



UNIVERSIDADE FEDERAL DA BAHIA
INSTITUTO DE MATEMÁTICA E ESTATÍSTICA
PROGRAMA DE PÓS-GRADUAÇÃO EM MATEMÁTICA
DISSERTAÇÃO DE MESTRADO



NOVOS MODELOS DE REGRESSÃO BINÁRIA USANDO
FUNÇÕES DE LIGAÇÃO SIMÉTRICAS E ASSIMÉTRICAS

TATIANA FELIX DA MATTA

Área de Concentração: ESTATÍSTICA

Salvador - Bahia
AGOSTO DE 2021

NOVOS MODELOS DE REGRESSÃO BINÁRIA USANDO FUNÇÕES DE LIGAÇÃO SIMÉTRICAS E ASSIMÉTRICAS

TATIANA FELIX DA MATTA

Dissertação de Mestrado apresentada ao Colegiado da Pós-Graduação em Matemática da Universidade Federal da Bahia (UFBa), como parte dos requisitos para obtenção do título de Mestre em Matemática. Área de concentração: Estatística.

Orientador: Prof. Dr. Paulo Henrique Ferreira da Silva

Coorientador: Prof. Dr. Anderson Luiz Ara Souza

Salvador - Bahia

Agosto de 2021

NOVOS MODELOS DE REGRESSÃO BINÁRIA USANDO FUNÇÕES DE LIGAÇÃO SIMÉTRICAS E ASSIMÉTRICAS

TATIANA FELIX DA MATTA

Dissertação de Mestrado apresentada ao Colegiado da Pós-Graduação em Matemática da Universidade Federal da Bahia (UFBA), como parte dos requisitos para a obtenção do título de Mestre em Matemática. Área de concentração: Estatística.

Banca Examinadora:

Prof. Dr. Paulo Henrique Ferreira da Silva (Orientador)
IME - UFBA

Prof. Dr. Paulo Jorge Canas Rodrigues
IME - UFBA

Prof. Dr. Francisco Louzada Neto
ICMC - USP

Dedico este trabalho à minha mãe Maria das Graças (*in memoriam*).

Agradecimentos

Agradeço a **Deus**, por me fortalecer a cada dia nos momentos difíceis que, por muitas vezes, pareciam não ter saída. Sem Ele não teria sabedoria para escolher os caminhos certos para conquistar meus objetivos.

Meu afetuoso agradecimento aos meus pais (*in memoriam*), especialmente à minha mãe Maria das Graças, que se dedicou por toda sua vida a me apoiar, me passando tranquilidade e confiança para que eu pudesse atingir meus objetivos. Mãe, serei eternamente grata pelo que proporcionastes na minha vida.

Ao meu avô Leonardo Felix (*in memoriam*), por todo amor, compreensão, esforço e por sempre me incentivar a seguir em frente quando encontrei dificuldades. E também por ter cuidado de mim com muita dedicação e carinho.

Aos meus irmãos de coração Adriana Campos e Luiz Carlos Junior, por sempre se fazerem presentes na minha vida, pela amizade, companheirismo e por estarem sempre preocupados com meu bem-estar. Às minhas tias e primos, em especial, aos meus primos Hugo Leonardo, Jaqueline, Marcelle, Mikael e Theodoro, que sempre torceram por mim, emanando energias positivas no meu caminho para que essa etapa se concretizasse em minha vida. Tenho muita sorte por ter pessoas tão especiais na minha vida.

Ao meu afilhado Paulo Victor que, através da sua simplicidade de criança, seus gestos e sorriso, me ajuda a renovar as esperanças em vários momentos. Como sempre te falo: “Dindo”, você é o melhor afilhado do mundo, o amor da minha vida!

Meus sinceros agradecimentos ao meu orientador, Prof. Dr. Paulo Henrique Ferreira, pela confiança depositada em mim, pela paciência, atenção, por ter compartilhado seu conhecimento e pelas palavras de apoio durante o desenvolvimento da dissertação. **Muito obrigada!** Ao meu coorientador, Prof. Dr. Anderson Ara, pelas sugestões que contribuíram para a realização deste trabalho.

Um agradecimento especial à minha amiga Laís Sacramento, minha dupla de estudo, de caminhada desde a graduação, por todas as chamadas, as conversas (até durante a madrugada), por sempre ter uma palavra de conforto e apoio. Foram momentos de dificuldade, alegria e ajuda mútua. A toda sua família pelo acolhimento, por me fazer sentir parte dela.

Aos meus amigos da graduação em Estatística, Adriano Cruz, Alexandro Teles e Aline Souza, pela amizade, pelo incentivo e por sempre estarem dispostos a me ajudar. Muito obrigado aos meus colegas de mestrado, pelos momentos de descontração durante os estudos, o que fizeram essa caminhada mais prazerosa.

Agradeço à Universidade Federal da Bahia, onde fiz minha vida acadêmica. Aos professores do Departamento de Estatística e professores da pós-graduação, pela oportunidade de assistir suas aulas e adquirir valiosos conhecimentos. Aos funcionários do IME, pela disponibilidade e ajuda.

Aos professores membros da banca, Prof. Dr. Francisco Louzada e Prof. Dr. Paulo Canas Rodrigues, por disponibilizarem seu tempo em avaliar este trabalho.

O presente trabalho foi desenvolvido com apoio da Coordenação de Aperfeiçoamento de Pessoal de Nível Superior – Brasil (CAPES) - Código de Financiamento 001.

*“Deus, se um dia eu perder as
esperanças, ajude-me a lembrar que os
teus planos são melhores que os
meus.”*

Chico Xavier

Resumo

Os modelos de regressão com variáveis respostas binárias (1 - ocorrência do evento de interesse ou “sucesso”, 0 - não ocorrência do evento de interesse ou “fracasso”) têm sido aplicados intensamente em diversas áreas do conhecimento, como saúde, finanças, indústria, entre outras. Tradicionalmente, o modelo mais usado na regressão binária tem sido o modelo de regressão logística. Contudo, ele utiliza a função de ligação *logit* (ou *logito*), a qual é uma função de ligação simétrica e pode não ser adequada em determinadas situações, por exemplo, quando uma das classes da variável resposta é desbalanceada em relação à outra (conjunto de dados desbalanceados). Este trabalho tem como objetivo principal apresentar novos modelos de regressão binária usando funções de ligação simétricas e assimétricas. A estimação dos parâmetros dos modelos descritos neste trabalho (a saber: modelos de regressão binária double Lindley, double Lindley assimétrica, potência double Lindley e reversa de potência double Lindley) é feita pelo método clássico da máxima verossimilhança. Para comparação e seleção do “melhor” modelo dentre as diferentes distribuições, são empregados critérios de informação (AIC e BIC) e medidas de avaliação da capacidade preditiva (AUC, acurácia balanceada, sensibilidade, especificidade, valores de predição positivo e negativo, *F1-Score*, coeficiente de correlação de Matthews, dentre outras). Os resultados da análise de dois conjuntos de dados reais, um sobre câncer de mama, obtido do UCI (*University of California, Irvine*) *Machine Learning Repository*, e outro referente a uma competição promovida pelo Banco Santander para a comunidade do *Kaggle*, mostram que os modelos com as funções de ligação propostas podem apresentar melhor ajuste e performance preditiva do que os modelos usando ligações tradicionais, como a *logito*.

Palavras-chave: Dados desbalanceados, distribuição double Lindley, ligações potência e reversa de potência, método de máxima verossimilhança, performance preditiva.

Abstract

Regression models with binary response variables (1 - occurrence of the event of interest or “success”, 0 - non-occurrence of the event of interest or “failure”) have been intensively applied in several areas of knowledge, such as health, finance, industry, among others. Traditionally, the most used model in binary regression has been the logistic regression model. However, it uses the logit link function, which is a symmetric link function and may not be suitable in certain situations, for example, when one of the classes of the response variable is disproportionate to the other (imbalanced data set). The main aim of this work is to present new binary regression models using symmetric and asymmetric link functions. The parameter estimation of the models (namely, the double Lindley, asymmetric double Lindley, power double Lindley, and reversal power double Lindley binary regression models) is performed with the classical maximum likelihood method. In order to compare and select the “best” model among the different distributions, information criteria (AIC and BIC) and measures of predictive performance (AUC, balanced accuracy, sensitivity, specificity, positive and negative predictive values, F1-Score, Matthews correlation coefficient, among others) are used. Through the analysis of two real data sets, one on breast cancer, obtained from the University of California, Irvine’s (UCI) Machine Learning Repository, and another on a competition promoted by Santander Bank for the Kaggle community, we show that models using the proposed link functions can provide a better fit and predictive ability than models using standard links, such as logit.

Keywords: Imbalanced data, double Lindley distribution, power and reversal power link functions, maximum likelihood method, predictive performance.

Sumário

1	Introdução	1
1.1	Objetivo	3
1.2	Organização do Trabalho	3
2	Modelos de Regressão para Dados Binários	4
2.1	Regressão Binária	4
2.2	Dados Desbalanceados	5
2.3	Estimação por Máxima Verossimilhança	6
2.4	Distribuições de Probabilidade Simétricas e Assimétricas	7
2.4.1	Distribuição Logística	8
2.4.2	Distribuição de Lindley	8
2.4.2.1	Distribuição Double Lindley	9
2.4.2.2	Distribuição Double Lindley Assimétrica	11
2.4.2.3	Distribuição Potência Double Lindley	14
2.4.2.4	Distribuição Reversa de Potência Double Lindley	15
2.5	Comparação de Modelos	18
2.5.1	Critério de Informação de Akaike	19
2.5.2	Critério de Informação Bayesiano	19
2.6	Métricas de Desempenho	20
2.7	Método de Validação	24
3	Estudos de Simulação	26
3.1	Recuperação de Parâmetros	26
3.2	<i>Misspecification</i>	29
3.3	Avaliação da Capacidade Preditiva	32
4	Aplicações	39
4.1	<i>Breast Cancer Wisconsin (Diagnostic) Data Set</i>	39
4.2	<i>Santander Customer Transaction Prediction Dataset</i>	43

5	Considerações Finais e Trabalhos Futuros	47
	Referências Bibliográficas	49

Capítulo 1

Introdução

Os problemas de classificação de dados são bastante comuns em diversas áreas do conhecimento, como indústria, saúde e finanças. Muitas vezes, a variável de interesse, denominada variável dependente ou resposta, é considerada binária e assume apenas uma de duas categorias (níveis ou classes) possíveis. Como, por exemplo, o *status* de um item produzido (bom estado ou defeituoso), a remissão de uma doença (sim ou não), o resultado de um tratamento (bom ou ruim), se um cliente se tornará inadimplente (sim ou não), dentre outros. Nestes casos, o problema de classificação é dito ser de classificação binária ou dicotômica. Frequentemente, um conjunto de variáveis que influenciam a resposta de interesse, chamadas de variáveis independentes ou explicativas (ou ainda, covariáveis), também está disponível. Essas variáveis podem ser tanto qualitativas (sexo, raça, estado civil, grau de escolaridade, etc.) quanto quantitativas (idade, altura, peso, renda, etc.).

Dentre as ferramentas de modelagem estatística que têm sido amplamente utilizadas no auxílio à tomada de decisões em situações como as descritas anteriormente (isto é, qual tratamento escolher, se irá conceder ou não o crédito solicitado, etc.), estão os modelos de regressão (ou classificação) binária.

Nos modelos de regressão em geral, o objetivo principal é descrever uma possível relação existente entre a variável resposta Y e as covariáveis X_1, \dots, X_p , com $p \geq 1$. Nos modelos de regressão binária, a variável resposta Y é dicotômica, isto é, permite somente dois resultados, aos quais atribuímos convencionalmente o valor 1 para a ocorrência do evento de interesse (“sucesso”), e o valor 0 para a não ocorrência do evento de interesse (“fracasso”). Estudos mais detalhados sobre a regressão binária podem ser encontrados em Cox & Snell [28] e Collett [26], dentre outros.

O modelo de regressão logística binária (ou simplesmente, modelo de regressão logística) é conhecido desde os anos 1950. No entanto, estudos iniciais sobre esse modelo foram publicados nos artigos de Verhulst [79, 80, 81]. Bliss [15], em seu artigo, estudou experimentos biológicos do tipo dose-resposta para doses fixas e respostas aleatórias que

refletiam a distribuição individual de níveis de tolerância. Além disso, aplicações do modelo de regressão logística em áreas como economia e pesquisa de mercado surgiram nos anos 1950 e 1960 (Farrell [33], Aitchison & Brown [2], Adam [1]). A partir da década de 1980, torna-se mais utilizado com o trabalho de Cox & Snell [28]. Esse modelo tem sido aplicado em várias áreas do conhecimento para a análise de dados binários, pois não precisa atender alguns pressupostos, como a igualdade das matrizes de covariância e a normalidade dos erros. Ele traz, ainda, a vantagem da facilidade de interpretação dos parâmetros, conforme apontam Hosmer *et al.* [43].

Na regressão logística, a função de ligação empregada é a *logit* (ou logito), que é uma função de ligação simétrica, resultante da função de distribuição acumulada (FDA) da distribuição logística padrão. Outras funções de ligação simétricas são comumente usadas na regressão binária, como por exemplo, a *cauchit* (ou cauchito) e a *probit* (ou probito), que resultam das FDAs das distribuições Cauchy padrão e normal padrão, respectivamente.

Considere n observações de uma variável aleatória independente (Y_1, Y_2, \dots, Y_n) com distribuição de Bernoulli, com probabilidade de sucesso μ_i e probabilidade de fracasso $1 - \mu_i$, para $i = 1, 2, \dots, n$. Para associar a probabilidade μ_i com as covariáveis X_{i1}, \dots, X_{ip} , para $i = 1, 2, \dots, n$ e $p \geq 1$, as funções de ligação mais empregadas são a logito e a probito (ambas simétricas); entretanto, elas podem não se legitimar caso os dados (isto é, as classes da variável resposta) sejam desbalanceados e levar a conclusões erradas. Neste contexto, Czado & Santner [29] pesquisaram os efeitos da má especificação na função de ligação, concluindo que supor a logito como função de ligação no ajuste quando a função de ligação dos dados é outra, traz viés na estimação dos parâmetros da regressão e nas probabilidades preditas.

De fato, quando lidamos com dados binários, é frequente a presença de uma variável resposta cujo “sucesso” é pouco provável de ocorrer, ou seja, temos um evento raro (amostra ou base de dados desbalanceados). Neste caso, os estimadores de máxima verossimilhança podem não fornecer resultados satisfatórios (boas estimativas) para os parâmetros (ou coeficientes) do modelo de regressão, nem produzir boas previsões.

Contudo, as abordagens discutidas anteriormente não são as únicas possíveis e, dada a importância de se identificar os eventos raros, é imprescindível que se entenda melhor como obter um modelo com bom poder preditivo frente às limitações de um banco de dados com tais características, seja na área da saúde, indústria, finanças ou em qualquer outra área.

1.1 Objetivo

O objetivo geral desta dissertação é desenvolver novos modelos de regressão binária usando funções de ligação simétricas e assimétricas ¹. A metodologia proposta (novos modelos de regressão binária, métodos de estimação e validação, medidas de comparação de modelos e avaliação da performance preditiva) é ilustrada em dois conjuntos de dados reais: o primeiro refere-se ao diagnóstico de câncer de mama (tumores benignos e malignos) em pacientes do Hospital Universitário de Wisconsin, nos Estados Unidos (*Breast Cancer Wisconsin Data Set*), e o segundo diz respeito a dados do Banco Santander (*Santander Customer Transaction*), no qual o objetivo era prever se um cliente iria efetivar uma operação financeira específica futuramente, independente da quantia transacionada. Também são realizados estudos de simulação para avaliar a performance do método de estimação considerado (máxima verossimilhança), a adequação das medidas de comparação/seleção de modelos utilizadas (AIC e BIC), assim como o desempenho preditivo dos novos modelos de classificação binária.

1.2 Organização do Trabalho

Este trabalho está estruturado em cinco capítulos, dispostos da seguinte forma. No Capítulo 2 é realizada uma breve revisão sobre alguns conceitos fundamentais, tais como regressão binária, funções de ligação simétricas e assimétricas, método de estimação por máxima verossimilhança, critérios para comparação de modelos e métricas de desempenho preditivo. Além disso, são apresentados neste capítulo os modelos de regressão binária propostos, a saber: modelos de regressão double Lindley, double Lindley assimétrica, e os novos modelos potência double Lindley e reversa de potência double Lindley. No Capítulo 3 são mostrados os resultados dos estudos de simulação (recuperação de parâmetros, *misspecification* e avaliação da capacidade preditiva) e a forma como foram realizados, considerando todos os modelos propostos. No Capítulo 4 são exibidos os principais resultados da aplicação da metodologia proposta a dois conjuntos de dados reais, sendo um da área de saúde (referente ao diagnóstico de câncer de mama em pacientes de um hospital norte-americano) e outro da área financeira (que trata sobre possíveis transações bancárias efetuadas por clientes de uma instituição financeira). As considerações finais deste trabalho estão no Capítulo 5.

¹Uma vez que informações a respeito de tais modelos não foram encontradas em pesquisas na literatura, acredita-se que todos eles sejam realmente inéditos.

Capítulo 2

Modelos de Regressão para Dados Binários

Neste capítulo são apresentados o modelo de regressão binária usual (regressão logística), assim como detalhes acerca das distribuições de probabilidade (double Lindley, double Lindley assimétrica, potência double Lindley e reversa de potência double Lindley) consideradas nos modelos propostos neste trabalho. São também descritos o procedimento de estimação clássica pelo método da máxima verossimilhança e as diferentes medidas de comparação de modelos, incluindo os critérios de informação (AIC e BIC) e as métricas de avaliação da performance preditiva.

2.1 Regressão Binária

Os modelos de regressão binária têm sido aplicados intensamente em várias áreas do conhecimento, tais como indústria (Pacagnella *et al.* [60]), saúde (Silva *et al.* [72]), finanças (Ritta *et al.* [30]), etc. Eles são indicados quando a variável resposta é dicotômica (= 1 para o evento ou resultado de interesse - “sucesso”; = 0 para o evento complementar - “fracasso”). A variável resposta está geralmente associada a outras variáveis, que podem ser discretas, contínuas ou categóricas. Sendo que a probabilidade de “sucesso” pode ser explicada por essas variáveis, denominadas variáveis explicativas ou covariáveis.

Seja $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^\top$ um vetor $n \times 1$ de variáveis aleatórias respostas, $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ik})^\top$ um vetor $k \times 1$ de covariáveis associadas a Y_i , $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_k)^\top$ um vetor $k \times 1$ de coeficientes de regressão associados às covariáveis e, finalmente, considere a probabilidade de “sucesso” $P(Y_i = 1) = \mu_i$ e a probabilidade de “fracasso” $P(Y_i = 0) = 1 - \mu_i$, para $i = 1, 2, \dots, n$. Logo, no modelo de regressão binária, Y_i segue uma distribuição de Bernoulli com parâmetro μ_i , conforme especificado a seguir:

$$Y_i \sim \text{Bernoulli}(\mu_i), \quad \text{para } i = 1, 2, \dots, n,$$

$$\mu_i = P(Y_i = 1) = F(\eta_i),$$

$$\eta_i = F^{-1}(\mu_i),$$

em que $F(\cdot)$ é a FDA correspondente, $F^{-1}(\cdot)$ é a função de ligação e $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta} = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}$.

Observe que a função de ligação $F^{-1}(\cdot)$ transforma o intervalo $(0, 1)$, isto é, o suporte de μ_i (a média de Y_i), para a linha dos reais. Assim, quando o preditor linear η , com valores nos reais, é avaliado na FDA $F(\cdot)$ (também referida como função de ligação inversa), os resultados obtidos são valores de probabilidade válidos, que estão entre 0 e 1. A função de ligação logito, por exemplo, é obtida da inversa da FDA da distribuição logística padrão. Assim, o modelo de regressão logística é determinado pela seguinte relação:

$$\log\left(\frac{\mu_i}{1 - \mu_i}\right) = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}.$$

2.2 Dados Desbalanceados

Na regressão binária com ligações simétricas, que são as mais comumente usadas, sendo elas originárias da inversa da FDA padrão de distribuições simétricas (por exemplo, normal, logística, Cauchy, t -Student e Laplace, dentre outras), um obstáculo difícil existe quando uma das classes da variável resposta é desbalanceada em relação à outra. Segundo Van der Paal [78], um conjunto de dados é considerado desbalanceado (ou com raridade relativa) quando a classe de interesse (evento) é consideravelmente menor do que a outra classe (não evento), ou seja, o número de observações referente a uma classe é desproporcional ao número de observações da outra classe.

O desbalanceamento dos dados influencia a estimação dos parâmetros do modelo de regressão binária, no que diz respeito ao vício, erro padrão e erro quadrático médio das estimativas obtidas, assim como as predições das métricas de desempenho, como por exemplo, a acurácia (ou taxa de acerto geral). Por isso, quando são considerados dados desbalanceados, é possível incidir em um erro, visto que o acerto geral pode estar próximo à totalidade de 100%, entretanto, o acerto na classe de interesse pode ficar próximo de zero.

Considerando modelos de classificação para dados desbalanceados, segundo Ali *et al.* [5], uma das abordagens que vem sendo utilizada é trabalhar com o nivelamento dos dados, que consiste em técnicas de reamostragem dos dados visando a redução do desbalanceamento. Uma dessas técnicas é o *undersampling* (ou subamostragem), que trata-se do processo de redução da quantidade de eventos em maior número, de forma que o conjunto de dados alcance uma quantidade de registros considerável para resolução

do problema em questão. Outro método é o *oversampling* (o superamostragem), que consiste no incremento da classe minoritária, considerando um processo de reamostragem da classe até que se alcance o resultado desejado. Com a utilização desta abordagem, pode-se realizar desde um processo de reamostragem aleatória simples do banco de dados até métodos mais elaborados, como o SMOTE (*Synthetic Minority Oversampling Technique*) proposto por Chawla *et al.* [24]. Outro método citado na literatura (que é baseado no SMOTE) é o ADASYN (*Adaptive Synthetic Sampling*), proposto por He *et al.* [41], que consiste na utilização de uma distribuição de densidade como um critério de decisão sobre a quantidade de dados sintéticos que devem ser gerados.

Quando o banco de dados possui grande número de observações, a utilização de *undersampling* com amostragem aleatória simples é muito comum no mercado de trabalho brasileiro, uma vez que é de fácil explicação e reduz o potencial de *overfitting* que, segundo Padhi *et al.* [61], ocorre quando uma função explica muito bem um determinado conjunto de observações, mas não possui a capacidade de generalização para outros bancos de dados com características similares.

Outra abordagem que vem sendo empregada, segundo Ali *et al.* [5], é uma metodologia de predição que, por natureza de sua criação, pode gerar resultados satisfatórios para dados desbalanceados, sem necessidade de realização de amostragens, como é o caso, por exemplo, do algoritmo SVM (*Support Vector Machines*). Entretanto, quando o desbalanceamento é alto, o desempenho do classificador SVM pode ser afetado (i.e., decrementado) significativamente (Tao *et al.* [75]; Imam *et al.* [46]).

Este trabalho traz algo que faz um contraponto às técnicas mencionadas anteriormente. O objetivo é construir novos modelos de classificação binária que possuam interpretabilidade, isto é, que garantam uma maior flexibilidade (em relação aos modelos estatísticos tradicionais, como o de regressão logística, por exemplo) preservando a sua facilidade de interpretação, o que não se encontra em alguns desses algoritmos de aprendizado de máquina.

2.3 Estimação por Máxima Verossimilhança

A estimação dos parâmetros do modelo de regressão binária é feita, usualmente, com a aplicação de distribuições de probabilidade na modelagem de variáveis aleatórias, com o objetivo de estimar quantidades populacionais desconhecidas. De acordo com Cordeiro & Demétrio [27], muitos métodos podem ser empregados para estimar os parâmetros β_1, \dots, β_k , incluindo o método dos mínimos quadrados, o Bayesiano e o da máxima verossimilhança (MV), o qual possui várias propriedades ótimas, tais como consistência e eficiência assintótica.

Neste trabalho, considera-se o método de MV para estimar os parâmetros lineares β_1, \dots, β_k do modelo de regressão binária. Com a suposição de independência dos valores de Y_i , para $i = 1, 2, \dots, n$, a função de verossimilhança para esses parâmetros é dada por:

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n \mu_i^{y_i} (1 - \mu_i)^{1-y_i}.$$

Como os valores que maximizam a função de verossimilhança acima são os mesmos que maximizam seu logaritmo, então pode-se escrevê-la da seguinte forma:

$$\begin{aligned} \ell(\boldsymbol{\beta}) &= \log(L(\boldsymbol{\beta})) = \sum_{i=1}^n \log(\mu_i^{y_i} (1 - \mu_i)^{1-y_i}) \\ &= \sum_{i=1}^n y_i \log(\mu_i) + \sum_{i=1}^n (1 - y_i) \log(1 - \mu_i). \end{aligned}$$

Quando as condições de regularidade estão satisfeitas, segundo Casella & Berger [23], o máximo global da função $\ell(\boldsymbol{\beta})$ é encontrado, unicamente, pelas soluções da expressão:

$$\frac{\partial \log(L(\boldsymbol{\beta}))}{\partial \boldsymbol{\beta}} = \frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \mathbf{0}.$$

De modo que o estimador de MV $\hat{\boldsymbol{\beta}}$ de $\boldsymbol{\beta}$ seja obtido pela solução do sistema de k equações que torna o vetor escore igual a zero, ou seja,

$$U(\hat{\boldsymbol{\beta}}) = \mathbf{0}, \quad (2.1)$$

em que:

$$U_j(\boldsymbol{\beta}) = \frac{\partial \ell(\boldsymbol{\beta})}{\partial \beta_j} = 0, \quad \text{para } j = 1, \dots, k.$$

Em geral, quando não existem soluções exatas (ou analíticas) para a equação (2.1), elas têm que ser obtidas numericamente por meio de processos iterativos, como por exemplo, Newton-Raphson e BFGS (Broyden-Fletcher-Goldfarb-Shanno). Neste trabalho, para a estimação dos parâmetros pelo método de MV foi utilizada a função $\text{maxLik}(\cdot, \text{method} = \text{"BFGS"})$ do pacote de mesmo nome (ver Henningsen & Toomet [42]) do programa estatístico R, versão 3.6.3 (R Core Team [76]).

2.4 Distribuições de Probabilidade Simétricas e Assimétricas

Nos últimos anos, vem crescendo o interesse de pesquisadores em estudar diferentes famílias de distribuições paramétricas simétricas, bem como suas versões assimétricas. De

acordo com Nitha & Krishnarani [51], uma razão para isto é que muitas das famílias de distribuições simétricas existentes não são eficazes em modelar os conjuntos de dados assimétricos e de caudas pesadas que surgem em várias situações da vida real. Com isso, aumentou o interesse em buscar novas famílias de distribuições simétricas e assimétricas.

Nas subseções seguintes são vistas algumas distribuições de probabilidade utilizadas como base para a formulação dos diferentes modelos de regressão binária.

2.4.1 Distribuição Logística

A distribuição logística foi proposta inicialmente para estudos de crescimento populacional humano (Balakrishnan [12]). Em teoria da probabilidade e estatística, a distribuição logística é classificada como sendo uma distribuição de probabilidade contínua. Então, sendo X uma variável aleatória contínua, diz-se que X tem distribuição logística com parâmetros de locação $\mu \in \mathbb{R}$ e de escala $s > 0$, se sua função densidade de probabilidade (FDP) é dada por:

$$f(x; \mu, s) = \frac{e^{-\frac{x-\mu}{s}}}{s \left(1 + e^{-\frac{x-\mu}{s}}\right)^2}, \quad -\infty < x < \infty.$$

E a FDA é dada por:

$$F(x; \mu, s) = \frac{1}{1 + e^{-\frac{x-\mu}{s}}}, \quad -\infty < x < \infty.$$

Se $\mu = 0$ e $s = 1$, a distribuição assume a sua forma padrão.

2.4.2 Distribuição de Lindley

Lindley [53, 54] introduziu uma nova família de distribuições contínuas para uma variável aleatória X com suporte nos reais não-negativos. Diz-se que uma variável aleatória X segue uma distribuição de Lindley com parâmetro $\theta > 0$ se sua FDP é dada por:

$$f(x; \theta) = \frac{\theta^2}{(\theta + 1)}(1 + x)e^{-\theta x}, \quad x \geq 0. \quad (2.2)$$

Diversas propriedades e aplicações desta distribuição foram estudadas por Ghitany *et al.* [37] e Al-Mutairi *et al.* [4]. Em (2.2), nota-se que a distribuição de Lindley é uma mistura das distribuições Exponencial(θ) e Gama(2, θ), sendo que os pesos dados a essa mistura são, respectivamente, γ e $(1 - \gamma)$, com $\gamma = \theta/(1 + \theta)$. Esses mesmos autores observaram que, embora a distribuição de Lindley seja semelhante à distribuição exponencial, ela pode ser usada como um modelo melhor do que a distribuição exponencial em algumas situações, devido ao fato de que, enquanto a distribuição exponencial tem taxa

de risco e função de vida residual média (MRLF, do inglês *mean residual life function*) constantes, a distribuição de Lindley possui taxa de risco crescente e MRLF decrescente.

Estudos sobre a distribuição de Lindley vêm crescendo e ganhando bastante espaço; várias extensões/generalizações dela podem ser encontradas na literatura estatística recente. Conforme exemplificado por Kumar & Jose [66], Nadarajah *et al.* [58] introduziram uma forma generalizada da distribuição de Lindley e mostraram sua superioridade sobre as distribuições gama, log-normal e a forma exponenciada de diferentes distribuições; Bakouch *et al.* [11] obtiveram uma forma estendida da distribuição de Lindley; Elbatal & Elgarhy [32] estudaram a distribuição Kumaraswamy quase Lindley; Gómez-Déniz *et al.* [38] consideraram a distribuição log-Lindley; Ashour & Eltehiwy [8] consideraram a distribuição potência Lindley exponenciada; Nedjar & Zeghdoudi [59]) introduziram a distribuição gama Lindley. Ademais, outra forma estendida da distribuição Lindley generalizada com aplicações a dados sobre tempos de vida foi considerada por Torabi *et al.* [77], enquanto certas propriedades e aplicações da distribuição Lindley ponderada generalizada foram discutidas por Ramos & Louzada [77].

A distribuição de Lindley oferece muitas vantagens ao ter seu suporte estendido para toda a reta real, visto que resulta num modelo mais flexível e competitivo do que muitas classes de distribuições simétricas com suporte em $(-\infty, \infty)$ (Nitha & Krishnarani [51]).

2.4.2.1 Distribuição Double Lindley

Tal distribuição foi proposta por Nitha & Krishnarani [51] e Kumar & Jose [66]. Seja X uma variável aleatória contínua com distribuição double Lindley (DL) de parâmetro $\theta > 0$, então sua FDP é dada por:

$$f(x; \theta) = \frac{\theta^2}{2(\theta + 1)}(1 + |x|)e^{-\theta|x|}, \quad -\infty < x < \infty. \quad (2.3)$$

Conforme observado por Nitha & Krishnarani [51], a FDP da variável aleatória com distribuição DL pode ser vista como uma mistura de duas FDPs: uma Laplace com média 0 e variância $2\theta^2$ e outra gama bilateral (*two-sided*) com parâmetro de forma 2 e parâmetro de escala θ , cujos pesos são, respectivamente, γ e $(1 - \gamma)$, com $\gamma = \theta/(1 + \theta)$.

Observa-se na Figura 2.1 a forma da FDP da distribuição DL(θ) para diferentes valores atribuídos a θ , onde fica evidente que a densidade é simétrica em relação a 0 (média e mediana da distribuição), tem natureza unimodal (para $\theta \geq 1$; neste caso, a moda é igual a 0) e bimodal (para $\theta < 1$; neste caso, as modas são iguais a $\pm(1 - 1/\theta)$), e fica mais elevada (ou “pontuda”) à medida que o o valor de θ aumenta.

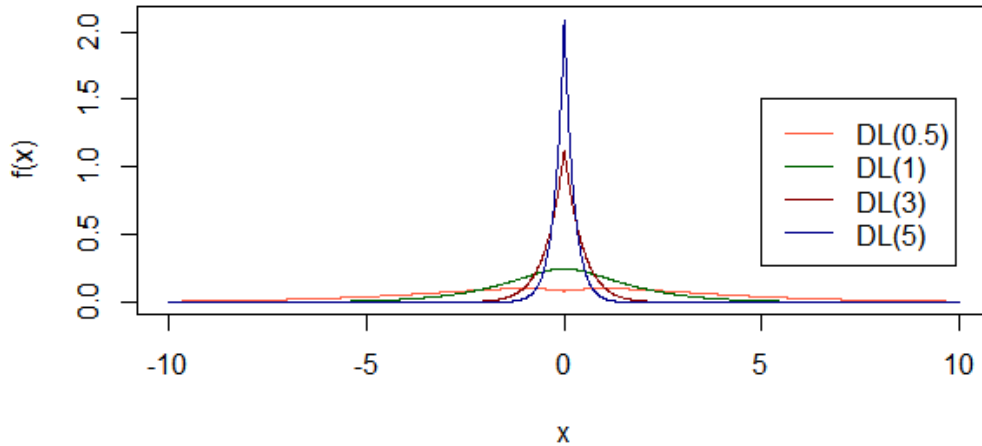


Figura 2.1: Função densidade de probabilidade da distribuição DL, para diferentes valores de θ .

Por sua vez, a FDA da distribuição DL(θ) é expressa por:

$$F(x; \theta) = \begin{cases} \frac{1}{2(\theta + 1)} [1 + \theta(1 - x)] e^{\theta x}, & \text{se } x \leq 0, \\ 1 - \frac{1}{2(\theta + 1)} [1 + \theta(1 + x)] e^{-\theta x}, & \text{se } x > 0, \end{cases} \quad (2.4)$$

que também pode ser escrita da seguinte forma:

$$F(x; \theta) = \frac{1}{2} + \frac{1}{2} \text{sgn}(x) \left\{ 1 - \frac{1}{\theta + 1} [1 + \theta(1 + |x|)] e^{-\theta|x|} \right\}, \quad -\infty < x < \infty, \quad (2.5)$$

em que $\text{sgn}(\cdot)$ denota a função sinal.

A Figura 2.2 mostra a FDA da distribuição DL(θ) para diferentes valores de θ , em que observa-se a sua forma não decrescente e contínua.

Para maiores detalhes, assim como outras propriedades (momentos, curtose, assimetria, etc.) da distribuição DL, recomenda-se ao leitor consultar os trabalhos de Nitha & Krishnarani [51] e Kumar & Jose [66].

Importante ressaltar que, para obter o modelo de regressão binária a partir da distribuição DL, deve-se considerá-la em sua forma padrão. Como ela possui um único parâmetro θ , basta então tomar $\theta = 1$. Assim, a forma padrão dessa distribuição é representada por:

$$F(x) = \frac{1}{2} + \frac{1}{2} \text{sgn}(x) \left\{ 1 - \left(1 + \frac{|x|}{2} \right) e^{-|x|} \right\}, \quad -\infty < x < \infty. \quad (2.6)$$

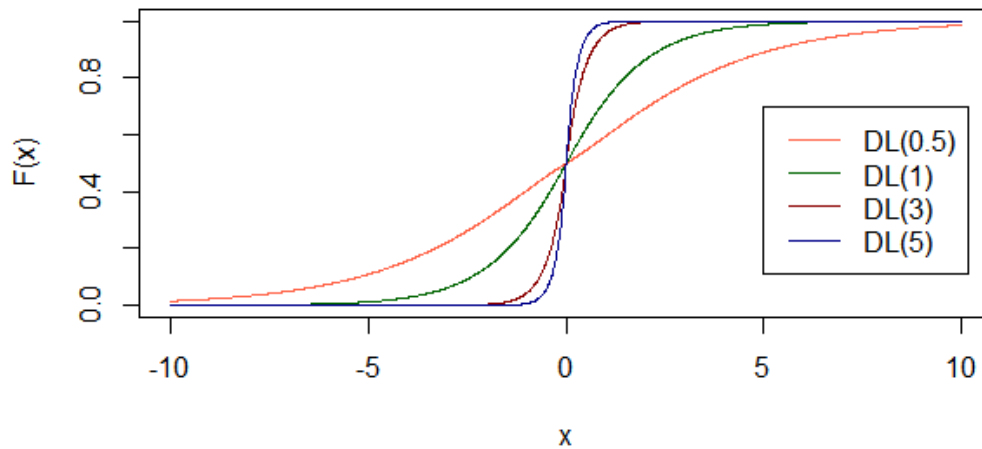


Figura 2.2: Função de distribuição acumulada da distribuição DL, para diferentes valores de θ .

Segundo os trabalhos de Bazán, Torres-Avilés, Suzuki & Louzada [64] e Silva, Anyosa & Bazán [73], para a construção de distribuições de probabilidade do tipo potência ou reversa de potência, parte-se de uma distribuição dita de base (ou basal), cuja FDP é unimodal, log-côncava e de suporte real. Essas propriedades são válidas para a forma padrão da distribuição DL. Portanto, ela pode ser empregada para o desenvolvimento de novas distribuições do tipo potência e reversa de potência. Isto será discutido em maiores detalhes nas Seções 2.4.2.3 e 2.19.

2.4.2.2 Distribuição Double Lindley Assimétrica

Como visto anteriormente, a distribuição DL pertence à família de distribuições simétricas. Logo, sua aplicabilidade, no contexto de regressão binária, pode estar restrita a situações em que os conjuntos de dados reais são balanceados ou com desbalanceamento pouco acentuado. Existem diferentes métodos de introdução de assimetria em uma família de distribuições simétricas. Nitha & Krishnarani [51] sugerem ao leitor consultar, por exemplo, os trabalhos de Ayebo & Kozubowski [9], Kotz *et al.* [49] e Azzalini [10], para algumas dessas propostas. E, para as aplicações das distribuições assim formadas, consultar, por exemplo, Julia & Vives-Rego [48] e Kozubowski & Podgorski [50].

Nesta subsubseção é apresentada uma versão assimétrica da distribuição DL, obtida por Nitha & Krishnarani [51] usando a ideia de fatores de escala inversa de Fernandez & Steel [34]. Neste método, um parâmetro novo é adicionado, que atua como um parâmetro

de assimetria na família de distribuições simétricas.

A FDP da distribuição double Lindley assimétrica (ADL) com parâmetros $\theta > 0$ e $\lambda > 0$, de acordo com Nitha & Krishnarani [51], é dada por:

$$f(x; \theta, \lambda) = \frac{\theta^2}{(\theta + 1)} \frac{\lambda}{(1 + \lambda^2)} \begin{cases} \left(1 - \frac{x}{\lambda}\right) e^{\frac{\theta x}{\lambda}}, & \text{se } x \leq 0, \\ (1 + \lambda x) e^{-\theta \lambda x}, & \text{se } x > 0. \end{cases} \quad (2.7)$$

Nota-se que, para $\lambda = 1$, a distribuição ADL(θ, λ) se reduz à distribuição DL(θ).

Os momentos não centrais de ordem r da distribuição ADL(θ, λ) podem ser calculados como segue (Nitha & Krishnarani [51]):

$$\begin{aligned} \mu'_r = \mathbb{E}[X^r] &= \frac{\theta^2}{(\theta + 1)} \left(\int_{-\infty}^0 x^r \left(1 - \frac{x}{\lambda}\right) e^{\frac{\theta x}{\lambda}} dx + \int_0^{\infty} x^r (1 + \lambda x) e^{-\theta \lambda x} dx \right) \\ &= \frac{\theta^{1-r}}{(\theta + 1)} \frac{\lambda}{(1 + \lambda^2)} \Gamma(r + 1) \left[1 + \left(\frac{r + 1}{\theta}\right) \right] \left[(-1)^r \lambda^{r+1} + \frac{1}{\lambda^{r+1}} \right]. \end{aligned}$$

Em particular, quando $r = 1$, tem-se a média dessa distribuição:

$$\mathbb{E}[X] = \frac{\theta + 2}{\theta} \left(\frac{1 - \lambda^2}{\lambda} \right). \quad (2.8)$$

Para $r = 2$, tem-se:

$$\mathbb{E}[X^2] = \frac{2(\theta + 3)}{\theta^2(\theta + 1)} \left(\frac{\lambda}{1 + \lambda^2} \right) \left(\lambda^3 + \frac{1}{\lambda^3} \right). \quad (2.9)$$

Logo, a partir dos resultados (2.8) e (2.9), pode-se calcular a variância da distribuição ADL através da seguinte expressão: $Var[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$.

Nas Figuras 2.3 e 2.4, observa-se as diferentes formas da FDP (2.7) quando são considerados diferentes valores de θ e λ . Na Figura 2.4, por exemplo, nota-se uma assimetria à direita (ou positiva) para os valores de $\lambda = \{0,5; 0,6\}$, enquanto que, para $\lambda = 3$, percebe-se uma assimetria à esquerda (ou negativa).

Finalmente, a FDA da distribuição ADL(θ, λ) é expressa por:

$$F(x; \theta, \lambda) = \begin{cases} \frac{\lambda^2 e^{\frac{\theta x}{\lambda}}}{(\theta + 1)(1 + \lambda^2)} \left[1 + \theta \left(1 - \frac{x}{\lambda}\right) \right], & \text{se } x \leq 0, \\ 1 - \frac{e^{-\theta \lambda x}}{(\theta + 1)(1 + \lambda^2)} [1 + \theta(1 + \lambda x)], & \text{se } x > 0. \end{cases} \quad (2.10)$$

Para obter o modelo de regressão binária a partir da distribuição ADL, de maneira similar à DL, adota-se $\theta = 1$. Já o parâmetro λ (assimetria) será estimado.

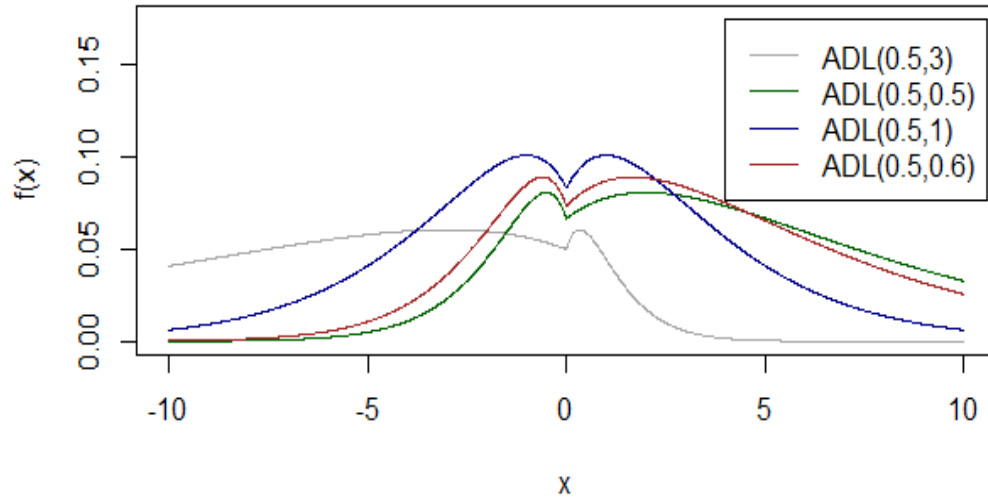


Figura 2.3: Função densidade de probabilidade da distribuição ADL, considerando $\theta = 0,5$ e $\lambda = \{0,5, 0,6, 1, 3\}$.

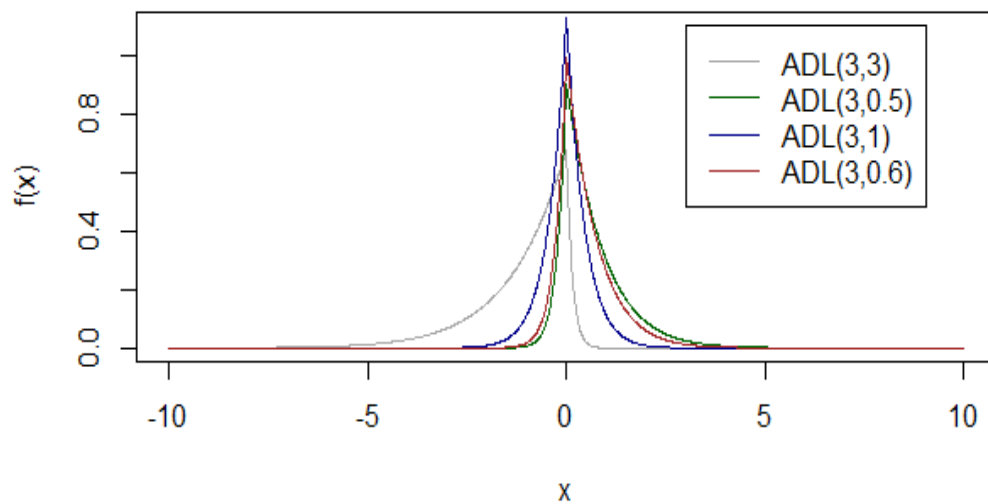


Figura 2.4: Função densidade de probabilidade da distribuição ADL, considerando $\theta = 3$ e $\lambda = \{0,5, 0,6, 1, 3\}$.

Então, a FDA da distribuição padrão ADL(1, λ) é expressa por:

$$F(x; \lambda) = \begin{cases} \frac{\lambda^2 e^{\frac{x}{\lambda}}}{2(1 + \lambda^2)} \left(2 - \frac{x}{\lambda}\right), & \text{se } x \leq 0, \\ 1 - \frac{e^{-\lambda x}}{2(1 + \lambda^2)} (2 + \lambda x), & \text{se } x > 0. \end{cases} \quad (2.11)$$

Para maiores detalhes e outras propriedades da distribuição ADL, ver Nitha & Krishnarani [51].

2.4.2.3 Distribuição Potência Double Lindley

Lemonte & Bazán [52] descrevem a composição de uma distribuição de probabilidade, que está fundamentada em considerar uma FDA contínua arbitrária e elevá-la a uma potência real positiva. Com isso, é apresentada uma nova FDA com um parâmetro de potência adicional, que é conhecida como distribuição de potência.

Diz-se que uma variável aleatória univariada X segue uma distribuição de potência com parâmetro de locação $\mu \in \mathbb{R}$, parâmetro de escala $\sigma > 0$ e parâmetro de forma $\lambda > 0$, se X tem FDP representada por:

$$f(x; \mu, \sigma, \lambda) = \frac{\lambda}{\sigma} g\left(\frac{x - \mu}{\sigma}\right) \left[G\left(\frac{x - \mu}{\sigma}\right)\right]^{\lambda-1}, \quad -\infty < x < \infty, \quad (2.12)$$

em que $g(\cdot)$ e $G(\cdot)$ são, respectivamente, a FDP e FDA padrão de qualquer distribuição univariada contínua com suporte na reta real (dita distribuição de base).

A FDA da distribuição de potência é definida como:

$$F(x; \mu, \sigma, \lambda) = G\left(\frac{x - \mu}{\sigma}\right)^\lambda, \quad -\infty < x < \infty. \quad (2.13)$$

As distribuições de potência são assimétricas à direita se $\lambda > 1$, e assimétricas à esquerda se $0 < \lambda < 1$.

Alguns modelos de regressão binária construídos com base em distribuições de potência, como por exemplo, potência logística, potência normal e potência Cauchy, dentre outras, foram explorados e investigados mais a fundo nos trabalhos de Anyosa [7], Huayanay [44], Bazán, Torres-Avilés, Suzuki & Louzada [64].

Quando analisa-se dados binários que apresentam certo grau de assimetria ou desbalanceamento entre as classes, as funções de ligação simétricas podem não ser úteis para ajustar esses dados, produzindo resultados insatisfatórios. Neste contexto, é desenvolvida aqui a versão de potência para a distribuição DL, que dá origem à distribuição potência double Lindley (PDL).

Então, é apresentada uma nova FDP com um parâmetro de potência adicional $\lambda > 0$, que é definida da forma:

$$f(x; \theta, \lambda) = \lambda \frac{\theta^2(1 + |x|)e^{-\theta|x|}}{2(\theta + 1)} \left(\frac{1}{2} + \frac{1}{2} \operatorname{sgn}(x) \left\{ 1 - \left[\frac{1 + \theta(1 + |x|)e^{-\theta|x|}}{\theta + 1} \right] \right\} \right)^{\lambda-1}, \quad (2.14)$$

para $x \in \mathbb{R}$.

Nas Figuras 2.5 e 2.6, observa-se as diferentes formas da FDP (2.14) quando são considerados diferentes valores de θ e λ .

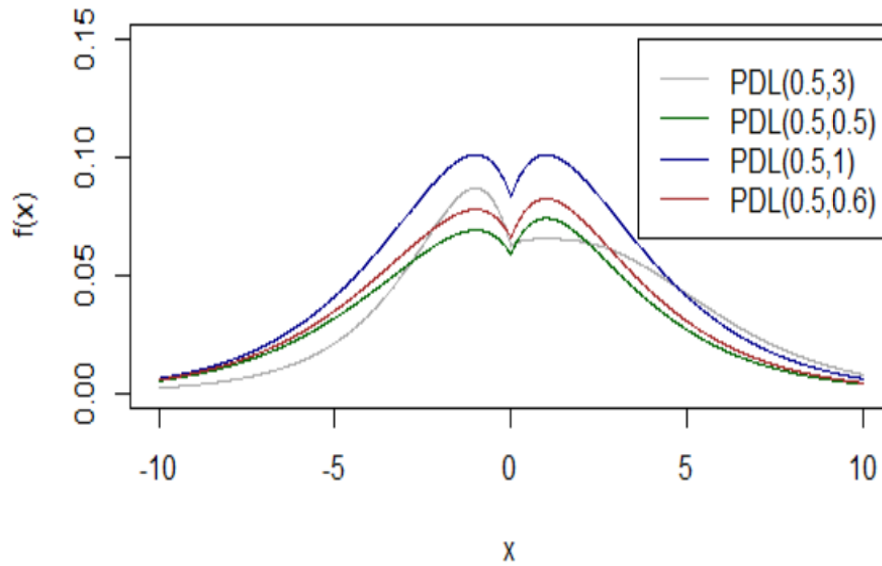


Figura 2.5: Função densidade de probabilidade da distribuição PDL, considerando $\theta = 0,5$ e $\lambda = \{0,5, 0,6, 1, 3\}$.

A FDA, por sua vez, é definida da forma:

$$F(x; \theta, \lambda) = \left(\frac{1}{2} + \frac{1}{2} \operatorname{sgn}(x) \left\{ 1 - \frac{1}{\theta + 1} [1 + \theta(1 + |x|)] e^{-\theta|x|} \right\} \right)^{\lambda}, \quad x \in \mathbb{R}. \quad (2.15)$$

Analogamente à distribuição ADL, para obter o modelo de regressão binária baseado na distribuição PDL, fixa-se $\theta = 1$, enquanto o parâmetro λ (assimetria) será estimado. Assim, a FDA padrão fica definida como segue:

$$F(x; \lambda) = \left(\frac{1}{2} + \frac{1}{2} \operatorname{sgn}(x) \left\{ 1 - \left(1 + \frac{|x|}{2} \right) e^{-|x|} \right\} \right)^{\lambda}, \quad x \in \mathbb{R}. \quad (2.16)$$

2.4.2.4 Distribuição Reversa de Potência Double Lindley

Para as distribuições reversa de potência, uma propriedade que precisa ser considerada quando são usadas diferentes funções de ligação é a de reversibilidade, segundo a

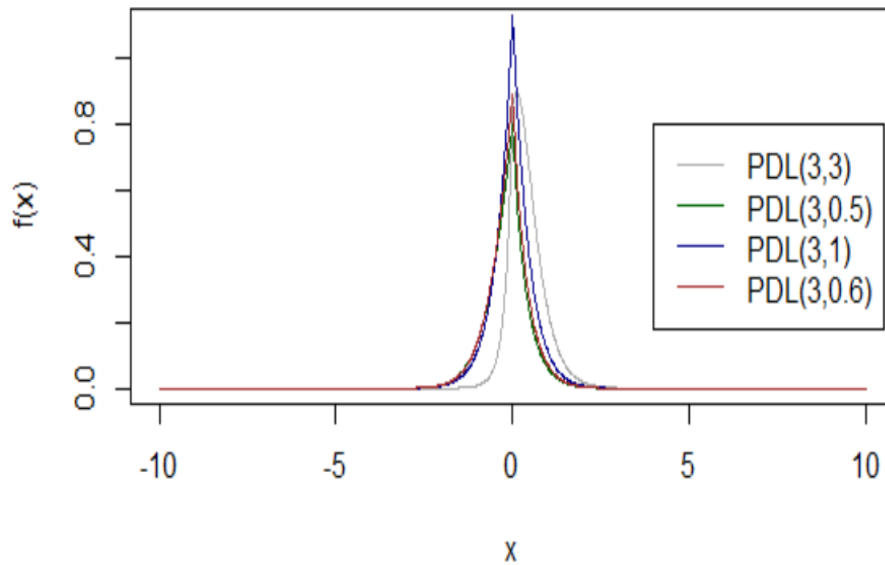


Figura 2.6: Função densidade de probabilidade da distribuição PDL, considerando $\theta = 3$ e $\lambda = \{0,5, 0,6, 1, 3\}$.

qual a distribuição de S satisfaz a propriedade de reversibilidade se $S \sim F(\cdot) \implies -S \sim G(\cdot) \equiv 1 - F(-\cdot)$. Neste caso, $G(\cdot)$ é chamada de distribuição reversa de $F(\cdot)$ (Bazán *et al.* [14]).

Como também descrito por Lemonte & Bazán [52], a FDP da distribuição reversa de potência pode ser representada pela expressão a seguir:

$$f(x; \mu, \sigma, \lambda) = \frac{\lambda}{\sigma} g\left(\frac{x - \mu}{\sigma}\right) \left[G\left(-\left(\frac{x - \mu}{\sigma}\right)\right)\right]^{\lambda-1}, \quad -\infty < x < \infty, \quad (2.17)$$

e a sua FDA pode ser expressa por:

$$F(x; \mu, \sigma, \lambda) = 1 - G\left(-\left(\frac{x - \mu}{\sigma}\right)\right)^{\lambda}, \quad -\infty < x < \infty, \quad (2.18)$$

sendo $\mu \in \mathbb{R}$, $\sigma > 0$ e $\lambda > 0$ os parâmetros de localização, escala e forma, respectivamente, e $G(\cdot)$ é a FDA padrão da distribuição de base.

Considerando a propriedade de reversibilidade, é possível propor outras distribuições reversa de potência. Neste contexto, como o objetivo do trabalho é propor distribuições para lidar com dados assimétricos, é desenvolvida aqui a versão reversa de potência para a distribuição DL, que dá origem à distribuição reversa de potência double Lindley (RPDL).

Logo, é apresentada uma nova FDP com um parâmetro de potência adicional

$\lambda > 0$, que é definida da forma:

$$f(x; \theta, \lambda) = \lambda \frac{\theta^2(1 + |x|)e^{-\theta|x|}}{2(\theta + 1)} \times \left(\frac{1}{2} + \frac{1}{2} \operatorname{sgn}(-x) \left\{ 1 - \left[\frac{1 + \theta(1 + |-x|)e^{-\theta|-x|}}{\theta + 1} \right] \right\} \right)^{\lambda - 1}, \quad (2.19)$$

para $x \in \mathbb{R}$.

Nas Figuras 2.7 e 2.8, observa-se as diferentes formas da FDP (2.19) quando são considerados diferentes valores de θ e λ .

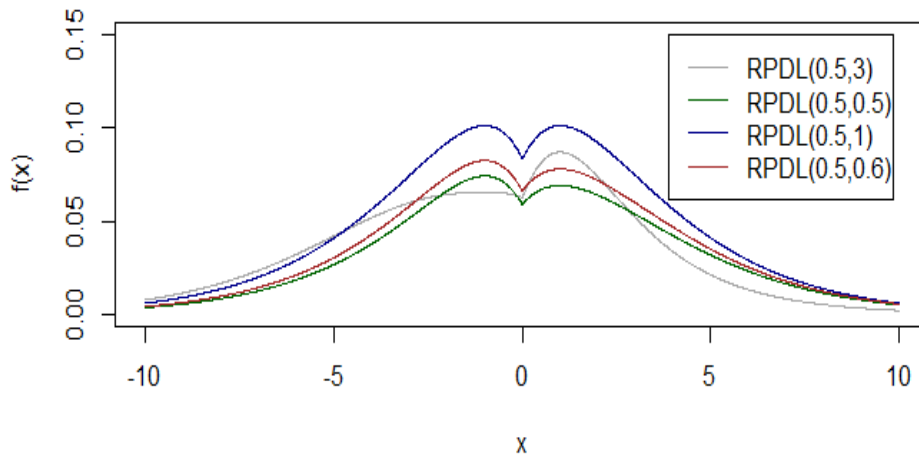


Figura 2.7: Função densidade de probabilidade da distribuição RPD L, considerando $\theta = 0,5$ e $\lambda = \{0,5, 0,6, 1, 3\}$.

A FDA, por sua vez, é definida como:

$$F(x; \theta, \lambda) = 1 - \left(\frac{1}{2} + \frac{1}{2} \operatorname{sgn}(-x) \left\{ 1 - \frac{1}{\theta + 1} [1 + \theta(1 + |-x|)] e^{-\theta|-x|} \right\} \right)^{\lambda}, \quad (2.20)$$

para $x \in \mathbb{R}$. Em sua forma padrão, isto é, para $\theta = 1$, tem-se:

$$F(x; \lambda) = 1 - \left(\frac{1}{2} + \frac{1}{2} \operatorname{sgn}(-x) \left\{ 1 - \left(1 + \frac{|-x|}{2} \right) e^{-|-x|} \right\} \right)^{\lambda}, \quad x \in \mathbb{R}. \quad (2.21)$$

Vale lembrar que a propriedade de reversibilidade indica que $F^{**}(x) + F^*(-x) = 1$, sendo $F^{**}(\cdot)$ a distribuição de interesse e $F^*(\cdot)$ a sua distribuição reversa. Considerando $F_P(\cdot)$ a FDA padrão da distribuição PDL, dada por (2.16), e $F_{RP}(\cdot)$ a FDA padrão da distribuição RPD L, dada por (2.21), então é fácil verificar que (Anyosa [7]):

- $F_P(\cdot)$ e $F_{RP}(\cdot)$ não são ponto-simétricas, uma vez que $F_P(-x) \neq 1 - F_P(x)$ ou $F_{RP}(-x) \neq 1 - F_{RP}(x)$ para $\lambda \neq 1$;

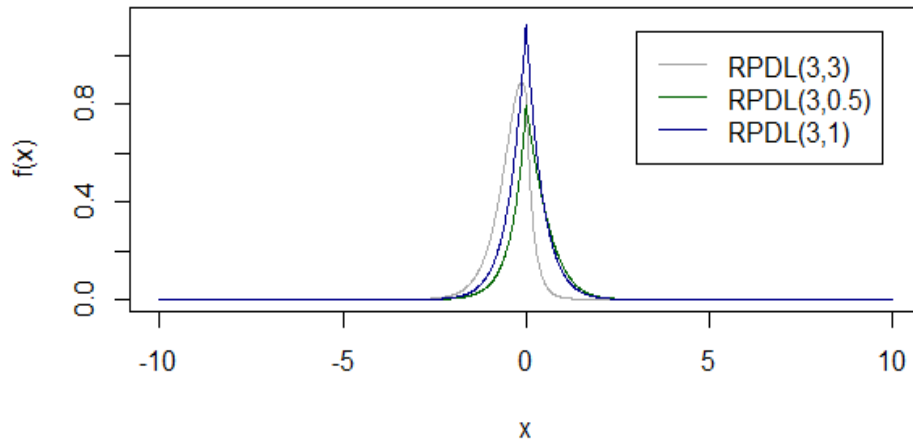


Figura 2.8: Função densidade de probabilidade da distribuição RPD L, considerando $\theta = 3$ e $\lambda = \{0,5, 1, 3\}$.

- $F_P(\pm x) + F_{RP}(\mp x) = 1$, o que significa que as duas distribuições são diferentes mas estão relacionadas porque uma é reversa da outra;
- Se $\lambda = 1$, então $F_P(x) = G(x) = F_{RP}(x)$, o que quer dizer que $G(\cdot)$ é um caso particular das duas distribuições ou distribuição de base;
- Se $X \sim F_P(\cdot)$, então $-X \sim F_{RP}(\cdot)$;
- Se $X \sim F_P(\cdot)$ ou $X \sim F_{RP}(\cdot)$, então X não é simétrica se $\lambda \neq 1$.

Particularmente, para os modelos de regressão binária com base nas distribuições ADL, PDL e RPD L (vide Seções 2.4.2.2, 2.4.2.3 e 2.19, respectivamente), utiliza-se a transformação $\lambda = \exp\{\lambda^*\}$ para que o parâmetro λ^* esteja definido nos reais; com isso, calcula-se o erro padrão de $\hat{\lambda}$ (estimador de MV de λ) a partir da aplicação do método delta, isto é, mediante o uso da função `deltamethod(·)` do pacote `msm` (Jackson [47]).

2.5 Comparação de Modelos

A comparação de modelos é fundamental quando são observados dois ou mais modelos ajustados descrevendo eventos semelhantes, com o intuito de selecionar aquele que represente melhor ou se aproxime mais da realidade. Tal comparação pode ser feita conforme descrito em Gelman *et al.* [36]. Após ajustar todos os modelos candidatos

propostos, a escolha é feita com base naquele que apresentar como resultado o melhor desempenho, ou seja, o menor valor dentre todos os critérios de informação.

Na literatura existem vários critérios de informação disponíveis, como por exemplo, o critério de informação de Akaike (AIC) (Akaike [3]), o critério de informação Bayesiano ou de Schwarz (BIC) (Schwarz [67]), o AIC corrigido (AICc) (Sugiura [74]; Hurvich & Tsai [45]), o AIC consistente (CAIC) (Bozdogan, [17]; Anderson *et al.* [6]), o critério de informação de Hannan-Quinn (HQIC) (Hannan & Quinn [40]), etc. Dentre eles, os mais conhecidos e utilizados são o AIC e o BIC. Esses dois critérios de comparação ou seleção de modelos são descritos nas subseções seguintes.

2.5.1 Critério de Informação de Akaike

Desenvolvido pelo estatístico Hirotugu Akaike com o nome de “um critério de informação” em 1971, e apresentado na literatura por Akaike [3], o AIC é uma medida relativa da qualidade do ajuste de um modelo estimado. Ele é um critério que avalia a qualidade do ajuste do modelo paramétrico estimado pelo método da máxima verossimilhança, e consiste no fato de que o viés tende ao número de parâmetros a serem estimados no modelo. O AIC é definido como:

$$\text{AIC} = -2 \log \left(L(\hat{\theta}) \right) + 2p,$$

sendo:

- $L(\hat{\theta})$ o valor máximo da função de verossimilhança, isto é, a função de verossimilhança avaliada no estimador (ou estimativa) de máxima verossimilhança $\hat{\theta}$ de θ ;
- p o número de parâmetros a serem estimados no modelo.

De acordo com o AIC, pode-se classificar vários modelos concorrentes ajustados a um mesmo conjunto de dados, com aqueles tendo os menores valores de AIC sendo considerados os melhores (Burnham & Anderson [21]). A partir do valor do AIC, pode-se inferir, por exemplo, que os principais modelos estão “empatados” e os restantes não seriam recomendados. Por isso, o valor do AIC torna-se útil quando são comparados diversos modelos. O modelo com o menor AIC é classificado como o “melhor” modelo, dentre todos os modelos comparados.

2.5.2 Critério de Informação Bayesiano

Também conhecido como Critério de Informação de Schwarz, o BIC é usado quando o interesse reside em selecionar apenas um dentro de um conjunto finito de modelos

estatísticos. A sua vantagem em relação a outros critérios, como por exemplo, o AIC, está na inclusão do tamanho amostral na penalização, diminuindo, assim, a chance de selecionar um modelo que incorra em *overfitting* (sobreajuste). Proposto por Schwarz [67], o BIC baseia-se, em partes, no máximo da função de verossimilhança do modelo, e é definido como:

$$\text{BIC} = -2 \log \left(L(\hat{\theta}) \right) + p \log(n),$$

em que:

- $L(\hat{\theta})$ é o máximo da função de verossimilhança;
- p é o número de parâmetros a serem estimados no modelo;
- n é o número de observações da amostra.

Então, o modelo com o menor BIC é classificado como o “melhor” modelo, dentre todos os modelos comparados.

2.6 Métricas de Desempenho

Nesta seção são apresentadas algumas métricas de desempenho para avaliar a capacidade preditiva dos modelos de regressão binária propostos. Existem inúmeras métricas diferentes e algumas funcionam melhor do que outras para um determinado tipo de problema; escolher uma medida adequada para avaliar o modelo é tão importante quanto escolher um bom modelo.

A maioria das métricas de desempenho é baseada na matriz de confusão, apresentada em uma tabela de contingência 2×2 , que é uma das formas mais simples de visualizar e estabelecer o cálculo dessas métricas.

A partir da Tabela 2.1 (matriz de confusão), tem-se que:

- Y equivale às respostas binárias observadas na amostra, isto é, $Y = 1$ representa “sucesso” e $Y = 0$ denota “fracasso”;
- \hat{Y} é a variável predita correspondente às respostas classificadas pelo modelo em análise, isto é, $\hat{Y} = 1$ indica “sucesso” e $\hat{Y} = 0$ representa “fracasso”;
- VP (Verdadeiro positivo) é quando a observação é classificada como “sucesso” e é “sucesso”;
- VN (Verdadeiro negativo) é quando a observação é classificada como “fracasso” e é “fracasso”;

- FP (Falso positivo) é quando a observação é classificada como “sucesso” e é “fracasso”;
- FN (Falso negativo) é quando a observação é classificada como “fracasso” e é “sucesso”.

Logo, tem-se que $VP + VN + FP + FN = n$, sendo n o tamanho da amostra.

Tabela 2.1: Matriz de confusão 2×2 (classificação binária).

	Predito (\hat{Y})		
	1	0	
Observado (Y)	1	VP	FN
	0	FP	VN

Desta forma, pode-se utilizar as seguintes métricas para avaliação da capacidade preditiva (no caso, valores dessas medidas próximos a 1 são indicativos de um “bom” modelo):

Sensibilidade (SEN): também conhecida como “revocação” (ou *recall*), corresponde à proporção dos verdadeiros positivos entre todas as observações que realmente são positivas no conjunto de dados, isto é,

$$SEN = \frac{VP}{VP + FN}.$$

Especificidade (SPE): é a proporção dos verdadeiros negativos entre todas as observações que realmente são negativas no conjunto de dados, ou seja,

$$SPE = \frac{VN}{VN + FP}.$$

Valor Preditivo Positivo (VPP): também conhecida como “precisão” (ou *precision*), tal métrica traz a informação da quantidade de observações classificadas como positivas que são realmente positivas, ou seja, a proporção de verdadeiros positivos em relação a todas as predições positivas:

$$VPP = \frac{VP}{VP + FP}.$$

Valor Preditivo Negativo (VPN): é a métrica que traz a informação da quantidade de observações classificadas como negativas que são realmente negativas, ou seja, a proporção de verdadeiros negativos em relação a todas as predições negativas:

$$VPN = \frac{VN}{VN + FN}.$$

Acurácia (ACC): é a métrica mais popular; ela representa a proporção de acertos do modelo, ou seja, é a fração de verdadeiros positivos e verdadeiros negativos em relação a todos os resultados possíveis:

$$ACC = \frac{VP + VN}{VP + VN + FP + FN}.$$

No entanto, quando as classes são desbalanceadas, utilizar a ACC não é adequado, pois tal medida poderia causar uma falsa impressão de bom desempenho, levando a tirar conclusões não satisfatórias.

Acurácia Balanceada ($ACCB$): a métrica anterior é geralmente usada na forma balanceada quando o problema a ser estudado envolve classificação desbalanceada; é dada pela média aritmética entre a fração de verdadeiros positivos e a fração de verdadeiros negativos:

$$ACCB = \frac{\frac{VP}{VP + FN} + \frac{VN}{VN + FP}}{2} = \frac{SEN + SPE}{2}.$$

F1-Score: é a métrica referente à média harmônica entre SEN e VPP :

$$F1 - Score = 2 \times \frac{SEN \times VPP}{SEN + VPP}.$$

Coefficiente de Correlação de Matthews (MCC): esta medida é usada para interpretar a classificação geral do modelo (Baldi *et al.* [13]). O MCC leva em consideração valores positivos e negativos, verdadeiros e falsos, e geralmente é considerado uma medida equilibrada que pode ser usada mesmo que as classes estudadas sejam desbalanceadas. Segundo Boughorbel [16], o MCC é um coeficiente de correlação entre as classificações binárias observadas e previstas, e retorna um valor entre -1 e $+1$. Sendo que um coeficiente de $+1$ representa uma predição (ou classificação) perfeita; quando igual a 0 , uma predição aleatória; e quando igual a -1 , indica que a predição é totalmente inversa. Tal coeficiente é definido como:

$$MCC = \frac{VP \times VN - FP \times FN}{\sqrt{(VP + FP)(VP + FN)(VN + FP)(VN + FN)}}.$$

Observa-se que o MCC não está definido quando pelo menos uma das somas ($VP + FP$, $VP + FN$, $VN + FP$ ou $VN + FN$) é igual a zero, como no caso de não existir valores previstos como positivos. Para essas situações, Burset & Guigó [22] definiram uma medida de correlação aproximada para compensar o problema declarado com o MCC . Maiores detalhes sobre a medida de correlação aproximada podem ser vistos também em Louzada, Ferreira & Diniz [56].

Uma ferramenta gráfica relevante é a curva característica de operação do receptor (ROC, do inglês *receiver operating characteristic curve*), a qual é descrita a seguir.

Curva ROC: é definida como um gráfico que indica o comportamento de um classificador binário em possibilidades diferentes do valor limite para a classificação. Para o eixo x (abscissa), tem-se a medida de $1 - SPE$, e, para o eixo y (ordenada), tem-se a medida de SEN . A curva ROC deve ser interpretada de forma que, se a curva estiver mais distante da diagonal principal, melhor o desempenho do modelo associado a ela. Para definir o melhor ponto de corte, identifica-se aquele que maximiza simultaneamente a SPE e a SEN da classificação, ou seja, o ponto de corte ótimo deve estar mais próximo do eixo superior esquerdo do gráfico. Um exemplo de curva ROC é visto na Figura 2.9.

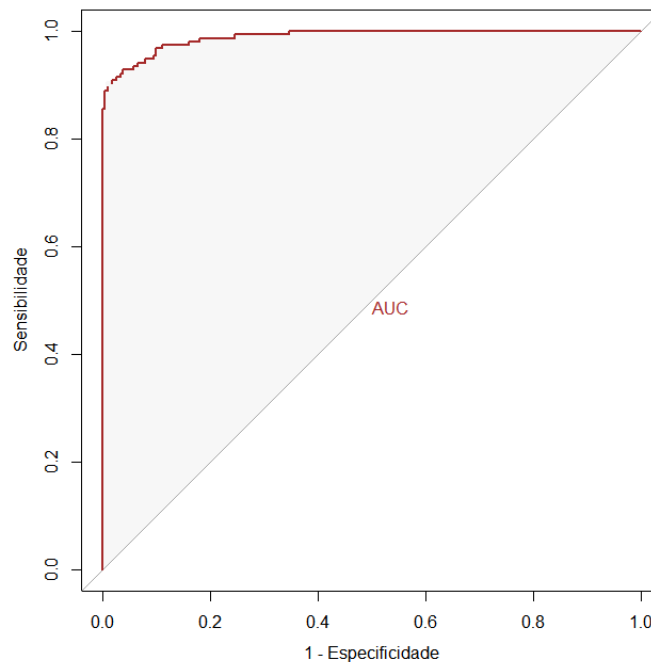


Figura 2.9: Exemplo de curva ROC.

Pode-se, ainda, comparar o desempenho dos classificadores através da área sob a curva ROC, ou AUC (do inglês *area under the curve*). Hanley & McNeil [39] indicaram que a AUC tem a propriedade importante de ser equivalente ao teste de Wilcoxon. Um modelo de classificação binária é considerado apropriado se o valor da AUC for próximo de 1; caso não seja apropriado, será próximo a 0,5 (Powers [62]).

Computacionalmente, para a implementação da curva ROC, cálculo da AUC e escolha do ponto de corte ótimo, foi utilizado o pacote `pROC` (Robin *et al.* [63]) do *software* R.

Outras duas medidas de desempenho importantes e úteis são descritas a seguir.

Área Sob a Curva Precisão \times Revocação ($AUCPR$): esta medida diz respeito à área sob a curva obtida com a comparação da precisão *versus* revocação (ou ainda, VPP

versus SEN), para diferentes pontos de corte de probabilidade. Assim como a *AUC*, tal medida é sugerida para a comparação do desempenho entre classificadores distintos, sendo mais recomendada que a *AUC* para o caso de bases de dados desbalanceadas (Saito & Rehmsmeier [65]), pois foca-se na classificação correta da classe positiva, enquanto que a *AUC* observa ambas as classes. No caso de um classificador aleatório, o valor obtido é a própria proporção de casos positivos nos dados.

Brier Score (*BS*): proposto por Glenn W. Brier em 1950 (Brier [20]), o Brier *Score* serve para avaliar (e comparar) a precisão das previsões probabilísticas. Pode ser definido como:

$$BS = \frac{1}{N} \sum_{t=1}^N (p_t - o_t)^2, \quad (2.22)$$

em que N é o número de instâncias, p_t é a probabilidade prevista da t -ésima instância pertencer à classe de interesse (evento), e o_t é o resultado real do evento na instância t (1 se ocorrer e 0 se não ocorrer). No caso, pontuações menores (mais próximas de zero) indicam melhores previsões, sendo que o mínimo que se espera de um bom modelo é que sua medida *BS* seja menor do que $1/4$. Por outro lado, são considerados modelos com previsões de má qualidade aqueles que apresentam medidas maiores do que $1/4$.

Vale ressaltar que a expressão (2.22) é apropriada somente para eventos binários. No caso de múltiplas categorias (ou classes) da variável resposta de interesse, deve-se usar a definição original dada por Brier [20]. Ainda segundo essa definição original, para problemas de classificação binária, tem-se $1/2$ como valor de referência para o *BS*, sendo desejáveis valores de $BS < 1/2$.

2.7 Método de Validação

Ainda sobre métricas de desempenho aceitáveis, o processo de construção de modelos preditivos deve passar por uma importante etapa de validação. O intuito é verificar se o modelo ajustado possui uma boa capacidade de generalização, ao ser aplicado a dados que possuem as mesmas características do conjunto a partir do qual ele foi desenvolvido. Um dos procedimentos utilizados para garantir essas validações é conhecido como *cross-validation* (ou validação cruzada), sendo uma das técnicas mais empregadas para avaliar a performance de um modelo preditivo.

Neste trabalho é considerado o método *holdout*, que consiste em dividir a base de dados original em duas partes, sendo uma delas utilizada para o desenvolvimento/construção do modelo, enquanto a outra é empregada para a avaliação da performance preditiva. Ou seja, a maior parcela dos dados é usada para a estimação do modelo (amostra treinamento), e a menor para verificação da adequabilidade/performance preditiva do modelo

(amostra teste). Na prática, divide-se em 70% e 30% para treinamento e teste, respectivamente.

Capítulo 3

Estudos de Simulação

Neste capítulo são apresentados os principais resultados de estudos de simulação dos tipos: recuperação de parâmetros (Seção 3.1), *misspecification* (Seção 3.2) e avaliação da capacidade preditiva (Seção 3.3), considerando os modelos apresentados (DL, ADL, PDL e RPDL), além do modelo de regressão binária tradicional (logística).

3.1 Recuperação de Parâmetros

Neste primeiro estudo de simulação, foram geradas $M = 1.000$ amostras (ou conjuntos de dados) de tamanhos $n = \{100, 200, 500, 1.000\}$ de cada um dos quatro modelos sugeridos. Para os modelos assimétricos (ADL, PDL e RPDL), considerou-se, ainda, três valores distintos para o parâmetro de assimetria λ , a fim de obter diferentes proporções de sucessos (uns) nas amostras geradas: 50% (dados balanceados), 25% (moderado grau de desbalanceamento entre as classes) e 10% (alto grau de desbalanceamento). Sendo assim, foram totalizados 40 cenários simulados.

Para todos os modelos, os coeficientes de regressão foram fixados com os valores: $\beta_0 = 0$ e $\beta_1 = 1$, assim como no trabalho desenvolvido por Bazán, Romeo & Rodrigues [64]. A covariável foi gerada considerando $X \sim \text{Uniforme}(-4, 4)$. A partir dessas especificações, os valores da variável resposta Y foram simulados de uma distribuição de Bernoulli com parâmetro $\mu_i = F(\beta_0 + \beta_1 x)$, sendo $F(\cdot)$ a FDA padrão do modelo correspondente.

Os parâmetros dos quatro modelos foram estimados pelo método da MV. Para avaliar a performance dos estimadores de MV, foram calculados o viés e a raiz do erro quadrático médio (REQM) das estimativas, como segue:

$$\text{Viés}(\hat{\omega}_j) = \frac{1}{M} \sum_{m=1}^M (\hat{\omega}_j^{(m)} - \omega_j) = \left(\frac{1}{M} \sum_{m=1}^M \hat{\omega}_j^{(m)} \right) - \omega_j,$$

$$\text{REQM}(\hat{\omega}_j) = \sqrt{\frac{1}{M} \sum_{m=1}^M (\hat{\omega}_j^{(m)} - \omega_j)^2},$$

em que $M = 1.000$ é o número de réplicas de Monte Carlo, e $\hat{\omega}_j^{(m)}$ é a estimativa de MV do parâmetro ω_j na m -ésima amostra simulada.

Os resultados dessas duas medidas são apresentados nas Figuras 3.1 - 3.10, para cada modelo proposto e grau de desbalanceamento entre as classes. Na Figura 3.1, a única referente ao modelo DL, observa-se que ambos o viés e a REQM das estimativas de MV dos parâmetros β_0 e β_1 apresentam valores baixos (próximos a zero) quando n cresce, o que era esperado (desejado).

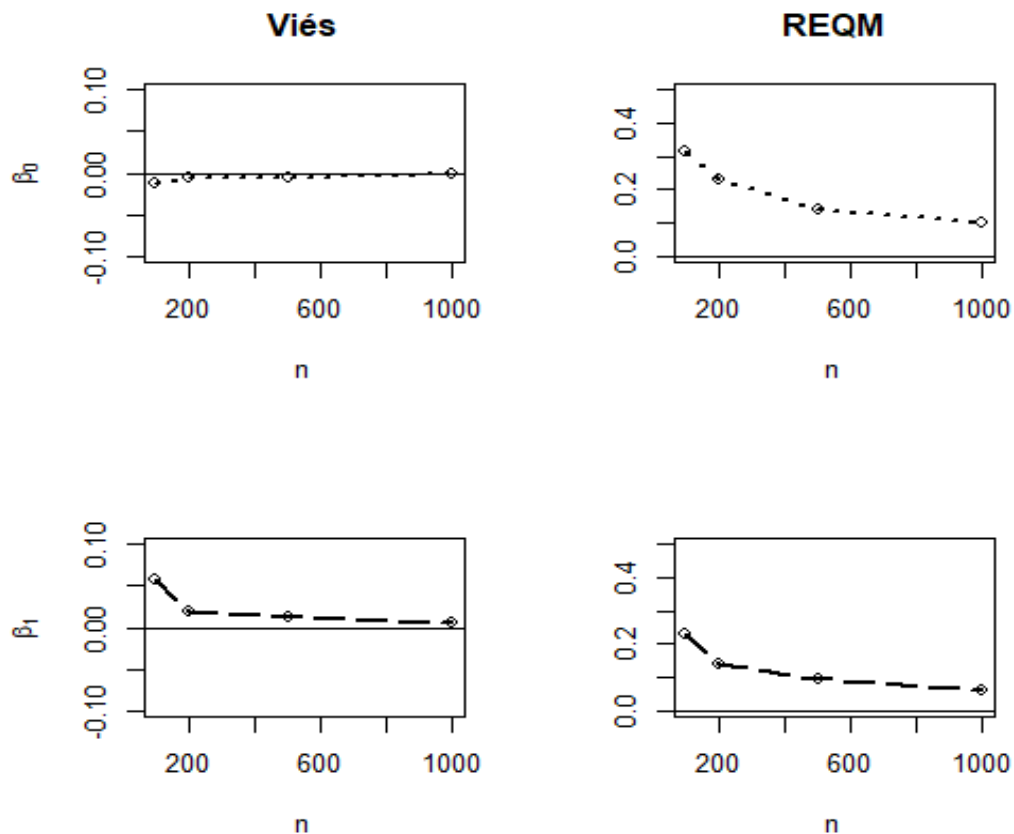


Figura 3.1: Viés e REQM das estimativas de MV dos parâmetros β_0 e β_1 do modelo DL, considerando diferentes tamanhos de amostra.

Para o modelo ADL, foram considerados os seguintes valores de λ : $\lambda = \exp\{0; 1,4; 2,4\}$, que correspondem, respectivamente, a 50%, 25% e 10% de uns nas amostras geradas. As Figuras 3.2, 3.3 e 3.4 ilustram os resultados obtidos para esses três casos, nas quais se observa que, de forma similar ao modelo DL, o viés e a REQM das estimativas de MV dos

parâmetros β_0 , β_1 e λ apresentam valores baixos quando n cresce, tornando as estimativas obtidas também melhores (isto é, mais acuradas e precisas) à medida que o tamanho amostral aumenta.

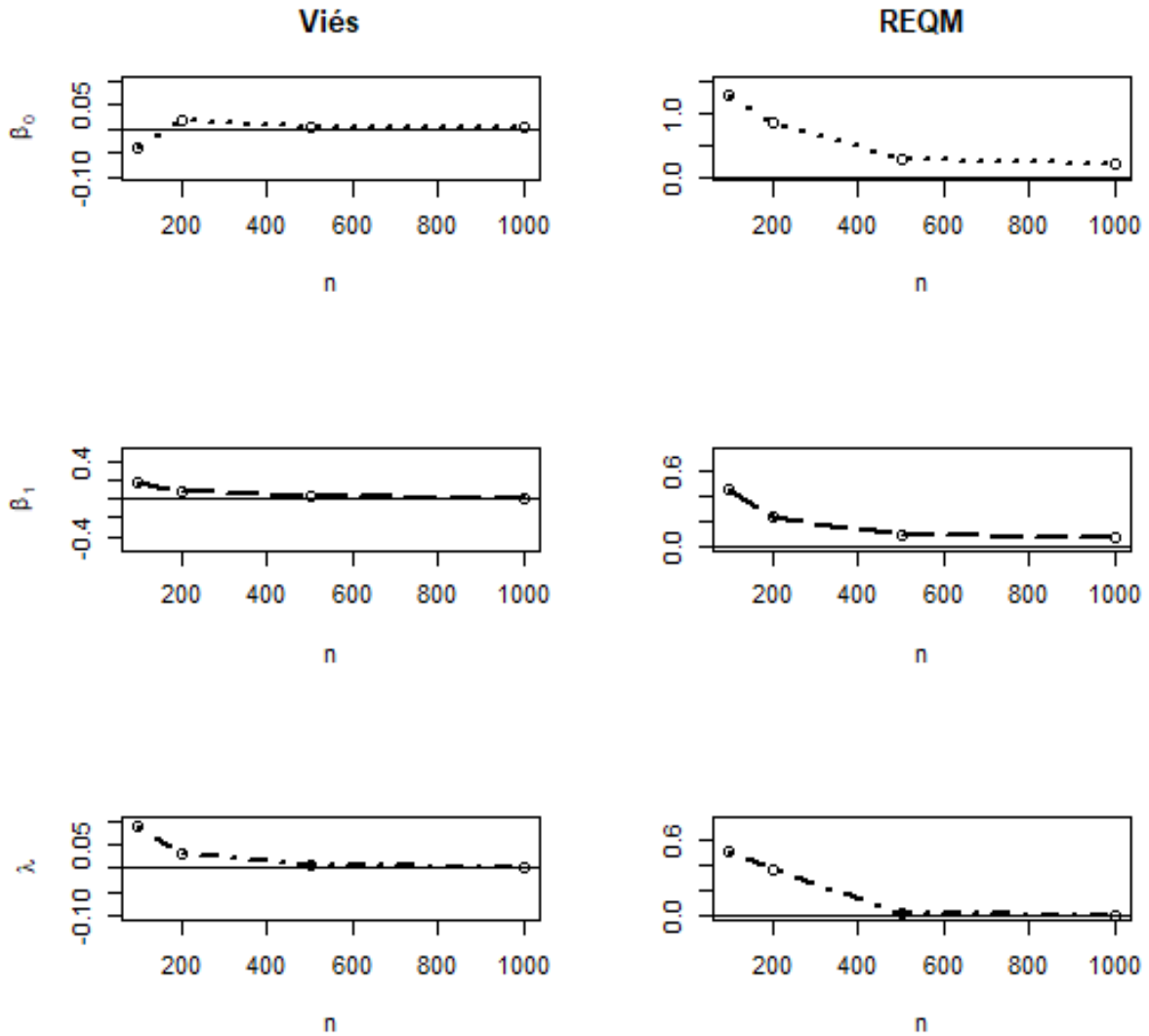


Figura 3.2: Viés e REQM das estimativas de MV dos parâmetros β_0 , β_1 e λ do modelo ADL, considerando diferentes tamanhos de amostra e $\lambda = \exp\{0\}$.

Assim como para o modelo ADL, considerou-se $\lambda = \exp\{0; 1,4; 2,4\}$ para o modelo PDL (ver Figuras 3.5, 3.6 e 3.7) e $\lambda = \exp\{0; -1,4; -2,4\}$ para o modelo RPD (ver Figuras 3.8, 3.9 e 3.10), a fim de obter 50%, 25% e 10% de uns nas amostras geradas, respectivamente. Novamente, são observados comportamentos similares para o viés e a REQM das estimativas produzidas, que tendem a assumir valores menores à medida que n aumenta.

Em resumo, os resultados obtidos demonstram uma recuperação adequada dos

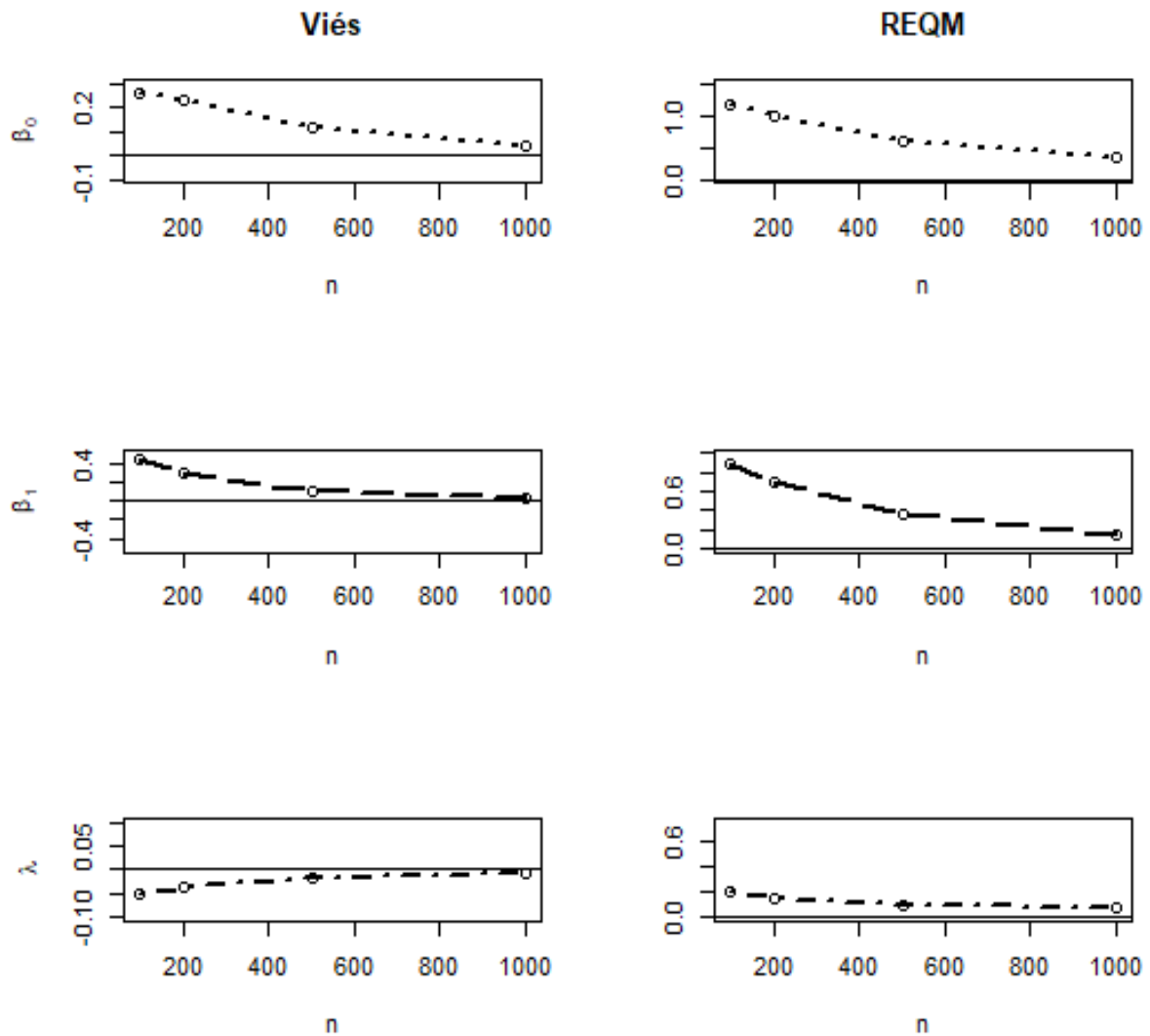


Figura 3.3: Viés e REQM das estimativas de MV dos parâmetros β_0 , β_1 e λ do modelo ADL, considerando diferentes tamanhos de amostra e $\lambda = \exp\{1,4\}$.

parâmetros dos quatro modelos aqui desenvolvidos, segundo o procedimento de estimação proposto (máxima verossimilhança).

3.2 Misspecification

O segundo estudo de simulação considerou $M = 1.000$ amostras (ou conjuntos de dados) de tamanho $n = 1.000$ de cada um dos modelos descritos, incluindo o modelo logístico. Foram fixados $\beta_0 = 0$, $\beta_1 = 1$ e $X \sim \text{Uniforme}(-4, 4)$, e também estabelecidos diferentes valores para o parâmetro λ , de modo a obter cerca de 50%, 25% e 10% de uns nas amostras geradas dos modelos assimétricos (ADL, PDL e RPD). Então, para

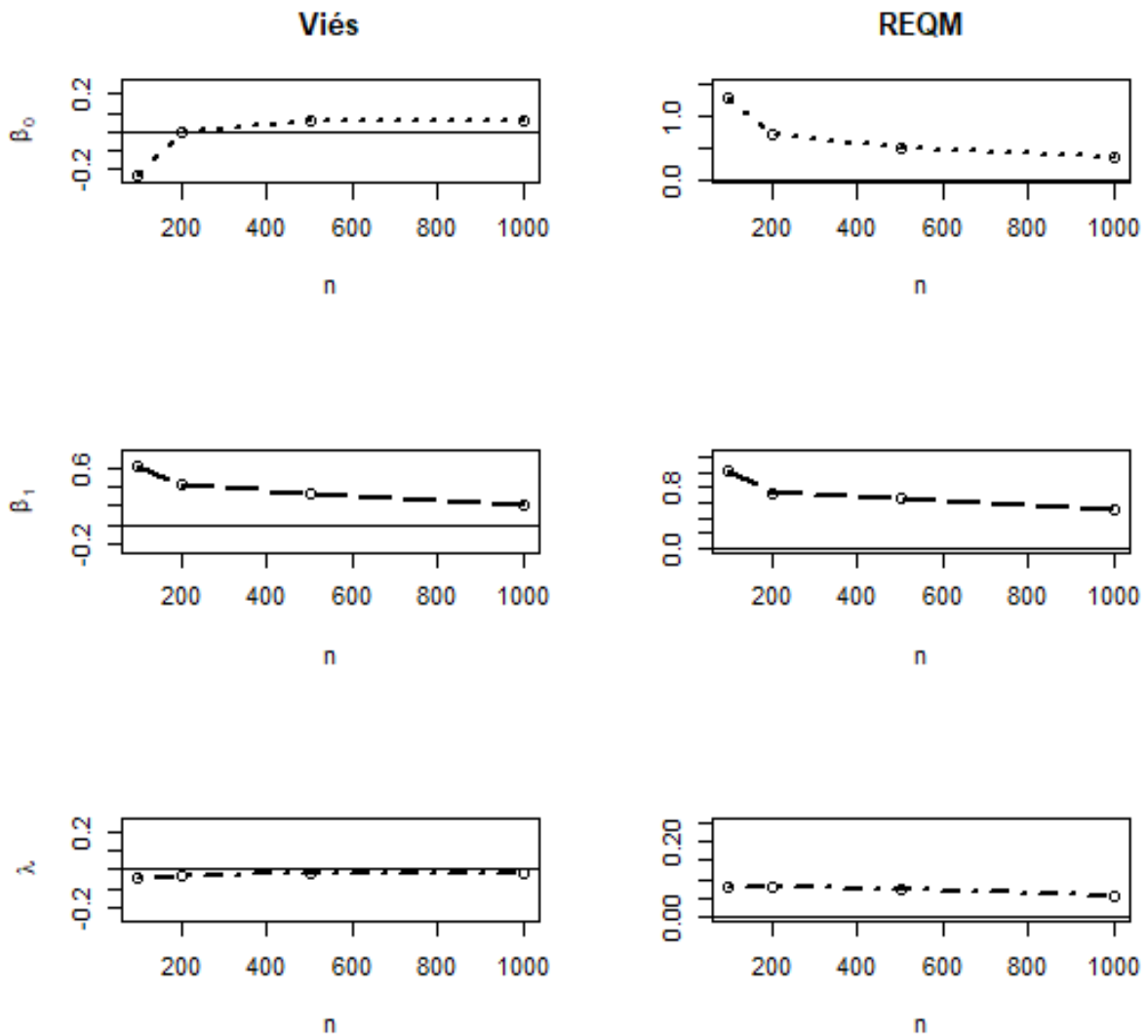


Figura 3.4: Viés e REQM das estimativas de MV dos parâmetros β_0 , β_1 e λ do modelo ADL, considerando diferentes tamanhos de amostra e $\lambda = \exp \{2,4\}$.

cada amostra simulada de cada um dos cinco modelos de regressão binária, foi empregado o seguinte processo: ajustou-se o modelo verdadeiro mais os outros quatro modelos candidatos e, em seguida, calculou-se os valores de AIC, BIC e viés relativo absoluto das probabilidades de evento (uns) estimadas para cada um deles, selecionando como o melhor modelo aquele que obteve os menores valores desses critérios.

A Tabela 3.1 mostra a proporção de vezes em que cada modelo forneceu o melhor ajuste, de acordo com os critérios AIC e BIC. De modo geral, observa-se uma boa conformidade desses critérios em discriminar os diferentes modelos estudados, sendo que ambos tiveram performances parecidas, porém com ligeira vantagem para o critério AIC nos resultados obtidos.

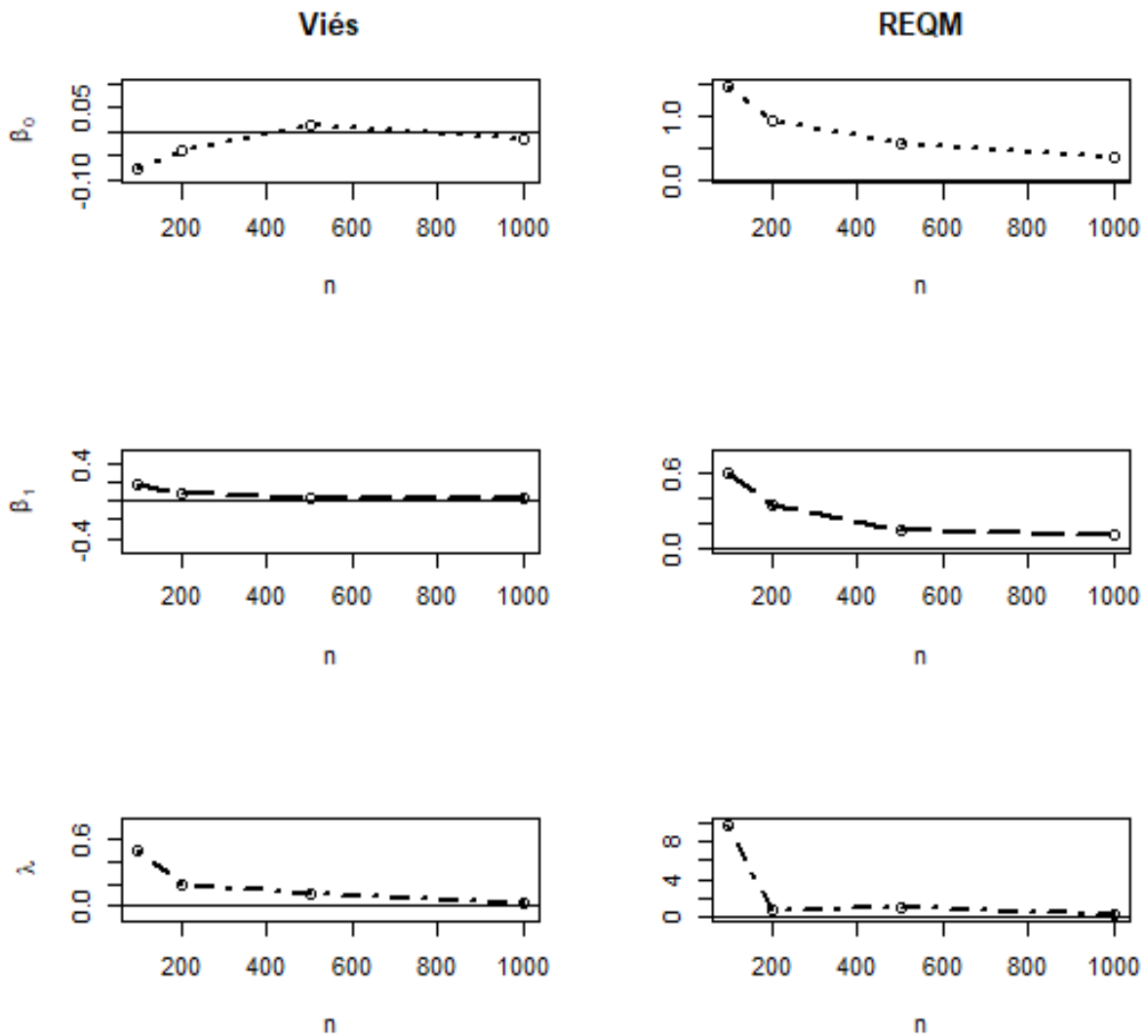


Figura 3.5: Viés e REQM das estimativas de MV dos parâmetros β_0 , β_1 e λ do modelo PDL, considerando diferentes tamanhos de amostra e $\lambda = \exp\{0\}$.

Na Tabela 3.2 são apresentados os resultados acerca do viés relativo absoluto na estimação das probabilidades de uns (“sucessos”). Observa-se que, nas situações de desbalanceamento entre as classes, os modelos assimétricos produzem as melhores estimativas das probabilidades de evento, superando, de maneira geral, os modelos logístico e DL, por apresentarem viés mais próximo a zero. Vale notar também que há uma certa superioridade do modelo DL sobre o modelo logístico nos casos em que o desbalanceamento é mais acentuado (10% de uns).

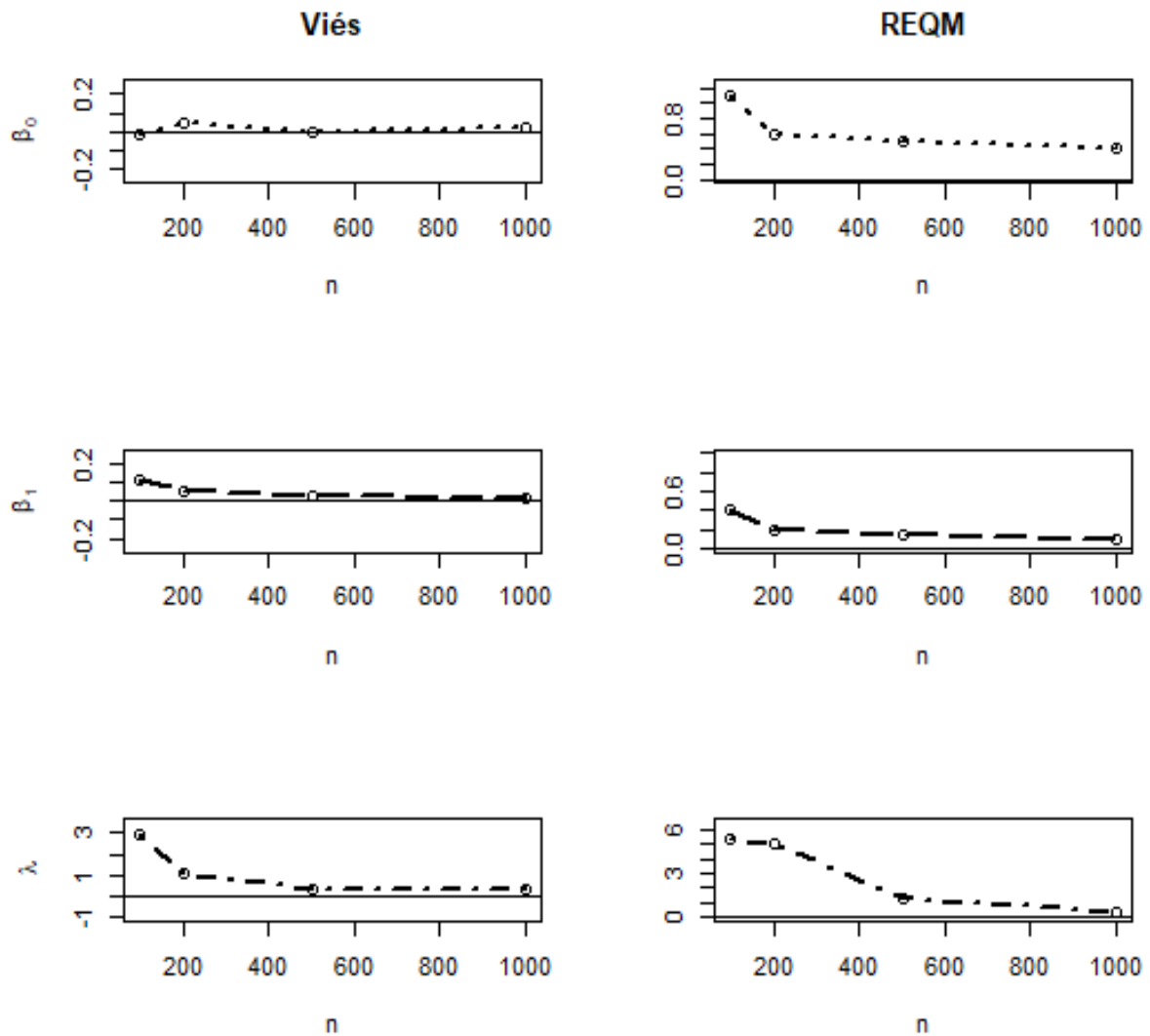


Figura 3.6: Viés e REQM das estimativas de MV dos parâmetros β_0 , β_1 e λ do modelo PDL, considerando diferentes tamanhos de amostra e $\lambda = \exp\{1,4\}$.

3.3 Avaliação da Capacidade Preditiva

Foi conduzido um terceiro estudo de simulação, com o propósito de avaliar a performance preditiva dos quatro modelos de classificação, em comparação com o modelo logístico (tradicional), segundo algumas das medidas de desempenho descritas na Seção 2.6 (SEN, VPP, *F1-Score*, ACCB e AUCPR). Baseando-se na geração proposta por Breiman [19], também utilizada por outros autores como Louzada *et al.* [56], Louzada & Ara [55] e Ferreira *et al.* [35], os valores das covariáveis para a classe de interesse (uns ou “sucessos”) foram simulados de uma distribuição normal multivariada (ou 6-variada) com vetor de médias $\boldsymbol{\mu}_S = (0, 0, 0, 0, 0, 0)^\top$ e matriz de covariâncias $4 \times I_6$, sendo I_6 a matriz

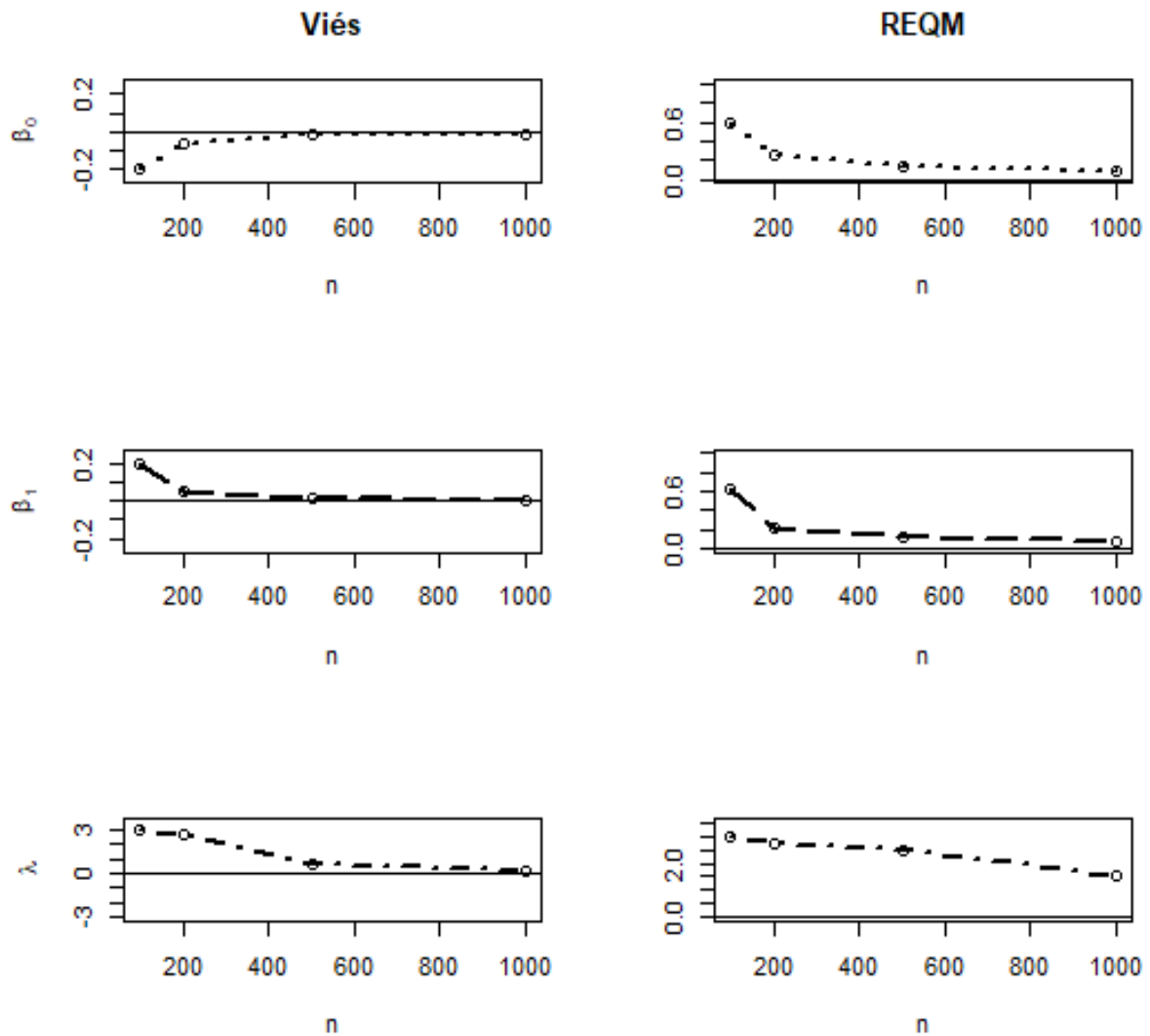


Figura 3.7: Viés e REQM das estimativas de MV dos parâmetros β_0 , β_1 e λ do modelo PDL, considerando diferentes tamanhos de amostra e $\lambda = \exp\{2,4\}$.

identidade de ordem 6. Já os valores das covariáveis para a outra classe (zeros ou “fracassos”) foram gerados de uma distribuição normal multivariada de dimensão 6 com vetor de médias igual a $\boldsymbol{\mu}_F = (1/\sqrt{6}, 1/\sqrt{6}, 1/\sqrt{6}, 1/\sqrt{6}, 1/\sqrt{6}, 1/\sqrt{6})^\top$ e matriz de covariâncias I_6 . Em seguida, a amostra final (constituída de observações dos dois grupos), com $n = 1.430$ e diferentes proporções de uns (50%, 25% e 10%), foi dividida aleatoriamente em 70% ($n = 1.000$) para treinamento e 30% ($n = 430$) para teste (validação *holdout*), verificando a performance preditiva dos cinco modelos candidatos na amostra teste. Este processo foi repetido $M = 1.000$ vezes.

Nas Tabelas 3.3 - 3.5 são apresentados os resultados acerca do desempenho preditivo médio dos cinco modelos estudados. Observa-se uma performance similar dos diferen-

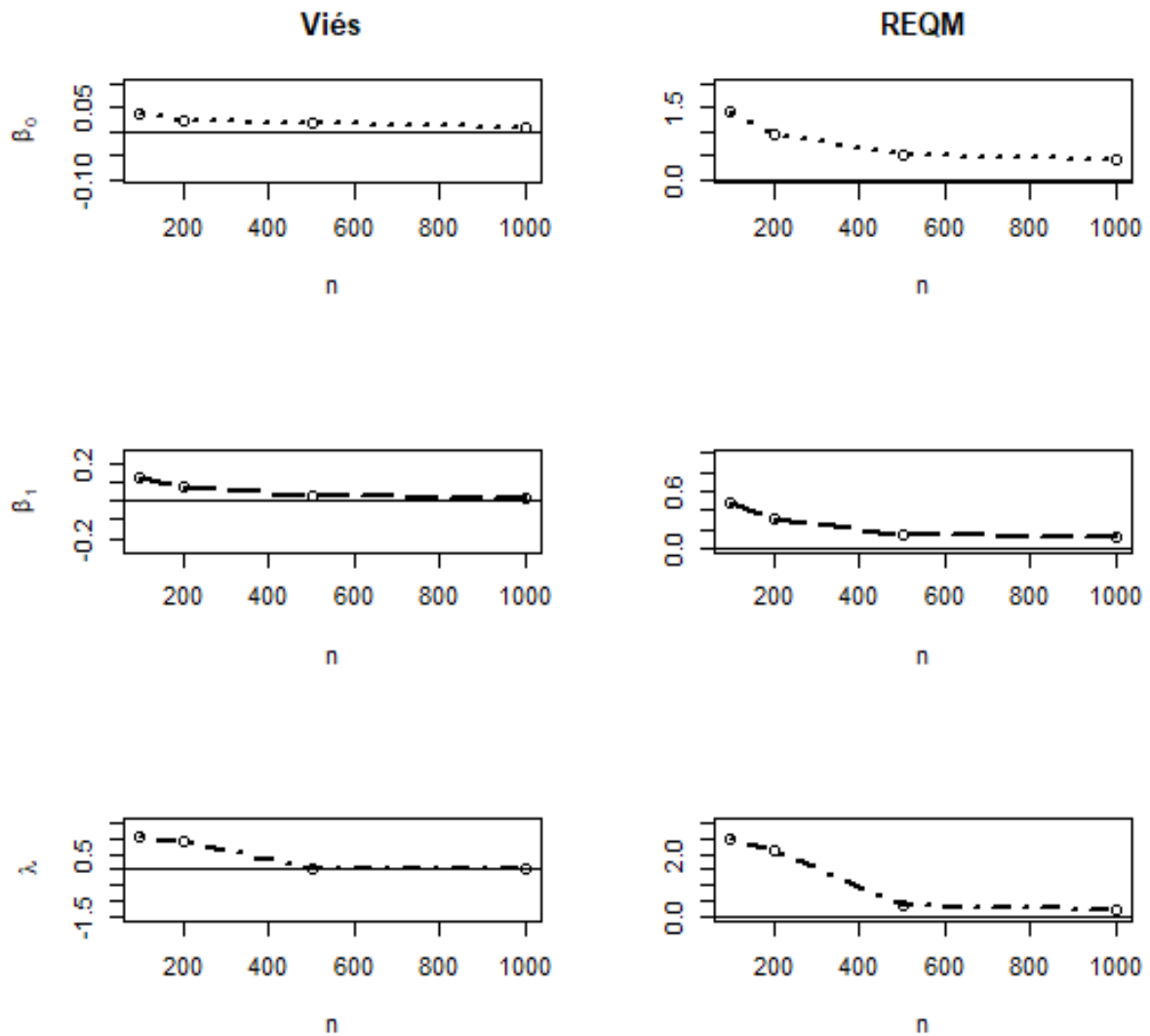


Figura 3.8: Viés e REQM das estimativas de MV dos parâmetros β_0 , β_1 e λ do modelo RPDL, considerando diferentes tamanhos de amostra e $\lambda = \exp\{0\}$.

tes modelos ajustados; em geral, com ligeira vantagem dos modelos novos sobre o modelo logístico, porém com diferenças observadas apenas na terceira ou quarta casa decimal. Nota-se, ainda, uma diminuição da capacidade preditiva de todos os modelos (segundo as métricas VPP e *F1-Score*, por exemplo) à medida que o grau de desbalanceamento entre as classes aumenta.

No entanto, vale ressaltar que, embora pequeno, tal ganho de performance pode fazer grande diferença na prática, por exemplo, evitando maiores prejuízos por parte das instituições financeiras ao aplicar esses novos modelos em problemas de *Credit Scoring*. Outro detalhe que é importante destacar, é que neste estudo de simulação, os pontos de corte (*cut-offs* ou *thresholds*) foram otimizados para cada modelo.

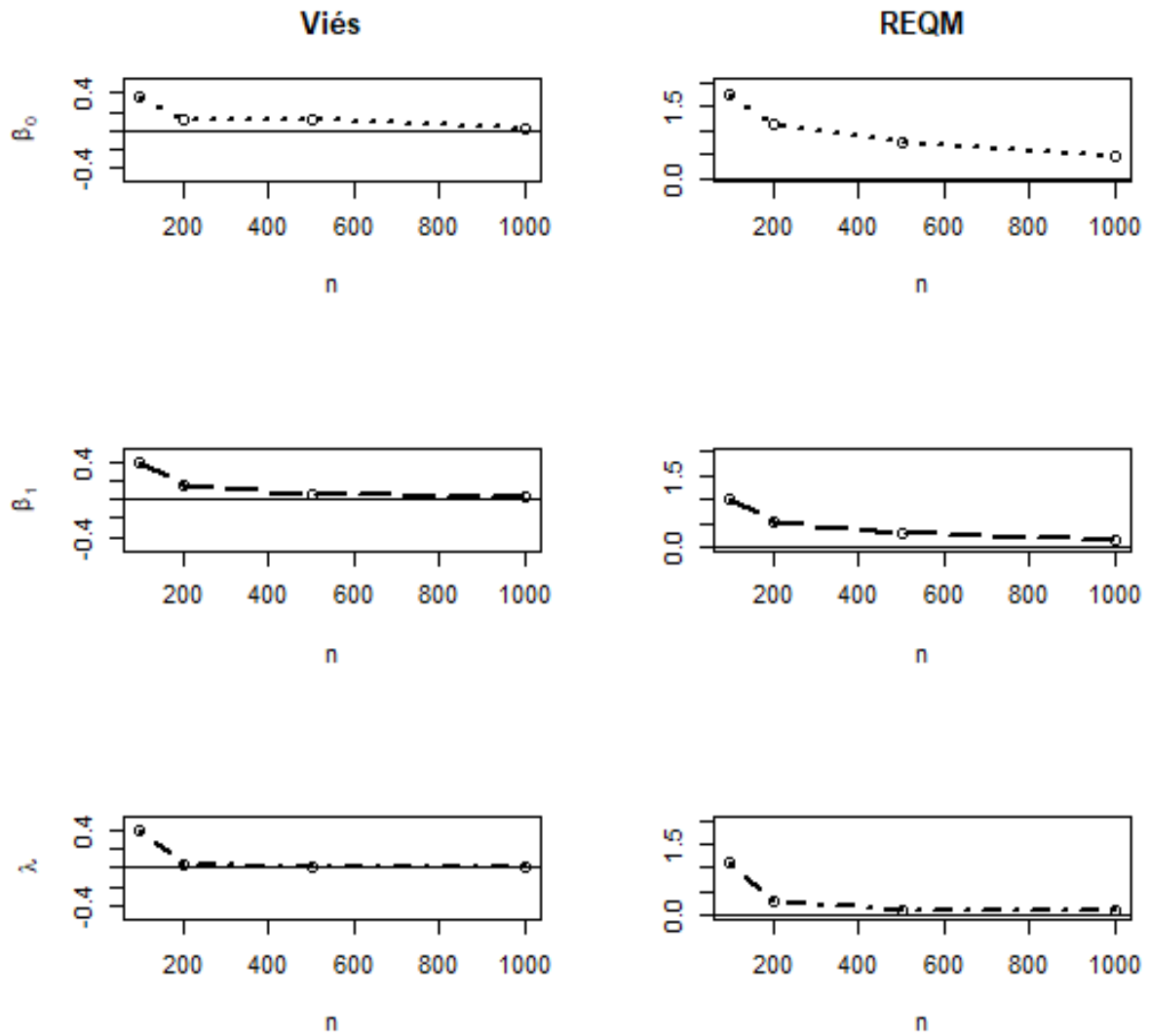


Figura 3.9: Viés e REQM das estimativas de MV dos parâmetros β_0 , β_1 e λ do modelo RPDL, considerando diferentes tamanhos de amostra e $\lambda = \exp\{-1,4\}$.

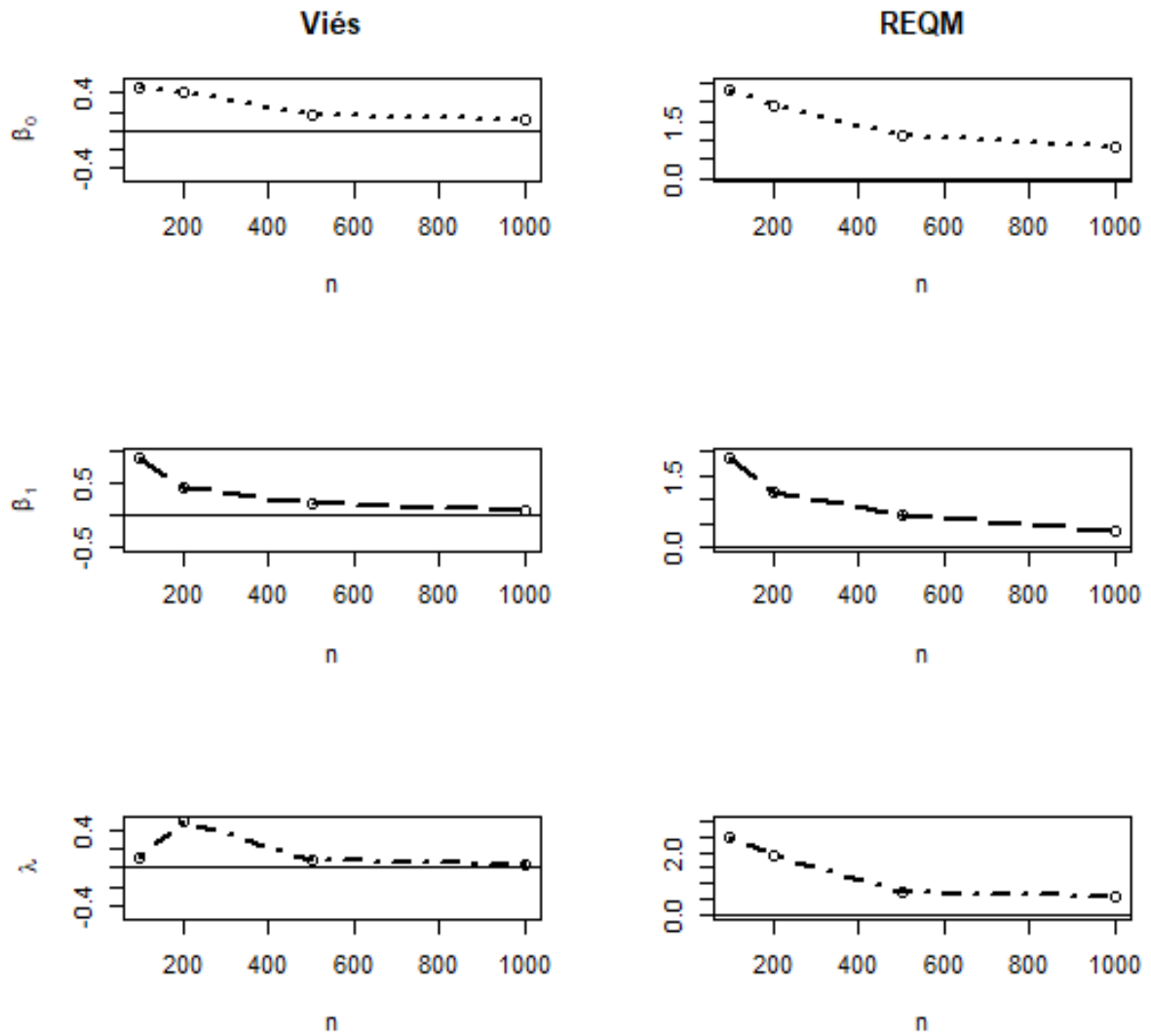


Figura 3.10: Viés e REQM das estimativas de MV dos parâmetros β_0 , β_1 e λ do modelo RPDL, considerando diferentes tamanhos de amostra e $\lambda = \exp\{-2,4\}$.

Tabela 3.1: Proporção de vezes em que cada modelo de regressão binária foi eleito o melhor segundo os critérios AIC e BIC.

Modelo Verdadeiro	Modelo Ajustado									
	Logístico		DL		ADL		PDL		RPDL	
	AIC	BIC	AIC	BIC	AIC	BIC	AIC	BIC	AIC	BIC
Logístico	0,502	0,584	0,340	0,408	0,045	0,003	0,048	0,001	0,065	0,004
DL	0,346	0,396	0,505	0,593	0,065	0,007	0,039	0,002	0,045	0,002
ADL ($\lambda = \exp\{0\}$)	0,377	0,422	0,487	0,572	0,065	0,001	0,034	0,003	0,034	0,002
ADL ($\lambda = \exp\{1,4\}$)	0,000	0,022	0,000	0,000	0,677	0,672	0,082	0,066	0,241	0,240
ADL ($\lambda = \exp\{2,4\}$)	0,000	0,056	0,000	0,000	0,890	0,868	0,028	0,001	0,082	0,075
PDL ($\lambda = \exp\{0\}$)	0,386	0,427	0,490	0,569	0,045	0,002	0,040	0,001	0,039	0,001
PDL ($\lambda = \exp\{1,4\}$)	0,046	0,314	0,010	0,016	0,452	0,423	0,453	0,211	0,039	0,036
PDL ($\lambda = \exp\{2,4\}$)	0,200	0,152	0,005	0,005	0,181	0,124	0,562	0,688	0,052	0,031
RPDL ($\lambda = \exp\{0\}$)	0,377	0,424	0,497	0,570	0,041	0,002	0,041	0,002	0,044	0,002
RPDL ($\lambda = \exp\{-1,4\}$)	0,167	0,613	0,009	0,056	0,401	0,197	0,074	0,006	0,349	0,128
RPDL ($\lambda = \exp\{-2,4\}$)	0,306	0,777	0,027	0,028	0,332	0,170	0,062	0,000	0,273	0,025

Tabela 3.2: Viés relativo absoluto na estimação das probabilidades de evento (uns).

Modelo Verdadeiro	Modelo Ajustado				
	Logístico	DL	ADL	PDL	RPDL
Logístico	0,1406	0,1517	0,1424	0,1484	0,1549
DL	0,0346	0,0246	0,0497	0,0372	0,0367
ADL ($\lambda = \exp\{0\}$)	0,0392	0,0283	0,0125	0,0139	0,0173
ADL ($\lambda = \exp\{1,4\}$)	8,0692	9,9745	0,7554	1,4404	0,8434
ADL ($\lambda = \exp\{2,4\}$)	7,4711	1,1899	0,1218	8,7562	1,8160
PDL ($\lambda = \exp\{0\}$)	0,1155	0,1192	0,0987	0,0897	0,1097
PDL ($\lambda = \exp\{1,4\}$)	6,4018	9,2803	0,4189	0,4504	1,3930
PDL ($\lambda = \exp\{2,4\}$)	8,2054	2,2249	1,7093	0,3471	0,6567
RPDL ($\lambda = \exp\{0\}$)	0,0825	0,0868	0,0541	0,0630	0,0531
RPDL ($\lambda = \exp\{-1,4\}$)	0,2159	0,2888	0,0588	0,2031	0,0521
RPDL ($\lambda = \exp\{-2,4\}$)	0,4160	0,5066	0,3068	0,0849	0,2806

Tabela 3.3: Desempenho preditivo médio dos modelos candidatos na amostra teste (50% de uns).

Modelo	Medida				
	SEN	VPP	F1-Score	ACCB	AUCPR
Logístico	0,8712	0,6184	0,7225	0,6672	0,5652
DL	0,8715	0,6184	0,7226	0,6672	0,5652
ADL	0,8728	0,6182	0,7229	0,6672	0,5653
PDL	0,8712	0,6184	0,7224	0,6671	0,5653
RPDL	0,8734	0,6182	0,7231	0,6674	0,5654

Tabela 3.4: Desempenho preditivo médio dos modelos candidatos na amostra teste (25% de uns).

Modelo	Medida				
	SEN	VPP	F1-Score	ACCB	AUCPR
Logístico	0,8635	0,3538	0,5008	0,6676	0,3116
DL	0,8632	0,3539	0,5009	0,6677	0,3115
ADL	0,8627	0,3541	0,5010	0,6678	0,3115
PDL	0,8630	0,3537	0,5007	0,6675	0,3116
RPDL	0,8633	0,3539	0,5009	0,6677	0,3113

Tabela 3.5: Desempenho preditivo médio dos modelos candidatos na amostra teste (10% de uns).

Modelo	Medida				
	SEN	VPP	F1-Score	ACCB	AUCPR
Logístico	0,8488	0,1538	0,2595	0,6648	0,1326
DL	0,8500	0,1538	0,2595	0,6651	0,1326
ADL	0,8498	0,1537	0,2594	0,6649	0,1325
PDL	0,8502	0,1538	0,2596	0,6652	0,1326
RPDL	0,8502	0,1538	0,2596	0,6652	0,1326

Capítulo 4

Aplicações

Neste capítulo são apresentadas aplicações da metodologia descrita no Capítulo 3 a dois conjuntos de dados reais, sendo um deles oriundo da área de saúde (Seção 4.1) e o outro da área financeira (Seção 4.2).

4.1 *Breast Cancer Wisconsin (Diagnostic) Data Set*

Este conjunto de dados é referente ao diagnóstico de câncer de mama, se o tumor detectado é maligno ou benigno, e foi doado no ano de 1995 para o repositório de *Machine Learning* da Universidade da Califórnia em Irvine, contendo 569 observações (ou instâncias) e 32 atributos (ou variáveis), a partir de amostras que chegam periodicamente quando o Dr. William H. Wolberg, do Hospital da Universidade de Wisconsin em Madison, relata seus casos clínicos. Essas características/variáveis descrevem o núcleo da célula presente na imagem obtida. Caso um nódulo seja encontrado, uma biópsia de aspiração com agulha fina é realizada, que utiliza uma agulha oca para extrair uma pequena amostra de células. Um médico clínico, em seguida, examina as células sob um microscópio para determinar se a massa é provável que seja maligna ou benigna.

O tumor maligno representa a classe positiva (“sucesso”) com 212 observações, e o tumor benigno representa a classe negativa (“fracasso”) com 357 observações. O banco de dados pode ser acessado pelo endereço: [https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)).

As variáveis independentes utilizadas nesta ilustração foram somente quatro variáveis contínuas (ver descrição na Tabela 4.1), e sem a presença de valores faltantes (*missing data*). Elas foram escolhidas por serem consideradas importantes para o diagnóstico da doença. A variável dependente, por sua vez, é uma resposta binária que indicará se o tumor foi diagnosticado como maligno (codificado como “1” pelo estudo) ou benigno (codificado como “0” pelo estudo).

Tabela 4.1: Variáveis consideradas na aplicação da metodologia proposta aos dados sobre câncer de mama.

Variável	Notação	Classificação
Diagnóstico	Y	Categórica 0-1
Compactação	X_1	Contínua
Suavidade	X_2	Contínua
Pior Raio	X_3	Contínua
Pior Dimensão	X_4	Contínua

Observa-se que, das 569 pacientes que foram diagnosticadas com câncer de mama, 37,26% tiveram diagnóstico de tumor maligno, enquanto que 62,74% tiveram diagnóstico de tumor benigno.

Os principais resultados de uma análise descritiva inicial para as variáveis escolhidas neste estudo, são apresentados na Tabela 4.2. Para verificar se há indícios de multicolinearidade, foi calculada a matriz de correlação entre as covariáveis (Figura 4.1); observa-se, no entanto, que nenhum dos coeficientes de correlação possui valor absoluto maior do que 0,70, por exemplo, o que representa, conforme indicado por Mukaka [57], que não existe uma correlação forte entre as variáveis. Ademais, nos gráficos da Figura 4.2, observa-se que as covariáveis referentes à compactação (X_1) e pior raio (X_3), aparentam ter maior possibilidade de prever o diagnóstico de um tumor maligno.

Tabela 4.2: Algumas estatísticas descritivas das variáveis independentes selecionadas. CV = coeficiente de variação.

Variável	Medidas-Resumo							
	Mínimo	Q_1	Mediana	Média	Q_3	Máximo	Desvio-Padrão	CV (%)
X_1	0,0193	0,0649	0,0926	0,1043	0,1304	0,3454	0,0529	50,61
X_2	0,0017	0,0052	0,0064	0,0070	0,0081	0,0311	0,0030	42,64
X_3	7,9300	13,0100	14,9700	16,2700	18,7900	36,0400	4,8332	29,70
X_4	0,0550	0,0715	0,0800	0,0839	0,0921	0,2075	0,0181	21,51

Os resultados da estimação para os novos modelos de regressão binária propostos, assim como para o modelo logístico usual, estão disponíveis na Tabela 4.3. Dentre outros, observa-se uma concordância do sinal do intercepto (negativo) e também dos coeficientes de regressão (todos positivos e de magnitudes parecidas para os diferentes modelos). Verifica-se também que as quatro covariáveis são relevantes nos modelos assimétricos, com exceção de X_1 (Compactação) para o modelo ADL. Quanto às estimativas dos erros padrões, estas foram calculadas a partir da raiz quadrada dos elementos da diagonal da inversa da matriz de informação observada; nota-se que os modelos que obtiveram

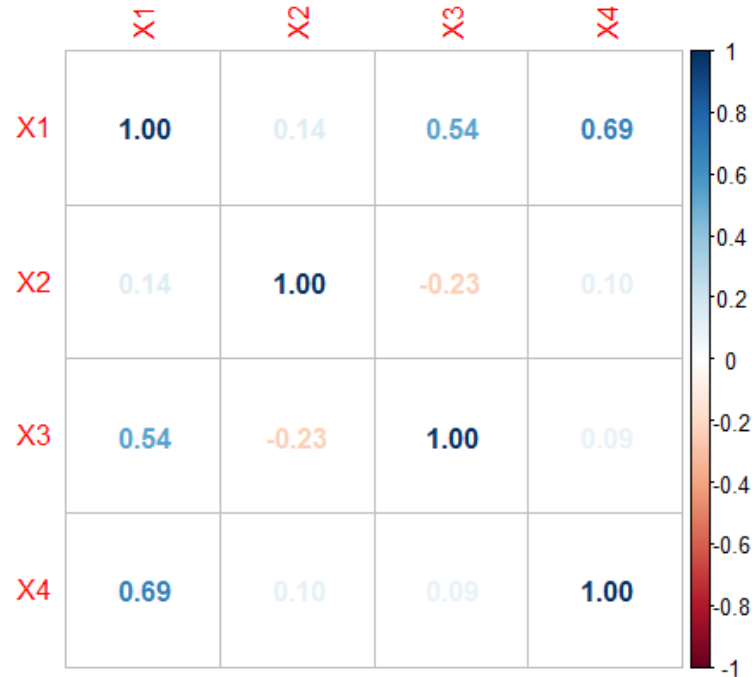


Figura 4.1: Matriz de correlação das covariáveis.

Tabela 4.3: Resultados da estimação para os modelos de regressão binária ajustados ao conjunto de dados sobre câncer de mama.

Parâmetro	Modelo									
	Logístico		DL		ADL		PDL		RPDL	
	Estimativa	Erro Padrão	Estimativa	Erro Padrão	Estimativa	Erro Padrão	Estimativa	Erro Padrão	Estimativa	Erro Padrão
β_0	-31,4110*	3,2099	-34,3974*	3,1321	-79,4003*	7,3827	-75,7806*	0,0350	-43,0790*	0,0053
β_1	6,2539	15,3341	-0,4306	8,3633	2,4894	3,3323	16,2654*	0,0884	7,9200*	0,0090
β_2	66,1407*	3,7874	186,8435*	7,2132	553,2222*	9,2832	134,3720*	0,1256	36,0431*	0,0070
β_3	1,3116*	0,1685	1,5165*	0,1819	3,2344*	0,3718	3,1878*	0,0357	1,1377*	0,0023
β_4	71,9795*	16,0975	92,1532*	4,9490	205,6775*	16,2234	190,7268*	0,0189	63,6264*	0,0065
λ	–	–	–	–	2,8855*	0,6976	0,2976*	0,1984	15,2446*	0,1941

*Significante ao nível de 10% (i.e., $\beta_j \neq 0$, para $j = 0, 1, 2, 3, 4$, e $\lambda \neq 1$).

melhores resultados, ou seja, menores valores para os erros padrões, foram os modelos PDL e RPDL. Pode-se ver também que no modelo ADL, por exemplo, $\hat{\lambda} = 2,8855 > 1$, o que indica assimetria positiva (ou à direita), sendo este resultado coerente à proporção de 1's na variável resposta (menor que 0,5).

De acordo com os valores dos critérios de informação (AIC e BIC), exibidos na Tabela 4.4, seleciona-se o modelo de regressão binária ADL como sendo o de “melhor” ajuste, por ter apresentado os menores valores de AIC e BIC. Porém, para não basear a seleção apenas nesses critérios de informação, é recomendável também investigar a capacidade preditiva desses modelos. Por isso, na sequência da análise foi avaliada a performance preditiva dos modelos ajustados.

A partir da Tabela 4.5, observa-se que os valores de AUC para os cinco modelos

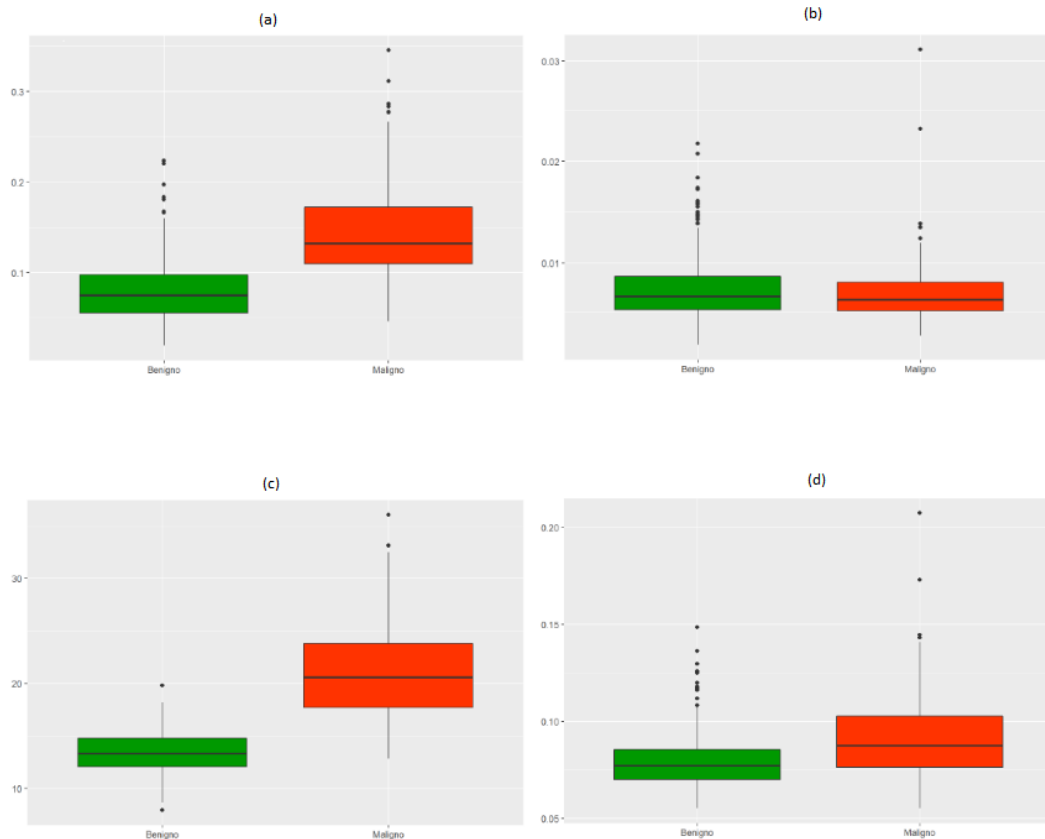


Figura 4.2: Diagnóstico de tumor (benigno ou maligno) segundo as variáveis: (a) Compactação, (b) Suavidade, (c) Pior Raio, e (d) Pior Dimensão.

candidatos são bastante parecidos entre si e próximos de 1, o que indica que todos eles são igualmente adequados. Repare ainda que os valores do ponto de corte ótimo para esses modelos não são muito parecidos entre si, sendo o ponto de corte ótimo do modelo RPDL o mais próximo de 0,5, e o do modelo logístico inferior ao dos demais modelos para o diagnóstico de tumor maligno.

Analisando-se os resultados mostrados na Tabela 4.6, observa-se que o modelo de regressão binária PDL apresentou melhor performance preditiva (em termos de ACC, MCC e *F1-Score*, por exemplo) para o diagnóstico de tumor maligno, quando avaliada na amostra teste (30%). Observa-se ainda que o modelo ADL, embora selecionado pelos critérios AIC e BIC como o de melhor ajuste, não demonstrou tão bom desempenho preditivo quando avaliado na menor parcela dos dados originais separada para este fim.

Portanto, com relação aos procedimentos apresentados, observa-se que os modelos de regressão binária com as funções de ligação aqui propostas apresentaram resultados promissores, com ajuste e performance preditiva melhores (neste caso, ligeiramente me-

Tabela 4.4: Valores de AIC e BIC para os modelos de regressão binária ajustados ao conjunto de dados sobre câncer de mama.

Critério	Modelo				
	Logístico	DL	ADL	PDL	RPDL
AIC	119,2507	115,7655	110,8239	116,3451	131,6976
BIC	139,1830	135,6978	134,7427	140,2638	155,6163

Tabela 4.5: Comparação do desempenho preditivo dos diferentes classificadores (amostra treinamento - 70%).

Modelo	Medida				
	AUC	Ponto de Corte Ótimo	SPE	SEN	
Logístico	0,9842	0,3263	0,9341	0,9500	
DL	0,9846	0,4695	0,9729	0,9143	
ADL	0,9845	0,3896	0,9806	0,9143	
PDL	0,9843	0,4160	0,9767	0,9143	
RPDL	0,9850	0,4799	0,9753	0,9032	

lhores, mas que pode fazer muita diferença na prática) do que os modelos usando ligações tradicionais, como a logito.

4.2 *Santander Customer Transaction Prediction Dataset*

O segundo conjunto de dados utilizado neste trabalho é referente a uma competição promovida pelo Banco Santander para a comunidade do Kaggle, cujo objetivo era prever se um cliente iria efetuar uma operação financeira específica no futuro, independentemente da quantia transacionada. A base de dados é composta de 200.000 observações (ou instâncias) e 200 variáveis preditoras contínuas, normalizadas, sem a presença de valores faltantes e anonimizadas. A variável resposta binária *Target*, se o cliente irá realizar uma transação ou não, está distribuída da seguinte forma: 10% de casos positivos (“sucessos”) e 90% de casos negativos (“fracassos”). Portanto, essa base é altamente desbalanceada, em uma proporção de 9 casos negativos para cada 1 positivo. O banco de dados pode ser acessado pelo endereço: <https://www.kaggle.com/lakshmi25npathi/santander-customer-transaction-prediction-dataset>.

Uma vez que existe um alto grau de desbalanceamento entre as classes (zeros e uns)

Tabela 4.6: Medidas de avaliação da performance preditiva para os modelos de regressão binária ajustados ao conjunto de dados sobre câncer de mama (amostra teste - 30%).

Modelo	Medida							
	ACC	SEN	SPE	VPP	VPN	MCC	F1-Score	ACCB
Logístico	0,9357	0,9167	0,9495	0,9296	0,9400	0,8679	0,9231	0,9331
DL	0,9415	0,9028	0,9697	0,9559	0,9320	0,8802	0,9286	0,9362
ADL	0,8538	0,6250	0,9907	0,9756	0,8154	0,6978	0,7619	0,8078
PDL	0,9474	0,9028	0,9798	0,9701	0,9327	0,8927	0,9446	0,9413
RPDL	0,9298	0,8906	0,9533	0,9194	0,9358	0,8495	0,9048	0,9219

da variável resposta, o uso de ligações simétricas poderia não ser adequado na modelagem de regressão binária do conjunto de dados em questão (Chen *et al.* [25]).

Para a seleção das variáveis (do total de 200) que mais contribuem para prever a realização ou não da transação específica, utilizou-se a técnica de *random forest*, ou floresta aleatória (para maiores detalhes, ver, por exemplo, Breiman [18] e Dellier [31]). Desta forma, as variáveis escolhidas como preditoras para os cinco modelos candidatos (logístico, DL, ADL, PDL e RPDL) foram sete, descritas na Tabela 4.7.

Tabela 4.7: Variáveis consideradas na aplicação da metodologia proposta aos dados do Banco Santander.

Variável	Notação	Classificação
<i>Target</i>	Y	Categórica 0-1
<i>var81</i>	X_1	Contínua
<i>var139</i>	X_2	Contínua
<i>var6</i>	X_3	Contínua
<i>var53</i>	X_4	Contínua
<i>var110</i>	X_5	Contínua
<i>var26</i>	X_6	Contínua
<i>var146</i>	X_7	Contínua

De acordo com os valores dos critérios de informação, apresentados na Tabela 4.8, seleciona-se o modelo de regressão binária DL como sendo o de “melhor” ajuste, por ter apresentado os (ligeiramente) menores valores de AIC e BIC.

Na Figura 4.3 são apresentadas as curvas ROC para os cinco modelos, em que os valores do AUC obtidos são próximos a 0,73. Isso indica que todos esses modelos de regressão binária são igualmente adequados.

Tabela 4.8: Valores de AIC e BIC para os modelos de regressão binária ajustados ao conjunto de dados do Banco Santander.

Critério	Modelo				
	Logístico	DL	ADL	PDL	RPDL
AIC	598,8691	598,6489	601,2960	600,8749	600,6522
BIC	638,1312	637,9109	645,4658	645,0447	644,8220

Verificando os resultados expostos na Tabela 4.9, nota-se que o modelo PDL apresentou o melhor desempenho preditivo (entre todas as medidas exploradas), para verificar se um cliente vai efetuar uma operação financeira específica. Superou, inclusive, o modelo DL que, por sua vez, performou melhor do que o modelo logístico. Tal avaliação foi feita na amostra teste (30%).

Tabela 4.9: Medidas de avaliação da performance preditiva para os modelos de regressão binária ajustados ao conjunto de dados do Banco Santander (amostra teste - 30%).

Modelo	Medida				
	SEN	VPP	F1-Score	ACCB	AUCPR
Logístico	0,5909	0,0850	0,1486	0,6239	0,1036
DL	0,6364	0,0864	0,1522	0,6368	0,1020
ADL	0,5455	0,0870	0,1500	0,6183	0,1001
PDL	0,6818	0,0904	0,1596	0,6559	0,1045
RPDL	0,5455	0,0789	0,1379	0,6012	0,1008

Considerando os resultados das medidas de avaliação da performance preditiva para os cinco modelos de regressão binária (Tabela 4.9), pode-se então selecionar o modelo PDL. Portanto, a probabilidade de que um cliente i irá efetuar uma operação financeira específica, de acordo com o modelo escolhido (PDL), é calculada por:

$$\hat{\mu}_i = F(\hat{\eta}_i) = \left(\frac{1}{2} + \frac{1}{2} \text{sgn}(\hat{\eta}_i) \left\{ 1 - \left(1 + \frac{|\hat{\eta}_i|}{2} \right) e^{-|\hat{\eta}_i|} \right\} \right)^{0,2313},$$

sendo:

$$\hat{\eta}_i = -4,2616 - 0,1157x_{i1} - 0,0417x_{i2} + 0,4724x_{i3} + 0,5474x_{i4} + 0,0579x_{i5} + 0,0697x_{i6} - 0,1592x_{i7}.$$

Finalmente, a regra de decisão (ou de classificação) para alocar um cliente i numa das duas classes (1 - realização da transação específica, 0 - não realização da transação específica) é:

$$\begin{aligned} i \in 1 \quad (\text{ou } \hat{Y}_i = 1) & \quad \text{se} \quad \hat{\mu}_i \geq 0,094, \\ i \in 0 \quad (\text{ou } \hat{Y}_i = 0) & \quad \text{se} \quad \hat{\mu}_i < 0,094. \end{aligned}$$

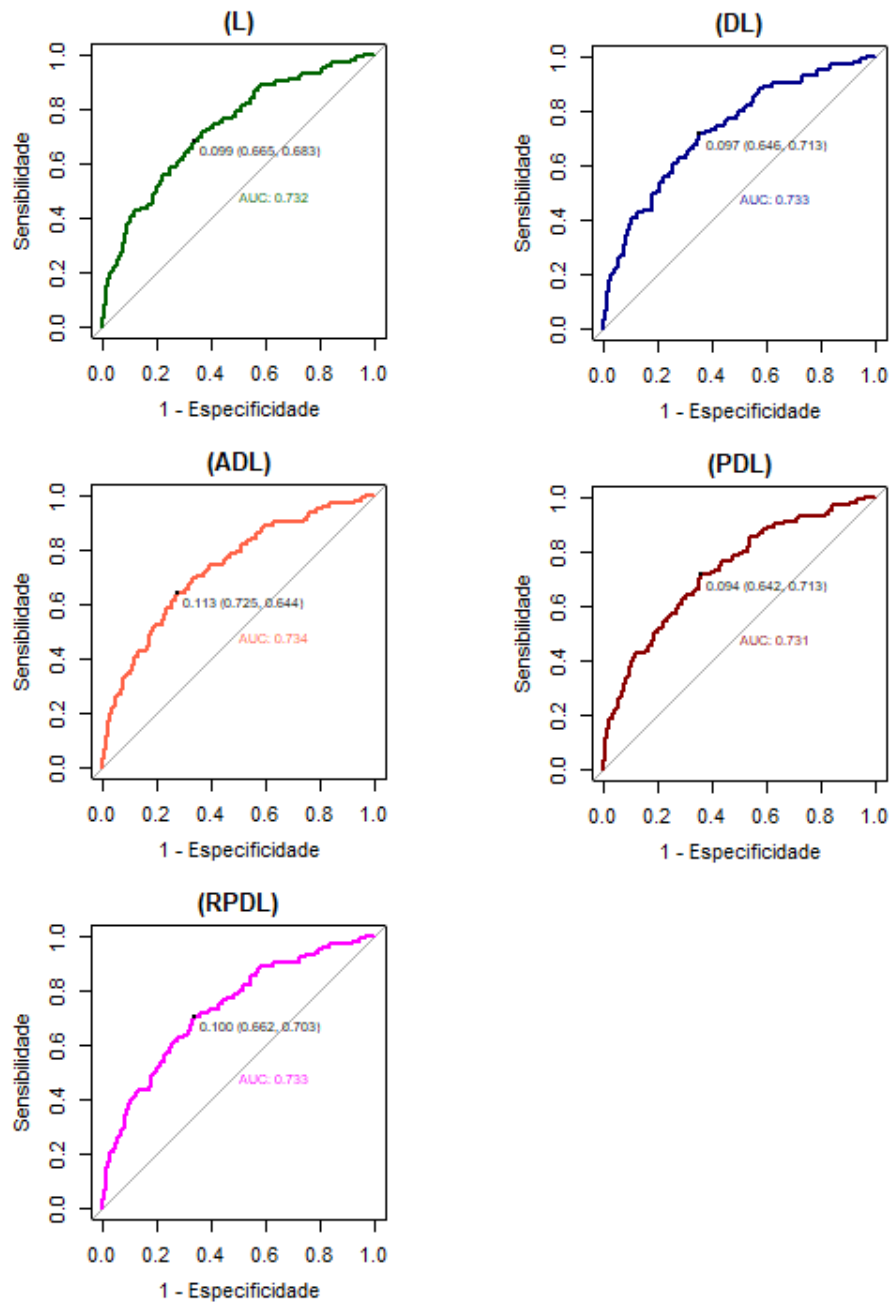


Figura 4.3: Curva ROC dos modelos de regressão binária logística, DL, ADL, PDL e RPDL, ajustados ao conjunto de dados do Banco Santander (amostra treinamento - 70%).

Capítulo 5

Considerações Finais e Trabalhos Futuros

Neste trabalho de dissertação, foi inicialmente realizada uma breve revisão dos principais métodos estatísticos e alguns métodos computacionais para lidar com problemas de classificação binária na presença de dados desbalanceados. No que tange à modelagem de regressão binária, foram consideradas algumas funções de ligação alternativas à logito (tradicional), dentre as quais, as ligações assimétricas do tipo potência e reversa de potência.

Dentre as principais contribuições deste trabalho, foram introduzidos novos modelos de regressão para a análise de dados binários, baseados nas distribuições double Lindley, double Lindley assimétrica, potência double Lindley e reversa de potência double Lindley (sendo as duas últimas distribuições também inéditas na literatura).

Estudos de simulação foram realizados com os objetivos, dentre outros, de avaliar a performance do método de estimação considerado (máxima verossimilhança), bem como a capacidade preditiva dos modelos de regressão binária propostos.

Por fim, dois bancos de dados reais foram utilizados para ilustrar a aplicabilidade dos modelos e métodos propostos. Dentre outros, verificou-se que os modelos de regressão binária com as funções de ligação potência e reversa de potência aqui propostas apresentaram resultados promissores, com ajuste e performance preditiva melhores do que os modelos de regressão binária usando ligações comuns (e.g., logito).

Como sugestão de desenvolvimentos futuros, pode-se propor classes gerais (e inéditas) de modelos de regressão para a análise de dados binários, as quais seriam compostas pelas versões double (assim como suas versões potência e reversa de potência) das distribuições de probabilidade obtidas da mistura de componentes exponencial e gama, como é o da distribuição de Lindley, considerada nesta dissertação, mas também de inúmeras distribuições introduzidas na literatura recente, como as de Akash [68], Shanker [69], Sujatha

[71], Ishita [70], etc. Também é de interesse explorar métodos de *oversampling* ou *undersampling* combinados com os modelos de regressão binária aqui propostos, bem como comparar o desempenho dos novos modelos com o de modelos de *Machine Learning*.

Além disso, seria importante buscar (e analisar) outros conjuntos de dados reais, oriundos de diferentes áreas do conhecimento, a fim de demonstrar maior ganho de performance dos modelos aqui propostos sobre os usuais. Também se pretende comparar os novos modelos com os modelos de regressão binária construídos com base em outras distribuições de potência e reversa de potência, como por exemplo, a potência logística, a potência normal, a potência Cauchy, a reversa de potência logística, a reversa de potência normal, e a reversa de potência Cauchy, que foram propostos por Bazán, Torres-Avilés, Suzuki & Louzada [64].

Referências Bibliográficas

- [1] Daniel Adam. *Les réactions du consommateur devant le prix: contribution aux études de comportement*, volume 15. Sedes, 1958.
- [2] John Aitchison and James AC Brown. The lognormal distribution with special reference to its uses in economics. 1957.
- [3] Hirotugu Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6):716–723, 1974.
- [4] DK Al-Mutairi, ME Ghitany, and Debasis Kundu. Inferences on stress-strength reliability from lindley distributions. *Communications in Statistics-Theory and Methods*, 42(8):1443–1463, 2013.
- [5] Aida Ali, Siti Mariyam Shamsuddin, Anca L Ralescu, et al. Classification with class imbalance problem: a review. *International Journal Of Advances In Soft Computing And Its Applications*, 7(3):176–204, 2015.
- [6] D. R. Anderson, K. P. Burnham, and G. C. White. Comparison of akaike information criterion and consistent akaike information criterion for model selection and statistical inference from capture-recapture studies. *Journal of Applied Statistics*, 25(2):263–282, 1998.
- [7] Susan Alicia Chumbimune Anyosa. *Regressão binária usando ligações potência e reversa de potência*. PhD thesis, Universidade de São Paulo, 2017.
- [8] Samir K Ashour and Mahmoud A Eltehiwy. Exponentiated power lindley distribution. *Journal of Advanced Research*, 6(6):895–905, 2015.
- [9] Abraham Ayebo and Tomasz J Kozubowski. An asymmetric generalization of gaussian and laplace laws. *Journal of Probability and Statistical Science*, 1(2):187–210, 2003.
- [10] Adelchi Azzalini. A class of distributions which includes the normal ones. *Scandinavian Journal of Statistics*, pages 171–178, 1985.

- [11] Hassan S Bakouch, Bander M Al-Zahrani, Ali A Al-Shomrani, Vitor AA Marchi, and Francisco Louzada. An extended lindley distribution. *Journal of the Korean Statistical Society*, 41:75–85, 2012.
- [12] Narayanaswamy Balakrishnan. *Handbook of the Logistic Distribution*. CRC Press, 1991.
- [13] Pierre Baldi, Søren Brunak, Yves Chauvin, Claus AF Andersen, and Henrik Nielsen. Assessing the accuracy of prediction algorithms for classification: an overview. *Bioinformatics*, 16(5):412–424, 2000.
- [14] Jorge Bazán, F. Torres-Avilés, Adriano Suzuki, and Francisco Louzada. Power and reversal power links for binary regressions: An application for motor insurance policyholders: J.I. bazán et al. *Applied Stochastic Models in Business and Industry*, 33, 11 2016.
- [15] Chester I Bliss. The method of probits. *Science*, 1934.
- [16] Sabri Boughorbel, Fethi Jarray, and Mohammed El-Anbari. Optimal classifier for imbalanced data using matthews correlation coefficient metric. *PloS One*, 12(6):e0177678, 2017.
- [17] Hamparsum Bozdogan. Model selection and akaike’s information criterion (aic): The general theory and its analytical extensions. *Psychometrika*, 52:345–370, 02 1987.
- [18] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [19] Leo Breiman et al. Arcing classifier (with discussion and a rejoinder by the author). *The annals of statistics*, 26(3):801–849, 1998.
- [20] G. Brier. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78:1–3, 1950.
- [21] Kenneth P Burnham and David R Anderson. A practical information-theoretic approach. *Model selection and multimodel inference, 2nd ed.* Springer, New York, 2, 2002.
- [22] Guigó R. Buset, M. Evaluation of gene structure prediction programs. *Genomics*, 34:353–367, 1996.
- [23] George Casella and Roger L Berger. *Statistical Inference*, volume 2. Duxbury Pacific Grove, CA, 2002.

- [24] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16:321–357, 2002.
- [25] Ming-Hui Chen, Dipak K Dey, and Qi-Man Shao. A new skewed link model for dichotomous quantal response data. *Journal of the American Statistical Association*, 94(448):1172–1186, 1999.
- [26] David Collett. *Modelling Binary Data*. CRC press, 2 edition, 2002.
- [27] Gauss Moutinho Cordeiro and Clarice GB Demétrio. Modelos lineares generalizados e extensões. *Piracicaba: USP*, 2008.
- [28] David Roxbee Cox and E Joyce Snell. *Analysis of Binary Data*, volume 32. CRC press, 1989.
- [29] Claudia Czado and Thomas J Santner. The effect of link misspecification on binary regression inference. *Journal of Statistical Planning and Inference*, 33(2):213–231, 1992.
- [30] Cleyton de Oliveira Ritta, Marcelo Christiano Gorla, and Nelso Hein. Modelo de regressão logística para análise de risco de crédito em uma instituição de microcrédito produtivo orientado. *Iberoamerican Journal of Industrial Engineering*, 7(13):103–122, 2015.
- [31] Fernando Dellier Antunes de Souza. Aplicações de deep learning em problemas de machine learning. trabalho de conclusão de curso do programa de mba em ciência de dados (icmc), 2020.
- [32] I Elbatal and M Elgarhy. Statistical properties of kumaraswamy quasi lindley distribution. *International Journal of Mathematics Trends and Technology*, 4(10):237–246, 2013.
- [33] Michael J Farrell. The demand for motor-cars in the united states. *Journal of the Royal Statistical Society. Series A (General)*, 117(2):171–201, 1954.
- [34] Carmen Fernández and Mark FJ Steel. On bayesian modeling of fat tails and skewness. *Journal of the American Statistical Association*, 93(441):359–371, 1998.
- [35] Paulo H Ferreira, Francisco Louzada, and Carlos Diniz. Credit scoring modeling with state-dependent sample selection: A comparison study with the usual logistic modeling. *Pesquisa Operacional*, 35:39–56, 2015.

- [36] Andrew Gelman, Jessica Hwang, and Aki Vehtari. Understanding predictive information criteria for bayesian models. *Statistics and Computing*, 24(6):997–1016, 2014.
- [37] Mohamed E Ghitany, Barbra Atieh, and Saralees Nadarajah. Lindley distribution and its application. *Mathematics and Computers in Simulation*, 78(4):493–506, 2008.
- [38] Emilio Gómez-Déniz, Miguel A Sordo, and Enrique Calderín-Ojeda. The log–lindley distribution as an alternative to the beta regression model with applications in insurance. *Insurance: Mathematics and Economics*, 54:49–57, 2014.
- [39] James A Hanley and Barbara J McNeil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 143(1):29–36, 1982.
- [40] E. J. Hannan and Barry G. Quinn. The determination of the order of an autoregression. *Journal of the Royal Statistical Society. Series B: Statistical Methodology*, 41(2):190–195, 1979.
- [41] Haibo He, Yang Bai, Edwardo A Garcia, and Shutao Li. Adasyn: Adaptive synthetic sampling approach for imbalanced learning. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, pages 1322–1328. IEEE, 2008.
- [42] Arne Henningsen and Ott Toomet. maxlik: A package for maximum likelihood estimation in r. *Computational Statistics*, 26(3):443–458, 2011.
- [43] David W Hosmer Jr, Stanley Lemeshow, and Rodney X Sturdivant. *Applied Logistic Regression*, volume 398. John Wiley & Sons, 2013.
- [44] Alex de La Cruz Huayanay et al. Modelos de regressão para resposta binária na presença de dados desbalanceados. 2019.
- [45] C. Hurvich and Chih-Ling Tsai. Regression and time series model selection in small samples. *Biometrika*, 76:297–307, 1989.
- [46] Tasadduq Imam, Kai Ming Ting, and Joarder Kamruzzaman. z-svm: An svm for improved classification of imbalanced data. In *Australasian Joint Conference on Artificial Intelligence*, pages 264–273. Springer, 2006.
- [47] Christopher H. Jackson. Multi-state models for panel data: The msm package for R. *Journal of Statistical Software*, 38(8):1–29, 2011.
- [48] Olga Julia and Josep Vives-Rego. Skew-laplace distribution in gram-negative bacterial axenic cultures: new insights into intrinsic cellular heterogeneity. *Microbiology*, 151(3):749–755, 2005.

- [49] Samuel Kotz, Tomasz Kozubowski, and Krzysztof Podgorski. *The Laplace distribution and generalizations: a revisit with applications to communications, economics, engineering, and finance*. Springer Science & Business Media, 2012.
- [50] Tomasz J Kozubowski and Krzysztof Podgórski. Asymmetric laplace laws and modeling financial data. *Mathematical and Computer Modelling*, 34(9-11):1003–1021, 2001.
- [51] Nitha KU and SD Krishnarani. A new family of heavy tailed symmetric distribution for modeling financial data. *Journal of Statistics Applications & Probability*, 6(3):577–586, 2017.
- [52] Artur J Lemonte and Jorge L Bazán. New links for binary regression: an application to coca cultivation in peru. *Test*, 27(3):597–617, 2018.
- [53] Dennis V Lindley. Fiducial distributions and bayes’ theorem. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 102–107, 1958.
- [54] DV Lindley. Introduction to probability and statistics from a bayesian viewpoint, part ii: Inference. *Camb. Univ. Press, New York*, 1965.
- [55] Francisco Louzada and Anderson Ara. Bagging k-dependence probabilistic networks: An alternative powerful fraud detection tool. *Expert Systems with Applications*, 39(14):11583–11592, 2012.
- [56] Francisco Louzada, Paulo H. Ferreira-Silva, and Carlos A.R. Diniz. On the impact of disproportional samples in credit scoring models: An application to a brazilian bank data. *Expert Systems with Applications*, 39(9):8071–8078, 2012.
- [57] Mavuto Mukaka. Statistics corner: A guide to appropriate use of correlation coefficient in medical research. *Malawi medical journal : the journal of Medical Association of Malawi*, 24:69–71, 09 2012.
- [58] Saralees Nadarajah, Hassan S Bakouch, and Rasool Tahmasbi. A generalized lindley distribution. *Sankhya B*, 73(2):331–359, 2011.
- [59] Sihem Nedjar and Halim Zeghdoudi. On gamma lindley distribution: Properties and simulations. *Journal of Computational and Applied Mathematics*, 298:167–174, 2016.
- [60] Antônio Carlos Pacagnella Júnior, Geciane Silveira Porto, Sérgio Kannebley Júnior, Sérgio Luís da Silva, and Carlos Alberto Grespan Bonacim. Obtenção de patentes na indústria do estado de são paulo: uma análise utilizando regressão logística. *Production*, 19(2):261–273, 2009.

- [61] Saswat Padhi, Todd Millstein, Aditya Nori, and Rahul Sharma. Overfitting in synthesis: Theory and practice. In *International Conference on Computer Aided Verification*, pages 315–334. Springer, 2019.
- [62] David Martin Powers. Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation. 2011.
- [63] Xavier Robin, Natacha Turck, Alexandre Hainard, Natalia Tiberti, Frédérique Lisacek, Jean-Charles Sanchez, and Markus Müller. proc: an open-source package for r and s+ to analyze and compare roc curves. *BMC Bioinformatics*, 12(1):1–8, 2011.
- [64] Jose Romeo, Jorge Bazán, and Josemar Rodrigues. Bayesian skew-probit regression for binary response data. *Brazilian Journal of Probability and Statistics*, 28:467–482, 06 2014.
- [65] Takaya Saito and Marc Rehmsmeier. The precision-recall plot is more informative than the roc plot when evaluating binary classifiers on imbalanced datasets. *PloS one*, 10:e0118432, 03 2015.
- [66] C Satheesh Kumar and Rosmi Jose. On double lindley distribution and some of its properties. *American Journal of Mathematical and Management Sciences*, 38(1):23–43, 2019.
- [67] Gideon Schwarz et al. Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464, 1978.
- [68] R Shanker. Akash distribution and its applications. *International Journal of Probability and Statistics*, 4(3):65–75, 2015.
- [69] R Shanker. Shanker distribution and its applications. *International Journal of Statistics and Applications*, 5(6):338–348, 2015.
- [70] R Shanker and KK Shukla. Ishita distribution and its applications. *Biometrics & Biostatistics International Journal*, 5(2):1–9, 2017.
- [71] Rama Shanker. Sujatha distribution and its applications. *Statistics in Transition. New Series*, 17(3):391–410, 2016.
- [72] Jaqueline Trentino Silva, Ediney Magalhães Junior, Mariano Martínez Espinosa, and Dayane de Carvalho Rodrigues. Domínios da qualidade de vida associados á percepção de saúde em idosos: comparação do modelo de regressão logística com o de regressão de poisson. *Sigmae*, 8(2):576–583, 2019.

- [73] S.; BAZÁN J. L. SILVA, A. N.; ANYOSA. Bayesian binary regression modeling for unbalanced data using new links. *Rev. Bras. Biom., Lavras*, 38:385–417, 2020.
- [74] N. Sugiura. Further analysts of the data by akaike' s information criterion and the finite corrections. *Communications in Statistics-theory and Methods*, 7:13–26, 1978.
- [75] X.-M Tao, Z.-J Tong, Yan Liu, and D.-D Fu. Svm classifier for unbalanced data based on combination of odr and bsmote. *Kongzhi yu Juece/Control and Decision*.
- [76] R Core Team. R: A language and environment for statistical computing, version 3.6. 3. vienna, austria; 2013.
- [77] H Torabi, M Falahati-Naeini, and NH Montazeri. An extended generalized lindley distribution and its applications to lifetime data. *Journal of Statistical Research of Iran JSRI*, 11(2):203–222, 2015.
- [78] Bart Van der Paal. A comparison of different methods for modelling rare events data. *PhD thesis*, 2014.
- [79] Pierre-François Verhulst. Notice sur la loi que la population suit dans son accroissement. *Corresp. Math. Phys.*, 10:113–126, 1838.
- [80] Pierre-François Verhulst. Resherches mathematiques sur la loi d'accroissement de la population. *Nouveaux memoires de l'academie royale des sciences*, 18:1–41, 1845.
- [81] Pierre-François Verhulst. Deuxième mémoire sur la loi d'accroissement de la population. *Mémoires de l'académie royale des sciences, des lettres et des beaux-arts de Belgique*, 20:1–32, 1847.