



FERNANDA GOMES DA SILVEIRA

**ABORDAGEM GEOMÉTRICA DO MÉTODO
DOS QUADRADOS MÍNIMOS PARCIAIS
COM UMA APLICAÇÃO A DADOS DE
SELEÇÃO GENÔMICA**

LAVRAS-MG

2014

FERNANDA GOMES DA SILVEIRA

**ABORDAGEM GEOMÉTRICA DO MÉTODO DOS QUADRADOS
MÍNIMOS PARCIAIS COM UMA APLICAÇÃO A DADOS DE 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

Coorientador

Dr. Fabyano Fonseca e Silva

LAVRAS-MG

2014

**Ficha Catalográfica Elaborada pela Coordenadoria de Produtos e
Serviços da Biblioteca Universitária da UFLA**

Silveira, Fernanda Gomes da.

Abordagem geométrica do método dos quadrados mínimos
parciais com uma aplicação a dados de seleção genômica /

Fernanda Gomes da Silveira. – Lavras : UFLA, 2014.

176 p. : il.

Tese (Doutorado) - Universidade Federal de Lavras, 2014.

Orientador: Lucas Monteiro Chaves.

Bibliografia.

1. Projeções. 2. Regressão. 3. Componentes principais. 4.
Suíno. I. Universidade Federal de Lavras. II. Título.

CDD – 519.536

FERNANDA GOMES DA SILVEIRA

**ABORDAGEM GEOMÉTRICA DO MÉTODO DOS QUADRADOS
MÍNIMOS PARCIAIS COM UMA APLICAÇÃO A DADOS DE 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 2014.

Dr. Daniel Furtado Ferreira	UFLA
Dr. Márcio Balestre	UFLA
Dr. André Luís da Costa Paiva	IFMG - Campus Bambuí
Dr. Antonio Policarpo Souza Carneiro	UFV
Dr. Fabyano Fonseca e Silva	UFV

Dr. Lucas Monteiro Chaves
Orientador

**LAVRAS-MG
2014**

Aos meus pais, Geraldo e Maria de Lurdes,
aos meus irmãos Fábio e Renata,
pelo amor incondicional.

DEDICO

AGRADECIMENTOS

A Deus, pelo dom da vida, por estar sempre me abençoando, me guiando e dando forças para vencer os obstáculos.

Aos meus pais, Geraldo e Maria de Lurdes pela oportunidade que me proporcionaram nos estudos, pelo apoio e pela compreensão nos momentos em que não pude estar presente e por serem sempre meu esteio e minha inspiração.

Aos meus irmãos, Fábio e Renata; aos cunhados, Renata e Norton; aos sobrinhos Lucas e Bernardo, pelo incentivo, força e carinho.

A toda minha querida família pelos constantes incentivos: avós, tios e tias, primos e primas.

À Universidade Federal de Lavras, pela oportunidade de aprimoramento acadêmico.

À Fundação de Amparo à Pesquisa do Estado de Minas Gerais (FAPEMIG) e à Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES), pelo apoio financeiro.

Ao Instituto Federal de Minas Gerais, campus Bambuí, pelo afastamento que tornou possível minha qualificação.

Ao professor Lucas Monteiro Chaves pela orientação, apoio e dedicação nestes quatro anos.

Ao professor Fabyano Fonseca e Silva pela co-orientação, amizade, confiança e por estar sempre bem disposto a ajudar.

Aos membros da banca examinadora, pelas críticas e sugestões a este trabalho.

Aos professores do Departamento de Ciências Exatas, pelos conhecimentos transmitidos.

Aos funcionários do Departamento de Ciências Exatas, pela prontidão em

ajudar, em especial a Josi, secretária da Pós-graduação.

Aos professores Paulo Sávio Lopes e Simone Eliza Facioni Guimaraes, responsáveis pela granja de melhoramento de suínos do DZO-UFV e Laboratório de Biotecnologia Animal da UFV, respectivamente, pela concessão dos dados utilizados nas aplicações deste trabalho.

Aos colegas do curso de Pós-graduação em Estatística e Experimentação Agropecuária, pelo convívio.

À amizade das companheiras de repúblicas: Adriana, Alcilene, Isabelle e Larissa.

À hospedagem e o carinho com que me acolheram em suas casas para que as viagens se tornassem menos exaustivas: Alice e João Paulo, Fátima e Diego, Élide e Edna, Luzia e Fábio, Alcilene (Elza, Júlia e Tatiane).

A todos que, direta ou indiretamente, contribuíram para a realização deste trabalho, muito obrigada!

“Quanto maior a dificuldade, tanto maior é o mérito em superá-la.”

Henry Ward Beecher

RESUMO

Quando, em uma regressão múltipla, tem-se relações lineares ou quase colinearidade entre as covariáveis ou, ainda, quando o número de covariáveis é maior que o número de observações, o método dos Quadrados Mínimos Ordinários pode não ser adequado. Neste contexto, o método dos Quadrados Mínimos Parciais (PLS) tem se mostrado eficiente. O método consiste em obter a redução de dimensionalidade, uma vez que a regressão é realizada em relação a componentes relevantes. O método é abordado na literatura principalmente sob dois aspectos: algorítmico e algébrico. Neste trabalho, uma abordagem geométrica, utilizando projeções ortogonais, é apresentada no sentido de explicitar todas as etapas teóricas e da construção do algoritmo PLS. No intuito de tornar o texto mais didático, a abordagem geométrica também foi aplicada na teoria da Regressão em Componentes Principais (PCR), uma vez que os métodos PLS e PCR são similares. Uma rotina foi desenvolvida usando o software R, visando também a explicitar o passo a passo da construção do algoritmo. Como em qualquer análise de redução de dimensionalidade, um passo importante na aplicação do método PLS é a determinação de um número ótimo de componentes. Para tal, foi apresentada a teoria de Graus de Liberdade e o método de Validação Cruzada. Os métodos PCR e PLS, além da regressão tradicional sem redução de dimensionalidade, foram aplicados em uma análise de seleção genômica em suínos considerando um painel de marcadores SNPs de baixa densidade e dois fenótipos relacionados com a qualidade da carne.

Palavras-chave: Abordagem Geométrica. Quadrados Mínimos Parciais. Regressão. Seleção Genômica Ampla.

ABSTRACT

When estimating the coefficients of a multiple regression model, if the vectors of predictors are highly correlated, meaning that one can be (almost) a linear combination of the others, or if the number of predictors is greater than the number of observations, the Ordinary Least Square method may be non-appropriate. In that case, the Partial Least Square (PLS) method has shown to be efficient. It consists of obtaining a reduction in dimension by restricting the regression to relevant components. It is usual the literature to be restricted to two main aspects: algorithmic and algebraic. In this work, a geometric approach, based in orthogonal projections, is used to explicit all the theory behind the PLS method as well as in the construction of the PLS algorithm. Aiming to make the text more didactic, the same approach is applied to the Principal Components Regression (PCR) method, since both PLS and PCR are similar. Also a step by step routine for the PLS algorithm was developed using the R software. As in any procedure of reduction of dimensionality, the determination of the optimal number of components is a key step. To do that, we have described and used the Degree of Freedom Method and the Cross Validation Method. Both the PLS and PCR methods, besides the usual regression with no dimensionality reduction, were applied to a genomic selection analysis of pigs, considering a panel of low density SNP markers and two phenotypes related to meat quality.

Keywords: Geometric Approach. Partial Least Squares. Regression. Genome Wide Selection.

LISTA DE FIGURAS

Capítulo 1	19
Figura 1 Transformação linear de \mathbb{R}^n em \mathbb{R}^m	24
Figura 2 Interpretação geométrica do produto interno	25
Figura 3 Projetor ortogonal de um vetor aleatório \mathbf{Y} sobre um subespaço W	26
Figura 4 Representação geométrica para o modelo linear de Gauss-Markov	32
Figura 5 Representação geométrica do vetor residual	33
Figura 6 Representação gráfica do núcleo e da imagem de \mathbf{X} e \mathbf{X}'	34
Figura 7 Modelo de Gauss-Markov quando \mathbf{X} não é injetiva	36
Figura 8 Representação geométrica da obtenção do componente principal	55
Figura 9 Representação geométrica da curva parametrizada pelo comprimento do arco sobre a esfera de raio unitário, centrada na origem do espaço \mathbb{R}^n	58
Figura 10 Representação geométrica da maximização da $\text{var}(\mathbf{X}\boldsymbol{\beta} \cdot \mathbf{Y})$ restrito a $\boldsymbol{\beta}'\boldsymbol{\beta} = 1$	59
Figura 11 Representação geométrica dos autovetores $\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2, \dots, \boldsymbol{\gamma}_p$ da matriz $\mathbf{X}'\mathbf{X}$	60
Figura 12 Subespaço W_m gerado pelos m -primeiros autovetores da matriz $\mathbf{X}'\mathbf{X}$	61
Figura 13 Projeção ortogonal de \mathbf{Y} no subespaço W_m definindo o parâmetro de regressão $\hat{\boldsymbol{\beta}}_{\text{PCR}}$	62
Figura 14 Representação da função injetiva f entre os conjuntos A e B e da função g como uma projeção na $\text{Im}f$	62
Figura 15 Representação da composição de uma função h sobrejetiva com uma função f injetiva, definindo uma bijeção $h \circ f$	63
Figura 16 Representação dos vetores $\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_m$ como imagem dos autovetores $\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2, \dots, \boldsymbol{\gamma}_m$ pela matriz \mathbf{X}	64
Figura 17 Representação da transformação linear $\boldsymbol{\Xi}_m : \mathbb{R}^m \rightarrow \mathbb{R}^n$ em que $\boldsymbol{\Xi}_m \mathbf{e}_i = \boldsymbol{\xi}_i$	65
Figura 18 Representação do parâmetro $\hat{\boldsymbol{\beta}}_{\text{PCR}}$ como a projeção ortogonal do vetor $\hat{\boldsymbol{\beta}}_{\text{OLS}}$ no subespaço gerado por $\{\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2, \dots, \boldsymbol{\gamma}_m\}$	68
Figura 19 Projeção do vetor de dados \mathbf{Y} no subespaço gerado pelos m primeiros componentes principais	69
Figura 20 Representação geométrica do algoritmo PLS Populacional com duas iterações	81
Figura 21 Representação de $\mathbf{X}'\mathbf{Y}$ como uma transformação linear de \mathbb{R}^k em \mathbb{R}^p	89

Figura 22	Representação geométrica das seqüências de vetores $\mathbf{u}_n = \mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{u}_{n-1}$, $\mathbf{c}_n = \mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{c}_{n-1}$, $\mathbf{t}_n = \mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{t}_{n-1}$ e $\mathbf{w}_n = \mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{w}_{n-1}$	95
Figura 23	Representação da matriz de covariáveis $\mathbf{X}_{n \times p}$ como uma transformação do espaço de parâmetros \mathbb{R}^p no espaço dos dados \mathbb{R}^n .	105
Figura 24	Representação das colunas de \mathbf{X} como imagem dos vetores canônicos \mathbf{e}_i	106
Figura 25	Representação das variáveis centradas através da projeção ortogonal no subespaço gerado pelo vetor $\mathbf{1}$	106
Figura 26	Projeção ortogonal do vetor \mathbf{U}_1 em cada um dos vetores \mathbf{V}_{1j} . .	107
Figura 27	Projeção ortogonal do vetor \mathbf{U}_1 em cada um dos vetores \mathbf{V}_{1j} . .	107
Figura 28	Construção do vetor \mathbf{T}_1 como média ponderada dos vetores \mathbf{U}_{1j}	108
Figura 29	Construção dos vetores \mathbf{U}_2 e \mathbf{V}_{2j} através da projeção dos vetores \mathbf{U}_1 e \mathbf{V}_{1j} , respectivamente, em \mathbf{T}_1	109
Figura 30	Construção dos vetores \mathbf{U}_{2j} como projeção do vetor \mathbf{U}_2 nos vetores \mathbf{V}_{2j}	110
Figura 31	Construção do vetor \mathbf{T}_2 como média ponderada dos vetores \mathbf{U}_{2j}	110
Figura 32	Componentes ortogonais \mathbf{T}_1 e \mathbf{T}_2	111
Figura 33	Representação geométrica dos componentes como imagem da transformação linear \mathbf{T} do \mathbb{R}^m para o \mathbb{R}^n	113
Figura 34	Projeção do vetor de dados \mathbf{Y} no subespaço $\text{Im}\mathbf{T}$ gerado pelos componentes	113
Figura 35	Representação geométrica das transformações lineares \mathbf{T} , \mathbf{R} e \mathbf{X} .	114
Figura 36	Vetor $\hat{\beta}_{\text{PLS}}$ como projeção oblíqua do vetor $\hat{\beta}_{\text{OLS}}$	116
Figura 37	Representação de uma curva gerada por valores dos graus liberdade, DoF, em função do número de componentes m	123
Figura 38	Representação de uma curva gerada por valores de erro quadrático médio em função do número de componentes m	124
Figura 39	Como a transformação \mathbf{L} deforma uma esfera em um elipsóide .	143
Figura 40	Representação geométrica do método das potências	145
Capítulo 2		152
Figura 1	Determinação do número ótimo de componentes na análise PLS utilizando a teoria de graus de liberdade (a) e validação cruzada (b) para a variável pH da carne suína aos 45 min após o abate . .	165
Figura 2	Determinação do número ótimo de componentes na análise PLS utilizando a teoria de graus de liberdade (a) e validação cruzada (b) para a variável pH da carne suína 24 horas após o abate . . .	165

Figura 3	Capacidades preditivas obtidas pelos métodos PLS, PCR e OLS (regressão múltipla tradicional) para as características de pH ₄₅ (a) e pH _u (b)	166
Figura 4	<i>Manhattan plot</i> (efeitos estimados dos SNPs ao longo das posições genômicas) para as características pH da carne suína aos 45 min (a) e 24 horas (b) após o abate	168
Figura 5	Densidade de QTLs reportados na região intermediária do cromossomo SSC4 proveniente da base de dados PigQTLdb.	170

LISTA DE TABELAS

Capítulo 1	19
Capítulo 2	152
Tabela 1 Identificação dos SNPs reportados como sendo os mais relevantes para os fenótipos pH_{45} (45 min após o abate) e pH_u (24 horas após o abate)	169

SUMÁRIO

	INTRODUÇÃO GERAL	15
	CAPÍTULO 1 Abordagem Geométrica da Regressão por Quadrados Mínimos Parciais	19
1	INTRODUÇÃO	21
2	REFERENCIAL TEÓRICO	23
2.1	Vetores Aleatórios	23
2.2	Regressão Linear Múltipla	29
2.2.1	As equações normais e os estimadores de quadrados mínimos	32
2.2.2	Propriedades dos estimadores de quadrados mínimos	37
2.3	Regressão Estatística	39
2.4	Componentes Principais	49
2.5	Regressão por Componentes Principais	49
2.6	Regressão por Quadrados Mínimos Parciais	50
3	METODOLOGIA	51
3.1	Componentes Principais	51
3.2	Regressão por Componentes Principais	51
3.3	Regressão por Quadrados Mínimos Parciais	52
3.4	Determinação do número ótimo de componentes	53
3.5	Exemplo Didático	54
4	RESULTADOS	55
4.1	Componentes Principais	55
4.2	Regressão por Componentes Principais: uma abordagem geométrica	59
4.3	Regressão por Quadrados Mínimos Parciais: uma abordagem geométrica	74
4.3.1	PLS Populacional	74
4.3.2	PLS Amostral	87
4.3.2.1	Componentes com covariância máxima	88
4.3.2.2	O algoritmo PLS Multivariado	92
4.3.2.3	Regressão em PLS	102
4.3.3	Uma alternativa ao algoritmo PLS	105
4.4	Determinação do número ótimo de componentes (variáveis latentes)	116
4.4.1	Graus de Liberdade e seleção de modelos	117
4.4.2	Validação Cruzada	124
4.5	Exemplo	125
5	CONCLUSÕES	136

	REFERÊNCIAS	137
	APÊNDICES	140
	CAPÍTULO 2 Quadrados Mínimos Parciais aplicado à seleção genômica para qualidade de carne em suínos	152
1	INTRODUÇÃO	154
2	MATERIAL E MÉTODOS	157
2.1	Descrição dos dados utilizados	157
2.2	Quadrados Mínimos Parciais - PLS	158
2.3	Número ótimo de componentes (variáveis latentes) no PLS	161
2.4	Capacidade Preditiva	162
3	RESULTADOS E DISCUSSÃO	164
4	CONCLUSÕES	171
	REFERÊNCIAS	172
	APÊNDICE	175

INTRODUÇÃO GERAL

O método mais utilizado para obter equações de predição é o método dos quadrados mínimos. Contudo, quando existem fortes relações lineares entre as covariáveis, ou quando o número de covariáveis é maior que o número de observações em uma regressão múltipla, tem-se um problema denominado multicolinearidade. Nesta situação, as estimativas dos coeficientes de regressão por quadrados mínimos, tendem a ser instáveis, geralmente com grandes erros-padrão, podendo resultar em inferências errôneas (GARTHWAITE, 1994).

Uma forma de contornar o problema da multicolinearidade é a aplicação de métodos de redução de dimensão nas covariáveis. Dentre estes métodos destaca-se o método dos Quadrados Mínimos Parciais (*Partial Least Squares* - PLS). Este método foi introduzido por Wold em 1975, sendo considerado útil para a construção de equações de predições em situações nas quais se tem um grande número de variáveis explicativas e um número relativamente pequeno de dados amostrais (HOSKULDSSON, 1988). Resumidamente, a ideia geral do PLS é formar componentes, isto é, combinações lineares, que capturem a maior quantidade de informação possível disposta nas variáveis explicativas X_1, \dots, X_p para prever as variáveis respostas Y_1, \dots, Y_k . O método PLS apresenta similaridades com o tradicional método de Regressão via Componentes Principais (*Principal Components Regression* - PCR). A maior diferença é dada pelo fato do PCR levar em consideração, na construção dos componentes, apenas as variáveis explicativas, enquanto que o PLS também leva em consideração as variáveis respostas (GARTHWAITE, 1994).

Muitos estudos são encontrados na literatura com aplicações utilizando o método PLS, mas são raras as referências em que se aborda a teoria do PLS,

em especial a abordagem geométrica. Assim, trabalhos contemplando esta teoria geométrica são úteis para o entendimento e desenvolvimento do método.

Como em qualquer análise de redução de dimensionalidade, um passo importante na execução do PLS é a determinação do número ótimo de componentes. No entanto, isso tem sido uma lacuna na aplicação dessa teoria. Metodologias como a teoria de graus de liberdade (*Degrees of Freedom* - DoF) e validação cruzada (*Cross Validation* - CV) foram propostas (KRÄMER; SUGIYAMA, 2011), porém, até o momento, não há relatos da comparação de tais metodologias na aplicação de dados.

A título de ilustração, considere a grande contribuição da genética molecular no melhoramento animal. A utilização direta das informações de DNA no processo de identificação de animais geneticamente superiores, denominado Seleção Genômica Ampla (*Genome Wide Selection* - GWS) consiste da análise de um grande número de marcadores de Polimorfismo de Nucleotídeo Único (*Single Nucleotide Polymorphisms* - SNP) amplamente distribuídos no genoma (MEUWISSEN; HAYES; GODDARD, 2001). Porém, a utilização dessas informações é um desafio, uma vez que, geralmente, o número de marcadores é muito maior que o número de animais genotipados (alta dimensionalidade) e tais marcadores são altamente correlacionados (multicolinearidade). Assim, o sucesso da seleção genômica ampla deve-se à escolha de metodologias que contornem essas adversidades. Dentre estas, o PLS apresenta-se como uma alternativa interessante que merece ser investigada.

Diante do exposto, objetivou-se, no capítulo 1, apresentar uma abordagem geométrica à teoria do PLS (populacional e amostral) e do PCR. E no capítulo 2, objetivou-se aplicar os métodos PCR e PLS, além da regressão tradicional sem redução de dimensionalidade, em uma análise de seleção genômica em suínos

considerando um painel de marcadores SNPs de baixa densidade e dois fenótipos relacionados com a qualidade da carne (pH medido aos 45 min e às 24 horas após o abate). Além disso, objetivou-se, ainda, testar dois diferentes métodos de determinação do número ótimo de componentes na análise PLS, a teoria de graus de liberdade e a validação cruzada, bem como sua influência na performance preditiva do método.

CAPÍTULO 1

Abordagem Geométrica da Regressão por Quadrados Mínimos Parciais

RESUMO

A abordagem geométrica do método dos Quadrados Mínimos Parciais (PLS) é natural e intuitiva. Esta abordagem explicita as relações existentes entre os métodos de regressão PLS, Quadrados Mínimos Ordinários (OLS) e Componentes Principais (PCR). Em termos de projeções ortogonais, são explicadas, neste capítulo, as etapas da construção dos algoritmos PLS, populacional e amostral, bem como as construções relativas à regressão em componentes principais. Uma rotina foi desenvolvida no software R e aplicada a um exemplo didático, explicitando o passo a passo do algoritmo PLS. A teoria dos Graus de Liberdade é apresentada como uma alternativa ao método da Validação Cruzada, uma vez que esta vem sendo proposta para determinar o número de componentes a ser utilizado na regressão.

Palavras-chave: Componentes Principais. Projeções. Redução de dimensionalidade. Vetores Aleatórios.

ABSTRACT

The geometric approach to the Partial Least Square (PLS) Method is natural and intuitive. It makes clear the similarities among the PLS, the Ordinary Least Square (OLS) and the Principal Components (PCR) regression methods. In this chapter we use the orthogonal projections to explain the step by step construction of the PLS algorithm, for sample and for population, as well as the PCR. A routine was developed in the software R and applied in a didactic example. The Degrees of Freedom theory is presented as an alternative to the Cross Validation method, since it has been proposed to choose the number of components to be used in the regression.

Keywords: Principal Components. Projections. Dimensionality Reduction. Random Vectors.

1 INTRODUÇÃO

A álgebra linear é talvez a área da matemática mais acessível. Uma razão que pode justificar tal fato é a sua natureza dupla, isto é, álgebra, com toda a sua abstração e elegância, e geometria, com todo o seu apelo intuitivo. Vetores, ângulos, subespaços vetoriais são conceitos geométricos e de fácil visualização. A álgebra é, também, a área da matemática com maior aplicabilidade em estatística, como modelos lineares e estatística multivariada. O uso de transformações lineares ortogonais, diagonalização de matrizes simétricas e diagonalização de formas quadráticas, são essenciais e de uso constante. Esses resultados de álgebra linear são encontrados em praticamente todos os livros utilizados nos cursos de graduação em matemática, engenharias e estatística.

A teoria da regressão é uma área ampla e muito importante da estatística, e, apesar de existirem muitas referências sobre o assunto, pouco se encontra na literatura sobre sua abordagem geométrica. O fato interessante a ser observado é que a abordagem feita pelos autores estatísticos, em geral, é totalmente algébrica, extremamente analítica e abstrata. Desse modo, a utilização de uma abordagem geométrica poderá dar uma contribuição ao entendimento de vários conceitos na teoria da regressão.

Dentre os métodos de regressão, a regressão via componentes principais e a regressão via método dos quadrados mínimos parciais, têm destaque pela sua aplicabilidade, em especial esta última, que vem sendo muito utilizada em áreas como quimiometria e genética. No entanto, apesar de ser recorrente o uso do método PLS em trabalhos encontrados na literatura, muito se discute sobre sua aplicação, mas muito pouco se encontra sobre sua teoria.

O objetivo deste capítulo é apresentar a teoria da regressão em compo-

mentes principais e o método dos quadrados mínimos parciais, sob o ponto de vista geométrico, e a teoria dos Graus de Liberdade e Validação Cruzada, que vêm sendo propostas para determinar o número de componentes, tomando como referências os artigos: Helland (1990), Hoskuldsson (1988), Garthwaite (1994), Phatak e Jong (1997) e Krämer e Sugiyama (2011).

2 REFERENCIAL TEÓRICO

2.1 Vetores Aleatórios

Quando se quer medir vários aspectos aleatórios de um fenômeno, ou mesmo fazer medidas repetidas de algum deles, tem-se a teoria das variáveis aleatórias multidimensionais, $\mathbf{Y} = (Y_1, \dots, Y_n)$, em que cada coordenada é uma variável aleatória unidimensional. O exemplo mais usual é quando se tem uma amostra aleatória simples $\mathbf{Y} = (Y_1, \dots, Y_n)$. Neste caso \mathbf{Y} representa medidas repetidas independentes de uma mesma grandeza populacional. Outra situação é quando se tem uma variável multidimensional propriamente dita, isto é, a medida de várias características de um fenômeno aleatório. Pode-se também ter uma amostra de tais variáveis. Portanto, as coordenadas podem estar relacionadas a valores de característica ou valores obtidos por repetição.

Os valores que $\mathbf{Y} = (Y_1, \dots, Y_n)$ assume podem ser vistos como vetores no espaço \mathbb{R}^n , e, portanto, podem também ser apropriadamente denominados de vetores aleatórios. Neste trabalho, as variáveis aleatórias multidimensionais serão vistas como vetores aleatórios e estudam-se suas projeções em subespaços vetoriais, utilizando-se uma abordagem inteiramente geométrica.

Quando uma base do espaço vetorial é escolhida será usada a notação \mathbb{R}^n . Caso o espaço vetorial seja abstrato, isto é, sem uma escolha particular de base, será utilizada a notação V . Uma vez fixadas bases no domínio e no contradomínio, uma transformação linear pode ser representada por uma matriz. Dessa forma, não será feita distinção entre os conceitos de transformações lineares e de matrizes.

Uma situação usual em estatística é a transformação linear de dados, como mudanças de escalas, por exemplo, ou o procedimento de se trabalhar com a mé-

dia amostral. Tais procedimentos, em termos de álgebra linear consistem em aplicar transformações lineares, fazer projeções ortogonais, calcular ângulos, etc. Da mesma forma, todos estes procedimentos podem ser aplicados em vetores aleatórios.

Considere transformações lineares: $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$, conforme na Figura 1.

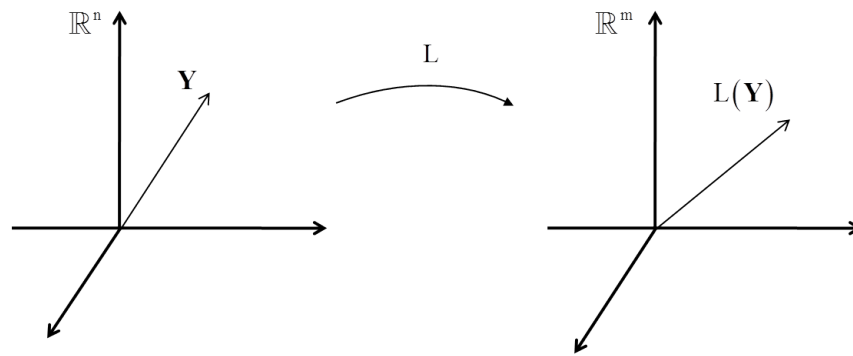


Figura 1 Transformação linear de \mathbb{R}^n em \mathbb{R}^m

Considere L como uma matriz $L_{m \times n}$ e L' sua transposta. Em termos de produto de matrizes, tem-se que a transformação linear pode ser dada por:

$$L_{m \times n} \mathbf{Y}_{n \times 1} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}.$$

O produto interno do vetor aleatório \mathbf{Y} com um vetor $\mathbf{x}' = (x_1, x_2, \dots, x_n)$ é dado por $\mathbf{x} \cdot \mathbf{Y} = x_1 Y_1 + \dots + x_n Y_n$. Em termos de matrizes, se \mathbf{x} é considerado

uma matriz $n \times 1$, ou seja, $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, então $\mathbf{x} \cdot \mathbf{Y}$ pode ser expresso como o produto de matrizes

$$\mathbf{x}'\mathbf{Y} = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}.$$

As seguintes notações serão usadas indistintamente: $\mathbf{x} \cdot \mathbf{Y} = \mathbf{x}'\mathbf{Y} = \langle \mathbf{x}, \mathbf{Y} \rangle$.

O produto interno é dado por $\mathbf{x} \cdot \mathbf{Y} = \|\mathbf{x}\| \|\mathbf{Y}\| \cos \theta$. Geometricamente, $\mathbf{x} \cdot \mathbf{Y}$ significa o tamanho do vetor \mathbf{x} vezes o tamanho do vetor aleatório obtido pela projeção do vetor \mathbf{Y} no subespaço unidimensional gerado pelo vetor \mathbf{x} (Figura 2).

$$\cos \theta = \frac{\|P_{\mathbf{x}}\mathbf{Y}\|}{\|\mathbf{Y}\|} \Rightarrow \|P_{\mathbf{x}}\mathbf{Y}\| = \|\mathbf{Y}\| \cos \theta,$$

$$\mathbf{x} \cdot \mathbf{Y} = \|\mathbf{x}\| \|\mathbf{Y}\| \cos \theta = \|\mathbf{x}\| \|P_{\mathbf{x}}\mathbf{Y}\|.$$

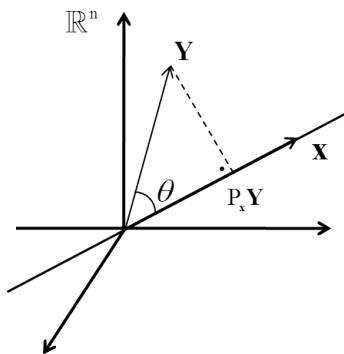


Figura 2 Interpretação geométrica do produto interno

A esperança de um vetor aleatório é o vetor dado pela esperança de cada componente, denominado $E(\mathbf{Y}) = \boldsymbol{\beta}$, em que $\boldsymbol{\beta}$ é um vetor $p \times 1$. Outra característica fundamental de um vetor aleatório é sua matriz de variâncias e covariâncias, $\text{cov}(\mathbf{Y}) = \boldsymbol{\Sigma} = E((\mathbf{Y} - E(\mathbf{Y}))(\mathbf{Y} - E(\mathbf{Y}))')$, em que $\boldsymbol{\Sigma}$ é uma matriz positiva definida $p \times p$.

Alguns teoremas importantes envolvendo projetores ortogonais de um vetor aleatório \mathbf{Y} são apresentados a seguir. Deve ser observado que, sendo \mathbf{Y} um vetor aleatório em um espaço vetorial V , a sua projeção ortogonal em um subespaço W , d -dimensional, de V , denotada por $P_W(\mathbf{Y})$, também é um vetor aleatório que depende, obviamente, de \mathbf{Y} (Figura 3). O vetor $P_W(\mathbf{Y})$ pode ser expresso, em relação a uma base ortogonal $\{\mathbf{u}_1, \dots, \mathbf{u}_d\}$ de W , por

$$P_W \mathbf{Y} = \left(\frac{\mathbf{Y} \cdot \mathbf{u}_1}{\mathbf{u}_1 \cdot \mathbf{u}_1} \right) \mathbf{u}_1 + \left(\frac{\mathbf{Y} \cdot \mathbf{u}_2}{\mathbf{u}_2 \cdot \mathbf{u}_2} \right) \mathbf{u}_2 + \dots + \left(\frac{\mathbf{Y} \cdot \mathbf{u}_d}{\mathbf{u}_d \cdot \mathbf{u}_d} \right) \mathbf{u}_d.$$

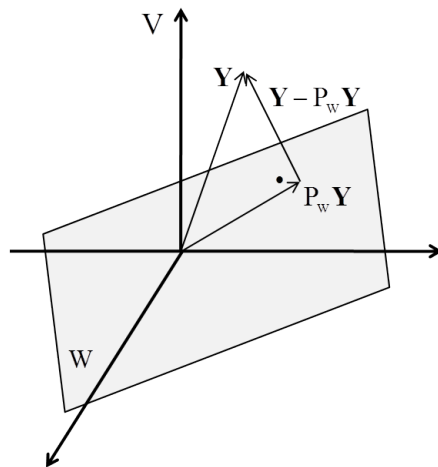


Figura 3 Projetor ortogonal de um vetor aleatório \mathbf{Y} sobre um subespaço W

Um fato importante sobre projetores é sua caracterização como matrizes. Uma matriz quadrada tal que $\mathbf{A}^2 = \mathbf{A}$ é chamada matriz de projeção ou proje-

tor. Tem-se que um projetor \mathbf{A} restrito à $\text{Im}\mathbf{A}$ é a identidade, ou seja, $\mathbf{A}(\mathbf{Ax}) = \mathbf{A}^2\mathbf{x} = \mathbf{Ax}$, $\forall \mathbf{x} \in \mathbb{R}^n$. $\mathbf{I} - \mathbf{A}$ também é um projetor, pois $(\mathbf{I} - \mathbf{A})^2 = \mathbf{I} - 2\mathbf{A} + \mathbf{A}^2 = \mathbf{I} - \mathbf{A}$, em que \mathbf{I} é a matriz identidade.

Tem-se que: $\text{Ker}\mathbf{A} = \text{Im}(\mathbf{I} - \mathbf{A})$ e $\text{Im}\mathbf{A} = \text{Ker}(\mathbf{I} - \mathbf{A})$.

Demonstração.

$$\mathbf{A}((\mathbf{I} - \mathbf{A})\mathbf{x}) = \mathbf{A}(\mathbf{x} - \mathbf{Ax}) = \mathbf{Ax} - \mathbf{A}^2\mathbf{x} = \mathbf{Ax} - \mathbf{Ax} = \mathbf{0}$$

$$\Rightarrow \text{Im}(\mathbf{I} - \mathbf{A}) \subset \text{Ker}\mathbf{A} \quad \text{e}$$

$$\mathbf{x} \in \text{Ker}\mathbf{A} \Rightarrow (\mathbf{I} - \mathbf{A})\mathbf{x} = \mathbf{x} - \mathbf{Ax} = \mathbf{x} - \mathbf{0} \Rightarrow \mathbf{x} \in \text{Im}(\mathbf{I} - \mathbf{A})$$

$$\Rightarrow \text{Ker}\mathbf{A} \subset \text{Im}(\mathbf{I} - \mathbf{A}),$$

logo, $\text{Ker}\mathbf{A} = \text{Im}(\mathbf{I} - \mathbf{A})$.

$$\mathbf{y} \in \text{Ker}(\mathbf{I} - \mathbf{A}) \Rightarrow \mathbf{0} = (\mathbf{I} - \mathbf{A})\mathbf{y} = \mathbf{y} - \mathbf{Ay} \Rightarrow \mathbf{Ay} = \mathbf{y} \Rightarrow \mathbf{y} \in \text{Im}\mathbf{A}$$

$$\Rightarrow \text{Ker}(\mathbf{I} - \mathbf{A}) \subset \text{Im}\mathbf{A} \quad \text{e}$$

$$\mathbf{y} \in \text{Im}\mathbf{A} \Rightarrow \mathbf{y} = \mathbf{Ax} \Rightarrow (\mathbf{I} - \mathbf{A})\mathbf{y} = (\mathbf{I} - \mathbf{A})\mathbf{Ax} = \mathbf{Ax} - \mathbf{A}^2\mathbf{x}$$

$$= \mathbf{Ax} - \mathbf{Ax} = \mathbf{0} \Rightarrow \mathbf{y} \in \text{Ker}(\mathbf{I} - \mathbf{A})$$

$$\Rightarrow \text{Im}\mathbf{A} \subset \text{Ker}(\mathbf{I} - \mathbf{A}),$$

logo, $\text{Im}\mathbf{A} = \text{Ker}(\mathbf{I} - \mathbf{A})$. □

Uma matriz de projeção \mathbf{A} é dita um projetor ortogonal se $\mathbf{Av} - \mathbf{v}$ é perpendicular ao subespaço $\text{Im}\mathbf{A}$.

Proposição 1. *Uma matriz de projeção $A_{n \times n}$ é simétrica se, e somente se, é um projetor ortogonal, isto é,*

$$\langle \mathbf{A}\mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{v}, \mathbf{A}\mathbf{w} \rangle \Leftrightarrow \langle \mathbf{A}\mathbf{v} - \mathbf{v}, \mathbf{A}\mathbf{w} \rangle = 0, \quad \forall \mathbf{v}, \mathbf{w} \in \mathbb{R}^n.$$

Demonstração. (\Rightarrow) Se \mathbf{A} é simétrica ($\mathbf{A}' = \mathbf{A}$), então:

$$\begin{aligned} \langle \mathbf{v} - \mathbf{A}\mathbf{v}, \mathbf{A}\mathbf{w} \rangle &= \langle \mathbf{v}, \mathbf{A}\mathbf{w} \rangle - \langle \mathbf{A}\mathbf{v}, \mathbf{A}\mathbf{w} \rangle = \langle \mathbf{v}, \mathbf{A}\mathbf{w} \rangle - \langle \mathbf{v}, \mathbf{A}^2\mathbf{w} \rangle \\ &= \langle \mathbf{v}, \mathbf{A}\mathbf{w} \rangle - \langle \mathbf{v}, \mathbf{A}\mathbf{w} \rangle = 0, \quad \forall \mathbf{v}, \mathbf{w} \in \mathbb{R}^n. \end{aligned}$$

Logo, \mathbf{A} é um projetor ortogonal.

(\Leftarrow) Se \mathbf{A} é um projetor ortogonal, $\mathbf{A}\mathbf{v} - \mathbf{v}$ é perpendicular à imagem de \mathbf{A} , isto é, $\langle \mathbf{A}\mathbf{v} - \mathbf{v}, \mathbf{A}\mathbf{w} \rangle = 0 \Rightarrow \langle \mathbf{A}\mathbf{v}, \mathbf{A}\mathbf{w} \rangle = \langle \mathbf{v}, \mathbf{A}\mathbf{w} \rangle$, da mesma forma $\mathbf{A}\mathbf{w} - \mathbf{w}$ é perpendicular a $\mathbf{A}\mathbf{v}$, isto é, $\langle \mathbf{A}\mathbf{w} - \mathbf{w}, \mathbf{A}\mathbf{v} \rangle = 0 \Rightarrow \langle \mathbf{A}\mathbf{w}, \mathbf{A}\mathbf{v} \rangle = \langle \mathbf{w}, \mathbf{A}\mathbf{v} \rangle$, portanto, $\langle \mathbf{A}\mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{v}, \mathbf{A}\mathbf{w} \rangle$. Logo, \mathbf{A} é simétrica.

□

Teorema 1. *Se $\mathbf{Y} \in V$ é vetor aleatório n -dimensional com $E(\mathbf{Y}) = \boldsymbol{\tau}$, $\text{cov}(\mathbf{Y}) = \boldsymbol{\Sigma}$ e W um subespaço d -dimensional de V . Então:*

- $E(\mathbf{x} \cdot \mathbf{Y}) = \mathbf{x} \cdot E(\mathbf{Y}) = \mathbf{x} \cdot \boldsymbol{\tau}$
- $\text{var}(\mathbf{x} \cdot \mathbf{Y}) = \mathbf{x}'\boldsymbol{\Sigma}\mathbf{x}$
- $\text{cov}(\mathbf{x} \cdot \mathbf{Y}, \mathbf{z} \cdot \mathbf{Y}) = \mathbf{x}'\boldsymbol{\Sigma}\mathbf{z}$
- Se \mathbf{x} é um autovetor de $\boldsymbol{\Sigma}$ relativo ao autovalor ξ então $\text{var}(\mathbf{x} \cdot \mathbf{Y}) = \|\mathbf{x}\|^2 \xi$.
- Sejam \mathbf{x}, \mathbf{z} autovetores relativos aos autovalores ξ e η .

Se $\xi = \eta$, então $\text{cov}(\mathbf{x} \cdot \mathbf{Y}, \mathbf{z} \cdot \mathbf{Y}) = (\mathbf{x} \cdot \mathbf{z})\xi$.

Se $\xi \neq \eta$, então $\text{cov}(\mathbf{x} \cdot \mathbf{Y}, \mathbf{z} \cdot \mathbf{Y}) = 0$.

- Se W é um subespaço d -dimensional do subespaço dos autovetores relativos ao autovalor ξ , tem-se:

$$E(P_W \mathbf{Y}) = P_W E(\mathbf{Y}) = P_W \boldsymbol{\tau}$$

$$E\left(\|P_W \mathbf{Y}\|^2\right) = \|P_W \boldsymbol{\tau}\|^2 + d\xi.$$

Para o caso em que $\text{cov}(\mathbf{Y}) = \sigma^2 \mathbf{I}$, tem-se:

- $\text{var}(\mathbf{x} \cdot \mathbf{Y}) = \sigma^2 \mathbf{x}' \mathbf{x}$
- $\text{cov}(\mathbf{x} \cdot \mathbf{Y}, \mathbf{z} \cdot \mathbf{Y}) = \sigma^2 \mathbf{x}' \mathbf{z}$
- $E\left(\|P_W \mathbf{Y}\|^2\right) = \|P_W \boldsymbol{\tau}\|^2 + d\sigma^2$

Corolário 1. Se \mathbf{x} e \mathbf{z} são vetores ortogonais então $\mathbf{x} \cdot \mathbf{Y}$ e $\mathbf{z} \cdot \mathbf{Y}$ são variáveis aleatórias não correlacionadas.

As demonstrações destes resultados encontram-se em Adão (2011).

2.2 Regressão Linear Múltipla

Considere observações (Y_1, Y_2, \dots, Y_n) não correlacionadas, tais que:

$$E(Y_i) = x_{i1}\beta_1 + x_{i2}\beta_2 + \dots + x_{ip}\beta_p \text{ e } \text{var}(Y_i) = \sigma^2,$$

para $i = 1, 2, \dots, n$, em que $\beta_1, \beta_2, \dots, \beta_p$ e σ^2 são parâmetros desconhecidos.

Tem-se, portanto, o sistema:

$$\begin{aligned} E(Y_1) &= x_{11}\beta_1 + x_{12}\beta_2 + \cdots + x_{1p}\beta_p \\ E(Y_2) &= x_{21}\beta_1 + x_{22}\beta_2 + \cdots + x_{2p}\beta_p \\ &\vdots \\ E(Y_n) &= x_{n1}\beta_1 + x_{n2}\beta_2 + \cdots + x_{np}\beta_p. \end{aligned}$$

Representando matricialmente este sistema de equações, tem-se:

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \quad \text{e} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}$$

$$\begin{bmatrix} E(Y_1) \\ E(Y_2) \\ \vdots \\ E(Y_n) \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix},$$

isto é, denominando o vetor de médias por $E(\mathbf{Y}) = \boldsymbol{\tau}$ o problema de se estimar o vetor $\boldsymbol{\tau}$ é equivalente a se obter o vetor $\boldsymbol{\beta}$ tal que $\boldsymbol{\tau} = \mathbf{X}\boldsymbol{\beta}$, em que $\mathbf{X} = (x_{ij})$ é uma matriz de coeficientes conhecidos que depende do delineamento subjacente à obtenção dos dados. Como casos extremos, tem-se o caso em que os y_i são observações independentes de uma mesma população, ou, utilizando a linguagem da estatística experimental, todos os y_i são respostas de um mesmo tratamento, neste caso, $E(\mathbf{Y}) = \boldsymbol{\tau}$ é um vetor com coordenadas iguais e a matriz \mathbf{X} é uma matriz $n \times 1$ com entradas iguais a 1. A outra situação ocorre quando os y_i são observações independentes de n populações diferentes ($p = n$), ou, utilizando a linguagem

da estatística experimental, os y_i são respostas relativas a n tratamentos diferentes. Neste caso, $E(\mathbf{Y}) = \boldsymbol{\tau}$ é um vetor com todas as coordenadas possivelmente diferentes e a matriz \mathbf{X} é uma matriz $n \times n$.

É claro que os valores observados no vetor de respostas \mathbf{Y} são, em razão da aleatoriedade que envolve todo experimento, diferentes dos valores esperados $E(\mathbf{Y})$. Note que os valores esperados não são conhecidos, mas sabe-se, por exemplo, que a repetição de um tratamento deveria ter a mesma resposta e tal fato não acontece em razão dos fatores aleatórios que ocorrem em todo experimento.

Desse modo, pode-se considerar uma variável aleatória ε , tal que

$$\varepsilon = \mathbf{Y} - E(\mathbf{Y}) \Rightarrow \mathbf{Y} = E(\mathbf{Y}) + \varepsilon.$$

O vetor ε é denominado vetor de erros ou vetor de resíduos, uma vez que mede a diferença entre o valor efetivamente observado e o valor desconhecido esperado. Considerando que os componentes do vetor de erros ε sejam independentes e igualmente distribuídos, tem-se que:

$$\begin{aligned} E(\varepsilon) &= E(\mathbf{Y} - E(\mathbf{Y})) = E(\mathbf{Y}) - E(E(\mathbf{Y})) = E(\mathbf{Y}) - E(\mathbf{Y}) = \mathbf{0}, \\ \text{cov}(\varepsilon) &= \text{cov}(\mathbf{Y} - E(\mathbf{Y})) = \text{cov}(\mathbf{Y}) = \sigma^2 \mathbf{I}_n, \end{aligned}$$

em que $\text{cov}(\cdot)$ é a matriz de variâncias e covariâncias e \mathbf{I}_n é a matriz identidade de ordem n .

O problema que se quer resolver é: como estimar os parâmetros desconhecidos β_j com base nas observações y_i ?

A equação $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \varepsilon$ é a chamada equação do modelo linear geral. A razão do nome é que o modelo é linear em seus parâmetros β_j . O modelo linear geral é a equação do modelo, sob as condições $E(\varepsilon) = \mathbf{0}$ e $\text{cov}(\varepsilon) = \sigma^2 \mathbf{I}_n$.

Geometricamente, o modelo linear de Gauss-Markov pode ser descrito como na Figura 4.

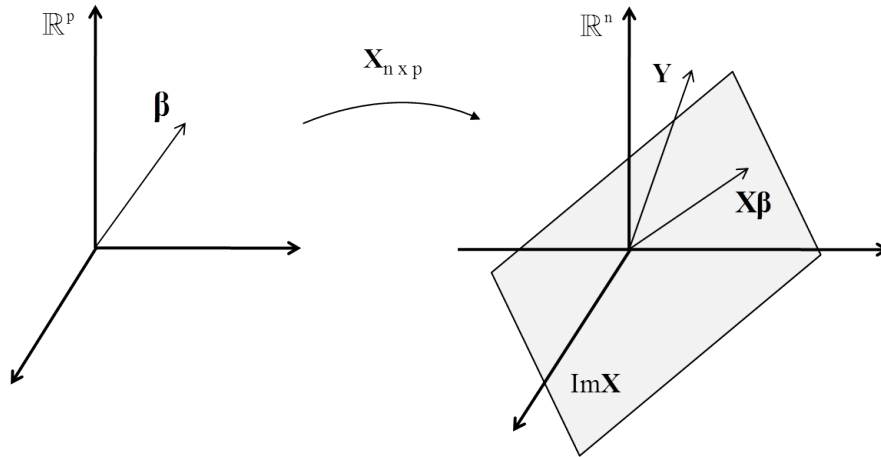


Figura 4 Representação geométrica para o modelo linear de Gauss-Markov

Seja \mathbb{R}^n o espaço dos dados, \mathbb{R}^p o espaço dos parâmetros, e $\text{Im}X = \{v = X\beta, \beta \in \mathbb{R}^p\}$ o subespaço formado pelos vetores que são imagem da transformação linear X , isto é, o subespaço vetorial gerado pelas colunas da matriz X . O subespaço dos vetores que são levados no vetor nulo por X , é denominado kernel da transformação linear X , $\text{Ker}X$. Em geral, n é maior que p , pois o número de dados geralmente é maior que o número de parâmetros.

2.2.1 As equações normais e os estimadores de quadrados mínimos

No modelo linear geral: $Y = X\beta + \varepsilon$ ou $\varepsilon = Y - X\beta$, o estimador de quadrados mínimos de β é o vetor $\hat{\beta}$ que minimiza a soma de quadrados dos erros, isto é, o valor do vetor de parâmetros que faz a norma do vetor ε , $\|\varepsilon\|$, ser a menor possível, ou de forma equivalente, que faz a norma ao quadrado, $\|\varepsilon\|^2$, a menor possível.

Tal problema pode ser resolvido por considerações geométricas de forma bastante simples, sem o uso de cálculo diferencial. A solução é o vetor obtido pela diferença entre o vetor observado \mathbf{Y} e a projeção deste vetor no subespaço $\text{Im}\mathbf{X}$ uma vez que a projeção ortogonal minimiza a distância de um ponto a um subespaço vetorial (Figura 5).

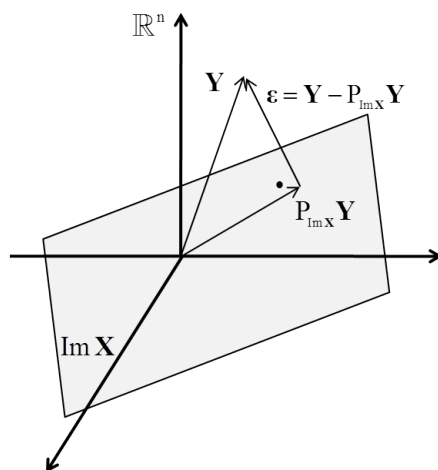


Figura 5 Representação geométrica do vetor residual

O problema, então, torna-se: como obter o projetor ortogonal no subespaço $\text{Im}\mathbf{X}$. Projetar ortogonalmente em um subespaço é equivalente a obter o subespaço complementar ortogonal deste subespaço. Tem-se, portanto, que se obter o subespaço $\text{Im}\mathbf{X}^\perp$. Uma maneira de se obter tal subespaço é trabalhar com a matriz transposta \mathbf{X}' .

A transposta de \mathbf{X} é uma transformação linear $\mathbf{X}' : \mathbb{R}^n \rightarrow \mathbb{R}^p$, tal que, para $\mathbf{x} \in \mathbb{R}^p$ e $\mathbf{y} \in \mathbb{R}^n$ quaisquer, tem-se $\langle \mathbf{X}\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{X}'\mathbf{y} \rangle$.

$\langle \mathbf{X}\mathbf{x}, \mathbf{y} \rangle = (\mathbf{X}\mathbf{x})'\mathbf{y} = \mathbf{x}'\mathbf{X}'\mathbf{y} = \langle \mathbf{x}, \mathbf{X}'\mathbf{y} \rangle$. Se $\mathbf{x} \in \text{Ker}\mathbf{X}$ então, para todo $\mathbf{y} \in \mathbb{R}^n$, $\langle \mathbf{X}\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{0}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{X}'\mathbf{y} \rangle = 0$ e, conseqüentemente, $\text{Ker}\mathbf{X}$ é perpendicular à $\text{Im}\mathbf{X}'$. Da mesma forma segue que $\text{Ker}\mathbf{X}'$ é perpendicular à $\text{Im}\mathbf{X}$.

Temos então a configuração geométrica essencial, mostrada na Figura 6.

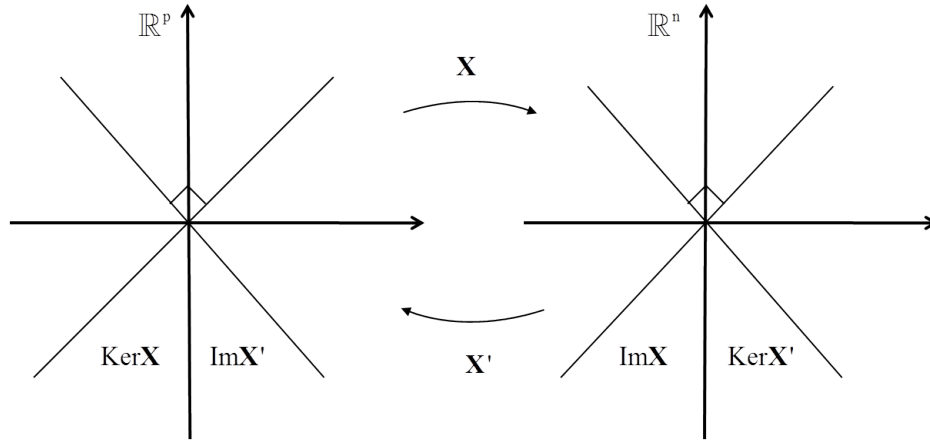


Figura 6 Representação gráfica do núcleo e da imagem de X e X'

Primeiramente vamos considerar o caso mais simples em que $\text{Ker}X = \{0\}$. Neste caso, a matriz $X'X$ é inversível e define um operador $X'X : \mathbb{R}^p \rightarrow \mathbb{R}^p$. A ideia é definir a matriz $P = X(X'X)^{-1}X'$. Como

$$P^2 = X(X'X)^{-1}X'X(X'X)^{-1}X' = X(X'X)^{-1}X' = P,$$

tem-se que P é um projetor, e como

$$P' = \left(X(X'X)^{-1}X' \right)' = X(X'X)^{-1}X' = P,$$

segue que P é um projetor ortogonal, e a projeção é no subespaço $\text{Im}X$.

Segue deste fato que o estimador de quadrados mínimos é dado pelo único vetor β que é levado por X no vetor $P_{\text{Im}X}Y = \left(X(X'X)^{-1}X' \right) Y$, isto é, o estimador é obtido resolvendo-se a equação $X\beta = X(X'X)^{-1}X'Y$. Novamente, utilizando que a transformação X é injetiva, a equação pode ser simplificada aplicando-se o operador X' em ambos os lados da igualdade obtendo-se

$\mathbf{X}'\mathbf{X}\beta = \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \Rightarrow \mathbf{X}'\mathbf{X}\beta = \mathbf{X}'\mathbf{Y}$, que são as chamadas equações normais do modelo linear, com solução da forma $\hat{\beta}_{\text{OLS}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$, em que OLS significa Quadrados Mínimos Ordinários (*Ordinary Least Square*), ou, simplesmente, $\hat{\beta}$.

Quando \mathbf{X} não é injetiva, o projetor ortogonal é dado por $\mathbf{P}_{\text{Im}\mathbf{X}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$ em que $(\mathbf{X}'\mathbf{X})^{-}$ é uma inversa generalizada de $\mathbf{X}'\mathbf{X}$. Geralmente se utiliza da inversa generalizada de Moore-Penrose. Note que fica então provado que o sistema de equações normais sempre admite solução. O problema agora é que quando \mathbf{X} não é injetiva existem infinitos vetores β em \mathbb{R}^p que são levados por \mathbf{X} em $\mathbf{P}_{\text{Im}\mathbf{X}}\mathbf{Y} = \left(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\right)\mathbf{Y}$, pois se $\mathbf{w} \in \text{Ker}\mathbf{X}$, $\mathbf{X}(\beta + \mathbf{w}) = \mathbf{X}\beta$, ou seja, como o núcleo de \mathbf{X} contém infinitos elementos, então existem infinitos vetores que são levados na mesma imagem e, portanto, têm-se infinitas estimativas de quadrados mínimos.

Mesmo neste caso de não injetividade de \mathbf{X} , isto é, mesmo quando \mathbf{X} não tem posto coluna completo, os vetores que são as soluções da equação $\mathbf{X}\beta = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}$ são os mesmos vetores que são soluções do sistema obtido aplicando-se \mathbf{X}' em ambos os lados da igualdade. Tal fato ocorre, pois como a $\text{Im}\mathbf{X}$ é perpendicular ao $\text{Ker}\mathbf{X}'$, \mathbf{X}' restrita à $\text{Im}\mathbf{X}$ é injetiva.

No caso de posto incompleto, escolhendo-se uma inversa generalizada de Quadrados Mínimos (RAO, 2002) fica definido uma pré-imagem, isto é, $\hat{\beta}$ também pode ser expresso explicitamente como função de \mathbf{Y} da forma (Figura 7):

$$\mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y} \Rightarrow \mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y} \Rightarrow \hat{\beta} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}.$$

Ocorre que, para o caso da matriz do delineamento \mathbf{X} não ser de posto completo, havendo infinitas estimativas (soluções) de quadrados mínimos para os parâmetros β , é um complicante na teoria. As propriedades e a construção geo-

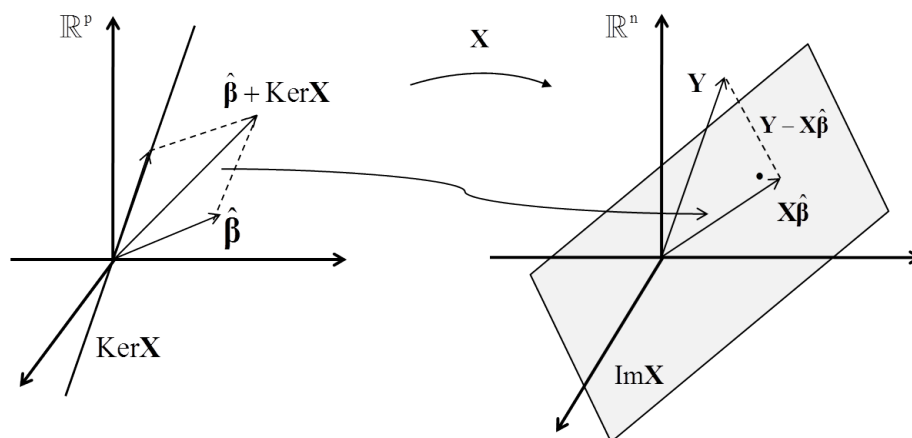


Figura 7 Modelo de Gauss-Markov quando \mathbf{X} não é injetiva

métrica das inversas generalizadas podem ser vistas em Guimarães (2010).

De uma maneira geral, se o interesse é estimar a combinação linear entre os efeitos dos tratamentos, tem-se $x_1\tau_1 + \dots + x_n\tau_n = \langle \mathbf{x}, \boldsymbol{\tau} \rangle$. Como $\langle \mathbf{x}, \boldsymbol{\tau} \rangle = \langle \mathbf{x}, \mathbf{X}\boldsymbol{\beta} \rangle = \langle \mathbf{X}'\mathbf{x}, \boldsymbol{\beta} \rangle$, uma combinação linear entre os efeitos dos tratamentos pode ser, de maneira equivalente, descrita como uma combinação linear dos parâmetros. Uma maneira de se estimar tais combinações lineares é tomar um estimador de quadrados mínimos de $\boldsymbol{\beta}$ e fazer o produto interno $\langle \hat{\boldsymbol{\beta}}, \mathbf{X}'\mathbf{x} \rangle$. Tal procedimento, para estar bem definido, não pode depender de qual estimativa de quadrados mínimos, isto é, $\langle \hat{\boldsymbol{\beta}}, \mathbf{X}'\mathbf{x} \rangle$ tem que ser o mesmo valor, qualquer que seja $\hat{\boldsymbol{\beta}}$. Para que isto ocorra, o vetor $\mathbf{X}'\mathbf{x}$ tem que ser ortogonal ao $\text{Ker}\mathbf{X}$. Como tal fato é verdadeiro, pois $\text{Im}\mathbf{X}'$ é ortogonal ao $\text{Ker}\mathbf{X}$, o procedimento pode sempre ser feito. De uma maneira geral, uma combinação linear $\langle \boldsymbol{\beta}, \mathbf{v} \rangle$ é estimável no sentido dos quadrados mínimos se \mathbf{v} for perpendicular ao espaço $\text{Ker}\mathbf{X}$. Caso contrário, $\langle \boldsymbol{\beta}, \mathbf{v} \rangle$ é dito não estimável, nome não muito adequado, pois não é estimável via quadrados mínimos. Mas a expressão quadrados mínimos relativa à combinação linear $\langle \boldsymbol{\beta}, \mathbf{v} \rangle$ também não é muito adequada, pois não foi feito ne-

nhum procedimento de se minimizar quadrados em relação ao valor $\langle \beta, \mathbf{v} \rangle$.

2.2.2 Propriedades dos estimadores de quadrados mínimos

No caso de posto completo, $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$. Assim:

- *Não viesado:*

$$\begin{aligned} E(\hat{\beta}) &= E\left((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}\right) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{Y}) \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\tau} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta}, \end{aligned}$$

e, portanto, é um estimador não viesado.

- $\text{cov}(\hat{\beta}) = \text{cov}\left((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}\right) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\text{cov}(\mathbf{Y})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$.

Se $\text{cov}(\mathbf{Y}) = \sigma^2\mathbf{I}$, então,

$$\text{cov}(\hat{\beta}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}.$$

No caso de posto incompleto, $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}$. Assim:

- *Viesado:*

$$\begin{aligned} E(\hat{\beta}) &= E\left((\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}\right) = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'E(\mathbf{Y}) \\ &= (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\boldsymbol{\tau} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{X}\boldsymbol{\beta}, \end{aligned}$$

logo, é um estimador viesado.

- $\text{cov}(\hat{\beta}) = \text{cov}\left((\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}\right) = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\text{cov}(\mathbf{Y})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}$.

Se $\text{cov}(\mathbf{Y}) = \sigma^2\mathbf{I}$, então,

$$\text{cov}(\hat{\boldsymbol{\beta}}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}.$$

- Seja a combinação linear dos parâmetros, $\langle \mathbf{r}, \boldsymbol{\beta} \rangle$. Se $\langle \mathbf{r}, \hat{\boldsymbol{\beta}} \rangle$ é estimador de quadrados mínimos de $\langle \mathbf{r}, \boldsymbol{\beta} \rangle$, então:

– No caso de posto completo e $\text{cov}(\mathbf{Y}) = \sigma^2\mathbf{I}$:

$$E(\langle \mathbf{r}, \hat{\boldsymbol{\beta}} \rangle) = \langle \mathbf{r}, E(\hat{\boldsymbol{\beta}}) \rangle = \langle \mathbf{r}, \boldsymbol{\beta} \rangle, \text{ que é não viesado, e}$$

$$\text{cov}(\langle \mathbf{r}, \hat{\boldsymbol{\beta}} \rangle) = \mathbf{r}'\text{cov}(\hat{\boldsymbol{\beta}})\mathbf{r} = \sigma^2\mathbf{r}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{r}.$$

– No caso de posto incompleto e $\text{cov}(\mathbf{Y}) = \sigma^2\mathbf{I}$:

$$E(\langle \mathbf{r}, \hat{\boldsymbol{\beta}} \rangle) = \langle \mathbf{r}, E(\hat{\boldsymbol{\beta}}) \rangle = \langle \mathbf{r}, (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} \rangle, \text{ que é viesado, e}$$

$$\text{cov}(\langle \mathbf{r}, \hat{\boldsymbol{\beta}} \rangle) = \mathbf{r}'\text{cov}(\hat{\boldsymbol{\beta}})\mathbf{r} = \sigma^2\mathbf{r}'(\mathbf{X}'\mathbf{X})^{-}\mathbf{r}.$$

Se $\langle \mathbf{1}, \boldsymbol{\beta} \rangle$ é outra combinação linear dos parâmetros, a covariância entre os estimadores é dada por:

– No caso de posto completo e $\text{cov}(\mathbf{Y}) = \sigma^2\mathbf{I}$:

$$\text{cov}(\langle \mathbf{r}, \hat{\boldsymbol{\beta}} \rangle, \langle \mathbf{1}, \hat{\boldsymbol{\beta}} \rangle) = \sigma^2\mathbf{r}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{1}.$$

– No caso de posto incompleto e $\text{cov}(\mathbf{Y}) = \sigma^2\mathbf{I}$:

$$\text{cov}(\langle \mathbf{r}, \hat{\boldsymbol{\beta}} \rangle, \langle \mathbf{1}, \hat{\boldsymbol{\beta}} \rangle) = \sigma^2\mathbf{r}'(\mathbf{X}'\mathbf{X})^{-}\mathbf{1}.$$

A vantagem de se estimar combinações lineares dos parâmetros utilizando o estimador de quadrados mínimos segue da proposição:

Proposição 2. *Entre todos os estimadores não viesados de $\langle \mathbf{r}, \boldsymbol{\beta} \rangle$ dados por combinações lineares dos dados, o de menor variância é $\langle \mathbf{r}, \hat{\boldsymbol{\beta}} \rangle$.*

Demonstração. Seja $\langle \mathbf{l}, \mathbf{Y} \rangle$ um estimador não viesado de $\langle \mathbf{r}, \boldsymbol{\beta} \rangle$. Então $E(\langle \mathbf{l}, \mathbf{Y} \rangle) = \langle \mathbf{r}, \boldsymbol{\beta} \rangle$. Mas $E(\langle \mathbf{l}, \mathbf{Y} \rangle) = \langle \mathbf{l}, E(\mathbf{Y}) \rangle = \langle \mathbf{l}, \mathbf{X}\boldsymbol{\beta} \rangle = \langle \mathbf{X}'\mathbf{l}, \boldsymbol{\beta} \rangle$ e, portanto, $\langle \mathbf{r}, \boldsymbol{\beta} \rangle = \langle \mathbf{X}'\mathbf{l}, \boldsymbol{\beta} \rangle$. Como $\boldsymbol{\beta}$ é desconhecido, esta igualdade tem que ser válida para todo $\boldsymbol{\beta}$ e, portanto, $\mathbf{r} = \mathbf{X}'\mathbf{l}$. Como $\text{var}(\mathbf{l}'\mathbf{Y}) = \|\mathbf{l}\|^2 \sigma^2$, temos que obter \mathbf{l} de norma mínima na pré-imagem de \mathbf{r} por \mathbf{X}' . Vamos projetar ortogonalmente \mathbf{l} em $\text{Im } \mathbf{X}$, isto é, $\mathbf{P}\mathbf{l}$. Logo $\mathbf{l} = \mathbf{P}\mathbf{l} + (\mathbf{I} - \mathbf{P})\mathbf{l} \Rightarrow \|\mathbf{l}\|^2 \geq \|\mathbf{P}\mathbf{l}\|^2$ e $\mathbf{X}'\mathbf{l} = \mathbf{X}'\mathbf{P}\mathbf{l}$.

Assim, $E(\langle \mathbf{l} - \mathbf{P}\mathbf{l}, \mathbf{Y} \rangle) = \langle \mathbf{l} - \mathbf{P}\mathbf{l}, E(\mathbf{Y}) \rangle = \langle \mathbf{l} - \mathbf{P}\mathbf{l}, \mathbf{X}\boldsymbol{\beta} \rangle = 0 \Rightarrow \langle \mathbf{l}, \mathbf{X}\boldsymbol{\beta} \rangle = \langle \mathbf{P}\mathbf{l}, \mathbf{X}\boldsymbol{\beta} \rangle$.

Como $\langle \mathbf{P}\mathbf{l}, \mathbf{Y} \rangle = \langle \mathbf{l}, \mathbf{P}\mathbf{Y} \rangle = \langle \mathbf{l}, \mathbf{X}\hat{\boldsymbol{\beta}} \rangle = \langle \mathbf{X}'\mathbf{l}, \hat{\boldsymbol{\beta}} \rangle = \langle \mathbf{r}, \hat{\boldsymbol{\beta}} \rangle$, o estimador de menor variância para $\langle \mathbf{r}, \boldsymbol{\beta} \rangle$ é a combinação linear dos dados expressa por $\langle \mathbf{r}, \hat{\boldsymbol{\beta}} \rangle$.

□

2.3 Regressão Estatística

Uma das situações mais comuns em estatística é quando se tem variáveis aleatórias Y_1, \dots, Y_k e X_1, \dots, X_p com uma distribuição conjunta

$$f_{Y_1, \dots, Y_k, X_1, \dots, X_p}(y_1, \dots, y_k, x_1, \dots, x_p)$$

e, por alguma razão, as variáveis X_1, \dots, X_p podem ser amostradas, mas as variáveis Y_i não. Após a observação das variáveis X_i , o conhecimento sobre as variáveis Y_i aumenta pois podemos calcular a função densidade de probabilidade condicional

$$f_{Y_1, \dots, Y_k | X_1, \dots, X_p}(y_1, \dots, y_k; x_1, \dots, x_p).$$

Como, em geral, a densidade de probabilidade conjunta é desconhecida, a densidade de probabilidade condicional não pode ser calculada. As variáveis observáveis X_i são denominadas variáveis preditoras e as variáveis Y_i de variáveis de interesse ou respostas.

Considere a situação mais simples em que se tem apenas uma variável de interesse Y . Seja

$$\begin{aligned} E(Y|X_1=x_1, \dots, X_p=x_p) \\ = \int \cdots \int y f_{Y|X_1, \dots, X_p}(y; x_1, \dots, x_p) dy = M(x_1, \dots, x_p). \end{aligned}$$

Pode-se, agora, considerar a variável aleatória

$$M(X_1, \dots, X_p) = E(Y|X_1, \dots, X_p).$$

Esta variável é a melhor variável aleatória para prever Y , no sentido de apresentar a menor esperança do quadrado do erro, isto é, se $f(X_1, \dots, X_p)$ é uma função qualquer das variáveis X_1, \dots, X_p , então:

Proposição 3. $E[(Y - f(X_1, \dots, X_p))^2] \geq E[(Y - M(X_1, \dots, X_p))^2]$.

Demonstração.

$$\begin{aligned} & E\left((Y - f(X_1, \dots, X_p))^2\right) \\ &= E\left((Y - M(X_1, \dots, X_p) + M(X_1, \dots, X_p) - f(X_1, \dots, X_p))^2\right) \\ &= E\left((Y - M(X_1, \dots, X_p))^2\right) + E\left((M(X_1, \dots, X_p) - f(X_1, \dots, X_p))^2\right) \\ &\quad + 2E\left((Y - M(X_1, \dots, X_p))(M(X_1, \dots, X_p) - f(X_1, \dots, X_p))\right). \end{aligned}$$

Tem-se que

$$E[(Y - M(X_1, \dots, X_m))(M(X_1, \dots, X_m) - f(X_1, \dots, X_m))] = 0.$$

Para tal, considere a notação simplificada $E[(Y - M)(M - f)]$:

$$\begin{aligned} & E[(Y - M)(M - f)] \\ &= \int \int (y - m)(m - f) f_{Y, \mathbf{X}}(y, \mathbf{x}) dy dx \\ &= \int \int (y - m)(m - f) f_{Y|\mathbf{X}}(y|\mathbf{x}) f_{\mathbf{X}}(\mathbf{x}) dy dx \\ &= \int \left(\int (y - m) f_{Y|\mathbf{X}}(y|\mathbf{x}) dy \right) (m - f) f_{\mathbf{X}}(\mathbf{x}) dx \\ &= \int \left(\int (y - E(y|\mathbf{x})) f_{Y|\mathbf{X}}(y|\mathbf{x}) dy \right) (m - f) f_{\mathbf{X}}(\mathbf{x}) dx \\ &= \int \left(\int y f_{Y|\mathbf{X}}(y|\mathbf{x}) dy - \int E(y|\mathbf{x}) f_{Y|\mathbf{X}}(y|\mathbf{x}) dy \right) (m - f) f_{\mathbf{X}}(\mathbf{x}) dx \\ &= \int \left(E(y|\mathbf{x}) - E(y|\mathbf{x}) \int f_{Y|\mathbf{X}}(y|\mathbf{x}) dy \right) (m - f) f_{\mathbf{X}}(\mathbf{x}) dx \\ &= \int (E(y|\mathbf{x}) - E(y|\mathbf{x})) (m - f) f_{\mathbf{X}}(\mathbf{x}) dx \\ &= \int (0) (m - f) f_{\mathbf{X}}(\mathbf{x}) dx = 0 \end{aligned}$$

$$\text{Logo, } E\left((Y - f(X_1, \dots, X_p))^2\right) \geq E\left((Y - M(X_1, \dots, X_p))^2\right).$$

□

A variável $E(Y|X_1, \dots, X_p) = M(X_1, \dots, X_p)$ é o preditor de menor erro quadrático médio, denominada função de regressão.

Exemplo : Caso Normal

Seja \mathbf{x} uma variável aleatória multidimensional e y uma variável aleatória unidimensional. Suponha que:

$$\begin{bmatrix} \mathbf{x} \\ y \end{bmatrix} \sim N \left(\begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \Sigma & \sigma \\ \sigma' & \sigma_y^2 \end{bmatrix} \right)$$

em que, $\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$, $\sigma = \text{cov}(\mathbf{x}, y) = \begin{bmatrix} \text{cov}(x_1, y) \\ \vdots \\ \text{cov}(x_p, y) \end{bmatrix}$ e $\Sigma = \text{cov}(\mathbf{x})$.

Dada uma matriz particionada $\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$, sua inversa é dada por (RAO, 2002):

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

ou ainda,

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

Assim, a inversa da matriz $\begin{bmatrix} \Sigma & \sigma \\ \sigma' & \sigma_y^2 \end{bmatrix}$ é:

$$\begin{bmatrix} \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2}\right)^{-1} & -\frac{\Sigma^{-1}\sigma}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} \\ -\frac{\sigma'}{\sigma_y^2} \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2}\right)^{-1} & \frac{1}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} \end{bmatrix} = \begin{bmatrix} \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2}\right)^{-1} & -\left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2}\right)^{-1} \frac{\sigma}{\sigma_y^2} \\ -\frac{\sigma'\Sigma^{-1}}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} & \frac{1}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} \end{bmatrix}.$$

Na distribuição conjunta das variáveis \mathbf{x} e y , o expoente da exponencial

é:

$$\begin{aligned}
& \begin{bmatrix} \mathbf{x}' & y' \end{bmatrix} \begin{bmatrix} \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} & -\frac{\Sigma^{-1}\sigma}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} \\ -\frac{\sigma'}{\sigma_y^2} \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} & \frac{1}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{x}' \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} & -y' \frac{\sigma'}{\sigma_y^2} \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} & -\mathbf{x}' \frac{\Sigma^{-1}\sigma}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} & +y' \frac{1}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix} \\
&= \left(\mathbf{x}' \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} - y' \frac{\sigma'}{\sigma_y^2} \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \right) \mathbf{x} + \left(-\mathbf{x}' \frac{\Sigma^{-1}\sigma}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} + y' \frac{1}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} \right) y \\
&= \mathbf{x}' \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \mathbf{x} - y' \frac{\sigma'}{\sigma_y^2} \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \mathbf{x} - \mathbf{x}' \frac{\Sigma^{-1}\sigma}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} y + y' \frac{1}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} y \\
&= \mathbf{x}' \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \mathbf{x} - y' \frac{\sigma'}{\sigma_y^2} \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \mathbf{x} - \mathbf{x}' \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \frac{\sigma}{\sigma_y^2} y + y' \frac{1}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} y \\
&= \mathbf{x}' \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \mathbf{x} - 2\mathbf{x}' \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \frac{\sigma}{\sigma_y^2} y + y' \frac{1}{\sigma_y^2 - \sigma'\Sigma^{-1}\sigma} y.
\end{aligned}$$

Fazendo,

$$\mathbf{C} = \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \text{ e } \mathbf{H} = \frac{\sigma}{\sigma_y^2},$$

a expressão anterior fica da forma:

$$\mathbf{x}'\mathbf{C}\mathbf{x} - 2\mathbf{x}'\mathbf{C}\mathbf{H}y + y'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})y.$$

Logo, a marginal da variável \mathbf{x} é dada por (Apêndice A):

$$\begin{aligned}
f_{\mathbf{x}}(\mathbf{x}) &\propto \exp \left(\mathbf{x}' \left(\mathbf{C} - \mathbf{C}\mathbf{H}(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})^{-1}\mathbf{H}'\mathbf{C} \right) \mathbf{x} \right) \\
&= \exp \left(\mathbf{x}' \left(\left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} - \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \frac{\sigma}{\sigma_y^2} (\sigma_y^2 - \sigma'\Sigma^{-1}\sigma) \frac{\sigma'}{\sigma_y^2} \left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} \right) \mathbf{x} \right) \\
&= \exp \left(\mathbf{x}' \left(\left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} - \frac{\Sigma^{-1}\sigma}{(\sigma_y^2 - \sigma'\Sigma^{-1}\sigma)} (\sigma_y^2 - \sigma'\Sigma^{-1}\sigma) \frac{\sigma'\Sigma^{-1}}{(\sigma_y^2 - \sigma'\Sigma^{-1}\sigma)} \right) \mathbf{x} \right) \\
&= \exp \left(\mathbf{x}' \left(\left(\Sigma - \frac{\sigma\sigma'}{\sigma_y^2} \right)^{-1} - \frac{\Sigma^{-1}\sigma\sigma'\Sigma^{-1}}{(\sigma_y^2 - \sigma'\Sigma^{-1}\sigma)} \right) \mathbf{x} \right)
\end{aligned}$$

e a distribuição condicional fica da forma:

$$\begin{aligned}
f(y|\mathbf{x}) &= \frac{f(\mathbf{y}, \mathbf{x})}{f(\mathbf{x})} \\
&\propto \frac{\exp\left(\mathbf{x}'\left(\boldsymbol{\Sigma} - \frac{\boldsymbol{\sigma}\boldsymbol{\sigma}'}{\sigma_y^2}\right)^{-1}\mathbf{x} - 2\mathbf{x}'\left(\boldsymbol{\Sigma} - \frac{\boldsymbol{\sigma}\boldsymbol{\sigma}'}{\sigma_y^2}\right)^{-1}\frac{\boldsymbol{\sigma}}{\sigma_y^2}y + y'\frac{1}{\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}}y\right)}{\exp\left(\mathbf{x}'\left(\left(\boldsymbol{\Sigma} - \frac{\boldsymbol{\sigma}\boldsymbol{\sigma}'}{\sigma_y^2}\right)^{-1} - \frac{\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}\boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}}{(\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma})}\right)\mathbf{x}\right)} \\
&= \exp\left(\left(\mathbf{x}'\left(\boldsymbol{\Sigma} - \frac{\boldsymbol{\sigma}\boldsymbol{\sigma}'}{\sigma_y^2}\right)^{-1}\mathbf{x} - 2\mathbf{x}'\left(\boldsymbol{\Sigma} - \frac{\boldsymbol{\sigma}\boldsymbol{\sigma}'}{\sigma_y^2}\right)^{-1}\frac{\boldsymbol{\sigma}}{\sigma_y^2}y + y'\frac{1}{\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}}y\right) - \left(\mathbf{x}'\left(\boldsymbol{\Sigma} - \frac{\boldsymbol{\sigma}\boldsymbol{\sigma}'}{\sigma_y^2}\right)^{-1}\mathbf{x} - \mathbf{x}'\frac{\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}\boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}}{(\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma})}\mathbf{x}\right)\right) \\
&= \exp\left(-2\mathbf{x}'\left(\boldsymbol{\Sigma} - \frac{\boldsymbol{\sigma}\boldsymbol{\sigma}'}{\sigma_y^2}\right)^{-1}\frac{\boldsymbol{\sigma}}{\sigma_y^2}y + y'\frac{1}{\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}}y + \mathbf{x}'\frac{\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}\boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}}{(\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma})}\mathbf{x}\right) \\
&= \exp\left(-2\mathbf{x}'\frac{\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}}{\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}}y + y'\frac{1}{\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}}y + \mathbf{x}'\frac{\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}\boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}}{\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}}\mathbf{x}\right).
\end{aligned}$$

Comparando a expressão anterior com a forma geral $(y - w)' \frac{1}{a} (y - w)$, em relação ao termo cruzado, em que w é a média da distribuição y , tem-se:

$$\begin{aligned}
2w' \frac{1}{a} y &= 2\mathbf{x}' \frac{\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}}{\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}} y \\
\Rightarrow w' &= \mathbf{x}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\sigma} \text{ e } a = \sigma_y^2 - \boldsymbol{\sigma}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\sigma}.
\end{aligned}$$

De onde segue que:

$$f(y|\mathbf{x}) \propto \exp\left(\left(y - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\mathbf{x}\right)' \frac{1}{\sigma_y^2 - \boldsymbol{\sigma}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}} (y - \mathbf{x}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma})\right).$$

Logo, $\hat{Y} = \mu_y + \boldsymbol{\beta}'(\mathbf{x} - \boldsymbol{\mu}_x)$ com $\boldsymbol{\beta} = \boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}$.

Portanto, para o caso normal, a função de regressão é linear.

Uma questão importante para distribuições em geral é obter, entre todas as funções lineares nas covariáveis, $\alpha + \boldsymbol{\beta}'\mathbf{x}$, aquela que minimiza o erro quadrático médio $E\left((y - (\alpha + \boldsymbol{\beta}'\mathbf{x}))^2\right)$. Assim:

$$\begin{aligned}
& E\left((y - (\alpha + \beta' \mathbf{x}))^2\right) \\
&= E\left(y^2 - 2(\alpha + \beta' \mathbf{x})y + (\alpha + \beta' \mathbf{x})^2\right) \\
&= E(y^2) - 2E((\alpha + \beta' \mathbf{x})y) + E((\alpha + \beta' \mathbf{x})^2) \\
&= E(y^2) - 2(E(\alpha y) + E(\beta' \mathbf{x}y)) + E((\alpha + \beta' \mathbf{x})^2) \\
&= E(y^2) - 2\alpha E(y) - 2E((\beta' \mathbf{x})y) + E(\alpha^2 + 2\alpha\beta' \mathbf{x} + (\beta' \mathbf{x})^2) \\
&= E(y^2) - 2\alpha E(y) - 2\beta' E(\mathbf{x}y) + \alpha^2 + 2\alpha\beta' E(\mathbf{x}) + E((\beta' \mathbf{x})^2) \\
&= E(y^2) - 2\alpha E(y) - 2\beta' E(\mathbf{x}y) + \alpha^2 + 2\alpha\beta' E(\mathbf{x}) + E((\beta' \mathbf{x})^2) \\
&\quad - \left(E(\beta' \mathbf{x})^2\right) + \left(E(\beta' \mathbf{x})^2\right) \\
&= E(y^2) - 2\alpha E(y) - 2\beta' E(\mathbf{x}y) + (\alpha + \beta' E(\mathbf{x}))^2 + \text{var}(\beta' \mathbf{x}) \\
&= E(y^2) - 2\alpha E(y) - 2\beta' E(\mathbf{x}y) + (\alpha + \beta' E(\mathbf{x}))^2 + \beta' \text{cov}(\mathbf{x})\beta.
\end{aligned}$$

Derivando em relação a α e igualando a zero, obtém-se:

$$\begin{aligned}
& -2E(y) + 2(\alpha + \beta' E(\mathbf{x})) = 0 \\
& \Rightarrow \alpha^* = E(y) - \beta' E(\mathbf{x}).
\end{aligned}$$

Derivando em relação a β e igualando a zero, obtém-se:

$$\begin{aligned}
& -2E(\mathbf{x}y) + 2(\alpha + \beta' E(\mathbf{x})) E(\mathbf{x}) + 2\text{cov}(\mathbf{x})\beta \\
&= -2E(\mathbf{x}y) + 2E(y) E(\mathbf{x}) + 2\text{cov}(\mathbf{x})\beta = 0
\end{aligned}$$

$$\text{cov}(\mathbf{x})\boldsymbol{\beta} = \text{E}(\mathbf{x}y) - \text{E}(y)\text{E}(\mathbf{x})$$

$$\text{cov}(\mathbf{x})\boldsymbol{\beta}^* = \text{cov}(y, \mathbf{x})$$

$$\boldsymbol{\beta}^* = (\text{cov}(\mathbf{x}))^{-1} \text{cov}(y, \mathbf{x})$$

$$\boldsymbol{\beta}^* = \boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma}.$$

Esta combinação linear, $\alpha^* + \boldsymbol{\beta}^{*\prime} \mathbf{x}$, tem covariância máxima com y , em relação a qualquer outra combinação linear $\alpha + \boldsymbol{\beta}' \mathbf{x}$, com $\boldsymbol{\beta}' \text{cov}(\mathbf{x}) \boldsymbol{\beta} = \boldsymbol{\beta}^{*\prime} \text{cov}(\mathbf{x}) \boldsymbol{\beta}^*$.

Demonstração.

$$\begin{aligned} \boldsymbol{\beta}^{*\prime} \text{cov}(\mathbf{x}) \boldsymbol{\beta}^* &= \left((\text{cov}(\mathbf{x}))^{-1} \text{cov}(y, \mathbf{x}) \right)' \text{cov}(\mathbf{x}) \left((\text{cov}(\mathbf{x}))^{-1} \text{cov}(y, \mathbf{x}) \right) \\ &= (\text{cov}(y, \mathbf{x}))' (\text{cov}(\mathbf{x}))^{-1} \text{cov}(\mathbf{x}) (\text{cov}(\mathbf{x}))^{-1} \text{cov}(y, \mathbf{x}) \\ &= (\text{cov}(y, \mathbf{x}))' (\text{cov}(\mathbf{x}))^{-1} \text{cov}(y, \mathbf{x}) \\ &= \boldsymbol{\sigma}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\sigma}, \end{aligned}$$

$$\text{e cov}(y, \alpha + \boldsymbol{\beta}' \mathbf{x}) = \text{cov}(y, \boldsymbol{\beta}' \mathbf{x}) = \boldsymbol{\beta}' \text{cov}(y, \mathbf{x}) = \boldsymbol{\beta}' \boldsymbol{\sigma}.$$

Assim, tem-se que maximizar $\boldsymbol{\beta}' \boldsymbol{\sigma}$, com a restrição $\boldsymbol{\beta}' \text{cov}(\mathbf{x}) \boldsymbol{\beta} = \boldsymbol{\beta}^{*\prime} \text{cov}(\mathbf{x}) \boldsymbol{\beta}^* = \boldsymbol{\sigma}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\sigma}$.

A lagrangiana é dada por:

$$H(\boldsymbol{\beta}, \lambda) = \boldsymbol{\beta}' \boldsymbol{\sigma} + \lambda \left(\boldsymbol{\beta}' \text{cov}(\mathbf{x}) \boldsymbol{\beta} - \boldsymbol{\sigma}' (\text{cov}(\mathbf{x}))^{-1} \boldsymbol{\sigma} \right).$$

Derivando e igualando a zero, obtém-se

$$\frac{dH}{d\lambda} = 0 \Rightarrow \beta' \text{cov}(\mathbf{x}) \beta = \sigma' (\text{cov}(\mathbf{x}))^{-1} \sigma \quad (1)$$

$$\frac{dH}{d\beta} = 0 \Rightarrow \sigma + \lambda (2 \text{cov}(\mathbf{x}) \beta) = 0. \quad (2)$$

Segue que

$$\sigma = -2\lambda \text{cov}(\mathbf{x}) \beta \Rightarrow \beta = -\frac{(\text{cov}(\mathbf{x}))^{-1} \sigma}{2\lambda}. \quad (3)$$

Substituindo (3) na equação (1), tem-se

$$\begin{aligned} & \left(-\frac{(\text{cov}(\mathbf{x}))^{-1} \sigma}{2\lambda} \right)' \text{cov}(\mathbf{x}) \left(-\frac{(\text{cov}(\mathbf{x}))^{-1} \sigma}{2\lambda} \right) = \sigma' (\text{cov}(\mathbf{x}))^{-1} \sigma \\ \Rightarrow & \frac{1}{4\lambda^2} \sigma' (\text{cov}(\mathbf{x}))^{-1} \sigma = \sigma' (\text{cov}(\mathbf{x}))^{-1} \sigma \\ \Rightarrow & \frac{1}{4\lambda^2} = 1 \Rightarrow 4\lambda^2 = 1 \Rightarrow \lambda = \pm \frac{1}{2}. \end{aligned}$$

Logo,

$$\begin{aligned} \beta &= -\frac{(\text{cov}(\mathbf{x}))^{-1} \sigma}{2\lambda} = -\frac{(\text{cov}(\mathbf{x}))^{-1} \sigma}{2 \left(\pm \frac{1}{2} \right)} \\ &= \pm (\text{cov}(\mathbf{x}))^{-1} \sigma = \pm \Sigma^{-1} \sigma = \pm \beta^*. \end{aligned}$$

O valor máximo ocorre quando o sinal é positivo, assim, $\beta = \beta^*$.

□

Segue, então, que o melhor preditor linear (*Best Linear Predictor* - BLP) de menor erro quadrático médio é dado por:

$$\begin{aligned}\alpha^* + \mathbf{x}'\boldsymbol{\beta}^* &= E(y) - \boldsymbol{\beta}^{*\prime}E(\mathbf{x}) + \boldsymbol{\beta}^{*\prime}\mathbf{x} \\ &= E(y) + \boldsymbol{\beta}^{*\prime}(\mathbf{x} - E(\mathbf{x})).\end{aligned}$$

Neste caso, o EQM do preditor é:

$$\begin{aligned}& E\left(\left(y - \left(E(y) + \boldsymbol{\beta}^{*\prime}(\mathbf{x} - E(\mathbf{x}))\right)\right)^2\right) \\ &= E\left(\left(y - E(y)\right)^2 - 2\left(y - E(y)\right)\boldsymbol{\beta}^{*\prime}(\mathbf{x} - E(\mathbf{x})) + \left(\boldsymbol{\beta}^{*\prime}(\mathbf{x} - E(\mathbf{x}))\right)^2\right) \\ &= E\left(y - E(y)\right)^2 - 2E\left(\left(y - E(y)\right)\boldsymbol{\beta}^{*\prime}(\mathbf{x} - E(\mathbf{x}))\right) + E\left(\boldsymbol{\beta}^{*\prime}(\mathbf{x} - E(\mathbf{x}))\right)^2 \\ &= \text{var}(y) - 2\boldsymbol{\beta}^{*\prime}\text{cov}(y, \mathbf{x}) + \boldsymbol{\beta}^{*\prime}\text{cov}(\mathbf{x})\boldsymbol{\beta}^* \\ &= \text{var}(y) - 2\left(\text{cov}(\mathbf{x})^{-1}\text{cov}(y, \mathbf{x})\right)'\text{cov}(y, \mathbf{x}) \\ &\quad + \left(\text{cov}(\mathbf{x})^{-1}\text{cov}(y, \mathbf{x})\right)'\text{cov}(\mathbf{x})\left(\text{cov}(\mathbf{x})^{-1}\text{cov}(y, \mathbf{x})\right) \\ &= \text{var}(y) - \left(\text{cov}(\mathbf{x})^{-1}\text{cov}(y, \mathbf{x})\right)'\text{cov}(y, \mathbf{x}) \\ &= \sigma_y^2 - (\boldsymbol{\Sigma}^{-1}\boldsymbol{\sigma})'\boldsymbol{\sigma}.\end{aligned}$$

Observação: $\boldsymbol{\beta}^*$ é um vetor paramétrico, não é uma estimativa. O interessante é que este vetor paramétrico depende das variâncias e covariâncias e não da forma específica da distribuição populacional $f_{y,\mathbf{x}}(y, \mathbf{x})$.

2.4 Componentes Principais

A análise de componentes principais (*Principal Component Analysis* - PCA) pode ser aplicada como um método de redução de dimensionalidade, visando a solucionar eventuais problemas de multicolinearidade da matriz \mathbf{X} e a necessidade de um número excessivo de amostras para a construção de um modelo viável, sem acarretar na perda significativa de informações presentes nos dados (OTTO, 1999).

O PCA pode ser feito com base em autovalores e autovetores de uma matriz de covariância. Os elementos dos autovetores são os coeficientes requeridos para a combinação linear das variáveis originais e são conhecidos como *loadings*. Os elementos individuais das novas variáveis são derivados das variáveis explicativas e são conhecidos como os *scores*. A variância total dos dados originais é dada pela soma de suas variâncias individuais e os autovalores representam a variância associada a cada componente principal.

2.5 Regressão por Componentes Principais

A regressão via componentes principais (*Principal Components Regression* - PCR) foi introduzida por Kendall (1957) e Hotelling (1957). Este método faz uso dos procedimentos empregados na análise de componentes principais visando a contornar os problemas apresentados pela regressão linear múltipla.

A metodologia do método PCR consiste, basicamente, em eliminar componentes que não contribuam na explicação da variância presente nos dados, o que reduz a dimensionalidade dos dados originais. Após a escolha do número ótimo de componentes a serem utilizados, os coeficientes podem ser determinados por

meio do método dos quadrados mínimos ordinários.

Como em qualquer regressão múltipla, a escolha do número de variáveis a serem incluídas no modelo é de extrema importância, tendo em vista evitar a perda de informações relevantes. Segundo Roggo *et al.* (2007), cada componente descreve uma fração da variação total contida nos dados, tornando possível a determinação do número ótimo de componentes a serem incluídos na regressão.

2.6 Regressão por Quadrados Mínimos Parciais

Quadrados Mínimos Parciais (*Partial Least Squares* - PLS) foi introduzido por Wold em 1975, sendo considerado útil para a construção de equações de predições em situações nas quais se tem um grande número de variáveis explicativas e um número relativamente pequeno de dados amostrais (HOSKULDSSON, 1988). Resumidamente, a ideia geral do PLS é formar componentes que capturem a maior quantidade de informação possível disposta nas variáveis explicativas (X_1, \dots, X_p) para prever as variáveis respostas (Y_1, \dots, Y_k). O método PLS apresenta similaridades com o método de Regressão via Componentes Principais, sendo a maior diferença entre eles dada pelo fato do PCR levar em consideração apenas as variáveis explicativas na construção dos componentes, enquanto que o PLS também leva em consideração as variáveis respostas.

3 METODOLOGIA

3.1 Componentes Principais

Foi feita uma abordagem da obtenção dos Componentes Principais por meio do cálculo diferencial, além do tradicional método multiplicadores de Lagrange. Geometricamente, objetivou-se encontrar a direção de maior variabilidade do vetor aleatório \mathbf{X} (covariáveis).

3.2 Regressão por Componentes Principais

A regressão em componentes principais é colocada neste trabalho em razão das semelhanças com o PLS. Na linha deste trabalho, a exposição da teoria do PCR foi bastante geométrica e, portanto, diferente da abordagem usual na literatura.

Os componentes são obtidos pelas direções no subespaço gerado pelas observações das covariáveis que fornecem a maior variabilidade dos dados. Na construção dos componentes, utiliza-se apenas as covariáveis, e estes componentes representam as novas variáveis a serem utilizadas na regressão. Os coeficientes da regressão podem ser estimados pelo método dos quadrados mínimos.

Geometricamente, foi demonstrado que os componentes são representados como imagem dos autovetores da matriz $\mathbf{X}'\mathbf{X}$. Obtendo o subespaço gerado por estes componentes, a regressão é feita pela projeção do vetor de dados neste subespaço. Também é demonstrado que o vetor de coeficientes estimado via Regressão em Componentes Principais, $\hat{\beta}_{\text{PCR}}$, é um estimador de encolhimento.

3.3 Regressão por Quadrados Mínimos Parciais

Uma ferramenta para resolver problemas de regressão nos quais as variáveis preditoras (ou covariáveis) apresentam um alto nível de colinearidade é o uso da teoria dos Quadrados Mínimos Parciais. Neste caso, a matriz das covariáveis \mathbf{X} pode ser de posto incompleto. Outra situação, muito comum em aplicações na área de química, ocorre quando o número de covariáveis é muito maior que o número de observações. Em ambas as situações, as equações normais para obtenção do estimador de quadrados mínimos admitem infinitas soluções. Para a escolha de uma estimativa em particular, um método possível é o uso de inversas generalizadas. No entanto, a escolha desta inversa é, de certa forma, arbitrária, gerando situações indesejáveis para as equações de predição obtidas. No método de regressão utilizando PLS, uma estimativa única para o parâmetro de regressão é obtida.

O método é semelhante ao de Regressão em Componentes Principais, mas enquanto este só depende da matriz das covariáveis \mathbf{X} , o método PLS depende também do vetor de dados \mathbf{Y} . A ideia básica é regredir o vetor de dados \mathbf{Y} em cada um dos vetores das covariáveis e, através de escolha de pesos, obter combinações lineares destas regressões. Desta forma, o vetor de dados \mathbf{Y} é explicado em termos de um número menor de vetores. A escolha destas combinações lineares é justificada estatisticamente, pela maximização de covariâncias.

A estimativa do vetor de parâmetros através do método PLS, $\hat{\beta}_{\text{PLS}}$, admite uma expressão como uma transformação linear da estimativa $\hat{\beta}_{\text{OLS}}$, invariante em relação ao $\hat{\beta}_{\text{OLS}}$ escolhido. O método possui uma natureza geométrica sendo obtido por meio de projeções oblíquas, diferentemente do método dos quadrados mínimos, que utiliza projeções ortogonais.

A literatura referente ao método PLS, apesar de abundante, é em sua mai-

oria, aplicada em problemas oriundos da Quimiometria e Genética, por exemplo. Poucas referências existem relativas à parte teórica do método. Em geral, os artigos se restringem a apresentar os algoritmos básicos utilizados pelo método. As referências mais acessíveis relacionadas à teoria geométrica e algébrica do PLS e à descrição dos algoritmos são os artigos Helland (1990), Hoskuldsson (1988), Garthwaite (1994) e Phatak e Jong (1997). Sendo assim, o texto que segue discute o método PLS e suas propriedades, essencialmente, sob o ponto de vista dos autores destes quatro artigos, apresentando, sempre que possível, uma abordagem geométrica. Além disso, é feito um paralelo entre os algoritmos para o PLS populacional e amostral.

3.4 Determinação do número ótimo de componentes

Uma questão a ser abordada é o número de componentes que deve ser utilizado na regressão em PLS e, para isto, é necessário utilizar alguns conceitos de seleção de modelos. Dentre os métodos propostos na literatura, são abordados os conceitos de Graus de Liberdade (*Degrees of Freedom* - DoF) e Validação Cruzada (*Cross Validation* - CV).

São apresentadas as definições do DoF para modelos de regressão linear e não linear, bem como a demonstração para uma cota inferior do DoF quando $m = 1$. Ainda, dois gráficos ilustram como o DoF e o Erro Quadrático Médio (EQM - utilizado na validação cruzada) se comportam em função do número de componentes. Os gráficos foram gerados no software R (R Core Team, 2013), por meio das funções *pls.model* e *pls.cv* do pacote *plsdo*, respectivamente, para DoF e CV.

3.5 Exemplo Didático

Foi elaborada uma rotina no software R (Apêndice D) com o objetivo de explicitar passo a passo os resultados do algoritmo PLS (HOSKULDSSON, 1988). Um exemplo didático foi apresentado para ilustrar tal rotina.

Inicialmente, gerou-se uma matriz de covariâncias, 5×5 , da forma

$$\begin{bmatrix} \Sigma & \sigma \\ \sigma' & \sigma_y^2 \end{bmatrix}, \text{ em que } \sigma = \text{cov}(\mathbf{X}, Y), \quad \Sigma = \text{cov}(\mathbf{X}) \text{ e } \sigma_y^2 = \text{var}(Y).$$

Utilizando esta matriz de covariâncias, gerou-se de uma Normal Multivariada, uma matriz de covariáveis \mathbf{X} , com 4 covariáveis, e um vetor para a variável de interesse Y , ambos com $n = 6$ observações.

Obtidos \mathbf{X} e Y , implementou-se o algoritmo PLS para a construção dos componentes. Para efeito de comparação, estimou-se o vetor de parâmetros $\hat{\beta}$ pela regressão por Quadrados Mínimos Ordinários, $\hat{\beta}_{OLS}$, que só depende de \mathbf{X} e de Y , e pela regressão por Quadrados Mínimos Parciais, $\hat{\beta}_{PLS}$, utilizando os m componentes obtidos no algoritmo, com $m = 1, 2, 3$ e 4 .

Uma vez obtidos os parâmetros da regressão, gerou-se um novo valor para cada uma das quatro covariáveis (\mathbf{X}_p), e um novo valor para a variável de interesse (Y_p), a fim de fazer predição. Com os novos valores de \mathbf{X}_p e com os vetores $\hat{\beta}_{OLS}$ e $\hat{\beta}_{PLS}$, encontrou-se o valor predito para Y_p ($\hat{Y}_p = \mathbf{X}_p \hat{\beta}$), denominados $\hat{Y}_{p,OLS}$ e $\hat{Y}_{p,PLS}$.

4 RESULTADOS

4.1 Componentes Principais

A teoria de componentes principais está relacionada ao problema básico: dado um vetor aleatório \mathbf{Z} , encontrar a direção em que este vetor apresenta a maior variação. Geometricamente, o problema pode ser descrito como: dada uma direção $\mathbf{a} \in \mathbb{R}^n$, $\|\mathbf{a}\| = 1$, projeta-se \mathbf{Z} na linha definida por \mathbf{a} (Figura 8) obtendo-se um vetor aleatório totalmente descrito por sua norma $\mathbf{a} \cdot \mathbf{Z} = \|\mathbf{Z}\| \cos \theta$ que define uma variável aleatória unidimensional, \mathbf{Y} . A variância desta é a variabilidade do vetor \mathbf{Z} na direção de \mathbf{a} . Nas condições do Teorema 1, tem-se:

$$\text{var}(\mathbf{Y}) = \text{var}(\mathbf{a} \cdot \mathbf{Z}) = \mathbf{a}' \boldsymbol{\Sigma} \mathbf{a}.$$

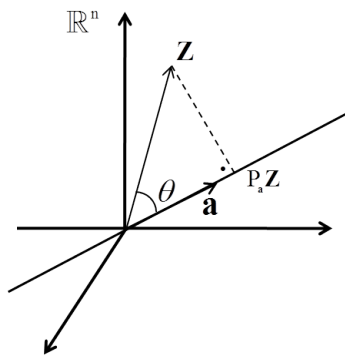


Figura 8 Representação geométrica da obtenção do componente principal

É necessário, então, maximizar $\text{var}(\mathbf{a} \cdot \mathbf{Z}) = \mathbf{a}' \boldsymbol{\Sigma} \mathbf{a}$ com a restrição $\|\mathbf{a}\| = 1$. Serão apresentados dois métodos de maximização: utilizando multiplicadores de Lagrange e Cálculo Diferencial.

Utilizando multiplicadores de Lagrange, tem-se:

$$\mathbf{H}(\mathbf{a}, \lambda) = \mathbf{a}'\Sigma\mathbf{a} - \lambda \left(\|\mathbf{a}\|^2 - 1 \right)$$

$$\frac{\partial \mathbf{H}}{\partial \mathbf{a}} = 2\Sigma\mathbf{a} - 2\lambda\mathbf{a} = 0 \Rightarrow \Sigma\mathbf{a} = \lambda\mathbf{a}.$$

Portanto, os pontos críticos, \mathbf{a} , são os autovetores da matriz Σ .

Uma outra demonstração pode ser feita através do cálculo diferencial: seja $a(t) = (a_1(t), \dots, a_n(t))$ uma curva parametrizada com $\|a(t)\| = 1$, isto é, uma curva sobre a esfera de raio unitário, centrada na origem do espaço \mathbb{R}^n . Vamos também supor que esta curva esteja parametrizada pelo comprimento do arco, isto é, $\left\| \frac{da}{dt}(t) \right\| = 1$. Definindo a função $f(a(t)) = a'(t) \Sigma a(t)$, o que se quer obter são os pontos críticos de f .

Derivando e igualando a zero,

$$\frac{df}{dt}(0) = \frac{df}{da}(a(0)) \frac{da}{dt}(0) = 0$$

Como $\frac{df}{da}(a(0)) = 2\Sigma a(0)$ tem-se

$$2(\Sigma a(0))' \frac{da}{dt}(0) = 0$$

e segue, portanto, que o vetor $\frac{da}{dt}(0)$ é perpendicular ao vetor $\Sigma a(0)$.

Como tal fato tem que ser verdade para qualquer curva que passa em $a(0)$ significa que $\Sigma a(0)$ é um vetor perpendicular ao espaço tangente da esfera em $a(0)$. Neste caso $\Sigma a(0)$ tem que ser um múltiplo de $a(0)$ e, portanto, $\Sigma a(0) = \lambda a(0)$. Assim, os pontos críticos de f são justamente os autovetores de Σ .

A natureza dos pontos críticos pode ser estudada obtendo-se a derivada

segunda:

$$\begin{aligned} f'(t) &= 2a'(t) \Sigma \frac{da}{dt}(t) \\ f''(t) &= 2 \frac{da'}{dt}(t) \Sigma \left(\frac{da}{dt}(t) \right) + 2a'(t) \Sigma \frac{d^2a}{dt^2}(t) \\ f''(0) &= 2 \frac{da'}{dt}(0) \Sigma \left(\frac{da}{dt}(0) \right) + 2a'(0) \Sigma \frac{d^2a}{dt^2}(0). \end{aligned}$$

Como $\left\| \frac{da}{dt}(t) \right\| = 1$ (parametrizado pelo comprimento do arco), tem-se:

$$1 = \left(\frac{da}{dt}(t) \right)' \left(\frac{da}{dt}(t) \right) \Rightarrow 0 = 2 \left(\frac{d^2a}{dt^2} \right)' \frac{da}{dt}(t).$$

Logo, $\frac{d^2a}{dt^2}(0)$ é paralelo a $a(0)$ e, portanto, também é um autovetor de Σ , isto é,

$$\Sigma \left(\frac{d^2a}{dt^2}(0) \right) = \lambda \left(\frac{d^2a}{dt^2}(0) \right) = \lambda k a(0),$$

em que k é negativo (note que $\frac{d^2a}{dt^2}(t)$ é a aceleração centrípeta). Assim,

$$\begin{aligned} f''(0) &= 2 \left(\frac{da}{dt}(0) \right)' \Sigma \left(\frac{da}{dt}(0) \right) + 2\lambda k \|a(0)\|^2 \\ &= 2 \left(\frac{da}{dt}(0) \right)' \Sigma \left(\frac{da}{dt}(0) \right) + 2\lambda k, \end{aligned}$$

e $\left(\frac{da}{dt}(0) \right)' \Sigma \left(\frac{da}{dt}(0) \right) \geq 0$ pois Σ é positiva semidefinida. Os autovalores de Σ são positivos pelo mesmo motivo. A curva $a(t)$ pode ser sempre pensada como

$$a(t) = \cos(t) \mathbf{i} + \sin(t) \mathbf{j} \Rightarrow \frac{d^2a}{dt^2}(t) = -\cos(t) \mathbf{i} - \sin(t) \mathbf{j} = -a(t).$$

Dessa forma, $k = -1$, e temos $f''(0) = 2 \left(\frac{da}{dt}(0) \right)' \Sigma \left(\frac{da}{dt}(0) \right) - \lambda$.

Portanto, $a(0)$ é um ponto de máximo se

$$\left. \frac{d^2 f}{dt^2}(a(t)) \right|_{t=0} < 0 \Rightarrow \left(\frac{da}{dt}(0) \right)' \Sigma \left(\frac{da}{dt}(0) \right) < \lambda,$$

para todo vetor $\frac{da}{dt}(0)$ no espaço tangente da esfera S^n em $a(0)$ (Figura 9).

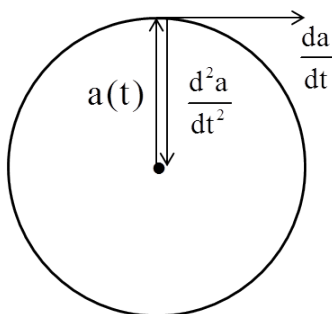


Figura 9 Representação geométrica da curva parametrizada pelo comprimento do arco sobre a esfera de raio unitário, centrada na origem do espaço \mathbb{R}^n

Obtidos os autovetores \mathbf{a} de Σ , as componentes principais \mathbf{Y} são obtidas como as combinações lineares $\mathbf{a} \cdot \mathbf{Z}$. Onde conclui-se que para o maior autovalor se tem o máximo, para o menor autovalor o mínimo, e os outros pontos críticos são selas.

4.2 Regressão por Componentes Principais: uma abordagem geométrica

Esta seção tem como referência o artigo Phatak e Jong (1997).

Para o problema de regressão em que se tem o modelo $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, a ideia básica é obter a direção, no subespaço $\text{Im}\mathbf{X}$, que fornece a maior variabilidade de \mathbf{Y} .

Neste caso, o problema de multiplicadores de Lagrange, agora, é mais complicado uma vez que se têm duas restrições ($\|\boldsymbol{\beta}\| = 1$ e subespaço $\text{Im}\mathbf{X}$). No entanto, pode-se modificar o problema para maximizar $\text{var}(\mathbf{X}\boldsymbol{\beta} \cdot \mathbf{Y})$ restrito a $\boldsymbol{\beta}'\boldsymbol{\beta} = 1$, isto é, maximizar $(\mathbf{X}\boldsymbol{\beta})' \boldsymbol{\Sigma} (\mathbf{X}\boldsymbol{\beta})$, restrito a $\|\boldsymbol{\beta}\| = 1$, ou seja, $\boldsymbol{\beta}'\mathbf{X}'\boldsymbol{\Sigma} \mathbf{X}\boldsymbol{\beta}$, restrito a $\|\boldsymbol{\beta}\| = 1$, como mostrado na Figura 10.

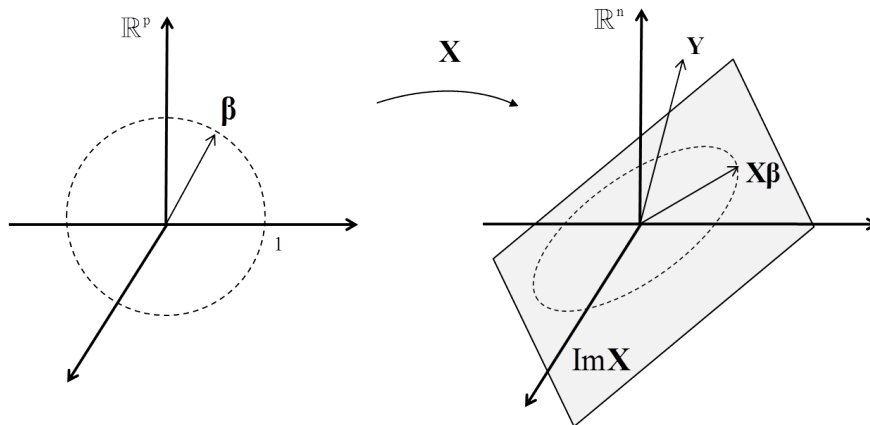


Figura 10 Representação geométrica da maximização da $\text{var}(\mathbf{X}\boldsymbol{\beta} \cdot \mathbf{Y})$ restrito a $\boldsymbol{\beta}'\boldsymbol{\beta} = 1$

Este problema é semelhante ao anteriormente feito, com a direção de maior variação dada agora no espaço paramétrico. As direções são dadas pelos autovetores de $\mathbf{X}'\boldsymbol{\Sigma}\mathbf{X}$. O caso mais simples ocorre quando $\boldsymbol{\Sigma} = \sigma^2\mathbf{I}$, e então, as direções são dadas pelos autovetores de $\mathbf{X}'\mathbf{X}$.

Sejam $\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2, \dots, \boldsymbol{\gamma}_p$ os autovetores e $\lambda_1, \lambda_2, \dots, \lambda_p$ os respectivos autova-

lores, da matriz $\mathbf{X}'\mathbf{X}$. Assim, $\mathbf{X}'\mathbf{X}\gamma_i = \lambda_i\gamma_i$ com $\gamma_i'\gamma_i = 1$ e $\gamma_i'\gamma_j = 0$ para $i \neq j : 1, \dots, p$ (Figura 11).

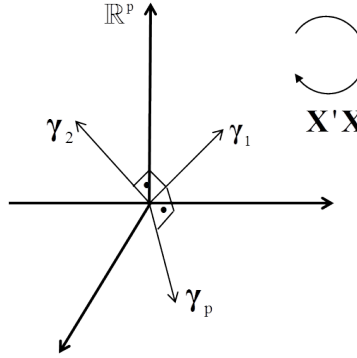


Figura 11 Representação geométrica dos autovetores $\gamma_1, \gamma_2, \dots, \gamma_p$ da matriz $\mathbf{X}'\mathbf{X}$

Note que,

$$\text{var}(\mathbf{X}\gamma_i \cdot \mathbf{Y}) = (\mathbf{X}\gamma_i)' (\sigma^2 \mathbf{I}) (\mathbf{X}\gamma_i) = \sigma^2 \gamma_i' \mathbf{X}'\mathbf{X}\gamma_i = \sigma^2 \gamma_i' (\lambda_i \gamma_i) = \sigma^2 \lambda_i.$$

Os vetores $\mathbf{X}\gamma_1, \mathbf{X}\gamma_2, \dots, \mathbf{X}\gamma_p$ são ortogonais, pois:

$$\mathbf{X}\mathbf{X}'(\mathbf{X}\gamma_i) = \mathbf{X}(\mathbf{X}'\mathbf{X})\gamma_i = \mathbf{X}\lambda_i\gamma_i = \lambda_i(\mathbf{X}\gamma_i).$$

Assim, os $\mathbf{X}\gamma_i$ são autovetores de $\mathbf{X}\mathbf{X}'$, portanto, ortogonais.

Ordenando os autovalores $\lambda_1 > \lambda_2 > \dots > \lambda_p$, definem-se como componentes principais as direções dadas por $\mathbf{X}\gamma_1, \mathbf{X}\gamma_2, \mathbf{X}\gamma_3, \dots$. Como estas direções apresentam uma ordem decrescente de variação do vetor aleatório \mathbf{Y} , elas “explicam” bem o vetor \mathbf{Y} . A ideia, então, é formar os subespaços (Figura 12):

$$W_1 = \text{span} \{ \mathbf{X}\gamma_1 \}$$

$$W_2 = \text{span} \{ \mathbf{X}\gamma_1, \mathbf{X}\gamma_2 \}$$

$$W_3 = \text{span} \{ \mathbf{X}\gamma_1, \mathbf{X}\gamma_2, \mathbf{X}\gamma_3 \}$$

$$\vdots$$

$$W_m = \text{span} \{ \mathbf{X}\gamma_1, \mathbf{X}\gamma_2, \mathbf{X}\gamma_3, \dots, \mathbf{X}\gamma_m \}$$

$$\text{Im}\mathbf{X} = \text{span} \{ \mathbf{X}\gamma_1, \mathbf{X}\gamma_2, \mathbf{X}\gamma_3, \dots, \mathbf{X}\gamma_p \}.$$

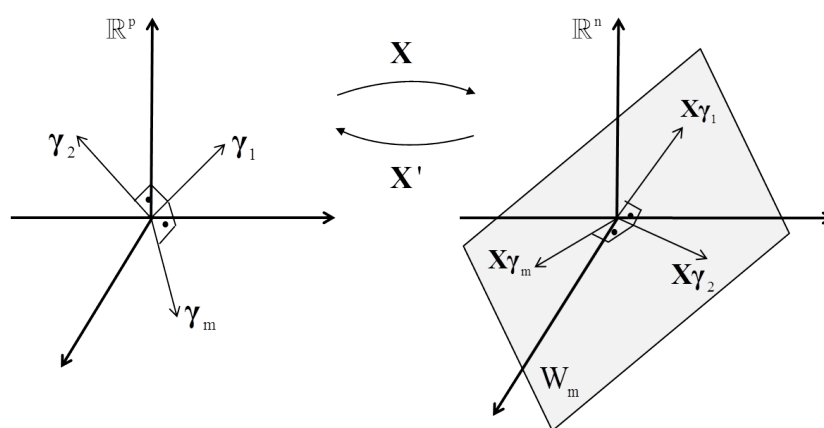


Figura 12 Subespaço W_m gerado pelos m -primeiros autovetores da matriz $\mathbf{X}'\mathbf{X}$

Como estas direções explicam a variabilidade de \mathbf{Y} , a ideia é projetar \mathbf{Y} ortogonalmente no subespaço W_m , isto é, regressar \mathbf{Y} aos m -primeiros componentes principais. A projeção define, então, o parâmetro de regressão denominado $\hat{\beta}_{\text{PCR}}$, cujas equações normais são dadas por $P_{W_m}\mathbf{Y} = \mathbf{X}\hat{\beta}_{\text{PCR}}$. Geometricamente, fica como representado na Figura 13.

O triângulo retângulo formado pelos vetores $P_{\text{Im}\mathbf{X}}\mathbf{Y}$, $P_{W_m}\mathbf{Y}$ e pelo vetor dado pela diferença entre eles, garante que $\|P_{W_m}\mathbf{Y}\| \leq \|P_{\text{Im}\mathbf{X}}\mathbf{Y}\|$ e, apesar deste fato não implicar em $\|\hat{\beta}_{\text{PCR}}\| < \|\hat{\beta}\|$, o procedimento da regressão em componentes principais é claramente conservador, no sentido que fornece estimativas menores para o vetor de parâmetros β . De fato, será demonstrado que $\|\hat{\beta}_{\text{PCR}}\| < \|\hat{\beta}\|$. Resta então, obter expressões matriciais para as equações nor-

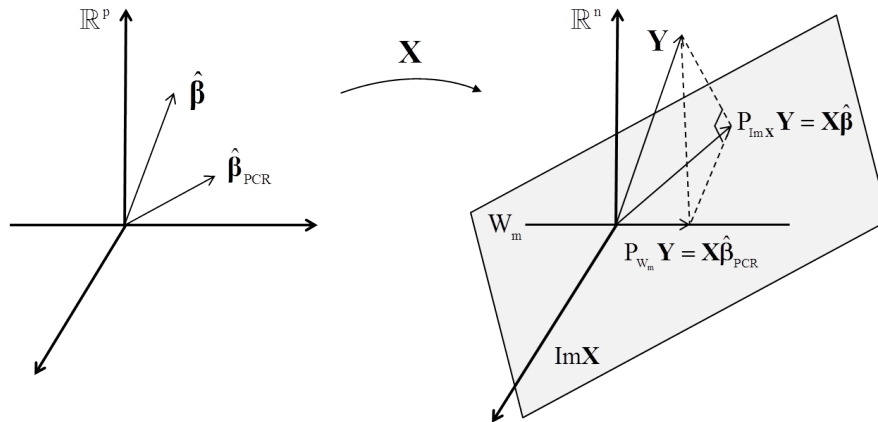


Figura 13 Projeção ortogonal de \mathbf{Y} no subespaço W_m definindo o parâmetro de regressão $\hat{\beta}_{\text{PCR}}$

mais, isto é, expressar o projetor P_{W_m} matricialmente.

A construção será heurística baseada no argumento: seja $f : A \rightarrow B$ uma função injetiva entre os conjuntos A e B . O que se quer é uma aplicação $g : B \rightarrow B$ que seja uma projeção em $\text{Im}f$, isto é, $g(f(x)) = f(x), \forall x \in A$ e $g(g(y)) = g(y), \forall y \in B$, como mostra a Figura 14.

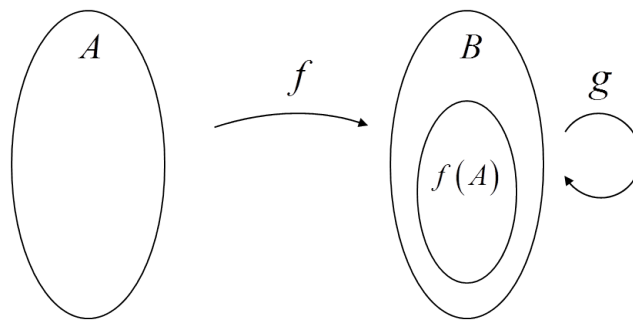


Figura 14 Representação da função injetiva f entre os conjuntos A e B e da função g como uma projeção na $\text{Im}f$

A ideia básica é que a composição de uma função injetiva com uma sobrejetiva define uma bijeção. Seja $h : B \rightarrow A$ uma função sobrejetiva tal que h

restrita à imagem de f seja injetiva, $h \circ f : A \rightarrow A$ é uma bijeção (Figura 15).

Define-se então,

$$g = f \circ (h \circ f)^{-1} \circ h. \quad (4)$$

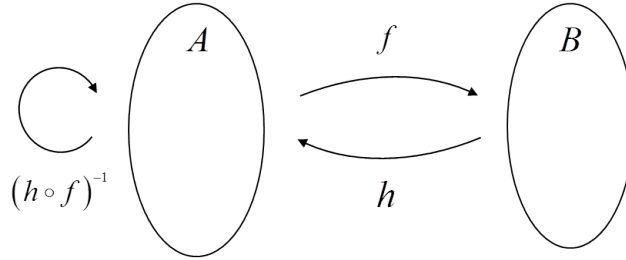


Figura 15 Representação da composição de uma função h sobrejetiva com uma função f injetiva, definindo uma bijeção $h \circ f$

$$\begin{aligned} g \circ g &= \left(f \circ (h \circ f)^{-1} \circ h \right) \circ \left(f \circ (h \circ f)^{-1} \circ h \right) \\ &= f \circ (h \circ f)^{-1} \circ (h \circ f) \circ (h \circ f)^{-1} \circ h \\ &= f \circ (h \circ f)^{-1} \circ h = g \quad e \end{aligned}$$

$$\begin{aligned} g(f(x)) &= f \circ (h \circ f)^{-1} \circ h(f(x)) \\ &= f \circ (h \circ f)^{-1} (h \circ f)(x) = f(x). \end{aligned}$$

Portanto, g é uma projeção em $\text{Im}f$, como procurado. Basta agora fazer a construção de g utilizando as funções lineares convenientes para se obter a projeção linear em W_m .

Sejam $\xi_1 = \mathbf{X}\gamma_1, \xi_2 = \mathbf{X}\gamma_2, \dots, \xi_m = \mathbf{X}\gamma_m$, como na Figura 16.

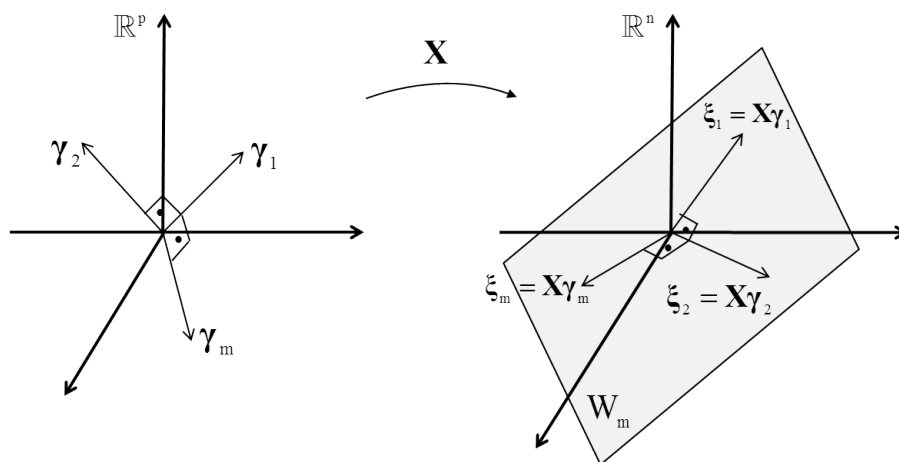


Figura 16 Representação dos vetores $\xi_1, \xi_2, \dots, \xi_m$ como imagem dos autovetores $\gamma_1, \gamma_2, \dots, \gamma_m$ pela matriz \mathbf{X}

Seja Ξ_m a matriz $n \times m$ cujas colunas são os vetores ξ_i ,

$$\Xi_m = [\xi_1, \xi_2, \dots, \xi_m].$$

A matriz Ξ_m define uma transformação linear $\Xi_m : \mathbb{R}^m \rightarrow \mathbb{R}^n$ tal que $\Xi_m e_i = \xi_i$, em que e_i são os vetores canônicos. Desta forma, Ξ_m é injetiva de posto m e a imagem de Ξ_m é o subespaço W_m (Figura 17).

Considere agora $\Xi'_m : \mathbb{R}^n \rightarrow \mathbb{R}^m$ a transformação linear definida pela matriz transposta de Ξ_m ,

$$\Xi'_m = \begin{bmatrix} \xi'_1 \\ \xi'_2 \\ \vdots \\ \xi'_m \end{bmatrix}.$$

Como o posto linha é igual ao posto coluna, tem-se que Ξ'_m é uma aplica-

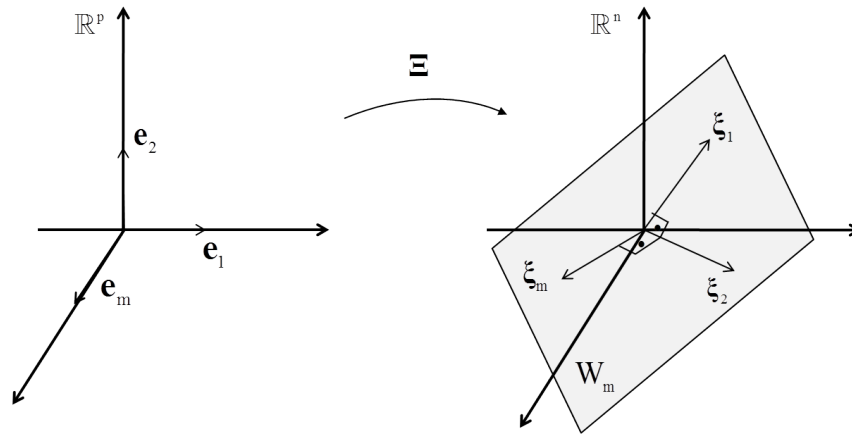


Figura 17 Representação da transformação linear $\Xi_m : \mathbb{R}^m \rightarrow \mathbb{R}^n$ em que $\Xi_m e_i = \xi_i$

ção sobrejetiva. Ξ'_m restrito à $\text{Im}X$ deve ser injetiva. De fato,

$$\begin{bmatrix} \xi'_1 \\ \xi'_2 \\ \vdots \\ \xi'_i \\ \vdots \\ \xi'_m \end{bmatrix} [\xi_i] = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \|\xi_i\|^2 \\ \vdots \\ 0 \end{bmatrix}.$$

Podemos, então, definir o projetor

$$P_{W_m} = \Xi_m (\Xi'_m \Xi_m)^{-1} \Xi'_m,$$

de acordo com a equação (4). De onde seguem as equações normais do estimador

de componentes principais, $\hat{\beta}_{\text{PCR}}$:

$$\mathbf{X}\hat{\beta}_{\text{PCR}} = P_{W_m} \mathbf{Y} = \mathbf{\Xi}_m \left(\mathbf{\Xi}'_m \mathbf{\Xi}_m \right)^{-1} \mathbf{\Xi}'_m \mathbf{Y}.$$

Pode-se obter uma relação entre os estimadores $\hat{\beta}_{\text{PCR}}$ e $\hat{\beta}_{\text{OLS}}$. Seja $\mathbf{\Gamma}_m$ a matriz $p \times m$ definida por $\mathbf{\Gamma}_m = [\gamma_1, \gamma_2, \dots, \gamma_m]$. Tem-se que:

$$\mathbf{X}\mathbf{\Gamma}_m = \mathbf{X} [\gamma_1, \gamma_2, \dots, \gamma_m] = [\mathbf{X}\gamma_1, \mathbf{X}\gamma_2, \dots, \mathbf{X}\gamma_m] = \mathbf{\Xi}_m.$$

Substituindo $\mathbf{\Xi}_m = \mathbf{X}\mathbf{\Gamma}_m$ na definição de $\hat{\beta}_{\text{PCR}}$,

$$\begin{aligned} \mathbf{X}\hat{\beta}_{\text{PCR}} &= \mathbf{X}\mathbf{\Gamma}_m \left(\mathbf{\Xi}'_m \mathbf{\Xi}_m \right)^{-1} \mathbf{\Gamma}'_m \mathbf{X}' \mathbf{Y} \\ &= \mathbf{X}\mathbf{\Gamma}_m \left(\mathbf{\Xi}'_m \mathbf{\Xi}_m \right)^{-1} \mathbf{\Gamma}'_m (\mathbf{X}' \mathbf{X}) (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} \\ &= \mathbf{X}\mathbf{\Gamma}_m \left(\mathbf{\Xi}'_m \mathbf{\Xi}_m \right)^{-1} (\mathbf{X}\mathbf{\Gamma}_m)' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} \\ &= \mathbf{X}\mathbf{\Gamma}_m \left(\mathbf{\Xi}'_m \mathbf{\Xi}_m \right)^{-1} (\mathbf{X}\mathbf{\Gamma}_m)' \mathbf{X} \hat{\beta}_{\text{OLS}}. \end{aligned}$$

Como os ξ_i são ortogonais, $\mathbf{\Xi}'_m \mathbf{\Xi}_m = \mathbf{\Lambda}_m$ é uma matriz diagonal

$$\mathbf{\Lambda}_m = \begin{bmatrix} \|\xi_1\|^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \|\xi_m\|^2 \end{bmatrix}.$$

Assim,

$$\begin{aligned} &\mathbf{X}\mathbf{\Gamma}_m \left(\mathbf{\Xi}'_m \mathbf{\Xi}_m \right)^{-1} (\mathbf{X}\mathbf{\Gamma}_m)' \\ &= \mathbf{X}\mathbf{\Gamma}_m \mathbf{\Lambda}_m^{-1} (\mathbf{X}\mathbf{\Gamma}_m)' \end{aligned}$$

$$\begin{aligned}
&= [\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_m] \begin{bmatrix} \|\boldsymbol{\xi}_1\|^{-2} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \|\boldsymbol{\xi}_m\|^{-2} \end{bmatrix} \begin{bmatrix} \boldsymbol{\xi}'_1 \\ \boldsymbol{\xi}'_2 \\ \vdots \\ \boldsymbol{\xi}'_m \end{bmatrix} \\
&= [\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_m] \begin{bmatrix} \boldsymbol{\xi}'_1 / \|\boldsymbol{\xi}_1\|^2 \\ \boldsymbol{\xi}'_2 / \|\boldsymbol{\xi}_2\|^2 \\ \vdots \\ \boldsymbol{\xi}'_m / \|\boldsymbol{\xi}_m\|^2 \end{bmatrix}.
\end{aligned}$$

Portanto, $\mathbf{X}\hat{\boldsymbol{\beta}}_{\text{PCR}} = [\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_m] \begin{bmatrix} \boldsymbol{\xi}'_1 / \|\boldsymbol{\xi}_1\|^2 \\ \boldsymbol{\xi}'_2 / \|\boldsymbol{\xi}_2\|^2 \\ \vdots \\ \boldsymbol{\xi}'_m / \|\boldsymbol{\xi}_m\|^2 \end{bmatrix} \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{OLS}}.$

Utilizando a decomposição de \mathbf{X} da forma:

$$\mathbf{X} = \boldsymbol{\xi}_1 \boldsymbol{\gamma}'_1 + \boldsymbol{\xi}_2 \boldsymbol{\gamma}'_2 + \dots + \boldsymbol{\xi}_p \boldsymbol{\gamma}'_p,$$

$$\begin{aligned}
&[\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_m] \begin{bmatrix} \boldsymbol{\xi}'_1 / \|\boldsymbol{\xi}_1\|^2 \\ \boldsymbol{\xi}'_2 / \|\boldsymbol{\xi}_2\|^2 \\ \vdots \\ \boldsymbol{\xi}'_m / \|\boldsymbol{\xi}_m\|^2 \end{bmatrix} (\boldsymbol{\xi}_1 \boldsymbol{\gamma}'_1 + \boldsymbol{\xi}_2 \boldsymbol{\gamma}'_2 + \dots + \boldsymbol{\xi}_p \boldsymbol{\gamma}'_p) \hat{\boldsymbol{\beta}}_{\text{OLS}} \\
&= [\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_m] \left\{ \sum_{i=1}^p \begin{bmatrix} \boldsymbol{\xi}'_1 / \|\boldsymbol{\xi}_1\|^2 \\ \boldsymbol{\xi}'_2 / \|\boldsymbol{\xi}_2\|^2 \\ \vdots \\ \boldsymbol{\xi}'_m / \|\boldsymbol{\xi}_m\|^2 \end{bmatrix} \boldsymbol{\xi}_i \boldsymbol{\gamma}'_i \right\} \hat{\boldsymbol{\beta}}_{\text{OLS}}
\end{aligned}$$

$$\begin{aligned}
&= [\xi_1, \xi_2, \dots, \xi_m] \left\{ \sum_{i=1}^p \frac{\xi_i'}{\|\xi_i\|^2} \xi_i \gamma_i' \right\} \hat{\beta}_{OLS} \\
&= [\xi_1, \xi_2, \dots, \xi_m] \Gamma_m' \hat{\beta}_{OLS} \\
&= \mathbf{X} \Gamma_m \Gamma_m' \hat{\beta}_{OLS}.
\end{aligned}$$

Logo, $\mathbf{X} \hat{\beta}_{PCR} = \mathbf{X} (\Gamma_m \Gamma_m') \hat{\beta}_{OLS} \Rightarrow \hat{\beta}_{PCR} = (\Gamma_m \Gamma_m') \hat{\beta}_{OLS}$. Como os vetores γ_i são ortonormais, $\Gamma_m' \Gamma_m = \mathbf{I}_m$ e $\hat{\beta}_{PCR} = \Gamma_m (\Gamma_m' \Gamma_m)^{-1} \Gamma_m' \hat{\beta}_{OLS}$, portanto, $\hat{\beta}_{PCR}$ é obtido como a projeção ortogonal do vetor $\hat{\beta}_{OLS}$ no subespaço gerado por $\{\gamma_1, \gamma_2, \dots, \gamma_m\}$, como observado na Figura 18.

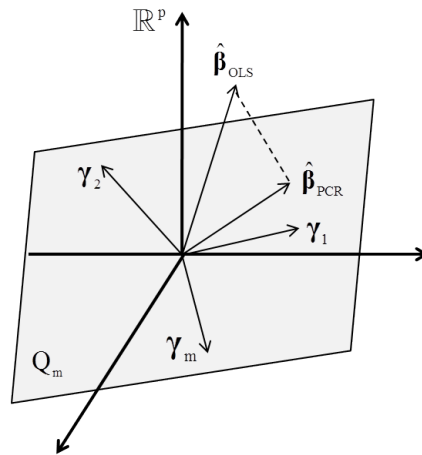


Figura 18 Representação do parâmetro $\hat{\beta}_{PCR}$ como a projeção ortogonal do vetor $\hat{\beta}_{OLS}$ no subespaço gerado por $\{\gamma_1, \gamma_2, \dots, \gamma_m\}$

Em particular, fica demonstrado que $\|\hat{\beta}_{PCR}\| < \|\hat{\beta}_{OLS}\|$, logo, $\hat{\beta}_{PCR}$ é um estimador de encolhimento.

Outra maneira mais geométrica de se ver esta relação é: considerando $\{\gamma_1, \gamma_2, \dots, \gamma_p\}$ como base no domínio e $\{\xi_1, \xi_2, \dots, \xi_p\}$ como base no contra-

domínio. A transformação linear \mathbf{X} se expressa como a matriz identidade

$$\mathbf{X} \left(\sum_{i=1}^p a_i \gamma_i \right) = \sum_{i=1}^p a_i \xi_i.$$

Projetar \mathbf{Y} ortogonalmente na $\text{Im}\mathbf{X}$ e depois projetar ortogonalmente em W_m é o mesmo que projetar ortogonalmente \mathbf{Y} em W_m . A demonstração desse fato segue da aplicação do Teorema de Pitágoras nos vários triângulos retângulos da Figura 19.

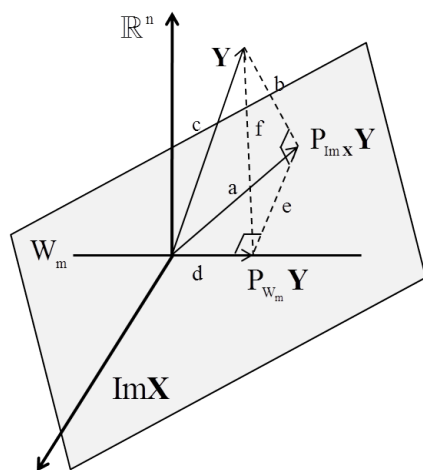


Figura 19 Projeção do vetor de dados \mathbf{Y} no subespaço gerado pelos m primeiros componentes principais

Portanto, se $\mathbf{X}\hat{\beta}_{\text{OLS}} = P_{\text{Im}\mathbf{X}}\mathbf{Y} \Rightarrow \mathbf{X}(\sum a_i \gamma_i) = \sum a_i \xi_i$, projetar $P_{\text{Im}\mathbf{X}}\mathbf{Y}$ em W_m é simplesmente tomar as m primeiras coordenadas, $P_{W_m}\mathbf{Y} = \sum_{i=1}^m a_i \xi_i$ e $\mathbf{X}\hat{\beta}_{\text{PCR}} = P_{W_m}\mathbf{Y}$ implica que $\hat{\beta}_{\text{PCR}} = \sum_{i=1}^m a_i \gamma_i$ e assim, $\hat{\beta}_{\text{PCR}}$ é a projeção de $\hat{\beta}_{\text{OLS}}$ no subespaço gerado por $\{\gamma_1, \gamma_2, \dots, \gamma_m\}$.

A matriz de variâncias e covariâncias de $\hat{\beta}_{\text{PCR}}$ é:

$$\begin{aligned}
\text{cov}(\hat{\beta}_{\text{PCR}}) &= \text{cov}(\mathbf{\Gamma}_m \mathbf{\Gamma}'_m \hat{\beta}_{\text{OLS}}) \\
&= \mathbf{\Gamma}_m \mathbf{\Gamma}'_m \text{cov}(\hat{\beta}_{\text{OLS}}) \mathbf{\Gamma}_m \mathbf{\Gamma}'_m \\
&= \sigma^2 \mathbf{\Gamma}_m \mathbf{\Gamma}'_m (\mathbf{X}'\mathbf{X})^{-1} \mathbf{\Gamma}_m \mathbf{\Gamma}'_m.
\end{aligned}$$

Matricialmente, tem-se:

$$\begin{aligned}
& [\gamma_1, \gamma_2, \dots, \gamma_m] \begin{bmatrix} \gamma'_1 \\ \gamma'_2 \\ \vdots \\ \gamma'_m \end{bmatrix} (\mathbf{X}'\mathbf{X})^{-1} [\gamma_1, \gamma_2, \dots, \gamma_m] \begin{bmatrix} \gamma'_1 \\ \gamma'_2 \\ \vdots \\ \gamma'_m \end{bmatrix} \\
&= [\gamma_1, \gamma_2, \dots, \gamma_m] \begin{bmatrix} \gamma'_1 \\ \gamma'_2 \\ \vdots \\ \gamma'_m \end{bmatrix} \begin{bmatrix} \frac{1}{\lambda_1} \gamma_1, \frac{1}{\lambda_2} \gamma_2, \dots, \frac{1}{\lambda_m} \gamma_m \end{bmatrix} \begin{bmatrix} \gamma'_1 \\ \gamma'_2 \\ \vdots \\ \gamma'_m \end{bmatrix} \\
&= [\gamma_1, \gamma_2, \dots, \gamma_m] \begin{bmatrix} 1/\lambda_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & 1/\lambda_m \end{bmatrix} \begin{bmatrix} \gamma'_1 \\ \gamma'_2 \\ \vdots \\ \gamma'_m \end{bmatrix} \\
&= \begin{bmatrix} \gamma_1 & \gamma_2 & \dots & \gamma_m \\ \lambda_1 & \lambda_2 & \dots & \lambda_m \end{bmatrix} \begin{bmatrix} \gamma'_1 \\ \gamma'_2 \\ \vdots \\ \gamma'_m \end{bmatrix}.
\end{aligned}$$

Essa matriz possui γ_i como autovetores, pois,

$$\begin{aligned} & \left(\begin{array}{c} \left[\frac{\gamma_1}{\lambda_1}, \frac{\gamma_2}{\lambda_2}, \dots, \frac{\gamma_m}{\lambda_m} \right] \begin{bmatrix} \gamma'_1 \\ \gamma'_2 \\ \vdots \\ \gamma'_m \end{bmatrix} \end{array} \right) (\gamma_i) \\ &= \left[\frac{\gamma_1}{\lambda_1}, \frac{\gamma_2}{\lambda_2}, \dots, \frac{\gamma_m}{\lambda_m} \right] \begin{pmatrix} \begin{bmatrix} \gamma'_1 \\ \gamma'_2 \\ \vdots \\ \gamma'_m \end{bmatrix} \\ (\gamma_i) \end{pmatrix} \\ &= \left[\frac{\gamma_1}{\lambda_1}, \frac{\gamma_2}{\lambda_2}, \dots, \frac{\gamma_m}{\lambda_m} \right] \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} = \frac{1}{\lambda_i} \gamma_i \end{aligned}$$

e, portanto, a variância total de $\hat{\beta}_{\text{PCR}}$ é $\sum_{i=1}^m \frac{1}{\lambda_i}$.

Para o caso em que $\text{cov}(\mathbf{Y}) = \Sigma$, o estimador de menor variância é dado pelo estimador de Gauss-Markov (PEREIRA, 2013):

$$\begin{aligned} \hat{\beta}_{\text{GM}} &= (\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1}\mathbf{X}'\Sigma^{-1}\mathbf{Y} \\ \mathbf{X}\hat{\beta}_{\text{GM}} &= \mathbf{X}(\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1}\mathbf{X}'\Sigma^{-1}\mathbf{Y}. \end{aligned}$$

Assim, $\hat{\mathbf{Y}}$ é dado pela projeção de \mathbf{Y} , em que $\mathbf{X}(\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1}\mathbf{X}'\Sigma^{-1}$ é o projetor.

No caso usual, em que $\text{cov}(\mathbf{Y}) = \sigma^2\mathbf{I}$, as projeções nos subespaços W_m que definem $\hat{\beta}_{\text{PCR},m}$ são ortogonais, e quando $m = p$, $\hat{\beta}_{\text{PCR}} = \hat{\beta}_{\text{OLS}}$. Portanto, no caso $\text{cov}(\mathbf{Y}) = \Sigma$, é natural utilizar a mesma projeção. O projetor $\mathbf{X}(\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1}\mathbf{X}'\Sigma^{-1}$ é idempotente, mas não é simétrico, logo, não é um pro-

jetor ortogonal. No entanto, em relação ao produto interno $\langle\langle \mathbf{x}, \mathbf{y} \rangle\rangle = \mathbf{x}'\boldsymbol{\Sigma}^{-1}\mathbf{y}$, esse projetor é simétrico. Assim, o projetor $\mathbf{X}(\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Sigma}^{-1}$ é ortogonal em relação ao produto interno $\langle\langle \mathbf{x}, \mathbf{y} \rangle\rangle = \mathbf{x}'\boldsymbol{\Sigma}^{-1}\mathbf{y}$.

Deseja-se então, maximizar a variância da projeção de \mathbf{Y} na direção de $\mathbf{X}\boldsymbol{\beta}$, isto é, maximizar $\text{var}(\langle\langle \mathbf{X}\boldsymbol{\beta}, \mathbf{Y} \rangle\rangle)$, restrito a $\boldsymbol{\beta}'\boldsymbol{\beta} = 1$:

$$\begin{aligned} \text{var}((\mathbf{X}\boldsymbol{\beta})'\boldsymbol{\Sigma}^{-1}\mathbf{Y}) &= \text{var}(\boldsymbol{\beta}'\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{Y}) \\ &= \boldsymbol{\beta}'\mathbf{X}'\boldsymbol{\Sigma}^{-1}\text{cov}(\mathbf{Y})\boldsymbol{\Sigma}^{-1}\mathbf{X}\boldsymbol{\beta} \\ &= \boldsymbol{\beta}'\mathbf{X}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\Sigma}\boldsymbol{\Sigma}^{-1}\mathbf{X}\boldsymbol{\beta} \\ &= \boldsymbol{\beta}'\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X}\boldsymbol{\beta}. \end{aligned}$$

Sejam $\boldsymbol{\gamma}_i$ autovetores de $\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X}$:

$$\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X}(\boldsymbol{\gamma}_i) = \lambda_i\boldsymbol{\gamma}_i \text{ e } \boldsymbol{\gamma}_i'\boldsymbol{\gamma}_i = 1.$$

Os vetores $\mathbf{X}\boldsymbol{\gamma}_i$ são ortogonais, pois

$$\begin{aligned} \langle\langle \mathbf{X}\boldsymbol{\gamma}_i, \mathbf{X}\boldsymbol{\gamma}_j \rangle\rangle &= (\mathbf{X}\boldsymbol{\gamma}_i)'\boldsymbol{\Sigma}^{-1}\mathbf{X}\boldsymbol{\gamma}_j = \boldsymbol{\gamma}_i'\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X}\boldsymbol{\gamma}_j \\ &= \boldsymbol{\gamma}_i'(\lambda_j\boldsymbol{\gamma}_j) = \lambda_j\boldsymbol{\gamma}_i'\boldsymbol{\gamma}_j = 0. \end{aligned}$$

Definindo o subespaço $W_m = \text{span}\{\mathbf{X}\boldsymbol{\gamma}_1, \mathbf{X}\boldsymbol{\gamma}_2, \mathbf{X}\boldsymbol{\gamma}_3, \dots, \mathbf{X}\boldsymbol{\gamma}_m\}$, tem-se:

$$\mathbf{X}\hat{\boldsymbol{\beta}}_{\text{PCR},m} = P_{W_m}\mathbf{Y},$$

em que P_{W_m} é a projeção ortogonal, em relação ao produto interno $\langle\langle \cdot, \cdot \rangle\rangle$, no subespaço gerado por $\mathbf{X}\boldsymbol{\gamma}_1, \mathbf{X}\boldsymbol{\gamma}_2, \mathbf{X}\boldsymbol{\gamma}_3, \dots, \mathbf{X}\boldsymbol{\gamma}_m$.

Da mesma forma, como feito para o caso em que $\text{cov}(\mathbf{Y}) = \sigma^2\mathbf{I}$, obtém-

se:

$$P_{W_m} = \Xi_m \left(\Xi_m' \Sigma^{-1} \Xi_m \right)^{-1} \Xi_m' \Sigma^{-1}.$$

Como $\Xi_m = \mathbf{X}\Gamma_m$, segue que

$$\begin{aligned} \mathbf{X}\hat{\beta}_{\text{PCR},m} &= \mathbf{X}\Gamma_m \left(\Gamma_m' \mathbf{X}' \Sigma^{-1} \mathbf{X} \Gamma_m \right)^{-1} \Gamma_m' \mathbf{X}' \Sigma^{-1} \mathbf{Y} \\ &= \mathbf{X}\Gamma_m \left(\Gamma_m' \mathbf{X}' \Sigma^{-1} \mathbf{X} \Gamma_m \right)^{-1} \Gamma_m' (\mathbf{X}' \Sigma^{-1} \mathbf{X}) (\mathbf{X}' \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}' \Sigma^{-1} \mathbf{Y} \\ &= \mathbf{X}\Gamma_m \left(\Gamma_m' \mathbf{X}' \Sigma^{-1} \mathbf{X} \Gamma_m \right)^{-1} \Gamma_m' (\mathbf{X}' \Sigma^{-1} \mathbf{X}) \hat{\beta}_{\text{GM}}. \end{aligned}$$

Logo,

$$\hat{\beta}_{\text{PCR},m} = \Gamma_m \left(\Gamma_m' \mathbf{X}' \Sigma^{-1} \mathbf{X} \Gamma_m \right)^{-1} \Gamma_m' (\mathbf{X}' \Sigma^{-1} \mathbf{X}) \hat{\beta}_{\text{GM}}.$$

Observe que, se $m = p$, Γ_p é uma matriz $p \times p$ ortogonal. Assim:

$$\begin{aligned} \hat{\beta}_{\text{PCR},p} &= \Gamma_p \left(\Gamma_p' \mathbf{X}' \Sigma^{-1} \mathbf{X} \Gamma_p \right)^{-1} \Gamma_p' (\mathbf{X}' \Sigma^{-1} \mathbf{X}) \hat{\beta}_{\text{GM}} \\ &= \Gamma_p (\Gamma_p)^{-1} (\mathbf{X}' \Sigma^{-1} \mathbf{X})^{-1} (\Gamma_p')^{-1} \Gamma_p' (\mathbf{X}' \Sigma^{-1} \mathbf{X}) \hat{\beta}_{\text{GM}} \\ &= \hat{\beta}_{\text{GM}}. \end{aligned}$$

4.3 Regressão por Quadrados Mínimos Parciais: uma abordagem geométrica

4.3.1 PLS Populacional

Esta seção tem como referência o artigo Helland (1990).

Seja uma situação em que se quer prever uma variável aleatória y a partir de p variáveis aleatórias x_1, \dots, x_p :

$$E \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix} \quad \text{e} \quad \text{cov} \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix} = \begin{bmatrix} \Sigma & \sigma \\ \sigma' & \sigma_y^2 \end{bmatrix},$$

$$\text{em que, } \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}, \quad \boldsymbol{\sigma} = \text{cov}(\mathbf{x}, y) = \begin{bmatrix} \text{cov}(x_1, y) \\ \vdots \\ \text{cov}(x_p, y) \end{bmatrix} \quad \text{e} \quad \Sigma = \text{cov}(\mathbf{x}).$$

Uma pergunta interessante: Qual é o vetor $\boldsymbol{\eta} \in \mathbb{R}^p$ que define a direção de maior covariância com a variável y ? Isto é: se $\mathbf{Z} = \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix}$, deseja-se o vetor $\boldsymbol{\eta}$ em \mathbb{R}^p , $\|\boldsymbol{\eta}\| = 1$, que, ao se projetar o vetor aleatório \mathbf{Z} na direção de $\boldsymbol{\eta}$, fornece a maior covariância com $\mathbf{e}_{p+1} \cdot \mathbf{Z} = y$. Tem-se, então, o problema de maximização

$$\max_{\|\boldsymbol{\eta}\|=1} \text{cov}(\mathbf{e}_{p+1} \cdot \mathbf{Z}, \boldsymbol{\eta}^* \cdot \mathbf{Z}), \quad \text{em que } \boldsymbol{\eta}^* = \begin{bmatrix} \boldsymbol{\eta} \\ 0 \end{bmatrix}. \quad \text{Utilizando o Teorema 1 sobre}$$

vetores aleatórios, tem-se:

$$\begin{aligned} \text{cov} \left(\mathbf{e}_{p+1} \cdot \mathbf{Z}, \begin{bmatrix} \boldsymbol{\eta} \\ 0 \end{bmatrix} \cdot \mathbf{Z} \right) &= \left\langle \begin{pmatrix} \mathbf{0} \\ 1 \end{pmatrix}, \begin{bmatrix} \boldsymbol{\Sigma} & \boldsymbol{\sigma} \\ \boldsymbol{\sigma}' & \sigma_y^2 \end{bmatrix} \begin{bmatrix} \boldsymbol{\eta} \\ 0 \end{bmatrix} \right\rangle \\ &= \left\langle \begin{bmatrix} \mathbf{0} \\ 1 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}\boldsymbol{\eta} \\ \boldsymbol{\sigma}'\boldsymbol{\eta} \end{bmatrix} \right\rangle = \boldsymbol{\sigma}'\boldsymbol{\eta} = \|\boldsymbol{\eta}\| \|\boldsymbol{\sigma}\| \cos \theta = \|\boldsymbol{\sigma}\| \cos \theta \end{aligned}$$

Para se ter covariância máxima, o ângulo deve ser zero ($\cos 0 = 1$). Logo, $\boldsymbol{\eta}$ deve estar na mesma direção de $\boldsymbol{\sigma}$. Assim, a direção que fornece a covariância máxima é dada pelo próprio vetor de covariâncias $\boldsymbol{\sigma}$. A projeção de \mathbf{Z} em $\boldsymbol{\sigma}$ é $\boldsymbol{\sigma}'\mathbf{Z} = \boldsymbol{\sigma}'\mathbf{x}$.

A ideia é, como uma primeira aproximação, substituir a variável aleatória de interesse y por um múltiplo da variável aleatória preditora $\boldsymbol{\sigma}'\mathbf{x}$. Usando como critério de predição o Erro Quadrático Médio (EQM), devemos minimizar a expressão $E \left[(y - \alpha (\boldsymbol{\sigma}'\mathbf{x}))^2 \right]$:

$$\begin{aligned} &E \left[(y - \alpha (\boldsymbol{\sigma}'\mathbf{x}))^2 \right] \\ &= E \left[y^2 \right] - 2\alpha E \left[y (\boldsymbol{\sigma}'\mathbf{x}) \right] + \alpha^2 E \left[(\boldsymbol{\sigma}'\mathbf{x})^2 \right] \\ &= \sigma_y^2 - 2\alpha \text{cov} (y, \boldsymbol{\sigma}'\mathbf{x}) + \alpha^2 \text{var} (\boldsymbol{\sigma}'\mathbf{x}) \end{aligned}$$

Derivando e igualando a zero, obtém-se:

$$\begin{aligned} -2\text{cov} (y, \boldsymbol{\sigma}'\mathbf{x}) + 2\alpha \text{var} (\boldsymbol{\sigma}'\mathbf{x}) &= 0 \\ \alpha &= \frac{\text{cov}(y, \boldsymbol{\sigma}'\mathbf{x})}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \end{aligned}$$

Note que, a escala da variável y é diferente das escalas das variáveis x_i , e com esta ponderação pela $\text{var} (\boldsymbol{\sigma}'\mathbf{x})$, este problema fica minimizado.

Assim, a variável que prediz y com menor erro quadrático médio é a va-

riável $\frac{\text{cov}(y, \boldsymbol{\sigma}'\mathbf{x})}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \boldsymbol{\sigma}'\mathbf{x}$.

Denominando $q_1 = \text{cov}\left(y, \frac{\boldsymbol{\sigma}'\mathbf{x}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}\right)$, tem-se:

$$\begin{aligned} & \text{cov}\left(y, \frac{\boldsymbol{\sigma}'\mathbf{x}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}\right) \\ &= \text{cov}\left(\mathbf{e}_{p+1} \cdot \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix}, \begin{bmatrix} \frac{\boldsymbol{\sigma}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \\ 0 \end{bmatrix} \cdot \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix}\right) \\ &= \left\langle \mathbf{e}_{p+1}, \begin{bmatrix} \boldsymbol{\Sigma} & \boldsymbol{\sigma} \\ \boldsymbol{\sigma}' & \sigma_y^2 \end{bmatrix} \begin{bmatrix} \frac{\boldsymbol{\sigma}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \\ 0 \end{bmatrix} \right\rangle \\ &= \frac{\boldsymbol{\sigma}'\boldsymbol{\sigma}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} = \frac{\boldsymbol{\sigma}'\boldsymbol{\sigma}}{\boldsymbol{\sigma}'\text{cov}(\mathbf{x})\boldsymbol{\sigma}} = \frac{\boldsymbol{\sigma}'\boldsymbol{\sigma}}{\boldsymbol{\sigma}'\boldsymbol{\Sigma}\boldsymbol{\sigma}}. \end{aligned}$$

Assim, a parte de y que não foi explicada por $\frac{\text{cov}(y, \boldsymbol{\sigma}'\mathbf{x})}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \boldsymbol{\sigma}'\mathbf{x}$, denominada por f_1 , é dada por:

$$f_1 = y - \frac{\text{cov}(y, \boldsymbol{\sigma}'\mathbf{x})}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \boldsymbol{\sigma}'\mathbf{x} = y - \frac{\boldsymbol{\sigma}'\boldsymbol{\sigma}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \boldsymbol{\sigma}'\mathbf{x} = y - \frac{\boldsymbol{\sigma}'\boldsymbol{\sigma}}{\boldsymbol{\sigma}'\boldsymbol{\Sigma}\boldsymbol{\sigma}} \boldsymbol{\sigma}'\mathbf{x}.$$

No sentido de se obter uma sequência de variáveis preditoras, considera-se como primeira variável preditora de y , apenas a variável $\frac{\boldsymbol{\sigma}'\mathbf{x}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}$. Assim, o vetor de correlações $\boldsymbol{\sigma} = \text{cov}(\mathbf{x}, y)$, renomeado por \mathbf{w}_1 , pode ser substituído pelo vetor $\text{cov}\left(\mathbf{x}, \frac{\boldsymbol{\sigma}'\mathbf{x}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}\right)$. Denomina-se este vetor por:

$$\mathbf{p}_1 = \frac{\text{cov}(\mathbf{x}, \boldsymbol{\sigma}'\mathbf{x})}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} = \begin{bmatrix} \frac{\text{cov}(x_1, \boldsymbol{\sigma}'\mathbf{x})}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \\ \vdots \\ \frac{\text{cov}(x_p, \boldsymbol{\sigma}'\mathbf{x})}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \end{bmatrix}.$$

A i -ésima coordenada do vetor \mathbf{p}_1 é dada por:

$$\begin{aligned}
 x_i &= \mathbf{e}_i \cdot \mathbf{Z} \\
 \boldsymbol{\sigma}'\mathbf{x} &= (\boldsymbol{\sigma}, 0) \cdot \mathbf{Z} \\
 \text{cov}(x_i, \boldsymbol{\sigma}'\mathbf{x}) &= \text{cov}(\mathbf{e}_i \cdot \mathbf{Z}, \boldsymbol{\sigma} \cdot \mathbf{Z}) \\
 &= \left\langle \mathbf{e}_i, \begin{bmatrix} \boldsymbol{\Sigma} & \boldsymbol{\sigma} \\ \boldsymbol{\sigma}' & \sigma_y^2 \end{bmatrix} \begin{bmatrix} \boldsymbol{\sigma} \\ 0 \end{bmatrix} \right\rangle = \left\langle \mathbf{e}_i, \begin{bmatrix} \boldsymbol{\Sigma}\boldsymbol{\sigma} \\ \boldsymbol{\sigma}'\boldsymbol{\sigma} \end{bmatrix} \right\rangle \\
 &= \sum_{j=1}^p \text{cov}(x_i, x_j) \sigma_j = \sum_{j=1}^p \text{cov}(x_i, x_j) \text{cov}(x_j, y) \\
 &\Rightarrow \text{cov}\left(x_i, \frac{\boldsymbol{\sigma}'\mathbf{x}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}\right) = \frac{1}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \sum_{j=1}^p \text{cov}(x_i, x_j) \text{cov}(x_j, y).
 \end{aligned}$$

Ou ainda,

$$\begin{aligned}
 \text{cov}\left(x_i, \frac{\boldsymbol{\sigma}'\mathbf{x}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}\right) &= \text{cov}\left(x_i, \sum_{j=1}^p x_j \frac{\text{cov}(x_j, y)}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}\right) \\
 &= \sum_{j=1}^p \text{cov}\left(x_i, x_j \frac{\text{cov}(x_j, y)}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}\right) = \sum_{j=1}^p \frac{\text{cov}(x_j, y)}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \text{cov}(x_i, x_j) \\
 &= \sum_{j=1}^p \frac{\text{cov}(x_i, x_j)}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \sigma_j = \frac{1}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} (\boldsymbol{\Sigma}\boldsymbol{\sigma})_i.
 \end{aligned}$$

$$\text{Logo, } \mathbf{p}_1 = \frac{\boldsymbol{\Sigma}\boldsymbol{\sigma}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}.$$

A componente de \mathbf{x} perpendicular ao vetor $\boldsymbol{\sigma}$ não foi utilizada na predição de y . Deve-se então, encontrar α tal que $\mathbf{x} - \alpha\mathbf{p}_1$ seja perpendicular à $\boldsymbol{\sigma}$:

$$\begin{aligned}
 \langle \boldsymbol{\sigma}, \mathbf{x} - \alpha\mathbf{p}_1 \rangle &= \langle \boldsymbol{\sigma}, \mathbf{x} \rangle - \alpha \langle \boldsymbol{\sigma}, \mathbf{p}_1 \rangle \\
 &\Rightarrow \langle \boldsymbol{\sigma}, \mathbf{x} \rangle - \alpha \left\langle \boldsymbol{\sigma}, \frac{\boldsymbol{\Sigma}\boldsymbol{\sigma}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} \right\rangle = 0 \\
 &\Rightarrow \langle \boldsymbol{\sigma}, \mathbf{x} \rangle - \alpha \frac{\boldsymbol{\sigma}'\boldsymbol{\Sigma}\boldsymbol{\sigma}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})} = 0 \\
 &\Rightarrow \boldsymbol{\sigma}'\mathbf{x} - \alpha = 0 \\
 &\Rightarrow \alpha = \boldsymbol{\sigma}'\mathbf{x}.
 \end{aligned}$$

Desta forma, denominando $t_1 = \boldsymbol{\sigma}'\mathbf{x}$, tem-se o vetor $\mathbf{e}_1 = \mathbf{x} - t_1\mathbf{p}_1$, $\mathbf{e}_1 \perp \boldsymbol{\sigma}$. Note que, a preditora de y pode ser reescrita como $\frac{t_1}{\text{var}(t_1)}$ e

$$\mathbf{e}_1 = \mathbf{x} - (\boldsymbol{\sigma}'\mathbf{x}) \frac{\boldsymbol{\Sigma}\boldsymbol{\sigma}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}.$$

O processo agora é repetido, substituindo-se y por sua parte não explicada f_1 e \mathbf{x} por sua componente ainda não utilizada \mathbf{e}_1 . A direção de maior covariância entre f_1 e os vetores aleatórios no subespaço perpendicular a $\boldsymbol{\sigma}$, é dada pelo vetor \mathbf{w}_2 :

$$\begin{aligned} \mathbf{w}_2 &= \text{cov}(\mathbf{e}_1, f_1) = \text{cov}\left(\mathbf{e}_1, y - \frac{\boldsymbol{\sigma}'\boldsymbol{\sigma}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}\boldsymbol{\sigma}'\mathbf{x}\right) \\ &= \text{cov}(\mathbf{e}_1, y) - \frac{\boldsymbol{\sigma}'\boldsymbol{\sigma}}{\text{var}(\boldsymbol{\sigma}'\mathbf{x})}\text{cov}(\mathbf{e}_1, \boldsymbol{\sigma}'\mathbf{x}) \stackrel{\mathbf{e}_1 \perp \boldsymbol{\sigma}}{=} \text{cov}(\mathbf{e}_1, y) \\ &= \text{cov}(\mathbf{x} - \mathbf{p}_1 t_1, y) = \text{cov}(\mathbf{x}, y) - \text{cov}(\mathbf{p}_1 t_1, y) \\ &= \boldsymbol{\sigma} - \mathbf{p}_1 \text{cov}(\boldsymbol{\sigma}'\mathbf{x}, y) = \boldsymbol{\sigma} - \frac{\boldsymbol{\Sigma}\boldsymbol{\sigma}}{\boldsymbol{\sigma}'\boldsymbol{\Sigma}\boldsymbol{\sigma}}\boldsymbol{\sigma}'\boldsymbol{\sigma}. \end{aligned}$$

Definindo $t_2 = \mathbf{e}_1'\mathbf{w}_2$, a variável preditora de f_1 será $\frac{t_2}{\text{var}(t_2)}$. Calculando-se a covariância entre f_1 e seu preditor $\frac{t_2}{\text{var}(t_2)}$, temos a grandeza q_2 , dada por

$$q_2 = \frac{\text{cov}(f_1, t_2)}{\text{var}(t_2)}.$$

Desta forma, a parte não explicada de f_1 será dada por $f_2 = f_1 - q_2 t_2$.

Define-se agora, o vetor dado pela covariância entre o vetor aleatório \mathbf{e}_1 e a variável preditora de f_1 , $\frac{t_2}{\text{var}(t_2)}$, denominado \mathbf{p}_2 :

$$\mathbf{p}_2 = \frac{\text{cov}(\mathbf{e}_1, t_2)}{\text{var}(t_2)}.$$

Novamente, a componente do vetor \mathbf{e}_1 que não foi utilizada para explicar

f_1 será dada pelo vetor $\mathbf{e}_2 = \mathbf{e}_1 - \mathbf{p}_2 t_2$.

Resumindo:

Para $a = 1$:

$$\mathbf{w}_1 = \text{cov}(\mathbf{x}, y) = \boldsymbol{\sigma}$$

$$t_1 = \mathbf{x}' \mathbf{w}_1 = \boldsymbol{\sigma}' \mathbf{x}$$

$$\mathbf{p}_1 = \text{cov}\left(\mathbf{x}, \frac{t_1}{\text{var}(t_1)}\right) = \text{cov}\left(\mathbf{x}, \frac{\boldsymbol{\sigma}' \mathbf{x}}{\text{var}(\boldsymbol{\sigma}' \mathbf{x})}\right) = \frac{\boldsymbol{\Sigma} \boldsymbol{\sigma}}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}}$$

$$q_1 = \text{cov}\left(y, \frac{t_1}{\text{var}(t_1)}\right) = \text{cov}\left(y, \frac{\boldsymbol{\sigma}' \mathbf{x}}{\text{var}(\boldsymbol{\sigma}' \mathbf{x})}\right) = \frac{\boldsymbol{\sigma}' \boldsymbol{\sigma}}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}}$$

$$\mathbf{e}_1 = \mathbf{x} - \mathbf{p}_1 t_1 = \mathbf{x} - \frac{\boldsymbol{\Sigma} \boldsymbol{\sigma}}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}} \boldsymbol{\sigma}' \mathbf{x}$$

$$f_1 = y - q_1 t_1 = y - \frac{\boldsymbol{\sigma}' \boldsymbol{\sigma}}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}} \boldsymbol{\sigma}' \mathbf{x}$$

Para $a = 2$:

$$t_2 = \mathbf{e}_1' \mathbf{w}_2 = \mathbf{x}' \left(\boldsymbol{\sigma} - \frac{\boldsymbol{\sigma}' \boldsymbol{\sigma}}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}} \boldsymbol{\Sigma} \boldsymbol{\sigma} \right) = \mathbf{x}' \boldsymbol{\sigma} - \frac{\boldsymbol{\sigma}' \boldsymbol{\sigma}}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}} \mathbf{x}' \boldsymbol{\Sigma} \boldsymbol{\sigma}$$

$$\text{var}(t_2) = \left(\boldsymbol{\sigma} - \frac{\|\boldsymbol{\sigma}\|^2}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}} \boldsymbol{\Sigma} \boldsymbol{\sigma} \right)' \boldsymbol{\Sigma} \left(\boldsymbol{\sigma} - \frac{\|\boldsymbol{\sigma}\|^2}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}} \boldsymbol{\Sigma} \boldsymbol{\sigma} \right)$$

$$\mathbf{w}_2 = \text{cov}(\mathbf{e}_1, f_1) = \boldsymbol{\sigma} - \frac{\boldsymbol{\Sigma} \boldsymbol{\sigma}}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}} \boldsymbol{\sigma}' \boldsymbol{\sigma}$$

$$\mathbf{p}_2 = \frac{\text{cov}(\mathbf{e}_1, t_2)}{\text{var}(t_2)} = \frac{\text{cov}(\mathbf{x}, t_2)}{\text{var}(t_2)}$$

$$= \text{cov}\left(\mathbf{x}, \mathbf{x}' \boldsymbol{\sigma} - \frac{\boldsymbol{\sigma}' \boldsymbol{\sigma}}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}} \mathbf{x}' \boldsymbol{\Sigma} \boldsymbol{\sigma}\right) = \boldsymbol{\Sigma} \boldsymbol{\sigma} - \frac{\boldsymbol{\sigma}' \boldsymbol{\sigma}}{\boldsymbol{\sigma}' \boldsymbol{\Sigma} \boldsymbol{\sigma}} \boldsymbol{\Sigma}^2 \boldsymbol{\sigma}$$

$$q_2 = \frac{\text{cov}(f_1, t_2)}{\text{var}(t_2)}$$

$$\mathbf{e}_2 = \mathbf{e}_1 - \mathbf{p}_2 t_2$$

$$f_2 = f_1 - q_2 t_2$$

Um algoritmo para a predição da variável aleatória de interesse y , em termos das covariáveis fica, então, definido.

Algoritmo para o PLS Populacional (HELLAND, 1990):

1. Defina os valores iniciais para \mathbf{x} residual (\mathbf{e}_a) e y residual (f_a) :

$$\mathbf{e}_0 = \mathbf{x} - \boldsymbol{\mu}_x$$

$$f_0 = y - \mu_y$$

Para $a = 1, 2, \dots$, execute os passos (2) – (4) abaixo:

2. Introduza *scores* t como combinações lineares do \mathbf{x} residual a partir do último passo; para fazer os *scores* mais relacionados com y , use covariâncias com y residual como coeficiente/pesos (\mathbf{w}) :

$$\mathbf{w}_a = \text{cov}(\mathbf{e}_{a-1}, f_{a-1})$$

$$t_a = \mathbf{e}_{a-1}' \mathbf{w}_a$$

3. Determine \mathbf{x} *loadings* (\mathbf{p}_a) e y *loading* (q_a) por quadrados mínimos:

$$\mathbf{p}_a = \frac{\text{cov}(\mathbf{e}_{a-1}, t_a)}{\text{var}(t_a)}$$

$$q_a = \frac{\text{cov}(f_{a-1}, t_a)}{\text{var}(t_a)}$$

4. Encontre novos resíduos:

$$\mathbf{e}_a = \mathbf{e}_{a-1} - \mathbf{p}_a t_a$$

$$f_a = f_{a-1} - q_a t_a$$

Deve-se observar que os vetores $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_A$ e $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_A$ são vetores populacionais. Além disso, q_a são grandezas populacionais, isto é, fixada uma determinada população, os números q_a ficam determinados. Portanto, as grandezas aleatórias são os vetores \mathbf{e}_a , os pesos t_a e as variáveis f_a . Duas iterações do algoritmo estão descritas na Figura 20.

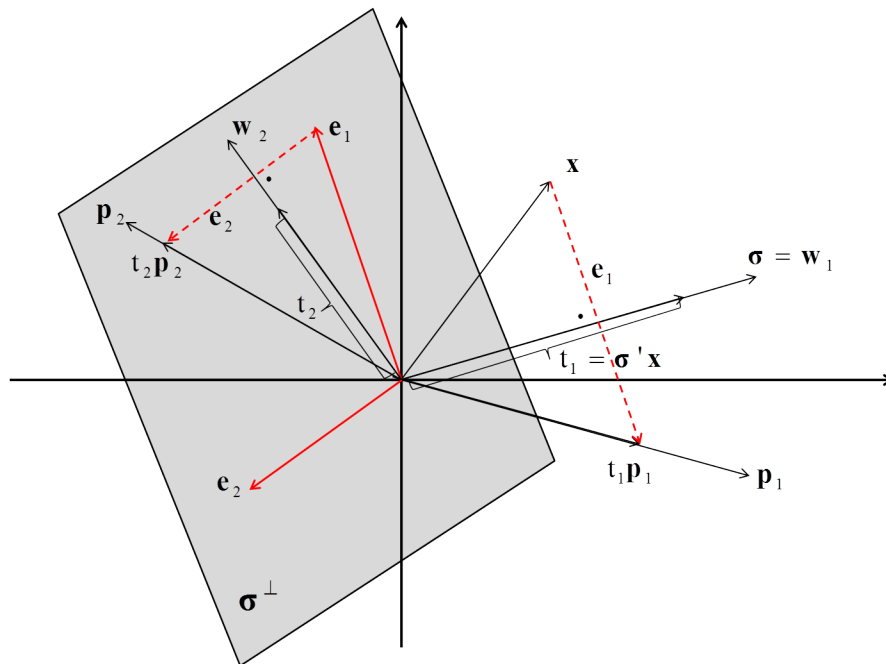


Figura 20 Representação geométrica do algoritmo PLS Populacional com duas iterações

Da construção geométrica do algoritmo, seguem os seguintes fatos fundamentais:

1. Os vetores $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_A$ são ortogonais.
2. As variáveis aleatórias t_1, t_2, \dots, t_A são não correlacionadas. Para o caso em que $\Sigma = \mathbf{I}$, este fato segue diretamente da construção geométrica, pois

as variáveis t_1, t_2, \dots, t_A são obtidas por projeções nos vetores \mathbf{w}_i , que são ortogonais.

3. O vetor aleatório \mathbf{e}_a é não correlacionado à variável t_a , pois t_a é obtida por uma projeção em \mathbf{w}_a e o vetor \mathbf{e}_a é perpendicular ao vetor \mathbf{w}_a . De fato, temos mais do que isso: como \mathbf{e}_a é perpendicular a \mathbf{w}_a , por construção, e \mathbf{w}_a é perpendicular a $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{a-1}$, temos que \mathbf{e}_a é perpendicular a todos eles, $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_a$, implicando então, que \mathbf{e}_a é não correlacionado às variáveis t_1, t_2, \dots, t_a . Da mesma forma, este resultado segue da construção geométrica se $\Sigma = \mathbf{I}$.

Estas propriedades podem ser demonstradas algebricamente (HOSKUL-DSSON, 1988).

As propriedades 2 e 3 permitem as seguintes simplificações no algoritmo:

•

$$\begin{aligned}
 \mathbf{w}_a &= \text{cov}(\mathbf{e}_{a-1}, f_{a-1}) = \text{cov}(\mathbf{e}_{a-1}, f_{a-2} - q_{a-1}t_{a-1}) \\
 &= \text{cov}(\mathbf{e}_{a-1}, f_{a-2}) - q_{a-1}\text{cov}(\mathbf{e}_{a-1}, t_{a-1}) \\
 &\stackrel{\text{P.3}}{=} \text{cov}(\mathbf{e}_{a-1}, f_{a-2}) = \text{cov}(\mathbf{e}_{a-1}, f_{a-3} - q_{a-2}t_{a-2}) \\
 &= \text{cov}(\mathbf{e}_{a-1}, f_{a-3}) - q_{a-2}\text{cov}(\mathbf{e}_{a-1}, t_{a-2}) \\
 &\stackrel{\text{P.3}}{=} \text{cov}(\mathbf{e}_{a-1}, f_{a-3}) = \dots = \text{cov}(\mathbf{e}_{a-1}, y)
 \end{aligned}$$

Logo, $\mathbf{w}_a = \text{cov}(\mathbf{e}_{a-1}, f_{a-1}) = \text{cov}(\mathbf{e}_{a-1}, y)$.

•

$$\begin{aligned}
\text{cov}(\mathbf{e}_{a-1}, t_a) &= \text{cov}(\mathbf{e}_{a-2} - \mathbf{p}_{a-1}t_{a-1}, t_a) \\
&= \text{cov}(\mathbf{e}_{a-2}, t_a) - \text{cov}(t_{a-1}, t_a) \mathbf{p}_{a-1} \\
&\stackrel{\text{P.2}}{=} \text{cov}(\mathbf{e}_{a-2}, t_a) = \text{cov}(\mathbf{e}_{a-3} - \mathbf{p}_{a-2}t_{a-2}, t_a) \\
&= \text{cov}(\mathbf{e}_{a-3}, t_a) - \text{cov}(t_{a-2}, t_a) \mathbf{p}_{a-2} \\
&\stackrel{\text{P.2}}{=} \text{cov}(\mathbf{e}_{a-3}, t_a) = \dots = \text{cov}(\mathbf{x}, t_a)
\end{aligned}$$

$$\text{Logo, } \mathbf{p}_a = \frac{\text{cov}(\mathbf{e}_{a-1}, t_a)}{\text{var}(t_a)} = \frac{\text{cov}(\mathbf{x}, t_a)}{\text{var}(t_a)}.$$

•

$$\begin{aligned}
\text{cov}(f_{a-1}, t_a) &= \text{cov}(f_{a-2} - q_{a-1}t_{a-1}, t_a) \\
&= \text{cov}(f_{a-2}, t_a) - q_{a-1} \text{cov}(t_{a-1}, t_a) \\
&\stackrel{\text{P.2}}{=} \text{cov}(f_{a-2}, t_a) = \text{cov}(f_{a-3} - q_{a-2}t_{a-2}, t_a) \\
&= \text{cov}(f_{a-3}, t_a) - q_{a-2} \text{cov}(t_{a-2}, t_a) \\
&\stackrel{\text{P.2}}{=} \text{cov}(f_{a-3}, t_a) = \dots = \text{cov}(y, t_a)
\end{aligned}$$

$$\text{Logo, } q_a = \frac{\text{cov}(f_{a-1}, t_a)}{\text{var}(t_a)} = \frac{\text{cov}(y, t_a)}{\text{var}(t_a)}.$$

Portanto, \mathbf{e}_{a-1} pode ser substituído por \mathbf{x} na definição de \mathbf{p}_a , e f_{a-1} pode ser substituído por y nas definições de \mathbf{w}_a e q_a , nos passos 2 e 3 do algoritmo.

Assim, tem-se que o algoritmo fica da forma:

$$\mathbf{e}_0 = \mathbf{x} - \boldsymbol{\mu}_x$$

$$f_0 = y - \mu_y$$

$$t_a = \mathbf{e}_{a-1}' \mathbf{w}_a$$

$$\mathbf{w}_a = \text{cov}(\mathbf{e}_{a-1}, y)$$

$$\mathbf{p}_a = \frac{\text{cov}(\mathbf{x}, t_a)}{\text{var}(t_a)}$$

$$q_a = \frac{\text{cov}(y, t_a)}{\text{var}(t_a)}$$

$$\mathbf{e}_a = \mathbf{e}_{a-1} - \mathbf{p}_a t_a$$

$$f_a = f_{a-1} - q_a t_a.$$

Apesar do algoritmo nessa versão ser mais simples, como o objetivo é a motivação para se obter o algoritmo para o caso amostral, a primeira versão é a mais útil.

No passo A do algoritmo, obtemos a representação:

$$\begin{aligned} & \mathbf{p}_1 t_1 + \mathbf{p}_2 t_2 + \mathbf{p}_3 t_3 + \cdots + \mathbf{p}_A t_A \\ &= (\mathbf{e}_0 - \mathbf{e}_1) + (\mathbf{e}_1 - \mathbf{e}_2) + (\mathbf{e}_2 - \mathbf{e}_3) + \cdots + (\mathbf{e}_{A-1} - \mathbf{e}_A) \\ &= \mathbf{e}_0 - \mathbf{e}_A = (\mathbf{x} - \boldsymbol{\mu}_x) - \mathbf{e}_A, \end{aligned}$$

logo, $\mathbf{x} = \boldsymbol{\mu}_x + \mathbf{p}_1 t_1 + \cdots + \mathbf{p}_A t_A + \mathbf{e}_A$.

$$\begin{aligned} & q_1 t_1 + q_2 t_2 + q_3 t_3 + \cdots + q_A t_A \\ &= (f_0 - f_1) + (f_1 - f_2) + (f_2 - f_3) + \cdots + (f_{A-1} - f_A) \\ &= f_0 - f_A = (y - \mu_y) - f_A, \end{aligned}$$

logo, $y = \mu_y + q_1 t_1 + \cdots + q_A t_A + f_A$.

Temos então, um preditor para a variável resposta y , dado por

$$\hat{y}_{A, \text{PLS}} = \mu_y + q_1 t_1 + \cdots + q_A t_A.$$

Este preditor é, de fato, um preditor linear, pois utilizando a segunda versão do

algoritmo, tem-se:

$$\hat{y}_{A,PLS} = \mu_y + \beta'_{A,PLS} (\mathbf{x} - \mu_x)$$

em que (HELLAND, 1990),

$$\beta_{A,PLS} = \mathbf{W}_A (\mathbf{W}_A' \Sigma \mathbf{W}_A)^{-1} \mathbf{W}_A' \sigma.$$

Explicitando essas matrizes:

$$\mathbf{W}_A = [\mathbf{w}_1, \dots, \mathbf{w}_A]$$

$$\Sigma (\mathbf{W}_A) = \Sigma [\mathbf{w}_1, \dots, \mathbf{w}_A] = [\Sigma \mathbf{w}_1, \dots, \Sigma \mathbf{w}_A]$$

$$\mathbf{W}_A' \Sigma \mathbf{W}_A = \begin{bmatrix} \mathbf{w}_1' \\ \vdots \\ \mathbf{w}_A' \end{bmatrix} [\Sigma \mathbf{w}_1, \dots, \Sigma \mathbf{w}_A] = \begin{bmatrix} \mathbf{w}_1' \Sigma \mathbf{w}_1 & \cdots & \mathbf{w}_1' \Sigma \mathbf{w}_A \\ \vdots & \ddots & \vdots \\ \mathbf{w}_A' \Sigma \mathbf{w}_1 & \cdots & \mathbf{w}_A' \Sigma \mathbf{w}_A \end{bmatrix}.$$

Casos particulares:

1º Caso: Se $\mathbf{w}_i' \Sigma \mathbf{w}_j = 0, \forall i \neq j$, temos uma matriz diagonal e poderemos calcular a inversa:

$$(\mathbf{W}_A' \Sigma \mathbf{W}_A)^{-1} = \begin{bmatrix} \frac{1}{\mathbf{w}_1' \Sigma \mathbf{w}_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{\mathbf{w}_2' \Sigma \mathbf{w}_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{\mathbf{w}_A' \Sigma \mathbf{w}_A} \end{bmatrix}$$

$$\mathbf{W}_A' \sigma = \begin{bmatrix} \mathbf{w}_1' \\ \vdots \\ \mathbf{w}_A' \end{bmatrix} \sigma = \begin{bmatrix} \mathbf{w}_1' \sigma \\ \vdots \\ \mathbf{w}_A' \sigma \end{bmatrix}$$

$$\begin{aligned}
(\mathbf{W}_A' \Sigma \mathbf{W}_A)^{-1} \mathbf{W}_A' \boldsymbol{\sigma} &= \begin{bmatrix} \frac{\mathbf{w}_1' \boldsymbol{\sigma}}{\mathbf{w}_1' \Sigma \mathbf{w}_1} \\ \vdots \\ \frac{\mathbf{w}_A' \boldsymbol{\sigma}}{\mathbf{w}_A' \Sigma \mathbf{w}_A} \end{bmatrix} \\
\mathbf{W}_A (\mathbf{W}_A' \Sigma \mathbf{W}_A)^{-1} \mathbf{W}_A' \boldsymbol{\sigma} &= [\mathbf{w}_1, \dots, \mathbf{w}_A] \begin{bmatrix} \frac{\mathbf{w}_1' \boldsymbol{\sigma}}{\mathbf{w}_1' \Sigma \mathbf{w}_1} \\ \vdots \\ \frac{\mathbf{w}_A' \boldsymbol{\sigma}}{\mathbf{w}_A' \Sigma \mathbf{w}_A} \end{bmatrix} \\
&= \left[\left(\frac{\mathbf{w}_1' \boldsymbol{\sigma}}{\mathbf{w}_1' \Sigma \mathbf{w}_1} \right) \mathbf{w}_1 + \dots + \left(\frac{\mathbf{w}_A' \boldsymbol{\sigma}}{\mathbf{w}_A' \Sigma \mathbf{w}_A} \right) \mathbf{w}_A \right].
\end{aligned}$$

2º Caso: Se o número de passos A do algoritmo é igual a n , $A = n$, a matriz \mathbf{W}_A é quadrada e inversível, assim,

$$\begin{aligned}
\beta_{A,PLS} &= \mathbf{W}_A (\mathbf{W}_A' \Sigma \mathbf{W}_A)^{-1} \mathbf{W}_A' \boldsymbol{\sigma}, \\
\beta_{n,PLS} &= \mathbf{W}_A \mathbf{W}_A^{-1} \Sigma^{-1} (\mathbf{W}_A')^{-1} \mathbf{W}_A' \boldsymbol{\sigma} = \Sigma^{-1} \boldsymbol{\sigma},
\end{aligned}$$

que é o mesmo coeficiente do melhor preditor linear, no caso em que $\begin{bmatrix} \mathbf{x} \\ y \end{bmatrix}$ tem distribuição Normal.

3º Caso: O preditor para apenas 1 passo, fica da forma:

$$\hat{\mathbf{Y}}_{A,PLS} = \mu_y + q_1 t_1 = \mu_y + \frac{\|\boldsymbol{\sigma}\|^2}{\boldsymbol{\sigma}' \Sigma \boldsymbol{\sigma}} \mathbf{x}' \boldsymbol{\sigma} = \mu_y + \frac{\|\boldsymbol{\sigma}\|^2}{\boldsymbol{\sigma}' \Sigma \boldsymbol{\sigma}} \boldsymbol{\sigma}' (\mathbf{x} - \boldsymbol{\mu}_x).$$

4.3.2 PLS Amostral

Esta seção tem como referência o artigo Hoskuldsson (1988).

Sejam duas matrizes $\mathbf{X}_{n \times p}$ e $\mathbf{Y}_{n \times k}$, centradas na média, em que p é o número de covariáveis, k o número de variáveis respostas e n o número de observações (repetições):

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} y_{11} & y_{12} & \cdots & y_{1k} \\ y_{21} & y_{22} & \cdots & y_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \cdots & y_{nk} \end{bmatrix}.$$

A matriz de covariâncias amostrais, a menos de uma constante, é dada por $\mathbf{X}'\mathbf{Y}$:

$$\begin{aligned} \mathbf{X}'\mathbf{Y} &= \begin{bmatrix} x_{11} & x_{21} & \cdots & x_{n1} \\ x_{12} & x_{22} & \cdots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1p} & x_{2p} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} y_{11} & y_{12} & \cdots & y_{1k} \\ y_{21} & y_{22} & \cdots & y_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \cdots & y_{nk} \end{bmatrix} \\ &= \begin{bmatrix} \hat{\text{cov}}(\mathbf{x}_1, \mathbf{y}_1) & \hat{\text{cov}}(\mathbf{x}_1, \mathbf{y}_2) & \cdots & \hat{\text{cov}}(\mathbf{x}_1, \mathbf{y}_k) \\ \hat{\text{cov}}(\mathbf{x}_2, \mathbf{y}_1) & \hat{\text{cov}}(\mathbf{x}_2, \mathbf{y}_2) & \cdots & \hat{\text{cov}}(\mathbf{x}_2, \mathbf{y}_k) \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\text{cov}}(\mathbf{x}_p, \mathbf{y}_1) & \hat{\text{cov}}(\mathbf{x}_p, \mathbf{y}_2) & \cdots & \hat{\text{cov}}(\mathbf{x}_p, \mathbf{y}_k) \end{bmatrix}. \end{aligned}$$

Observe que as colunas da matriz $\mathbf{X}'\mathbf{Y}$, a menos de uma constante, estimam os vetores $\boldsymbol{\sigma}_j = \text{cov}(\mathbf{X}, \mathbf{Y}_j)$ com $j = 1, \dots, k$.

4.3.2.1 Componentes com covariância máxima

Considere dois componentes f e g definidos pelas colunas das matrizes \mathbf{X} e \mathbf{Y} :

$$\mathbf{f} = \mathbf{X}\mathbf{d}, \quad \|\mathbf{d}\| = 1$$

$$\mathbf{g} = \mathbf{Y}\mathbf{e}, \quad \|\mathbf{e}\| = 1.$$

Note que,

$$\mathbf{f} = \mathbf{X}\mathbf{d} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_p \end{bmatrix} = d_1 \begin{bmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{n1} \end{bmatrix} + \cdots + d_p \begin{bmatrix} x_{1p} \\ x_{2p} \\ \vdots \\ x_{np} \end{bmatrix},$$

logo, $\mathbf{f} = \mathbf{X}\mathbf{d}$ é uma combinação linear das colunas de \mathbf{X} , isto é, é um vetor definido como uma média ponderada das realizações de todas as covariáveis. Da mesma forma, $\mathbf{g} = \mathbf{Y}\mathbf{e}$ é uma combinação linear das colunas de \mathbf{Y} .

A covariância amostral entre os dois componentes é dada por:

$$\text{cov}(\mathbf{f}, \mathbf{g}) = \frac{\mathbf{f}'\mathbf{g}}{n}.$$

Uma boa escolha para \mathbf{f} e \mathbf{g} será aquela que fornecer covariância máxima. Aqui, a covariância deve ser máxima em valor absoluto, visto que a covariância pode ser positiva ou negativa.

Para obter os vetores que realizam essa maximização da covariância, é necessário fazer a decomposição em valores singulares (Apêndice B) da matriz de covariâncias amostrais $\mathbf{X}'\mathbf{Y}$. Para tal, considere $\mathbf{X}'\mathbf{Y}$ como uma transformação linear de \mathbb{R}^k em \mathbb{R}^p , $k < p$. Existe uma base $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ em \mathbb{R}^k , e uma base

$\{z_1, z_2, \dots, z_p\}$ em \mathbb{R}^p , ambas ortonormais, tais que $\mathbf{X}'\mathbf{Y}(v_i) = a_i z_i$, sendo a ordem dos vetores nas bases definida por $a_1^2 \geq a_2^2 \geq \dots \geq a_k^2$, conforme mostra a Figura 21. Desta forma, tem-se a igualdade de matrizes $\mathbf{X}'\mathbf{Y} = \sum_{i=1}^k a_i z_i v_i'$, pois

$$\left(\sum_{i=1}^k a_i z_i v_i' \right) (v_j) = \sum_{i=1}^k (a_i z_i v_i') (v_j) = \sum_{i=1}^k (a_i z_i) (v_i' v_j) = a_j z_j = \mathbf{X}'\mathbf{Y}(v_j).$$

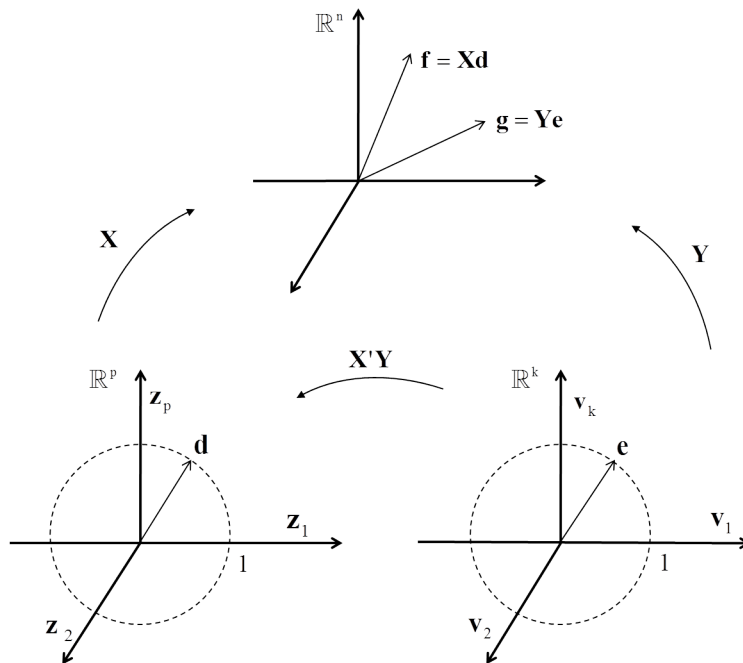


Figura 21 Representação de $\mathbf{X}'\mathbf{Y}$ como uma transformação linear de \mathbb{R}^k em \mathbb{R}^p

O problema da maximização, $\max(\mathbf{f}'\mathbf{g})$, $\|\mathbf{d}\| = \|\mathbf{e}\| = 1$, pode ser colocado da forma:

$$\begin{aligned}
\max(\mathbf{f}'\mathbf{g}) &= \max((\mathbf{X}\mathbf{d})'(\mathbf{Y}\mathbf{e})) = \max(\mathbf{d}'\mathbf{X}'\mathbf{Y}\mathbf{e}) \\
&= \max(\mathbf{d}'(\mathbf{X}'\mathbf{Y})\mathbf{e}) = \max\left(\mathbf{d}'\left(\sum_{i=1}^k a_i \mathbf{z}_i \mathbf{v}_i'\right)\mathbf{e}\right) \\
&= \max\left(\sum_{i=1}^k a_i (\mathbf{d}'(\mathbf{z}_i \mathbf{v}_i')\mathbf{e})\right) = \max\left(\sum_{i=1}^k a_i (\mathbf{d}'\mathbf{z}_i)(\mathbf{v}_i'\mathbf{e})\right).
\end{aligned}$$

Assim, os vetores \mathbf{d} e \mathbf{e} que realizam essa maximização são os vetores $\mathbf{d} = \mathbf{z}_1$ e $\mathbf{e} = \mathbf{v}_1$, pois,

$$\left(\sum_{i=1}^k a_i (\mathbf{z}_1'\mathbf{z}_i)(\mathbf{v}_1'\mathbf{v}_i)\right) = a_1.$$

Demonstração. Objetiva-se maximizar (ou minimizar) a expressão:

$$\begin{aligned}
\max((\mathbf{X}\mathbf{d})'(\mathbf{Y}\mathbf{e})) &= \max(\mathbf{d}'(\mathbf{X}'\mathbf{Y})\mathbf{e}) \\
&= \max\left(\mathbf{d}'\left(\sum_{i=1}^k a_i \mathbf{z}_i \mathbf{v}_i'\right)\mathbf{e}\right) \\
&= \max\left(\left(\sum_{j=1}^k d_j \mathbf{z}_j\right)' \left(\sum_{i=1}^k a_i \mathbf{z}_i \mathbf{v}_i'\right) \left(\sum_{s=1}^k e_s \mathbf{v}_s\right)\right) \\
&= \max \sum_{j,i,s} d_j a_i e_s \mathbf{z}_j'(\mathbf{z}_i \mathbf{v}_i') \mathbf{v}_s \\
&= \max \sum_{j,i,s} d_j a_i e_s (\mathbf{z}_j' \mathbf{z}_i) (\mathbf{v}_i' \mathbf{v}_s) \\
&= \max \sum_i d_i a_i e_i (\mathbf{z}_i' \mathbf{z}_i) (\mathbf{v}_i' \mathbf{v}_i) \\
&= \max_{\|\mathbf{d}\|=\|\mathbf{e}\|=1} \sum_i d_i a_i e_i
\end{aligned}$$

Assim, obtém-se a função lagrangiana com duas restrições, $\|\mathbf{d}\| = \|\mathbf{e}\| = 1$:

$$H(d_1, \dots, d_k, e_1, \dots, e_k, \lambda, \alpha) = \sum d_i a_i e_i + \lambda \left(\sum d_i^2 - 1 \right) + \alpha \left(\sum e_i^2 - 1 \right).$$

De onde segue o sistema de equações:

$$\begin{cases} \frac{\partial H}{\partial d_j} = a_j e_j + 2\lambda d_j = 0 \\ \frac{\partial H}{\partial e_1} = a_1 d_1 + 2\alpha e_1 = 0 \\ \sum d_i^2 = 1 \\ \sum e_i^2 = 1 \end{cases}$$

Multiplicando a 1ª equação por d_j e a 2ª equação por e_j , tem-se:

$$\begin{aligned} d_j a_j e_j + 2\lambda d_j^2 &= 0 \\ e_j a_j d_j + 2\alpha e_j^2 &= 0 \end{aligned}$$

Igualando as duas equações, segue que:

$$\lambda d_j^2 = \alpha e_j^2 \Rightarrow \lambda \sum d_j^2 = \alpha \sum e_j^2 \Rightarrow \lambda = \alpha \Rightarrow d_j = |e_j|.$$

Assim, volta-se ao problema da maximização, mas agora com apenas uma restrição:

$$\max_{\|\mathbf{d}\|=1} \sum a_i d_i^2.$$

Dada a função lagrangiana:

$$H(d_1, \dots, d_k, \lambda) = \sum a_i d_i^2 + \lambda \left(\sum d_i^2 - 1 \right),$$

segue que

$$\frac{\partial H}{\partial d_j} = 2a_j d_j + 2\lambda d_j = 0.$$

Logo, $a_j d_j = \lambda d_j$.

Supondo todos os a_i diferentes entre si, $a_i \neq a_j$ para $i \neq j$, a solução consiste em todos os d_j serem zero, exceto um deles que deve ser 1, pois $\|\mathbf{d}\| = 1$. Assim, o máximo ocorre para $i = 1$, pois a_1 é maior que os outros. Logo,

$$\max_{\|\mathbf{d}\|=\|\mathbf{e}\|=1} (\mathbf{d}' (\mathbf{X}'\mathbf{Y}) \mathbf{e}) = a_1.$$

□

Portanto, foram obtidos os vetores de correlação máxima $\mathbf{f} = \mathbf{X}\mathbf{z}_1$ e $\mathbf{g} = \mathbf{Y}\mathbf{v}_1$. O algoritmo PLS é justamente um processo iterativo para se obter outros pares de componentes que também maximizam a covariância com mais restrições (em ordem decrescente de magnitude).

4.3.2.2 O algoritmo PLS Multivariado

No sentido de melhor explicar os passos do algoritmo, seu desenvolvimento será paralelamente relacionado ao algoritmo PLS Populacional, anteriormente descrito. Com o objetivo de simplificar a notação, todo vetor da forma $\frac{\mathbf{a}}{\mathbf{a}'\mathbf{a}}$ receberá o mesmo nome \mathbf{a} .

O algoritmo inicia-se da seguinte forma: tome o vetor \mathbf{u} como a primeira coluna da matriz de respostas \mathbf{Y} , isto é, \mathbf{y}_1 . Faça $\frac{\mathbf{u}}{\mathbf{u}'\mathbf{u}}$ e seja $\mathbf{w} = \frac{\mathbf{X}'\mathbf{u}}{\mathbf{u}'\mathbf{u}}$. Logo, \mathbf{w} é uma combinação linear das linhas da matriz \mathbf{X} , ou seja, é um vetor de covariâncias

amostrais entre a primeira variável resposta e as covariáveis:

$$\begin{aligned} \mathbf{w} = \mathbf{X}'\mathbf{u} &= \begin{bmatrix} x_{11} & x_{21} & \cdots & x_{n1} \\ x_{12} & x_{22} & \cdots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1p} & x_{2p} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} y_{11} \\ y_{21} \\ \vdots \\ y_{n1} \end{bmatrix} \\ &= y_{11} \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1p} \end{bmatrix} + \cdots + y_{n1} \begin{bmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{np} \end{bmatrix} = \begin{bmatrix} \text{côv}(\mathbf{x}_1, \mathbf{y}_1) \\ \text{côv}(\mathbf{x}_2, \mathbf{y}_1) \\ \vdots \\ \text{côv}(\mathbf{x}_p, \mathbf{y}_1) \end{bmatrix}. \end{aligned}$$

Observe que $\mathbf{w} = \frac{\mathbf{X}'\mathbf{u}}{\mathbf{u}'\mathbf{u}}$ é o equivalente amostral do vetor $\mathbf{w}_1 = \boldsymbol{\sigma} = \text{cov}(\mathbf{x}, y)$, no algoritmo PLS populacional univariado, com $y = y_1$.

Normalize \mathbf{w} e seja $\mathbf{t} = \mathbf{X}\mathbf{w}$:

$$\mathbf{t} = \mathbf{X}\mathbf{w} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} \text{côv}(\mathbf{x}_1, \mathbf{y}_1) \\ \text{côv}(\mathbf{x}_2, \mathbf{y}_1) \\ \vdots \\ \text{côv}(\mathbf{x}_p, \mathbf{y}_1) \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^p x_{1i} \text{côv}(\mathbf{x}_i, \mathbf{y}_1) \\ \sum_{i=1}^p x_{2i} \text{côv}(\mathbf{x}_i, \mathbf{y}_1) \\ \vdots \\ \sum_{i=1}^p x_{ni} \text{côv}(\mathbf{x}_i, \mathbf{y}_1) \end{bmatrix}.$$

Veja que, a i -ésima entrada de \mathbf{t} corresponde a uma estimativa da variável $t = \boldsymbol{\sigma}'\mathbf{x} = (\text{cov}(\mathbf{x}, y))'\mathbf{x}$, com $y = y_1$.

Em seguida, faça $\mathbf{c} = \frac{\mathbf{Y}'\mathbf{t}}{\mathbf{t}'\mathbf{t}}$:

$$\begin{aligned} \mathbf{c} = \mathbf{Y}'\mathbf{t} &= \begin{bmatrix} y_{11} & y_{21} & \cdots & y_{n1} \\ y_{12} & y_{22} & \cdots & y_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ y_{1k} & y_{2k} & \cdots & y_{nk} \end{bmatrix} \begin{bmatrix} \sum_{i=1}^p x_{1i} \hat{\text{cov}}(x_i, \mathbf{y}_1) \\ \sum_{i=1}^p x_{2i} \hat{\text{cov}}(x_i, \mathbf{y}_1) \\ \vdots \\ \sum_{i=1}^p x_{ni} \hat{\text{cov}}(x_i, \mathbf{y}_1) \end{bmatrix} \\ &= \begin{bmatrix} \sum_{j=1}^n \sum_{i=1}^p y_{j1} x_{ji} \hat{\text{cov}}(x_i, \mathbf{y}_1) \\ \sum_{j=1}^n \sum_{i=1}^p y_{j2} x_{ji} \hat{\text{cov}}(x_i, \mathbf{y}_1) \\ \vdots \\ \sum_{j=1}^n \sum_{i=1}^p y_{jk} x_{ji} \hat{\text{cov}}(x_i, \mathbf{y}_1) \end{bmatrix} . \end{aligned}$$

Fazendo correspondência com o PLS populacional, a i -ésima entrada do vetor \mathbf{c} é uma estimativa do parâmetro $q_{1i} = \text{cov}(y, \boldsymbol{\sigma}'\mathbf{x}) = \text{cov}(y, t_1)$, com $y = y_1$.

Normalize \mathbf{c} e tome $\mathbf{u} = \mathbf{Y}\mathbf{c}$:

$$\mathbf{u} = \mathbf{Y}\mathbf{c} = \begin{bmatrix} y_{11} & y_{12} & \cdots & y_{1k} \\ y_{21} & y_{22} & \cdots & y_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \cdots & y_{nk} \end{bmatrix} \begin{bmatrix} \sum_{j=1}^n \sum_{i=1}^p y_{j1} x_{ji} \hat{\text{cov}}(x_i, \mathbf{y}_1) \\ \sum_{j=1}^n \sum_{i=1}^p y_{j2} x_{ji} \hat{\text{cov}}(x_i, \mathbf{y}_1) \\ \vdots \\ \sum_{j=1}^n \sum_{i=1}^p y_{jk} x_{ji} \hat{\text{cov}}(x_i, \mathbf{y}_1) \end{bmatrix}$$

$$= \begin{bmatrix} \sum_{l=1}^k \sum_{j=1}^n \sum_{i=1}^p y_{1l} y_{j1} x_{ji} \hat{c} \hat{v}(x_i, \mathbf{y}_1) \\ \sum_{l=1}^k \sum_{j=1}^n \sum_{i=1}^p y_{2l} y_{j1} x_{ji} \hat{c} \hat{v}(x_i, \mathbf{y}_1) \\ \vdots \\ \sum_{l=1}^k \sum_{j=1}^n \sum_{i=1}^p y_{nl} y_{j1} x_{ji} \hat{c} \hat{v}(x_i, \mathbf{y}_1) \end{bmatrix}.$$

Este vetor não tem uma interpretação em termos do PLS populacional.

Nesta etapa, o algoritmo define um *loop* obtendo seqüências de vetores $\mathbf{u}_n = \mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{u}_{n-1}$, $\mathbf{c}_n = \mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{c}_{n-1}$, $\mathbf{t}_n = \mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{t}_{n-1}$ e $\mathbf{w}_n = \mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{w}_{n-1}$. Essas seqüências são convergentes em razão de uma série de propriedades de álgebra linear, relativas ao diagrama da Figura 22.

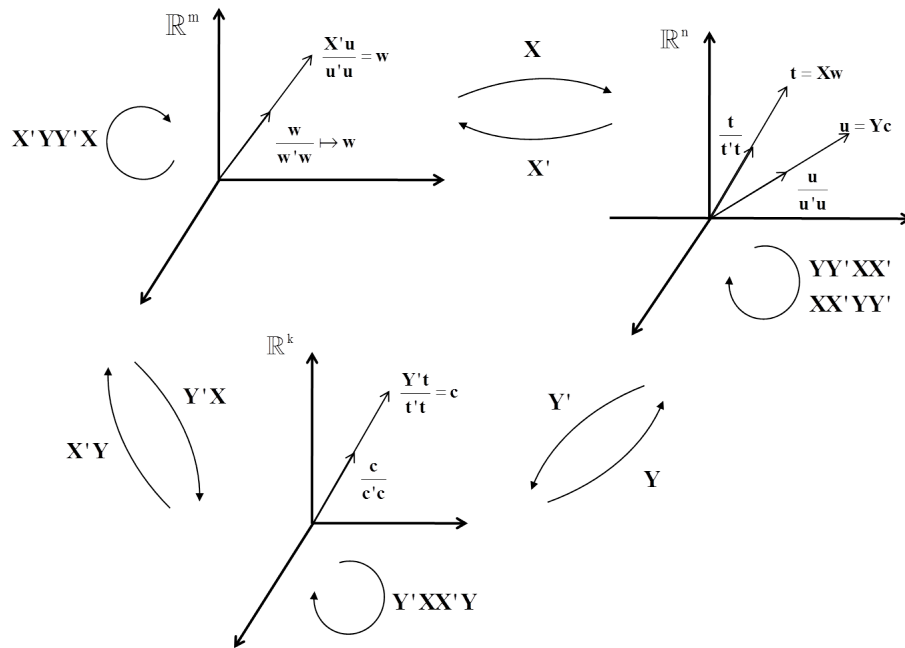


Figura 22 Representação geométrica das seqüências de vetores $\mathbf{u}_n = \mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{u}_{n-1}$, $\mathbf{c}_n = \mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{c}_{n-1}$, $\mathbf{t}_n = \mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{t}_{n-1}$ e $\mathbf{w}_n = \mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{w}_{n-1}$

Observe que:

$$\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X} = \left(\sum_{i=1}^k a_i \mathbf{z}_i \mathbf{v}_i' \right) \left(\sum_{j=1}^k a_j \mathbf{v}_j \mathbf{z}_j' \right) = \sum_{i=1}^k \sum_{j=1}^k a_i a_j \mathbf{z}_i \mathbf{v}_i' \mathbf{v}_j \mathbf{z}_j' = \sum_{i=1}^k a_i^2 \mathbf{z}_i \mathbf{z}_i' .$$

A partir dessa decomposição, os autovalores de $\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X}$ são $a_1^2 \geq a_2^2 \geq \dots \geq a_k^2$, com respectivos autovetores $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k$, pois:

$$\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X} (\mathbf{z}_j) = \sum_{i=1}^k a_i^2 \mathbf{z}_i \mathbf{z}_i' \mathbf{z}_j = a_j^2 \mathbf{z}_j .$$

- As matrizes $\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X}$, $\mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{X}'$, $\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y}$ e $\mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{Y}'$ possuem os mesmos autovalores, pois:

se \mathbf{r} é um autovetor da matriz $\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X}$, então

$$\mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{Y}' (\mathbf{X}\mathbf{r}) = \mathbf{X} (\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X}) \mathbf{r} = \mathbf{X} (\lambda \mathbf{r}) = \lambda (\mathbf{X}\mathbf{r}) .$$

Se \mathbf{r} é um autovetor da matriz $\mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{Y}'$, então

$$\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y} (\mathbf{Y}'\mathbf{r}) = \mathbf{Y}' (\mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{Y}') \mathbf{r} = \mathbf{Y}' (\lambda \mathbf{r}) = \lambda (\mathbf{Y}'\mathbf{r}) .$$

Se \mathbf{r} é um autovetor da matriz $\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y}$, então

$$\mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{X}' (\mathbf{Y}\mathbf{r}) = \mathbf{Y} (\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y}) \mathbf{r} = \mathbf{Y} (\lambda \mathbf{r}) = \lambda (\mathbf{Y}\mathbf{r}) .$$

Portanto, o maior autovalor para todas as matrizes é dado por a_1^2 .

- Aplicando o método das potências (Apêndice C) para o cálculo do autovetor relativo ao maior autovalor às sequências $\mathbf{u}_n = \mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{u}_{n-1}$, $\mathbf{c}_n = \mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{c}_{n-1}$, $\mathbf{t}_n = \mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{t}_{n-1}$, e $\mathbf{w}_n = \mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{w}_{n-1}$, obtém-se

os vetores limites das sequências $\mathbf{u} = \lim \mathbf{u}_n$, $\mathbf{c} = \lim \mathbf{c}_n$, $\mathbf{t} = \lim \mathbf{t}_n$ e $\mathbf{w} = \lim \mathbf{w}_n$. Assim:

$$\begin{aligned} \mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{u} &= a_1^2\mathbf{u} \\ \mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{c} &= a_1^2\mathbf{c} \\ \mathbf{X}\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{t} &= a_1^2\mathbf{t} \\ \mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{X}\mathbf{w} &= a_1^2\mathbf{w}. \end{aligned}$$

Ocorrendo a convergência pela aplicação de algum critério de parada, ficam definidos os vetores \mathbf{w} e \mathbf{c} , relativos ao maior autovalor a_1^2 . Pelo problema de maximização das componentes com maior covariância, os vetores relativos ao maior autovalor são dados explicitamente a partir da decomposição em valores singulares $\mathbf{X}'\mathbf{Y}$, isto é, $\mathbf{w} = \mathbf{z}_1$ e $\mathbf{c} = \mathbf{v}_1$. Portanto, as componentes com maior covariância são dadas por $\mathbf{t} = \mathbf{X}\mathbf{w}$ e $\mathbf{u} = \mathbf{Y}\mathbf{c}$.

O próximo passo do algoritmo é definido por: primeiramente se calcula a covariância amostral entre as componentes $\mathbf{t} = \mathbf{X}\mathbf{w}$ e $\mathbf{u} = \mathbf{Y}\mathbf{c}$, dada por $b = \frac{\mathbf{u}'\mathbf{t}}{\mathbf{t}'\mathbf{t}}$.

Como o vetor \mathbf{t} é a combinação linear das colunas de \mathbf{X} que possui maior covariância com uma combinação linear das colunas de \mathbf{Y} , a ideia agora é obter uma nova matriz \mathbf{X} de tal forma que as colunas desta nova matriz \mathbf{X} sejam ortogonais ao vetor \mathbf{t} . A projeção ortogonal no subespaço ortogonal ao vetor \mathbf{t} é dada por $\mathbf{I} - \frac{\mathbf{t}\mathbf{t}'}{\mathbf{t}'\mathbf{t}}$. Veja:

$$\left(\mathbf{I} - \frac{\mathbf{t}\mathbf{t}'}{\mathbf{t}'\mathbf{t}}\right) \left(\mathbf{I} - \frac{\mathbf{t}\mathbf{t}'}{\mathbf{t}'\mathbf{t}}\right) = \mathbf{I} - \frac{\mathbf{t}\mathbf{t}'}{\mathbf{t}'\mathbf{t}} - \frac{\mathbf{t}\mathbf{t}'}{\mathbf{t}'\mathbf{t}} + \frac{\mathbf{t}\mathbf{t}'\mathbf{t}\mathbf{t}'}{\mathbf{t}'\mathbf{t}} = \mathbf{I} - \frac{\mathbf{t}\mathbf{t}'}{\mathbf{t}'\mathbf{t}}.$$

Portanto, a nova matriz \mathbf{X} é dada por:

$$\left(\mathbf{I} - \frac{\mathbf{t}\mathbf{t}'}{\mathbf{t}'\mathbf{t}} \right) \mathbf{X}.$$

Assim,

$$\left(\mathbf{I} - \frac{\mathbf{t}\mathbf{t}'}{\mathbf{t}'\mathbf{t}} \right) \mathbf{X} = \mathbf{X} - \frac{\mathbf{t}\mathbf{t}'}{\mathbf{t}'\mathbf{t}} \mathbf{X} = \mathbf{X} - \mathbf{t} \frac{\mathbf{t}'\mathbf{X}}{\mathbf{t}'\mathbf{t}} = \mathbf{X} - \mathbf{t}\mathbf{p}',$$

em que $\mathbf{p} = \frac{\mathbf{X}'\mathbf{t}}{\mathbf{t}'\mathbf{t}}$.

O vetor \mathbf{p} representa o vetor de covariâncias amostrais entre as covariáveis dadas pelas colunas de \mathbf{X} e o vetor de estimativas das variáveis predictoras \mathbf{t} , e a nova matriz \mathbf{X} representa as componentes das covariáveis que não foram utilizadas na predição, sendo denominada matriz residual e que será utilizada na próxima iteração do algoritmo.

Uma construção semelhante é feita para a matriz \mathbf{Y} . O melhor preditor

das variáveis respostas $\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_k \end{bmatrix}$ é o vetor $\text{cov}(\mathbf{t}, \mathbf{y}) = \begin{bmatrix} \text{cov}(\mathbf{t}, \mathbf{y}_1) \\ \text{cov}(\mathbf{t}, \mathbf{y}_2) \\ \vdots \\ \text{cov}(\mathbf{t}, \mathbf{y}_k) \end{bmatrix}$, que é

estimado pelo vetor \mathbf{c} , em que,

$$\mathbf{c}_{k \times 1} = \frac{\mathbf{Y}'\mathbf{t}}{\mathbf{t}'\mathbf{t}} = \begin{bmatrix} \mathbf{Y}'_1 \\ \mathbf{Y}'_2 \\ \vdots \\ \mathbf{Y}'_k \end{bmatrix}_{k \times n} \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} \frac{\mathbf{Y}'_1 \mathbf{t}}{\mathbf{t}'\mathbf{t}} \\ \frac{\mathbf{Y}'_2 \mathbf{t}}{\mathbf{t}'\mathbf{t}} \\ \vdots \\ \frac{\mathbf{Y}'_k \mathbf{t}}{\mathbf{t}'\mathbf{t}} \end{bmatrix}_{k \times 1}.$$

Segue que,

$$\mathbf{tc}' = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}_{n \times 1} \begin{bmatrix} \frac{\mathbf{Y}'_1 \mathbf{t}}{\mathbf{t}'\mathbf{t}}, \frac{\mathbf{Y}'_2 \mathbf{t}}{\mathbf{t}'\mathbf{t}}, \dots, \frac{\mathbf{Y}'_k \mathbf{t}}{\mathbf{t}'\mathbf{t}} \end{bmatrix}_{1 \times k} = \begin{bmatrix} \frac{\mathbf{Y}'_1 \mathbf{t}}{\mathbf{t}'\mathbf{t}} \mathbf{t}, \frac{\mathbf{Y}'_2 \mathbf{t}}{\mathbf{t}'\mathbf{t}} \mathbf{t}, \dots, \frac{\mathbf{Y}'_k \mathbf{t}}{\mathbf{t}'\mathbf{t}} \mathbf{t} \end{bmatrix}_{n \times k}.$$

As colunas da matriz \mathbf{tc}' representam a parte que o vetor \mathbf{t} explica de \mathbf{Y} . Tem-se, portanto, duas fontes de explicação da variável \mathbf{Y} , a saber, a parte de \mathbf{t} que explica \mathbf{u} , que é uma combinação linear das colunas de \mathbf{Y} , e a parte de \mathbf{t} que explica \mathbf{Y} . Portanto, a quantidade de informação que o vetor \mathbf{t} tem do vetor \mathbf{u} e dos vetores coluna da matriz \mathbf{Y} é obtida pelo produto $b \frac{\mathbf{t}'\mathbf{Y}_i}{\mathbf{t}'\mathbf{t}}$.

Assim, ponderando a quantidade de informação pelo produto, a parte não explicada de \mathbf{Y} é a nova matriz \mathbf{Y} , dada pela diferença:

$$\mathbf{Y} - \mathbf{bt} \left(\frac{\mathbf{t}\mathbf{Y}'}{\mathbf{t}'\mathbf{t}} \right) = \mathbf{Y} - \mathbf{btc}'.$$

Obtidas as novas matrizes \mathbf{X} e \mathbf{Y} , o processo é repetido, ficando o algoritmo da forma:

Algoritmo para o PLS Amostral (HOSKULDSSON, 1988):

1. Seja \mathbf{u} a primeira coluna de \mathbf{Y}
2. $\mathbf{w} = \frac{\mathbf{X}'\mathbf{u}}{\mathbf{u}'\mathbf{u}}$
3. Normalize \mathbf{w}
4. $\mathbf{t} = \mathbf{X}\mathbf{w}$
5. $\mathbf{c} = \frac{\mathbf{Y}'\mathbf{t}}{\mathbf{t}'\mathbf{t}}$

6. Normalize \mathbf{c}
7. $\mathbf{u} = \mathbf{Y}\mathbf{c}$
8. Se houver convergência vá para o passo 9, caso contrário, volte ao passo 2
9. \mathbf{X} -loadings: $\mathbf{p} = \frac{\mathbf{X}'\mathbf{t}}{\mathbf{t}'\mathbf{t}}$
10. \mathbf{Y} -loadings: $\mathbf{q} = \frac{\mathbf{Y}'\mathbf{u}}{\mathbf{u}'\mathbf{u}}$
11. Regressão (\mathbf{u} sobre \mathbf{t}): $b = \frac{\mathbf{u}'\mathbf{t}}{\mathbf{t}'\mathbf{t}}$
12. Matrizes Residuais: $\mathbf{X} \leftarrow \mathbf{X} - \mathbf{t}\mathbf{p}'$ e $\mathbf{Y} \leftarrow \mathbf{Y} - \mathbf{b}\mathbf{t}\mathbf{c}'$

Realizadas A iterações deste algoritmo, obtêm-se os vetores:

$$\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_A,$$

$$\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_A,$$

$$\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_A,$$

$$\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_A,$$

$$\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_A.$$

- É importante observar que da construção do algoritmo segue que a sequência dos vetores de componentes $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_A$, formam uma sequência de vetores ortogonais.
- Os vetores $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_A$, são vetores mutuamente ortogonais.

Demonstração. Do algoritmo, sabe-se que:

$$\begin{aligned}
 \mathbf{X}_2 &= \left(\mathbf{I} - \frac{\mathbf{t}_1 \mathbf{t}'_1}{\mathbf{t}'_1 \mathbf{t}_1} \right) \mathbf{X}_1 \\
 \mathbf{X}_3 &= \left(\mathbf{I} - \frac{\mathbf{t}_2 \mathbf{t}'_2}{\mathbf{t}'_2 \mathbf{t}_2} \right) \mathbf{X}_2 = \left(\mathbf{I} - \frac{\mathbf{t}_2 \mathbf{t}'_2}{\mathbf{t}'_2 \mathbf{t}_2} \right) \left(\mathbf{I} - \frac{\mathbf{t}_1 \mathbf{t}'_1}{\mathbf{t}'_1 \mathbf{t}_1} \right) \mathbf{X}_1 \\
 &\vdots \\
 \mathbf{X}_j &= \left(\mathbf{I} - \frac{\mathbf{t}_{j-1} \mathbf{t}'_{j-1}}{\mathbf{t}'_{j-1} \mathbf{t}_{j-1}} \right) \mathbf{X}_{j-1} \\
 &= \left(\mathbf{I} - \frac{\mathbf{t}_{j-1} \mathbf{t}'_{j-1}}{\mathbf{t}'_{j-1} \mathbf{t}_{j-1}} \right) \left(\mathbf{I} - \frac{\mathbf{t}_{j-2} \mathbf{t}'_{j-2}}{\mathbf{t}'_{j-2} \mathbf{t}_{j-2}} \right) \dots \left(\mathbf{I} - \frac{\mathbf{t}_1 \mathbf{t}'_1}{\mathbf{t}'_1 \mathbf{t}_1} \right) \mathbf{X}_1.
 \end{aligned}$$

Assim, para $i < j$, tem-se:

$$\mathbf{X}_j = \mathbf{Z} \left(\mathbf{I} - \frac{\mathbf{t}_i \mathbf{t}'_i}{\mathbf{t}'_i \mathbf{t}_i} \right) \mathbf{X}_i,$$

em que \mathbf{Z} é o produto de matrizes $\left(\mathbf{I} - \frac{\mathbf{t}_k \mathbf{t}'_k}{\mathbf{t}'_k \mathbf{t}_k} \right)$, $k = i + 1, \dots, j$.

De onde segue que, para $i < j$:

$$\begin{aligned}
 \mathbf{X}_j \mathbf{w}_i &= \mathbf{Z} \left(\mathbf{X}_i - \frac{\mathbf{t}_i \mathbf{t}'_i}{\mathbf{t}'_i \mathbf{t}_i} \mathbf{X}_i \right) \mathbf{w}_i = \mathbf{Z} \left(\mathbf{X}_i \mathbf{w}_i - \frac{\mathbf{t}_i \mathbf{t}'_i}{\mathbf{t}'_i \mathbf{t}_i} \mathbf{X}_i \mathbf{w}_i \right) \\
 &= \mathbf{Z} \left(\mathbf{t}_i - \frac{\mathbf{t}_i \mathbf{t}'_i}{\mathbf{t}'_i \mathbf{t}_i} \mathbf{t}_i \right) = \mathbf{Z} (\mathbf{t}_i - \mathbf{t}_i) = 0.
 \end{aligned}$$

Sabendo que $\mathbf{X}'_j \mathbf{Y}_j \mathbf{Y}'_j \mathbf{X}_j \mathbf{w}_j = a_j \mathbf{w}_j$, obtém-se:

$$\mathbf{w}'_j \mathbf{w}_i = \left(\frac{\mathbf{X}'_j \mathbf{Y}_j \mathbf{Y}'_j \mathbf{X}_j \mathbf{w}_j}{a_j} \right)' \mathbf{w}_i = \frac{\mathbf{w}'_j \mathbf{X}'_j \mathbf{Y}_j \mathbf{Y}'_j \mathbf{X}_j}{a_j} \mathbf{w}_i$$

$$= \frac{\mathbf{w}_j' \mathbf{X}_j' \mathbf{Y}_j \mathbf{Y}_j'}{a_j} (\mathbf{X}_j \mathbf{w}_i) = 0.$$

□

- Os vetores \mathbf{w}_i são ortogonais aos vetores \mathbf{p}_j , para $i < j$.

Demonstração.

$$\mathbf{w}_i' \mathbf{p}_j = \mathbf{w}_i' \frac{\mathbf{X}_j' \mathbf{t}_j}{\mathbf{t}_j' \mathbf{t}_j} = \mathbf{w}_i' \mathbf{X}_j' \frac{\mathbf{t}_j}{\mathbf{t}_j' \mathbf{t}_j} = (\mathbf{X}_j \mathbf{w}_i)' \frac{\mathbf{t}_j}{\mathbf{t}_j' \mathbf{t}_j} = 0.$$

□

- A sequência de vetores $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_A$, também é uma sequência de vetores ortogonais.

4.3.2.3 Regressão em PLS

Uma das questões mais importantes em análise de regressão é reduzir o número de variáveis independentes. Uma situação típica em que é vantajosa a seleção de variáveis é que a variância dos coeficientes de regressão aumenta quando novas variáveis são introduzidas no modelo (HOSKULDSSON, 1988). A regressão PLS pode ser vista como um bom método de análise de regressão porque os componentes são selecionados de modo que eles descrevem as variáveis respostas. A principal característica é que o método PLS é capaz de reduzir o número de variáveis, em termos de componentes, durante o processo de estimação. Vamos considerar aqui, a regressão via método PLS.

Um modelo de regressão linear pode ser escrito como:

$$\mathbf{Y}_{n \times k} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times k} + \mathbf{e}_{n \times k},$$

em que $\boldsymbol{\beta}$ é o vetor dos coeficientes de regressão, \mathbf{X} é a matriz do delineamento e \mathbf{e} é a matriz residual.

Sendo obtidos os vetores $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_A$, em A iterações, a matriz do delineamento \mathbf{X} pode ser decomposta da forma

$$\mathbf{X}_{n \times p} = \sum_{i=1}^A \mathbf{t}_i \mathbf{p}_i' + \mathbf{X}_0,$$

onde \mathbf{X}_0 é a matriz residual do processo.

Note que se o número de iterações A é igual ao posto de \mathbf{X} , então $\mathbf{X}_0 = 0$.

Desprezando a matriz residual \mathbf{X}_0 , tem-se

$$\mathbf{X}_{n \times p} = \sum_{i=1}^A \mathbf{t}_i \mathbf{p}_i' = \mathbf{T}_{n \times A} \mathbf{P}'_{A \times p}$$

com $\mathbf{T}_{n \times A} = (\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_A)$ em que \mathbf{t}_i é a i -ésima coluna de \mathbf{T} e $\mathbf{P}_{n \times A} = (\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_A)$ em que \mathbf{p}_i é a i -ésima coluna de \mathbf{P} .

Como os vetores \mathbf{p}_i são ortonormais, multiplica-se a igualdade $\mathbf{X}_{n \times p} = \mathbf{T}_{n \times A} \mathbf{P}'_{A \times p}$ por $\mathbf{P}_{p \times A}$, obtendo

$$\mathbf{X}_{n \times p} \mathbf{P}_{p \times A} = \mathbf{T}_{n \times A} \mathbf{P}'_{A \times p} \mathbf{P}_{p \times A} = \mathbf{T}_{n \times A} \mathbf{I}_{A \times A} = \mathbf{T}_{n \times A}.$$

Observe que \mathbf{P} é uma inversa generalizada à direita de \mathbf{P}' . De fato, qualquer inversa generalizada poderia ser utilizada (JONG, 1993).

Assim, $\mathbf{T} = \mathbf{XP}$, e a equação de regressão pode ser escrita como

$$\mathbf{Y}_{n \times k} = \mathbf{X}_{n \times p} \mathbf{P}_{p \times A} \mathbf{P}'_{A \times p} \boldsymbol{\beta}_{p \times k} + \mathbf{e}_{n \times k} = \mathbf{T}_{n \times A} (\boldsymbol{\beta}_{\text{PLS},A})_{A \times k} + \mathbf{e}_{n \times k},$$

com \mathbf{T} a nova matriz de covariáveis, denominada matriz de componentes PLS, e $\boldsymbol{\beta}_{\text{PLS},A} = \mathbf{P}'\boldsymbol{\beta}$ a nova matriz de parâmetros de regressão, denominada matriz de parâmetros PLS.

Com os novos parâmetros PLS, temos um problema de regressão usual, porém, agora, com o número de parâmetros muito menor que o problema original. O vetor de parâmetros, $\boldsymbol{\beta}_{\text{PLS},A}$, é estimado pelo método dos quadrados mínimos, isto é,

$$\left(\hat{\boldsymbol{\beta}}_{\text{PLS},A} \right)_{A \times k} = (\mathbf{T}'\mathbf{T})_{A \times A}^{-1} \mathbf{T}'_{A \times n} \mathbf{Y}_{n \times k}. \quad (5)$$

Obtidas as estimativas dos coeficientes de regressão, a equação de predição para novos valores \mathbf{x}_0 das covariáveis, retornando às variáveis originais não centradas, fica da forma:

$$\hat{\mathbf{y}} - \hat{\boldsymbol{\mu}}_{\mathbf{y}} = (\mathbf{t} - \hat{\boldsymbol{\mu}}_{\mathbf{t}})' \hat{\boldsymbol{\beta}}_{\text{PLS},A} = (\mathbf{x}_0 - \hat{\boldsymbol{\mu}}_{\mathbf{x}})' \mathbf{P} \hat{\boldsymbol{\beta}}_{\text{PLS},A}. \quad (6)$$

4.3.3 Uma alternativa ao algoritmo PLS

Esta seção tem como referências os artigos Garthwaite (1994) e Phatak e Jong (1997), e será utilizada a mesma notação. A descrição deste algoritmo é feita para apenas uma resposta, $k = 1$.

Como anteriormente exposto, a matriz de covariáveis $\mathbf{X}_{n \times p}$, é considerada como uma transformação do espaço dos parâmetros \mathbb{R}^p no espaço dos dados \mathbb{R}^n , e o vetor da variável resposta é $\mathbf{Y}_{n \times 1}$, conforme Figura 23.

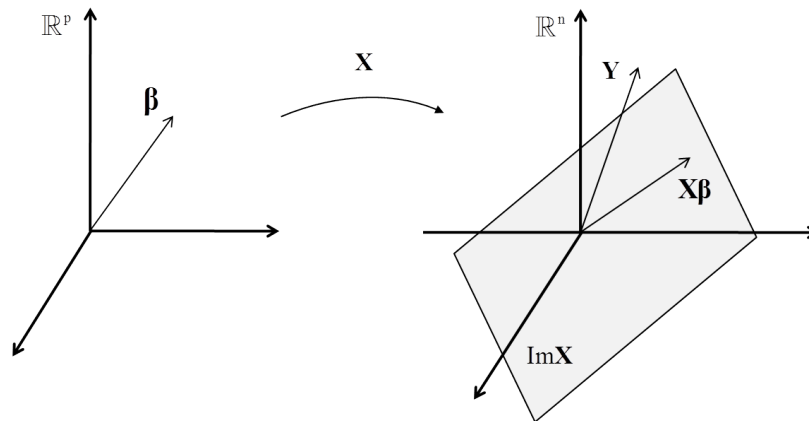


Figura 23 Representação da matriz de covariáveis $\mathbf{X}_{n \times p}$ como uma transformação do espaço de parâmetros \mathbb{R}^p no espaço dos dados \mathbb{R}^n

A matriz \mathbf{X} aplicada ao vetor canônico $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)$ tem como imagem a i -ésima coluna de \mathbf{X} , denotada por \mathbf{X}_i (Figura 24).

Como é usual, utilizam-se variáveis centradas na média:

$$\mathbf{U}_1 = \mathbf{Y} - \bar{\mathbf{Y}} \quad \text{e} \quad \mathbf{V}_{1j} = \mathbf{X}_j - \bar{\mathbf{X}}_j \quad \text{para } j = 1, \dots, p,$$

sendo $\bar{\mathbf{Y}}' = (\bar{y}_1, \bar{y}_2, \dots, \bar{y}_n)'$ e $\mathbf{1}' = (1, 1, \dots, 1)'$, vetores linha n -dimensionais. Observe, como na Figura 25, que os vetores $\mathbf{Y} - \bar{\mathbf{Y}}$ e $\mathbf{X}_j - \bar{\mathbf{X}}_j$ são ortogonais ao

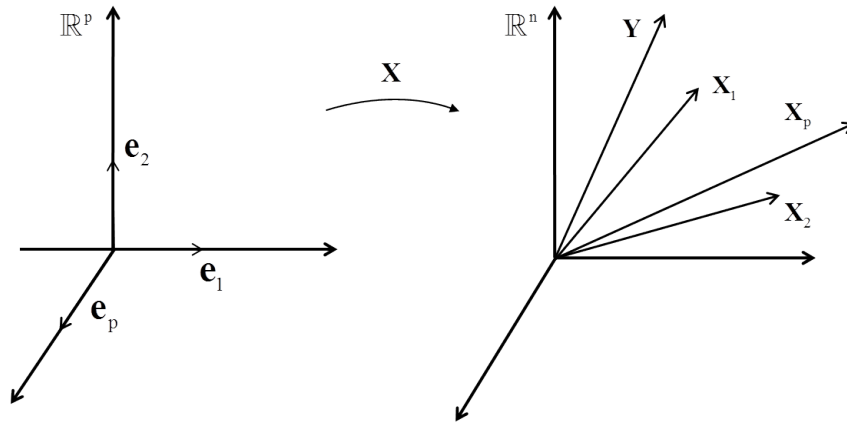


Figura 24 Representação das colunas de X como imagem dos vetores canônicos e_i

subespaço gerado pelo vetor $\mathbf{1}$.

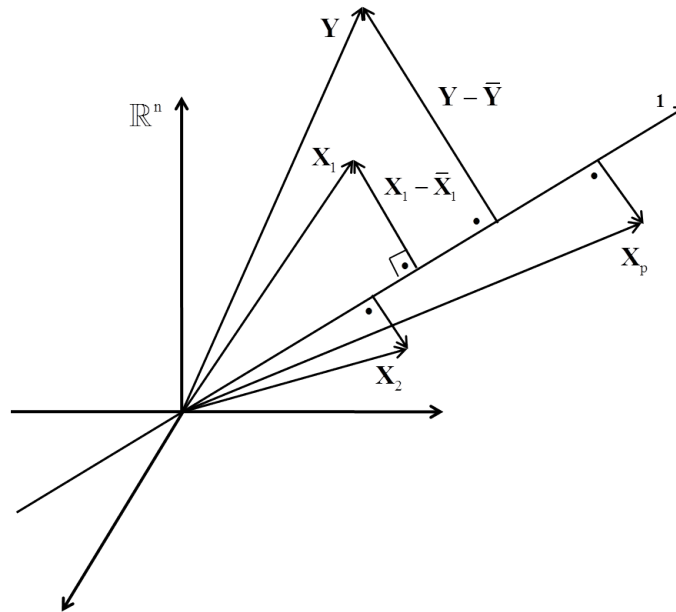


Figura 25 Representação das variáveis centradas através da projeção ortogonal no subespaço gerado pelo vetor $\mathbf{1}$

A ideia é regressir o vetor \mathbf{U}_1 em cada um dos vetores \mathbf{V}_{1j} (Figuras 26 e 27). Tal procedimento é uma forma de medir o quanto de informação o vetor da covariável \mathbf{X}_j explica o vetor de dados, obtendo-se o vetor $\mathbf{U}_{1j} = \frac{\mathbf{u}'_1 \mathbf{v}_{1j}}{\mathbf{v}'_{1j} \mathbf{v}_{1j}} \mathbf{v}_{1j}$.

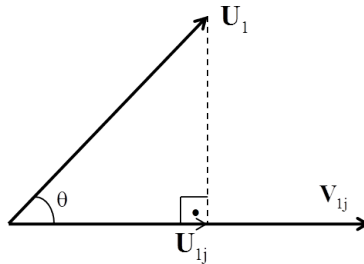


Figura 26 Projeção ortogonal do vetor \mathbf{U}_1 em cada um dos vetores \mathbf{V}_{1j}

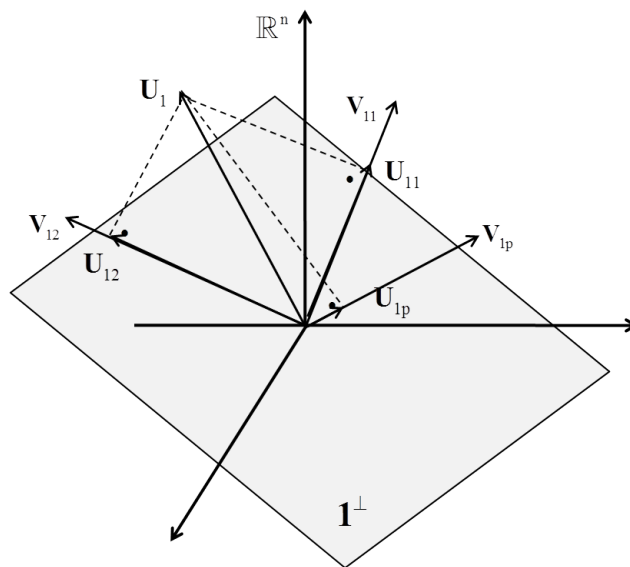


Figura 27 Projeção ortogonal do vetor \mathbf{U}_1 em cada um dos vetores \mathbf{V}_{1j}

Uma maneira de se coletarem todas estas informações sobre o vetor de dados é através de uma média ponderada dos vetores \mathbf{U}_{1j} , utilizando pesos apro-

prriadamente escolhidos, $w_{1j} \geq 0$ e $\sum_{j=1}^p w_{1j} = 1$, como na Figura 28.

$$\mathbf{T}_1 = \sum_{j=1}^p w_{1j} \mathbf{U}_{1j}. \quad (7)$$

Uma escolha usual para os pesos é a média simples, $w_{1j} = 1/p$.

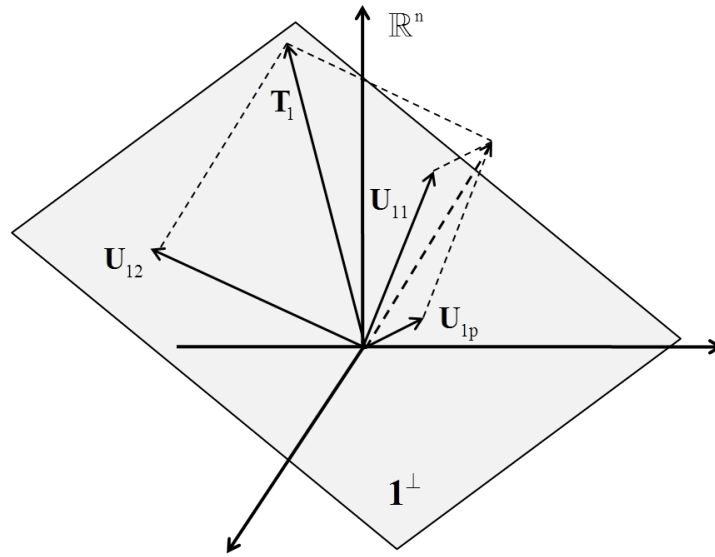


Figura 28 Construção do vetor \mathbf{T}_1 como média ponderada dos vetores \mathbf{U}_{1j}

A quantidade de informação que \mathbf{T}_1 contém do vetor de dados \mathbf{U}_1 é obtida pela regressão do vetor \mathbf{U}_1 em \mathbf{T}_1 , $P_{\mathbf{T}_1} \mathbf{U}_1 = \frac{\mathbf{u}'_1 \mathbf{t}_1}{\mathbf{t}'_1 \mathbf{t}_1} \mathbf{t}_1$. A parte não explicada é definida pelo vetor de resíduo $\mathbf{U}_2 = \mathbf{U}_1 - P_{\mathbf{T}_1} \mathbf{U}_1$. Novamente, regride-se \mathbf{V}_{1j} em \mathbf{T}_1 , $P_{\mathbf{T}_1} \mathbf{V}_{1j} = \frac{\mathbf{v}'_{1j} \mathbf{t}_1}{\mathbf{t}'_1 \mathbf{t}_1} \mathbf{t}_1$, e toma-se os vetores de resíduos $\mathbf{V}_{2j} = \mathbf{V}_{1j} - P_{\mathbf{T}_1} \mathbf{V}_{1j}$, que são uma medida do quanto o vetor \mathbf{T}_1 não explica as covariáveis, conforme Figura 29.

Procede-se a regressão do vetor \mathbf{U}_2 nos vetores \mathbf{V}_{2j} , que são denotados por $\mathbf{U}_{2j} = P_{\mathbf{V}_{2j}} \mathbf{U}_2 = b_{2j} \mathbf{V}_{2j} = \frac{\mathbf{U}'_2 \mathbf{V}_{2j}}{\mathbf{V}'_{2j} \mathbf{V}_{2j}} \mathbf{V}_{2j}$ (Figura 30).

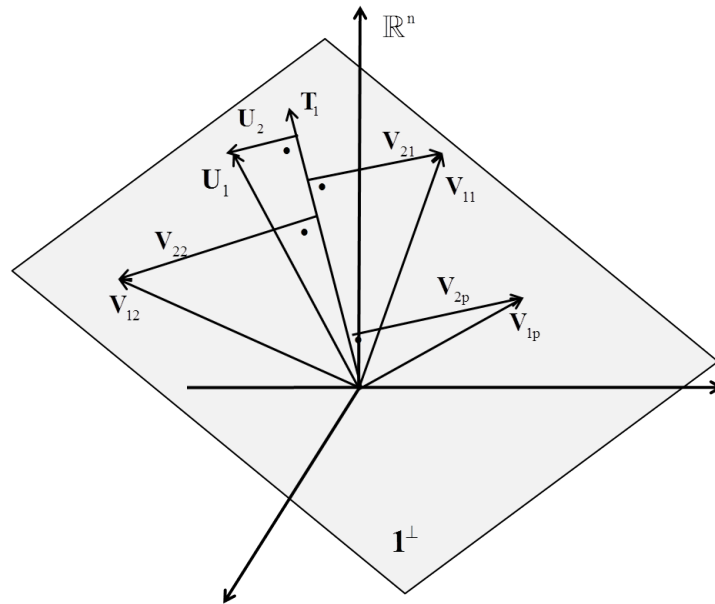


Figura 29 Construção dos vetores \mathbf{U}_2 e \mathbf{V}_{2j} através da projeção dos vetores \mathbf{U}_1 e \mathbf{V}_{1j} , respectivamente, em \mathbf{T}_1

Utilizando pesos w_{2j} , encontramos o vetor \mathbf{T}_2 fazendo a média ponderada dos vetores \mathbf{U}_{2j} (Figura 31).

Observe que, por construção, os vetores \mathbf{T}_1 e \mathbf{T}_2 são ortogonais (Figura 32).

O procedimento pode ser repetido para a construção de outros vetores $\mathbf{T}_3, \mathbf{T}_4, \dots, \mathbf{T}_m$.

Suponha que \mathbf{T}_i ($i \geq 1$) tenha sido construído a partir das variáveis \mathbf{U}_i e \mathbf{V}_{ij} com $j = 1, \dots, p$. Para obter o componente \mathbf{T}_{i+1} , as variáveis $\mathbf{V}_{(i+1)j}$ e $\mathbf{U}_{(i+1)}$ devem ser determinadas. Com este propósito, é feita uma regressão entre \mathbf{V}_{ij} e \mathbf{T}_i e $\mathbf{V}_{(i+1)j}$ fica definido por:

$$\mathbf{V}_{(i+1)j} = \mathbf{V}_{ij} - \frac{\mathbf{v}'_{ij} \mathbf{t}_i}{\mathbf{t}'_i \mathbf{t}_i} \mathbf{t}_i,$$

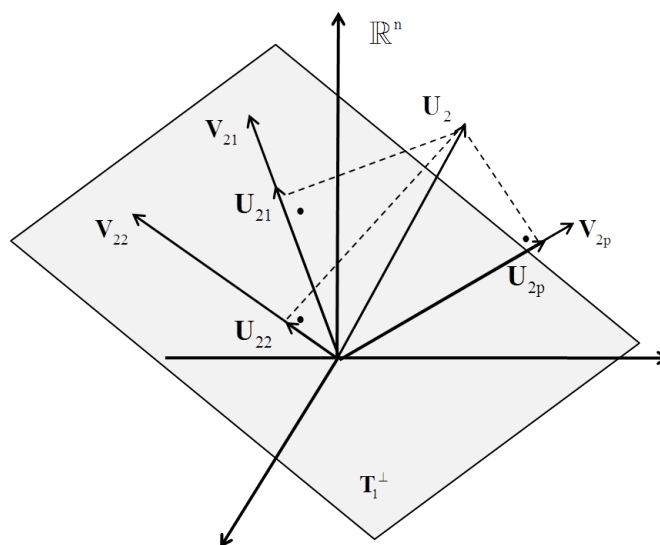


Figura 30 Construção dos vetores U_{2j} como projeção do vetor U_2 nos vetores V_{2j}

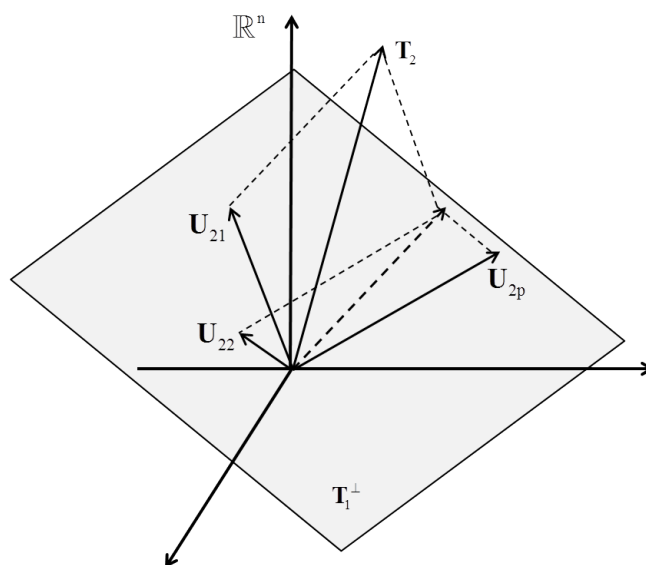


Figura 31 Construção do vetor T_2 como média ponderada dos vetores U_{2j}

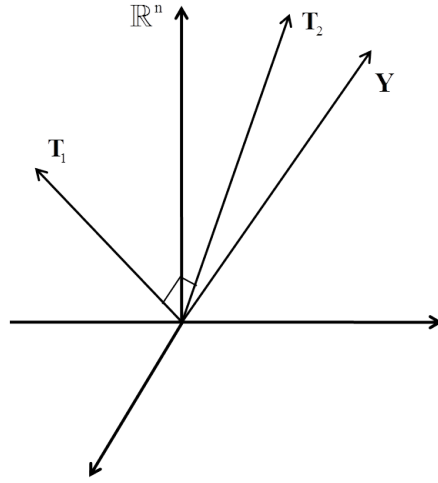


Figura 32 Componentes ortogonais \mathbf{T}_1 e \mathbf{T}_2

sendo \mathbf{t}_i o vetor de valores de \mathbf{T}_i , $\mathbf{v}_{(i+1)j}$ os resíduos da regressão e $\frac{\mathbf{v}'_{ij}\mathbf{t}_i}{\mathbf{t}'_i\mathbf{t}_i}$, o coeficiente da regressão entre \mathbf{V}_{ij} e \mathbf{T}_i .

De forma análoga, $\mathbf{U}_{(i+1)} = \mathbf{U}_i - \frac{\mathbf{u}'_i\mathbf{t}_i}{\mathbf{t}'_i\mathbf{t}_i}\mathbf{t}_i$ e $\mathbf{u}_{(i+1)}$ são os resíduos da regressão entre \mathbf{U}_i e \mathbf{T}_i . Assim, a j -ésima regressão entre \mathbf{U}_{i+1} e $\mathbf{V}_{(i+1)j}$ resulta num coeficiente de regressão dado por:

$$b_{(i+1)j} = \frac{\mathbf{u}'_{(i+1)}\mathbf{v}_{(i+1)j}}{\mathbf{v}'_{(i+1)j}\mathbf{v}_{(i+1)j}}.$$

Analogamente à equação 7, define-se $\mathbf{T}_{(i+1)}$ como sendo:

$$\mathbf{T}_{(i+1)} = \sum_{j=1}^p w_{(i+1)j} b_{(i+1)j} \mathbf{V}_{(i+1)j}.$$

Os vetores $\mathbf{T}_1, \mathbf{T}_2, \dots, \mathbf{T}_m$ são ortogonais e, portanto, como vetores aleatórios, são não correlacionados (no caso em que $\text{cov}(\mathbf{Y}) = \sigma^2\mathbf{I}$).

Construídos os vetores $\mathbf{T}_1, \mathbf{T}_2, \dots, \mathbf{T}_m$ fica definido na $\text{Im}\mathbf{X}$, um subespaço

ço m -dimensional, $\text{Im}\mathbf{T}$. Projeta-se ortogonalmente o vetor \mathbf{Y} no subespaço $\text{Im}\mathbf{T}$ obtendo-se $\hat{\mathbf{Y}}_{\text{PLS}} = P_{\text{Im}\mathbf{T}}\mathbf{Y}$. O processo de estimação do $\hat{\boldsymbol{\beta}}_{\text{PLS}}$ é o usual utilizando o método dos quadrados mínimos, isto é, $\hat{\boldsymbol{\beta}}_{\text{PLS}}$ é obtido como uma solução particular, bem determinada, das equações normais $\mathbf{X}\hat{\boldsymbol{\beta}}_{\text{PLS}} = P_{\text{Im}\mathbf{T}}\mathbf{Y}$. Uma vez obtida a estimativa $\hat{\boldsymbol{\beta}}_{\text{PLS}}$, a equação de predição fica da forma:

$$\mathbf{y} - \bar{y} = \hat{\boldsymbol{\beta}}'_{\text{PLS}} \begin{bmatrix} t_1 \\ \vdots \\ t_m \end{bmatrix}.$$

Para obter a expressão matricial das equações normais que definem o $\hat{\boldsymbol{\beta}}_{\text{PLS}}$ e também obter uma relação linear entre $\hat{\boldsymbol{\beta}}_{\text{PLS}}$ e $\hat{\boldsymbol{\beta}}_{\text{OLS}}$, foi utilizada uma formalização do método em termos de matrizes, isto é, transformações lineares.

Suponha que m componentes t_1, t_2, \dots, t_m tenham sido construídos. Considere a matriz $\mathbf{T}_{n \times m}$, cujas colunas são definidas pelos vetores \mathbf{t}_i , $\mathbf{T} = [\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_m]$. \mathbf{T} define uma transformação linear do \mathbb{R}^m para o \mathbb{R}^n , tal que $\mathbf{T}\mathbf{e}_i = \mathbf{t}_i$, conforme Figura 33.

A imagem de \mathbf{T} , $\text{Im}\mathbf{T}$, é um subespaço da imagem de \mathbf{X} gerado pelos vetores $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_m$. Portanto, a projeção ortogonal sobre a $\text{Im}\mathbf{T}$, é dada pela matriz $\mathbf{T}(\mathbf{T}'\mathbf{T})^{-1}\mathbf{T}'$, e a projeção do vetor de dados \mathbf{Y} na $\text{Im}\mathbf{T}$ é dada por:

$$\hat{\mathbf{y}}_{\text{PLS}}^m = \mathbf{T}_m(\mathbf{T}'_m \mathbf{T}_m)^{-1} \mathbf{T}'_m \mathbf{y}, \quad (8)$$

conforme Figura 34.

Seja o vetor \mathbf{r}_i definido como uma solução qualquer da equação $\mathbf{X}\mathbf{r}_i = \mathbf{t}_i$, $i=1, \dots, m$. Da mesma forma, pode-se definir a matriz $\mathbf{R}_{p \times m}$, cujas colunas são os vetores \mathbf{r}_i , $i=1, \dots, m$, $\mathbf{R}_m = [\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_m]$. Assim, $\mathbf{R}\mathbf{e}_i = \mathbf{r}_i$.

Como $\mathbf{T}\mathbf{e}_i = \mathbf{t}_i$ e $\mathbf{X}\mathbf{R}\mathbf{e}_i = \mathbf{t}_i$, tem-se que $\mathbf{T}\mathbf{e}_i = \mathbf{X}\mathbf{R}\mathbf{e}_i$. Logo, $\mathbf{T} = \mathbf{X}\mathbf{R}$

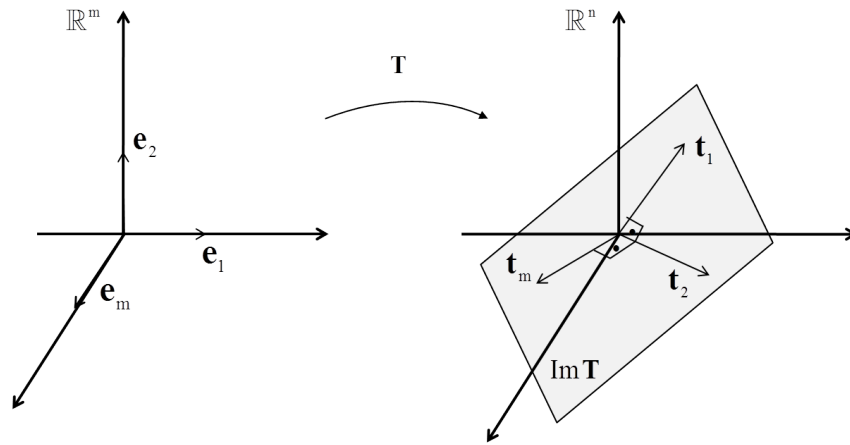


Figura 33 Representação geométrica dos componentes como imagem da transformação linear \mathbf{T} do \mathbb{R}^m para o \mathbb{R}^n

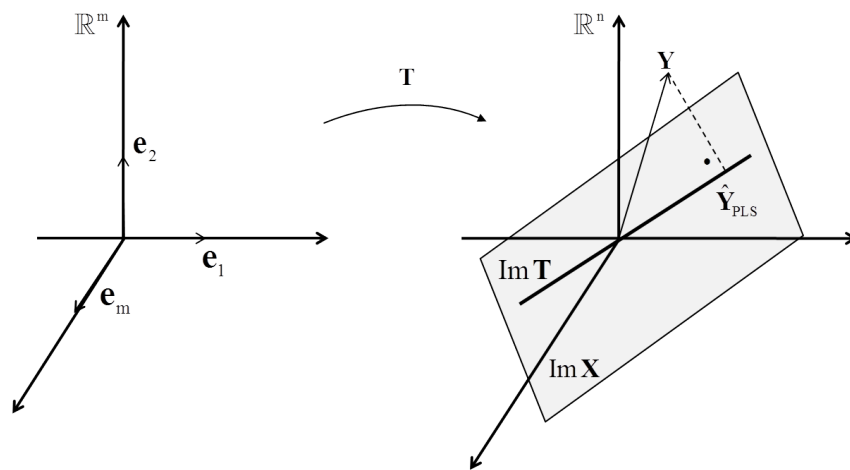


Figura 34 Projeção do vetor de dados \mathbf{Y} no subespaço $\text{Im}\mathbf{T}$ gerado pelos componentes

e obtém-se o diagrama comutativo da Figura 35.

A matriz \mathbf{T} é claramente uma aplicação linear injetiva, pois, suas colunas são linearmente independentes. Assim, restrita à imagem de \mathbf{R} , segue que \mathbf{X} também é uma aplicação injetiva. Vetorialmente, $\text{Im}\mathbf{R} \cap \text{Ker}\mathbf{X} = \{\mathbf{0}\}$.

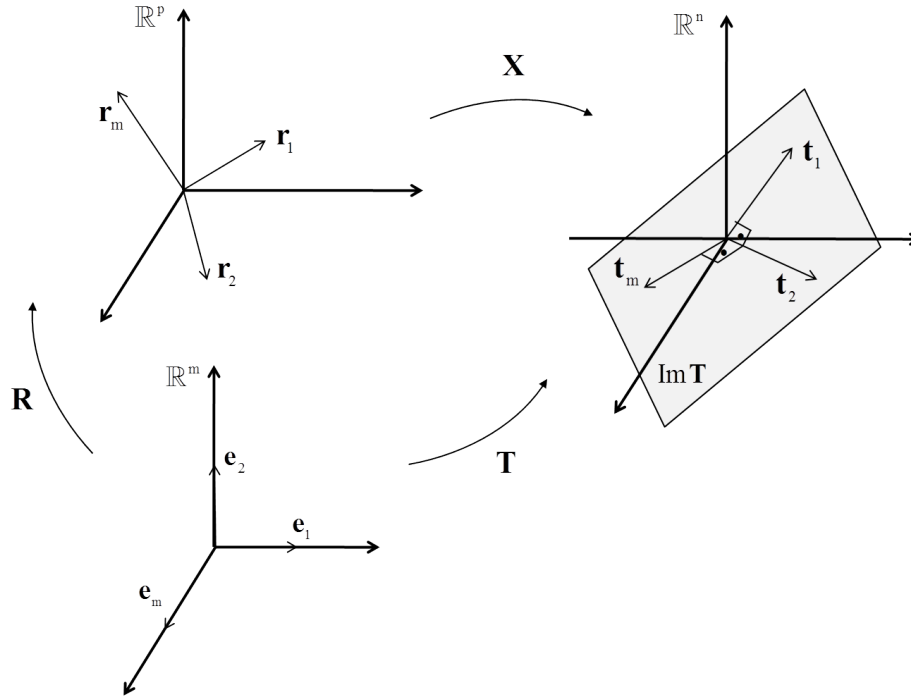


Figura 35 Representação geométrica das transformações lineares \mathbf{T} , \mathbf{R} e \mathbf{X}

Substituindo \mathbf{T} por \mathbf{XR} e $\mathbf{X}'\mathbf{y}$ por $\mathbf{X}'\mathbf{X}\hat{\beta}_{OLS}$, na equação (8), obtemos:

$$\hat{\mathbf{y}}_{PLS}^m = \mathbf{XR}_m \left(\mathbf{R}_m' \mathbf{X}' \mathbf{X} \mathbf{R}_m \right)^{-1} \mathbf{R}_m' \mathbf{X}' \mathbf{X} \hat{\beta}_{OLS}.$$

Observe que a substituição de $\mathbf{X}'\mathbf{y}$ por $\mathbf{X}'\mathbf{X}\hat{\beta}_{OLS}$ é verdadeira em relação a qualquer $\hat{\beta}_{OLS}$ escolhido.

Utilizando o fato de \mathbf{X} ser injetiva quando restrita à imagem de \mathbf{R} , temos que $\hat{\beta}_{PLS}^m$ está bem definido e é dado por:

$$\hat{\beta}_{PLS}^m = \mathbf{R}_m \left(\mathbf{R}_m' \mathbf{X}' \mathbf{X} \mathbf{R}_m \right)^{-1} \mathbf{R}_m' \mathbf{X}' \mathbf{X} \hat{\beta}_{OLS}. \quad (9)$$

A escolha particular dos vetores \mathbf{r}_i , não afeta a estimativa de $\hat{\beta}_{PLS}^m$, pois

se \mathbf{A} é uma transformação linear inversível, $\mathbf{A} : \mathbb{R}^m \rightarrow \mathbb{R}^m$, a mudança de coordenadas $\tilde{\mathbf{R}} = \mathbf{R}\mathbf{A}$ é tal que, substituindo na equação (9), é obtida a mesma estimativa:

$$\begin{aligned}\hat{\beta}_{\text{PLS}}^m &= \tilde{\mathbf{R}}_m \left(\tilde{\mathbf{R}}_m' \mathbf{X}' \mathbf{X} \tilde{\mathbf{R}}_m \right)^{-1} \tilde{\mathbf{R}}_m' \mathbf{X}' \mathbf{X} \hat{\beta}_{\text{OLS}} \\ &= \mathbf{R}_m \mathbf{A} \left(\mathbf{A}' \mathbf{R}_m' \mathbf{X}' \mathbf{X} \mathbf{R}_m \mathbf{A} \right)^{-1} \mathbf{A}' \mathbf{R}_m' \mathbf{X}' \mathbf{X} \hat{\beta}_{\text{OLS}} \\ &= \mathbf{R}_m \mathbf{A} (\mathbf{A}')^{-1} \left(\mathbf{R}_m' \mathbf{X}' \mathbf{X} \mathbf{R}_m \right)^{-1} (\mathbf{A})^{-1} \mathbf{A}' \mathbf{R}_m' \mathbf{X}' \mathbf{X} \hat{\beta}_{\text{OLS}} \\ &= \mathbf{R}_m \left(\mathbf{R}_m' \mathbf{X}' \mathbf{X} \mathbf{R}_m \right)^{-1} \mathbf{R}_m' \mathbf{X}' \mathbf{X} \hat{\beta}_{\text{OLS}}.\end{aligned}$$

Este fato reflete que apenas o subespaço gerado pelos vetores \mathbf{r}_i é importante, e não uma escolha particular da base.

É possível obter uma relação linear entre $\hat{\beta}_{\text{PLS}}$ e $\hat{\beta}_{\text{OLS}}$ que possui uma interpretação geométrica muito interessante. Observe que:

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

Logo, essa transformação é um projetor. Como ela não é simétrica, esse projetor não é um projetor ortogonal e, sim, um projetor oblíquo, portanto, como todo projetor, a projeção se dá ao longo do kernel (Figura 36).

Como as matrizes \mathbf{R} e $\mathbf{X}'\mathbf{X}$ são injetivas, $\mathbf{R}_m' \mathbf{X}' \mathbf{X} \mathbf{R}_m$ é inversível e essas transformações não têm kernel, portanto, o kernel da transformação

$$\mathbf{R}_m \left(\mathbf{R}_m' \mathbf{X}' \mathbf{X} \mathbf{R}_m \right)^{-1} \mathbf{R}_m' \mathbf{X}' \mathbf{X},$$

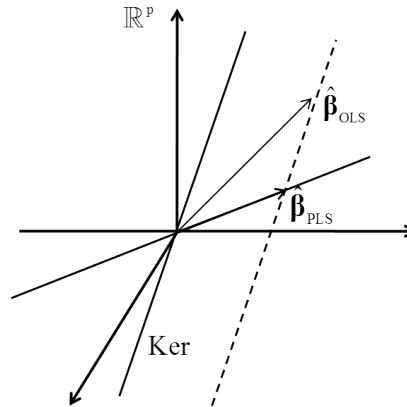


Figura 36 Vetor $\hat{\beta}_{PLS}$ como projeção oblíqua do vetor $\hat{\beta}_{OLS}$

é exatamente o kernel de \mathbf{R}_m' . Note também que, o kernel de \mathbf{X} está contido no kernel de $\mathbf{R}_m' \mathbf{X}' \mathbf{X}$.

4.4 Determinação do número ótimo de componentes (variáveis latentes)

Um fator fundamental na utilização do PLS é a escolha do número ótimo de componentes a serem utilizados no modelo. Não existe uma regra formal para determinar esse número, porém, na literatura, um critério de decisão utilizado para a metodologia PCR é adotar uma porcentagem da variação total explicada que se deseja obter e usar o número de componentes que atinja esse valor pré-determinado. Essa abordagem empírica vem sendo usada na prática, geralmente fixando uma porcentagem de 70% (AZEVEDO, 2012). Porém, por se tratar de um enfoque empírico, alternativas provenientes de derivações teóricas ou de recursos computacionais intensivos são mais adequadas, como por exemplo as teorias Graus de Liberdade e Validação Cruzada.

4.4.1 Graus de Liberdade e seleção de modelos

Esta seção tem como referência o artigo Krämer e Sugiyama (2011).

Em regressão, o Grau de Liberdade quantifica a complexidade intrínseca de um modelo. Para modelos de Regressão Linear, em que os valores ajustados são funções lineares dos dados, \mathbf{y} , isto é, $\hat{\mathbf{y}}_\lambda = \mathbf{H}_\lambda \mathbf{y}$ com $\mathbf{H}_\lambda \in \mathbb{R}^{n \times n}$, em que \mathbf{H} é uma matriz de projeção (*hat-matrix*), o grau de liberdade é definido como o traço da matriz de projeção:

$$\text{DoF}(\lambda) = \text{tr}(\mathbf{H}_\lambda). \quad (10)$$

Observe que o traço de \mathbf{H} é igual à dimensão do subespaço definido pelo modelo (em geral, o subespaço gerado pelos vetores colunas da matriz de covariáveis). No caso em que a matriz de covariáveis é de posto completo, o DoF é exatamente igual ao número de parâmetros. Como exemplo, essa definição aplica-se à Regressão de Componentes Principais.

Na regressão PLS, tem-se, conforme a equação (8), que

$$\hat{\mathbf{y}}_{\text{PLS}}^m = \bar{\mathbf{y}} + \mathbf{T}_m (\mathbf{T}_m' \mathbf{T}_m)^{-1} \mathbf{T}_m' \mathbf{y} = \bar{\mathbf{y}} + \mathbf{P}_{\mathbf{T}} \mathbf{y}.$$

Como as componentes dependem não somente das covariáveis, mas também dos dados, a regressão PLS não é linear. Portanto, a definição (10) não pode ser aplicada. É necessário, então, usar uma generalização proposta por Efron (2004).

Definição 4.1. Faça \hat{f}_λ ser uma estimativa da verdadeira função de regressão f , parametrizada por λ . Defina o vetor de valores ajustados como $\hat{\mathbf{y}}_\lambda = \left(\hat{f}_\lambda(\mathbf{x}_1), \dots, \hat{f}_\lambda(\mathbf{x}_n) \right)'$. Os graus de liberdade são

$$\text{DoF}(\lambda) = E \left(\text{tr} \left(\frac{\partial \hat{\mathbf{y}}_\lambda}{\partial \mathbf{y}} \right) \right).$$

Assumindo \mathbf{X} fixada e a esperança E tomada com relação a y_1, \dots, y_n .

Neste sentido, o grau de liberdade é uma medida da sensibilidade dos valores ajustados como função dos dados.

Para o caso PLS, o DoF fica definido como:

$$\text{DoF}(m) = 1 + E \left(\text{tr} \left(\frac{\partial \mathbf{P}_{\mathbf{T}\mathbf{Y}}}{\partial \mathbf{y}} \right) \right).$$

A constante 1 corresponde à estimação do intercepto, uma vez que a matriz \mathbf{X} foi considerada centrada, e m é o número de componentes utilizados no modelo.

Um estimador não viesado para o Grau de Liberdade da regressão PLS com m componentes latentes, $\mathbf{T} = (\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_m)$, é dado por:

$$\hat{\text{DoF}}(m) = 1 + \text{tr} \left(\frac{\partial \mathbf{P}_{\mathbf{T}\mathbf{Y}}}{\partial \mathbf{y}} \right).$$

Uma conjectura importante, afirma que $\text{DoF}(m) \geq m + 1$ (FRANK; FRIEDMAN, 1993; MARTENS; NAES, 1989). Para exemplificar, será calculado explicitamente uma cota inferior para o DoF quando $m = 1$.

Seja $\mathbf{S} = \frac{1}{n-1} \mathbf{X}'\mathbf{X} \in \mathbb{R}^{p \times p}$ a matriz de covariância amostral do grupo \mathbf{X} , e $\mathbf{s} = \frac{1}{n-1} \mathbf{X}'\mathbf{y} \in \mathbb{R}^p$, a covariância amostral entre os grupos \mathbf{X} e \mathbf{Y} .

Teorema 2. Se o maior autovalor λ_{\max} da matriz de covariância amostral \mathbf{S} satisfaz

$$\lambda_{\max} \leq \frac{1}{2} \text{tr}(\mathbf{S}),$$

então

$$\text{DôF}(m=1) \geq 1 + \frac{\text{tr}(\mathbf{S})}{\lambda_{\max}}.$$

Demonstração. Sabendo que o primeiro componente é definido por $\mathbf{t}_1 = \mathbf{X}\mathbf{s}$, tem-se:

$$\frac{\mathbf{t}_1}{\|\mathbf{t}_1\|} = \frac{\mathbf{t}_1}{\sqrt{\mathbf{t}_1' \mathbf{t}_1}} = \frac{\mathbf{X}\mathbf{s}}{\sqrt{(\mathbf{X}\mathbf{s})' \mathbf{X}\mathbf{s}}} = \frac{\mathbf{X}\mathbf{s}}{\sqrt{\mathbf{s}' \mathbf{X}' \mathbf{X} \mathbf{s}}} = \frac{\mathbf{X}\mathbf{s}}{\sqrt{\mathbf{s}' \mathbf{S} \mathbf{s}}},$$

de onde segue que a projeção de \mathbf{y} no primeiro componente \mathbf{t}_1 , é dada por:

$$\begin{aligned} \left(\mathbf{y} \cdot \frac{\mathbf{t}_1}{\|\mathbf{t}_1\|} \right) \frac{\mathbf{t}_1}{\|\mathbf{t}_1\|} &= \left(\mathbf{y} \cdot \frac{\mathbf{X}\mathbf{s}}{\sqrt{\mathbf{s}' \mathbf{S} \mathbf{s}}} \right) \frac{\mathbf{X}\mathbf{s}}{\sqrt{\mathbf{s}' \mathbf{S} \mathbf{s}}} = \frac{(\mathbf{y}' \mathbf{X}) \mathbf{s}}{\sqrt{\mathbf{s}' \mathbf{S} \mathbf{s}}} \frac{\mathbf{X}\mathbf{s}}{\sqrt{\mathbf{s}' \mathbf{S} \mathbf{s}}} \\ &= \frac{\mathbf{s}' \mathbf{s}}{\sqrt{\mathbf{s}' \mathbf{S} \mathbf{s}}} \frac{\mathbf{X}\mathbf{s}}{\sqrt{\mathbf{s}' \mathbf{S} \mathbf{s}}} = \frac{\mathbf{s}' \mathbf{s}}{\mathbf{s}' \mathbf{S} \mathbf{s}} \mathbf{X}\mathbf{s}. \end{aligned}$$

Assim,

$$\hat{\mathbf{y}}_1 = \bar{\mathbf{y}} + \frac{\mathbf{s}' \mathbf{s}}{\mathbf{s}' \mathbf{S} \mathbf{s}} \mathbf{X}\mathbf{s}. \quad (11)$$

Calculando a derivada de (11) em relação a \mathbf{y} e sendo a derivada de uma forma quadrática dada por $\frac{\partial}{\partial \mathbf{y}} (\mathbf{y}' \mathbf{A} \mathbf{y}) = 2 (\mathbf{A} \mathbf{y})'$, tem-se:

$$\frac{\partial}{\partial \mathbf{y}} (\bar{\mathbf{y}}) = \begin{bmatrix} 1/n & \dots & 1/n \\ \vdots & \ddots & \vdots \\ 1/n & \dots & 1/n \end{bmatrix} \Rightarrow \text{tr} \left(\frac{\partial}{\partial \mathbf{y}} (\bar{\mathbf{y}}) \right) = \frac{1}{n} n = 1, \quad (12)$$

e,

$$\begin{aligned}
\frac{\partial P_{t_1} \mathbf{y}}{\partial \mathbf{y}} &= \frac{\partial}{\partial \mathbf{y}} \left(\frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \mathbf{X}\mathbf{s} \right) \\
&= \frac{\partial}{\partial \mathbf{y}} \left(\frac{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}} \mathbf{X}\mathbf{X}'\mathbf{y} \right) \\
&= \left(\frac{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}} \right) \frac{\partial}{\partial \mathbf{y}} (\mathbf{X}\mathbf{X}'\mathbf{y}) + \frac{\partial}{\partial \mathbf{y}} \left(\frac{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}} \right) (\mathbf{X}\mathbf{X}'\mathbf{y}) \\
&= \left(\frac{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}} \right) \mathbf{X}\mathbf{X}' + \left(\frac{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})(2\mathbf{X}\mathbf{X}'\mathbf{y})' - (\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y})(2\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})'}{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})^2} \right) \mathbf{X}\mathbf{X}'\mathbf{y} \\
&= \left(\frac{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}} \right) \mathbf{X}\mathbf{X}' + \left(\frac{2(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})\mathbf{y}'\mathbf{X}\mathbf{X}' - 2(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y})\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'}{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})^2} \right) \mathbf{X}\mathbf{X}'\mathbf{y} \\
&= \left(\frac{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}} \right) \mathbf{X}\mathbf{X}' + 2 \left(\frac{\mathbf{y}'\mathbf{X}\mathbf{X}'}{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})} - \frac{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y})\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'}{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})^2} \right) \mathbf{X}\mathbf{X}'\mathbf{y} \\
&= \left(\frac{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}} \right) \mathbf{X}\mathbf{X}' + 2 \left(\frac{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}}{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})} - \frac{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y})\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}}{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})^2} \right) \\
&= \left(\frac{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}} \right) \mathbf{X}\mathbf{X}' + 2 \left(1 - \frac{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{y})}{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})} \frac{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})}{(\mathbf{y}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y})} \right) \\
&= \frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \mathbf{X}\mathbf{X}' + 2 \left(1 - \frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \frac{\mathbf{s}'\mathbf{S}^2\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \right).
\end{aligned}$$

Usando $tr(\mathbf{AB}) = tr(\mathbf{BA})$, segue que:

$$\begin{aligned}
\text{DôF} (m = 1) &= 1 + tr \left(\frac{\partial P_{\mathbf{T}\mathbf{Y}}}{\partial \mathbf{y}} \right) \\
&= 1 + tr \left(\frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \mathbf{X}\mathbf{X}' + 2 \left(1 - \frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \frac{\mathbf{s}'\mathbf{S}^2\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \right) \right) \\
&= 1 + \frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} tr(\mathbf{X}'\mathbf{X}) + 2 - 2 \frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \frac{\mathbf{s}'\mathbf{S}^2\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \\
&= 3 + \frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \left(tr(\mathbf{S}) - 2 \frac{\mathbf{s}'\mathbf{S}^2\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \right).
\end{aligned}$$

Note que a constante 1 da fórmula é exatamente a derivada do vetor $\bar{\mathbf{y}}$, calculada em (12).

Afirmção: $\frac{\mathbf{s}'\mathbf{S}^2\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \leq \lambda_{\max}$.

Seja \mathbf{A} uma matriz ortogonal, $\mathbf{A}'\mathbf{A} = \mathbf{A}\mathbf{A}' = \mathbf{I}$, que diagonaliza \mathbf{S} :

$$\mathbf{A}\mathbf{S}\mathbf{A}' = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix} = \mathbf{\Lambda}.$$

Tem-se que \mathbf{A} também diagonaliza \mathbf{S}^2 . Note que

$$\mathbf{S}\mathbf{v} = \lambda\mathbf{v} \Rightarrow \mathbf{S}\mathbf{S}\mathbf{v} = \lambda\mathbf{S}\mathbf{v} \Rightarrow \mathbf{S}^2\mathbf{v} = \lambda^2\mathbf{v},$$

logo,

$$\begin{aligned} \mathbf{A}\mathbf{S}^2\mathbf{A}' &= \mathbf{A}\mathbf{S}\mathbf{A}'\mathbf{A}\mathbf{S}\mathbf{A}' = \mathbf{\Lambda}\mathbf{\Lambda} \\ &= \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix} \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix} \\ &= \begin{pmatrix} \lambda_1^2 & & \\ & \ddots & \\ & & \lambda_p^2 \end{pmatrix}. \end{aligned}$$

Fazendo $\mathbf{s} = \mathbf{A}'\mathbf{w}$, em que $\mathbf{w} = \mathbf{A}\mathbf{s}$, tem-se:

$$\begin{aligned} \frac{\mathbf{s}'\mathbf{S}^2\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} &= \frac{\mathbf{w}'\mathbf{A}\mathbf{S}^2\mathbf{A}'\mathbf{w}}{\mathbf{w}'\mathbf{A}\mathbf{S}\mathbf{A}'\mathbf{w}} = \frac{\mathbf{w}'\mathbf{\Lambda}^2\mathbf{w}}{\mathbf{w}'\mathbf{\Lambda}\mathbf{w}} \\ &= \frac{\sum \lambda_i^2 \mathbf{w}_i^2}{\sum \lambda_i \mathbf{w}_i^2} = \frac{\lambda_{\max}^2 \sum \left(\frac{\lambda_i}{\lambda_{\max}}\right)^2 \mathbf{w}_i^2}{\lambda_{\max} \sum \frac{\lambda_i}{\lambda_{\max}} \mathbf{w}_i^2}. \end{aligned}$$

Como $\frac{\lambda_i}{\lambda_{\max}} \leq 1 \Rightarrow \left(\frac{\lambda_i}{\lambda_{\max}}\right)^2 \leq \frac{\lambda_i}{\lambda_{\max}}$, segue que

$$\frac{\mathbf{s}'\mathbf{S}^2\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \leq \lambda_{\max}.$$

Portanto,

$$\text{tr}(\mathbf{S}) - 2\frac{\mathbf{s}'\mathbf{S}^2\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \geq \text{tr}(\mathbf{S}) - 2\lambda_{\max}.$$

Além disso, $\frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \geq \frac{1}{\lambda_{\max}}$.

Novamente, fazendo $\mathbf{s} = \mathbf{A}'\mathbf{w}$,

$$\begin{aligned} \frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} &= \frac{\mathbf{w}'\mathbf{A}\mathbf{A}'\mathbf{w}}{\mathbf{w}'\mathbf{A}\mathbf{S}\mathbf{A}'\mathbf{w}} = \frac{\mathbf{w}'\mathbf{w}}{\mathbf{w}'\mathbf{\Lambda}\mathbf{w}} \\ &= \frac{\sum \mathbf{w}_i^2}{\sum \lambda_i \mathbf{w}_i^2} = \frac{1}{\lambda_{\max}} \frac{\sum \mathbf{w}_i^2}{\sum \frac{\lambda_i}{\lambda_{\max}} \mathbf{w}_i^2}. \end{aligned}$$

Como $\frac{\lambda_i}{\lambda_{\max}} \leq 1$, segue que:

$$\frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} \geq \frac{1}{\lambda_{\max}}.$$

Assim,

$$\begin{aligned} \hat{\text{DoF}}(m=1) &\geq 3 + \frac{\mathbf{s}'\mathbf{s}}{\mathbf{s}'\mathbf{S}\mathbf{s}} (\text{tr}(\mathbf{S}) - 2\lambda_{\max}) \\ &\geq 3 + \frac{1}{\lambda_{\max}} (\text{tr}(\mathbf{S}) - 2\lambda_{\max}) \\ &= 1 + \frac{\text{tr}(\mathbf{S})}{\lambda_{\max}}. \end{aligned}$$

□

Observe que $\text{tr}(\mathbf{S})$ é igual a soma dos autovalores e, portanto, $\frac{\text{tr}(\mathbf{S})}{\lambda_{\max}} > 1$. Logo, $\hat{\text{DoF}}(m=1) \geq 2$ e, $\text{DoF}(m=1) \geq 2$, mostrando que a conjectura é

verdadeira para $m = 1$.

Pode-se utilizar o conceito de Grau de Liberdade para determinar um número adequado de componentes a ser utilizado na regressão, determinando o DoF para cada valor de m . Se para um acréscimo em m tem-se um acréscimo considerável no valor do DoF, esse valor de m é justificado. Porém, se para um acréscimo em m , o acréscimo no valor do DoF é pequeno, o princípio da parcimônia recomenda que esse acréscimo na complexidade do modelo é desnecessário. Traçando-se o gráfico DoF em função do número de componentes, o número adequado de componentes ocorre quando a curva tende a se estabilizar em um valor constante, como ilustrado na Figura 37.

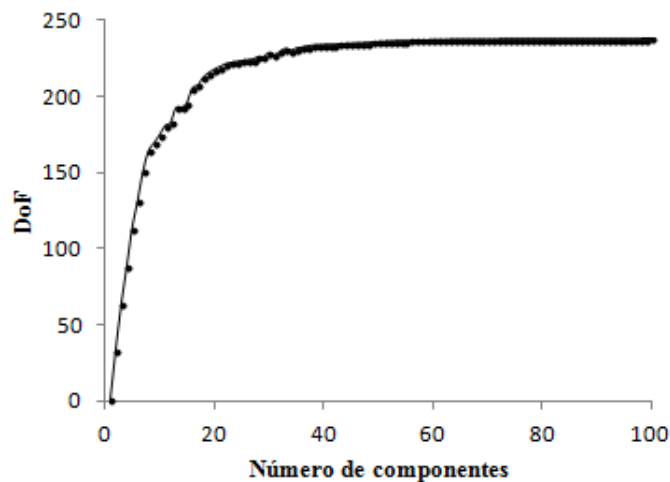


Figura 37 Representação de uma curva gerada por valores dos graus liberdade, DoF, em função do número de componentes m

4.4.2 Validação Cruzada

Uma forma mais prática de se determinar o número de componentes, é por meio da validação cruzada (PÉREZ-CABAL *et al.*, 2012). Este método consiste em dividir o conjunto de dados originais em N subconjuntos (*folds*), e realizar N análises, de forma que em cada uma delas um dos subconjuntos é retirado a fim de ser utilizado como recurso para validar a análise realizada. Assim, os valores preditos pelas equações estimadas em cada análise podem ser diretamente comparados com os valores observados que foram retirados. No contexto de PLSR, o número ótimo de componentes será aquele que fornecer o menor erro quadrático entre o conjunto de valores observados e preditos, como ilustra a Figura 38.

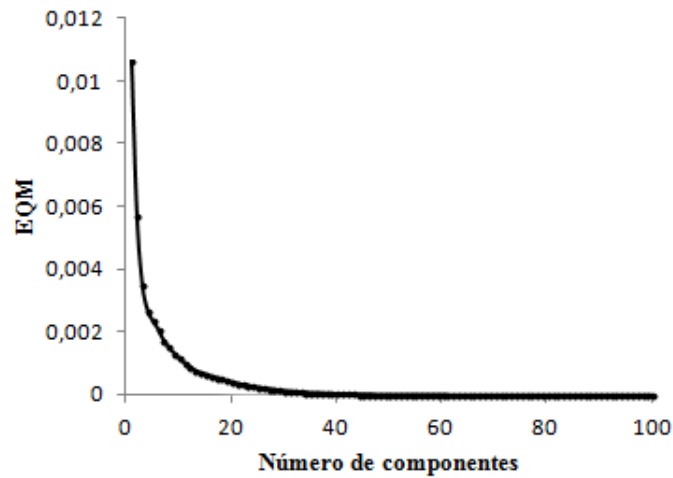


Figura 38 Representação de uma curva gerada por valores de erro quadrático médio em função do número de componentes m

4.5 Exemplo

Matriz de covariâncias simulada conforme a seção 3.5:

```
6.2810455  8.194268  0.3211009  5.693020  20.811639
8.1942679  10.856983  1.0961611  7.626992  27.924481
0.3211009  1.096161  4.5260950  7.049471  6.085769
5.6930197  7.626992  7.0494713  25.766833  26.599556
20.8116386  27.924481  6.0857694  26.599556  77.622607
```

Autovalores e autovetores da matriz de covariâncias:

\$values

```
1.043224e+02  1.758224e+01  2.610955e+00  5.375069e-01
4.210175e-04
```

\$vectors

```
-0.227191  0.178466 -0.292529 -0.457806  0.788272
-0.304403  0.223537 -0.140214 -0.695880 -0.594524
-0.080458 -0.327128  0.851982 -0.371038  0.151556
-0.342879 -0.865085 -0.362872 -0.015131 -0.046416
-0.855384  0.250587  0.192913  0.410200  0.006555
```

Matriz de correlação:

```
1.00000000 0.9922922 0.06022318 0.4475032 0.9425322
0.99229225 1.0000000 0.15637163 0.4560035 0.9619144
0.06022318 0.1563716 1.00000000 0.6527760 0.3246830
0.44750325 0.4560035 0.65277602 1.0000000 0.5947710
0.94253219 0.9619144 0.32468300 0.5947710 1.0000000
```

Matriz de covariáveis $X_{6 \times 4}$:

```
0.9391754 2.0814841 4.776478 8.597991
2.1972739 4.1631504 5.466290 6.725786
3.5575236 5.4117197 3.282917 6.600031
-0.8903883 0.1328435 5.472369 4.877636
-0.5096611 0.2326990 3.388711 0.957349
1.9463446 3.6372080 5.319583 7.467699
```

Vetor da variável resposta $Y_{6 \times 1}$:

```
7.820750
11.335945
16.909518
3.731732
-1.821228
9.446128
```

Vetor de parâmetros estimado via OLS, $\hat{\beta}_{OLS}$:

-42.710500
33.456902
-9.939255
3.077167

Vetor de respostas estimado via OLS, \hat{Y}_{OLS} :

8.510179
11.804963
16.795406
3.091492
-1.182066
8.667023

Utilizando 1 componente: $m = 1$

Vetor dos componentes:

9.723363
9.813350
9.514253
6.064935
2.297160
10.018880

Vetor de parâmetros estimado via PLS, $\hat{\beta}_{PLS}$, em função dos componentes:

1.088373

Vetor de parâmetros estimado via PLS, $\hat{\beta}_{PLS}$, em função das covariáveis originais:

0.2557785

0.4469298

0.5175329

0.8071838

Vetor de respostas estimado via PLS, \hat{Y}_{PLS} , em função dos componentes:

10.582641

10.680581

10.355052

6.600909

2.500166

10.904274

Vetor de respostas estimado via PLS, \hat{Y}_{PLS} , em função das covariáveis originais:

10.582641

10.680581

10.355052

6.600909

2.500166

10.904274

Utilizando 2 componentes: $m = 2$

Matriz dos componentes:

```
9.723363 -1.1645658
9.813350  0.5490979
9.514253  3.2962070
6.064935 -3.6298704
2.297160 -2.0393287
10.018880  0.1271250
```

Vetor de parâmetros estimado via PLS, $\hat{\beta}_{\text{PLS}}$, em função dos componentes:

```
1.088373
1.480342
```

Vetor de parâmetros estimado via PLS, $\hat{\beta}_{\text{PLS}}$, em função das covariáveis originais:

```
1.0619346
1.3092326
-0.3592859
0.6364624
```

Vetor de respostas estimado via PLS, \hat{Y}_{PLS} , em função dos componentes:

8.8586862
11.4934339
15.2345645
1.2274611
-0.5187373
11.0924624

Vetor de respostas estimado via PLS, \hat{Y}_{PLS} , em função das covariáveis originais:

7.4786659
10.1006417
13.8842227
0.3666752
-0.8447693
9.6704998

Utilizando 3 componentes: $m = 3$

Matriz dos componentes:

9.723363	-1.1645658	-1.52574743
9.813350	0.5490979	0.79694695
9.514253	3.2962070	0.28010446
6.064935	-3.6298704	0.13097361
2.297160	-2.0393287	1.30889216
10.018880	0.1271250	0.05475768

Vetor de parâmetros estimado via PLS, $\hat{\beta}_{\text{PLS}}$, em função dos componentes:

1.08837257
1.48034155
0.09640233

Vetor de parâmetros estimado via PLS, $\hat{\beta}_{\text{PLS}}$, em função das covariáveis originais:

1.0341975
1.3699283
-0.3146987
0.5830576

Vetor de respostas estimado via PLS, \hat{Y}_{PLS} , em função dos componentes:

8.711601
11.570261
15.261567
1.240087
-0.392557
11.097741

Vetor de respostas estimado via PLS, \hat{Y}_{PLS} , em função das covariáveis originais:

7.3327488
10.1769183
13.9079179
0.3829437
-0.7165427
9.6756510

Utilizando 4 componentes: $m = 4$

Matriz dos componentes:

9.723363	-1.1645658	-1.52574743	-0.004186004
9.813350	0.5490979	0.79694695	0.004877645
9.514253	3.2962070	0.28010446	0.031876755
6.064935	-3.6298704	0.13097361	0.038476500
2.297160	-2.0393287	1.30889216	-0.016407845
10.018880	0.1271250	0.05475768	-0.050515983

Vetor de parâmetros estimado via PLS, $\hat{\beta}_{\text{PLS}}$, em função dos componentes:

1.08837257
1.48034155
0.09640233
48.11779953

Vetor de parâmetros estimado via PLS, $\hat{\beta}_{\text{PLS}}$, em função das covariáveis originais:

-35.14784
16.25217
-22.32000
17.91708

Vetor de respostas estimado via PLS, \hat{Y}_{PLS} , em função dos componentes:

8.510179
11.804963
16.795406
3.091492
-1.182066
8.667023

Vetor de respostas estimado via PLS, \hat{Y}_{PLS} , em função das covariáveis originais:

48.258491
-11.070420
7.891424
-1.296101
-36.787790
5.768917

Novos valores de \mathbf{X}_p e Y_p gerados a fim de fazer predição:

X_p

3.289261 5.040320 3.607594 8.119808

Y_p

12.87691

Valores preditos para Y_p (\hat{Y}_p) via OLS e PLS:

OLS 17.23269

PLS 10.81509 - 1 componente

PLS 15.09155 - 2 componentes

PLS 15.01617 - 3 componentes

PLS 37.04379 - 4 componentes

5 CONCLUSÕES

A abordagem geométrica apresentou-se como uma maneira eficiente e intuitiva para explicitar os passos teóricos da teoria e do algoritmo do método dos Quadrados Mínimos Parciais. Além disso, permite uma abordagem unificada das teorias de regressão em Quadrados Mínimos Ordinários, Componentes Principais e Quadrados Mínimos Parciais.

REFERÊNCIAS

- ADÃO, A. S. **Introdução à teoria geométrica dos delineamentos experimentais**. 2011. 78 f. Dissertação (Mestrado em Estatística e Experimentação Agropecuária) - Universidade Federal de Lavras, Lavras, MG, 2011.
- AZEVEDO, C. F. **Métodos de redução de dimensionalidade aplicados na seleção genômica para características de carcaça em suínos**. 2012. 50 f. Dissertação (Mestrado em Estatística Aplicada e Biometria) - Universidade Federal de Viçosa, Viçosa, MG, 2012.
- EFRON, B. The Estimation of Prediction Error: covariance penalties and cross-validation. **Journal of the American Statistical Association**, [s.l.], v. 99, n. 467, p. 619-632, set. 2004.
- FERREIRA, D. F. **Estatística Multivariada**. Lavras: Ed. UFLA, 2008. 662 p.
- FRANK, I.; FRIEDMAN, J. A Statistical View of Some Chemometrics Regression Tools. **Technometrics**, [s.l.], v. 35, n. 2, p. 109-135, maio 1993.
- GARTHWAITE, P. H. An Interpretation of Partial Least Squares. **Journal of the American Statistical Association**, [s.l.], v. 89, n. 425, p. 122-127, mar. 1994.
- GUIMARÃES, P. H. S. **Uma abordagem geométrica da Teoria de Inversas Generalizadas**. 2010. 67 f. Dissertação (Mestrado em Estatística e Experimentação Agropecuária) - Universidade Federal de Lavras, Lavras, MG, 2010.
- HELLAND, I. S. Partial Least Squares Regression and Statistical Models. **Scandinavian Journal of Statistics**, [s.l.], v. 17, n. 2, p. 97-114, 1990.
- HOSKULDSSON, A. PLS Regression Methods. **Journal of Chemometrics**, [s.l.], v. 2, p. 211-228, 1988.

HOTELLING, H. The relations of the newer multivariate statistical methods to factor analysis. **British Journal of Mathematical and Statistical Psychology**, [s.l.], v. 10, n. 2, p. 69-79, nov. 1957.

JONG, S. SIMPLS: an alternative approach to partial least squares regression. **Chemometrics and Intelligent Laboratory Systems**, [s.l.], v. 18, p. 251-263, 1993.

KALMAN, D. A Singularly Valuable Decomposition: the SVD of a Matrix. **The College Mathematics Journal**, [s.l.], v. 27, n. 1, p. 2-23, jan. 1996.

KENDALL, M. G. **A Course in Multivariate Analysis**. London: Griffin, 1957.

KRÄMER, N.; SUGIYAMA, M. The Degrees of Freedom of Partial Least Squares Regression. **Journal of the American Statistical Association**, [s.l.], v. 106, n. 494, p. 697-705, fev. 2011.

MARTENS, H.; NAES, T. **Multivariate Calibration**. New York: Wiley, 1989.

MEUWISSEN, T. H. E.; HAYES, B. J.; GODDARD, M. E. Prediction of total genetic value using genome wide densemarker maps. **Genetics**, [s.l.], v. 157, p. 1819-1829, abr. 2001.

OTTO, M. **Chemometrics**. Weinheim: Wiley, 1999. 328 p.

PEREIRA, L. S. **Abordagem geométrica à teoria dos modelos de Gauss-Markov**. 2013. 130 f. Dissertação (Mestrado em Estatística e Experimentação Agropecuária) - Universidade Federal de Lavras, Lavras, MG, 2013.

PÉREZ-CABAL, M. A. *et al.* Accuracy of genome-enabled prediction in a dairy cattle population using different cross-validation layouts. **Frontiers in Genetics - Livestock Genomics**, [s.l.], v. 3, p. 1-7, fev. 2012.

PHATAK, A.; JONG, S. The geometry of partial least squares. **Journal of the Chemometrics**, [s.l.], v. 11, n. 4, p. 311-338, jul. 1997.

R DEVELOPMENT CORE TEAM. **R**: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, 2013. Disponível em: <<http://www.R-project.org>>. Acesso em: 14 nov. 2013.

RAO, C. R. **Linear Statistical Inference and Its Applications**. New York: Wiley, 2002. 625 p.

ROGGO, Y. *et al.* A review of near infrared spectroscopy and chemometrics in pharmaceutical technologies. **Journal of Pharmaceutical and Biomedical Analysis**, [s.l.], v. 44, n. 3, p. 683-700, jul. 2007.

WOLD, H. Path models with Latent Variables: the NIPALS Approach. In: H. B. *et al.* (Ed.). **Quantitative Sociology**: international Perspectives on Mathematical and Statistical Modeling. [New York]: Academic Press, 1975. p. 307-357.

APÊNDICES

APÊNDICE A - Completamento de Quadrados

Deseja-se o completamento de quadrados em y para o seguinte termo:

$$\mathbf{x}'\mathbf{C}\mathbf{x} - 2\mathbf{x}'\mathbf{C}\mathbf{H}\mathbf{y} + \mathbf{y}'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})\mathbf{y}.$$

Ou seja, deseja-se determinar w tal que:

$$(\mathbf{y} - w)'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})(\mathbf{y} - w) + g(\mathbf{x}) = \mathbf{x}'\mathbf{C}\mathbf{x} - 2\mathbf{x}'\mathbf{C}\mathbf{H}\mathbf{y} + \mathbf{y}'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})\mathbf{y},$$

em que $g(\mathbf{x})$ é uma função que não depende de y .

Tem-se que:

$$\begin{aligned} & (\mathbf{y} - w)'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})(\mathbf{y} - w) \\ &= \mathbf{y}'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})\mathbf{y} - 2w'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})\mathbf{y} + w'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})w. \end{aligned}$$

Observe que, $\mathbf{x}'\mathbf{C}\mathbf{H}\mathbf{y} = w'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})\mathbf{y}$. Logo,

$$w = (\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})^{-1}\mathbf{H}'\mathbf{C}\mathbf{x}.$$

Portanto,

$$\begin{aligned} & \mathbf{x}'\mathbf{C}\mathbf{x} - 2\mathbf{x}'\mathbf{C}\mathbf{H}\mathbf{y} + \mathbf{y}'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})\mathbf{y} \\ &= \mathbf{x}'\mathbf{C}\mathbf{x} + \mathbf{y}'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})\mathbf{y} - 2w'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})\mathbf{y} \\ &= \mathbf{x}'\mathbf{C}\mathbf{x} + \mathbf{y}'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})\mathbf{y} - 2w'(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})\mathbf{y} \end{aligned}$$

$$\begin{aligned}
& + \mathbf{w}' (\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B}) \mathbf{w} - \mathbf{w}' (\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B}) \mathbf{w} \\
= & (\mathbf{y} - \mathbf{w})' (\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B}) (\mathbf{y} - \mathbf{w}) + \mathbf{x}'\mathbf{C}\mathbf{x} - \mathbf{w}' (\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B}) \mathbf{w}.
\end{aligned}$$

Assim, a distribuição marginal de \mathbf{X} é dada por:

$$\begin{aligned}
f_{\mathbf{X}}(\mathbf{x}) & \propto \exp(\mathbf{x}'\mathbf{C}\mathbf{x} - \mathbf{w}' (\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B}) \mathbf{w}) \\
& \propto \exp\left(\mathbf{x}'\mathbf{C}\mathbf{x} - \mathbf{x}'\mathbf{C}\mathbf{H}(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})^{-1} (\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B}) \times \right. \\
& \quad \left. (\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})^{-1} \mathbf{H}'\mathbf{C}\mathbf{x}\right) \\
& \propto \exp\left(\mathbf{x}'\mathbf{C}\mathbf{x} - \mathbf{x}'\mathbf{C}\mathbf{H}(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})^{-1} \mathbf{H}'\mathbf{C}\mathbf{x}\right) \\
& \propto \exp\left(\mathbf{x}' \left(\mathbf{C} - \mathbf{C}\mathbf{H}(\mathbf{H}'\mathbf{C}\mathbf{H} + \mathbf{B})^{-1} \mathbf{H}'\mathbf{C}\right) \mathbf{x}\right).
\end{aligned}$$

APÊNDICE B - Decomposição em Valores Singulares

Teorema 3. (*Decomposição Espectral ou de Jordan*). Seja \mathbf{L} uma transformação linear; $\mathbf{L} : \mathbb{R}^n \rightarrow \mathbb{R}^n$, simétrica. Então, existe uma base ortonormal de autovetores $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$, tal que $\mathbf{L}(\mathbf{v}_i) = \lambda_i \mathbf{v}_i$, com λ_i autovalores de \mathbf{L} . De forma equivalente, para toda matriz simétrica $\mathbf{L}_{n \times n}$ existe uma decomposição da forma $\mathbf{L} = \mathbf{V}\mathbf{D}\mathbf{V}'$, em que \mathbf{D} é a matriz diagonal formada pelos autovalores λ_i de \mathbf{L} , e $\mathbf{V}_{n \times n} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n]$ é a matriz ortogonal formada pelos autovetores de \mathbf{L} em suas colunas.

Esse teorema pode ser generalizado da forma:

Teorema 4. (*Decomposição em valores singulares*). Seja \mathbf{L} uma transformação linear; $\mathbf{L} : \mathbb{R}^m \rightarrow \mathbb{R}^n$, de posto r . Então existe uma base ortonormal de \mathbb{R}^n , $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$, e uma base ortonormal de \mathbb{R}^m , $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m\}$, tal que $\mathbf{L}(\mathbf{v}_i) = \sigma_i \mathbf{u}_i$ e $\mathbf{L}'(\mathbf{u}_i) = \sigma_i \mathbf{v}_i$, $i = 1, \dots, r$ e $\sigma_i = 0$, para $i \geq r + 1$. De forma equivalente, toda matriz $\mathbf{L}_{n \times m}$ de posto r pode ser decomposta em $\mathbf{L} = \mathbf{V}\mathbf{D}\mathbf{U}'$, em que $\mathbf{D}_{r \times r}$ é a matriz diagonal formada pelos valores singulares σ_i de \mathbf{L} , e $\mathbf{V}_{n \times r} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r]$ e $\mathbf{U}_{m \times r} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r]$ são matrizes ortonormais por coluna.

O teorema da decomposição em valores singulares admite a seguinte interpretação geométrica (KALMAN, 1996):

Seja $\mathbf{L}_{n \times m}$, de posto r , e seja \mathbf{v} um vetor na esfera unitária de \mathbb{R}^n , isto é, $\mathbf{v} = \sum_{i=1}^n x_i \mathbf{v}_i$ e $\|\mathbf{v}\|^2 = x_1^2 + \dots + x_n^2 = 1$. Então, existem bases tais que

$\mathbf{L}(\mathbf{v}_i) = \sigma_i \mathbf{u}_i$. Logo, para $i = 1, \dots, r$,

$$\mathbf{L}\mathbf{v} = \mathbf{L}\left(\sum_{i=1}^r x_i \mathbf{v}_i\right) = \sum_{i=1}^r x_i \mathbf{L}(\mathbf{v}_i) = \sum_{i=1}^r x_i \sigma_i \mathbf{u}_i.$$

Fazendo $y_i = x_i \sigma_i$, segue que:

$$\frac{y_1^2}{\sigma_1^2} + \dots + \frac{y_r^2}{\sigma_r^2} = \frac{x_1^2 \sigma_1^2}{\sigma_1^2} + \dots + \frac{x_r^2 \sigma_r^2}{\sigma_r^2} = x_1^2 + \dots + x_r^2 \leq 1.$$

Portanto, a imagem pela transformação \mathbf{L} da esfera unitária no \mathbb{R}^n é um elipsóide, como pode ser observado pela Figura 39. Este elipsóide será sólido se o posto da transformação for menor que n .

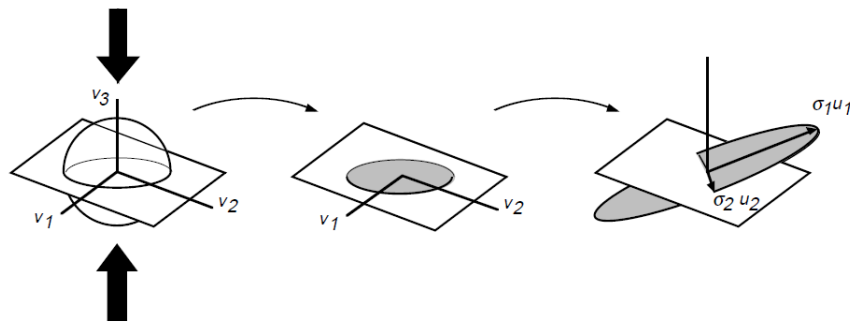


Figura 39 Como a transformação \mathbf{L} deforma uma esfera em um elipsóide

APÊNDICE C - Método das Potências

O Método das Potências, ou *Power Method*, é um método para obtenção dos autovalores e autovetores de uma matriz. Sejam $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ os autovetores de uma matriz $\mathbf{A}_{n \times n}$, associados aos autovalores $\lambda_1, \lambda_2, \dots, \lambda_n$, respectivamente, em que $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. Como os autovetores formam uma base ortonormal, pode-se obter um vetor arbitrário $\mathbf{v}^{(0)}$ como uma combinação linear dos vetores desta base:

$$\mathbf{v}^{(0)} = c_1 \mathbf{x}_1 + c_2 \mathbf{x}_2 + \dots + c_n \mathbf{x}_n, \quad (13)$$

em que c_1, c_2, \dots, c_n são constantes reais.

Pré-multiplicando por \mathbf{A} ambos os lados da igualdade (13), e utilizando que $\mathbf{A}\mathbf{x}_i = \lambda_i \mathbf{x}_i$, tem-se:

$$\begin{aligned} \mathbf{A}\mathbf{v}^{(0)} &= \mathbf{A}(c_1 \mathbf{x}_1 + c_2 \mathbf{x}_2 + \dots + c_n \mathbf{x}_n) \\ &= c_1 \mathbf{A}\mathbf{x}_1 + c_2 \mathbf{A}\mathbf{x}_2 + \dots + c_n \mathbf{A}\mathbf{x}_n \\ &= c_1 \lambda_1 \mathbf{x}_1 + c_2 \lambda_2 \mathbf{x}_2 + \dots + c_n \lambda_n \mathbf{x}_n. \end{aligned}$$

Repetindo o processo k vezes, obtém-se:

$$\begin{aligned} \mathbf{A}^2 \mathbf{v}^{(0)} &= c_1 \lambda_1^2 \mathbf{x}_1 + c_2 \lambda_2^2 \mathbf{x}_2 + \dots + c_n \lambda_n^2 \mathbf{x}_n \\ &\vdots \\ \mathbf{A}^k \mathbf{v}^{(0)} &= c_1 \lambda_1^k \mathbf{x}_1 + c_2 \lambda_2^k \mathbf{x}_2 + \dots + c_n \lambda_n^k \mathbf{x}_n \\ &= \lambda_1^k \left(c_1 \mathbf{x}_1 + c_2 \frac{\lambda_2^k}{\lambda_1^k} \mathbf{x}_2 + \dots + c_n \frac{\lambda_n^k}{\lambda_1^k} \mathbf{x}_n \right). \end{aligned}$$

Como λ_1 é o autovalor dominante, $\left(\frac{\lambda_i}{\lambda_1}\right)^k \rightarrow 0$ quando $k \rightarrow \infty$, logo,

$$\mathbf{A}^k \mathbf{v}^{(0)} \approx c_1 \lambda_1^k \mathbf{x}_1,$$

para k suficientemente grande (Figura 40).

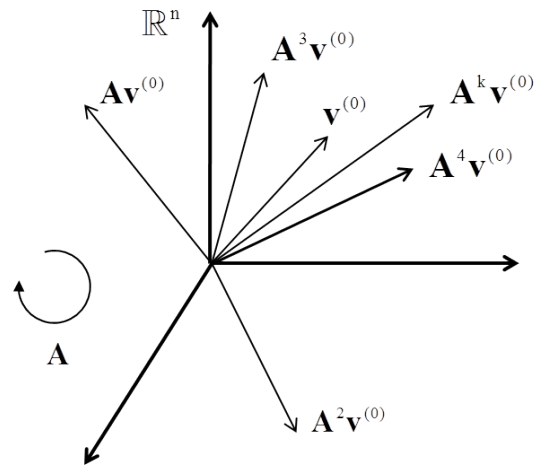


Figura 40 Representação geométrica do método das potências

O método das potências normaliza os produtos $\mathbf{A} \mathbf{v}^{(k-1)}$ para evitar *overflow* ou *underflow* (FERREIRA, 2008).

APÊNDICE D - Rotina Computacional para Algoritmo PLS

Rotina utilizada no Exemplo mostrando o passo a passo do Algoritmo Amostral PLS:

```

library(MBESS)
library(MASS)
library(Matrix)

# gerando uma matriz de covariâncias

RandomSigma <- function(p = 5, df = 10)
{
  T <- diag(sqrt(rgamma(c(rep(1, p)), df/2, 1/2)))
  T[lower.tri(T)] <- 0.98*rgamma((p * (p - 1)/2),
                                df/2, 1/2)

  Sigma <- T %*% t(T)
  return(Sigma)
}

p <- 5
df <- 2
Sigma <- RandomSigma(p, df)
Sigma
eigen(Sigma)
Rho <- diag(diag(Sigma)^-0.5) %*% Sigma %*%
          diag(diag(Sigma)^-0.5)

```

Rho

```
nxy=nrow(Sigma)
```

```
nxy
```

```
# vetor de médias populacional de X e Y
```

```
mpXY=as.matrix(c(1,2,3,4,5))
```

```
mpXY
```

```
# gerando X e Y de uma normal multivariada
```

```
XY=mvrnorm(6,mpXY,Sigma)
```

```
XY
```

```
n = dim(XY) [1]
```

```
n
```

```
X = as.matrix(XY[1:n,1:4])
```

```
X
```

```
m = dim(X) [2]
```

```
m
```

```
X1=sum(X[1:n,1])/n # média da covariável X1
```

```
X1
```

```
X2=sum(X[1:n,2])/n # média da covariável X2
```

```
X2
```

```
X3=sum(X[1:n,3])/n # média da covariável X3
```

```
X3
```

```
X4=sum(X[1:n,4])/n # média da covariável X4
```

```
X4

maX=(c(X1,X2,X3,X4)) # vetor de médias amostrais de X
maX

Y = as.matrix(XY[1:n,5])
Y
k = dim(Y)[2]
k

maY=sum(Y[1:n,1])/n # vetor de médias amostrais de Y
maY

# simulação PLS

nc = 4 # variando o n° de componentes

# os vetores armazenados
T = U = matrix(0,n,nc)
C = matrix(0,k,nc)
P = W = matrix(0,m,nc)
Cont = matrix(0,1,nc)

# O algoritmo:

contnumbcomp = 1 # contador do n° componentes
```

```

u = as.matrix(Y[,1])      # passo 1 - o chute inicial
while (contnumbcomp<=nc){
  flag = 0                # indicador de convergência
  cont = 0
  while (flag<1){
    cont = cont + 1
    w = t(X)%*%u/((t(u)%*%u)[1,1])      # passo 2
    w = w/sqrt((t(w)%*%w)[1,1])         # passo 3
    t = X%*%w                            # passo 4
    c = t(Y)%*%t/((t(t)%*%t)[1,1])     # passo 5
    c = c/sqrt((t(c)%*%c)[1,1])        # passo 6
    if (max(sqrt((u - Y%*%c)^2)) > 0.0001) u = Y%*%c
    else flag = 1                       # passos 7 e 8
  }
  p = t(X)%*%t/((t(t)%*%t)[1,1])      # passo 9
  q = t(Y)%*%u/((t(u)%*%u)[1,1])     # passo 10
  Cont[1,contnumbcomp] = cont
  W[1:m,contnumbcomp] = w
  T[1:n,contnumbcomp] = t
  C[1:k,contnumbcomp] = c
  U[1:n,contnumbcomp] = u
  P[1:m,contnumbcomp] = w
  b = (t(u)%*%t)[1]/((t(t)%*%t)[1])   # passo 11
  X = X - t%*%t(p)                    # passo 12
  Y = Y - b*(t%*%t(c))                # passo 12
  contnumbcomp = contnumbcomp + 1

```

```
}

# chamando as matrizes de dados X e Y novamente:
X = as.matrix(XY[1:n,1:4])
X
Y = as.matrix(XY[1:n,5])
Y

# estimando via OLS
beta_ols=solve(t(X)**X)**t(X)**Y
beta_ols
y_ols=X**beta_ols
y_ols

# estimando via PLS
beta_pls_T=solve(t(T)**T)**t(T)**Y
beta_pls_T
beta_pls_X=P**beta_pls_T
beta_pls_X

y_pls=T**beta_pls_T
y_pls

y_pls=X**P**beta_pls_T
y_pls
```

```
# gerando novas matrizes de dados X e Y para predição
XYp=mvrnorm(1,mpXY,Sigma)
XYp
Xp = XYp[1:4]
Xpcentrado=Xp-maX

Yp = XYp[5]
Yp

# Eq. de predição:  $y = y_m + (x-x_m)\beta$ 
y_pls=maY+Xpcentrado**P**beta_pls_T
y_pls
y_ols=maY+Xpcentrado**beta_ols
y_ols

# verificando ortonormalidade de T

t1 = as.matrix(T[1:n,1])
t2 = as.matrix(T[1:n,2])
t3 = as.matrix(T[1:n,3])

t(t1)**t2
t(t1)**t3
t(t2)**t3
```

CAPÍTULO 2

Quadrados Mínimos Parciais aplicado à seleção genômica para qualidade de carne em suínos

RESUMO

A principal contribuição da genética molecular no melhoramento animal é a utilização direta das informações de DNA no processo de identificação de animais geneticamente superiores. No âmbito da seleção genômica, uma vez que o número de marcadores é geralmente muito maior que o número de animais genotipados e tais marcadores são altamente correlacionados, métodos estatísticos baseados na redução de dimensionalidade apresentam grande aplicabilidade. Dentre esses métodos, destacam-se a Regressão em Componentes Principais (PCR) e os Quadrados Mínimos Parciais (PLS). Para a aplicação de tais métodos, a determinação do número ótimo de componentes a ser utilizado ainda se caracteriza como uma questão relevante, no entanto, metodologias como a teoria de graus de liberdade (DoF) e a validação cruzada (CV) têm sido propostas. Diante do exposto, objetivou-se aplicar os métodos PCR e PLS, além da regressão tradicional, sem redução de dimensionalidade, em uma análise de seleção genômica em suínos, considerando um painel de marcadores SNPs de baixa densidade e dois fenótipos relacionados à qualidade da carne (pH medido aos 45 min e às 24 horas após o abate). Objetivou-se, ainda, testar as teorias DoF e CV na determinação do número ótimo de componentes na análise PLS, e sua influência na performance preditiva do método. Os resultados mostraram que a metodologia PLS é eficiente para a seleção genômica por possibilitar previsões satisfatórias para o pH da carne suína utilizando apenas informações genotípicas, além de identificar uma região relevante no cromossomo 4. Os métodos DoF e CV foram compatíveis para a determinação do número ótimo de componentes na análise PLS, sendo este de eficiência superior ao PCR e à regressão múltipla tradicional.

Palavras-chave: Graus de Liberdade. Marcadores SNPs. pH. Validação Cruzada.

ABSTRACT

The main contribution of molecular genetics in animal breeding is the direct use of DNA information in the identification process of genetically superior animals. Within the genomic selection, since the number of markers is generally much larger than the number of genotyped animals and such markers are highly correlated, statistical methods based on dimensionality reduction have wide applicability. Among these methods the Principal Components Regression (PCR) and Partial Least Squares (PLS) are highlighted. To further the application of such methods, determining the optimal number of components to be used is characterized as a relevant issue, however, methodologies such as the theory of degrees of freedom (DoF) and cross-validation (CV) have been proposed. Given these facts, we sought to apply PCR and PLS methods, beyond traditional regression without dimensionality reduction, in an analysis of genomic selection in pigs considering a panel of low-density SNP markers and two phenotypes related to meat quality (pH measured at 45 minutes and at 24 hours after slaughter). We still aimed to test the DoF and CV theories to determine the optimal number of components in the PLS analysis, and its influence on the predictive performance of the method. The results showed that the PLS method is efficient for genomic selection for enabling satisfactory predictions for swine meat pH using only genotypic information, and for identifying an important region on chromosome 4. The DoF and CV methods were compatible for determining the optimal number of components in the PLS analysis, and this method demonstrated itself as being more efficient than PCR and traditional multiple regression.

Keywords: Degrees of Freedom. SNPs Markers. pH. Cross Validation.

1 INTRODUÇÃO

A partir do início do século XXI, os avanços biotecnológicos na área de automação do processo de genotipagem permitiram o desenvolvimento de novas classes de marcadores moleculares (JENKINS; GIBSON, 2002), dentre os quais se destacam os SNPs (*Single Nucleotide Polymorphisms*). Diante da abundância destes marcadores, Meuwissen; Hayes; Goddard (2001) idealizaram a seleção genômica ampla (*Genome Wide Selection - GWS*). A principal contribuição da genética molecular no melhoramento animal é a utilização direta das informações de DNA no processo de identificação de animais geneticamente superiores.

Devido à alta densidade dos marcadores SNPs no genoma, é possível assumir que alguns deles estejam em desequilíbrio de ligação com locos de características quantitativas (*Quantitative Trait Loci - QTL*), possibilitando sua utilização direta na estimação do valor genético genômico de indivíduos sujeitos a seleção, inclusive de indivíduos que ainda não foram fenotipados. No âmbito da seleção genômica, uma vez que o número de marcadores é geralmente muito maior que o número de animais genotipados (alta dimensionalidade) e tais marcadores são altamente correlacionados (multicolinearidade devida ao desequilíbrio de ligação), métodos estatísticos baseados na redução de dimensionalidade apresentam grande aplicabilidade.

Dentre estes métodos destacam-se a regressão em Componentes Principais (*Principal Components Regression - PCR*) e os Quadrados Mínimos Parciais (*Partial Least Squares - PLS*), os quais são recomendados para situações em que se tem mais covariáveis do que observações (HOSKULDSSON, 1988) e também alta correlação entre as covariáveis, uma vez que garantem que a correlação entre qualquer par de variáveis latentes, ou componentes (combinações lineares das

covariáveis), seja igual a zero. A principal diferença entre estes dois métodos é que o PCR leva em consideração apenas as variáveis explicativas na construção dos componentes, enquanto que o PLS também leva em consideração as variáveis respostas (GARTHWAITE, 1994).

Embora o PLS apresente uma vantagem teórica sobre o PCR por considerar a variável resposta no processo de formação dos componentes, a determinação do número ótimo de tais componentes ainda se caracteriza como um problema relevante para a aplicação do referido método, uma vez que diferentes números de componentes podem levar a diferentes resultados. Metodologias como a teoria de graus de liberdade (*Degrees of Freedom - DoF*), a qual é baseada em uma sofisticada teoria estatística (KRÄMER; SUGIYAMA, 2011); e a validação cruzada (*Cross Validation - CV*), a qual é de caráter empírico, têm sido propostas. Porém, até o momento, comparações entre essas metodologias sob o enfoque de predição genômica ainda não foram reportadas na literatura e merecem ser averiguadas.

Atualmente, fenótipos relacionados com a qualidade da carne são relevantes para a seleção genômica em suínos e, dentre estes, destacam-se as medidas de pH, como o pH inicial (45 minutos após o abate) e o pH último (24 horas após abate). A importância de tais fenótipos diz respeito ao fato dos mesmos estarem diretamente relacionados com a retenção de água, maciez, suculência e aparência da carne, que por sua vez relacionam-se com a aceitabilidade e a palatabilidade (BENEVENTO JUNIOR, 2001). Além disso, como são fenótipos avaliados após o abate, a identificação de indivíduos superiores destinados a seleção pode ser realizada precocemente por meio de painéis de marcadores SNPs identificados a partir de análises estatísticas específicas como PCR e PLS.

Diante do exposto, objetivou-se aplicar os métodos PCR e PLS, além da regressão tradicional sem redução de dimensionalidade, em uma análise de sele-

ção genômica em suínos, considerando um painel de marcadores SNPs de baixa densidade e dois fenótipos relacionados com a qualidade da carne (pH medido aos 45 min e às 24 horas após o abate). Além disso, objetivou-se, ainda, testar dois diferentes métodos de determinação do número ótimo de componentes na análise PLS, a teoria de graus de liberdade e a validação cruzada, bem como sua influência na performance preditiva do método.

2 MATERIAL E MÉTODOS

2.1 Descrição dos dados utilizados

Os dados utilizados no presente estudo são provenientes da Granja de Melhoramento de Suínos do Departamento de Zootecnia (DZO) da Universidade Federal de Viçosa (UFV), em Viçosa, Minas Gerais, Brasil. Neste experimento, a população F_2 foi composta de 345 suínos provenientes do cruzamento de dois varrões da raça local brasileira Piau, com 18 fêmeas de linhagem desenvolvida na UFV, pelo acasalamento de animais de linha comercial (Landrace x Large White x Pietrain).

Os detalhes dos procedimentos utilizados, cuja extração do DNA foi realizada no Laboratório de Biotecnologia Animal do Departamento de Zootecnia da Universidade Federal de Viçosa, podem ser encontrados em Peixoto *et al.* (2006). A genotipagem foi realizada via tecnologia Golden Gate/Vera Code R, no Laboratório de Genética Animal (LGA), Embrapa Recursos Genéticos e Biotecnologia (CENARGEN), Brasília, DF, conforme descrito por Hidalgo *et al.* (2013). Os marcadores SNPs utilizados estão distribuídos da seguinte forma nos cromossomos da espécie *Sus scrofa domesticus*: SSC1 (56), SSC4 (54), SSC7 (59), SSC8 (31), SSC17 (25) e SSCX (12), totalizando, assim, 237 marcadores.

Os dados fenotípicos (características de qualidade de carne) de 345 indivíduos foram mensurados após o abate (realizado aproximadamente aos 105 dias de idade dos animais) e, dentre estas características, optou-se por aquelas relacionadas com o pH da carne *post-mortem* (pH aos 45 minutos, pH_{45} , e às 24 horas, pH_{24}). Estas medidas de pH foram realizadas pela inserção de um eletrodo de vidro (DIGIMED, DME-CV1), acoplado a um pHmetro DIGIMED DM-20, previa-

mente calibrado, no músculo *Longissimus dorsi*, retirado da região imediatamente posterior à última costela do animal.

2.2 Quadrados Mínimos Parciais - PLS

A metodologia do PLS consiste em formar componentes t_i , ($i = 1, \dots, m$), que capturem a maior quantidade de informação possível disposta nas variáveis explicativas X_1, \dots, X_p , neste caso, os marcadores SNPs (os quais assumem os valores 0, 1 e 2, respectivamente aos genótipos aa, aA e AA), visando a prever a variável dependente Y , neste caso, os dois fenótipos considerados (pH_{45} e pH_u). Sob esse enfoque, a equação de regressão é expressa por:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 t_1 + \hat{\beta}_2 t_2 + \dots + \hat{\beta}_m t_m, \quad (1)$$

em que t_i é o vetor coluna que compõe a matriz de componentes \mathbf{T} e $\hat{\beta}_i$ é a estimativa dos coeficientes da regressão entre \mathbf{Y} e \mathbf{T} , $\forall i = 1, \dots, m$.

A correlação de qualquer par de componentes é igual a 0, isto é: $cor(t_i, t_{i'}) = 0, \forall i \neq i'$. Assim, o método PLS reduz o número de termos na equação de regressão, uma vez que o número de componentes na equação (1) geralmente é menor que o número de variáveis \mathbf{X} .

Para simplificar os cálculos, utilizam-se variáveis centradas de \mathbf{Y} e \mathbf{X}_j denotadas respectivamente por \mathbf{U}_1 e \mathbf{V}_{1j} , em que:

$$\mathbf{U}_1 = \mathbf{Y} - \bar{Y}\mathbf{1} \text{ e } \mathbf{V}_{1j} = \mathbf{X}_j - \bar{X}_j\mathbf{1} \text{ para } j = 1, \dots, p,$$

sendo seus vetores de valores dados por $\mathbf{u}_1 = \mathbf{y} - \bar{y}\mathbf{1}$ e $\mathbf{v}_{1j} = \mathbf{x}_j - \bar{x}_j\mathbf{1}$, em que $\mathbf{1}$ é um vetor coluna n -dimensional, $\mathbf{1}' = (1, 1, \dots, 1)'$.

Os componentes são determinados de forma sequencial, sendo que primeiramente é realizada uma regressão de U_1 em função de V_{11} , V_{12} e, assim por diante, até V_{1p} . Tais regressões são independentes, portanto, no desenvolvimento do método PLS, as correlações entre os V_{1j} são ignoradas.

As equações de regressão resultantes do método dos quadrados mínimos ordinários para $j = 1, \dots, p$, são iguais a:

$$U_{1j} = \frac{\mathbf{u}'_1 \mathbf{v}_{1j}}{\mathbf{v}'_{1j} \mathbf{v}_{1j}} \mathbf{v}_{1j}.$$

Segundo Garthwaite (1994), geralmente define-se T_1 como uma média ponderada, dada por:

$$T_1 = \sum_{j=1}^p w_{1j} U_{1j}, \quad (2)$$

em que w_{1j} é um peso apropriadamente escolhido, com $w_{1j} \geq 0$ e $\sum_{j=1}^p w_{1j} = 1$.

Como o componente T_1 é uma média ponderada dos U_{1j} , T_1 não contém todas as informações existentes em X_j . A informação de X_j ausente em T_1 pode ser estimada pelos resíduos obtidos na regressão entre V_{1j} e T_1 . De modo análogo, a variabilidade em Y que não está sendo explicada por T_1 pode ser estimada pelos resíduos obtidos na regressão entre U_1 e T_1 . Estes resíduos são denotados por V_{2j} para V_{1j} e por U_2 para U_1 .

O segundo componente, T_2 , é construído do mesmo modo que T_1 , porém substituindo U_1 e V_{1j} por U_2 e V_{2j} , respectivamente. Desta forma, o procedimento se estende analogamente para os componentes T_2, \dots, T_m .

Generalizando, suponha que T_i seja construído a partir das variáveis U_i e V_{ij} com $j = 1, \dots, p$. Para obter o componente T_{i+1} , as variáveis $V_{(i+1)j}$, e $U_{(i+1)j}$, devem ser determinadas. Com este propósito, é feita uma regressão entre

\mathbf{V}_{ij} e \mathbf{T}_i e, assim, $\mathbf{V}_{(i+1)j}$ é definido por:

$$\mathbf{V}_{(i+1)j} = \mathbf{V}_{ij} - \frac{\mathbf{v}'_{ij} \mathbf{t}_i}{\mathbf{t}'_i \mathbf{t}_i} \mathbf{t}_i$$

sendo \mathbf{t}_i o vetor de valores de \mathbf{T}_i , $\mathbf{V}_{(i+1)j}$ os resíduos da regressão e $\frac{\mathbf{v}'_{ij} \mathbf{t}_i}{\mathbf{t}'_i \mathbf{t}_i}$ o coeficiente da regressão entre \mathbf{V}_{ij} e \mathbf{T}_i .

De forma análoga, $\mathbf{U}_{(i+1)} = \mathbf{U}_i - \frac{\mathbf{u}_i \mathbf{t}'_i}{\mathbf{t}'_i \mathbf{t}_i} \mathbf{t}_i$ e $\mathbf{u}_{(i+1)}$ são os resíduos da regressão entre \mathbf{U}_i e \mathbf{T}_i . Assim, a j -ésima regressão entre \mathbf{U}_{i+1} e $\mathbf{V}_{(i+1)j}$ resulta num coeficiente de regressão dado por:

$$b_{(i+1)j} = \frac{\mathbf{u}'_{(i+1)} \mathbf{v}_{(i+1)j}}{\mathbf{v}'_{(i+1)j} \mathbf{v}_{(i+1)j}}$$

Analogamente à equação 2, define-se $\mathbf{T}_{(i+1)}$ como sendo:

$$\mathbf{T}_{(i+1)} = \sum_{j=1}^p w_{(i+1)j} b_{(i+1)j} \mathbf{V}_{(i+1)j}$$

O método é repetido a fim de obter $\mathbf{T}_1, \mathbf{T}_2, \dots, \mathbf{T}_m$ e, após a obtenção dos m componentes, os coeficientes da equação de regressão (1) são determinados por meio do método dos quadrados mínimos ordinários.

Conforme já relatado, a principal característica do método PLS é que a correlação entre qualquer par de componentes é igual a 0, e isso se deve ao fato dos resíduos provenientes da regressão não serem correlacionados com o regressor, ou seja, $\mathbf{V}_{(i+1)j}$ é não correlacionado com \mathbf{T}_i . Assim, como cada componente $\mathbf{T}_{(i+1)}, \dots, \mathbf{T}_m$ é uma combinação linear de $\mathbf{V}_{(i+1)j}$, os componentes são não correlacionados com \mathbf{T}_i . Essa característica é de suma importância para a seleção genômica ampla, pois o método PLS torna-se uma alternativa eficiente para

se obter preditores apesar da multicolinearidade presente nos dados de marcadores (covariáveis \mathbf{X}_j). Além disso, os coeficientes da equação de regressão (1) podem ser estimados por uma simples regressão feita entre \mathbf{Y} e \mathbf{T}_i . Ainda, adicionando-se componentes à equação, os componentes anteriores não têm seus coeficientes alterados. Deve-se ressaltar que, após os coeficientes da equação (1) serem estimados, é possível expressar o modelo em relação a \mathbf{X}_j , ao invés dos componentes \mathbf{T}_i , ou seja, expressar o modelo em termos de suas variáveis originais e com interpretação biológica.

A necessidade dessa descrição detalhada do método PLS é justificada pelo fato desta teoria ser pouco difundida na área de genética e melhoramento. Por outro lado, a teoria do PCR é amplamente utilizada e já dispõe de vasta literatura abordando aspectos teóricos (ver Capítulo 1) e aplicações na referida área (AZEVEDO *et al.*, 2013), portanto, sua descrição aprofundada não é destacada no presente trabalho. Quanto ao método da regressão múltipla tradicional, isto é, sem redução de dimensionalidade, vale ressaltar que em situações nas quais o número de marcadores é maior que o número de observações fenotípicas, a aplicação da mesma só é possível via utilização de inversas generalizadas, como a inversa de Moore-Penrose.

2.3 Número ótimo de componentes (variáveis latentes) no PLS

Os graus de liberdade (DoF) quantificam a complexidade intrínseca de um método de regressão (VAN DER VOET, 1999). Segundo Krämer e Sugiyama (2011), a complexidade da análise PLS depende da colinearidade das variáveis predictoras, de forma que, quanto maior a colinearidade, menor é a complexidade e, portanto, menor o DoF (menor o número de componentes). Estes autores apresen-

tam uma estimativa não-viesada dos DoF para PLS com m componentes latentes, $\mathbf{T} = (\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_m)$, a qual é dada por:

$$\text{DôF}(m) = 1 + \text{tr} \left(\frac{\partial \mathbf{P}_T \mathbf{Y}}{\partial \mathbf{y}} \right). \quad (3)$$

O termo constante 1 corresponde à estimativa do intercepto β_0 , que consome um grau de liberdade. Para calcular o traço da derivada em (3) de forma explícita é necessário usar um algoritmo baseado na decomposição ortonormal de Lanczos (LANCZOS, 1950) de \mathbf{X} . Para y e m fixos, a decomposição é única (ELDÉN, 2004).

Uma forma mais prática de se determinar o número de componentes é por meio da validação cruzada. Esse método consiste em dividir o conjunto de dados originais em N subconjuntos (*folds*), e realizar N análises, de forma que em cada uma delas um dos subconjuntos é retirado a fim de ser utilizado como recurso para validar a análise realizada. Assim, os valores preditos pelas equações estimadas em cada análise podem ser diretamente comparados com os valores retirados (observados). No contexto de PLS, o número ótimo de componentes será aquele que fornecer a maior correlação entre o conjunto de valores observados e preditos. Embora esse método seja teoricamente mais simples, o mesmo é mais complexo sob o ponto de vista computacional.

2.4 Capacidade Preditiva

A população original de 345 indivíduos foi fracionada em duas diferentes populações, ($N = 2$): população de treinamento e população de validação. A fim de separar os grupos levando em consideração o menor parentesco, utilizou-se o programa RENUMF90 para prover informações de parentesco que permitiram

compor dois grupos com maior parentesco dentro de cada grupo e menor parentesco entre cada grupo. A população de treinamento, contendo 175 indivíduos, foi utilizada para estimar os efeitos dos marcadores SNPs, considerando os métodos PLS, PCR e regressão múltipla tradicional (sem redução de dimensionalidade). A população de validação, contendo 170 indivíduos, não foi utilizada para estimar os efeitos de marcadores, mas sim, para avaliar o potencial preditivo de cada um dos métodos. Dessa forma, os valores preditos para os fenótipos da população de validação (\hat{Y}_p) foram obtidos da seguinte forma: $\hat{Y}_p = X_p \hat{\beta}$, sendo X_p os genótipos dos marcadores SNPs dos indivíduos da população de validação e $\hat{\beta}$ o vetor de estimativas dos efeitos de marcadores provenientes da população de treinamento.

A correlação entre o vetor de fenótipos predito e observado na população de validação é denominada de capacidade preditiva. Sob o ponto de vista prático, essa grandeza permite calcular a eficiência da seleção dos indivíduos da população de validação utilizando apenas informações genotípicas (X_p) e um vetor de estimativas ($\hat{\beta}$) proveniente da população de treinamento. Em outras palavras, a capacidade preditiva quantifica a habilidade de se predizer o mérito genético dos indivíduos da população de validação sem a necessidade de se utilizar fenótipos desta população.

Considerando os fenótipos pH_{45} e pH_u do presente estudo, nota-se que, realmente, a seleção genômica constitui uma ferramenta importante para o melhoramento genético de características de qualidade de carne em suínos, pois se a capacidade preditiva for alta, é possível predizer a qualidade de carne e, assim, identificar os animais superiores, sem a necessidade de abate.

Em posse da escolha do melhor método (PLS, PCR e regressão múltipla tradicional) por meio da capacidade preditiva, os efeitos dos marcadores (em valor absoluto) estimados por este método foram utilizados para a construção dos

gráficos *Manhattan plot*, nos quais cada ponto representa um marcador SNP com sua exata localização no genoma (eixo X mostrando sua localização no cromossomo) e a magnitude de seu efeito (eixo Y). Estes gráficos são importantes para identificar possíveis regiões cromossômicas (QTLs) diretamente associadas com as características de interesse.

Todas as análises foram realizadas no software R (R DEVELOPMENT CORE TEAM, 2013) por meio das funções *pcr*, *pls.model* e *pls.cv* do pacote *pls-dof*.

3 RESULTADOS E DISCUSSÃO

As Figuras 1 e 2 mostram, respectivamente, os resultados da análise de determinação do número ótimo de componentes no PLS para as variáveis pH_{45} e pH_{11} . Por meio destas figuras nota-se que ambos os métodos, DoF e CV, proporcionaram o mesmo número ótimo, o qual pode ser resumido em 30 componentes. Esse valor foi escolhido devido ao fato de se observar a estabilização das curvas a partir deste número de componentes. Ao utilizar números de componentes acima deste valor, tem-se um aumento na complexidade do modelo, mas sua performance permanece inalterada.

Os resultados obtidos na análise de determinação do número de componentes concordam com aqueles mostrados por Krämer e Sugiyama (2011), os quais simularam vários cenários envolvendo diferentes números de componentes e observaram alta concordância entre os métodos DoF e CV quanto à escolha do número ótimo de componentes. Porém, da mesma forma que constatado no presente trabalho, o método CV, por necessitar de maior demanda computacional, é mais oneroso em relação ao tempo de computação e, portanto, em termos práti-

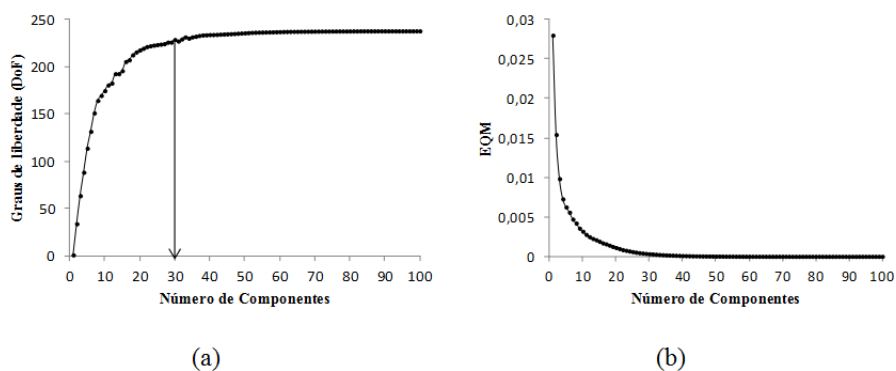


Figura 1 Determinação do número ótimo de componentes na análise PLS utilizando a teoria de graus de liberdade (a) e validação cruzada (b) para a variável pH da carne suína aos 45 min após o abate

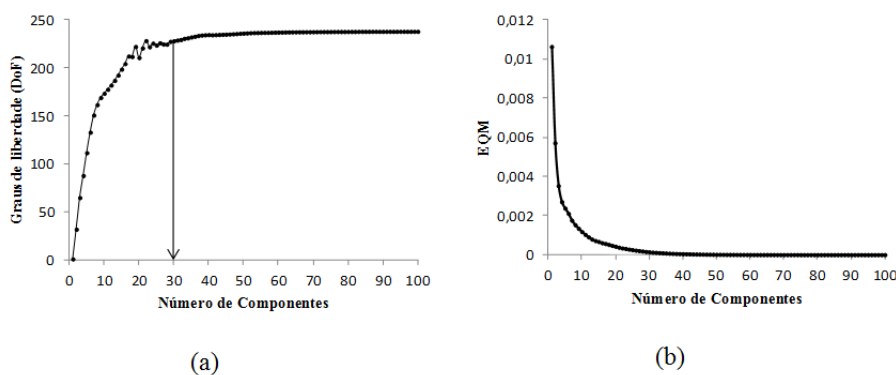


Figura 2 Determinação do número ótimo de componentes na análise PLS utilizando a teoria de graus de liberdade (a) e validação cruzada (b) para a variável pH da carne suína 24 horas após o abate

cos, recomenda-se o método DoF. Além disso, sob o ponto de vista estatístico, tal método deve ser exaltado devido ao seu maior embasamento teórico em relação ao CV que, por sua vez, pode ser caracterizado como um método empírico baseado em esforço computacional.

As capacidades preditivas obtidas pelos métodos de redução dimensional

(PLS e PCR) e Quadrados Mínimos Ordinários (*Ordinary Least Square* - OLS), para as características de pH_{45} e pH_u , são apresentadas na Figura 3. Devido ao fato da análise de número ótimo de componentes no PLS ter mostrado resultados consistentes (DoF e CV indicando 30 componentes), esse mesmo número também foi utilizado na análise PCR a fim de proporcionar uma situação na qual tais métodos sejam diretamente comparáveis.

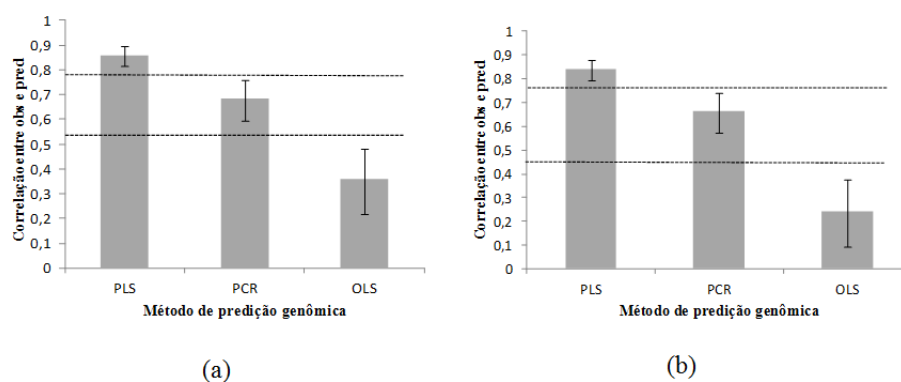


Figura 3 Capacidades preditivas obtidas pelos métodos PLS, PCR e OLS (regressão múltipla tradicional) para as características de pH_{45} (a) e pH_u (b)

Nota-se na Figura 3 que o método PLS mostrou-se mais eficiente para ambas as características, seguido pelo método PCR e, por fim, o método OLS. Nota-se, ainda, que essas diferenças na capacidade preditiva são estatisticamente significativas devido ao fato dos intervalos de confiança de 95% para a correlação não se sobreporem.

De forma geral, embora os métodos PLS e PCR sejam semelhantes no que diz respeito à condição de correlação nula entre os componentes e à possibilidade de serem empregados em situações envolvendo número de variáveis maior que o número de observações, como na presente situação (175 observações individuais e 237 marcadores), observou-se uma nítida vantagem do método PLS em relação a

sua performance preditiva. Certamente, essa vantagem deve-se ao fato do método PLS levar em consideração a variável resposta no processo de composição dos componentes, enquanto o método PCR considera apenas as próprias covariáveis. O fraco desempenho do método OLS pode estar associado à não correção dos problemas de multicolinearidade e à utilização de uma matriz inversa generalizada (Moore-Penrose), no caso $n < p$, a qual pode, em alguns casos, levar a problemas de sobreparametrização (*overfitting*) se comparada com as inversas clássicas utilizadas nos métodos PLS e PCR.

Ainda em relação a Figura 3, vale ressaltar que os valores de capacidade preditiva obtidos pelo método PLS para os fenótipos pH_{45} e pH_u , os quais foram respectivamente 0,85 e 0,84, realmente comprovam a eficiência da seleção para a análise de características de qualidade de carne em suínos. Até o momento, não há relatos na literatura de valores de acurácia de seleção genômica para pH de carne suína, fato este que impossibilita a comparação direta dos resultados obtidos no presente trabalho com outros relatos científicos. Porém, de forma geral, estes valores acima de 80% permitem inferir que seja possível identificar animais geneticamente superiores para pH da carne sem a necessidade de abatê-los, ou seja, utilizando apenas suas informações genótípicas. Tal fato representa um grande avanço tecnológico na área de melhoramento de suínos, pois além de reduzir o intervalo de geração por não necessitar que o animal atinja a idade de abate para inferir sobre seu mérito genético, também é possível evitar que bons animais sejam abatidos e, assim, não destinados à reprodução, por não se conhecer de antemão o mérito de tais animais.

Além da análise preditiva inerente à seleção genômica, também é de interesse a identificação de marcadores SNPs mais relevantes para cada característica, pois nas regiões cromossômicas destes marcadores é possível inferir sobre a pre-

sença de QTLs e/ou genes de efeito maior que podem vir a ser utilizados para fins de seleção. Para tanto, faz-se necessário identificar tais marcadores e suas posições específicas em cada cromossomo, o que pode ser facilmente implementado por análises gráficas como os *Manhattan plots* apresentados na Figura 4.

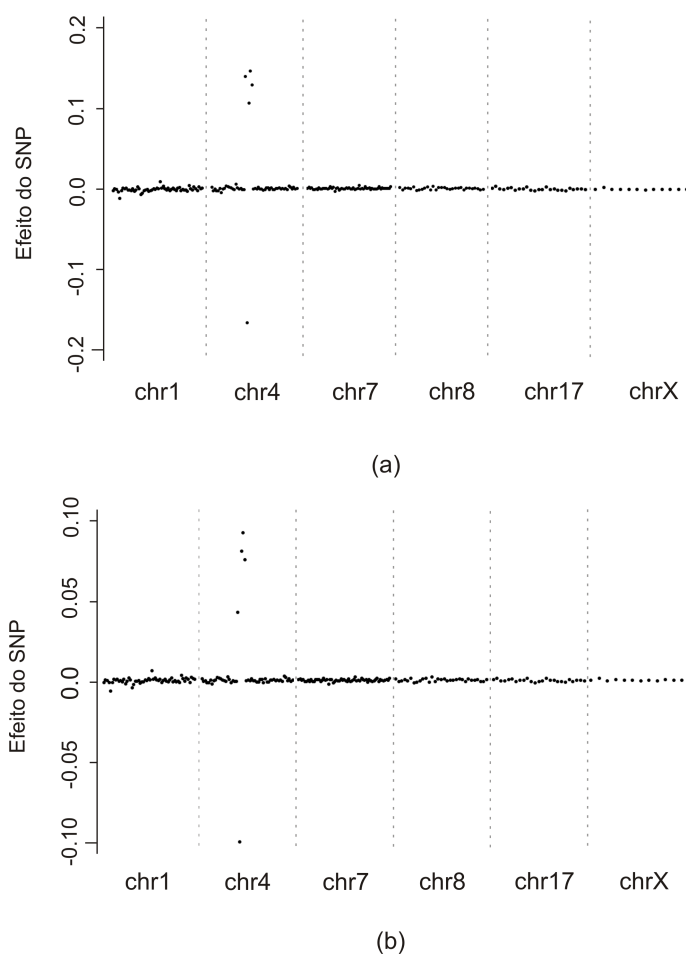


Figura 4 *Manhattan plot* (efeitos estimados dos SNPs ao longo das posições genômicas) para as características pH da carne suína aos 45 min (a) e 24 horas (b) após o abate

De acordo com a Figura 4, nota-se que para ambos os fenótipos (pH_{45} e

Tabela 1 Identificação dos SNPs reportados como sendo os mais relevantes para os fenótipos pH_{45} (45 min após o abate) e pH_u (24 horas após o abate)

Fenótipo	SNP	chr	pos_Mbp	Efeito
pH_{45}	ALGA0026103	4	75,55577	-0,17937
	ALGA0026237	4	80,01721	0,157341
	ALGA0026100	4	75,53339	0,150047
	ALGA0026241	4	80,13745	0,138221
	ALGA0026109	4	75,57379	0,114538
pH_u	ALGA0026103	4	75,55577	-0,11968
	ALGA0026237	4	80,01721	0,108745
	ALGA0026109	4	75,57379	0,095448
	ALGA0026241	4	80,13745	0,089195
	ALGA0026100	4	75,53339	0,050324

pH_u) detectou-se uma região sugestiva de QTL no cromossomo 4, pois é visível a alta influência dos SNPs localizados nessa região sobre os fenótipos em questão. Por meio destes gráficos é possível concluir que a mesma região cromossômica controla o pH da carne suína em diferentes tempos. Para um maior detalhamento dessa região, tem-se informações complementares na Tabela 1. Nessa, observa-se que os SNPs mais relevantes para pH_{45} são os mesmos para pH_u , os quais situam-se em torno da região entre 75 e 80 Mbp do cromossomo SSC4.

Essa região realmente apresenta influência sobre o pH da carne suína, uma vez que várias outras pesquisas têm reportado QTLs nessa mesma região cromossômica. Tal afirmação é comprovada por meio da Figura 5, proveniente da base de dados PigQTLdb (NATIONAL ANIMAL GENOME RESEARCH PROGRAM, 2014), na qual verifica-se alta densidade de QTLs relatados por várias pesquisas científicas. Dentre essas, destacam-se os resultados encontrados por Stratz *et al.* (2012) e Ma *et al.* (2013), que também utilizaram populações F_2 envolvendo linhas comerciais e identificaram QTLs para pH_{45} nas posições 79,5 e 88,26 Mbp,

respectivamente. Também destacam-se os resultados obtidos por Wimmers *et al.* (2006) e Ponsuksili *et al.* (2010), os quais reportaram QTLs para pH_u nas regiões 82 e 67 Mbp, respectivamente. Nesses trabalhos foram utilizadas diferentes populações derivadas da raça Duroc.

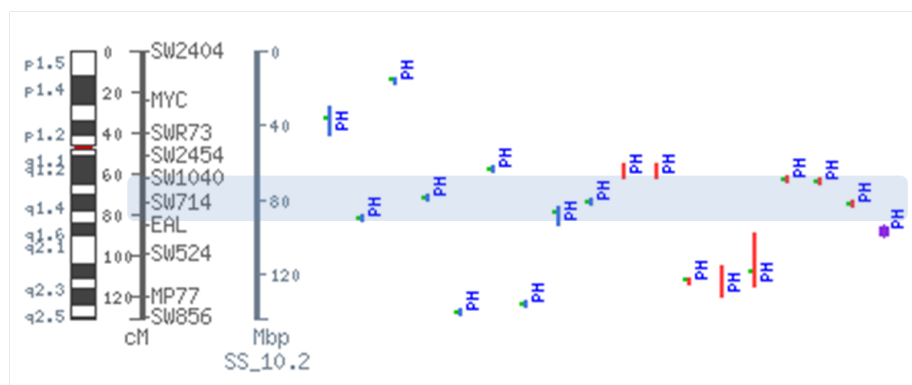


Figura 5 Densidade de QTLs reportados na região intermediária do cromossomo SSC4 proveniente da base de dados PigQTLdb.

Em resumo, embora a comprovação da importância desta região do SSC4 para o pH da carne seja interessante para pesquisas científicas na área de genética suína, ainda são necessárias pesquisas complementares a fim de identificar os genes de efeito maior nesta região. Dessa forma, uma vez constatada a presença desses genes e, compreendidos seus mecanismos de atuação, estes podem ser diretamente empregados com objetivos de seleção por meio de genotipagens específicas para os mesmos, as quais permitem inferir sobre o pH da carne, simplesmente por meio dos genótipos observados para estes genes.

4 CONCLUSÕES

- A metodologia Quadrados Mínimos Parciais é eficiente para a seleção genômica por possibilitar predições satisfatórias (capacidade preditiva acima de 80%) do pH da carne suína, utilizando apenas informações genótípicas baseadas em marcadores SNPs.
- Os métodos Graus de Liberdade e Validação Cruzada foram equivalentes na determinação do número ótimo de componentes na análise por Quadrados Mínimos Parciais.
- O método da Regressão em Componentes Principais apresenta vantagens sobre a regressão múltipla tradicional em relação à predição genômica, porém, este é de eficiência inferior ao método Quadrados Mínimos Parciais.
- A análise Quadrados Mínimos Parciais possibilitou identificar uma região (em torno de 70 e 80 Mbp) relevante no cromossomo 4 para o controle genético do pH da carne suína.

REFERÊNCIAS

- AZEVEDO, C. F. *et al.* Regressão via componentes independentes aplicada à seleção genômica para características de carcaça em suínos. **Pesquisa Agropecuária Brasileira**, Brasília, v. 48, n. 6, p. 619-626, jun. 2013.
- BENEVENUTO JUNIOR, A. A. **Avaliação de rendimento de carcaça e de qualidade da carne de suínos comerciais, nativos e cruzados**. 2001. 94 f. Dissertação (Mestrado em Ciência e Tecnologia de Alimentos) - Universidade Federal de Viçosa, Viçosa, MG, 2001.
- ELDÉN, L. Partial least-squares vs. Lanczos bidiagonalization—I: analysis of a projection method for multiple regression. **Computational Statistics and Data Analysis**, [s.l.], v. 46, n. 1, p. 11-31, 2004.
- GARTHWAITE, P. H. An Interpretation of Partial Least Squares. **Journal of the American Statistical Association**, [s.l.], v. 89, n. 425, p. 122-127, mar. 1994.
- HIDALGO, A. M. *et al.* Fine mapping and single nucleotide polymorphism effects estimation on pig chromosomes 1, 4, 7, 8, 17 and X. **Genetics and Molecular Biology**, [s.l.], v. 36, n. 4, p. 511-519, dez. 2013.
- HOSKULDSSON, A. PLS Regression Methods. **Journal of Chemometrics**, [s.l.], v. 2, p. 211-228, 1988.
- JENKINS, S.; GIBSON, N. High-throughput SNP genotyping. **Comparative and Functional Genomics**, [s.l.], v. 3, n. 1, p. 57-66, fev. 2002.
- KRÄMER, N.; SUGIYAMA, M. The Degrees of Freedom of Partial Least Squares Regression. **Journal of the American Statistical Association**, [s.l.], v. 106, n. 494, p. 697-705, fev. 2011.
- LANCZOS, C. An Iteration Method for the Solution of the Eigenvalue Problem of Linear Differential and Integral Operators. **Journal of Research of the**

National Bureau of Standards, [s.l.], v. 45, n. 4, p. 225-280, out. 1950.

MA, J. *et al.* Genome-Wide Association Study of Meat Quality Traits in a White Duroc \times Erhualian F2 Intercross and Chinese Sutai Pigs. **PLOS ONE**, [s.l.], v. 8, n. 5, p. 1-11, maio 2013.

MEUWISSEN, T. H. E.; HAYES, B. J.; GODDARD, M. E. Prediction of total genetic value using genome wide dense marker maps. **Genetics**, [s.l.], v. 157, p. 1819-1829, abr. 2001.

NATIONAL ANIMAL GENOME RESEARCH PROGRAM. **PigQTLdb**. Disponível em: <<http://www.animalgenome.org/QTLdb/>>. Acesso em: 8 jan. 2014.

PEIXOTO, J. O. *et al.* Associations of leptin gene polymorphisms with production traits in pigs. **Journal of Animal Breeding and Genetics**, [s.l.], v. 123, n. 6, p. 378-383, dez. 2006.

PONSUKSILI, S. *et al.* Identification of expression QTL (eQTL) of genes expressed in porcine *M. longissimus dorsi* and associated with meat quality traits. **BMC Genomics**, [s.l.], v. 11, n. 572, p. 1-14, 2010.

R DEVELOPMENT CORE TEAM. **R**: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, 2013. Disponível em: <<http://www.R-project.org/>>. Acesso em: 14 nov. 2013.

STRATZ, P. *et al.* A two-step approach to map quantitative trait loci for meat quality in connected porcine F_2 crosses considering main and epistatic effects. **Animal Genetics**, [s.l.], v. 44, p. 14-23, 2012.

VAN DER VOET, H. Pseudo-degrees of freedom for complex predictive models: the example of partial least squares. **Journal of Chemometrics**, [s.l.], v. 13, n. 3, p. 195-208, jul. 1999.

WIMMERS, K. *et al.* QTL for microstructural and biophysical muscle properties and body composition in pigs. **BMC Genetics**, [s.l.], v. 7, n. 15, p. 1-14, mar.

2006.

APÊNDICE

APÊNDICE A - Rotina Computacional usada na análise dos dados

```
memory.limit(size = 100000 )
dados=read.csv("all_gws_2013.csv",h=T,sep=";")

snp=as.matrix(dados[,-(1:13)])
ph45=as.matrix(dados[,2])

snp0=as.matrix(snp[-(176:345),])
ph450=as.matrix(ph45[-(176:345)])

snp1=as.matrix(snp[(176:345),])
ph451=as.matrix(ph45[(176:345)])

library("plsdo")

# Regressão Múltipla Tradicional

beta_ols=ginv(t(snp0)%*%snp0)%*%t(snp0)%*%ph450

gbv_ols=snp1%*%beta_ols
cor.test(gbv_ols,ph451)
```

```
# Regressão em Componentes Principais

pcr=pcr(snp0,ph450,m=30)

coef=pcr$coefficients[,30] #manhattan-plot
coef=as.matrix(coef)
max(coef)

gbv_pcr=snp1%*%coef
cor.test(gbv_pcr,ph451)

# Quadrados Mínimos Parciais: DoF

pls=pls.model(snp0,ph450,100,compute.DoF=TRUE,
              compute.jacobian=TRUE)
a=information.criteria(pls$RSS, pls$DoF, pls$yhat,
                       pls$sigmahat, nrow(snp),criterion="bic")
plot(a$DoF)
as.matrix(a$DoF)

coef=pls$coefficients[,30] #manhattan-plot
coef=as.matrix(coef)
max(coef)

gbv_dof=snp1%*%coef
cor.test(gbv_dof,ph451)
```

```
# Quadrados Mínimos Parciais: CV  
  
pls_cv=pls.cv(snp,ph45,m=100)  
plot(pls_cv$cv.error)
```