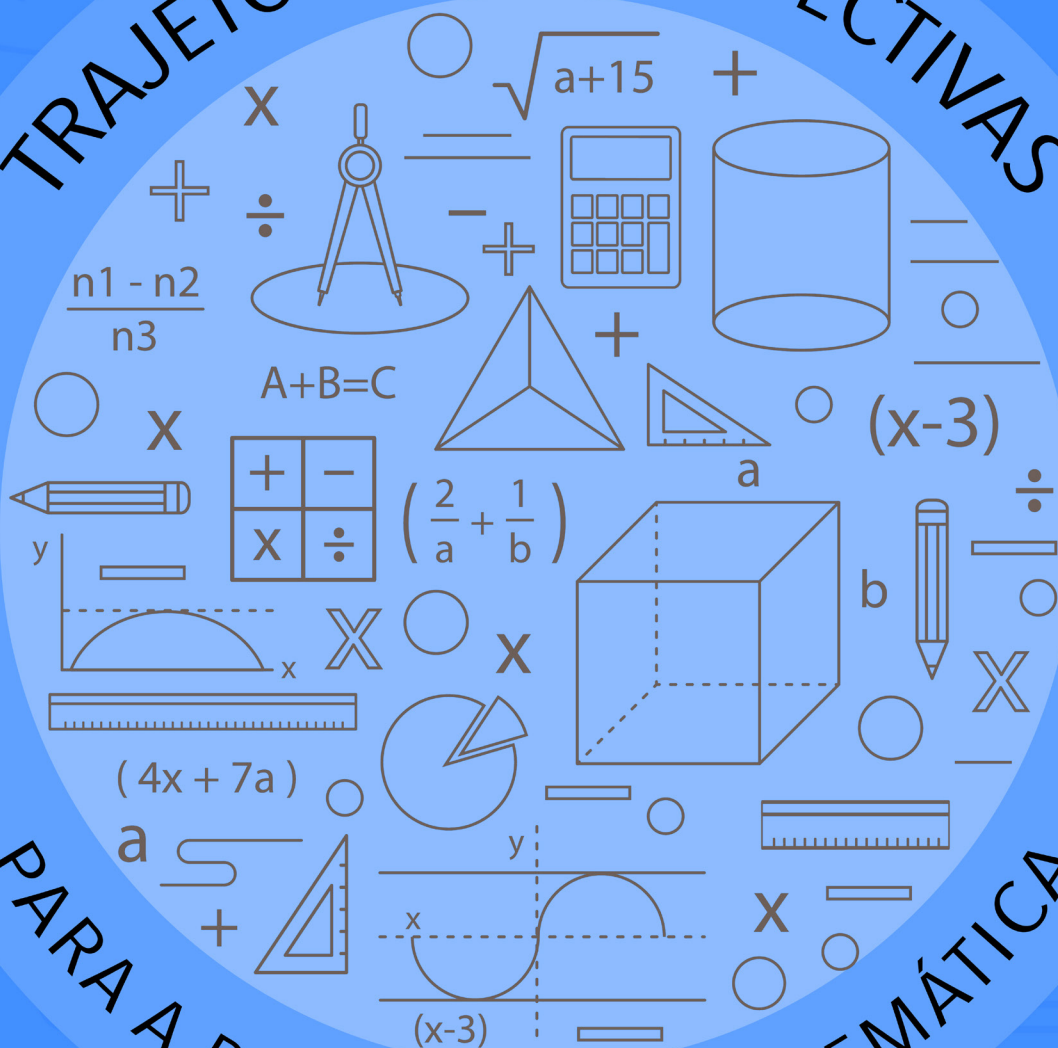


ANIELE DOMINGAS PIMENTEL SILVA  
(Organizadora)

# TRAJETÓRIAS E PERSPECTIVAS



# PARA A PESQUISA EM MATEMÁTICA



**Editora chefe**

Profª Drª Antonella Carvalho de Oliveira

**Editora executiva**

Natalia Oliveira

**Assistente editorial**

Flávia Roberta Barão

**Bibliotecária**

Janaina Ramos

**Projeto gráfico**

Bruno Oliveira

Camila Alves de Cremo

Luiza Alves Batista

**Imagens da capa**

iStock

**Edição de arte**

Luiza Alves Batista

2023 by Atena Editora

Copyright © Atena Editora

Copyright do texto © 2023 Os autores

Copyright da edição © 2023 Atena Editora

Direitos para esta edição cedidos à Atena Editora pelos autores.

Open access publication by Atena Editora



Todo o conteúdo deste livro está licenciado sob uma Licença de Atribuição *Creative Commons*. Atribuição-Não-Comercial-NãoDerivativos 4.0 Internacional (CC BY-NC-ND 4.0).

O conteúdo dos artigos e seus dados em sua forma, correção e confiabilidade são de responsabilidade exclusiva dos autores, inclusive não representam necessariamente a posição oficial da Atena Editora. Permitido o *download* da obra e o compartilhamento desde que sejam atribuídos créditos aos autores, mas sem a possibilidade de alterá-la de nenhuma forma ou utilizá-la para fins comerciais.

Todos os manuscritos foram previamente submetidos à avaliação cega pelos pares, membros do Conselho Editorial desta Editora, tendo sido aprovados para a publicação com base em critérios de neutralidade e imparcialidade acadêmica.

A Atena Editora é comprometida em garantir a integridade editorial em todas as etapas do processo de publicação, evitando plágio, dados ou resultados fraudulentos e impedindo que interesses financeiros comprometam os padrões éticos da publicação. Situações suspeitas de má conduta científica serão investigadas sob o mais alto padrão de rigor acadêmico e ético.

**Conselho Editorial****Ciências Exatas e da Terra e Engenharias**

Prof. Dr. Adélio Alcino Sampaio Castro Machado – Universidade do Porto

Profª Drª Alana Maria Cerqueira de Oliveira – Instituto Federal do Acre

Profª Drª Ana Grasielle Dionísio Corrêa – Universidade Presbiteriana Mackenzie

Profª Drª Ana Paula Florêncio Aires – Universidade de Trás-os-Montes e Alto Douro

Prof. Dr. Carlos Eduardo Sanches de Andrade – Universidade Federal de Goiás

Profª Drª Carmen Lúcia Voigt – Universidade Norte do Paraná

Prof. Dr. Cleiseano Emanuel da Silva Paniagua – Instituto Federal de Educação, Ciência e Tecnologia de Goiás

Prof. Dr. Douglas Gonçalves da Silva – Universidade Estadual do Sudoeste da Bahia  
Prof. Dr. Eloi Rufato Junior – Universidade Tecnológica Federal do Paraná  
Prof<sup>o</sup> Dr<sup>o</sup> Érica de Melo Azevedo – Instituto Federal do Rio de Janeiro  
Prof. Dr. Fabrício Menezes Ramos – Instituto Federal do Pará  
Prof<sup>o</sup> Dr<sup>o</sup> Glécilla Colombelli de Souza Nunes – Universidade Estadual de Maringá  
Prof<sup>o</sup> Dr<sup>o</sup> Iara Margolis Ribeiro – Universidade Federal de Pernambuco  
Prof<sup>o</sup> Dra. Jéssica Verger Nardeli – Universidade Estadual Paulista Júlio de Mesquita Filho  
Prof. Dr. Juliano Bitencourt Campos – Universidade do Extremo Sul Catarinense  
Prof. Dr. Juliano Carlo Rufino de Freitas – Universidade Federal de Campina Grande  
Prof<sup>o</sup> Dr<sup>o</sup> Luciana do Nascimento Mendes – Instituto Federal de Educação, Ciência e Tecnologia do Rio Grande do Norte  
Prof. Dr. Marcelo Marques – Universidade Estadual de Maringá  
Prof. Dr. Marco Aurélio Kistemann Junior – Universidade Federal de Juiz de Fora  
Prof<sup>o</sup> Dr<sup>o</sup> Maria José de Holanda Leite – Universidade Federal de Alagoas  
Prof. Dr. Miguel Adriano Inácio – Instituto Nacional de Pesquisas Espaciais  
Prof. Dr. Milson dos Santos Barbosa – Universidade Tiradentes  
Prof<sup>o</sup> Dr<sup>o</sup> Natiéli Piovesan – Instituto Federal do Rio Grande do Norte  
Prof<sup>o</sup> Dr<sup>o</sup> Neiva Maria de Almeida – Universidade Federal da Paraíba  
Prof. Dr. Nilzo Ivo Ladwig – Universidade do Extremo Sul Catarinense  
Prof<sup>o</sup> Dr<sup>o</sup> Priscila Tessmer Scaglioni – Universidade Federal de Pelotas  
Prof<sup>o</sup> Dr Ramiro Picoli Nippes – Universidade Estadual de Maringá  
Prof<sup>o</sup> Dr<sup>o</sup> Regina Célia da Silva Barros Allil – Universidade Federal do Rio de Janeiro  
Prof. Dr. Sidney Gonçalo de Lima – Universidade Federal do Piauí  
Prof. Dr. Takeshy Tachizawa – Faculdade de Campo Limpo Paulista

## Trajetórias e perspectivas para a pesquisa em matemática 2

**Diagramação:** Camila Alves de Cremo  
**Correção:** Yaidy Paola Martinez  
**Indexação:** Amanda Kelly da Costa Veiga  
**Revisão:** Os autores  
**Organizadora:** Aniele Domingas Pimentel Silva

<b>Dados Internacionais de Catalogação na Publicação (CIP)</b>	
T768	Trajетórias e perspectivas para a pesquisa em matemática 2 / Organizadora Aniele Domingas Pimentel Silva. – Ponta Grossa - PR: Atena, 202
	Formato: PDF Requisitos de sistema: Adobe Acrobat Reader Modo de acesso: World Wide Web Inclui bibliografia ISBN 978-65-258-1050-8 DOI: <a href="https://doi.org/10.22533/at.ed.508231502">https://doi.org/10.22533/at.ed.508231502</a>
	1. Matemática. I. Silva, Aniele Domingas Pimentel (Organizadora). II. Título.
	CDD 510
<b>Elaborado por Bibliotecária Janaina Ramos – CRB-8/9166</b>	

**Atena Editora**  
Ponta Grossa – Paraná – Brasil  
Telefone: +55 (42) 3323-5493  
[www.atenaeditora.com.br](http://www.atenaeditora.com.br)  
contato@atenaeditora.com.br

## DECLARAÇÃO DOS AUTORES

Os autores desta obra: 1. Atestam não possuir qualquer interesse comercial que constitua um conflito de interesses em relação ao artigo científico publicado; 2. Declaram que participaram ativamente da construção dos respectivos manuscritos, preferencialmente na: a) Concepção do estudo, e/ou aquisição de dados, e/ou análise e interpretação de dados; b) Elaboração do artigo ou revisão com vistas a tornar o material intelectualmente relevante; c) Aprovação final do manuscrito para submissão.; 3. Certificam que os artigos científicos publicados estão completamente isentos de dados e/ou resultados fraudulentos; 4. Confirmam a citação e a referência correta de todos os dados e de interpretações de dados de outras pesquisas; 5. Reconhecem terem informado todas as fontes de financiamento recebidas para a consecução da pesquisa; 6. Autorizam a edição da obra, que incluem os registros de ficha catalográfica, ISBN, DOI e demais indexadores, projeto visual e criação de capa, diagramação de miolo, assim como lançamento e divulgação da mesma conforme critérios da Atena Editora.

## DECLARAÇÃO DA EDITORA

A Atena Editora declara, para os devidos fins de direito, que: 1. A presente publicação constitui apenas transferência temporária dos direitos autorais, direito sobre a publicação, inclusive não constitui responsabilidade solidária na criação dos manuscritos publicados, nos termos previstos na Lei sobre direitos autorais (Lei 9610/98), no art. 184 do Código Penal e no art. 927 do Código Civil; 2. Autoriza e incentiva os autores a assinarem contratos com repositórios institucionais, com fins exclusivos de divulgação da obra, desde que com o devido reconhecimento de autoria e edição e sem qualquer finalidade comercial; 3. Todos os e-book são *open access*, *desta forma* não os comercializa em seu site, sites parceiros, plataformas de *e-commerce*, ou qualquer outro meio virtual ou físico, portanto, está isenta de repasses de direitos autorais aos autores; 4. Todos os membros do conselho editorial são doutores e vinculados a instituições de ensino superior públicas, conforme recomendação da CAPES para obtenção do Qualis livro; 5. Não cede, comercializa ou autoriza a utilização dos nomes e e-mails dos autores, bem como nenhum outro dado dos mesmos, para qualquer finalidade que não o escopo da divulgação desta obra.

A coleção “Trajetórias e perspectivas para a pesquisa em matemática 2” tem como foco criar espaços de discussão científica através dos diversificados trabalhos que a compõem. A coletânea abordará trabalhos, pesquisas com relatos de experiências e a matemática no campo interdisciplinar.







O objetivo principal é divulgar algumas pesquisas desenvolvidas por várias instituições de ensino superior do país, cujo eixo central dos trabalhos estão relacionados a metodologias de ensino, tendências em educação matemática e formação de professores. Nesse sentido, observa-se o avanço de pesquisas no campo da educação matemática, visando buscar maneiras que possam tornar a matemática mais atrativa e significativa aos alunos.

Os diversos temas discutidos nesse volume mostram que o conhecimento acadêmico é fundamental, propõe diálogo e reflexão para todos aqueles que tem interesse em conhecer e/ou melhorar sua prática pedagógica e ter um material disponível que permita o contato com essas pesquisas é extremamente relevante.

Deste modo a obra “Trajetórias e perspectivas para a pesquisa em matemática 2” apresenta resultados de pesquisas que foram satisfatórias e que podem aguçar a curiosidade e inspirar os leitores, por isso a importância de espaços como este de divulgação científica.

Aniele Domingas Pimentel Silva




<b>CAPÍTULO 1 .....</b>	<b>1</b>
AS CONTRIBUIÇÕES DO JOGO BATALHA CARTESIANA NO ENSINO E APRENDIZAGEM DE LOCALIZAÇÃO E IDENTIFICAÇÃO DE PONTOS NO PLANO CARTESIANO	
Phablo da Silva Medrado Mateus de Souza Galvão Lucília Batista Dantas Pereira	
 <a href="https://doi.org/10.22533/at.ed.5082315021">https://doi.org/10.22533/at.ed.5082315021</a>	
<b>CAPÍTULO 2 .....</b>	<b>20</b>
COMPREENDENDO A FUNÇÃO AFIM POR MEIO DA MODELAGEM MATEMÁTICA	
Joás Mariano da Silva Júnior Lucília Batista Dantas Pereira	
 <a href="https://doi.org/10.22533/at.ed.5082315022">https://doi.org/10.22533/at.ed.5082315022</a>	
<b>CAPÍTULO 3 .....</b>	<b>37</b>
ENSINO DE FUNÇÕES TRIGONOMÉTRICAS: AS POTENCIALIDADES DE ENSINO COM O GEOGEBRA	
Carlos Alberto Regis	
 <a href="https://doi.org/10.22533/at.ed.5082315023">https://doi.org/10.22533/at.ed.5082315023</a>	
<b>CAPÍTULO 4 .....</b>	<b>44</b>
CONTRIBUIÇÕES DOS OBSTÁCULOS EPISTEMOLÓGICOS DE BACHELARD NO ENSINO DE MATEMÁTICA	
Eduardo Sabel Cristiane Aparecida dos Santos	
 <a href="https://doi.org/10.22533/at.ed.5082315024">https://doi.org/10.22533/at.ed.5082315024</a>	
<b>CAPÍTULO 5 .....</b>	<b>56</b>
ENSINO DE ÁLGEBRA E A LINGUAGEM MATEMÁTICA: E AGORA, TEM LETRAS NA MATEMÁTICA?	
Heloisa Magalhães Barreto Joyce Jaquelinne Caetano	
 <a href="https://doi.org/10.22533/at.ed.5082315025">https://doi.org/10.22533/at.ed.5082315025</a>	
<b>CAPÍTULO 6 .....</b>	<b>68</b>
IDENTIDADE DE SER PROFESSOR NA PERCEPÇÃO DE PROFESSORES DE MATEMÁTICA EM FORMAÇÃO	
Paula Ledoux Tadeu Oliver Gonçalves	
 <a href="https://doi.org/10.22533/at.ed.5082315026">https://doi.org/10.22533/at.ed.5082315026</a>	
<b>CAPÍTULO 7 .....</b>	<b>87</b>
MATEMÁTICA PARA ENSINAR AS OPERAÇÕES BÁSICAS: INVESTIGANDO	

**O MANUAL PEDAGÓGICO DE IRENE DE ALBUQUERQUE DE 1964**


Karina Zolia Jacomelli-Alves

Eduardo Sabel

Eliandra Moraes Pires


 <https://doi.org/10.22533/at.ed.5082315027>**CAPÍTULO 8 ..... 98****TEORIA DE CONJUNTOS E BANCO DE DADOS RELACIONAIS: UMA ABORDAGEM A PARTIR DO USO DE UMA SEQUÊNCIA DIDÁTICA ADAPTATIVA**

Edilaine Jesus da Rocha

 <https://doi.org/10.22533/at.ed.5082315028>**CAPÍTULO 9 ..... 111****DESENVOLVIMENTO DO PENSAMENTO COMPUTACIONAL: UMA PROPOSTA DE ENSINO PARA ESTUDANTES QUE APRESENTAM DISCALCULIA**


Maria Luísa Visinoni Kotrybala

Joyce Jaquelinne Caetano

 <https://doi.org/10.22533/at.ed.5082315029>**CAPÍTULO 10..... 125****MÉTODOS PARA MAPEAMENTO DE QTL ATRAVÉS DE MARCADORES TIPO SNP: UMA COMPARAÇÃO**

Lara Midena João

Daiane Aparecida Zuanetti

 <https://doi.org/10.22533/at.ed.50823150210>**SOBRE A ORGANIZADORA ..... 141****ÍNDICE REMISSIVO ..... 142**

## MÉTODOS PARA MAPEAMENTO DE QTL ATRAVÉS DE MARCADORES TIPO SNP: UMA COMPARAÇÃO

Data de aceite: 01/02/2023

**Lara Midena João**

**Daiane Aparecida Zuanetti**

**RESUMO:** O mapeamento de regiões no genoma associadas a traços quantitativos (QTLs) através de marcadores genéticos do tipo SNP tem sido um dos problemas centrais em Genética e Biologia Molecular e vários métodos de detecção e identificação de QTLs tem sido propostos na literatura. Neste trabalho, três diferentes metodologias foram aplicadas e comparadas nos dados GAW17 quanto ao seus desempenhos em identificar corretamente SNPs relevantes e reguladores de um traço quantitativo. São elas: regressão linear simples com e sem a correção de Bonferroni no nível de significância, LASSO e SPLS. A fim de comparar o desempenho dessas metodologias, utilizamos a sensibilidade e a especificidade como métricas e notamos que o LASSO e a regressão linear simples com nível de significância de 5% apresentam os melhores resultados, uma vez que equilibram valores relativamente altos de sensibilidade e especificidade. O LASSO, por sua vez, também identifica SNPs

influentes mais raros. Para realizar esse estudo, utilizamos algoritmos que já estão implementados em pacotes estatísticos, tais como R, e que estão disponíveis para utilização.

### 1 | INTRODUÇÃO

Um dos problemas centrais na Genética e Biologia Molecular é a detecção e identificação de regiões no genoma associadas a traços quantitativos (fenótipos) de seres vivos. Essas regiões genômicas são geralmente conhecidas como regiões de traços quantitativos (do inglês *quantitative trait loci*, QTLs) e suas posições e efeitos sobre o fenótipo de interesse são estimados através de marcadores genéticos, mais comumente do tipo SNP (do inglês *single nucleotide polymorphism*), dos indivíduos.

Para identificar a(s) região(ões) genômica(s) causadora(s) ou promotora(s) (os QTLs) do fenótipo de interesse, milhares ou milhões de SNPs são genotipados em amostras compostas de centenas ou milhares de indivíduos. Os genótipos dos

SNPs são, então, vistos como covariáveis que podem afetar o fenótipo, considerado como a variável resposta. O fenótipo, quando se trata de uma variável contínua, é geralmente modelado com uma função linear dos efeitos aditivos e de dominância do genótipo dos SNPs e/ou suas interações de segunda, terceira ou maior ordem.

Muitos métodos tem sido propostos e estudados para identificar e selecionar os SNPs mais associados ao fenótipo de interesse. Os métodos mais usados são baseados na estimação do modelo de regressão linear simples entre o genótipo de cada SNP e o fenótipo em estudo.

Os SNPs mais significativos são escolhidos via teste de hipóteses baseado na significância do efeito do SNP no modelo estimado. As vantagens dessa abordagem simples são baixo tempo de processamento computacional, facilidade de uso e de interpretação dos resultados. No entanto, ela geralmente apresenta baixo poder, não permite que os efeitos dos SNPs sobre o fenótipo sejam conjuntamente estimados e também não considera a estrutura de associação existente entre eles (Yazdani, 2014, Feng et al., 2012, Oliveira, 2015).

Métodos alternativos aos modelos de regressão linear simples são as metodologias que permitem analisar os SNPs conjuntamente, alguns deles identificados como métodos de aprendizado de máquina. Entre eles se destacam: florestas aleatórias (Breiman, 2001, Oliveira, 2015), algoritmos genéticos (Goldberg, 1989, Oliveira, 2015), LASSO (do inglês *least absolute shrinkage and selection operator*, Tibshirani, 1996) que pode ser visto como um modelo de regressão Bayesiano com função de verossimilhança normal e distribuição *a priori* exponencial dupla para cada coeficiente (Park e Casella, 2008), e métodos de regressão por mínimos quadrados parciais – SPLS (do inglês *sparse partial least squares*, Chun, 2010). Todos esses métodos nos permitem trabalhar com dados cujo número de covariáveis é muito superior ao número de indivíduos da amostra e nos quais existe covariáveis altamente correlacionadas. Essas características estão presentes nos dados de mapeamento de QTLs através de marcadores SNPs e precisam ser consideradas.

Nesse trabalho, portanto, estudamos e aplicamos três diferentes metodologias nos dados GAW (Genetic Analysis Workshop) 17 (Almasy et al., 2011) para detectar e identificar SNPs associados ao fenótipo de interesse, principalmente SNPs com variantes raras e baixa frequência do alelo menor. Também discutimos as vantagens e desvantagens de cada método e elencamos os algoritmos que já estavam implementados em pacotes estatísticos, tais como R, e disponíveis para utilização.

O texto está organizado como segue. Na Seção 2 descrevemos as três metodologias que são estudadas e aplicadas para a identificação de SNPs relevantes e como o desempenho delas é analisado e comparado. A Seção 3 apresenta o conjunto de dados GAW17 que é estudado e os resultados obtidos. Finalmente, a Seção 4 traz conclusões e discussões.

## 2 I METODOLOGIAS PARA SELEÇÃO DE SNPs INFLUENTES

Nessa seção, descrevemos os métodos que foram utilizados para detectar e identificar SNPs associados ao fenótipo de interesse  $Y$ . São eles: o valor-p do teste de significância associado ao efeito de cada SNP em um modelo de regressão linear simples, o LASSO e o SPLS. Além disso, a fim de comparar qual desses métodos é o melhor para selecionar os SNPs influentes, empregamos as seguintes métricas: especificidade ( $E$ ) e sensibilidade ( $S$ ) (Lee et al., 2011) que são descritas na Seção 2.4.

### 2.1 Regressão linear simples

A fim de identificar e selecionar os SNPs mais associados ao fenótipo de interesse, um método muito utilizado é baseado na estimação do modelo de regressão linear simples entre o genótipo de cada SNP e o fenótipo em estudo. Nesse modelo o genótipo de cada SNP é tido como variável preditora e o fenótipo é considerado como variável resposta.

Seja  $Y = (Y_1, Y_2, \dots, Y_n)$  um traço quantitativo a ser observado em  $n$  indivíduos. O fenótipo  $Y_i$  do  $i$ -ésimo indivíduo pode ser modelado pelo seguinte modelo de regressão linear simples:

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad (1)$$

em que  $\beta_0$  é o valor esperado do fenótipo quando o genótipo do SNP vale zero,  $\beta_1$  é o efeito aditivo que representa o acréscimo no valor esperado do fenótipo quando o genótipo do SNP vale 1, ou seja, quando o indivíduo é homocigoto dominante,  $x_i$  é o genótipo do SNP, considerado na específica análise, do  $i$ -ésimo indivíduo codificado como  $-1, 0$  ou  $1$  para  $aa, Aa$  ou  $AA$ , respectivamente,  $i = 1, \dots, n$ ,  $\epsilon_i \sim \text{Normal}(0, \sigma^2)$  é o erro aleatório e  $\epsilon_i$  e  $\epsilon_{i'}$  são supostamente independentes para  $i \neq i'$ .

Os SNPs significativos são escolhidos via teste de hipóteses baseado na significância de cada SNP no modelo estimado. Esse teste de hipóteses tem como objetivo, além de verificar se a média do fenótipo dos indivíduos varia linearmente em função do genótipo do SNP, identificar quais marcadores genéticos são influentes. As hipóteses associadas ao teste são:

$H_0: \beta_1 = 0$  (a média do fenótipo não varia linearmente em função do genótipo do SNP, ou seja, o efeito do SNP não é significativo) contra

$H_1: \beta_1 \neq 0$  (a média do fenótipo varia linearmente em função do genótipo do SNP, ou seja, o efeito do SNP é significativo).

Para testar essas hipóteses, utilizamos a estatística teste  $F = \frac{QMReg}{QMRes}$ , em que  $QMReg = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$  e  $QMRes = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}$ , sendo  $\hat{y}_i$  o valor do fenótipo do  $i$ -ésimo indivíduo predito e  $\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$ . Sob  $H_0$ , essa estatística segue uma distribuição F de Snedecor com 1 grau de liberdade no numerador e  $n-2$  graus de liberdade no denominador.

Para rejeitarmos ou não a hipótese nula, utilizamos o valor-p definido como a probabilidade de obtermos um valor da estatística do teste igual ou mais desfavorável à hipótese nula ( $H_0$ ) que o valor observado quando  $H_0$  é verdadeira. Mais detalhes sobre esse teste de hipóteses são encontrados em Morettin e Bussab (2017).

Um SNP é considerado significativo, ou seja, ele afeta o fenótipo, quando  $\beta_1 \neq 0$ . Fixado um nível de significância  $\alpha$ , rejeitamos  $H_0$  se o valor-p é menor que  $\alpha$  e não rejeitamos  $H_0$ , caso contrário.

Nesse estudo, cada teste de hipóteses é realizado para diferentes níveis de significância ( $\alpha$ ), são eles: 5% e 10%, ambos com e sem a correção de Bonferroni para múltiplos testes parciais. Essa correção tem como finalidade reduzir a probabilidade de concluir simultânea e erroneamente que muitos SNPs são significativos e, conseqüentemente, identificar muitos SNPs falso positivos. O método é executado do seguinte modo: para que o nível de significância na decisão global (considerando conjuntamente todos os testes realizados, um para cada SNP) seja  $\alpha$ , o nível de significância de cada teste deve ser o resultado da divisão de  $\alpha$  pelo número de testes realizados (Abdi, 2007). Quando muitos testes devem ser feitos, como nesse caso em que testamos o efeito de milhares de SNPs, o nível de significância de cada teste é tão pequeno que o seu poder também cai drasticamente e, assim, podemos não selecionar muitos SNPs verdadeiramente significativos.

## 2.2 LASSO

Com o propósito de identificar e selecionar os SNPs que influenciam o fenótipo de interesse, o LASSO (do inglês *least absolute shrinkage and selection operator*) também tem sido uma das metodologias mais utilizadas. Ele se trata de uma regressão linear múltipla que seleciona as variáveis preditoras utilizando uma restrição através da qual a soma do valor absoluto dos coeficientes da regressão seja menor que uma constante fixada. Desse modo, esse modelo permite analisar os SNPs conjuntamente para cada cromossomo.

Nesse modelo, o genótipo de cada SNP é tido como variável preditora e o fenótipo é considerado como variável resposta. Seja  $Y = (Y_1, Y_2, \dots, Y_n)$  um traço quantitativo a ser observado em  $n$  indivíduos e  $x_i = (x_{i1}, \dots, x_{ip})$  o genótipo de  $p$  SNPs para o  $i$ -ésimo indivíduo. O fenótipo  $Y_i$  do  $i$ -ésimo indivíduo pode ser modelado pelo seguinte modelo de regressão múltipla:

$$Y_i = \beta_0 + \sum_{k=1}^p \beta_k x_{ki} + \epsilon_i, \quad (2)$$

em que  $\beta_0$  é o valor esperado do fenótipo quando o genótipo de todos os SNPs valem zero,  $\beta_k$  é o efeito aditivo do  $k$ -ésimo SNP no valor esperado do fenótipo,  $x_{ki}$  é o genótipo do  $k$ -ésimo SNP do  $i$ -ésimo indivíduo codificado como  $-1, 0$  ou  $1$  para  $aa, Aa$  ou  $AA$ , respectivamente,  $k = 1, \dots, p$  e  $i = 1, \dots, n$ ,  $\epsilon_i$  é o erro aleatório do  $i$ -ésimo indivíduo e  $\epsilon_i$  e  $\epsilon_j$  são supostamente independentes para  $i \neq j$ .

O LASSO é um método que seleciona as variáveis preditoras mais influentes através da adição de uma restrição na fórmula dos mínimos quadrados. Essa restrição é da forma:  $\sum_{k=1}^p |\beta_k| \leq c$ , em que  $c = c(\lambda)$ . Assim, por conta dessa limitação, temos que algumas das estimativas desses coeficientes ( $\beta_k$ s) serão aproximadamente 0. Desse modo, SNPs associados a um coeficiente com estimativa diferente de zero são selecionados como relevantes e SNPs com estimativa zero são descartados da análise.

Na sua forma lagrangiana, a metodologia LASSO busca

$$(\hat{\beta}_0, \hat{\beta}) = \arg \min_{\beta_0, \beta \in \mathbb{R}^p} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{k=1}^p \beta_k x_{ki} \right)^2 + \lambda \sum_{k=1}^p |\beta_k|, \quad (3)$$

em que  $\lambda \geq 0$  é um parâmetro de regularização. Além disso, quanto maior o valor de  $\lambda$ , mais restrita é a limitação e, portanto, menor a quantidade de SNPs significativos. Consequentemente, menos complexo é o modelo. Ademais, se  $\sum_{k=1}^p |\tilde{\beta}_k| \leq c(\lambda)$ , em que  $\tilde{\beta}_k$ s são as estimativas do método de mínimos quadrados, então as estimativas dos coeficientes do LASSO também serão os  $\tilde{\beta}_k$ s, pois não haverá restrição para as estimativas dos parâmetros (Rodrigues, 2018).

De acordo com Rodrigues (2018) temos que, se  $\lambda = 0$ , então o estimador obtido pelo modelo LASSO é idêntico ao obtido pelo método de mínimos quadrados. Ainda, para cada valor de  $\lambda$ , obtemos diferentes valores de  $\hat{\beta}$ .

Desse modo, a fim de obtermos o melhor modelo, é necessário escolher um valor para  $\lambda$  que não assuma o valor 0 nem valores muito grandes, uma vez que, no primeiro caso, todas as covariáveis seriam selecionadas, e, no segundo caso, provavelmente não seriam selecionadas covariáveis e o modelo possuiria apenas o intercepto  $\beta_0$  (Rodrigues, 2018). Assim, o valor de  $\lambda$  escolhido, para cada cromossomo, para compor o modelo é aquele que possui o menor erro de predição dentre todos os valores testados para o  $\lambda$ .

Esses valores do parâmetro de regularização foram definidos através da função “cv.glmnet” (do pacote “glmnet” do software R) que realiza a validação cruzada  $k$ -fold para obtê-los. Essa validação cruzada é um método no qual a amostra é dividida em  $l$  partes (subamostras) mutuamente exclusivas e aproximadamente de mesmo tamanho. O modelo é estimado com  $l - 1$  partes e testado na única parte que não é utilizada para estimar o modelo. Para cada valor de  $\lambda$ , são propostos  $l$  modelos e, para cada um, é calculado o erro de predição ( $EP$ ) como

$$EP = \sum_i \left( y_i - \hat{y}_i \right)^2, \quad (4)$$

em que  $y_i$  e  $\hat{y}_i$  são, respectivamente, o valor observado e o valor predito, pelo modelo de regressão múltipla, do fenótipo dos indivíduos que fazem parte da subamostra não utilizada na estimação.

Esse procedimento é repetido  $l$  vezes, alterando, a cada vez, as  $l - 1$  subamostras

sobre as quais o modelo é estimado e, conseqüentemente, a subamostra na qual ele é testado. Portanto, para cada valor de  $\lambda$ , temos  $l$  erros de predição e calculamos a média desses EPs. O valor de  $l$  utilizado para análise foi 10, que é o valor padrão da função no R e o  $\lambda$  escolhido, para o específico cromossomo, é aquele que apresenta a menor média de EP. Para mais detalhes, ver Lôca e Zuanetti (2021).

## 2.3 SPLS

Outra metodologia que tem sido muito utilizada para identificar e selecionar SNPs que possuem efeito significativo para o fenótipo de interesse é o SPLS (do inglês *sparse partial least squares regression*). Assim como o LASSO, o SPLS se trata de uma regressão linear múltipla que realiza a seleção das variáveis através de uma restrição na soma do valor absoluto dos coeficientes.

Essa metodologia consiste na redução de dimensão dos dados através da inclusão de variáveis latentes (combinações lineares das variáveis originais) e seleção de quais dessas variáveis são relevantes no modelo. Logo, esse método realiza a colapsagem de SNPs que possuem alta correlação em variáveis ocultas e verifica quais delas são realmente significativas para explicar a variabilidade do fenótipo em questão, de modo a analisar conjuntamente os SNPs para cada cromossomo.

Seja  $Y = (Y_1, Y_2, \dots, Y_n)$  um traço quantitativo a ser observado em  $n$  indivíduos. Supondo que  $p$  SNPs são considerados, sendo  $W_{p \times M}$  a matriz de carga fatorial (matriz com os coeficientes do genótipo de cada SNP em um cromossomo) de  $M$  vetores  $(w_1, w_2, \dots, w_M)$  de tamanho  $p$  e sendo  $X_{n \times p}$  a matriz que, para um cromossomo, contém o genótipo dos SNPs do  $i$ -ésimo indivíduo,  $i = 1, \dots, n$ , codificados em  $-1, 0$  ou  $1$  para  $aa, Aa$  ou  $AA$ , respectivamente. Dessa maneira,  $XW$  contém, para cada indivíduo, grande parte das informações originais dos SNPs colapsadas em  $M < \min(n, p)$  combinações lineares de seus genótipos. Cada uma dessas  $M$  combinações lineares do genótipo dos SNPs são denominadas componentes latentes, cujos valores não são observados no conjunto de dados, mas preditos durante o processo de ajuste do modelo.

De maneira similar ao LASSO, o SPLS impõe uma restrição em cada  $w_m$  para que haja uma solução esparsa, ou seja, para que muitos SNPs sejam identificados como não significativos e poucos tenham cargas fatoriais diferente de zero dentro de cada componente latente. Logo, a solução para cada estimativa de  $w_m$  é dada por:

$$\hat{w}_m = \max_{w_m} w_m^t Q w_m \quad (5)$$

em que  $Q = X^t Y Y^t X$  e sujeita a  $w_m^t w_m = 1$  e  $w_m^t S_{XX} \hat{w}_h = 0$ ,  $h = 1, \dots, m - 1$  e  $m = 1, \dots, M - 1$ , em que  $S_{XX}$  representa a matriz de covariância amostral dos SNPs. Desse modo, a solução  $\hat{w}_m$  traz como resultado uma regressão de mínimos quadrados parcial, em que diferentes componentes latentes explicam diferentes partes da variabilidade do fenótipo.

Apesar disso, uma solução esparsa não é alcançada e  $w_{mk}$  geralmente é diferente



de zero para muitos SNPs. Chung e Keles (2010) propõe, então, uma formulação impondo uma penalização em um vetor de cargas substituto ( $c$ ) e não no vetor de cargas fatoriais original ( $w$ ), enquanto mantém  $w$  e  $c$  próximos um do outro. A formulação é:

$$\tilde{w}_m, \tilde{c}_m = \arg \min_{w_m, c_m} \left\{ -d w_m^t Q w_m + (1-d)(c_m - w_m)^t Q (c_m - w_m) + \lambda_1 \sum_{k=1}^p |c_{mk}| + \lambda_2 \sum_{k=1}^p c_{mk}^2 \right\}$$

para  $0 < d \leq \frac{1}{2}$ , tal que  $w_m^t w_m = 1$  e  $w_m^t S X X^t \hat{w}_m = 0$  para  $h=1, \dots, m$ ;  $m=1, \dots, M-1$ , e em que  $Q = X^t Y Y^t X$ .

Enquanto a penalidade  $\lambda_1$  promove a esparsidade, a penalidade  $\lambda_2$  evita que a matriz  $Q$  seja singular (não inversível). Para  $Y$  univariado, como neste estudo, e assumindo  $\lambda_2 = \infty$ , a solução não depende de  $d$  e resulta no limite

$$\tilde{w}_m = \left( |\hat{w}_m| - \eta \max_{1 \leq k \leq p} |\hat{w}_{mk}| \right) \mathbb{I}_{(|\hat{w}_m| \geq \eta \max_{1 \leq k \leq p} |\hat{w}_{mk}|)} \text{sign}(\hat{w}_m), \quad (6)$$

em que  $0 \leq \eta \leq 1, \dots, \lambda_1$  e  $\hat{w}_m$  é a solução da regressão múltipla parcial encontrada na Equação (5) e  $\tilde{w}_m$  é a solução parcial e esparsa da Equação (6), como mostra Feng et al. (2012).

Logo, para cada cromossomo são encontrados valores de  $M$  e  $\eta$  que minimizam o erro de predição ( $EP$ ). Esses valores foram definidos via validação cruzada através da função “cv.spls”(do pacote “spls”do software R) de modo que a amostra foi dividida em  $l = 3$  partes mutuamente exclusivas e aproximadamente de mesmo tamanho. O modelo é estimado com  $l - 1$  partes e testado na única parte que não é utilizada para estimar o modelo. Para cada valor de  $M$  e  $\eta$  são ajustados  $l$  modelos e, para cada um, é calculado o erro de predição ( $EP$ ) da seguinte maneira

$$EP = \sum_i \left( y_i - \hat{y}_i \right)^2,$$

em que  $y_i$  e  $\hat{y}_i$  são, respectivamente, o valor observado e o valor predito, pelo modelo de regressão múltipla, do fenótipo dos indivíduos que fazem parte da subamostra não utilizada na estimação.

Esse procedimento é repetido  $l$  vezes, alterando, em cada momento, as  $l-1$  subamostras em que o modelo é estimado e, conseqüentemente, a subamostra na qual ele é testado. Portanto, para cada valor de  $M$  e  $\eta$ , temos  $l$  erros de predição e calculamos a média desses  $EP$ s. Os valores de  $M$  e  $\eta$  escolhidos, para o específico cromossomo, são aqueles que apresentam a menor média de  $EP$ . Os SNPs relevantes são os que apresentam cargas fatoriais diferentes de zero nos  $M$  componentes latentes que explicam a variabilidade do fenótipo. Para mais detalhes, ver Feng et al. (2012).

## 2.4 Medidas de desempenho

A fim de medir e comparar o desempenho das diferentes metodologias testadas na seleção de SNPs influentes, adotamos duas métricas denominadas especificidade e sensibilidade (Lee et al., 2011).

A sensibilidade é calculada como a proporção de estimativas de  $\beta$  diferentes de 0 dentre os verdadeiros elementos  $\beta$  que são não nulos. Já a especificidade é a razão da quantidade de estimativas de  $\beta$  iguais a 0 dentre os verdadeiros elementos  $\beta$  que são nulos. De maneira mais direta, o numerador da especificidade pode ser calculado como a diferença entre a quantidade de SNPs que não são verdadeiramente significativos e a quantidade de SNPs que foram erroneamente identificados como significativos e o denominador dessa fração é a quantidade de SNPs que realmente não são significativos.

Uma seleção de variáveis ideal ocorre quando tanto a sensibilidade, quanto a especificidade são iguais a um. Todavia, geralmente há uma compensação entre essas medidas, ou seja, quando a especificidade tende a 1, a sensibilidade tende a 0, e vice-versa.

## 3 | ANALISANDO O GAW 17

O banco de dados utilizado nesse estudo é denominado *Genetic Analysis Workshop 17* (GAW 17) e é formado por dados simulados, a partir de características reais da população, de 697 indivíduos sem parentesco, sendo eles 327 homens e 370 mulheres.

Como base para a simulação de uma doença comum e complexa, de fenótipos quantitativos e dos fatores de risco (os marcadores SNPs) foram utilizados os dados reais contidos no *1000 Genomes Project*, que consideram variações genéticas de vários grupos de populações humanas, sendo eles: Europa, Leste Asiático, do sul da Ásia, África Ocidental e de índios americanos.

Um fenótipo binário que indica presença ou não de uma doença e três fenótipos quantitativos contínuos:  $Q_1$ ,  $Q_2$  e  $Q_3$  foram simulados. Para mais detalhes, ver Lóca e Zuanetti (2021) e Almasy et al. (2011).

Originalmente, o banco de dados nos fornecia as informações dos genótipos dos marcadores SNPs através de 16 possíveis pares de bases nitrogenadas, sendo eles: A/A, T/T, C/C, G/G, A/T, A/C, A/G, T/A, T/C, T/G, C/A, C/T, C/G, G/A, G/T, G/C. Todavia, a fim de possibilitarmos esse estudo, classificamos essas informações do seguinte modo:

- C/C ou G/G: homocigoto dominante codificado como 1;
- A/T, A/C, A/G, T/A, T/C, T/G, C/A, C/T, C/G, G/A, G/T, G/C: heterocigoto codificado como 0; e
- A/A ou T/T: homocigoto recessivo codificado como -1.

Nesse estudo, analisamos o fenótipo  $Q_1$  que originalmente é impactado por 39

SNPs em 9 genes. Nosso objetivo é verificar o desempenho das metodologias estudadas em identificar e selecionar os 39 SNPs relevantes entre os 24487 SNPs disponíveis para análise. O boxplot do fenótipo estudado  $Q_1$  está disponível na Figura 1.

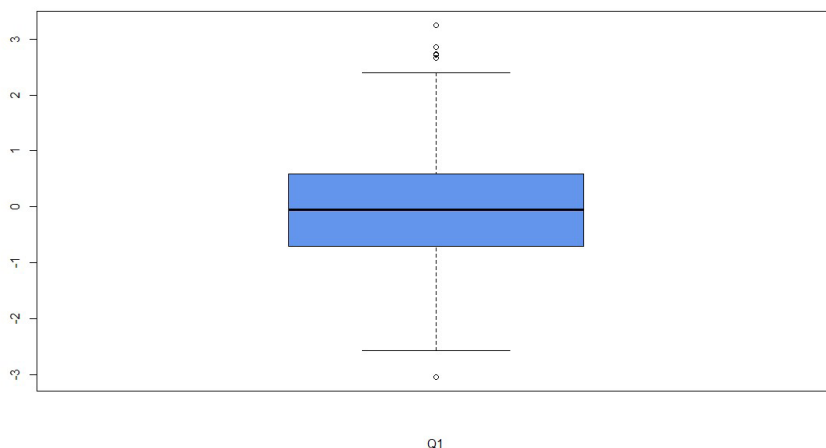


Figura 1: Boxplot do fenótipo  $Q_1$ .

Notamos, através da Figura 1 que o fenótipo  $Q_1$  é uma variável contínua que assume valores no intervalo  $(-3.0, 3.5)$  e possui 5 outliers, sendo que 4 deles ultrapassam o limite superior e 1 ultrapassa o limite inferior. Além disso,  $Q_1$  possui comportamento simétrico e 50% dos valores observados estão entre  $-1$  e  $1$ .

Ademais, de acordo com a simulação realizada, os cromossomos que possuem SNPs significativos são os cromossomos 1, 4, 5, 6, 13, 14 e 19. Esses SNPs que realmente influenciam o fenótipo  $Q_1$  são listados a seguir de acordo com o cromossomo em que eles se encontram (de modo que a notação CaSb identifica o SNP b presente no cromossomo a):

- Cromossomo 1: C1S6533, C1S6537, C1S6540, C1S6542, C1S6561, C1S3181 e C1S3182;
- Cromossomo 4: C4S1861, C4S1873, C4S1874, C4S1877, C4S1878, C4S1879, C4S1884, C4S1887, C4S1889, C4S1890 e C4S4935;
- Cromossomo 5: C5S5133 e C5S5156;
- Cromossomo 6: C6S2981;
- Cromossomo 13: C13S320, C13S399, C13S431, C13S479, C13S505, C13S514, C13S522, C13S523, C13S524, C13S547 e C13S567;
- Cromossomo 14: C14S1718, C14S1729, C14S1734 e C14S1736; e
- Cromossomo 19: C19S4799, C19S4815 e C19S4831.

Vale ressaltar que apenas 7 entre os 39 SNPs verdadeiramente relevantes possuem genótipos diferentes em mais que 1% do total de indivíduos, ou seja, grande parte dos SNPs são variantes raras e dificilmente identificados por metodologias de seleção de variáveis tradicionais.

### 3.1 Resultados

Nessa seção aplicamos os três métodos discutidos para identificar e selecionar os SNPs relevantes no fenótipo simulado  $Q_1$ , disponível nos dados GAW17. Mostramos os resultados obtidos e comparamos as metodologias em termos de sensibilidade e especificidade na seleção dos SNPs.

### 3.2 Regressão linear simples

Na primeira metodologia, foram realizados testes de hipóteses para verificar a significância do coeficiente de regressão associado ao genótipo de cada SNP em um modelo de regressão linear simples. Consideramos os níveis de significância  $\alpha$  de 0.10, 0.05 e suas correções de Bonferroni. Os resultados obtidos nesses testes para cada cromossomo autossômico e para cada  $\alpha$  são expostos a seguir e comparados com os 39 SNPs que são realmente significativos para  $Q_1$  de acordo com Almasy et al. (2011).

Utilizando o nível de significância de 10%, o modelo identificou que todos os cromossomos autossômicos possuem SNPs significativos, detectando, no total, 3575 SNPs como significativos. Além disso, dos 39 SNPs que realmente tem efeito significativo sobre  $Q_1$ , apenas 16 foram identificados pelo modelo. Desse modo, 3559 SNPs foram identificados incorretamente como significativos pela metodologia utilizada. Para o nível de significância de 10%, a sensibilidade ( $S$ ) e a especificidade ( $E$ ) valem, respectivamente:

$$S = \frac{16}{39} = 0.410 \text{ e } E = \frac{20889}{24448} = 0.854$$

Os valores acima possuem a seguinte interpretação: 41% dos SNPs que são realmente significativos são identificados pelo modelo e 85.4% dos SNPs que não possuem efeito importante sobre o fenótipo não são identificados pelo modelo.

Aplicando agora o nível de significância de 5%, o modelo identificou que todos os cromossomos autossômicos possuem SNPs significativos, em um total de 2177 SNPs relevantes. Comparando esse valor com a quantidade de SNPs identificados com  $\alpha = 10\%$ , observamos uma redução na quantidade de SNPs selecionados, uma vez que, quanto menor o  $\alpha$ , mais restritiva é a condição para considerar a hipótese  $H_1$  como sendo verdadeira, e, sendo assim, menor a quantidade de SNPs falsos positivos. Dos 39 SNPs significativos, apenas 15 são identificados pelo modelo, ou seja, há 24 SNPs que são significativos e não foram identificados.

A sensibilidade e a especificidade ao nível de significância de 5% valem, respectivamente:

$$S = \frac{15}{39} = 0.385 \text{ e } E = \frac{22286}{24448} = 0.912$$

Ao fazer uso do nível de significância de 10% com correção de Bonferroni, foi apontado que os cromossomos 1, 2, 7, 8, 9, 10, 12, 13, 14, 15, 16, 17, 19, 20 e 22 possuem, ao todo, 42 SNPs que afetam o fenótipo  $Q_1$ . Entre os 42 SNPs selecionados, apenas 4 SNPs no cromossomo 13 foram corretamente identificados como significativos. São eles: C13S431, C13S522, C13S523 e C13S524. Os SNPs identificados como significativos pelo modelo foram: C1S1683, C1S11539, C2S2355, C7S1598, C7S1632, C7S4110, C8S2899, C9S4121, C10S103, C12S703, C12S704, C12S706, C12S707, C12S708, C12S711, C12S718, C12S719, C12S792, C12S831, C12S971, C12S2028, C12S2798, C13S431, C13S522, C13S523, C13S524, C14S1137, C15S1580, C16S3109, C17S3017, C17S4607, C19S1163, C19S3512, C19S3829, C19S5085, C19S5879, C20S667, C22S116, C22S1507, C22S1508, C22S1540 e C22S1901.

Além disso, ao calcularmos a sensibilidade e a especificidade ao nível de significância de 10% com correção de Bonferroni, temos que essas são, respectivamente, iguais a 0.102 e 0.998. Portanto, 10.2% dos SNPs que são realmente importantes para  $Q_1$  são identificados pelo modelo e 99.8% dos SNPs que de fato não têm efeito significativo sobre o fenótipo não são identificados pelo modelo de regressão linear simples.

Ao utilizar o nível de significância de 5% com correção de Bonferroni, os cromossomos 1, 2, 7, 8, 12, 13, 14, 15, 17, 19, 20 e 22 foram caracterizados como detentores de SNPs significativos, possuindo, ao todo, 27 SNPs significativos. Entre os 27 SNPs selecionados, apenas 2 SNPs no cromossomo 13 foram corretamente identificados como significativos, são eles: C13S522 e C13S523. A seguir estão listados quais SNPs o modelo identificou como significativos: C1S11539, C2S2355, C7S1598, C7S4110, C8S2899, C12S704, C12S707, C12S711, C12S719, C12S831, C12S971, C12S2028, C12S2798, C13S522, C13S523, C14S1137, C15S1580, C17S3017, C17S4607, C19S3512, C19S3829, C19S5085, C19S5879, C20S667, C22S1507, C22S1508 e C22S1901.

Ao calcular a sensibilidade e especificidade ao nível de significância de 5% corrigido por Bonferroni, a fim de comparar qual seria o melhor nível de significância para o modelo de regressão linear simples, temos que a proporção de estimativas de  $\beta$  diferentes de 0 dentre os SNPs que são realmente significativos e a proporção de estimativas de  $\beta$  iguais a zero dentre os SNPs que na realidade não são significativos são, respectivamente: 0.051 e 0.999.

Vale destacar que os 4 SNPs verdadeiro positivos identificados praticamente em todos os níveis de significância são os SNPs que simultaneamente apresentam *MAF* (frequência do menor alelo) maior que 1% e coeficientes de regressão maior que 0.60 na simulação do  $Q_1$ . Enquanto que SNPs com baixíssima frequência do alelo menor (maior parte dos SNPs realmente influentes) ou SNPs com efeito mais baixo no fenótipo não foram selecionados por essas metodologias. SNPs com coeficientes maiores que 1 na simulação,

tais como: C4S4935, C6S2981, C4S1889, e C4S1877, mas com  $MAF < 0.2\%$  não foram identificados por essa metodologia.

Apesar de ser falso positivo, o SNP que foi identificado no cromossomo 14 pelos 4 níveis de significância está próximo a 4 SNPs relevantes, assim como os 4 SNPs identificados no cromossomo 19 estão próximos a 3 SNPs influentes.

### 3.3 LASSO

Os resultados obtidos através do LASSO para cada cromossomo são apresentados a seguir e comparados com os SNPs que realmente influenciam  $Q_1$  segundo Almasy et al. (2011). A sensibilidade e especificidade atingidas por esse método também são expostas a seguir.

O modelo identificou que todos os cromossomos possuem SNPs significativos, detectando, ao todo, 799 SNPs como significativos. Ademais, dos 39 SNPs que são realmente significativos para  $Q_1$ , apenas 9 foram identificados pelo modelo, são eles: C4S1877, C4S1884, C4S1889, C6S2981, C13S320, C13S431, C13S522, C13S523 e C14S1734. Desse modo, 790 SNPs foram erroneamente identificados como relevantes pela metodologia utilizada. A sensibilidade ( $S$ ) e a especificidade ( $E$ ) valem, respectivamente:

$$S = \frac{9}{39} = 0.231 \text{ e } E = \frac{23658}{24448} = 0.968 \quad (7)$$

Os valores observados acima possuem a seguinte interpretação: 23.1% dos SNPs que são realmente significativos são identificados pelo modelo e 96.8% dos SNPs que não são relevantes para o fenótipo não foram identificados pelo modelo. Notamos que o LASSO identificou 3 SNPs verdadeiros do cromossomo 13 que também foram identificados pelo valor-p do teste de hipóteses do modelo de regressão simples e outros SNPs que, apesar de apresentarem  $MAF$  bem baixo, apresentam grande efeito sobre o fenótipo, tais como C6S2981 e C4S1889. Outros 4 SNPs com pequeno  $MAF$  e razoável efeito sobre o fenótipo também foram selecionados.

### 3.4 SPLS

Para conseguir determinar os valores ótimos de  $M$  e  $\eta$  para cada cromossomo usando a validação cruzada no SPLS, foi necessário excluir da base e não analisar SNPs que possuíam menos que 5 indivíduos com variação nos cromossomos 2, 4, 6, 7, 14, 16, 17, 19, 20, 21 e 22; 6 indivíduos com variação nos cromossomos 1, 3, 5, 9, 10, 11, 13 e 18; 7 indivíduos com variação nos cromossomos 12 e 15; e 8 indivíduos com variação no cromossomo 8. Essa remoção foi necessária pois, ao dividir a base original em subamostras, esses SNPs muito raros não apresentavam variabilidade em alguma ou algumas das subamostras e o algoritmo não funcionava. Ademais, na validação cruzada o valor de  $M$  variava de 1 a 20, porém para todos os cromossomos foi selecionado  $M$  igual a 1 como o valor ótimo de componentes latentes.

A seguir, estão listados por cromossomo a quantidade de SNPs que o cromossomo possuía, a quantidade que foi filtrada (devido a quantidade mínima de variações necessária) e o número de SNPs que foram identificados como significativos, respectivamente: Cromossomo 1: 2237, 738, 252; cromossomo 2: 1599, 595, 7; cromossomo 3: 1211, 349, 1; cromossomo 4: 944, 304, 8; cromossomo 5: 1074, 361, 47; cromossomo 6: 1425, 550, 45; cromossomo 7: 1063, 417, 4; cromossomo 8: 982, 333, 17; cromossomo 9: 1166, 517, 34; cromossomo 10: 1396, 549, 423; cromossomo 11: 2102, 959, 3; cromossomo 12: 1435, 526, 102; cromossomo 13: 425, 135, 2; cromossomo 14: 795, 262, 1; cromossomo 15: 933, 268, 1; cromossomo 16: 844, 308, 14; cromossomo 17: 1223, 470, 2; cromossomo 18: 634, 204, 4; cromossomo 19: 1649, 792, 164; cromossomo 20: 591, 239, 1; cromossomo 21: 251, 86, 21 e cromossomo 22: 508, 175, 10.

Vale destacar que, apesar de excluirmos apenas SNPs muito raros, 9137 (37.3% do total) não foram considerados por essa metodologia.

O modelo identificou que todos os cromossomos possuem SNPs significativos e, apesar de descartar uma quantidade grande de SNPs raros, detectou, ao todo, 1163 SNPs como significativos. Dos 39 SNPs que realmente possuem efeito sobre o fenótipo  $Q_1$ , apenas 4 foram identificados pelo modelo, são eles: C1S6533, C4S1884, C13S522 e C13S523, sendo que os dois SNPs do cromossomo 13 também foram identificados pelas duas metodologias anteriores. Logo, 1159 SNPs foram incorretamente identificados como significativos pela metodologia implementada. A sensibilidade ( $S$ ) e a especificidade ( $E$ ) valem, respectivamente:

$$S = \frac{4}{39} = 0.102 \text{ e } E = \frac{23289}{24448} = 0.952 \quad (8)$$

Assim, 10.2% dos SNPs que são realmente significativos são identificados pelo modelo e 95.2% dos SNPs que não são relevantes para o fenótipo não foram identificados pelo modelo.

Vale destacar que o SPLS apenas identificou os SNPs verdadeiro positivos que apresentam  $MAF > 1\%$  e moderado efeito sobre o fenótipo e selecionou uma grande quantidade de falso positivos, mesmo que quase 40% dos SNPs não tenham sido analisados.

### 3.5 Comparações dos resultados obtidos pelos métodos empregados

Com o propósito de verificar o desempenho das metodologias para selecionar os SNPs que influenciam o fenótipo  $Q_1$ , comparamos o número de SNPs identificados como significativos e os valores de sensibilidade e especificidade, todos advindos de cada um dos métodos estudados anteriormente. Os valores são mostrados na Tabela 1.

	Regressão Linear Simples				LASSO	SPLS
	$\alpha = 10\%$	$\alpha = 5\%$	$\alpha_{BF} = 10\%$	$\alpha_{BF} = 5\%$		
sensibilidade	0.410	0.385	0.102	0.051	0.231	0.102
especificidade	0.854	0.912	0.998	0.999	0.968	0.952
n° de SNPs identificados como significativos	3575	2177	42	27	799	1163

Tabela 1: Valores de sensibilidade, especificidade e número de SNPs identificados como significativos de cada um dos métodos estudados. Aqui  $\alpha_{BF} = 10\%$  representa o nível de significância de 10% corrigido por Bonferroni.

Através da Tabela 1, notamos que a regressão linear simples com nível de significância de 5% corrigido por Bonferroni possui a melhor especificidade dentre os métodos comparados, pois seu valor é praticamente 1. Em contrapartida, a sua sensibilidade é baixa, indicando que esse modelo não considera como significativo grande parte dos SNPs que realmente influenciam o fenótipo. Notamos ainda que a regressão linear simples com nível de significância de 10% é a metodologia que possui maior sensibilidade dentre as utilizadas, embora seu valor para a especificidade seja o mais baixo de todos e isso se deve à grande quantidade de SNPs falso positivos.

Ao analisarmos os SNPs que são realmente significativos para o fenótipo  $Q_1$ , percebemos que os SNPs C13S522 e C13S523 foram identificados por todos os métodos utilizados e o SNP C4S1884 foi identificado pelo LASSO, pela regressão linear simples aos níveis de significância de 5 e 10% e pelo SPLS. Por sua vez, o SNP C1S6533 foi identificado como significativo pelo SPLS e pela regressão linear simples LASSO  $\alpha$  de 5% e de 10%. A identificação desses SNPs por mais de um método pode ocorrer devido à maior variação dos seus genótipos na base de dados ou devido ao elevado valor do efeito desses marcadores, os quais valem, respectivamente, 0.623466, 0.653351, 0.318125 e 0.589734, de modo que se distanciam de zero (caso em que o SNP não é considerado relevante para o fenótipo em questão).

Comparado com o SPLS e com a regressão a nível 5%, 10%, 10% Bonferroni e 5% Bonferroni, o LASSO é o método que apresenta simultaneamente bons valores de sensibilidade e especificidade, além de ter identificado (mesmo sem aumentar consideravelmente o número de falso positivos) SNPs influentes com baixo *MAF*. O SPLS, por sua vez, apresentou o pior desempenho no geral, com baixa sensibilidade e alto número de SNPs falso positivos, mesmo não considerando quase 40% dos SNPs.

A escolha entre utilizar LASSO ou regressão a um nível 5% para esse conjunto de dados ou conjuntos que tenham características parecidas depende de quão tolerante está o pesquisador em aumentar o número de SNPs falso positivos para identificar mais SNPs verdadeiro positivos.



## 4 | CONCLUSÃO E DISCUSSÃO

Nesse trabalho, estudamos, aplicamos e avaliamos o desempenho de três diferentes metodologias para identificação de SNPs relevantes no fenótipo  $Q_1$  dos dados GAW17.

Ao analisarmos os desempenhos das metodologias SPLS, LASSO e regressão linear simples ao nível de significância de 5% e 10%, ambos com e sem a correção de Bonferroni, todas utilizadas para selecionar os SNPs que influenciam o fenótipo de interesse, notamos que o LASSO além de apresentar um equilíbrio entre a sensibilidade e a especificidade, ou seja, identifica relativamente bem os SNPs que realmente possuem efeito sobre  $Q_1$  e não identifica aqueles que não são significativos, também identifica SNPs que afetam  $Q_1$  e que são muito raros (possuem  $MAF < 1\%$ ), são eles: C13S320, C4S1877, C4S1889 e C6S2981.

Os níveis de significância 5% e 10% corrigido por Bonferroni do teste de hipóteses em um modelo de regressão linear simples também apresentaram performance razoável. O SPLS, apesar de reduzir dimensão e selecionar variáveis simultaneamente, teve o pior desempenho no geral, além de não conseguir analisar SNPs muito raros.

Esse trabalho foi parcialmente financiado pelo Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) através da Bolsa PIBIC.

## REFERÊNCIAS

Abdi, H. (2007), 'Bonferroni and Šidák corrections for multiple comparisons', *Encyclopedia of Measurement and Statistics* **3**, 103–107. 4

Almasy, L., Dyer, T. D., Peralta, J. M., Kent, J. W., Charlesworth, J. C., Curran, J. E. e Blangero, J. (2011), Genetic Analysis Workshop 17 mini-exome simulation, in 'BMC Proceedings', Vol. 5, BioMed Central, p. S2. 2, 8, 10, 11

Breiman, L. (2001), 'Random forests', *Machine Learning* **45**(1), 5–32. 2

Chun, H. e Keles, S. (2010), 'Sparse partial least squares regression for simultaneous dimension reduction and variable selection', *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **72**(1), 3–25.

Chung, D. e Keles, S. (2010), 'Sparse partial least squares classification for high dimensional data', *Statistical Applications in Genetics and Molecular Biology* **9**(1). 6

Feng, Z. Z., Yang, X., Subedi, S. e McNicholas, P. D. (2012), 'The lasso and sparse least squares regression methods for SNP selection in predicting quantitative traits', *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)* **9**(2), 629–636. 2, 7

Goldberg, D. E. (1989), 'Genetic algorithm in search optimization and machine learning', *Addison Wesley*. 2

Ióca, M. P. e Zuanetti, D. A. (2021), 'Selection of SNP markers: Analyzing GAW17 data using different methodologies', *Brazilian Journal of Biometrics* **39**(1), 71–88. 5, 8

Lee, D., Lee, W., Lee, Y. e Pawitan, Y. (2011), 'Sparse partial least-squares regression and its applications to high-throughput data analysis', *Chemometrics and Intelligent Laboratory Systems* **109**(1), 1–8. 2, 7

Morettin, P. A. e Bussab, W. O. (2017), *Estatística básica*, Saraiva Educação SA. 3

Oliveira, F. C. (2015), 'Um método para seleção de atributos em dados genômicos', *Tese de Doutorado. Programa de Pós-graduação em Modelagem Computacional da Universidade Federal de Juiz de Fora* . 2

Park, T. e Casella, G. (2008), 'The Bayesian lasso', *Journal of the American Statistical Association* **103**(482), 681–686. 2

Rodrigues, K. A. S. (2018), Lasso clássico e bayesiano, Technical report, IME-USP. 5

Tibshirani, R. (1996), 'Regression shrinkage and selection via the lasso', *Journal of the Royal Statistical Society: Series B (Methodological)* **58**(1), 267–288.

Yazdani, A. (2014), 'Statistical approaches in genome-wide association studies', *Tese de*

*Doutorado. Dipartimento di Scienze Statistiche - Scuola de Dottorato di Ricerca in Scienze Statistiche* . 2

**ANIELE DOMINGAS PIMENTEL SILVA** - Possui graduação em Licenciatura Plena em Matemática e especialização em Educação Matemática pela Universidade Federal do Pará - UFPA. Mestre em Educação pela Universidade Federal do Oeste do Pará - UFOPA, na linha de pesquisa “Práticas Educativas, Linguagens e Tecnologias”, com ênfase em Modelagem Matemática e Tecnologias, atualmente é doutoranda pelo Programa de Pós-Graduação em Educação na Amazônia PGEDA/EDUCANORTE (2022) da Universidade Federal do Oeste do Pará – UFOPA, atuando na linha de pesquisa “Educação na Amazônia: formação do educador, práxis pedagógica e currículo”. Integra o Grupo de Estudos e Pesquisas em Educação Matemática e Interdisciplinaridade na Amazônia - GEPEIMAZ. Tem experiência como professora de matemática na educação básica pela Secretaria Municipal de Educação de Santarém-PA, professora colaboradora na UFOPA no programa PARFOR nos cursos de Licenciatura integrada em Matemática e Física e professora substituta no Instituto Federal do Amapá – IFAP em turmas do ensino médio integrado e de ensino superior.

**A**

Álgebra 53, 56, 57, 58, 60, 61, 65, 67, 98, 99, 101, 103, 105, 109

**B**

Banco de dados relacionais 98, 99, 100, 101, 103, 109

**C**

Conta de energia elétrica 20, 22, 24, 27, 29, 30, 31, 33, 34, 35, 36

**D**

Desenvolvimento cognitivo 3, 4, 12, 38

Discalculia 111, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124

**E**

Educação Matemática 1, 2, 18, 19, 20, 21, 23, 36, 43, 45, 52, 66, 67, 68, 88, 90, 92, 110, 116, 117, 123, 141

Ensino/aprendizagem 1, 17

Ensino de funções 37, 39

Ensino de Matemática 44, 46, 47, 50, 54, 57, 66, 87, 90, 121

Erros 5, 6, 9, 10, 12, 16, 17, 18, 46, 68, 69, 74, 75, 81, 82, 83, 95, 113, 117, 130, 131

Experiência 3, 48, 49, 50, 53, 54, 56, 61, 69, 71, 77, 79, 80, 84, 85, 90, 98, 107, 141

**F**

Ferramenta de ensino 13, 14, 16

Formação 2, 23, 24, 26, 39, 40, 42, 47, 51, 55, 68, 69, 70, 72, 73, 74, 75, 76, 77, 78, 79, 81, 83, 84, 85, 86, 88, 90, 91, 116, 141

Função afim 20, 22, 24, 27, 28, 30, 31, 33, 34, 35, 36

**G**

Geometria dinâmica 37, 38, 39

**I**

Identidade 68, 69, 70, 71, 72, 73, 74, 75, 76, 79, 80, 81, 82, 83, 84, 85, 86

**J**

Jogo Batalha Cartesiana 1, 8, 9, 10, 17

Jogos matemáticos 1, 2, 3, 13, 114, 123

**L**

LASSO 125, 126, 127, 128, 129, 130, 136, 138, 139, 140

Linguagem matemática 43, 56, 57, 58, 59, 60, 65, 66, 113

**M**

Manual pedagógico 87, 89, 91, 92, 96

Matemática 1, 2, 3, 4, 7, 13, 14, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 30, 31, 33, 35, 36, 37, 38, 43, 44, 45, 46, 47, 48, 49, 50, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 64, 65, 66, 67, 68, 69, 73, 74, 75, 76, 77, 78, 81, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 94, 95, 96, 97, 99, 108, 109, 110, 111, 113, 114, 116, 117, 120, 121, 122, 123, 124, 141

Matemática a ensinar 87, 91, 94, 96

Matemática para ensinar 87, 88, 89, 90, 91, 92, 94, 95, 96, 97

Material dourado 56, 61, 62, 63, 65, 66, 67

Metodologia de ensino 20, 26, 27

Modelagem Matemática 2, 20, 22, 23, 24, 25, 26, 27, 30, 33, 35, 36, 141

**O**

Obstáculos epistemológicos 44, 45, 47, 48, 49, 50, 51, 52, 53, 54, 55

Operações básicas 87, 88, 89, 90, 91, 92, 96, 97, 113

**P**

Pensamento computacional 26, 111, 112, 113, 115, 116, 118, 119, 122, 123, 124

Plano cartesiano 1, 2, 3, 7, 8, 10, 12, 15, 17, 18, 31, 35, 37, 39

Prática 25, 33, 43, 49, 55, 58, 61, 65, 69, 70, 78, 79, 80, 83, 84, 85, 91, 93, 95, 100, 110, 118, 123

Produtos notáveis 56, 58, 61, 62, 63, 65, 66

**R**

Rupturas do conhecimento 44, 46

**S**

Seleção de variáveis 132, 134

Sequência de atividades 36, 37, 38, 42

Sequência didática adaptativa 98, 99

SPLS 125, 126, 127, 130, 131, 136, 137, 138, 139

**T**

Técnico em informática 98, 109

Tecnologia educacional 37

Tendências em educação Matemática 18, 36

Teoria dos conjuntos 98, 99, 102, 103, 105, 109

Teste de significância 127

Trigonometria 37, 38, 39

**V**

Variantes raras 126, 134



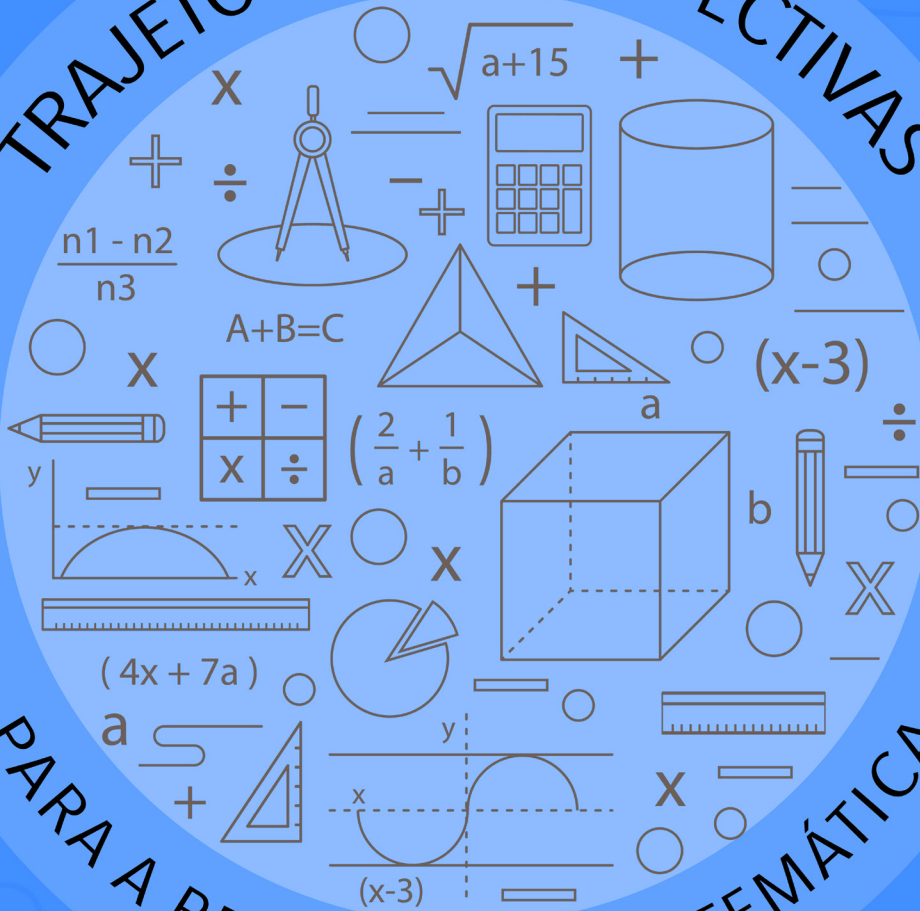
www.atenaeditora.com.br

contato@atenaeditora.com.br

@atenaeditora

www.facebook.com/atenaeditora.com.br

# TRAJETÓRIAS E PERSPECTIVAS



# PARA A PESQUISA EM MATEMÁTICA