LASSO. **Journal of the Royal Statistical Society**, v. 58, No.1, p. 267-288, 1996.

ZOU, H; HASTIE, Y. Regularization and variable selection via the elastic net. **Journal of the Royal Statistical Society**, v. 67, No.2, p. 301-320, 2005.

## APÊNDICE

### Código 1

```
library(glmnet)
z1=runif(100, min = 0, max = 20)
z2=runif(100, min = 0, max = 20)
x1=z1+rnorm(100,mean=0,sd=1/16)
x2=-z1+rnorm(100,mean=0,sd=1/16)
x3=z1+rnorm(100,mean=0,sd=1/16)
x4=z2+rnorm(100,mean=0,sd=1/16)
x5=-z2+rnorm(100,mean=0,sd=1/16)
x6=z2+rnorm(100,mean=0,sd=(1/16))
u<-cbind(x1,x2,x3,x4,x5,x6)
y=rnorm(100,mean=z1+0.1*z2,sd=1)
reg=glmnet(u,y,alpha=0.5)
plot(reg,label=TRUE)#ELASTICNET
reg=glmnet(u,y,alpha=1)
plot(reg,label=TRUE)#LASSO
```

### Código 2

```
library(glmnet)
rm(list=ls())
x1<-rnorm(100)
x2<-rnorm(100,mean=x1,sd=.01)
x3<-rnorm(100)
x4<-rnorm(100,mean=x3,sd=.008)
x5<-rnorm(100,mean=(x4+x3),sd=.005)
```



```

u<-cbind(x1,x2,x3,x4,x5)
#####
y<-rnorm(100,mean=0.01*x1+x2+2*x3-x4+2*x5)
reg=glmnet(u,y,alpha=0.5)
plot(reg,label=TRUE)
#####
y<-rnorm(100,mean=0.01*x1+x2+2*x3-10*x4+2*x5)
reg=glmnet(u,y,alpha=0.5)
plot(reg,label=TRUE)

```

### **Código 3**

```

install.packages("lasso2")
library(lasso2)#dados do Prostate Cancer
library(glmnet)
rm(list=ls())
prost=data(Prostate)
ls(Prostate) #conteudo do Prostate Data
X = as.matrix(Prostate[setdiff(colnames(Prostate),
"lpsa")])
#o comando anterior retira "lpsa" dos dados e monta
#a matriz X com os demais
y<-Prostate$lpsa
reg=glmnet(X,y,alpha=0.5)
plot(reg,label=TRUE)

```

### **Código 4**

```

library(glmnet)

```

```
install.packages("copula")
library(copula)
rm(list=ls())
X<-cbind(floor(rCopula(n = 100,copula =
normalCopula(param = .8,dim = 3))*3),
floor(rCopula(n = 100,copula =
normalCopula(param = .82,dim = 3))*3),
floor(rCopula(n = 100,copula =
normalCopula(param = .75,dim = 2))*3))
#o comando anterior cria um grupo de
#correlação 0.8 nas três primeiras
#covariáveis, 0.82 nas três seguintes
#e 0.75 nas duas últimas.
x1<-X[,1]
x2<-X[,2]
x3<-X[,3]
x4<-X[,4]
x5<-X[,5]
x6<-X[,6]
x7<-X[,7]
x8<-X[,8]
X<-cbind(x1,x2,x3,x4,x5,x6,x7,x8)
y=130*x1+120.1*x2+150*x3-15*x4+0.2*x5-0.15*x6
+60*x7-59*x8+rnorm(100,sd=0.1)
reg=glmnet(X,y,alpha=0.5)
plot(reg,label=TRUE)
```

```
cvreg=cv.glmnet(X,y,alpha=0.5)
coef(cvreg,s="lambda.min")
print(reg)
```

### **Código 5**

```
rm(list=ls())
library(glmnet)
gendad=read.table("C:/Doutorado/DoutFazendo/Defesa
/GenGWS.txt",header = TRUE)
names(gendad)
C=gendad[,64:300]
colunas=c()
for(i in 1:237){
  colunas[i]=paste("c",i,sep="")
}
colnames(C)=colunas
y=gendad[,11]#a coluna 11 eh o pcarc
j=90
k=100
X=C[,j:k]
X<-X[-c(which(is.na(y))),]
y=y[-c(which(is.na(y)))]
X=as.matrix(X)
reg=glmnet(X,y,alpha=0.5)#nlambda=500
plot(reg,label=TRUE)
```

**Código 6**

```
rm(list=ls())
library(glmnet)
gendad=read.table("C:/Doutorado/DoutFazendo/Defesa
/GenGWS.txt",header = TRUE)
names(gendad)
C=gendad[,64:300]
colunas=c()
for(i in 1:237){
  colunas[i]=paste("c",i,sep="")
}
colnames(C)=colunas
y=gendad[,5]#a coluna 5 eh o pcarc
j=1
k=237#ultimo=237
X=C[,j:k]
X=as.matrix(X)
reg=glmnet(X,y,alpha=0.5)#nlambda=500
plot(reg,label=TRUE)
coef(reg,s=1.2)
colSums(coef(reg,s=1.2) != 0)#conta quantos coef
#são não nulos.
a=c(36,82,129,144,148,162,190,206,208,211,219,229)
#a linha anterior cria um vetor com
#os 12 marcadores mais relevantes.
X=cbind(C[,a])
```

```

X=as.matrix(X)
y=gendad[,5]#a coluna 5 eh o pcarc
reg=glmnet(X,y,alpha=0.5)#nlambda=500
plot(reg,label=TRUE)

```

### **Código 7**

```

rm(list=ls())
library(glmnet)
gendad=read.table("C:/Doutorado/posdefesa/Defesa/
Codigos/GenGWS.txt",header = TRUE)
names(gendad)
C=gendad[,64:300]
colunas=c()
for(i in 1:237){
  colunas[i]=paste("c",i,sep="")
}
colnames(C)=colunas
y=gendad[,52]#a coluna 52 eh o pH45
j=1
k=237
X=C[,j:k]
X<-X[-c(which(is.na(y))),]
y=y[-c(which(is.na(y)))]
X=as.matrix(X)
reg=glmnet(X,y,alpha=0.5)
plot(reg,label=TRUE)
colSums(coef(reg,s=0.051) != 0)#conta quantos coef

```

```

#são não nulos.
coef(reg, s=0.051)
a=c(11,13,49,54,64,86,126,133,154,157,186,199,234)
C=gendad[,64:300]
colunas=c()
for(i in 1:237){
  colunas[i]=paste("c",i,sep="")
}
colnames(C)=colunas
y=gendad[,52]#a coluna 52 eh o pH45
X=cbind(C[,a])
X<-X[-c(which(is.na(y))),]
y=y[-c(which(is.na(y)))]
X=as.matrix(X)
reg=glmnet(X,y,alpha=0.5)
plot(reg,label=TRUE)

```

### **Código 8**

```

rm(list=ls())
library(glmnet)
gendad=read.table("C:/Doutorado/posdefesa/Defesa/
Codigos/GenGWS.txt",header = TRUE)
names(gendad)
C=gendad[,64:300]
colunas=c()
for(i in 1:237){
  colunas[i]=paste("c",i,sep="")
}

```

```

}
colnames(C)=colunas
y=gendad[,52]
j=57
k=110
X=C[,j:k]
X<-X[-c(which(is.na(y))),]
y=y[-c(which(is.na(y)))]
X=as.matrix(X)
reg=glmnet(X,y,alpha=0.5)
plot(reg,label=TRUE)
colSums(coef(reg,s=0.031) != 0)#conta quantos coef
#são não nulos.
coef(reg,s=0.031)
a=c(64,67,69,72,76,81,86,87,101)
C=gendad[,64:300]
colunas=c()
for(i in 1:237){
  colunas[i]=paste("c",i,sep="")
}
colnames(C)=colunas
y=gendad[,52]
X=cbind(C[,a])
X<-X[-c(which(is.na(y))),]
y=y[-c(which(is.na(y)))]
X=as.matrix(X)

```

```
reg=glmnet(X,y,alpha=0.5)  
plot(reg,label=TRUE)
```