

UNIVERSIDADE FEDERAL DE ALFENAS

Luiz Otávio de Oliveira Pala

Revisitando a estimação do Coeficiente de Determinação

ALFENAS, MG

2019

LUIZ OTÁVIO DE OLIVEIRA PALA

Revisitando a estimação do Coeficiente de Determinação

Dissertação apresentada ao Programa de Pós-Graduação em Estatística Aplicada e Biometria, Área de concentração em Estatística Aplicada e Biometria da Universidade Federal de Alfenas, MG, como parte dos requisitos para a obtenção do título de Mestre.

Linha de Pesquisa: Estatística Aplicada e Biometria.

Orientador: Prof. Dr. Eric Batista Ferreira.

Coorientador: Prof. Dr. Davi Butturi-Gomes.

ALFENAS, MG

2019

Dados Internacionais de Catalogação-na-Publicação (CIP)
Sistema de Bibliotecas da Universidade Federal de Alfenas

P153r Pala, Luiz Otávio de Oliveira.
Revisitando a estimação de coeficiente de determinação / Luiz Otávio de Oliveira Pala -- Alfenas/MG, 2019.
112 f.: il. –

Orientador: Eric Batista Ferreira.
Dissertação (Mestrado em Estatística Aplicada e Biometria) -
Universidade Federal de Alfenas, 2019.
Bibliografia.

1. Intervalos de Confiança. 2. Estatística como Assunto. 3. Método de Monte Carlo. I. Ferreira, Eric Batista. II. Título.

CDD-519.54



LUIZ OTÁVIO DE OLIVEIRA PALA

“REVISITANDO A ESTIMAÇÃO DO COEFICIENTE DE DETERMINAÇÃO”

A Banca Examinadora, abaixo assinada, aprova a Dissertação apresentada como parte dos requisitos para a obtenção do título de Mestre em Estatística Aplicada e Biometria pela Universidade Federal de Alfenas. Área de Concentração: Estatística Aplicada e Biometria

Aprovado em: 10 de julho de 2019.

Prof. Dr. Eric Batista Ferreira

Instituição: UNIFAL-MG

Assinatura: Eric Batista Ferreira

Profa. Dra. Gislene Araújo Pereira

Instituição: UNIFAL-MG

Assinatura: Gislene Araújo Pereira

Prof. Dr. Washington Santos da Silva

Instituição: IFMG

Assinatura: Washington Santos da Silva

Prof. Dr. Marcelino Alves Rosa de Pascoa

Instituição: ESALQ/USP

Assinatura: Marcelino Alves Rosa de Pascoa

AGRADECIMENTOS

Inicio os meus agradecimentos aos meus pais Maria e Luiz e ao meu irmão Mateus, que sempre me apoiaram e me incentivaram nos estudos. Agradeço também à Priscila e ao Daniel, este último que chegou em nossa família em um momento especial, trazendo muita alegria.

Aos meus amigos Aline, Ana Flávia, Lara, Tatiane, Thaisa e Rodolfo pela amizade e apoio de cada um durante o mestrado. Amo vocês.

Agradeço também a Universidade Federal de Alfenas e ao Programa de Pós Graduação em Estatística Aplicada e Biometria. Ao meu orientador Prof. Dr. Eric Batista Ferreira, coorientador Prof. Dr. Davi Butturi-Gomes e a minha orientadora de graduação Prof. Dr^a. Gislene Araujo Pereira. Obrigado pela confiança, atenção e conselhos.

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

RESUMO

O coeficiente de determinação (R^2) é uma métrica muito utilizada para a análise de qualidade de ajuste de modelos lineares. Este coeficiente assume valores no intervalo entre 0 e 1, de modo que quanto mais próximo de 1, maior parte da variação da variável resposta está sendo explicada pelo modelo. Há outras métricas com o mesmo objetivo, como o coeficiente de determinação ajustado, o erro absoluto e erro quadrático médio, por exemplo. Mesmo sendo muito utilizado, o R^2 é tratado com cautela na literatura, pois este pode ser viesado em modelos com poucas observações ou inflacionado quando se acrescentam covariáveis ao modelo. Neste sentido, autores sugerem tratá-lo como uma estatística que estima um parâmetro populacional (ρ^2), sendo este entendido como a qualidade de ajuste que um modelo possuiria se as infinitas observações do fenômeno viessem a ser coletadas. Desta forma, sendo ρ^2 um parâmetro e R^2 um estimador pontual, é natural pensar em estimação intervalar e testes de hipóteses para possibilitar a tomada de decisão sobre a adequação do modelo candidato ao fenômeno no qual deseja-se descrever. Entretanto, essa questão inferencial ainda não é considerada fechada na literatura, pois autores discutem distribuições de probabilidade para a modelagem deste em diferentes cenários e regiões do espaço paramétrico. Desta forma, este trabalho estudou a estimação do coeficiente de determinação paramétrico (ρ^2) a partir de cinco estimadores intervalares paramétricos. Para compará-los, foi realizado um estudo de simulação Monte Carlo, computando precisão e acurácia em diferentes cenários compostos pela combinação do número de covariáveis do modelo (k), tamanho amostral (n) e o verdadeiro valor paramétrico (ρ^2). Em conjunto a isso, elaborou-se um índice de desempenho de estimação intervalar que valoriza simultaneamente precisão e acurácia com importâncias relativas previamente fixadas. Os resultados permitiram a recomendação do melhor estimador para cada região do espaço paramétrico. Com isso, verificou-se que os estimadores propostos apresentaram qualidade similar aos indicados na literatura ao longo do espaço paramétrico. Por fim, foi construído um pacote R, possibilitando que o usuário estime de forma intervalar a qualidade do ajuste utilizando o estimador com melhor desempenho.

Palavras-Chave: Estatística como Assunto. Intervalos de Confiança. Método de Monte Carlo.

ABSTRACT

The coefficient of determination (R^2) is a widely used metric to analyze the quality of adjustment of linear models. This coefficient assumes values in the range between 0 and 1, so that the closer to 1, most of the capacity variation is being explained by the model. R^2 is treated with caution in the literature, as it can be biased in models with few observations or inflated when covariates are added to the model. In this sense, authors suggest treating it as a statistic that estimates a population parameter (ρ^2), which is understood as the quality of fit that a model would have if the infinite observations of the phenomenon were to be collected. Thus, we study the estimation of the parametric coefficient of determination (ρ^2) from five parametric interval estimators. For comparison, a Monte Carlo simulation study was performed, computing precision and accuracy in the different combinations to model the number of model variables (k), sample size (n) and the value of the parameter (ρ^2). The results allowed the recommendation of the best estimator for each region of the parametric space. Thus, it was found that the proposed estimators presented similar quality to those indicated in the literature in the parametric space. Finally, an R package was built, allowing the user to estimate the quality of the model using the best performing estimator.

Keywords: Statistics as a subject. Confidence intervals. Monte Carlo method.

LISTA DE FIGURAS

Figura 1 –	Comportamento da distribuição de probabilidade Beta para diferentes combinações de parâmetros a e b	18
Figura 2 –	Algumas relações e transformações entre distribuições de probabilidade associadas à distribuição Beta	20
Figura 3 –	Comportamento esperado de estimadores considerando o nível de precisão e acurácia	25
Figura 4 –	Comportamento da distribuição de M_1 considerando $n = 20$, $k = 1$ e $\rho^2 \in \{0,1; 0,5; 0,9\}$	36
Figura 5 –	Comportamento da distribuição de M_2 considerando $n = 20$, $k = 1$ e $\rho^2 \in \{0,1; 0,5; 0,9\}$	37
Figura 6 –	Comportamento da distribuição de P^* considerando $n = 20$ e $\rho^2 \in \{0,1; 0,5; 0,9\}$	38
Figura 7 –	Histogramas das amostras de R^2 para diferentes combinações de n . Sendo (a): $n = 15$, (b): $n = 50$ e (c): $n = 100$	40
Figura 8 –	Diagrama com os principais estudos que analisaram distribuições para o coeficiente de determinação ou que avaliaram computacionalmente intervalos de confiança	42
Figura 9 –	Taxa de acurácia referente ao modelo de regressão onde $k = 1$ e $n = 15$	45
Figura 10 –	Taxa de acurácia referente ao modelo de regressão onde $k = 1$ e $n = 50$	45
Figura 11 –	Taxa de acurácia referente ao modelo de regressão onde $k = 1$ e $n = 100$	46
Figura 12 –	Taxa de acurácia referente ao modelo de regressão onde $k = 8$ e $n = 15$	47
Figura 13 –	Taxa de acurácia referente ao modelo de regressão onde $k = 8$ e $n = 50$	47
Figura 14 –	Taxa de acurácia referente ao modelo de regressão onde $k = 8$ e $n = 100$	48
Figura 15 –	Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 15$	49
Figura 16 –	Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 50$	50
Figura 17 –	Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 100$	50

Figura 18 – Precisão dos intervalos referente ao modelo de regressão onde $k = 8$ e $n = 15$	51
Figura 19 – Precisão dos intervalos referente ao modelo de regressão onde $k = 8$ e $n = 50$	52
Figura 20 – Precisão dos intervalos referente ao modelo de regressão onde $k = 8$ e $n = 100$	52
Figura 21 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 1$ e $n = 15$	54
Figura 22 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 1$ e $n = 50$	54
Figura 23 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 1$ e $n = 100$	55
Figura 24 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 8$ e $n = 15$	56
Figura 25 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 8$ e $n = 50$	56
Figura 26 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 8$ e $n = 100$	57
Figura 27 – Dias transcorridos até o alcance máximo de amadurecimento em relação ao período de exposição ao 1-MCP	58
Figura 28 – Modelos ajustados ao experimento relativo ao amadurecimento de bananas do tipo maçã	59
Figura 29 – Comportamento entre a idade das matrizes e o ganho de peso da progênie no período de 1 a 7 dias após o nascimento	69
Figura 30 – Diagrama da função que recomenda o estimador a ser utilizado a partir de uma suavização dos cenários simulados, considerando o maior índice $\tau_{\alpha,i}$	73
Figura 31 – Estimador recomendado pela função sugestão disponibilizada no pacote ICR2 para o modelo linear	76
Figura 32 – Estimador recomendado pela função sugestão disponibilizada no pacote ICR2 para o modelo quadrático	76

Figura 33 – Estimador recomendado pela função simulação disponibilizada no pacote ICR2 para o modelo linear	77
Figura 34 – Modelos ajustados ao experimento relativo ao ganho de peso de filhotes em função da idade das matrizes	78
Figura 35 – Taxa de acurácia referente ao modelo de regressão onde $k = 2$ e $n = 15$	84
Figura 36 – Taxa de acurácia referente ao modelo de regressão onde $k = 2$ e $n = 50$	85
Figura 37 – Taxa de acurácia referente ao modelo de regressão onde $k = 2$ e $n = 100$	85
Figura 38 – Taxa de acurácia referente ao modelo de regressão onde $k = 3$ e $n = 15$	86
Figura 39 – Taxa de acurácia referente ao modelo de regressão onde $k = 3$ e $n = 50$	86
Figura 40 – Taxa de acurácia referente ao modelo de regressão onde $k = 3$ e $n = 100$	87
Figura 41 – Taxa de acurácia referente ao modelo de regressão onde $k = 4$ e $n = 15$	87
Figura 42 – Taxa de acurácia referente ao modelo de regressão onde $k = 4$ e $n = 50$	88
Figura 43 – Taxa de acurácia referente ao modelo de regressão onde $k = 4$ e $n = 100$	88
Figura 44 – Taxa de acurácia referente ao modelo de regressão onde $k = 5$ e $n = 15$	89
Figura 45 – Taxa de acurácia referente ao modelo de regressão onde $k = 5$ e $n = 50$	89
Figura 46 – Taxa de acurácia referente ao modelo de regressão onde $k = 5$ e $n = 100$	90
Figura 47 – Taxa de acurácia referente ao modelo de regressão onde $k = 6$ e $n = 15$	90
Figura 48 – Taxa de acurácia referente ao modelo de regressão onde $k = 6$ e $n = 50$	91
Figura 49 – Taxa de acurácia referente ao modelo de regressão onde $k = 6$ e $n = 100$	91
Figura 50 – Taxa de acurácia referente ao modelo de regressão onde $k = 7$ e $n = 15$	92
Figura 51 – Taxa de acurácia referente ao modelo de regressão onde $k = 7$ e $n = 50$	92
Figura 52 – Taxa de acurácia referente ao modelo de regressão onde $k = 7$ e $n = 100$	93
Figura 53 – Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 15$	94
Figura 54 – Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 50$	94
Figura 55 – Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 100$	95
Figura 56 – Precisão dos intervalos referente ao modelo de regressão onde $k = 3$ e $n = 15$	95
Figura 57 – Precisão dos intervalos referente ao modelo de regressão onde $k = 3$ e $n = 50$	96

Figura 58 – Precisão dos intervalos referente ao modelo de regressão onde $k = 3$ e $n = 100$	96
Figura 59 – Precisão dos intervalos referente ao modelo de regressão onde $k = 4$ e $n = 15$	97
Figura 60 – Precisão dos intervalos referente ao modelo de regressão onde $k = 4$ e $n = 50$	97
Figura 61 – Precisão dos intervalos referente ao modelo de regressão onde $k = 4$ e $n = 100$	98
Figura 62 – Precisão dos intervalos referente ao modelo de regressão onde $k = 5$ e $n = 15$	98
Figura 63 – Precisão dos intervalos referente ao modelo de regressão onde $k = 5$ e $n = 50$	99
Figura 64 – Precisão dos intervalos referente ao modelo de regressão onde $k = 5$ e $n = 100$	99
Figura 65 – Precisão dos intervalos referente ao modelo de regressão onde $k = 6$ e $n = 15$	100
Figura 66 – Precisão dos intervalos referente ao modelo de regressão onde $k = 6$ e $n = 50$	100
Figura 67 – Precisão dos intervalos referente ao modelo de regressão onde $k = 6$ e $n = 100$	101
Figura 68 – Precisão dos intervalos referente ao modelo de regressão onde $k = 7$ e $n = 15$	101
Figura 69 – Precisão dos intervalos referente ao modelo de regressão onde $k = 7$ e $n = 50$	102
Figura 70 – Precisão dos intervalos referente ao modelo de regressão onde $k = 7$ e $n = 100$	102
Figura 71 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 2$ e $n = 15$	103
Figura 72 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 2$ e $n = 50$	104
Figura 73 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 2$ e $n = 100$	104

Figura 74 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 3$ e $n = 15$	105
Figura 75 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 3$ e $n = 50$	105
Figura 76 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 3$ e $n = 100$	106
Figura 77 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 4$ e $n = 15$	106
Figura 78 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 4$ e $n = 50$	107
Figura 79 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 4$ e $n = 100$	107
Figura 80 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 5$ e $n = 15$	108
Figura 81 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 5$ e $n = 50$	108
Figura 82 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 5$ e $n = 100$	109
Figura 83 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 6$ e $n = 15$	109
Figura 84 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 6$ e $n = 50$	110
Figura 85 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 6$ e $n = 100$	110
Figura 86 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 7$ e $n = 15$	111
Figura 87 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 7$ e $n = 50$	111
Figura 88 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 7$ e $n = 100$	112

SUMÁRIO

	CAPÍTULO 1 - INTRODUÇÃO GERAL	15
1	REFERENCIAL TEÓRICO	16
1.1	O MODELO DE REGRESSÃO LINEAR E O COEFICIENTE DE DETERMINAÇÃO	16
1.2	A DISTRIBUIÇÃO DE PROBABILIDADE BETA	18
1.3	DISTRIBUIÇÕES DO COEFICIENTE DE DETERMINAÇÃO	20
1.4	ESTUDOS DE SIMULAÇÃO DE MONTE CARLO	23
1.5	AVALIAÇÃO DE ESTIMADORES INTERVALARES	24
	REFERÊNCIAS	29
	CAPÍTULO 2 - ESTIMAÇÃO INTERVALAR DO COEFICIENTE DE DETERMINAÇÃO: UMA APLICAÇÃO NA QUALIDADE DE AJUSTE DE MODELOS LINEARES	31
	RESUMO	31
1	INTRODUÇÃO	32
2	MATERIAL E MÉTODOS	34
2.1	ESTIMADORES INTERVALARES UTILIZADOS E PROPOSTOS	34
2.1.1	Estimador M_1	35
2.1.2	Estimador M_2	36
2.1.3	Estimador P^*	37
2.2	GERAÇÃO DOS DADOS E CENÁRIOS AVALIADOS	39
2.3	ÍNDICE DE DESEMPENHO DE ESTIMAÇÃO INTERVALAR (τ_α)	41
2.4	APLICAÇÃO: DURABILIDADE PÓS COLHEITA DE BANANAS DO TIPO MAÇÃ	42
3	RESULTADOS E DISCUSSÃO	44
3.1	AVALIAÇÃO E COMPARAÇÃO DOS ESTIMADORES	44
3.1.1	Taxas de acurácia	44
3.1.2	Níveis de precisão	49
3.1.3	Índices de desempenho de estimação intervalar	53
4	CONCLUSÃO	61
	REFERÊNCIAS	63
	CAPÍTULO 3 - O PACOTE ICR2: INTERVALOS DE CONFIANÇA PARA O COEFICIENTE DE DETERMINAÇÃO	65
	RESUMO	65
1	INTRODUÇÃO	66
2	EXEMPLO	68
3	IMPLEMENTAÇÃO	69
3.1	DISTRIBUIÇÕES DOS ESTIMADORES	69
3.1.1	Distribuição de W	69
3.1.2	Distribuição de E^*	70
3.1.3	Distribuição de M_1	71
3.1.4	Distribuição de M_2	71
3.1.5	Distribuição de P^*	72
3.2	FUNÇÕES PARA A CONSTRUÇÃO DE INTERVALOS DE CONFIANÇA	72
4	APLICAÇÃO AO EXEMPLO	75

5	CONCLUSÃO	79
	REFERÊNCIAS	81
	CAPÍTULO 4 - CONSIDERAÇÕES FINAIS	83
	APÊNDICES	84

CAPÍTULO 1 - INTRODUÇÃO GERAL

Com o desenvolver da ciência, os modelos teóricos são criados ou reformulados com o objetivo de representar com maior precisão os fenômenos ou adequar-se aos contextos da atualidade. Modelos que estão presentes em diversas áreas do conhecimento e possibilitam o levantamento de hipóteses e suas respectivas respostas acerca dos problemas de estudos de pesquisadores (SHOU et al., 2015).

No campo científico, autores como Gilbert e Osborne (1980), Montgomery, Peck e Vining (2006) e, mais recentemente, Frigg e Nguyen (2017) discutiram a importância e aplicabilidade de modelos. Esses autores exemplificaram os modelos atômicos, econométricos, climáticos e estatísticos. Sendo estes últimos apontados por Montgomery, Peck e Vining (2006) como os mais utilizados na ciência, contribuindo para delinear o relacionamento entre variáveis.

O método de regressão tem como propósito determinar a condição de causa e efeito ou o impacto de uma variável sobre as outras a partir de uma relação funcional estabelecida, possibilitando a realização de previsões da variável resposta (UYANIK; GÜLER, 2013). Esses modelos funcionais podem ser lineares ou não e são aproximações do verdadeiro relacionamento entre as variáveis analisadas pelo pesquisador.

Entretanto, dada as diversas relações funcionais que podem ser estabelecidas, métricas de qualidade de ajuste foram propostas na literatura para analisar o nível de adequabilidade do modelo ao fenômeno. Como exemplo dessas métricas, tem-se o Coeficiente de Determinação (R^2) e o Coeficiente de Determinação Ajustado (R_a^2), que são usuais nos modelos de regressão lineares, além do Erro Médio (EM), Erro Quadrático Médio (EQM) e Erro Médio Absoluto (EMA).

A avaliação do ajuste possibilita um dos pontos importantes da análise de regressão: a seleção de modelos; contexto tratado por Hamid et al. (2018) como um dos principais pontos da modelagem estatística. Neste sentido, as métricas de qualidade unem-se aos critérios BIC, AIC e a estatística de *Mallows*, permitindo a escolha de modelos parcimoniosos e com maiores taxas de acurácia, por exemplo.

No entanto, o R^2 é tratado por autores como uma métrica que deve ser analisada com certa precaução. Conforme Di Bucchianico (2008), este coeficiente está diretamente relacionado ao modelo estimado, entretanto, valores próximos a um podem não indicar adequabilidade do modelo ao fenômeno. Além de Di Bucchianico (2008), Montgomery, Peck e Vining

(2006) afirmaram que um valor alto de R^2 pode ser ilusório, dado que este pode ser inflacionado acrescentando covariáveis ao modelo.

Dada essa situação, autores sugerem analisar o coeficiente como uma estatística que estima a qualidade de ajuste de uma regressão, o que pode ser visto em Cramer (1987), Carroodus e Giles (1992) e Quinino, Reis e Bessegato (2012). Assim, é possível inferir a respeito de um coeficiente de determinação populacional (ρ^2) e tomar decisões sobre a qualidade do modelo.

Entretanto, nota-se que a estimação de ρ^2 não é um assunto esgotado na literatura. Estudos apresentam e discutem distribuições para os estimadores de ρ^2 sob diversas parametrizações e aspectos, como amostras pequenas e perturbações nos termos de erro da regressão. Em consequência deste fato, pode-se notar a presença de diferentes alternativas para a realização de inferências no espaço paramétrico de ρ^2 , dando abertura para novas propostas inferenciais.

Neste sentido, o objetivo deste trabalho é seccionado em dois artigos, dispostos nos Capítulos 2 e 3, sendo eles:

a) Capítulo 2: Estudar empiricamente a estimação intervalar de ρ^2 por meio cinco estimadores paramétricos, sendo três deles propostos por este trabalho. Além disso, comparar medidas de desempenho dos estimadores nos cenários analisados e construir intervalos de confiança para a qualidade de ajuste de um experimento, ilustrando a aplicação dos estimadores;

b) Capítulo 3: Construir um pacote em linguagem R, possibilitando a construção de intervalos de confiança para ρ^2 em cenários específicos do usuário, utilizando do estimador que apresente melhor qualidade e possibilite a seleção de modelos.

1 REFERENCIAL TEÓRICO

O objetivo desta seção é apresentar o coeficiente de determinação e as respectivas distribuições propostas por estudos para a estatística. Por fim, apresentar os aspectos da simulação computacional que possibilitam avaliar estimadores intervalares e testes de hipóteses.

1.1 O MODELO DE REGRESSÃO LINEAR E O COEFICIENTE DE DETERMINAÇÃO

Conforme Su, Yan e Tsai (2012), a metodologia de regressão desempenha um papel de grande importância na modelagem estatística. Em particular, os modelos lineares tornaram-se populares pela clareza quanto a forma de interpretação de seus parâmetros, unindo-se ao detalhamento teórico bem fundamentado na literatura.

De modo geral, o modelo de regressão linear pode ser denotado matricialmente da seguinte forma:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{I}),$$

sendo:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_{12} & \dots & x_{1p} \\ 1 & x_{22} & \dots & x_{2p} \\ 1 & x_{32} & \dots & x_{3p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n2} & \dots & x_{np} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}, \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{bmatrix}.$$

Sendo \mathbf{y} o vetor com a variável resposta; \mathbf{X} a matriz de delineamento e $\boldsymbol{\epsilon}$ o vetor de erros. Sob pressuposições de que a variável resposta seja função linear das covariáveis; $E(\epsilon_i) = 0$; os erros são homoscedásticos, não correlacionados e normalmente distribuídos (HOFFMANN, 2016).

A estimação dos parâmetros do modelo pode ser dada de diversas formas como mínimos quadrados, máxima verossimilhança e inferência Bayesiana. O método de mínimos quadrados é o mais utilizado na literatura, correspondendo a minimizar a distância da resposta observada com os valores preditos (SU; YAN; TSAI, 2012). Conforme Hoffmann (2016), o estimador de

mínimos quadrados é dado por:

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

sendo \mathbf{b} um estimador não tendencioso e de variância mínima de β conforme o Teorema de *Gauss-Markov*.

Posteriormente, pode ser de interesse do pesquisador a realização da análise de variância do modelo e a obtenção de algumas estatísticas. Neste contexto, define-se a Soma de Quadrados Totais (SQT) do modelo como:

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

que pode ser decomposta pela Soma de Quadrados da Regressão (SQR) e Soma de Quadrados dos Erros (SQE), respectivamente:

$$SQT = (\mathbf{X}\beta - \bar{\mathbf{y}})'(\mathbf{X}\beta - \bar{\mathbf{y}}) + \mathbf{e}'\mathbf{e}.$$

Com essas quantidades, pode-se obter o coeficiente de determinação amostral R^2 do modelo pela razão entre SQR e SQT, sendo interpretado como a proporção da variação total da resposta explicada pela regressão (SU; YAN; TSAI, 2012). Conforme Cramer (1987), tal medida representa a qualidade de ajuste do modelo, dada matricialmente por:

$$R^2 = \frac{\mathbf{b}'\mathbf{X}'\mathbf{X}\mathbf{b}}{\mathbf{b}'\mathbf{X}'\mathbf{X}\mathbf{b} + \mathbf{e}'\mathbf{e}}, \quad (1.1)$$

sendo \mathbf{b} o vetor de parâmetros estimados do modelo, \mathbf{X} a matriz de delineamento e \mathbf{e} o vetor de erros associado. Tal métrica assume valores no intervalo $0 \leq R^2 \leq 1$, de modo que quanto mais próximo de 1, maior parte da variação da variável resposta está sendo explicada.

Já foi apontado que o R^2 é uma métrica que pode apresentar problemas. Um destes foi indicado por Cramer (1987), afirmando que o coeficiente é viesado em pequenas amostras. Não sendo recomendável seu cálculo em modelos com menos de cinquenta observações, tamanho amostral este que pode ser considerado amplo para algumas áreas do conhecimento. Fato que também foi apontado por Quinino, Reis e Bessegato (2012), evidenciando que o valor de R^2 pode ser 'ilusório' em pequenas amostras.

Dada as afirmações de que o R^2 pode ser viesado e levar o pesquisador a tomar decisões equivocadas em certos cenários, autores como Weatherburn (1949), Cramer (1987) e Quinino,

Reis e Bessegato (2012) discutiram distribuições probabilísticas por meio de reparametrizações da distribuição de probabilidade Beta, possibilitando a realização de inferências a respeito de ρ^2 .

1.2 A DISTRIBUIÇÃO DE PROBABILIDADE BETA

A Beta é uma distribuição de probabilidade contínua, comumente utilizada para a modelagem de proporções, sendo muito flexível conforme a indexação de seus dois parâmetros (FERRARI; CRIBARI-NETO, 2004). Essa distribuição apresenta função de densidade de probabilidade dada por:

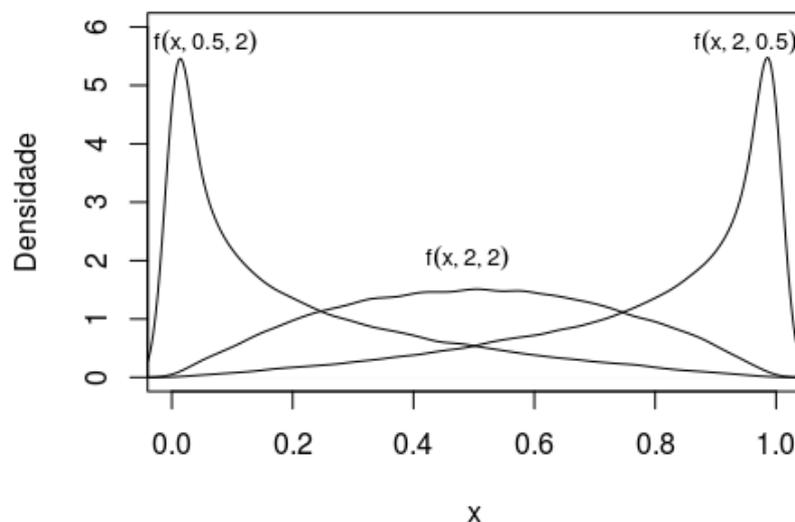
$$f(x; a, b) = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1} \mathbf{I}_{(0,1)}(x),$$

definida no espaço paramétrico $\Phi = \{(a, b) | a > 0, b > 0\}$, onde:

$$B(a, b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx,$$

corresponde a função beta. Como exemplo do comportamento da distribuição, a Figura 1 apresenta diferentes combinações dos parâmetros da distribuição.

Figura 1 – Comportamento da distribuição de probabilidade Beta para diferentes combinações de parâmetros a e b



Fonte: Do autor.

Conforme Mood, Graybill e Boes (1974), a função geradora de momentos da distribuição não tem forma trivial. Entretanto, alguns de seus momentos ou funções deles são acessíveis, como a média e a variância, dadas por:

$$E(X) = \frac{a}{(a+b)},$$

$$V(X) = \frac{ab}{(a+b+1)(a+b)^2}.$$

A existência da moda da distribuição depende do valor de seus parâmetros, de modo que sendo $a > 1$ e $b > 1$, a moda é dada por:

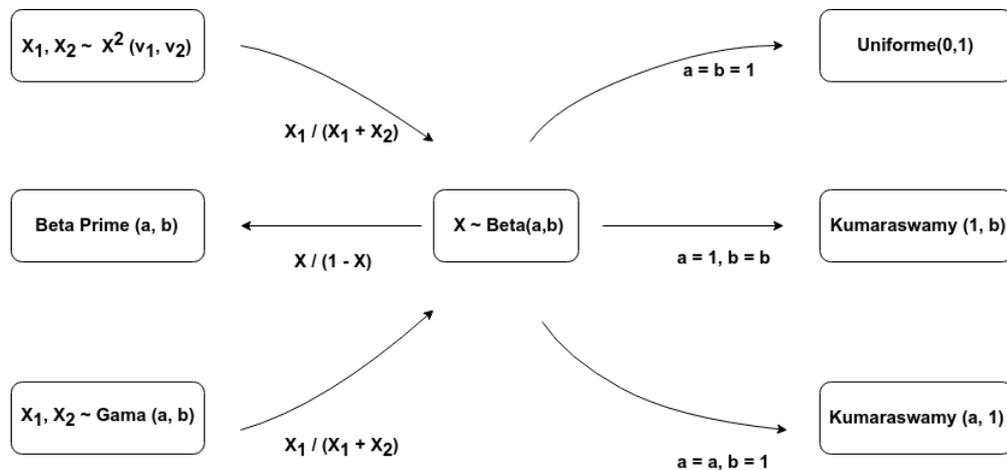
$$M_o(X) = \frac{a-1}{a+b-2},$$

e por outra perspectiva, se a ou $b < 1$, a moda da distribuição estará em uma das bordas (JOHNSON; BERVERLIN, 2013). Em oposição à moda, a mediana não apresenta uma forma algébrica fechada. O estudo de Kerman (2011) retratou uma aproximação para esta em pares a e b maiores que 1. Tal aproximação apresentou erro decrescente, ou seja, aproximando-se da mediana da distribuição a medida em que o valor dos parâmetros aumentaram.

Dado o domínio da distribuição, ela é comumente utilizada na modelagem de fenômenos que são mensurados no intervalo $(0,1)$. Como exemplo de aplicações, Ferrari e Cribari-Neto (2004) estabeleceram um modelo de regressão para fenômenos onde a variável resposta segue uma distribuição Beta, a partir de uma parametrização considerando a média e a dispersão. Quinino, Reis e Bessegato (2012), exploraram a distribuição para analisar o coeficiente de determinação de uma regressão, um tópico a ser explorado na próxima seção. Arnold e Ghosh (2017) analisaram o comportamento e algumas propriedades para a construção das distribuições Beta e Kumaraswamy bivariadas.

Além da forma mais conhecida da distribuição Beta, há variações na literatura. Uma destas corresponde a distribuição Beta *prime*, também denominada como Beta do segundo tipo, que apresenta uma cauda longa (TULUPYEV et al., 2013). Por fim, autores como Weatherburn (1946), Casella e Berger (2001) e Song (2005) discutiram algumas relações e/ou transformações a partir da distribuição Beta, como apresentado na Figura 2. É interessante notar que se X_1 e X_2 seguem distribuição χ^2 com v_1 e v_2 graus de liberdade, respectivamente, então a distribuição resultante será Beta com parâmetros $\frac{v_1}{2}$ e $\frac{v_2}{2}$ (WEATHERBURN, 1949). Ressalta-se que

Figura 2 – Algumas relações e transformações entre distribuições de probabilidade associadas à distribuição Beta



Fonte: Elaborado a partir de Weatherburn (1949), Casella e Berger (2001) e Song (2005).

algumas destas relações são usuais em contextos de inferência Bayesiana e modelos estatísticos em eventos discretos (JOHNSON, 2013).

1.3 DISTRIBUIÇÕES DO COEFICIENTE DE DETERMINAÇÃO

No ano de 1949, Weatherburn discutiu a distribuição do coeficiente de correlação (r) e do coeficiente de correlação ao quadrado. De modo que, a partir de uma transformação de variáveis obteve-se a distribuição amostral para um cenário específico do espaço paramétrico ($\rho^2 = 0$). Como resultado, mostrou-se que em pares de variáveis X e Y não correlacionadas, com respectivas médias e desvios padrão, o coeficiente de determinação segue uma distribuição Beta com parâmetros $a = \frac{1}{2}$ e $b = \frac{(n-2)}{2}$, onde n corresponde ao número de observações no modelo.

Expandindo o cenário apresentado por Weatherburn (1949), Foster, Smith e Whaley (1997) analisaram a distribuição quando k covariáveis compõem o modelo de regressão, sob a hipótese nula que o vetor de parâmetros é nulo ($\beta = 0$), implicando na condição de que $\rho^2 = 0$.

Com isso, o estimador do coeficiente seria distribuído da seguinte forma:

$$R_W^2 \sim \text{Beta} \left(\frac{k}{2}, \frac{n-k-1}{2} \right). \quad (1.2)$$

Essa distribuição amostral foi utilizada por Quinino, Reis e Bessegato (2012) para a realização de inferências em um modelo de regressão. Isso possibilitou a tomada de decisão a respeito do coeficiente de determinação populacional, conforme o seguinte par de hipóteses:

$$\begin{cases} H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0 \\ H_1 : B_i \neq 0, \text{ para algum } i = \{1, 2, \dots, k\} \end{cases} \iff \begin{cases} H_0 : \rho^2 = 0 \\ H_1 : \rho^2 > 0 \end{cases}.$$

Como resultado, tomou-se uma decisão a respeito da qualidade de ajuste em uma aplicação relacionando a velocidade de um automóvel em função de seu peso e potência.

Cramer (1987) obteve a distribuição exata de R^2 a partir de uma transformação considerando a situação de que $\beta \neq 0$, ou seja, um cenário mais geral do que o proposto por Weatherburn (1949). A transformação considerada foi dada por uma função (G) da seguinte forma:

$$G = \frac{R^2}{1 - R^2} = \frac{(b' X' X b) / \sigma^2}{(e' e) / \sigma^2},$$

sendo:

$$\frac{(b' X' X b)}{\sigma^2} \sim \chi_{(\lambda, k)}^2, \quad (1.3)$$

e,

$$\frac{e' e}{\sigma^2} \sim \chi_{(0, n-k)}^2. \quad (1.4)$$

Sendo λ o parâmetro de não centralidade da distribuição e b o vetor de parâmetros estimados. Para encontrar a distribuição dada pela transformação proposta, foram utilizadas fórmulas de recorrência como apresentado em Johson, Kotz e Balakrishnan (1995). Desta forma, a distribuição exata é dada por um produto entre as distribuições Beta e Poisson:

$$f(r) = \sum_{j=0}^{\infty} W(j) \frac{1}{B(u+j, v-u)} r^{u+j-1} (1-r)^{v-u-1}, \quad (1.5)$$

em que:

$$W(j) = \frac{e^{-\frac{\lambda}{2}} \left(\frac{\lambda}{2}\right)^j}{j!},$$

$$u = \frac{1}{2}(k),$$

e,

$$v = \frac{1}{2}(n - 1).$$

O estudo de Cramer (1987) relacionou o coeficiente de determinação amostral R^2 com ρ^2 , mostrando que R^2 possui uma condição assintótica conforme o tamanho amostral do modelo aumenta. Assim, ρ^2 seria estimado a partir dos parâmetros do modelo de regressão da seguinte forma:

$$\rho^2 = \frac{\beta' X' X \beta}{\beta' X' X \beta + n\sigma^2},$$

Desta forma, ρ^2 é incorporado à distribuição exata através de λ , dado por:

$$\lambda = n \left(\frac{\rho^2}{1 - \rho^2} \right). \quad (1.6)$$

Cramer (1987) também derivou a média e variância do coeficiente a partir da distribuição exata, indicando que, assim como R^2 converge para o seu valor populacional, o mesmo pode acontecer para a mediana e para a moda da distribuição.

Ohtani (1994) analisou a distribuição de R^2 e de R_a^2 , a partir de funções de risco em duas situações: i) quando variáveis importantes ao fenômeno são omitidas e ii) quando variáveis desprezíveis são mantidas. Como resultado, obteve-se que ambos coeficientes são subestimados na situação de omissão de variáveis relevantes e são sobrestimados na segunda situação. Ressalta-se que a distribuição utilizada por este estudo está relacionada com a distribuição proposta por Cramer (1987).

Por fim, Ohtani e Tanizaki (2004) apresentaram a distribuição exata em um modelo de regressão onde os erros seguem distribuição t multivariada. O estudo foi realizado via simulação com variações entre o número de covariáveis e o tamanho amostral. Foram avaliados computacionalmente cenários com $n \in \{10, 20, 30, 40\}$ e $k \in \{3, 4, 5, 6, 7, 8\}$. Posteriormente, mostrou-se que o R^2 é seriamente viesado em amostras pequenas e que os intervalos de confiança para cenários nas proximidades de zero são imprecisos.

1.4 ESTUDOS DE SIMULAÇÃO DE MONTE CARLO

Em estudos de simulação computacional relativos à avaliação de testes de hipóteses e construção de intervalos de confiança, comumente avaliam-se cenários com diferentes combinações de parâmetros e tamanhos amostrais. Com isso, é possível identificar o desempenho destes sob diversas perspectivas, possibilitando a indicação daqueles com maiores taxas de poder, por exemplo. Esses estudos comumente utilizam de métodos de Monte Carlo e simulação estocástica a partir da geração de números pseudoaleatórios, que podem apresentar certa distribuição de probabilidade.

A simulação de Monte Carlo tem por objetivo aproximar um valor esperado por sua média de resultados por meio da realização de repetições de um experimento, respaldado pela teoria dos grandes números (KORN; KORN; KROISANDT, 2010). Essa metodologia tem sido utilizada em diversos estudos para a solução de problemas numéricos e geração de variáveis aleatórias em contextos estatísticos, de engenharia e financeiros (PAULA, 2014).

Ionides et al. (2017) apontaram que os métodos de Monte Carlo possibilitam a avaliação de intervalos de confiança e de testes de hipóteses em estudos com modelos mais complexos e em cenários onde a função de verossimilhança pode ser intratável. Flowers-Cano et al. (2018) utilizaram de métodos de simulação para a comparação de quatro formas de construção de intervalos de confiança no contexto de precipitação, recomendando estimadores em diferentes cenários.

O uso da simulação computacional para a obtenção de intervalos de confiança pode apresentar um bom desempenho, principalmente em situações onde não é possível a obtenção de intervalos exatos por meio da quantidade pivotal. Com isso, a verdadeira probabilidade de cobertura de um intervalo construído por computação intensiva é uma variável aleatória que estará em torno da taxa de cobertura que seria obtida pelo método exato (ALMEIDA; SILVA, 2015).

Shawiesh, Banik e Kibria (2011) computaram o desempenho de intervalos de confiança de alguns estimadores. Foram utilizadas diferentes distribuições de probabilidade (normal, qui-quadrado e log-normal) e tamanhos amostrais $n \in \{5, 10, 15, 20, 40, 50, 70, 100\}$, obtendo regiões de confiança $(1 - \alpha)$ e possibilitando a comparação das taxas de cobertura de cada estimador.

1.5 AVALIAÇÃO DE ESTIMADORES INTERVALARES

No processo de avaliação de um estimador intervalar uma opção é verificar a probabilidade de cobertura gerada por cada cenário (SHAWIESH; BANIK; KIBRIA, 2011). A probabilidade de cobertura indica o percentual de vezes no qual o intervalo estimado contém o real valor do parâmetro. Conceito que também pode ser entendido como uma taxa de acurácia do estimador, ou seja, o percentual de vezes no qual o intervalo contém o real parâmetro.

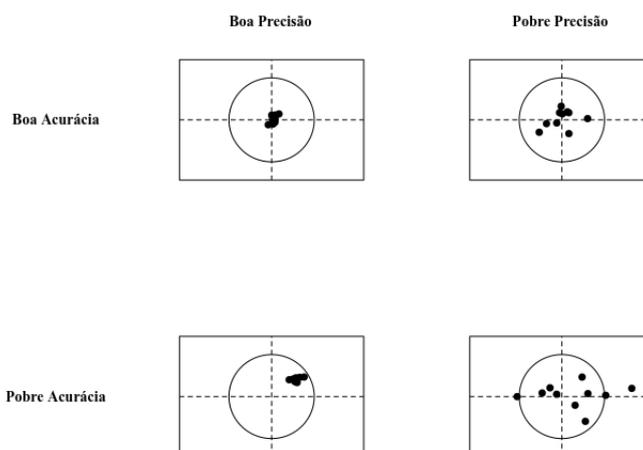
Mood, Graybill e Boes (1974) apontaram dois problemas no processo de construção de intervalos de confiança. O primeiro relacionado com a escolha da metodologia para a construção dos intervalos, e o segundo relacionado com a escolha de um intervalo satisfatório dentre todos possíveis. Na literatura, Patino e Ferreira (2015) indicaram que estudos almejam intervalos com menor comprimento, em consequência da maior precisão.

Dessa maneira, estudos unem dois critérios para a seleção de intervalos: a probabilidade de cobertura e o tamanho médio do intervalo. Isso pode ser visto no estudo de Scacabarozzi e Diniz (2007) na comparação de intervalos para o parâmetro da distribuição de Poisson, na avaliação de estimadores para proporções de Carari et al. (2010) e na análise do desempenho dos intervalos gerados pela distribuição generalizada de valores extremos (GEV) no estudo da duração máxima de secas em Butturi-Gomes, Beijo e Avelar (2018).

No entanto, uma problemática associada ao utilizar dois critérios de seleção é a de que um estimador pode ter uma boa cobertura e um tamanho de intervalo amplo, ou o inverso dessa situação. Fato que pode não trazer muita informação ao pesquisador, sendo necessário a comparação destes de forma separada. Tal forma foi realizada por Carari et al. (2010), analisando as taxas de acurácia seguido do comprimento médio para a indicação dos estimadores com melhor desempenho em cada cenário simulado.

De modo a exemplificar a condição entre precisão e acurácia de estimadores, a Figura 3 ilustra o comportamento entre o nível de precisão e de acurácia.

Figura 3 – Comportamento esperado de estimadores considerando o nível de precisão e acurácia



Fonte: Do autor.

É válido notar que a acurácia relaciona-se com a exatidão do estimador e a precisão relaciona-se com a menor dispersão. Note que um estimador com boa precisão e acurácia tem estimativas mais próximas do 'alvo', ou seja, o centro do arco. Por outro lado, um estimador com boa acurácia e pobre precisão apresenta intervalos de confiança reduzidos, mas sua 'mira' é prejudicada, sendo deslocada do centro do arco.

REFERÊNCIAS

- ALMEIDA, G; SILVA, I. Intervalos de confiança via simulação Monte Carlo : o estado da arte. **Revista da Estatística da Universidade Federal de Ouro Preto**, Ouro Preto, v. 4, p. 21-43, 2015. Disponível em: <http://www.repositorio.ufop.br/handle/123456789/6939>. Acesso em: 1 jan. 2019.
- ARNOLD, B; GHOSH, I. Bivariate Beta and Kumaraswamy models developed using the Arnold-ng bivariate Beta distribution. **REVSTAT**, Lisboa, v. 14, n. 2, p. 223-250, 2017. Disponível em: <https://www.ine.pt/revstat/pdf/rs170204.pdf>. Acesso em: 19 fev. 2019.
- BUTTURI-GOMES, D; BEIJO, L; AVELAR, F. On Modeling the maximum duration of dry spells: a simulation study under a Bayesian approach. **Theoretical and applied climatology**, Vienna, v. 137, n. 1-2, p. 1337-1346, 2018. Disponível em: <https://link.springer.com/article/10.1007%2Fs00704-018-2684-1>. Acesso em: 10 abr. 2018.
- CARARI, M. et al. Estimação de diferenças entre duas proporções binomiais via bootstrap. **Revista Brasileira de biometria**, São Paulo, v. 28, n. 3, p. 112-134, 2010.
- CARRODUS, M; GILES, D. The exact distribution of R2 when the regression disturbances are autocorrelated. **Economics letters**, v. 38, n. 4, p. 375-380, 1992. Disponível em: <https://www.sciencedirect.com/science/article/pii/016517659290021P>. Acesso em: 21 set. 2018.
- CASELLA, G; BERGER, R. **Statistical inference**. 2. ed. Pacific Grove: Duxbury Press, 2001.
- CRAMER, J.S. Mean and variance of R2 in small and moderate samples. **Journal of econometrics**, v. 35, ed. 2-3, p. 253-266, 1987. Disponível em: <https://www.sciencedirect.com/science/article/pii/0304407687900273>. Acesso em: 21 jun. 2018.
- DI BUCCHIANICO, A. Coefficient of determination (R2). In: RUGGERI, F; KENETT, R; FALTIN, F. **Encyclopedia of statistics in quality and reliability**, [s.l.]: John Wiley & Sons, 2008. p. 1-2.
- FERRARI, S; CRIBARI-NETO, F. Beta regression for modelling rates and proportions. **Journal of applied statistics**, v. 31, p. 799-815, 2004. Disponível em: <https://www.tandfonline.com/doi/abs/10.1080/0266476042000214501>. Acesso em: 1 dez. 2018.
- FLOWERS-CANO, R et al. Comparison of Bootstrap confidence intervals using monte carlo simulations. **Water**, v. 10, n. 2, 2018. Disponível em: <https://www.mdpi.com/2073-4441/10/2/166>. Acesso em: 11 dez. 2018.
- FOSTER, D; SMITH, T; WHALEY, R. Assessing goodness-of-fit of asset pricing models: The distribution of the maximal R2. **The journal of finance**, v. 52, n. 2, p. 591-607, 1997. Disponível em: <https://www.jstor.org/stable/2329491>. Acesso em: 23 jun. 2018.

FRIGG, R; NGUEYN, J. Models and representation. In: MAGNANI, L; BERTOLOTTI, T. **Springer handbook of model-based science**. Berlin: Springer, 2017. cap. 3, p. 51-96.

GILBERT, J. K; OSBORNE, R. J. The use of models in science and science teaching. **European journal of science education**, v. 2, n. 1, p. 3-13, 1980. Disponível em: <https://www.tandfonline.com/doi/abs/10.1080/0140528800020103>. Acesso em: 15 jul. 2018.

HAMID, H. et al. Investigating the power of goodness-of-fit tests for multinomial logistic regression. **Communications in statistics**, v. 47, n. 4, p. 1039-1055, 2018. Disponível em: <https://www.tandfonline.com/doi/abs/10.1080/03610918.2017.1303727>. Acesso em: 16 jul. 2018.

HOFFMANN, R. **Análise de regressão: uma introdução à econometria**. 5. ed. Piracicaba: Hucitec, 2016. Disponível em: <http://www.livrosabertos.sibi.usp.br/portaldelivrosUSP/catalog/book/73>. Acesso em: 16 jul. 2018.

IONIDES, E. Monte Carlo profile confidence intervals. **Journal of the royal society interface**, v. 14, n. 132, 2017. Disponível em: <https://royalsocietypublishing.org/doi/full/10.1098/rsif.2017.0126>. Acesso em: 18 jul. 2018.

JOHNSON, P; BEVERLIN, M. **Beta distribution**. 2013. Disponível em: <http://pj.freefaculty.org/guides/stat/Distributions/DistributionWriteups/Beta/Beta.pdf>. Acesso em: 08 nov. 2018.

JOHNSON, R. Applications of the Beta Distribution Part 1: Transformation Group Approach. **ArXiv.org**, 2013. Disponível em: <https://arxiv.org/abs/1307.6437>. Acesso em: 8 nov. 2018.

JOHNSON, N.; KOTZ, S.; BALAKRISHNAN, N. **Continuous univariate distributions**. 2. ed. [s.l.]: John Wiley & Sons, 1994. v. 1. Disponível em: <https://www.wiley.com/en-us/Continuous+Univariate+Distributions%2C+Volume+1%2C+2nd+Edition-p-9780471584957>. Acesso em: 11 nov. 2018.

KERMAN, J. A closed-form approximation for the median of the beta distribution. **arXiv.org**, 2011. Disponível em: <https://arxiv.org/abs/1111.0433>. Acesso em: 15 dez. 2018.

KORN, R; KORN, e; KROISANDT, G. **Monte Carlo methods and models in finance and insurance**. Londres: Chapman And Hall, 2010.

MONTGOMERY, D; PECK, E; VINING, G. **Introduction to Linear regression analysis**. 4. ed. Hoboken: John Wiley & Sons, 2006.

MOOD, A; GRAYBILL, F; BOES, D. **Introduction to theory of Statistics**. Estados Unidos: McGraw-Hill. 1974.

OHTANI, K. The density functions of R² and, and their risk performance under asymmetric loss in misspecified linear regression models. **Economic modelling**, v. 11, n. 4, 1994. Disponível em: <https://www.sciencedirect.com/science/article/pii/0264999394900035>. Acesso

em: 4 jun. 2018.

OHTANI, K; TANIZAKI, H. Exact distributions of R^2 and adjusted R^2 in a linear regression model with multivariate t error terms. **Journal of the Japan statistical society**, v. 34, n. 1, p. 101-109, 2004. Disponível em: <http://www2.econ.osaka-u.ac.jp/tanizaki/cv/papers/distr2.pdf>. Acesso em: 4 jun. 2018.

PATINO, C; FERREIRA, J. Intervalos de confiança: uma ferramenta útil para estimar o tamanho do efeito no mundo real. **Jornal Brasileiro de pneumologia**, v. 41, n. 6, 2015.

PAULA, Renato de. **Método de monte carlo e aplicações**. 2014. 83 f. Monografia (Especialização) - Universidade Federal Fluminense, Volta Redonda, 2014. Disponível em: <<https://app.uff.br/riuff/bitstream/1/4180/1/RenatoRicardoDePaula%202014-2.PDF>>. Acesso em: 19 jul. 2018.

QUININO, R.; REIS, E.; BESSEGATO, L. Using the coefficient of determination R^2 to test the significance of multiple linear regression. **Teaching statistics**, v. 35, n. 2, p. 84-88, 2012. Disponível em: <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-9639.2012.00525.x>. Acesso em: 19 ago. 2018.

SCACABAROZI, F; DINIZ, C. Uma comparação entre intervalos de credibilidade e o intervalo de confiança clássico para o parâmetro da distribuição de Poisson. In: SIMPÓSIO NACIONAL DE PROBABILIDADE E ESTATÍSTICA, 19., 2007, São Paulo. **Anais...** São Paulo: 19º Simpósio Nacional de Probabilidade e Estatística, 2007. p. 1 - 6. Disponível em: <<http://www2.ime.unicamp.br/sinape/sites/default/files/artigo2.pdf>>. Acesso em: 21 ago. 2018.

SHAWIESH, M; BANIK, S; KIBRIA, B. A simulation study on some confidence intervals for the population standard deviation. **SORT**, v. 35, n. 2, p. 83-201, 2011. Disponível em: <https://core.ac.uk/download/pdf/81625873.pdf>. Acesso em: 9 dez. 2018.

SHOU, W. et al. Theory, models and biology. **ELIFE**, v. 4, 2015. Disponível em: <https://elifesciences.org/articles/07158>. Acesso em: 1 jan. 2019.

SONG, W. Relationships among some univariate distributions. **IIE Transactions**, v. 37, n. 7, p. 651-656, 2005. Disponível em: <https://www.tandfonline.com/doi/abs/10.1080/07408170590948512?journalCode=uiie20>. Acesso em: 2 fev. 2019.

SU, X; YAN, X; TSAI, C. Linear regression. **WIREs**, [S. l.], v. 4, n. 3, p. 275-294, 2012. Disponível em: <https://dl.acm.org/citation.cfm?id=3163603>. Acesso em: 10 fev. 2019.

UYANIK, G; GÜLER, N. A study on multiple linear regression analysis. **Procedia**, v. 106, n. 10, p. 234-240, 2013. Disponível em: <https://www.sciencedirect.com/science/article/pii/S1877042813046429>. Acesso em: 5 mar. 2018.

TULUPYEV, A. Beta prime regression with application to risky behavior frequency

screening. **Stat. med.**, v. 32, p. 4044–4056, 2013. Disponível em:
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3789864/>. Acesso em: 5 mar. 2019.

WEATHERBURN, C. **A first course in mathematical statistics**. Cambridge: University Press, 1949. v. 32.

CAPÍTULO 2 - ESTIMAÇÃO INTERVALAR DO COEFICIENTE DE DETERMINAÇÃO: UMA APLICAÇÃO NA QUALIDADE DE AJUSTE DE MODELOS LINEARES

Resumo

A qualidade de ajuste de modelos de regressão é um tópico de grande importância na modelagem estatística. Para a análise dessa qualidade, a literatura geralmente utiliza de medidas como o coeficiente de determinação (R^2), o coeficiente de determinação ajustado (R_a^2) e algumas métricas baseadas em erros. O R^2 é popularmente utilizado pela rápida interpretação como a proporção da variação da variável resposta explicada pelo conjunto de preditores. Entretanto, é sabido que este pode conter alguns problemas, como o inflacionamento em modelos com elevado número de covariáveis não relevantes ou em modelos com reduzido tamanho amostral. De modo a reduzir o problema em modelos muito complexos, a literatura estabeleceu o R_a^2 , penalizando a inclusão de covariáveis não relevantes ao modelo. Porém, essa penalização pode fazer com que o R_a^2 seja negativo, dificultando a interpretação. Dado o contexto, estudos comumente têm tomado o R^2 como uma estatística que estima um parâmetro ρ^2 , sendo este associado a qualidade de ajuste que um modelo possui de forma populacional. Tal forma possibilita a construção de intervalos de confiança e de testes de hipóteses acerca da qualidade do modelo. Entretanto, a questão inferencial é um assunto não encerrado na literatura, pois os estimadores propostos até então consideram diferentes regiões do espaço paramétrico e podem gerar intervalos largos em alguns cenários. Desta forma, este trabalho tem como objetivo estudar, de forma intervalar, os estimadores paramétricos apresentados pela literatura e propor outros três estimadores baseados na distribuição de probabilidade Beta, que é comumente atribuída a modelagem de R^2 . Para isso, será realizado um estudo empírico de simulação em cenários estabelecidos, analisando medidas de desempenhos dos estimadores. Posteriormente, de modo a ilustrar a utilização dos estimadores, serão construídos intervalos de confiança para analisar a qualidade de ajuste de um modelo de regressão no contexto da estatística experimental. Como resultado, notou-se que a recomendação de um estimador com melhor desempenho em todo espaço paramétrico não é trivial, fazendo com que as recomendações sejam dadas em três regiões do espaço.

Palavras-Chave: Coeficiente de Determinação. Estimação Paramétrica. Regressão.

1 INTRODUÇÃO

A demanda por qualidade de ajuste em regressão tem possibilitado a criação de novos modelos e inferências em métricas que determinam a adequabilidade do modelo ao fenômeno. Uma destas corresponde ao coeficiente de determinação (R^2), que atua mensurando a proporção da variabilidade da variável resposta explicada pelo conjunto de preditores (ZHANG, 2017), sendo essa uma das medidas mais populares na análise da qualidade de modelos lineares (PI-EPHO, 2018).

Considerando o contexto matricial de um modelo de regressão, o R^2 pode ser obtido pela seguinte operação:

$$R^2 = \frac{\mathbf{b}'\mathbf{X}'\mathbf{X}\mathbf{b}}{\mathbf{b}'\mathbf{X}'\mathbf{X}\mathbf{b} + \mathbf{e}'\mathbf{e}},$$

sendo \mathbf{b} o vetor de parâmetros, \mathbf{X} a matriz de delineamento e \mathbf{e} o vetor de erros. Maiores detalhes podem ser obtidos em Cramer (1987). A partir da interpretação de proporção explicada, é evidente que a literatura busque modelos em que $R^2 \rightarrow 1$. Entretanto, estudos acerca de R^2 indicam que essa situação pode ser ilusória, pois este pode ser inflacionado em modelos com reduzido tamanho amostral ou com elevado número de covariáveis (CRAMER, 1987), (OHTANI; TANIZAKI, 2004), (QUININO; REIS; BESSEGATO, 2012).

De modo a penalizar o R^2 em modelos com elevado número de covariáveis, a literatura propôs uma versão modificada, denominada como coeficiente de determinação ajustado (R_a^2), dado por:

$$R_a^2 = 1 - (1 - R^2) \left(\frac{n - 1}{n - (k + 1)} \right),$$

sendo n o número de observações e k o respectivo número de covariáveis da regressão (DRAPER; SMITH, 1998). Com essa modificação, a fragilidade do aumento de R^2 com o acréscimo de covariáveis ao modelo é superada, fazendo com que o R_a^2 possa aumentar ou diminuir sob essa situação. Entretanto, em casos onde o modelo contém muitas covariáveis e reduzido tamanho amostral, o R_a^2 pode ser negativo, perdendo o sentido de interpretação como proporção explicada (QUININO; REIS; BESSEGATO, 2011).

Mesmo com as precauções apontadas na literatura estatística, estudos aplicados utilizam essas duas medidas (R^2 , R_a^2) para avaliar a qualidade de modelos, como visto em Yogesha et al. (2018), Bineesh et al. (2018) e Song, Deng e Ren (2019). Em conjunto a isso, novos

estudos acerca da expansão do R^2 para a classe de modelos lineares generalizados (MLG) estão sendo discutidos em Zhang (2017), Nakagawa, Johnson e Schielzeth (2017) e Piepho (2018). No contexto de MLG, Piepho (2018) apontou que o coeficiente de determinação de um modelo linear generalizado (R^2_{mlg}) comporta-se, em alguns casos, como o usual R^2 . Entretanto, o problema de inflacionamento incluindo covariáveis irrelevantes ao modelo ainda é persistente (ZHANG, 2018).

Em meio a essas situações, estudos sugerem a realização de inferências em relação ao parâmetro ρ^2 , estimado por R^2 , representando a qualidade que um modelo qualquer possuiria se as infinitas observações do fenômeno fossem coletadas. O que faz como natural a construção de intervalos de confiança e de testes de hipóteses que visem melhorar a incerteza a respeito a qualidade do modelo.

No entanto, a questão inferencial não aparenta ser um assunto encerrado na literatura, admitindo diferentes estimadores em regiões do espaço paramétrico. Mais precisamente, para o espaço onde $\rho^2 = 0$, tratado em Weatherburn (1949), indicando que R^2 segue distribuição Beta com parâmetros $\frac{k}{2}$ e $\frac{n-k-1}{2}$. E, de forma complementar da abordada por Weatherburn (1949), Cramer (1987) explorou a região de $\rho^2 \neq 0$, a partir da distribuição exata, disponível em 1.5.

A partir de um estudo de simulação, Ohtani e Tanizaki (2004) apontaram intervalos de confiança para ρ^2 são largos em regiões do espaço paramétrico. Essa situação estimula a inserção de novos estimadores, de modo a aperfeiçoar as inferências e permitir tomadas de decisão mais precisas em relação a qualidade do modelo. Sendo a qualidade de ajuste fundamental em diversas áreas do conhecimento como química, alimentos e ciência e tecnologia, permitindo o desenvolvimento e a otimização de produtos e processos (GRANATO; CALADO, 2014).

Sendo assim, este trabalho tem por objetivo estudar empiricamente a estimação intervalar de ρ^2 a partir de cinco estimadores paramétricos. Sendo dois deles apresentados pela literatura e três propostos nesse trabalho. Para isso, as medidas de acurácia, precisão e um índice de desempenho serão computadas em cenários estabelecidos pela simulação, permitindo a recomendação dos estimadores com melhor desempenho.

Posteriormente, a estimação intervalar da qualidade de ajuste de um modelo de regressão será ilustrada em um experimento sensorial, relativo a durabilidade da vida pós colheita de bananas do tipo maçã expostas ao 1-metilciclopropeno. Possibilitando, desta forma, a escolha de um período de exposição ao produto que resulte, em média, a maior durabilidade do fruto nas prateleiras de supermercado.

2 MATERIAL E MÉTODOS

2.1 ESTIMADORES INTERVALARES UTILIZADOS E PROPOSTOS

Inicialmente, optou-se pela utilização de dois estimadores já fundamentados na literatura. O primeiro estimador foi denominado como W e esteve associado a distribuição analisada por Weatherburn (1949) e Foster, Smith e Whaley (1997). O segundo estimador considerado foi baseado na distribuição exata do coeficiente de determinação, tratado em Cramer (1987). Este último foi denominado como E^* ¹.

Em conjunto aos apontados anteriormente, foram considerados três estimadores baseados na distribuição de probabilidade Beta. Tal proposta foi fundamentada pela condição de Cramer (1987) de que assim como o valor esperado do coeficiente converge para ρ^2 , há indícios de que o mesmo ocorra para a mediana e moda. Desta forma, reparametrizou-se por três formas a distribuição Beta, onde em duas delas a parametrização foi atribuída ao parâmetro a da distribuição de W .

Para justificar a reparametrização de a é necessário lembrar que ambos parâmetros da Beta dão forma a distribuição, com as seguintes características apontadas por Owen (2008): a) a distribuição torna-se assimétrica à esquerda a medida em que o parâmetro b é maior que a e o contrário acontece quando b é menor que a ; e b) quando $a > 1$ e $b > 1$ a distribuição é unimodal e quando $a = b = 1$ tem-se um caso particular da distribuição uniforme.

No caso de W , tem-se que $a = \frac{k}{2}$ e $b = \frac{n-k-1}{2}$, resultando em uma distribuição assimétrica à esquerda sempre que $k < \frac{n-1}{2}$, pois $b > a$. E essa condição sempre será verificada, pois a literatura recomenda $n \gg k$, justificando a proposta de W para a inferência em $\rho^2 = 0$. Por outro lado, indica uma maior probabilidade de um modelo qualquer ter qualidade baixa, e essa situação só será modificada com o aumento de k , o que pode não ser interessante.

De modo a estender essa condição, foi necessário a reparametrização de a em W , possibilitando o deslocamento da assimetria da distribuição e removendo o fato da necessidade do acréscimo de número de covariáveis para o aumento da qualidade do modelo. Esses estimadores estão apresentados nas seções 2.1.1, 2.1.2 e 2.1.3.

¹Adotou-se como E^* em função de que E é comumente associado ao operador de esperanças.

2.1.1 Estimador M_1

Para a proposta do estimador M_1 , será considerado a existência da moda da distribuição Beta, fazendo com que os parâmetros $a > 1$ e $b > 1$. Sendo X uma variável aleatória com distribuição de probabilidade Beta, tem-se que a moda de X é dada por:

$$M_o(X) = \frac{a - 1}{a + b - 2},$$

assumindo $\rho^2 = M_o(X)$:

$$\rho^2 = \frac{a - 1}{a + b - 2},$$

A partir disso, o parâmetro a foi isolado da seguinte forma:

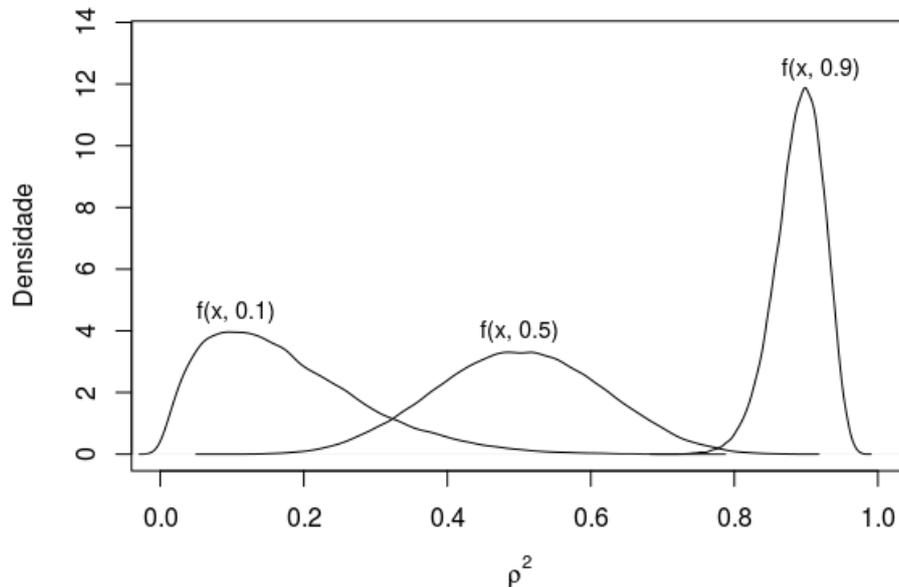
$$a = \frac{b\rho^2 - 2\rho^2 + 1}{1 - \rho^2}.$$

Mantendo o parâmetro b dado pela distribuição de W e reparametrizando a , tem-se que a distribuição do estimador M_1 é dada por:

$$R_{M_1}^2 \sim Beta\left(\frac{b\rho^2 - 2\rho^2 + 1}{1 - \rho^2}, \frac{n - k - 1}{2}\right). \quad (2.1)$$

em que $(1 - \rho^2) \neq 0$. Essa distribuição será unimodal se $\frac{b\rho^2 - 2\rho^2 + 1}{1 - \rho^2} > 1$ e $\frac{n - k - 1}{2} > 1$. Para isso $n > k + 3$. De modo a ilustrar o comportamento da distribuição, a Figura 4 denota o cenário onde $n = 20$, $k = 1$ e $\rho^2 \in \{0,1; 0,5; 0,9\}$.

Figura 4 – Comportamento da distribuição de M_1 considerando $n = 20$, $k = 1$ e $\rho^2 \in \{0,1; 0,5; 0,9\}$



Fonte: Do autor.

Note que a moda da distribuição aproxima-se do valor paramétrico fixado de ρ^2 . Por fim, dada a condição de que em modelos de regressão a quantidade ρ^2 é desconhecida, este será substituído pela estimativa pontual, ou seja, R^2 . Desta forma o intervalo conforme M_1 será dado por $M_1 = [q_1; q_2]$, sendo q_1 o quantil superior $1 - \frac{\alpha}{2}$ e q_2 o quantil $\frac{\alpha}{2}$ da distribuição beta parametrizada em 2.1.1.

2.1.2 Estimador M_2

O estimador M_2 foi proposto assumindo ρ^2 como a média da distribuição Beta. Considerando X uma variável aleatória com distribuição Beta, tem-se que a média é dada por:

$$E[X] = \frac{a}{a+b},$$

sob restrição de que $a+b \neq 0$. Atribuindo $\rho^2 = E[X]$ e isolando o parâmetro a , tem-se que:

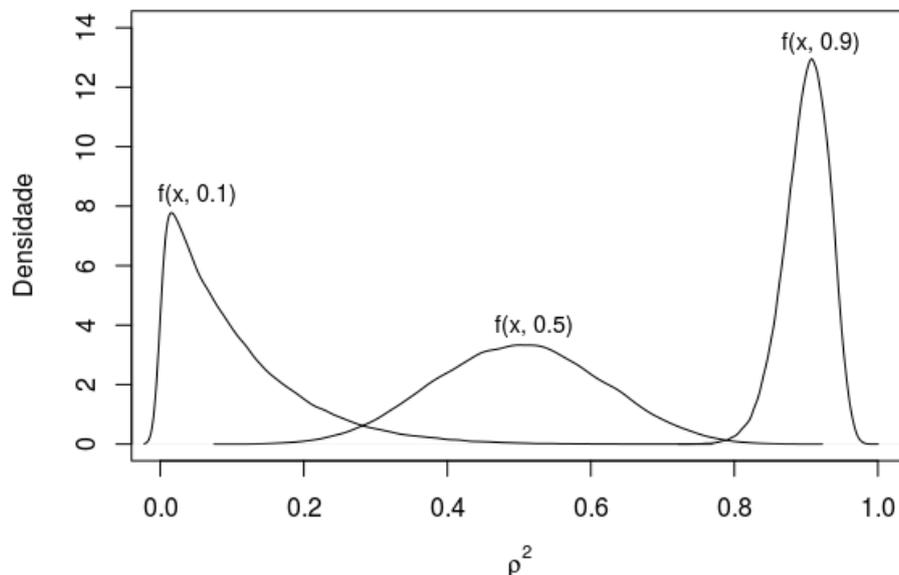
$$a = \frac{-b\rho^2}{\rho^2 - 1},$$

na restrição de que $\rho^2 - 1 \neq 0$. Então, a distribuição do estimador M_2 será dada por:

$$R_{M_2}^2 \sim \text{Beta} \left(\frac{-b\rho^2}{\rho^2 - 1}, \frac{n - k - 1}{2} \right).$$

Para a existência da distribuição deste estimador, é necessário que os parâmetros sejam maiores que zero. E como consequência deste fato, $n > k + 1$. A Figura 5 ilustra o comportamento da distribuição considerando $n = 20$, $k = 1$ e $\rho^2 \in \{0,1; 0,5; 0,9\}$.

Figura 5 – Comportamento da distribuição de M_2 considerando $n = 20$, $k = 1$ e $\rho^2 \in \{0,1; 0,5; 0,9\}$



Fonte: Do autor.

Novamente, em função de ρ^2 ser desconhecido, este foi substituído pelo valor estimado R^2 . O intervalo conforme M_2 será dado por $M_2 = [q_1; q_2]$, sendo q_1 o quantil superior $1 - \frac{\alpha}{2}$ e q_2 o quantil $\frac{\alpha}{2}$ da distribuição beta parametrizada em 2.1.2.

2.1.3 Estimador P^*

O último estimador considerado por este trabalho foi denominado P^* ², relacionando a esperança da soma de quadrados da regressão e soma de quadrados dos erros. Para essa

²Adotou-se como P^* em função de que a letra P é comumente associada à probabilidade.

proposta, supõem-se que a qualidade de ajuste de um modelo não dependa do número de covariáveis, alterando os dois parâmetros da distribuição do estimador W . Como apontado anteriormente, R^2 é calculado pela razão entre SQR e SQT. Então, o parâmetro a foi reparametrizado da seguinte forma:

$$\rho^2 = \frac{E(SQR)}{E(SQR + SQE)} = \frac{k}{k + (n - k)} = \frac{k}{n}.$$

Assumindo que o parâmetro $a = k$, pode-se escrever $a = n\rho^2$. Por outro lado, atribuindo b como a $E(SQE)$, então:

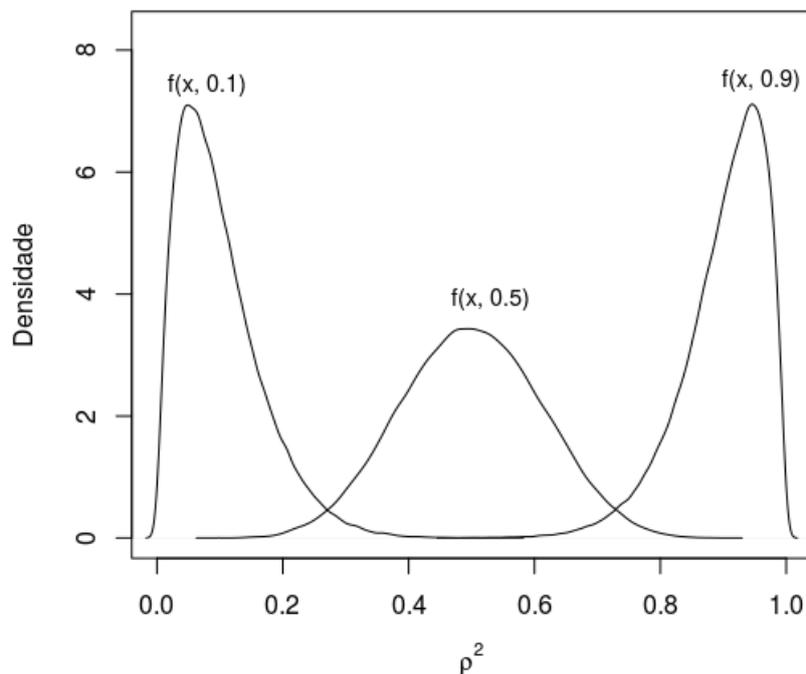
$$b = E(SQE) = n - a.$$

Por consequência: $b = n - n\rho^2$. Assim, a distribuição do estimador P^* será dada por:

$$R^2 \sim \text{Beta}(\rho^2 n, n - n\rho^2). \quad (2.2)$$

Sob a condição de existência da distribuição $\rho^2 n > 0$ e $n - n\rho^2 > 0$. Dada as características, a Figura 7 ilustra o comportamento da distribuição para $n = 20$ e $\rho^2 \in \{0,1; 0,5; 0,9\}$.

Figura 6 – Comportamento da distribuição de P^* considerando $n = 20$ e $\rho^2 \in \{0,1; 0,5; 0,9\}$



Novamente, em função do desconhecimento de ρ^2 , este será substituído pela estimativa pontual R^2 . O intervalo conforme P^* será dado por $P^* = [q_1; q_2]$, sendo q_1 o quantil superior $1 - \frac{\alpha}{2}$ e q_2 o quantil $\frac{\alpha}{2}$ da distribuição beta parametrizada em 2.1.3.

2.2 GERAÇÃO DOS DADOS E CENÁRIOS AVALIADOS

Para este estudo foram gerados mil valores aleatórios de R^2 para analisar modelos de regressão com até oito covariáveis. Considerou-se este número limite de covariáveis pois são os mais usuais na literatura. A simulação foi realizada combinando valores de n , k e ρ^2 e tomando o operador de esperanças das distribuições χ^2 a partir da Equação 1.1. A forma de amostragem foi dada pelo Algoritmo 1:

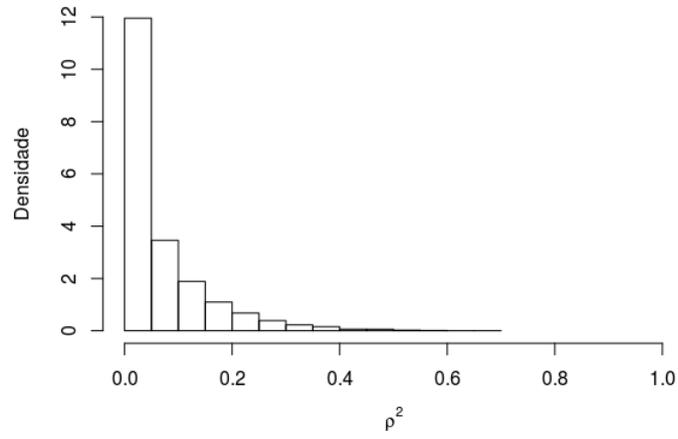
Algoritmo 1 Passos para a simulação de R^2

- 1: Fixar um valor de $\rho^2 \in \{0; 0,1; 0,2; 0,3; 0,4; 0,5; 0,6; 0,7; 0,8; 0,9\}$.
 - 2: Fixar um tamanho amostral, definido por $n \in \{15, 50, 100\}$.
 - 3: Calcular o parâmetro de não centralidade λ , dado por: $\lambda = n \left(\frac{\rho^2}{1-\rho^2} \right)$.
 - 4: Fixar o número de covariáveis $k \in \{1, 2, 3, 4, 5, 6, 7, 8\}$.
 - 5: Sortear um valor da distribuição $\chi^2_{\{\lambda, k\}}$ e um valor da distribuição $\chi^2_{\{0, n-k\}}$ que correspondem a distribuição de SQR e SQT, respectivamente, conforme 1.3 e 1.4.
 - 6: Tomar a esperança por meio da Equação 1.1. Sendo este resultado o valor esperado de uma realização das combinações de n , k e λ analisados.
-

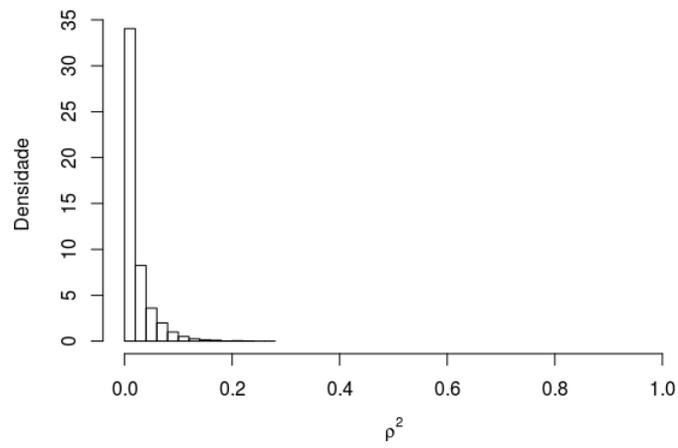
A partir disso, comparou-se o padrão dos valores amostrados com os cenários estabelecidos em Cramer (1987). A Figura 7 apresenta um cenário simulado considerando $k = 2$, $\rho^2 = 0$ e $n \in \{15, 50, 100\}$.

Com a simulação de valores amostrais, analisou-se o desempenho dos estimadores intervalares a um nível de significância $\alpha = 0,05$ com base nas distribuições apresentadas na seção 2.1. As taxas de acurácia e as precisões de cada estimador foram avaliadas. De modo que a precisão foi definida como o complementar da média de comprimento dos intervalos construídos.

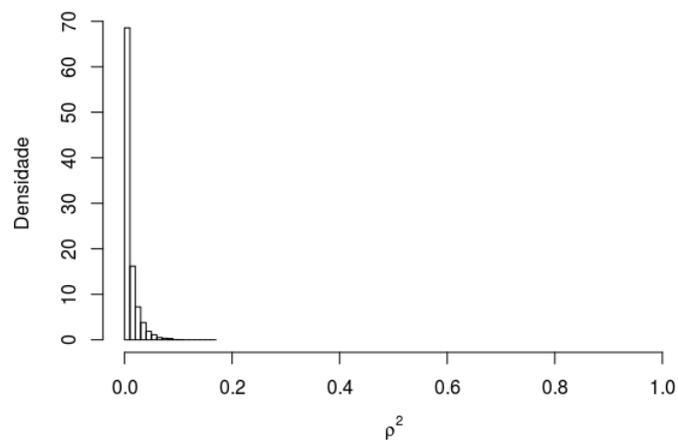
Figura 7 – Histogramas das amostras de R^2 para diferentes combinações de n . Sendo (a): $n = 15$, (b): $n = 50$ e (c): $n = 100$



(a)



(b)



(c)

2.3 ÍNDICE DE DESEMPENHO DE ESTIMAÇÃO INTERVALAR (τ_α)

Devido a dificuldade na escolha do estimador que gere intervalos de confiança com maior qualidade, optou-se pela elaboração de um índice. Esse índice resulta da combinação entre a taxa de acurácia e o tamanho médio dos intervalos gerados em cada cenário analisado.

O índice foi proposto da seguinte forma:

$$\tau_{\alpha,i} = \left(\frac{A_{(i)}}{\text{Max}(A_\omega)} \right) p_1 + \left(\frac{\text{Min}(M_\omega)}{M_{(i)}} \right) (1 - p_1), \quad i = 1, 2, 3, 4, 5.$$

onde $A_{(i)}$ e $M_{(i)}$ corresponde a taxa de acurácia e a média do comprimento dos intervalos do estimador i , respectivamente. $\text{Max}(A_\omega)$ corresponde à acurácia máxima dentre os estimadores analisados e, $\text{Min}(M_\omega)$ representa o valor mínimo de comprimento médio dos intervalos no cenário analisado.

O índice é ponderado por p_1 e $1 - p_1$, correspondendo ao peso atribuído à acurácia e precisão, respectivamente, sendo $p_1 + (1 - p_1) = 1$. Assim, o pesquisador pode optar pela penalização da medida de maior interesse. Entretanto, entende-se que essa penalização não é trivial, pois há contextos onde a acurácia pode ser de maior interesse que a precisão, por exemplo.

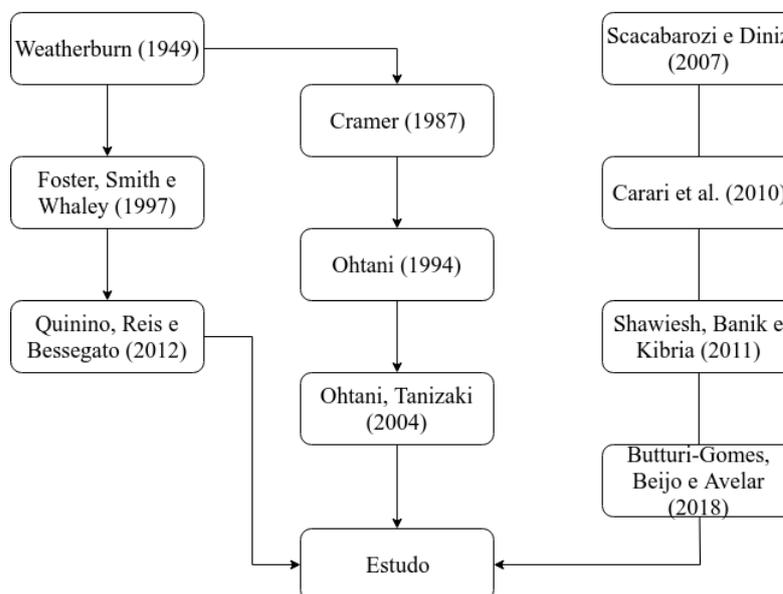
Devido ao contexto deste trabalho, atribuiu-se $p_1 = 0.5$, pois entende-se que no contexto de avaliação dos estimadores, as taxas de acurácia e os tamanhos médios são de grande importância e trazem informação, não sendo possível definir se uma delas deve ser mais penalizada.

Note que o índice é definido no espaço entre zero e um. Então, um estimador receberá um $\tau_\alpha = 1$ se, e somente se, este possuir a maior taxa de acurácia e a menor média de comprimento dos intervalos. E, os demais, serão ranqueados conforme a qualidade do melhor estimador no cenário analisado.

Após a simulação das amostras de R^2 descritas na seção 2.2, foram computados os intervalos de confiança para os estimadores W , E^* , M_2 , M_1 , P^* , bem como a acurácia, o comprimento médio e o índice de qualidade.

Portanto, este trabalho relaciona estudos que abordaram estimadores sobre o coeficiente de determinação com uma abordagem de avaliação computacional de intervalos de confiança, como apresentado na Figura 8.

Figura 8 – Diagrama com os principais estudos que analisaram distribuições para o coeficiente de determinação ou que avaliaram computacionalmente intervalos de confiança



Fonte: Do autor.

Como forma de interpretação da Figura 8, este trabalho engloba a distribuição de Weatherburn (1949) que foi ampliada por Foster, Smith e Whaley (1997) e aplicada por Quinino, Reis e Bessegato (2011). Além de utilizar a distribuição exata apresentada por Cramer (1987), também discutida por Ohtani (1994) e Ohtani, Tanizaki (2004). E, por fim, no contexto de simulação computacional, foram utilizadas das formas de avaliação de estimadores intervalares trabalhadas em Scacabarozi e Diniz (2007), Carari et al. (2010), Shawiesh, Banik e Kibria (2011) e Butturi-Gomes, Beijo e Avelar (2018).

2.4 APLICAÇÃO: DURABILIDADE PÓS COLHEITA DE BANANAS DO TIPO MAÇÃ

Neste estudo, serão utilizados dados provenientes do experimento realizado por Pinheiro (2007) no Departamento de Ciências dos Alimentos da Universidade Federal de Lavras. O experimento analisou a conservação de bananas do tipo maçã aplicando técnicas que prolonguem a vida pós colheita dos frutos, conservando seus atributos. Uma das metodologias utilizadas

no estudo foi a aplicação do 1-Metilciclopropeno (1-MCP), que é um produto utilizado para o prolongamento pós colheita e para a estabilidade de vida de vegetais.

Frutos completamente verdes foram colhidos na cidade de Lavras/MG e transportados para o Laboratório de Fisiologia da Universidade um dia após a colheita. Após isso, foram separados em buquês contendo 4 frutos cada. Em seguida, os buquês foram inseridos em caixas de isopor de 100 litros onde a aplicação de 1-MCP na concentração de 50nLL^{-1} na formulação em pó com 0,14% de ingrediente ativo foi realizada.

Os buquês ficaram expostos ao 1-MCP pelo período de 0, 6, 9, 12 ou 24 horas. Após a exposição, eles foram removidos e armazenados sob temperatura ambiente de 25°C . O experimento contou com 200 observações, das quais foram analisadas as seguintes variáveis: a) dias G7: Número de dias transcorridos onde os frutos alcançaram o grau máximo de amadurecimento, ou seja, bananas amarelas com manchas marrons, em cada período de exposição; e b) 1-MCP: Tempo de exposição dos frutos ao 1-MCP (em horas).

3 RESULTADOS E DISCUSSÃO

3.1 AVALIAÇÃO E COMPARAÇÃO DOS ESTIMADORES

Nas seções 3.1.1, 3.1.2 e 3.1.3 estão apresentados os resultados referentes a avaliação dos estimadores nos cenários definidos na seção 2. Para uma melhor discussão, foram exibidos inicialmente os gráficos relativos as taxas de acurácia, precisões e, por fim, os índices de qualidade. Ressalta-se que os resultados dos estimadores W e E^* estão apresentados nas cores azul e vermelho, possibilitando melhores comparações em relação aos propostos por este trabalho.

3.1.1 Taxas de acurácia

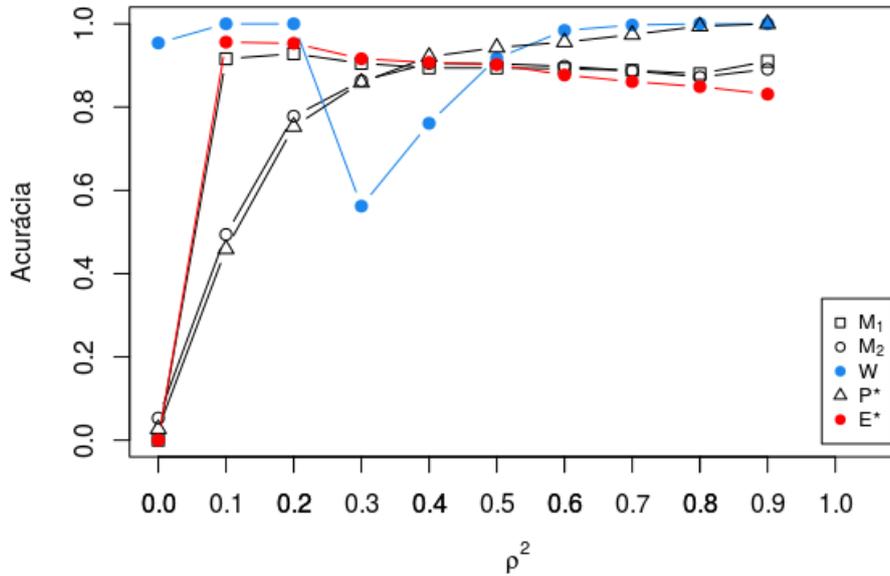
Considerando o modelo de regressão linear simples com quinze observações, apresentado na Figura 9, tem-se que somente o estimador W apresentou acurácia superior a 80% no cenário onde $\rho^2 = 0$, sendo diferente dos demais que apresentaram acurácia próxima a zero. A partir de $\rho^2 \geq 0,1$, os demais estimadores elevaram a acurácia de seus intervalos, superando a taxa de 80% no espaço de $\rho^2 \geq 0,5$. Houve também uma queda repentina da acurácia de W na região de $\rho^2 = 0,3$, atingindo um nível inferior a 60%.

Mantendo o modelo linear simples e aumentando o número de observações para 50, a Figura 10 aponta uma melhoria nas taxas de acurácia de todos os estimadores analisados no espaço $\rho^2 > 0,1$. Todos os intervalos gerados por cada estimador contiveram o real valor do parâmetro em, no mínimo, 70% das vezes. Na região de $\rho^2 = 0$, somente o estimador W manteve taxa de acurácia superior a 90% e os demais apresentaram uma taxa próxima a zero.

Com o acréscimo do número de observações no modelo linear simples para cem, tem-se que os estimadores mantiveram taxa de acurácia superior a 70% na região de $\rho^2 \geq 0,1$, o que pode ser verificado na Figura 11. Ressalta-se que o estimador W manteve a maior acurácia em grande parte deste cenário, apresentado maiores taxas de acurácia no espaço $\rho^2 \in \{0; 0,2; 0,3; 0,4; 0,5; 0,6\}$. Neste cenário, notou-se novamente uma queda da acurácia de W , porém na proximidade de $\rho^2 = 0,1$.

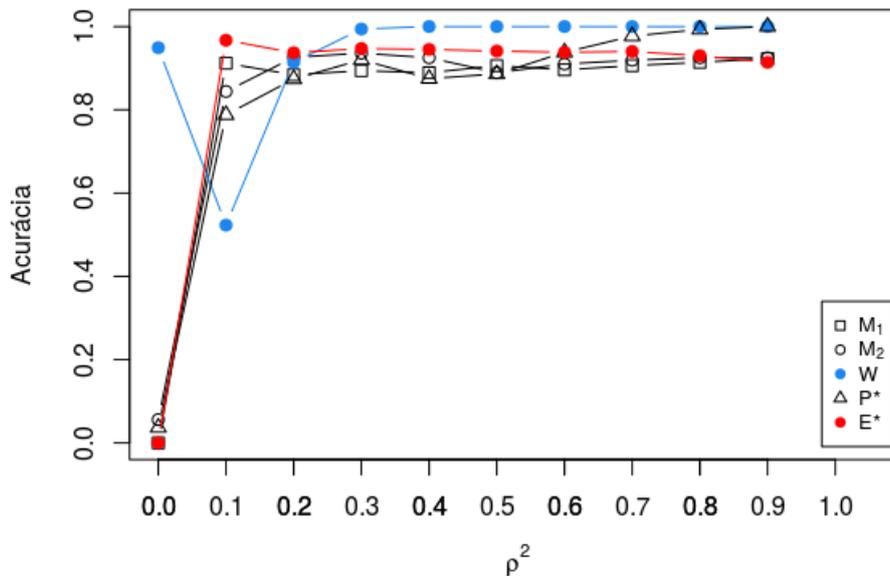
Para o modelo de regressão linear simples analisado anteriormente, houve uma melhoria

Figura 9 – Taxa de acurácia referente ao modelo de regressão onde $k = 1$ e $n = 15$



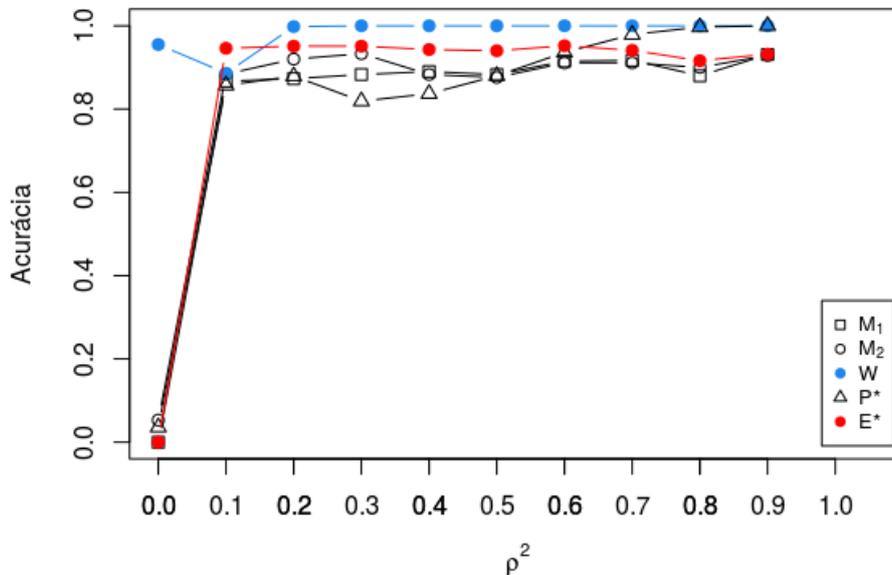
Fonte: Do autor.

Figura 10 – Taxa de acurácia referente ao modelo de regressão onde $k = 1$ e $n = 50$



Fonte: Do autor.

das taxas de acurácia de todos os estimadores com o aumento do tamanho amostral. Entretanto, os estimadores, com exceção de W , apresentaram taxas de acurácia próximas a zero no espaço onde $\rho^2 = 0$. É importante destacar que os estimadores P^* e W superaram os demais na região

Figura 11 – Taxa de acurácia referente ao modelo de regressão onde $k = 1$ e $n = 100$ 

Fonte: Do autor.

de $\rho^2 \geq 0,7$, ou seja, nos modelos com alta explicação do fenômeno.

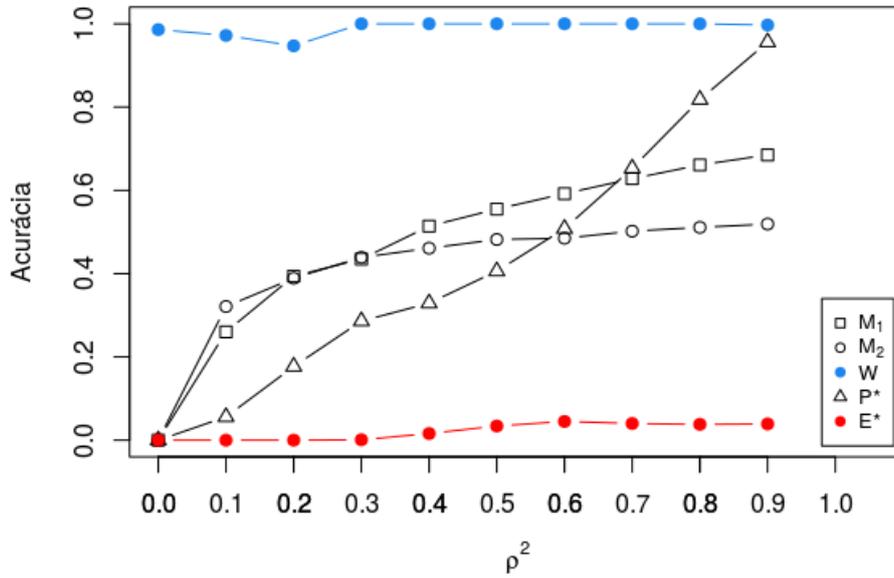
A Figura 12 indica o cenário de um modelo com oito covariáveis e quinze observações. Neste cenário, houve um domínio das taxas de acurácia de W em todo o espaço. Os demais estimadores apresentaram novamente baixa taxa em $\rho^2 = 0$. Mas, os estimadores M_1 , M_2 e P^* aumentaram esta na medida em que o valor de ρ^2 aumentou, com destaque para P^* que aproximou-se de W em $\rho^2 = 0,9$. Por outro lado, houve uma taxa de acurácia inferior a 10% pelo estimador E^* em todo espaço.

Com o aumento do tamanho amostral para cinquenta no modelo com oito covariáveis, houve uma melhoria das taxas de acurácia de todos os estimadores. Neste cenário, W novamente manteve certo domínio em relação aos demais. Entretanto, P^* apresentou taxas de acurácia próximas a W na região de $\rho^2 \geq 0,9$.

Considerando o modelo com cem observações e oito covariáveis, com resultado expresso na Figura 14, verificou-se a melhora do estimador E^* , aproximando-se dos demais estimadores. Novamente, houve um domínio da taxa de acurácia de W em grande parte do espaço, igualando-se a P^* no cenário de alto ρ^2 .

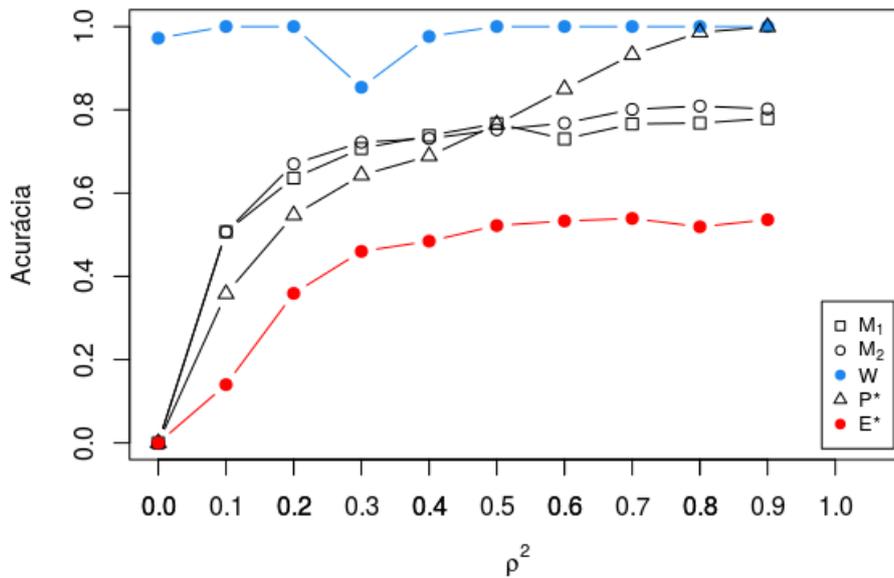
Neste sentido, é notória a melhoria na acurácia dos estimadores a medida em que n cresce na região de $\rho^2 \geq 0,1$. Essa condição assintótica indica certa consistência dos estimado-

Figura 12 – Taxa de acurácia referente ao modelo de regressão onde $k = 8$ e $n = 15$



Fonte: Do autor.

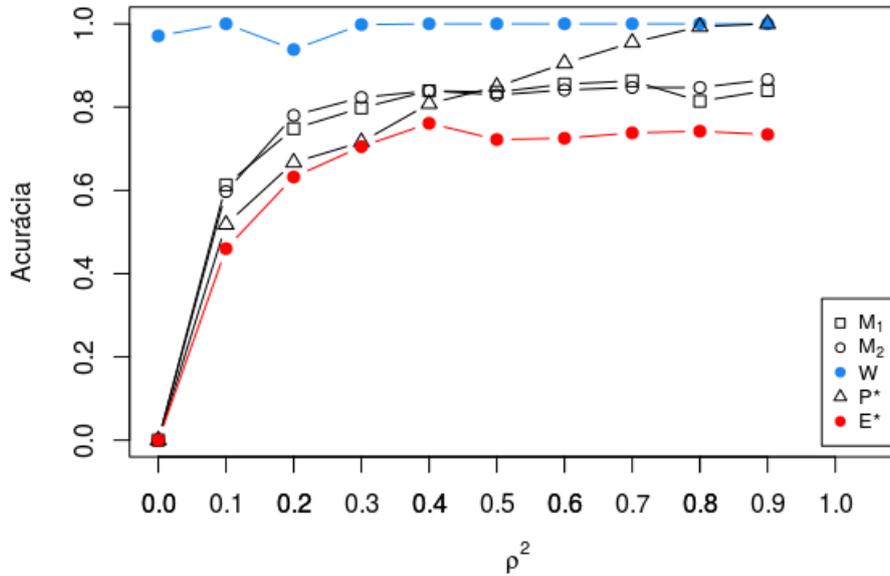
Figura 13 – Taxa de acurácia referente ao modelo de regressão onde $k = 8$ e $n = 50$



Fonte: Do autor.

res analisados, aumentando a probabilidade de conter o real parâmetro.

No cenário onde $\rho^2 = 0$, os estimadores M_2 , M_1 , P^* e E^* , tiveram baixa taxa de acurácia, condição esta que não melhorou com o aumento do tamanho amostral. Neste cenário, a baixa

Figura 14 – Taxa de acurácia referente ao modelo de regressão onde $k = 8$ e $n = 100$ 

Fonte: Do autor.

acurácia do estimador E^* pode ser justificada pela sua definição, de que a distribuição exata considera o cenário onde $\beta \neq 0$, ou seja, $\rho^2 \neq 0$. E, por outro lado, a alta acurácia de W também é justificada pela sua própria definição, baseada na condição de que $\rho^2 = 0$.

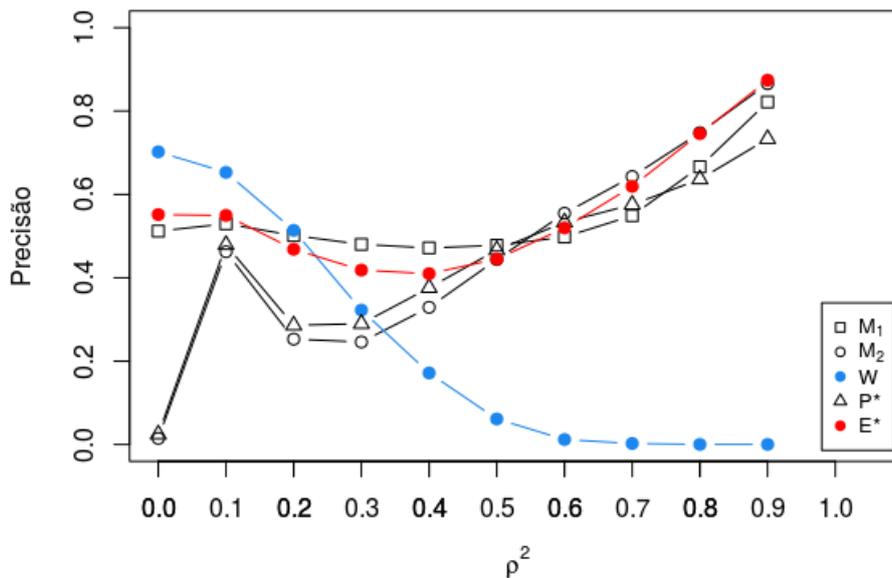
Além disso, notou-se uma situação de queda na acurácia do estimador W nas proximidades de $\rho^2 = 0,3$ em grande parte dos cenários avaliados. Há princípio não há uma justificativa para este fenômeno, abrindo a possibilidade de uma exploração computacional em mais pontos de ρ^2 na localidade.

Dada essas condições, os demais cenários encontram-se no Apêndice 5 e apresentaram comportamento similar aos tratados anteriormente. Dando continuidade à análise dos estimadores, a próxima seção apresenta o nível de precisão dos intervalos gerados. Como apontado, pode-se objetivar intervalos de menor comprimento e, conseqüentemente, maior precisão, aumentando o nível de informação a respeito do real valor do parâmetro analisado.

3.1.2 Níveis de precisão

A Figura 15 expressa a precisão dos intervalos computados para o modelo linear simples com quinze observações, ou seja, $k = 1$ e $n = 15$. Para este cenário, W apresentou maior precisão em $\rho^2 = 0$. Entretanto, este perdeu precisão a medida em que ρ^2 aproximou-se de 1, reduzindo o nível de informação. Os demais estimadores apresentaram baixa precisão em $\rho^2 = 0$. Porém, houve um aumento da precisão a medida em que o valor de ρ^2 cresceu. Na região de $\rho^2 > 0,8$, o estimador E^* apresentou a maior precisão quando comparado aos demais.

Figura 15 – Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 15$

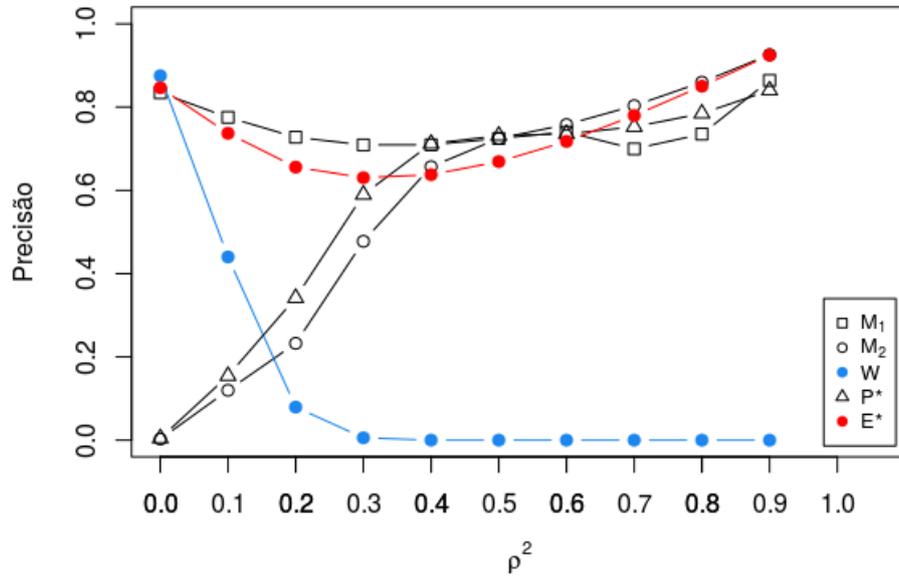


Fonte: Do autor.

Elevando o tamanho amostral do modelo linear simples para cinquenta observações, ou seja, $k = 1$ e $n = 50$, verificou-se na Figura 16 que o nível de informação dos intervalos gerados em W é comprimido na região de ρ^2 superior a 0,2. Em $\rho^2 = 0$, ocorreu uma melhoria na precisão dos estimadores M_2 e E^* , aproximando da precisão de W .

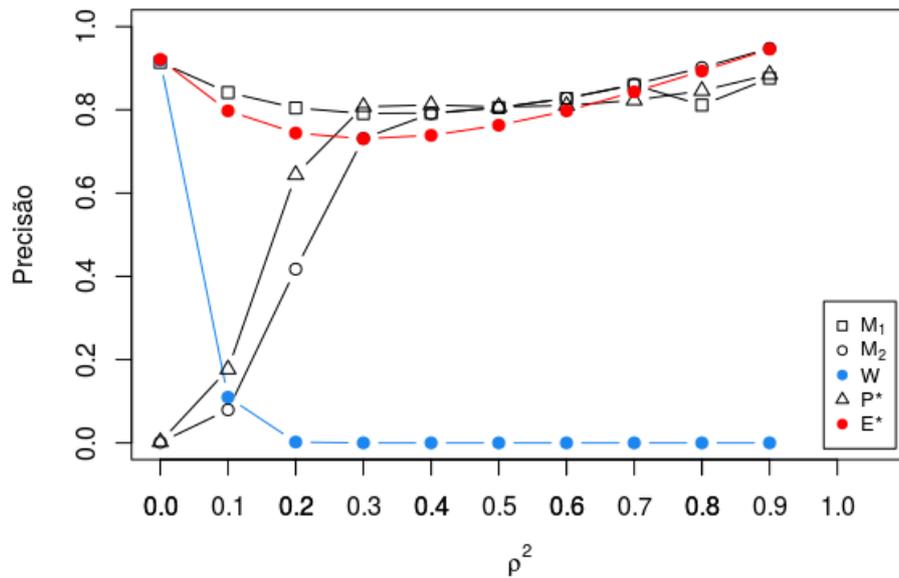
Expandindo o número de observações do modelo linear simples para cem, ou seja $k = 1$ e $n = 100$, tem-se que W apresentou uma precisão superior a 0,8 em $\rho^2 = 0$. Os estimadores E^* e M_1 aproximaram de W em $\rho^2 = 0$, e foram os mais bem comportados em todo o espaço. Além disso, P^* e M_2 perderam precisão no cenário superior a 0,7, o que pode ser verificado na Figura 17.

Figura 16 – Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 50$



Fonte: Do autor.

Figura 17 – Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 100$

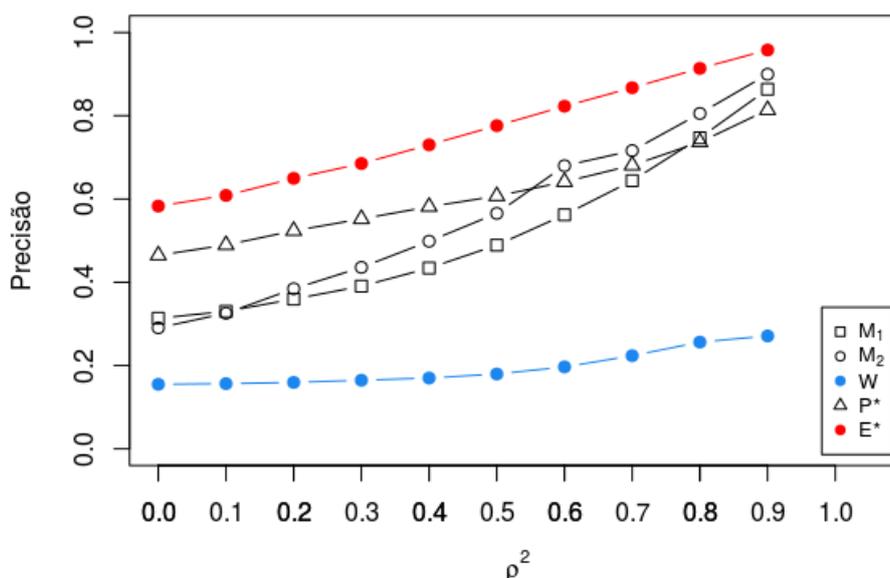


Fonte: Do autor.

Estendendo o grau do modelo para oito covariáveis e quinze observações, verificou-se por meio da Figura 18 que a precisão de W foi inferior a 0,3 em todo o espaço. Os demais estimadores aumentaram a precisão dos intervalos a medida em que o valor do coeficiente au-

mentou. Destaca-se neste cenário o comportamento do estimador E^* , com maior precisão em todo o espaço.

Figura 18 – Precisão dos intervalos referente ao modelo de regressão onde $k = 8$ e $n = 15$



Fonte: Do autor.

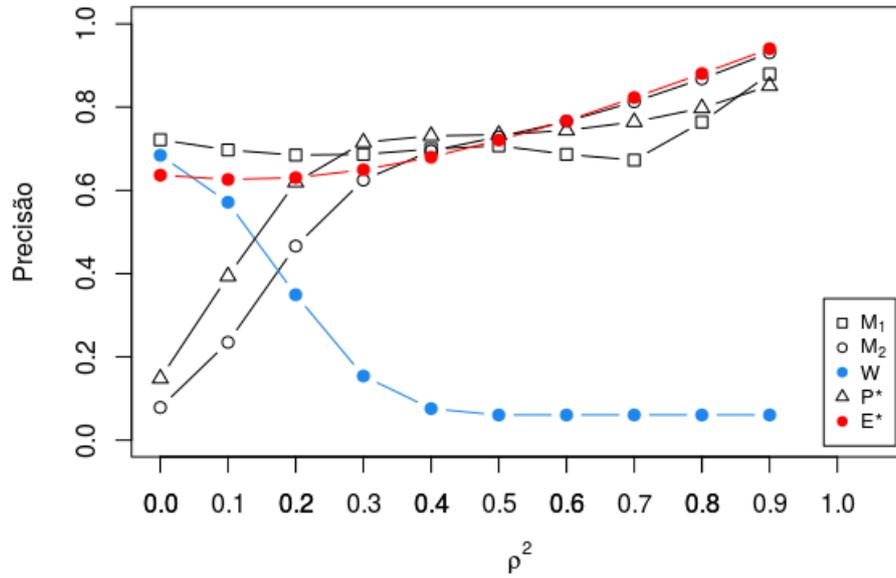
Com o aumento do número de observações do modelo com oito covariáveis para cinquenta, houve melhoria na precisão dos intervalos gerados pelos estimadores E^* e M_1 , apresentados na Figura 19. Nesse cenário, W reduziu sua precisão após $\rho^2 = 0$, atingindo precisão inferior a 20%.

A estabilidade do comprimento médio dos intervalos construídos pelos estimadores E^* e M_1 com o aumento do tamanho amostral foi verificada novamente com a ampliação do número de observações para cem no modelo com oito covariáveis, apresentados na Figura 20. W novamente perdeu o nível de informação a medida em que o valor de ρ^2 aproximou-se de 1.

Considerando os cenários analisados, pode-se notar que a boa acurácia não implica diretamente em um bom nível de informação. Esse fato pode ser exemplificado pelo estimador W , que apresentou um elevado nível de acurácia mas com pouca precisão em certos cenários, sendo plausível apenas na região em que $\rho^2 = 0$. Os estimadores E^* e M_1 apresentaram uma precisão mais equilibrada com o aumento do tamanho amostral, fato que também pode ser verificado nos demais cenários expressos no Apêndice 5.

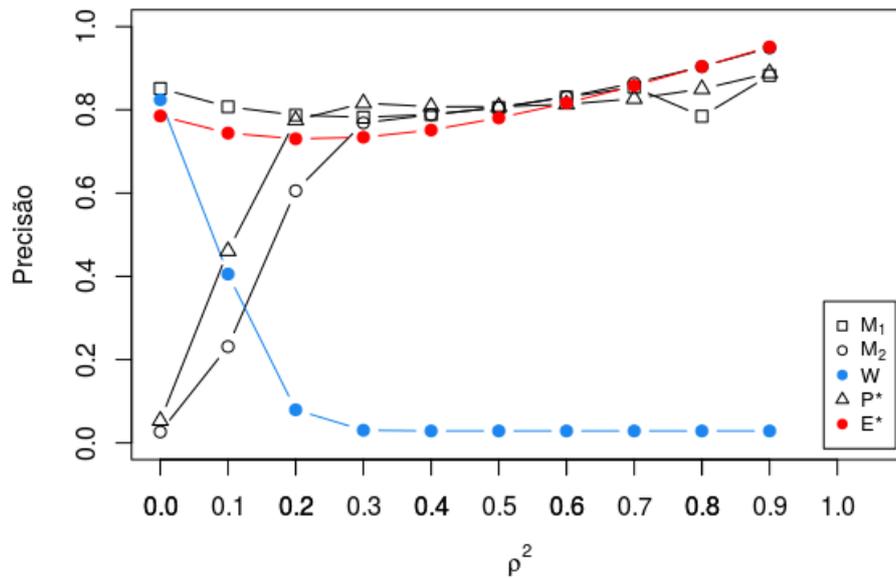
Com essa problemática, a proposta do índice de desempenho de estimação intervalar

Figura 19 – Precisão dos intervalos referente ao modelo de regressão onde $k = 8$ e $n = 50$



Fonte: Do autor.

Figura 20 – Precisão dos intervalos referente ao modelo de regressão onde $k = 8$ e $n = 100$



Fonte: Do autor.

pode ser útil para a análise conjunta da acurácia e da precisão em cada cenário. Essa situação possibilitou a escolha de um estimador com maior informação dentre os demais. A Seção 3.1.3 apresenta os índices computados para os cenários anteriores e dispõe dos demais cenários

avaliados no Apêndice 5.

3.1.3 Índices de desempenho de estimação intervalar

Como apontado na seção 2, o índice atuou como uma ponderação entre a precisão e a taxa de acurácia em cada cenário avaliado, possibilitando a escolha de um estimador mais adequado.

Iniciando novamente com o modelo de regressão linear simples com quinze observações, verificou-se na Figura 21 que W possuiu um maior nível de qualidade dentre os demais estimadores no cenário de $\rho^2 \in \{0; 0,1; 0,2\}$. Porém este perdeu qualidade a medida em que o valor de ρ^2 aumentou. Os estimadores E^* e M_1 possuíram baixa qualidade na região de $\rho^2 = 0$, mas ganharam qualidade na região superior a 0,1. Em $\rho^2 > 0,7$, os estimadores E^* e M_2 apresentaram melhor qualidade, com índices muito próximos, sendo indicados nesse cenário.

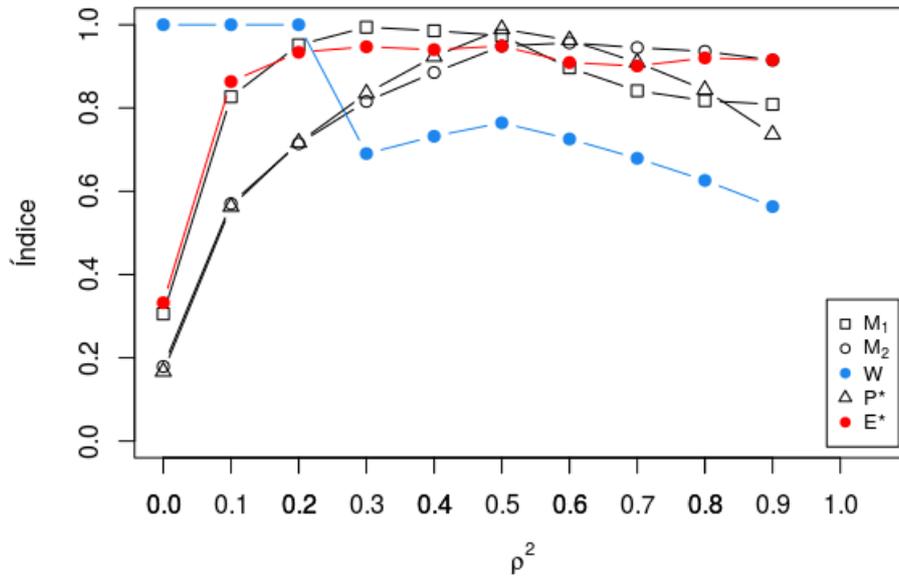
Apesar disso, notou-se também que nenhum estimador analisado recebeu um índice de desempenho máximo na região de $\rho^2 \geq 0,3$. Isso indicou que o estimador que obteve a melhor qualidade dentre os demais não possuiu a melhor taxa de acurácia e o menor comprimento médio.

Com o aumento do tamanho amostral do modelo linear simples para cinquenta observações, W perdeu qualidade conforme ρ^2 aumentou. Os estimadores M_2 e E^* apresentaram comportamento similar ao anteriormente relatado, melhorando a qualidade na região de elevado coeficiente. Os estimadores P^* e M_1 ganharam qualidade até a região de $\rho^2 = 0,5$ e foram penalizados na região de elevado coeficiente, expresso na Figura 22. Nesse cenário, somente W recebeu um índice máximo em $\rho^2 = 0$, ou seja, nas demais regiões os estimadores de melhor qualidade não obtiveram a melhor acurácia e o menor comprimento médio.

Com o aumento da informação dado pelo tamanho amostral para cem observações, o comportamento anterior foi reforçado. Os estimadores E^* e M_2 ganharam qualidade na região de elevado coeficiente, mas apresentaram baixa qualidade na região oposta. Entretanto, o estimador W exprimiu indícios de utilidade no cenário onde $\rho^2 = 0$. Ressalta-se que a perda de qualidade nos intervalos M_1 e P^* na região de $\rho^2 > 0,7$ pode não ser interessante, dado que muitas vezes objetiva-se realizar intervalos nesta região, fato ilustrado na Figura 23.

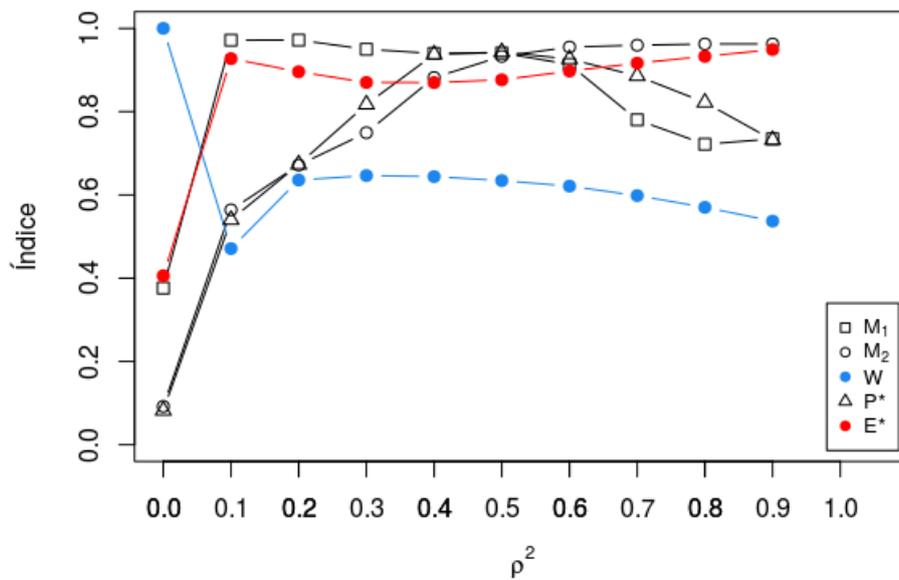
Elevando o número de covariáveis do modelo para $k = 8$ com um tamanho amostral

Figura 21 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 1$ e $n = 15$



Fonte: Do autor.

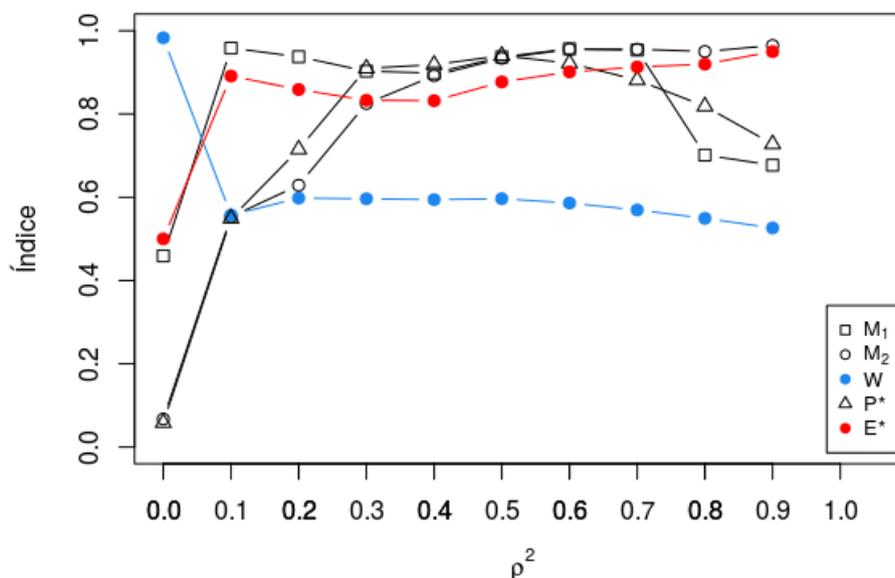
Figura 22 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 1$ e $n = 50$



Fonte: Do autor.

pequeno de quinze observações, não muito usual na literatura pelo baixo grau de liberdade para a estimação dos parâmetros, obteve-se os seguintes índices apresentados na Figura 24. Nesse

Figura 23 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 1$ e $n = 100$



Fonte: Do autor.

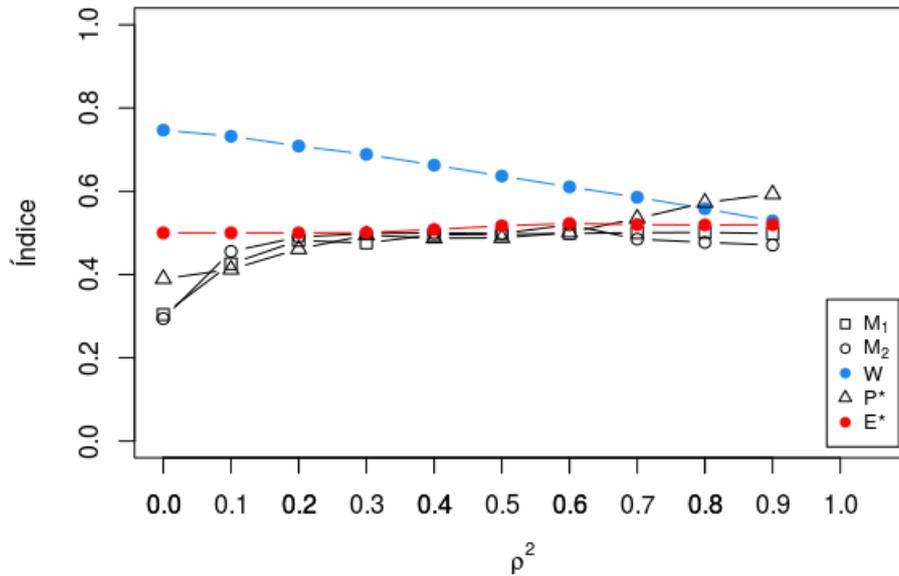
cenário, W perdeu novamente a qualidade conforme ρ^2 aumentou e os demais estimadores possuíram comportamento similar. Além disso, os índices foram inferiores a 0,8 e nenhum estimador recebeu um índice máximo.

Mantendo fixo o número de covariáveis em oito e aumentando o tamanho amostral n para 50, constatou-se por meio da Figura 25 que houve uma melhoria da qualidade dos estimadores M_1 , M_2 , P^* e E^* em relação ao cenário tratado anteriormente. Na região de $\rho^2 = 0,8$ e $\rho^2 = 0,9$, o estimador M_2 superou os demais, com maior taxa de acurácia e menor média de comprimento dos intervalos.

Considerando o aumento do tamanho amostral para $n = 100$ e $k = 8$, novamente notou-se a melhora na qualidade do estimador E^* , apresentado na Figura 26. O estimador M_2 possuiu da maior qualidade nos cenários de $\rho^2 > 0,7$, comportamento contrário ao de M_1 e P^* , que perderam qualidade nessa região. Ressalta-se que os demais cenários simulados estão disponíveis no Apêndice 5. Esses resultados mostraram que os estimadores M_2 e E^* melhoraram a qualidade a medida em que aumentou-se o tamanho amostral na região de ρ^2 superior a 0,7.

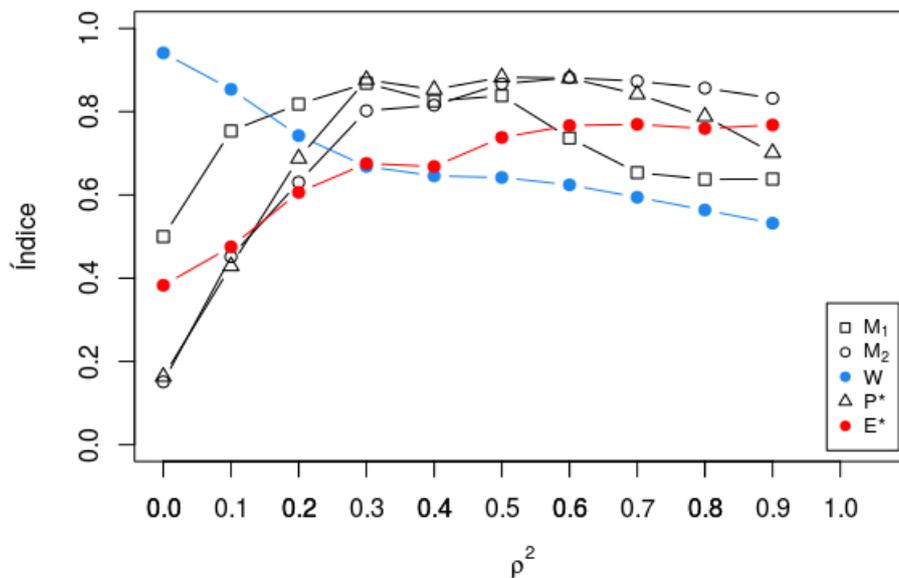
Com as medidas adotadas por este trabalho para a análise da qualidade dos estimadores intervalares, recomenda-se a utilização do estimador W para o cenário de nulidade do coeficiente de determinação, ou seja, $\rho^2 = 0$. Essa recomendação foi possível pelo fato de que este

Figura 24 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 8$ e $n = 15$



Fonte: Do autor.

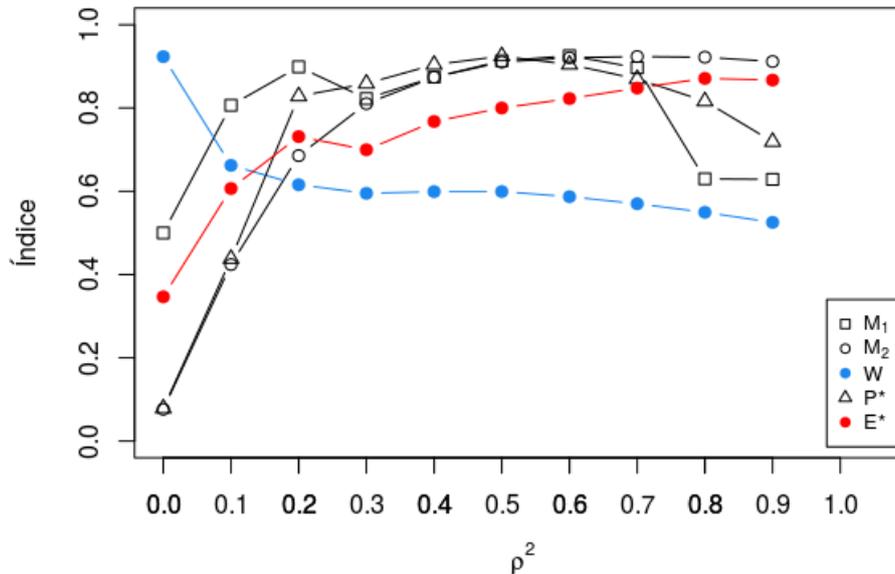
Figura 25 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 8$ e $n = 50$



Fonte: Do autor.

apresentou o melhor índice de desempenho em todos os cenários avaliados, com as maiores taxas de acurácia e maior precisão dos intervalos gerados em 54,16% das vezes. O desempenho

Figura 26 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 8$ e $n = 100$



Fonte: Do autor.

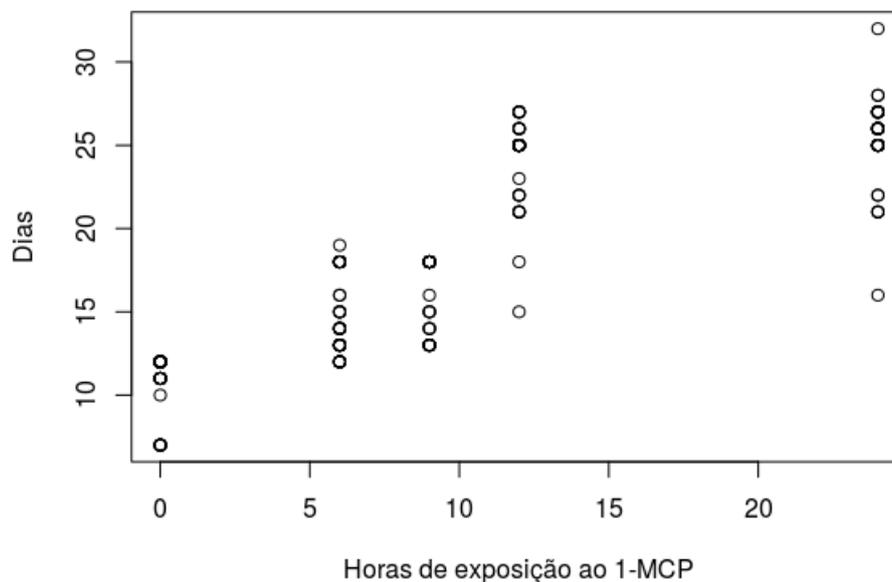
deste estimador em $\rho^2 = 0$ pode ser esperado na literatura, dado que este foi proposto sob tal condição.

Para o espaço de $0,1 \leq \rho^2 \leq 0,5$, o estimador M_1 apresentou um comportamento satisfatório, com índices que ampliaram-se com o aumento do tamanho amostral. Como exemplo deste fato, em tamanhos amostrais iguais a 100, este estimador obteve senão o melhor, índice superior a 0,8 em relação ao estimador recomendado do cenário. Tal circunstância leva a recomendação deste.

Na região de $0,5 < \rho^2 \leq 0,9$, recomenda-se a construção de intervalos de confiança pelos estimadores M_2 ou E^* . Estes apresentaram índices de qualidade que desenvolveram-se com o aumento do tamanho amostral, principalmente na região onde $\rho^2 \geq 0,7$, fato que não ocorreu com os estimadores P^* e M_1 . É válido ressaltar que em grande parte dos cenários analisados, nenhum estimador recebeu um índice de desempenho máximo. As exceções deste contexto foram dadas apenas em $\rho^2 = 0$.

Tendo em vista a análise da qualidade de ajuste do experimento relacionado ao prolongamento de vida pós colheita de bananas, observou-se o seguinte comportamento da variável Dias G7, apresentado na Figura 27.

Figura 27 – Dias transcorridos até o alcance máximo de amadurecimento em relação ao período de exposição ao 1-MCP



Fonte: Do autor.

Note na Figura 27 que os frutos expostos ao 1-MCP por 0 horas atingiram o grau máximo de amadurecimento em uma média de 10,9 dias. Por outro lado, os frutos expostos ao 1-MCP por 24 horas atingiram amadurecimento máximo em uma média de 25,35 dias.

Dada a situação, dois modelos de regressão foram ajustados com o objetivo de analisar o efeito do tempo de exposição ao 1-MCP no número de dias para o alcance máximo de amadurecimento. O modelo linear foi o primeiro a ser considerado, com resultados apresentados na Tabela 1.

Tabela 1 – Coeficientes estimados para a análise do efeito de exposição ao 1-MCP no número de dias transcorridos para o alcance máximo de amadurecimento das bananas

Preditor	Estimativa	Erro Padrão	$P_r > t $
Intercepto	11,80	0,3822	< 0,01
1-MCP	0,637	0,0290	< 0,01

Fonte: Do autor.

Analisando os resultados apresentados na Tabela 1, observa-se que a exposição ao 1-MCP tem efeito positivo e significativo no número de dias, ou seja, bananas expostas ao 1-MCP

tendem a apresentar maior vida de prateleira quando comparadas as não expostas. O modelo linear apresentou um $R^2 = 0,6981$, indicando que o modelo explica aproximadamente 69,8% da variação da variável resposta.

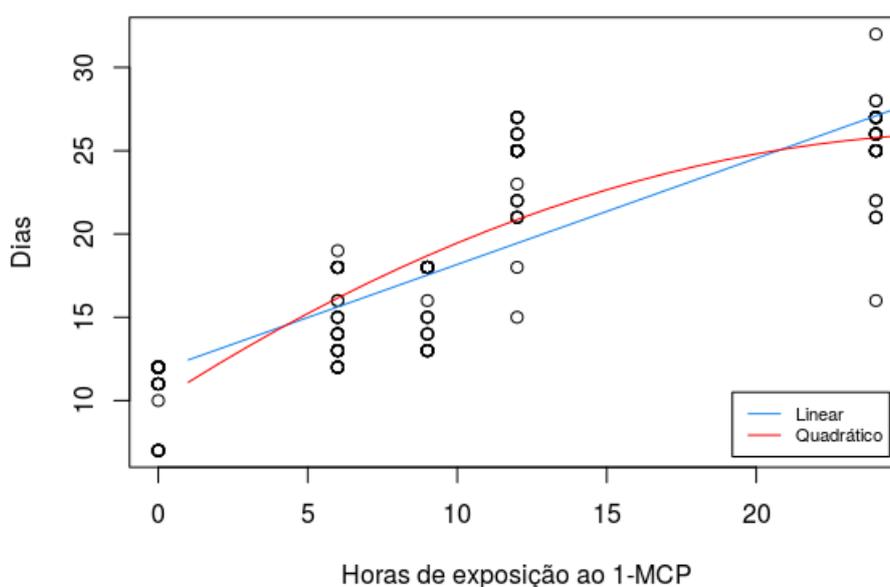
Posteriormente, o modelo quadrático foi o segundo a ser considerado, analisando o efeito do tempo de exposição e do tempo de exposição ao quadrado no número de dias para o alcance do amadurecimento máximo, com estimativas apresentadas na Tabela 2. O modelo quadrático apresentou de um $R^2 = 0,7461$. Como forma de ilustração, a Figura 28 denota o comportamento dos dois modelos considerados.

Tabela 2 – Coeficientes estimados para o modelo quadrático, analisando o efeito de exposição ao 1-MCP no número de dias transcorridos para o alcance máximo de amadurecimento de bananas

Preditor	Estimativa	Erro Padrão	$P_r > t $
Intercepto	9,957	0,463	< 0,01
1-MCP	1,158	0,089	< 0,01
1-MCP ²	-0,020	0,003	< 0,01

Fonte: Do autor.

Figura 28 – Modelos ajustados ao experimento relativo ao amadurecimento de bananas do tipo maçã



Fonte: Do autor.

Após o cálculo de R^2 nos modelos, os seguintes cenários foram simulados: i) $k = 1$, $n = 200$ e $\rho^2 = 0,6981$ e ii) $k = 2$, $n = 200$ e $\rho^2 = 0,7461$, que correspondem ao modelo linear e quadrático, respectivamente. Considerando ambos cenários, tem-se que M_2 com o melhor desempenho conforme o índice $(\tau_{\alpha,i})$, com os seguintes limites de confiança para o modelo linear:

$$IC(\rho^2, 95\%) = [0,6473 - 0,7465],$$

e para o modelo quadrático:

$$IC(\rho^2, 95\%) = [0,7016 - 0,7881].$$

A partir da construção dos intervalos para a qualidade de ajuste, pode-se notar que estes apresentaram interseção, não admitindo evidências de que os modelos tenham qualidades diferentes. Desta forma, pode-se optar pelo modelo mais parcimonioso, ou seja, o linear, ou considerar mais um critério disponível na literatura que possibilite a tomada de decisão. É válido ressaltar que o estimador utilizado para a construção dos intervalos em ambos cenários também foi o recomendado pela simulação realizada neste trabalho.

4 CONCLUSÃO

Considerando a questão de que a qualidade de ajuste de modelos é um dos grandes pontos da modelagem estatística, pode-se notar que mesmo com a popularidade do uso do R^2 , esta apresenta problemas e recomendações que devem ser considerados. A questão inferencial possibilita a tomada de decisão acerca da qualidade do modelo utilizado, a partir do momento em que trata-se R^2 como uma estatística que estima o parâmetro ρ^2 .

A literatura estatística oferece dois estimadores paramétricos que foram explorados em duas regiões do espaço de ρ^2 de forma complementar. Fato que possibilita a exploração do comportamento destes estimadores sob todo o espaço paramétrico e estimula o inserção de novos estimadores à literatura.

Considerando os estimadores W e E^* , notou-se que W apresentou o melhor desempenho na região na qual foi proposto mas perdeu qualidade do espaço paramétrico complementar. Por outro lado E^* possuiu um comportamento satisfatório, mas não apresentou a melhor qualidade em todo o espaço no qual foi construído.

Em relação aos estimadores propostos, notou-se que a reparametrização de um dos estimadores possibilitou a melhoria da acurácia e da respectiva qualidade dos intervalos construídos em alguns cenários. Como exemplo, tem-se o estimador M_2 , que apresentou uma qualidade razoável, aproximando-se da qualidade de E^* , com uma forma algébrica mais simples.

No processo de análise de todos os estimadores, não houve um predomínio da qualidade de um estimador em todo o espaço nos cenários simulados. Fato que incentivou a recomendar o de melhor qualidade em três regiões do espaço, sob a condição de que o tamanho amostral do modelo seja elevado. A importância de um tamanho amostral suficientemente grande identificada, retoma aos resultados analisados por Cramer (1987) da necessidade de pelo menos cinquenta observações para o cálculo do coeficiente, melhorando a taxa de acurácia e a precisão dos intervalos construídos.

Por fim, a utilização dos estimadores em um caso aplicado a um experimento apontou que a qualidade do modelo pode ser analisada por meio de um intervalo de confiança para ρ^2 , dado que estes comumente compreendem de um reduzido tamanho amostral. Como consequência, foi possível também a seleção de um modelo mais parcimonioso.

REFERÊNCIAS

BINEESH, K. et al. Length–weight relationships of eight deep-sea fish species collected from the southwest coast of India. **Journal of applied ichthyology**, v. 34, n. 5, p. 1220-1222, 2018. Disponível em: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/jai.13745>. Acesso em: 22 nov. 2018.

CRAMER, J.S. Mean and variance of R^2 in small and moderate samples. **Journal of econometrics**, v. 35, ed. 2-3, p. 253-266, 1987. Disponível em: <https://www.sciencedirect.com/science/article/pii/0304407687900273>. Acesso em: 21 jun. 2018.

DRAPER, N; SMITH, H. **Applied regression analysis**. 3. ed. rev. New York: John Wiley Sons, 1998. v. 326.

FOSTER, D; SMITH, T; WHALEY, R. Assessing goodness-of-fit of asset pricing models: The distribution of the maximal R^2 . **The journal of finance**, v. 52, n. 2, p. 591-607, 1997. Disponível em: <https://www.jstor.org/stable/2329491>. Acesso em: 23 jun. 2018.

GRANATO, D; CALADO, V. The use and importance of design of experiments (DOE) in process modelling in food science and technology. In: ARES, Gastón. **Mathematical and statistical methods in food science and technology**. 3. ed. New York: John Wiley & Sons, 1998.

NAKAGAWA, S; JOHNSON, PCD; SCHIELZETH, H. The coefficient of determination R^2 and intra-class correlation coefficient from generalized linear mixed-effects models revisited and expanded. **Royal society**, v. 14, n. 134, 2017. Disponível em: https://royalsocietypublishing.org/doi/full/10.1098/rsif.2017.0213?url_ver=Z39.88-2003rfr_id=ori:rid:crossref.org:rfr_dat=cr_pub%3dpubmed. Acesso em: 9 dez. 2018.

OHTANI, K; TANIZAKI, H. Exact distributions of R^2 and adjusted R^2 in a linear regression model with multivariate t error terms. **Journal of the Japan statistical society**, v. 34, n. 1, p. 101-109, 2004. Disponível em: <http://www2.econ.osaka-u.ac.jp/tanizaki/cv/papers/distr2.pdf>. Acesso em: 4 jun. 2018.

OHTANI, K. The density functions of R^2 and, and their risk performance under asymmetric loss in misspecified linear regression models. **Economic modelling**, v. 11, n. 4, 1994. Disponível em: <https://www.sciencedirect.com/science/article/pii/0264999394900035>. Acesso em: 4 jun. 2018.

OWEN, C. **Parameter estimation for the beta distribution**. 2008. 96 f. Dissertação (Mestrado) - Curso de Master Of Science, Brigham Young University, Provo, 2008. Disponível em: <https://scholarsarchive.byu.edu/cgi/viewcontent.cgi?article=2613context=etd>. Acesso em: 28 jan. 2018.

PIEPHO, H. A coefficient of determination (R^2) for generalized linear mixed models. **Biometrical journal**, v. 61, n. 4, p. 860-872, 2019. Disponível em:

<https://onlinelibrary.wiley.com/doi/10.1002/bimj.201800270>. Acesso em: 14 jan. 2019.

QUININO, R.; REIS, E.; BESSEGATO, L. Using the coefficient of determination R^2 to test the significance of multiple linear regression. **Teaching statistics**, v. 35, n. 2, p. 84-88, 2012. Disponível em: <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-9639.2012.00525.x>. Acesso em: 19 ago. 2018.

SONG, Z; DENG, Q; REN, Z. Correlation and principal component regression analysis for studying air quality and meteorological elements in Wuhan, China. **Environmental progress & sustainable energy**, 2019. Disponível em: <https://aiche.onlinelibrary.wiley.com/doi/abs/10.1002/ep.13278>. Acesso em: 19 jan. 2019.

WEATHERBURN, C. **A first course in mathematical statistics**. Cambridge: University Press, 1949. v. 32.

YOGESHA, S. et al. Simultaneous quantification of n-butylthiophosphoric triamide and dicyandiamide in urea formulation by liquid chromatography-tandem mass spectrometry. **Journal of separation science**, v. 42, n. 2, p. 484-490, 2019. Disponível em: <https://www.ncbi.nlm.nih.gov/pubmed/30450719>. Acesso em: 17 fev. 2019.

ZHANG, D. A Coefficient of determination for generalized linear models. **The American statistician**, v. 71, n. 4, p. 310-316, 2017. Disponível em: <https://amstat.tandfonline.com/doi/abs/10.1080/00031305.2016.1256839.XUiHCstKicw>. Acesso em: 23 fev. 2019.

CAPÍTULO 3 - O PACOTE ICR2: INTERVALOS DE CONFIANÇA PARA O COEFICIENTE DE DETERMINAÇÃO

Resumo

Com a extensa utilização de medidas de qualidade de ajuste em regressão, nota-se que o coeficiente de determinação (R^2) tem sido utilizado em grande parte dos estudos aplicados. Estudos estatísticos apontaram que este possui algumas fragilidades, mas possibilita uma análise e interpretação da qualidade do modelo proposto. Dado esse contexto, estudos buscam a realização de inferências acerca da qualidade populacional de um modelo (ρ^2) que é estimada pelo coeficiente de determinação amostral (R^2). Na literatura, há diferentes estimadores que podem ser utilizados para a construção de intervalos de confiança para ρ^2 , com diferentes recomendações em regiões do espaço paramétrico. Com isso, este trabalho tem por objetivo propor e apresentar o pacote ICR2, desenvolvido em linguagem R, que possibilita a construção de intervalos de confiança para a qualidade de ajuste de modelos de regressão que detenham de melhor qualidade de estimação para o cenário do usuário. Como forma de ilustração, será apresentado e discutido um exemplo analisando o efeito da idade de matrizes no crescimento de progênes durante os sete dias iniciais após o nascimento.

Palavras-Chave: Estimação Intervalar; R; Regressão.

1 INTRODUÇÃO

A extensa utilização do coeficiente de determinação (R^2) para a análise de qualidade de ajuste de modelos de regressão possibilitou a exploração de suas principais características e problemas pela literatura estatística, como visto em Bucchianico (2014) e Quinino Reis e Bessegato (2012).

Esse contexto possibilitou a inserção do coeficiente de determinação ajustado R_a^2 , atuando como um penalizador de R^2 conforme o número de covariáveis compreendidas pelo modelo. Essa penalização pode fazer com que o R_a^2 exceda o espaço paramétrico de R^2 , perdendo a facilidade de interpretação como a proporção da variabilidade da variável resposta explicada pelo modelo. No entanto, R_a^2 também é aceito e empregado pela literatura científica.

Mesmo com as particularidades, estudos acerca de R^2 buscam a ampliação deste para a análise de qualidade de modelos lineares generalizados e modelos mistos em Piepho (2019) e Zhang (2017). E, por outro lado, estudos destinam-se a realizar inferências acerca da qualidade do modelo, tratando R^2 como um estimador pontual de ρ^2 , baseado na distribuição Beta, fazendo com que ρ^2 seja um parâmetro respectivo a qualidade populacional do modelo.

A forma inferencial em respeito a ρ^2 é disposta por dois estimadores que atuam de forma complementar no espaço paramétrico de ρ^2 , isto é, $\rho^2 = 0$ e $\rho^2 \neq 0$, W e E^* , respectivamente. Esses dois estimadores foram propostos em Weatherburn (1949) e Cramer (1987), mas podem gerar intervalos largos e com baixa acurácia nos cenários que ultrapassem a região onde esses estimadores foram propostos.

Essa perspectiva proporciona a implantação de novos estimadores que possam aprimorar a decisão acerca da qualidade de um modelo. Com isso, três estimadores foram propostos no Capítulo 2, baseados na moda (M_1), média (M_2) e esperança da soma de quadrados da regressão (P^*). Posteriormente, os cinco estimadores foram avaliados de forma computacional e comparados quanto ao desempenho em acurácia e precisão. Entretanto, os resultados não possibilitaram evidenciar um estimador com melhor desempenho em todo o espaço paramétrico. De modo que cenários com diferentes combinações do número de covariáveis (k) e tamanho amostral (n) podem apresentar diferentes resultados dos explorados na simulação.

Para englobar essas particularidades, este trabalho tem por objetivo propor e apresentar o pacote ICR2, que possibilita a construção de intervalos de confiança para ρ^2 a partir das recomendações apontadas no capítulo anterior ou por meio de simulação do respectivo cenário

do usuário. A implementação do pacote será realizada em linguagem *R*, que é um ambiente livre para computação estatística (R CORE TEAM, 2019).

Na literatura, implementações de pacotes são comuns em *R* e exploram a utilização de exemplos e possibilitam a comunicabilidade entre desenvolvedores e pesquisadores. Como exemplo destes pacotes, tem-se o estudo de Sturtz, Ligges e Gelman (2005), apresentando o pacote *R2WinBUGS* destinado à análise Bayesiana e o estudo de Kuznetsova, Brockhoff, Christensen (2017) apresentando o pacote *lmerTest*.

2 EXEMPLO

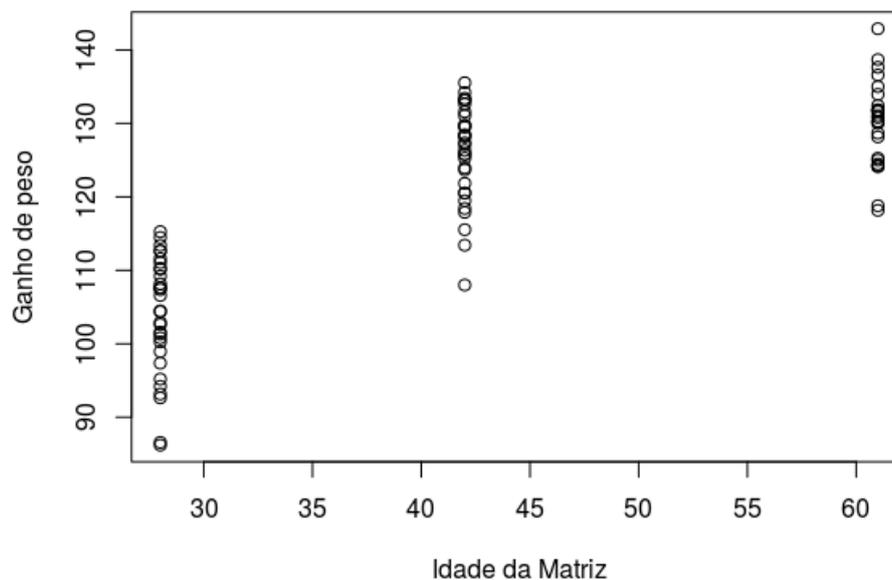
No contexto de melhoramento animal, o Brasil desenvolve pesquisas que possibilitam a avaliação genética e a seleção de animais que visem o aumento da produtividade. Essas pesquisas comumente baseiam-se em ajustes de modelos estatísticos e na aplicação da biologia molecular (LOBO, R; BITTNECOURT, T; BATISTA, L; 2010). Como exemplo, pode-se notar a crescente produção de aves no país, abrindo demanda por materiais genéticos de elevada qualidade, principalmente de aves poedeiras (ALMEIDA E SILVA, 2009).

Dada a situação, o exemplo a ser ilustrado corresponde a um experimento utilizando matrizes de galinhas. Foram utilizadas matrizes com idade de 28, 42 e 61 semanas, sendo 31 com 28 semanas, 31 com 42 semanas e 26 com 61 semanas. Como objetivo, analisou-se o desenvolvimento da progênie, através do ganho de peso em gramas, durante o período de 1 a 7 dias a partir do nascimento. O experimento contou com 88 observações, sendo desenvolvido por Moreno (2019).

A partir da realização do experimento, verificou-se o seguinte comportamento entre a idade das matrizes e o ganho de peso da progênie pela Figura 29. Note que os filhotes provenientes de matrizes com 28 semanas alcançaram, em média, 103,89g. Para as matrizes com 42 semanas, a média de peso dos filhotes foi 125,79g e, por fim, os filhotes de matrizes com 61 semanas apresentaram média de 129,9g no período.

Uma das hipóteses do estudo é analisar o efeito da idade das matrizes no ganho de peso da progênie durante o período analisado, ou seja, 1 a 7 dias após o nascimento. Fato que possibilita a seleção de um intervalo de idade que gere, em média, progênies que detenham de um maior ganho de peso nos primeiros dias de vida. Maiores detalhes podem ser vistos em Moreno (2019).

Figura 29 – Comportamento entre a idade das matrizes e o ganho de peso da progênie no período de 1 a 7 dias após o nascimento



Fonte: Do autor.

3 IMPLEMENTAÇÃO

Nesta seção estão apresentadas as distribuições dos estimadores a serem utilizadas pelo pacote e as respectivas funções usuais para a modelagem, como a geração de valores aleatórios, abordadas na subseção 3.1. A seção 3.2 apresenta as funções criadas para a construção de intervalos de confiança utilizando as distribuições dos estimadores tratados em 3.1.

3.1 DISTRIBUIÇÕES DOS ESTIMADORES

3.1.1 Distribuição de W

Tratado por Weatherburn (1949) para o cenário onde $\rho^2 = 0$, implicando na condição de que o vetor de parâmetros associado às covariáveis da regressão é nulo. Conforme o estudo, o

estimador tem distribuição Beta, dada da seguinte forma:

$$f(x; n, k) = \frac{1}{B\left(\frac{k}{2}, \frac{n-k-1}{2}\right)} x^{\frac{k-2}{2}} (1-x)^{\frac{n-k-3}{2}} I_{(0,1)}(x).$$

A distribuição de W pode ser acessada e é denominada como `distribuiçãoW`. As especificações e argumentos podem ser acessadas com o comando ajuda do *software*. A função `distribuiçãoW` contém as seguintes subfunções apresentadas na Tabela 3.

Tabela 3 – Subfunções disponíveis no pacote ICR2 para a distribuição do estimador W

Nome	Uso	Retorno
dW	dW(x, n, k)	O valor da densidade no ponto x.
qW	qW(p, n, k)	O quantil associado a probabilidade p.
rW	rW(m, n, k)	Gera m valores aleatórios da distribuição.
pW	pW(q, n, k)	O valor da integral de 0 a q.

Fonte: Do autor.

3.1.2 Distribuição de E^*

Corresponde a distribuição do estimador E^* apresentada por Cramer (1987), para o cenário onde $\rho^2 \neq 0$, dada por:

$$f(r) = \sum_{j=0}^{\infty} W(j) \frac{1}{B(u+j, v-u)} r^{u+j-1} (1-r)^{v-u-1} I_{(0,1)}(r).$$

em que:

$$W(j) = \frac{e^{-\frac{\lambda}{2}} \left(\frac{\lambda}{2}\right)^j}{j!}, \quad u = \frac{1}{2}(k), \quad v = \frac{1}{2}(n-1).$$

Atribuída como `distribuiçãoE`, a função contém as seguintes subfunções apresentadas na Tabela 4.

Tabela 4 – Subfunções disponíveis no pacote ICR2 para a distribuição do estimador E^*

Nome	Uso	Retorno
dE	dE(r, rho2, n, k)	O valor da densidade no ponto r.
qE	qE(x, rho2, n, k)	O quantil associado a prob. x.
rE	rE(m, rho2, n, k)	Gera m valores aleatórios.
pE	pE(r, rho2, n, k)	O valor da integral de 0 a r.

Fonte: Do autor.

3.1.3 Distribuição de M_1

O estimador M_1 foi construído baseado na distribuição do estimador W , através de uma reparametrização pela moda da distribuição Beta. Maiores detalhes podem ser vistos no capítulo anterior. Com isso, M_1 possui a seguinte distribuição:

$$f(x; n, k, \rho^2) = \frac{1}{B\left(\frac{\rho^2(n-k-5)+2}{2(1-\rho^2)}, \frac{n-k-1}{2}\right)} x^{\frac{\rho^2(n-k-5)+2}{2(1-\rho^2)}} (1-x)^{\frac{n-k-3}{2}} \mathbf{I}_{(0,1)}(x).$$

Denominada como `distribuiçãoM1`, a função contém as seguintes subfunções apresentadas na Tabela 5.

Tabela 5 – Subfunções disponíveis no pacote ICR2 para a distribuição do estimador M_1

Nome	Uso	Retorno
dM1	dM1 (x, n, k, rho2)	O valor da densidade no ponto x.
qM1	qM1 (p, n, k, rho2)	O quantil associado a prob. p.
rM1	rM1 (m, n, k, rho2)	Gera m valores aleatórios.
pM1	pM1 (q, n, k, rho2)	O valor da integral de 0 a q.

Fonte: Do autor.

3.1.4 Distribuição de M_2

O estimador M_2 foi baseado na distribuição do estimador W , através de uma reparametrização pela média da distribuição Beta. Com isso, a distribuição de M_2 é dada por:

$$f(x; n, k, \rho^2) = \frac{1}{B\left(\frac{-\rho^2(n-k-1)}{2(\rho^2-1)}, \frac{n-k-1}{2}\right)} x^{\frac{-\rho^2(n-k-1)}{2(\rho^2-1)}} (1-x)^{\frac{n-k-3}{2}} \mathbf{I}_{(0,1)}(x).$$

Denominada como `distribuiçãoM2`, a função contém as seguintes subfunções apresentadas na Tabela 6.

Tabela 6 – Subfunções disponíveis no pacote ICR2 para a distribuição do estimador M_2

Nome	Uso	Retorno
dM2	dM2(x, n, k, rho2)	O valor da densidade no ponto x.
qM2	qM2(p, n, k, rho2)	O quantil associado a prob. p.
rM2	rM2(m, n, k, rho2)	Gera m valores aleatórios.
pM2	pM2(q, n, k, rho2)	O valor da integral de 0 a q.

Fonte: Do autor.

3.1.5 Distribuição de P^*

O estimador P^* foi proposto a partir de relações entre a esperança das soma de quadrados da regressão e soma de quadrados totais, com distribuição dada por:

$$f(x; n, \rho^2) = \frac{1}{B(n\rho^2, n - n\rho^2)} x^{n\rho^2} (1 - x)^{n - n\rho^2 - 1} \mathbf{I}_{(0,1)}(x).$$

Denominada como distribuição $\circ P$, a função contém as seguintes subfunções apresentadas na Tabela 7.

Tabela 7 – Subfunções disponíveis no pacote ICR2 para a distribuição do estimador P^*

Nome	Uso	Retorno
dP	dP(x, n, rho2)	O valor da densidade no ponto x.
qP	qP(p, n, rho2)	O quantil associado a probabilidade p.
rP	rP(m, n, rho2)	Gera m valores aleatórios.
pP	pP(q, n, rho2)	O valor da integral de 0 a q.

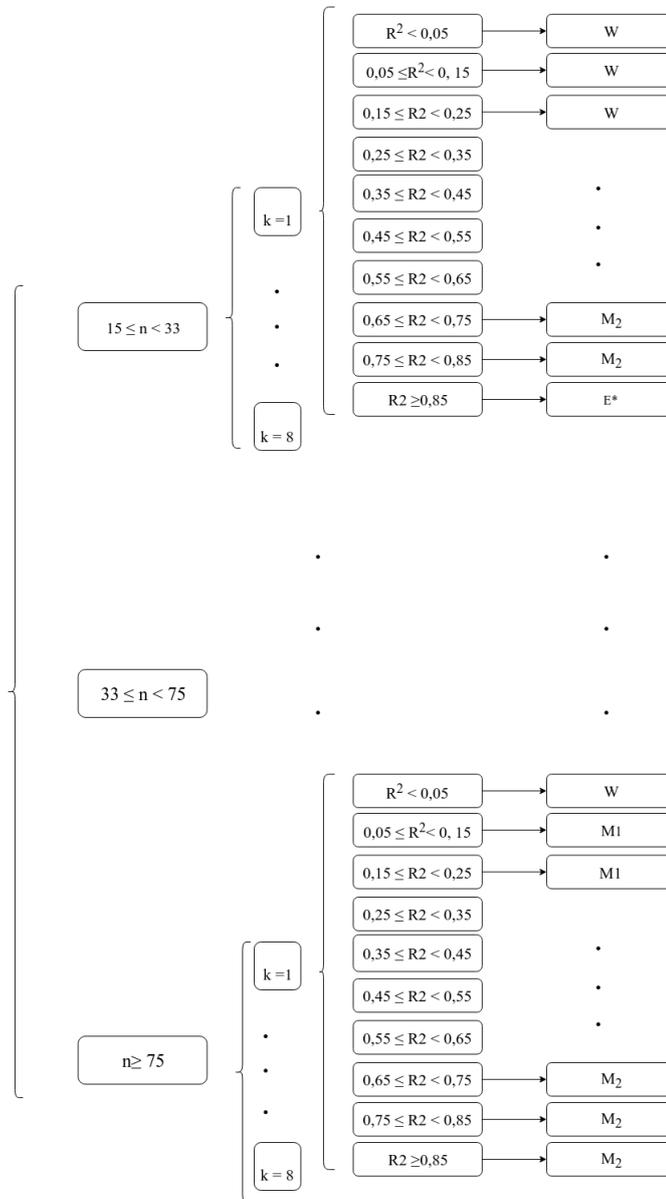
Fonte: Do autor.

3.2 FUNÇÕES PARA A CONSTRUÇÃO DE INTERVALOS DE CONFIANÇA

Como forma de disponibilizar os resultados apresentados no capítulo 2, foi implementado uma função que sugere o estimador maior índice $\tau_{\alpha,i}$ nas proximidades do cenário do usuário, ou seja, realizando uma suavização e permitindo amparar mais cenários, denotado na Figura 30.

Essa função foi denominada como `sEstimador`, e possui três argumentos, sendo eles: o tamanho amostral do modelo (n), o número de covariáveis (k) e o valor de ρ^2 , como apresentado na Tabela 8. É importante ressaltar que para a recomendação do estimador e construção do

Figura 30 – Diagrama da função que recomenda o estimador a ser utilizado a partir de uma suavização dos cenários simulados, considerando o maior índice $\tau_{\alpha,i}$



Fonte: Do autor.

respectivo intervalo, utiliza-se $\alpha = 0,05$ e $p_1 = 0,5$. Além disso, a função admite que o comportamento verificado em modelos com até oito covariáveis ocorrerá em modelos onde $k > 8$.

Caso o usuário opte por outro nível de confiança ou deseje realizar sua própria simulação, também é disponibilizada uma função denominada $sICR2$. Essa função computa a taxa de acurácia, o comprimento médio e o índice de desempenho de estimação intervalar no cená-

Tabela 8 – Função disponível no pacote ICR2 para a função `sEstimador`, que sugere o estimador a ser utilizado e retorna o intervalo de confiança a um nível $\alpha = 5\%$

Nome	Uso	Retorno
<code>sEstimador</code>	<code>sEstimador(n, k, rho2)</code>	Retorna o estimador e o intervalo com maior índice de qualidade na proximidade de ρ^2 , n e k .

Fonte: Do autor.

rio analisado. Para isso, o usuário deve escolher o número de simulações Monte Carlo (MC), o limite superior do somatório associado a distribuição do estimador E^* (JS), o peso atribuído à acurácia (p_1) e o o nível α adotado. O que pode ser visto na Tabela 9.

Tabela 9 – Função disponível no pacote ICR2 para a função `sICR2` para realizar a simulação Monte Carlo

Nome	Uso	Retorno
<code>sICR2</code>	<code>sICR2(rho2, n, k, MC = 1000, p1 = 0.5, alpha0 = 0.05, JS = 1000)</code>	Retorna o estimador com o melhor desempenho no cenário e o intervalo.

Fonte: Do autor.

É válido ressaltar que a função `sICR2` pode consumir um tempo computacional elevado em consequência da integração numérica da distribuição E^* . Como exemplo do tempo computacional gasto, a simulação para o cenário onde $k = 4$, $n = 61$ e $R^2 = 0,85$, durou 22,18 minutos. Entretanto, alguns cenários podem ultrapassar 50 minutos no processador Intel® Core™ i3-4005U (1.70GHz).

4 APLICAÇÃO AO EXEMPLO

Considerando o exemplo apresentado em 2, esta seção aplica as funções `sICR2` e `sEstimador` e analisa os resultados. Inicialmente, o pacote e os dados de exemplo contidos no pacote serão solicitados da seguinte forma:

```
> library(ICR2)
> experimento
> head(experimento)
```

	idade	g_peso
1	28	95.19
2	28	101.32
3	28	102.89
4	28	102.63
5	28	107.37
6	28	98.95

Para modelar os dados do experimento, serão utilizados dois modelos de regressão. O primeiro considerando a variável ganho de peso sendo explicada pela idade das matrizes e o segundo considerando o ganho de peso sendo explicado pela idade das matrizes e a idade das matrizes ao quadrado, da seguinte forma:

```
mod1 <-lm(g_peso ~ idade, data = experimento)
mod2 <-lm(g_peso ~ idade + I(idade ^ 2 ), data = experimento)
```

Após a estimação dos coeficientes, tem-se como necessidade a análise das pressuposições dos resíduos do modelo por parte do usuário. Posteriormente, pode-se consultar as estimativas do coeficiente de determinação amostral de cada modelo com o seguinte comando:

```
>summary(mod1)$r.squared
[1] 0.5889934
>summary(mod2)$r.squared
[1] 0.7400878
```

Note que o modelo quadrático apresentou um maior coeficiente de determinação quando comparado ao modelo linear, explicando 74% da variação da variável resposta. Dada a situação, os intervalos de confiança para a qualidade de ajuste serão construídos inicialmente utilizando a função que considera os resultados da simulação com um nível $\alpha = 0,05$. Para isso, é necessário conhecer o número de observações do modelo (nesse caso $n = 88$) e o número de covariáveis (k

= 1 em mod1 e $k = 2$ em mod2). Note que é utilizado a estimativa pontual R^2 em ρ^2 . Assim, considerando a função `sEstimador` para o primeiro modelo, tem-se:

```
> sEstimador(n = 88, k = 1, rho2 = 0.5889934)
```

Com a função, a seguinte tabela será apresentada conforme a Figura 31.

Figura 31 – Estimador recomendado pela função sugestão disponibilizada no pacote ICR2 para o modelo linear

```
Considerando o cenário, tem-se:
=====
Estimador  rho2    Lim_Inf Lim_Sup Confiança
-----
M1         0.5889934  0.492   0.68    0.95
-----
```

Fonte: Do autor.

Considerando o primeiro modelo, note que o intervalo de confiança foi construído com o estimador M_1 , com um comprimento de 0,1877 e confiança de 95%. Dado por:

$$IC(\rho^2, 95\%) = [0,4920 - 0,6800].$$

Para o modelo quadrático, o comando para a função é dado por:

```
> sEstimador(n = 88, k = 2, rho2 = 0.7400878)
```

retornando os seguintes resultados apresentados na Figura 32.

Figura 32 – Estimador recomendado pela função sugestão disponibilizada no pacote ICR2 para o modelo quadrático

```
Considerando o cenário, tem-se:
=====
Estimador  rho2    Lim_Inf Lim_Sup Confiança
-----
M2         0.7400878  0.67    0.804   0.95
-----
```

Fonte: Do autor.

Para esse cenário, o intervalo de confiança foi construído conforme o estimador M_2 , com comprimento de 0,1337, confiança de 95% e os seguintes limites:

$$IC(\rho^2, 95\%) = [0,6700 - 0,8040].$$

Note que os intervalos construídos interseccionam-se a um nível $\alpha = 0,05$, não havendo evidências de que um dos modelos tenha melhor qualidade.

Por outro lado, caso o usuário deseje realizar a simulação, é necessário utilizar a função `sICR2`. Para este exemplo serão considerados os argumentos conforme o recomendado pela função com o seguinte comando:

```
> sICR2(rho2 = 0.588934, n = 88, k = 1, MC = 1000, p1 = 0.5,
alpha0 = 0.05, JS = 1000)
```

A simulação será iniciada com esse comando e o processo poderá ser acompanhado na barra de progresso. Para este cenário, foi gasto um tempo computacional de 42 minutos, retornando os seguintes resultados apresentados na Figura 33.

Figura 33 – Estimador recomendado pela função simulação disponibilizada no pacote ICR2 para o modelo linear

```

Parâmetros de qualidade dos estimadores no cenário avaliado
=====

```

	M1	M2	W	P*	E*
Índice	0.955	0.953	0.593	0.927	0.896
Acurácia	0.913	0.906	1	0.938	0.944
Comprimento	0.186	0.186	1	0.203	0.219

```

-----
Intervalo considerando o estimador com maior Índice
=====

```

Estimador	rho2	Lim_Sup.	Lim_Inf.	Confiança
M1	0.588934	0.492	0.68	0.95

```

-----

```

Fonte: Do autor.

Os resultados da simulação indicaram que o estimador com melhor índice de qualidade de estimação intervalar foi o M_1 , com acurácia de 89,60% e comprimento médio de 0,1852. Considerando este estimador, o intervalo será:

$$IC(\rho^2, 95\%) = [0,4920 - 0,6800].$$

De forma análoga à anterior, realizando a simulação para o modelo quadrático:

```
> sICR2(rho2 = 0.7400878, n = 88, k = 2, MC = 1000, p1 = 0.5,
alpha0 = 0.05, JS = 1000)
```

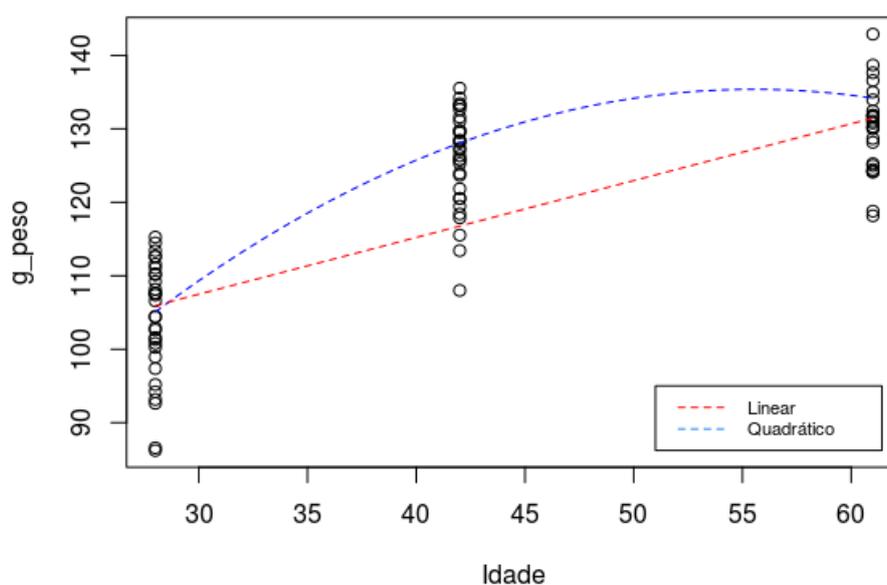
Para esse caso, o tempo computacional gasto foi de 45,36 minutos. O resultado apontou que o estimador M_2 possui melhor qualidade neste cenário, com acurácia de 88,9% e comprimento médio de 0,1309. Assim, os limites de confiança são dados por:

$$IC(\rho^2, 95\%) = [0,6704 - 0,8041].$$

É válido ressaltar que o tempo de simulação gasto pode ser inconveniente para o usuário. Entretanto, note que os resultados apontados pela função que estabelece os intervalos a partir dos resultados já simulados foram equivalentes com um esforço computacional muito menor. Com isso, o usuário pode escolher entre as duas funções para a análise da qualidade de ajuste dos modelos construídos.

Novamente, como os intervalos interseccionam-se, utilizou do critério de *Akaike* (AIC) para a tomada de decisão. O modelo quadrático apresentou um critério menor (AIC = 595,14) quando comparado ao linear (AIC = 633,47), contribuindo para a escolha do com menor AIC. A Figura 34 apresenta os modelos ajustados, sendo possível indicar que conforme o modelo quadrático que as matrizes com idade de 54,15 semanas apresentaram filhotes com maior média de ganho de peso nos 7 primeiros dias pós nascimento.

Figura 34 – Modelos ajustados ao experimento relativo ao ganho de peso de filhotes em função da idade das matrizes



Fonte: Do autor.

5 CONCLUSÃO

Implementações de pacotes em linguagem *R* têm se mostrado comuns na literatura estatística, possibilitando a interação entre os resultados de pesquisas e novos estudos. Para o caso pacote ICR2, há uma viabilização da construção de intervalos de confiança para a qualidade de ajuste de modelos de regressão, utilizando estimadores fundamentados na literatura ou propostos via simulação computacional.

A pluralidade de estimadores presentes é consequência de que a qualidade de estimação pode variar conforme o espaço paramétrico de ρ^2 , o que pode ser um entrave ao usuário quanto ao estimador a ser utilizado no cenário analisado. Desta forma, as funções disponibilizadas pelo pacote fazem com que o algoritmo decida em relação ao estimador com melhor desempenho, considerando as atribuições definidas pelo usuário, como o número de covariáveis e de observações.

Dado o esforço computacional que pode ser consumido por uma das funções, o pacote considerou também uma função que sugere o estimador a ser utilizado a partir de uma suavização dos cenários que já foram simulados, mantendo fixo o número de simulações Monte Carlo e o nível α . Com isso, há uma grande redução do tempo computacional. Por outro lado, o pacote também disponibilizou uma função que permite ao usuário a realização de uma simulação específica, caso este opte por outro nível de significância, por exemplo.

Assim, com a disponibilidade do pacote é possível introduzir a utilização dos estimadores à estudos aplicados de modo a avaliar a qualidade de ajuste, permitindo a análise e seleção de modelos. Como o caso do estudo de exemplo ilustrado pelo pacote, que analisou o efeito da idade da matriz no ganho de peso da progênie, assinalando que a qualidade dos modelos considerados eram iguais, permitindo a tomada de decisão por parte do usuário na escolha do modelo.

REFERÊNCIAS

ALMEIDA E SILVA, M. Evolução do melhoramento genético de aves no Brasil. **Ceres**, v. 56, n. 3, p. 437-445, 2009. Disponível em: <http://www.redalyc.org/pdf/3052/305226808008.pdf>. Acesso em: 2 fev. 2019.

DI BUCCHIANICO, A. Coefficient of determination (R^2). In: RUGGERI, F; KENETT, R; FALTIN, F. **Encyclopedia of statistics in quality and reliability**, [s.l.]: John Wiley & Sons, 2008. p. 1-2.

CRAMER, J.S. Mean and variance of R^2 in small and moderate samples. **Journal of econometrics**, v. 35, ed. 2-3, p. 253-266, 1987. Disponível em: <https://www.sciencedirect.com/science/article/pii/0304407687900273>. Acesso em: 21 jun. 2018.

LÔBO, R; BITTNECOURT, T; PINTO, L. Progresso científico em melhoramento animal no Brasil na primeira década do século XXI. **R. Bras. zootec.**, v. 39, p. 437-445, 2010. Disponível em: <http://www.scielo.br/scielo.php?script=sciarttextpid=S1516-35982010001300025>. Acesso em: 10 fev. 2019.

MORENO, F. **Efeito da idade da matriz e peso dos ovos no desempenho da progênie**. 2019. Dissertação (Mestrado em zootecnia) - Universidade Federal do Paraná, Curitiba, 2019. Disponível em: <https://acervodigital.ufpr.br/bitstream/handle/1884/61424/R%20-%20D%20-%20FILIPE%20AUGUSTO%20MORENO.pdf?sequence=1&isAllowed=y>. Acesso em: 17 fev. 2019.

KUZNETSOVA, A; BROCKHOFF, P.; CHRISTENSEN, R. LmerTest package: tests in linear mixed effects models. **Journal of statistical software**, v. 82, n. 13, 2017. Disponível em: <https://www.jstatsoft.org/article/view/v082i13>. Acesso em: 22 fev. 2019.

PIEPHO, H. A coefficient of determination (R^2) for generalized linear mixed models. **Biometrical journal**, v. 61, n. 4, p. 860-872, 2019. Disponível em: <https://onlinelibrary.wiley.com/doi/10.1002/bimj.201800270>. Acesso em: 14 jan. 2019.

QUININO, R.; REIS, E.; BESSEGATO, L. Using the coefficient of determination R^2 to test the significance of multiple linear regression. **Teaching statistics**, v. 35, n. 2, p. 84-88, 2012. Disponível em: <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-9639.2012.00525.x>. Acesso em: 19 ago. 2018.

R Core Team. **R: a language and environment for statistical computing**. Vienna, Austria, 2016. Disponível em: <https://www.R-project.org/>.

STURTZ, S; LIGGES, U; GELMAN, A. R2WinBUGS: a package for running WinBUGS from R. **Journal of statistical software**, v. 12, n. 3, 2005. Disponível em: <https://www.jstatsoft.org/article/view/v012i03/v12i03.pdf>. Acesso em: 2 mar. 2019.

WEATHERBURN, C. **A first course in mathematical statistics**. Cambridge: University Press, 1949. v. 32.

ZHANG, D. A coefficient of determination for generalized linear models. **The American statistician**, v. 71, n. 4, p. 310-316, 2017. Disponível em: <https://amstat.tandfonline.com/doi/abs/10.1080/00031305.2016.1256839.XUiHCstKicw>. Acesso em: 23 fev. 2019.

CONSIDERAÇÕES FINAIS

Neste trabalho, foi analisado o comportamento de dois estimadores intervalares indicados na literatura e propôs-se outros três, sendo dois deles construídos a partir da reparametrização do estimador W , disponível em Weatherburn (1949) e Foster, Smith e Whaley (1997). Para isto, avaliaram-se as taxas de acurácia e as médias de comprimento dos intervalos construídos por cada estimador a um nível de confiança de 95%, medidas também aplicadas por Carari et al. (2010) e Scacabarozzi e Diniz (2007) para a comparação de estimadores.

Até então os estudos que analisaram estimadores intervalares exploraram as taxas de acurácia e as médias de comprimento de forma separada, propiciando recomendações que podem ser opostas, ou seja, um estimador com maior acurácia e com pior média de comprimento dentre os demais. Assim, entendeu-se como necessário a elaboração de um índice, dado como uma média ponderada entre essas duas medidas, possibilitando a recomendação de um estimador mais parcimonioso dentre os demais.

Com os cenários analisados, os estimadores propostos neste trabalho apresentaram qualidade aproximada ou superior aos fundamentados na literatura em regiões do espaço avaliado. Em decorrência deste fato, foi possível a recomendação em três regiões, sendo elas: i) $\rho^2 = 0$, ii) $0,1 \leq \rho^2 \leq 0,5$ e iii) $0,5 > \rho^2 \leq 0,9$. Na primeira região, entendeu-se como válido a recomendação do estimador W , que apresentou os melhores índices de qualidade e as maiores taxas de acurácia nos cenários simulados. Na segunda região, o estimador M_1 foi recomendado, apresentando uma melhoria de seus índices sob condição assintótica e aproximando-se do estimador de maior qualidade quando este não apresentou o melhor índice. Na terceira região, os estimadores M_2 e E^* foram recomendados, apresentando melhoria dos índices de qualidade com o aumento do tamanho amostral, principalmente na região superior a 0,7. Essas considerações conduzem ao fato de que esse assunto não é esgotado na literatura, sendo possível a melhoria da qualidade de inferência a partir de reparametrizações dos estimadores já existentes e com características menos complexas do que a distribuição exata tratada em Cramer (1987), que pode apresentar problemas de integração numérica com o aumento do tamanho amostral e incerteza quanto ao limite superior do somatório infinito que pondera a distribuição.

Dada as recomendações justificadas anteriormente, entendeu-se como fundamental da expansão de cenários que ajustem-se a contextos particulares. Com isso, o pacote ICR2 viabilizou essa ampliação permitindo a estimação intervalar da qualidade de ajuste em outros cenários,

apontando os estimadores recomendados e indicando o estimador de melhor qualidade. Por fim, a disponibilidade do pacote torna acessível a proposição e a implementação de novos testes à literatura, de modo que o índice de desempenho aponte sempre aquele com maior qualidade. Fato que possibilita o aprimoramento da literatura e da tomada de decisão a respeito da qualidade de ajuste de modelos.

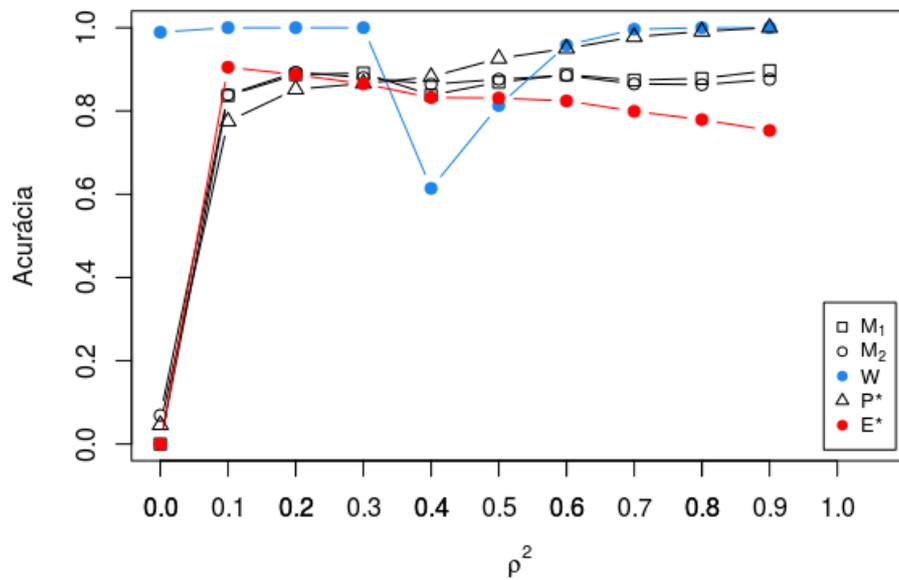
Como sugestão de trabalhos futuros é necessário explorar o impacto da ponderação p_1 no índice de desempenho de estimação intervalar no processo de escolha de um estimador. Um segundo ponto é analisar o desempenho dos estimadores em cenários com reduzido tamanho amostral, como o caso da regressão na análise de variância, comumente utilizada no contexto experimental. Além de relacionar outras distribuições aos estimadores como a distribuição de Kumaraswamy, por exemplo.

APÊNDICES

APÊNDICE A

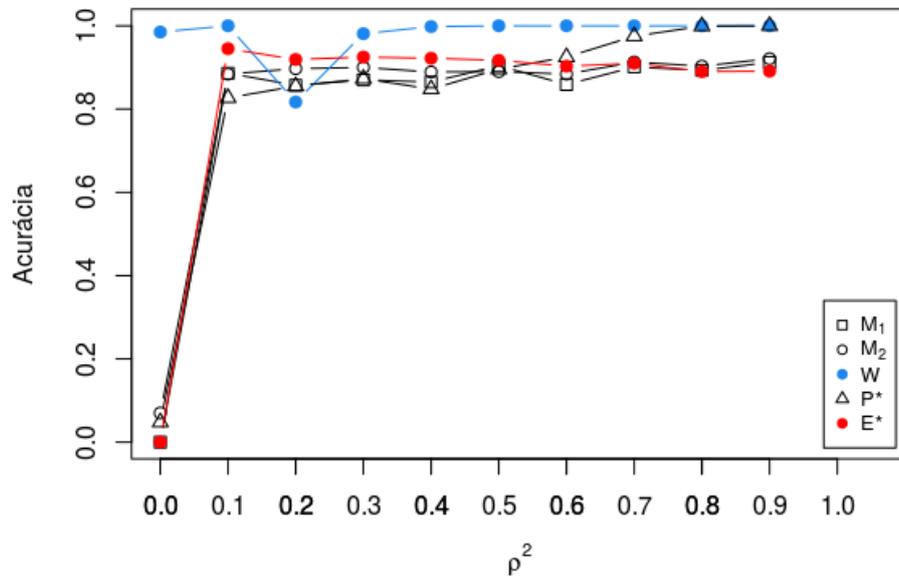
Nessa seção serão disponibilizadas as taxas de acurácia para os demais cenários analisados.

Figura 35 – Taxa de acurácia referente ao modelo de regressão onde $k = 2$ e $n = 15$



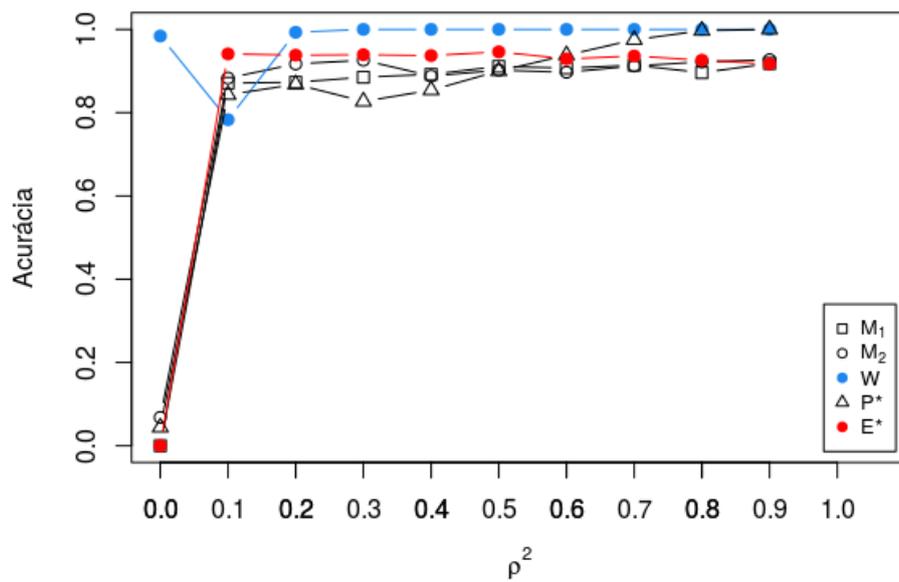
Fonte: Do autor.

Figura 36 – Taxa de acurácia referente ao modelo de regressão onde $k = 2$ e $n = 50$



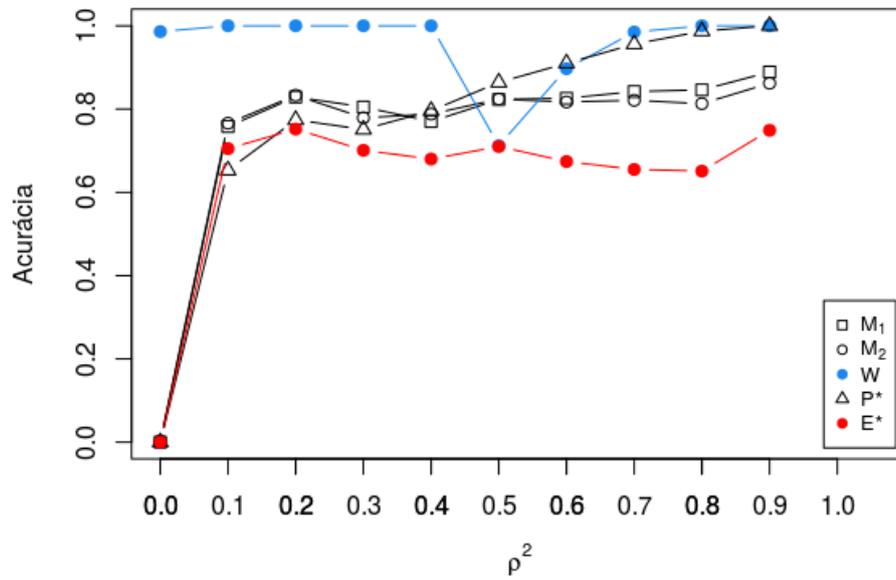
Fonte: Do autor.

Figura 37 – Taxa de acurácia referente ao modelo de regressão onde $k = 2$ e $n = 100$



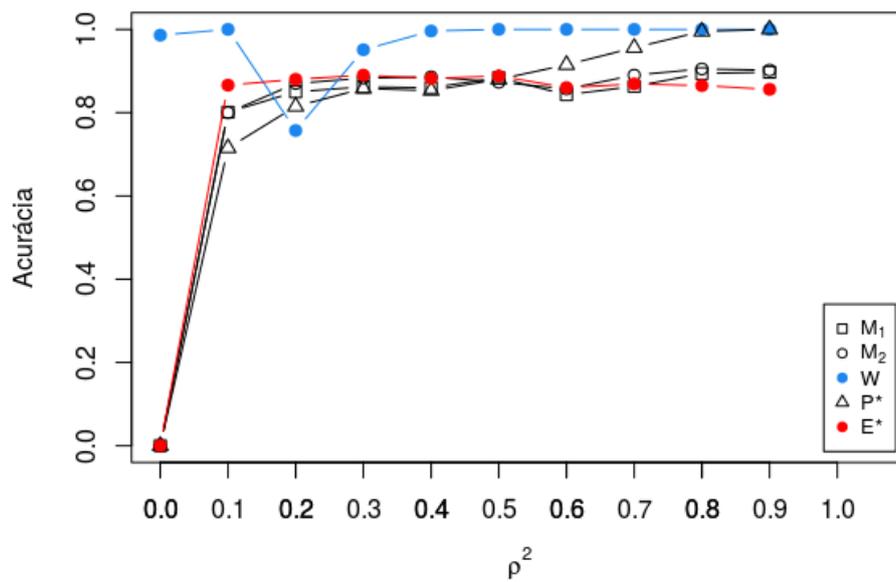
Fonte: Do autor.

Figura 38 – Taxa de acurácia referente ao modelo de regressão onde $k = 3$ e $n = 15$



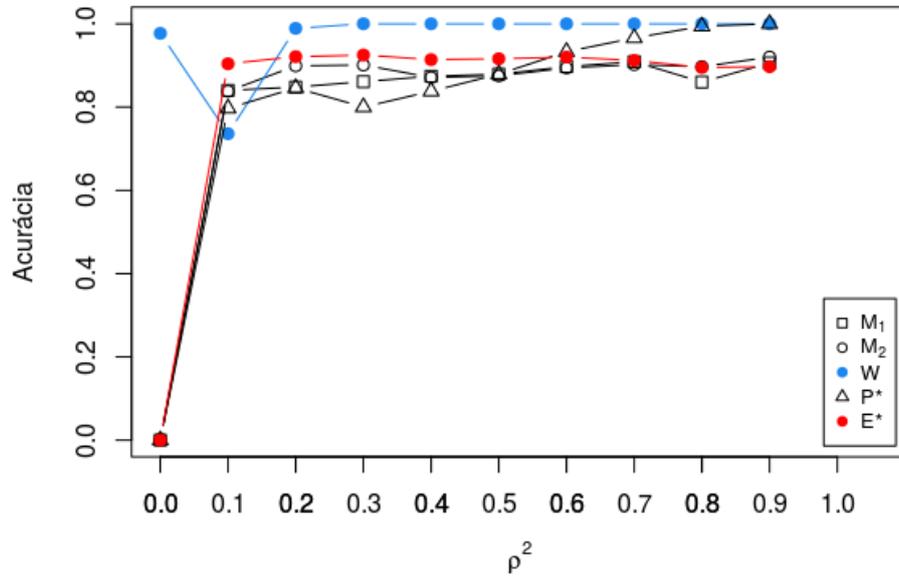
Fonte: Do autor.

Figura 39 – Taxa de acurácia referente ao modelo de regressão onde $k = 3$ e $n = 50$



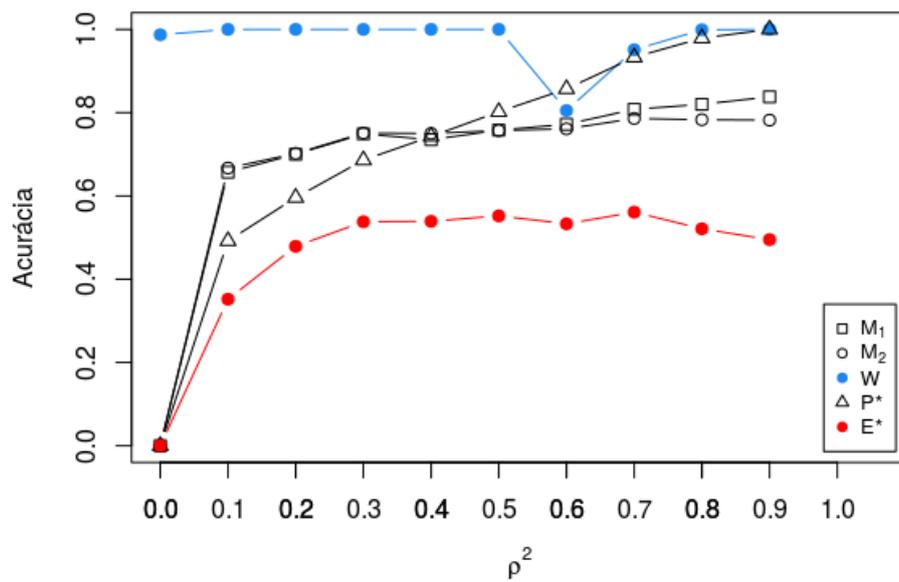
Fonte: Do autor.

Figura 40 – Taxa de acurácia referente ao modelo de regressão onde $k = 3$ e $n = 100$



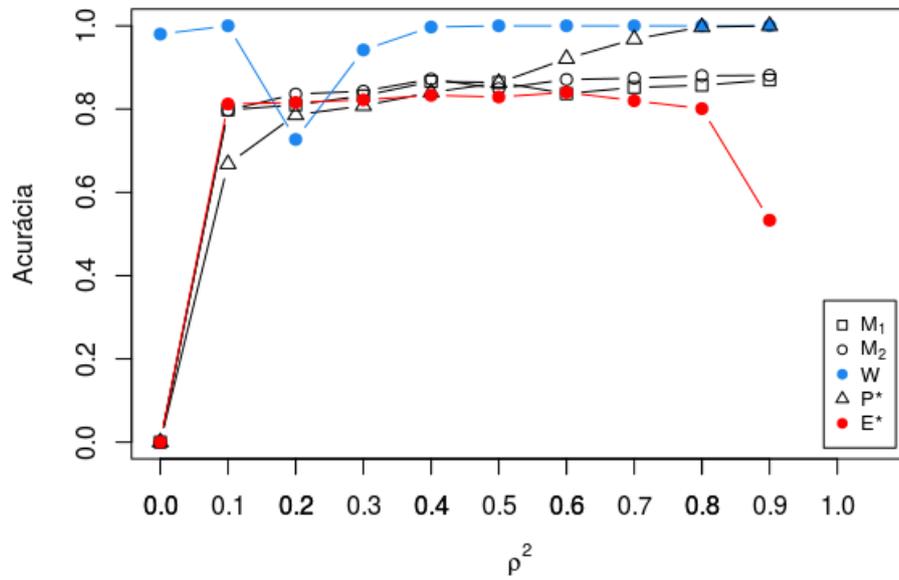
Fonte: Do autor.

Figura 41 – Taxa de acurácia referente ao modelo de regressão onde $k = 4$ e $n = 15$



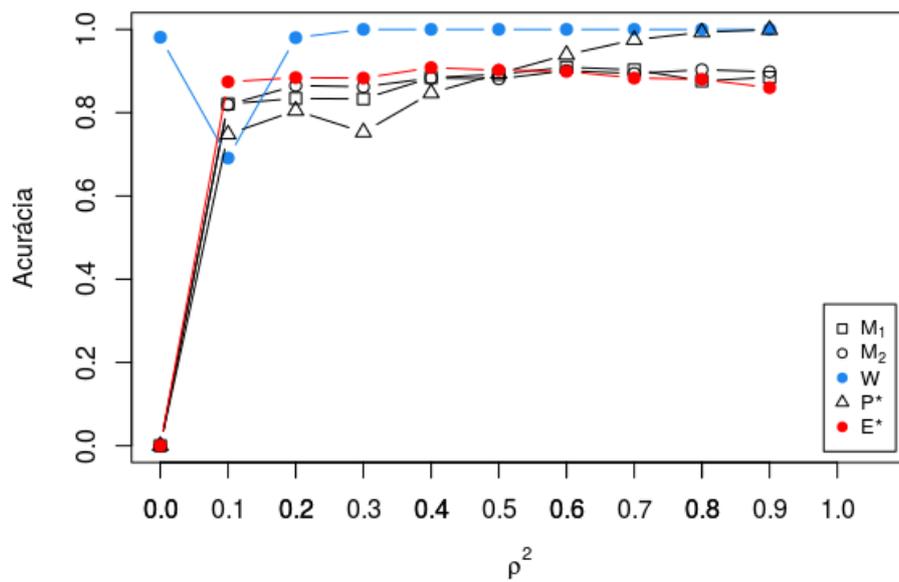
Fonte: Do autor.

Figura 42 – Taxa de acurácia referente ao modelo de regressão onde $k = 4$ e $n = 50$



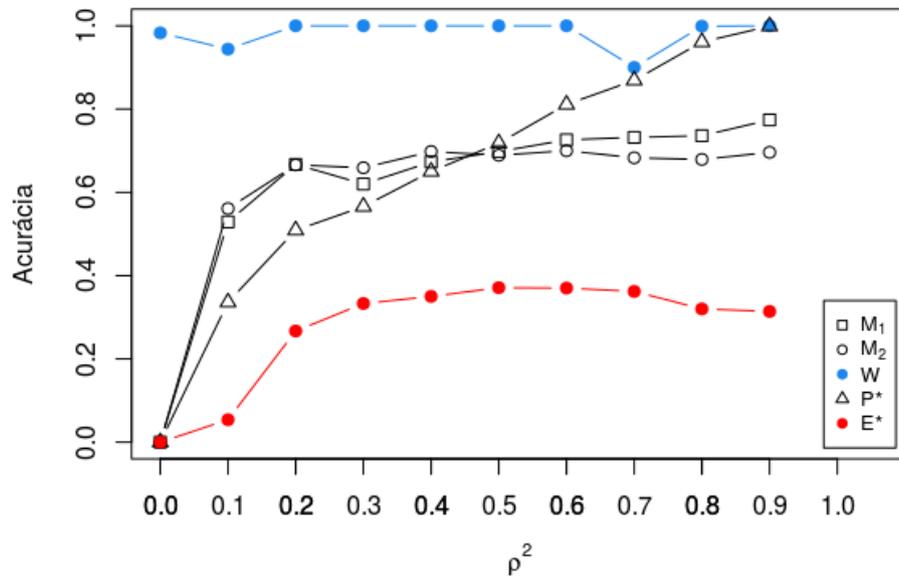
Fonte: Do autor.

Figura 43 – Taxa de acurácia referente ao modelo de regressão onde $k = 4$ e $n = 100$



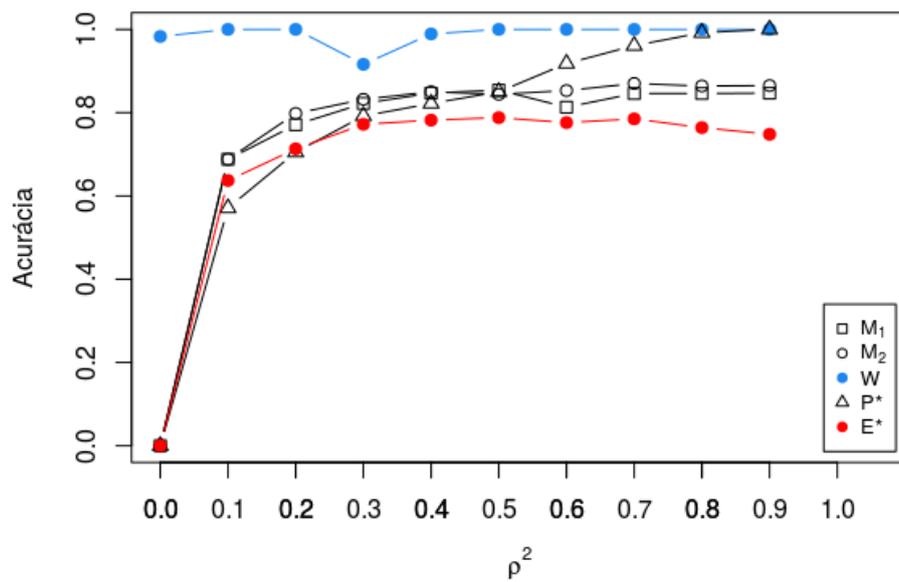
Fonte: Do autor.

Figura 44 – Taxa de acurácia referente ao modelo de regressão onde $k = 5$ e $n = 15$



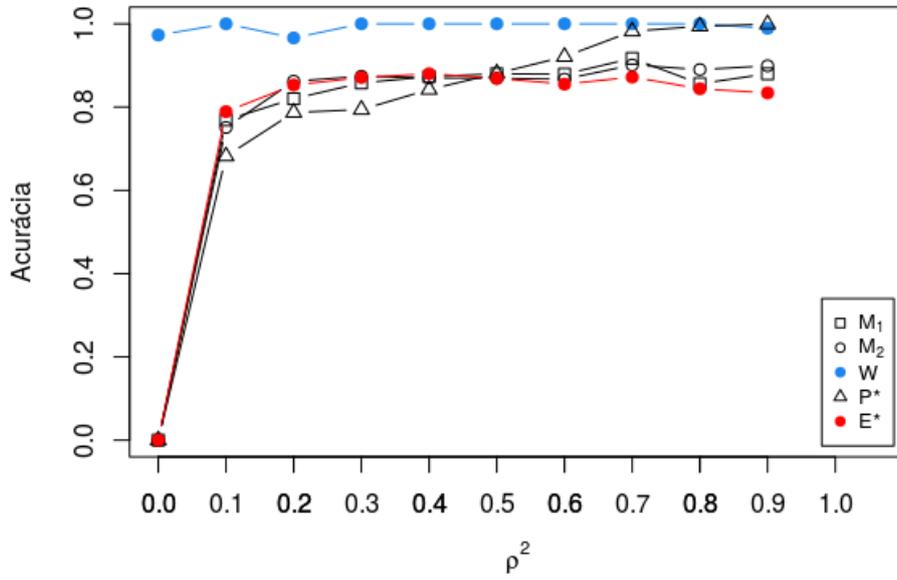
Fonte: Do autor.

Figura 45 – Taxa de acurácia referente ao modelo de regressão onde $k = 5$ e $n = 50$



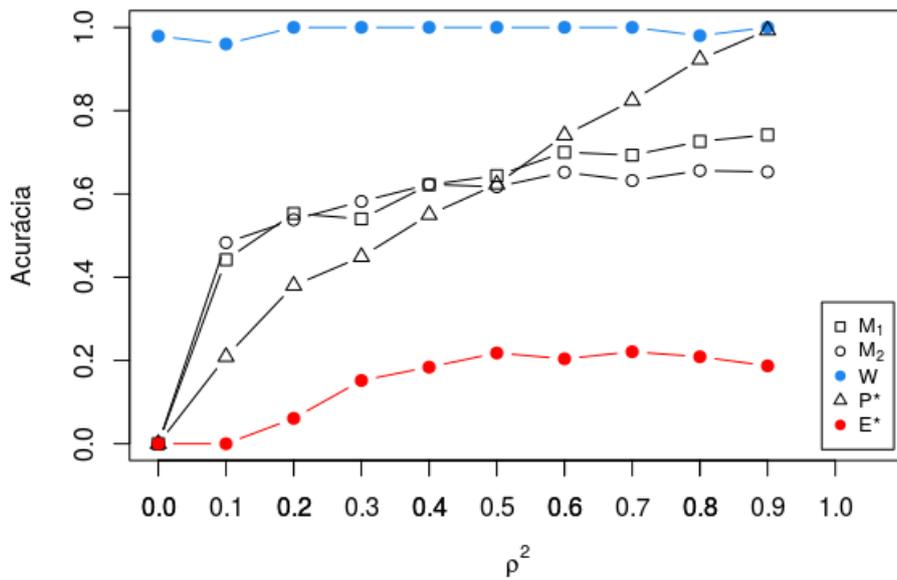
Fonte: Do autor.

Figura 46 – Taxa de acurácia referente ao modelo de regressão onde $k = 5$ e $n = 100$



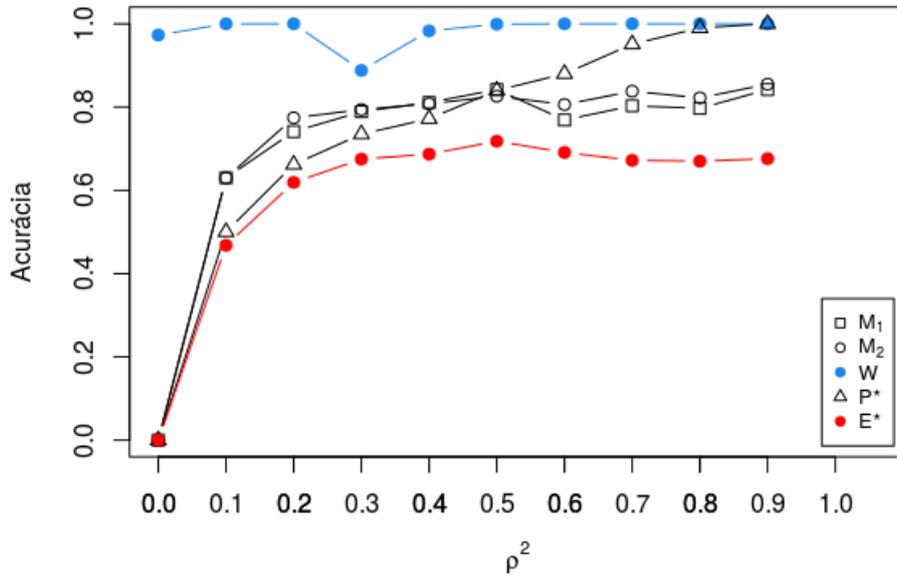
Fonte: Do autor.

Figura 47 – Taxa de acurácia referente ao modelo de regressão onde $k = 6$ e $n = 15$



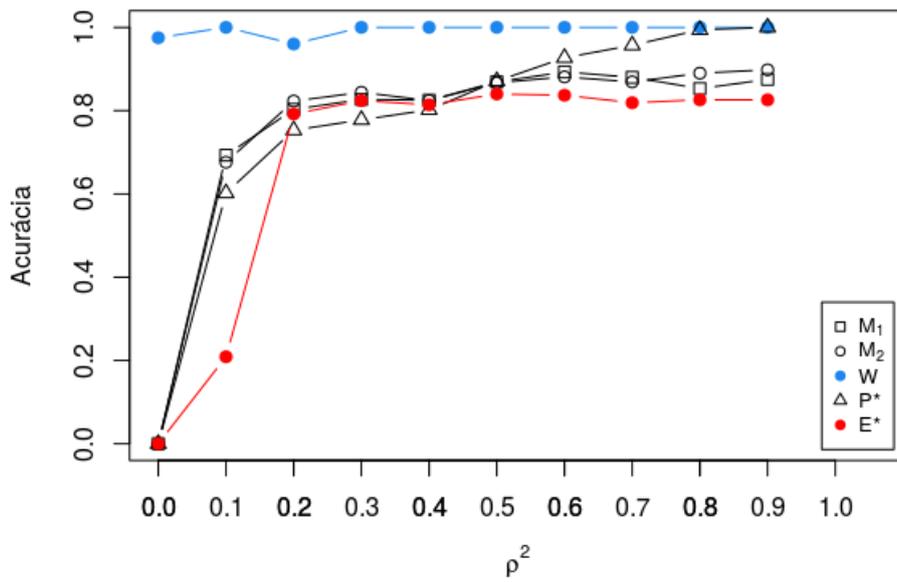
Fonte: Do autor.

Figura 48 – Taxa de acurácia referente ao modelo de regressão onde $k = 6$ e $n = 50$



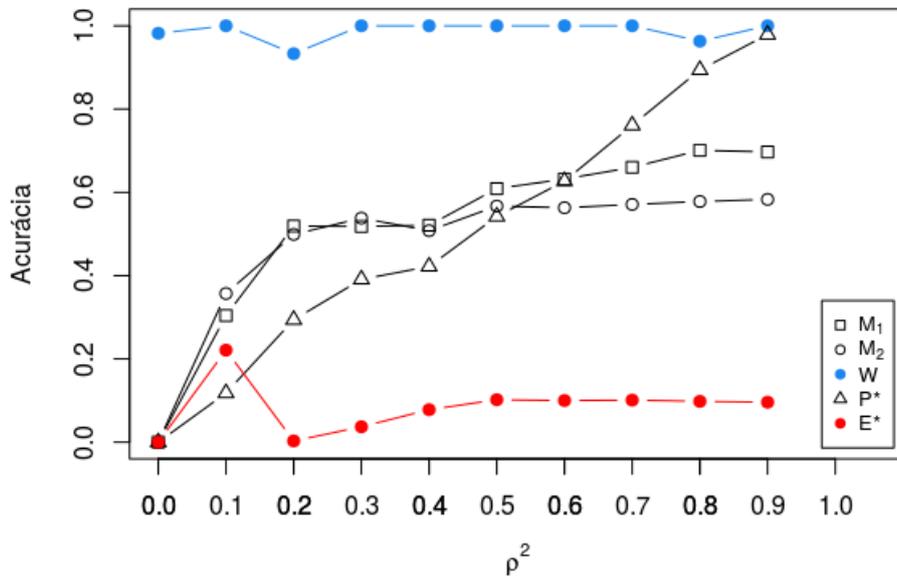
Fonte: Do autor.

Figura 49 – Taxa de acurácia referente ao modelo de regressão onde $k = 6$ e $n = 100$



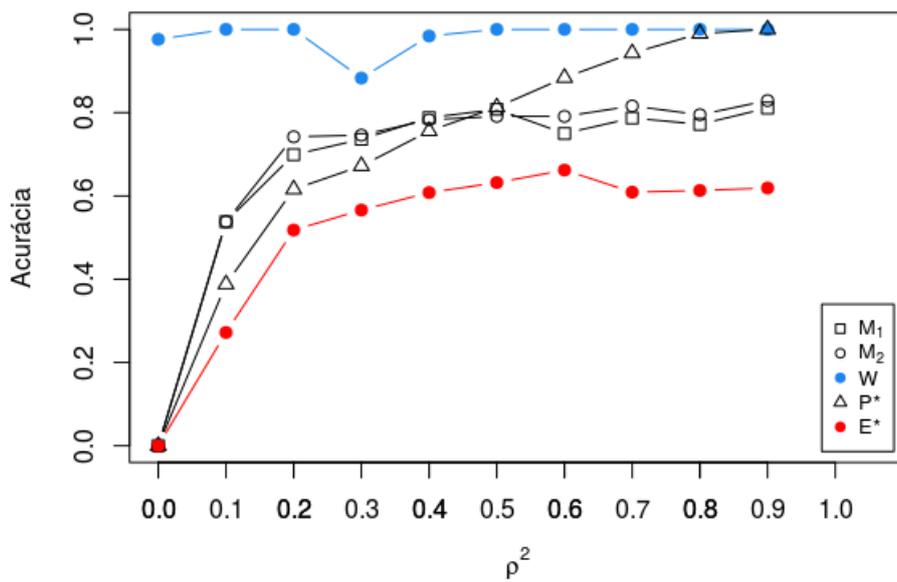
Fonte: Do autor.

Figura 50 – Taxa de acurácia referente ao modelo de regressão onde $k = 7$ e $n = 15$



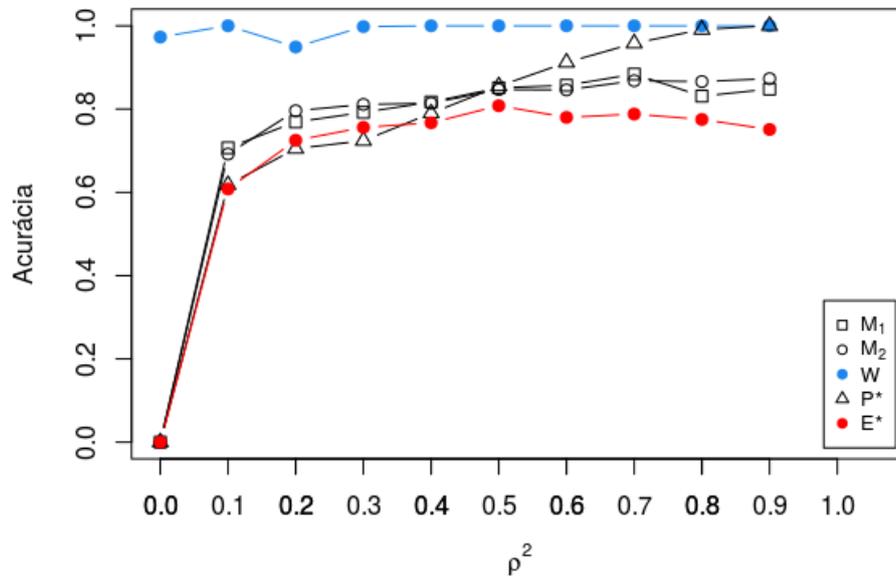
Fonte: Do autor.

Figura 51 – Taxa de acurácia referente ao modelo de regressão onde $k = 7$ e $n = 50$



Fonte: Do autor.

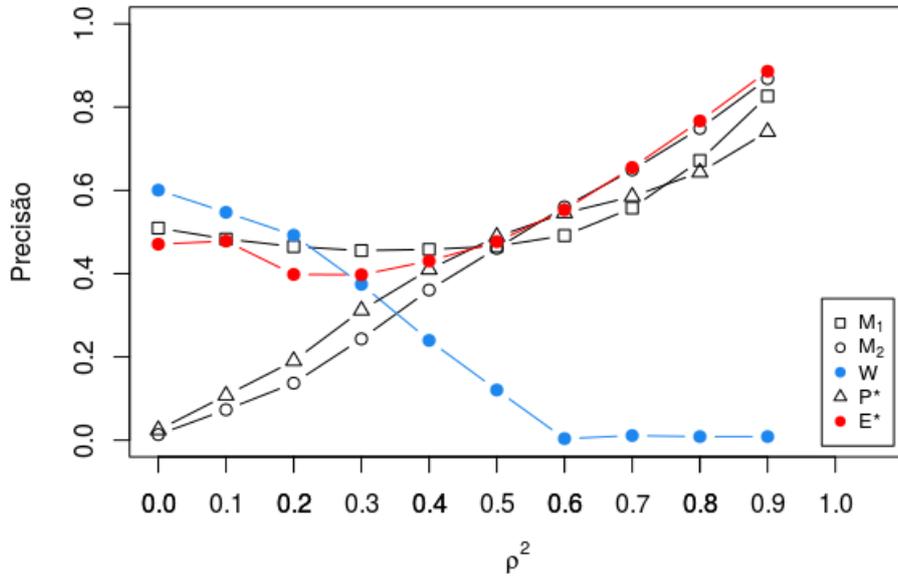
Figura 52 – Taxa de acurácia referente ao modelo de regressão onde $k = 7$ e $n = 100$



Fonte: Do autor.

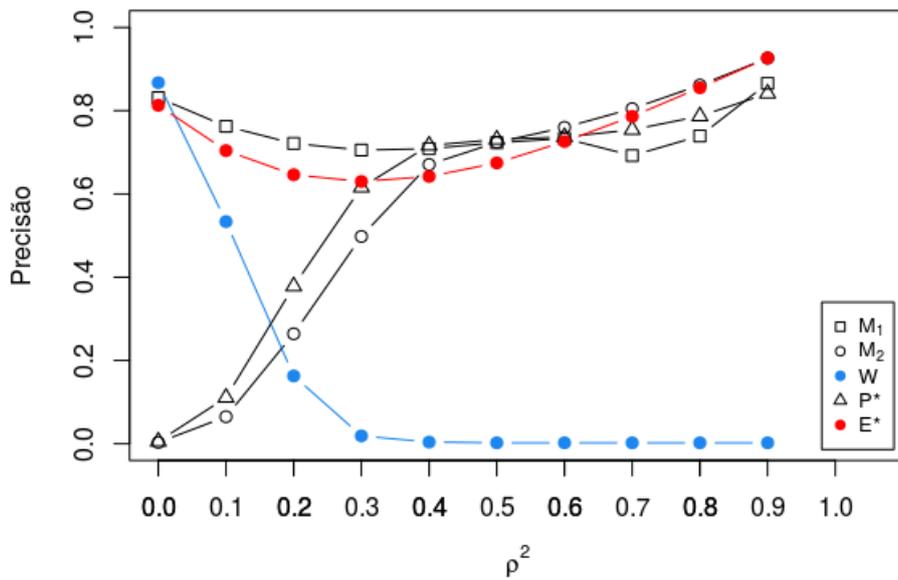
APÊNDICE B

Figura 53 – Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 15$



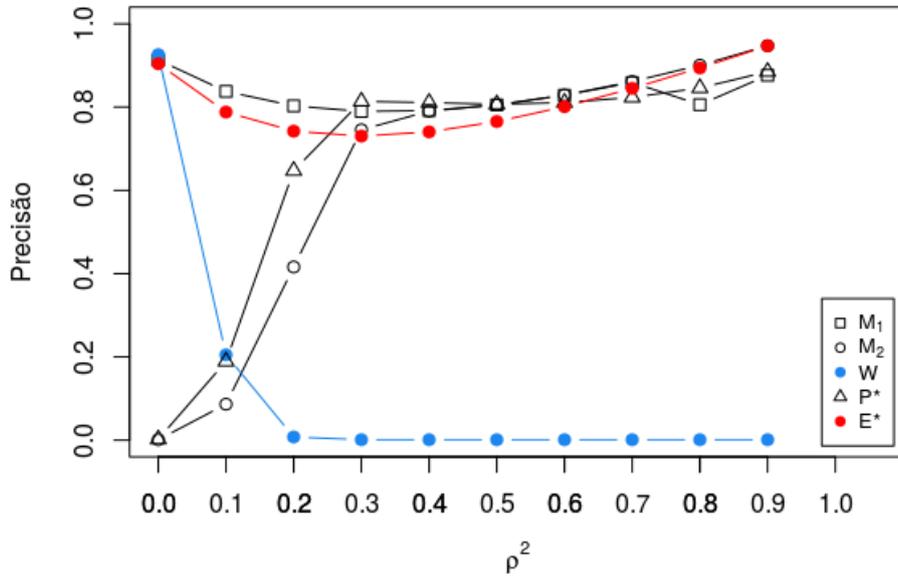
Fonte: Do autor.

Figura 54 – Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 50$



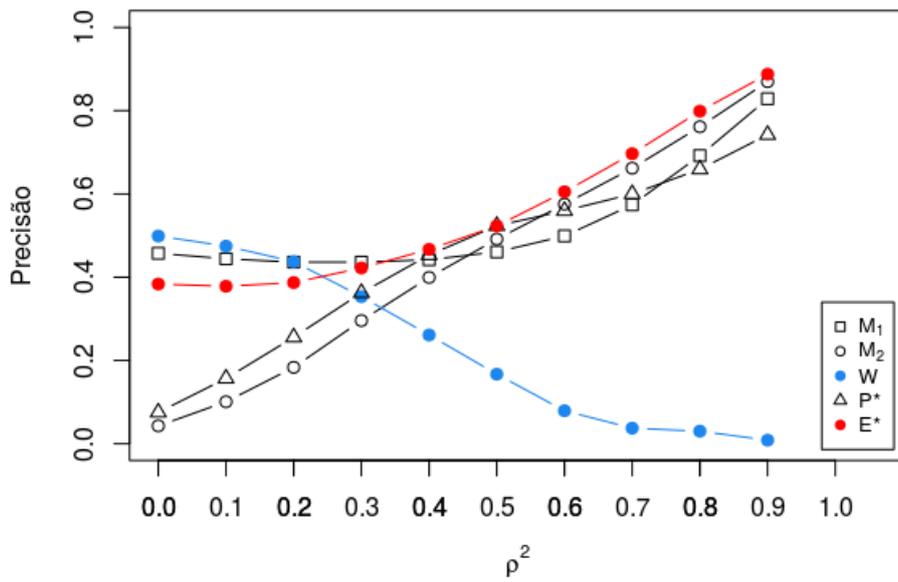
Fonte: Do autor.

Figura 55 – Precisão dos intervalos referente ao modelo de regressão onde $k = 1$ e $n = 100$



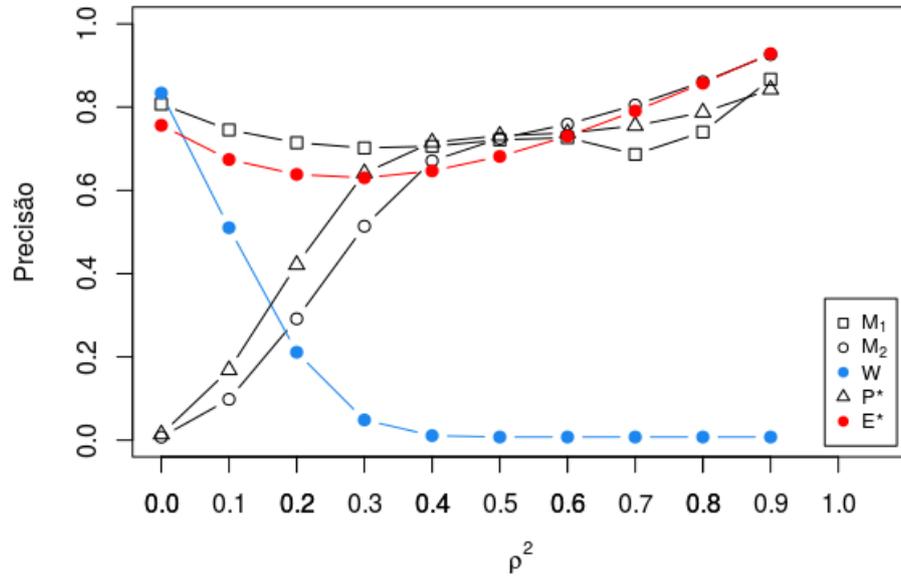
Fonte: Do autor.

Figura 56 – Precisão dos intervalos referente ao modelo de regressão onde $k = 3$ e $n = 15$



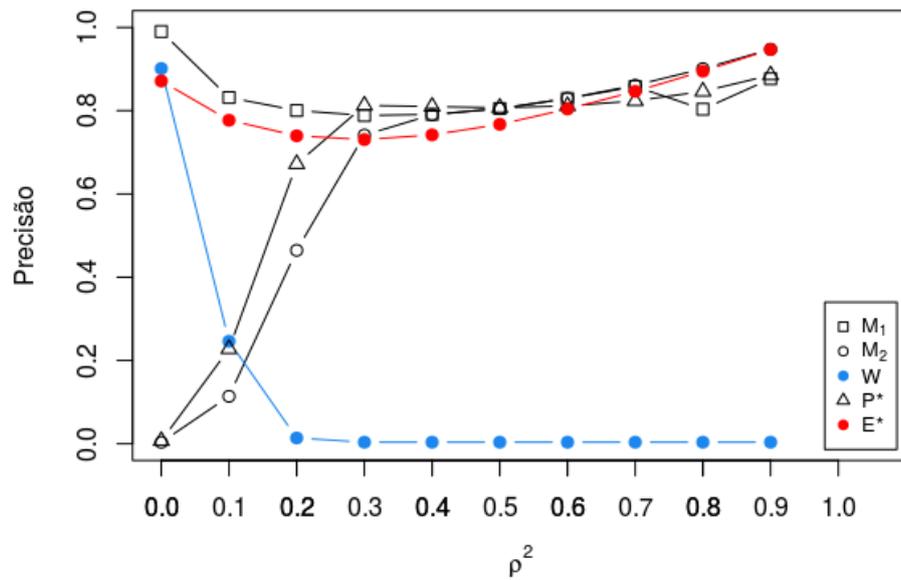
Fonte: Do autor.

Figura 57 – Precisão dos intervalos referente ao modelo de regressão onde $k = 3$ e $n = 50$



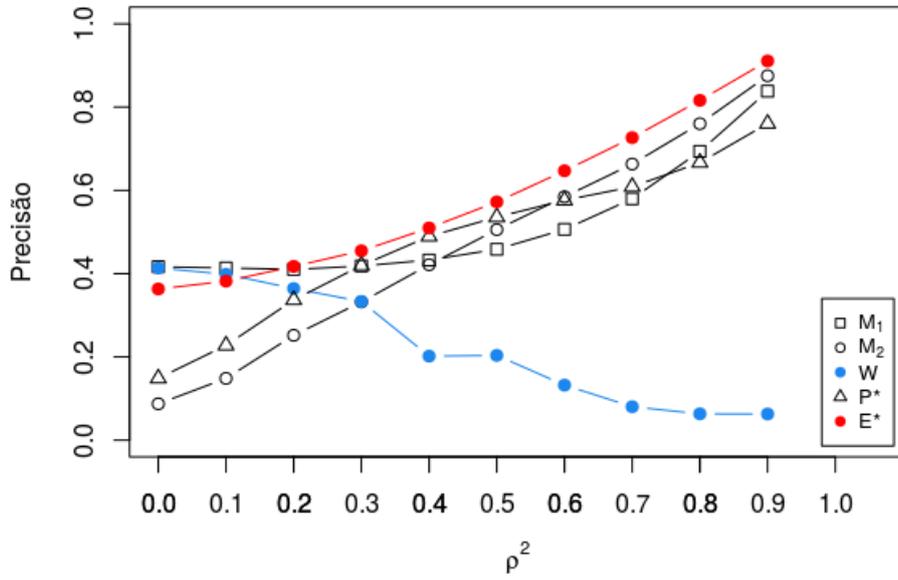
Fonte: Do autor.

Figura 58 – Precisão dos intervalos referente ao modelo de regressão onde $k = 3$ e $n = 100$



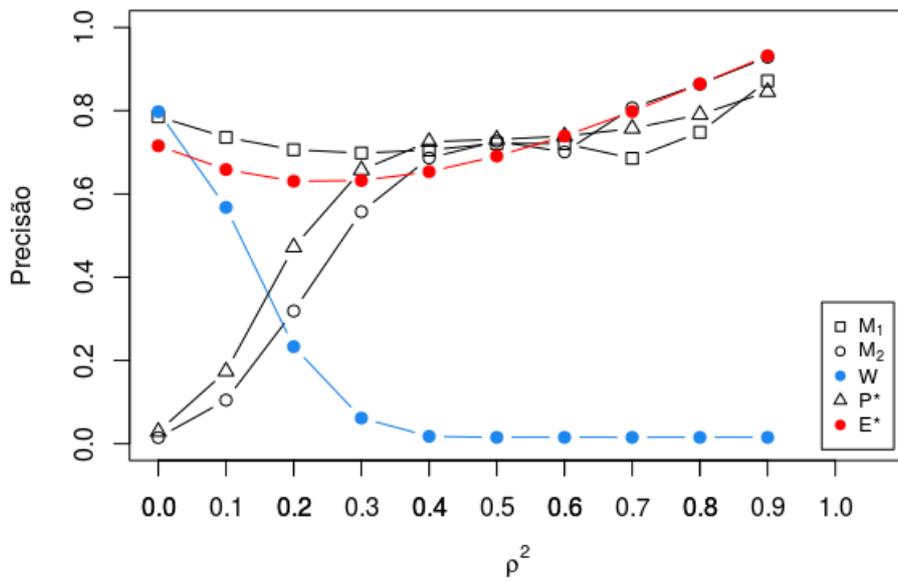
Fonte: Do autor.

Figura 59 – Precisão dos intervalos referente ao modelo de regressão onde $k = 4$ e $n = 15$



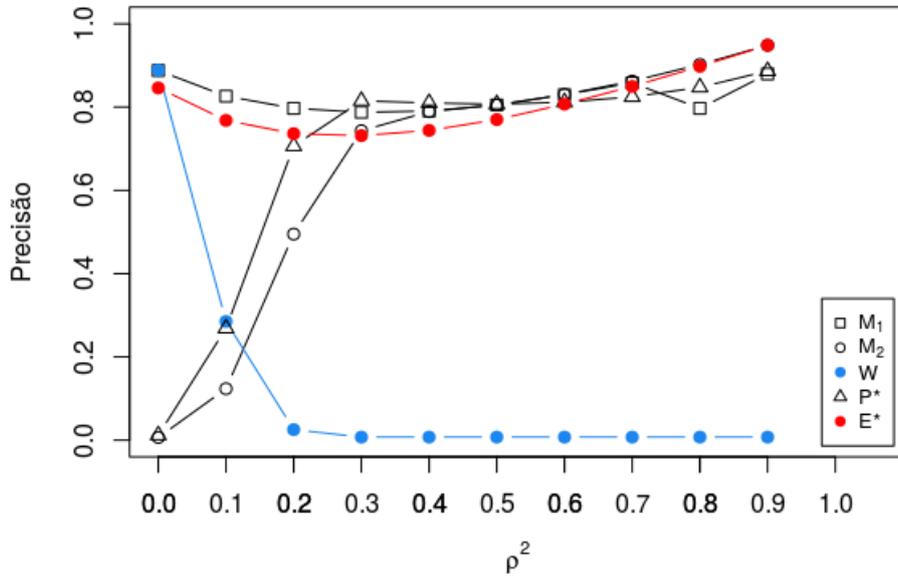
Fonte: Do autor.

Figura 60 – Precisão dos intervalos referente ao modelo de regressão onde $k = 4$ e $n = 50$



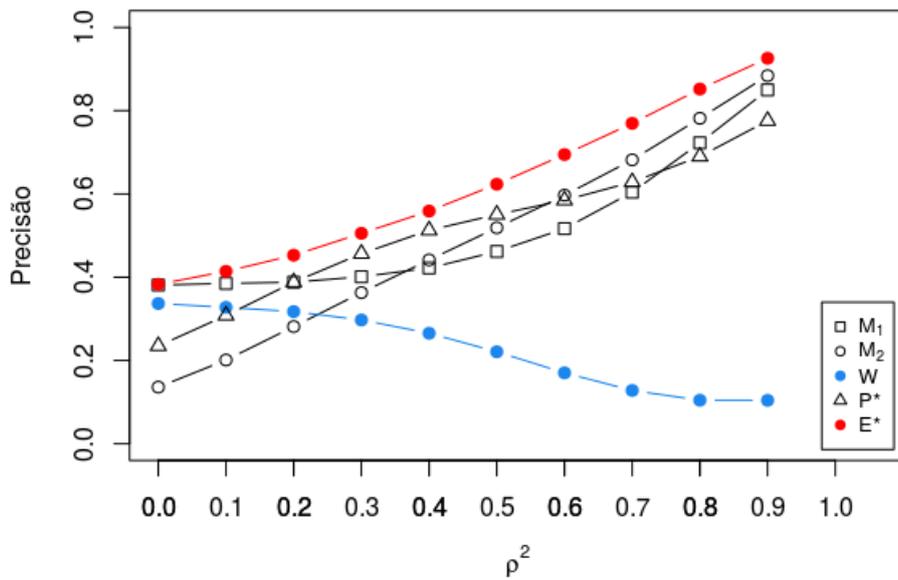
Fonte: Do autor.

Figura 61 – Precisão dos intervalos referente ao modelo de regressão onde $k = 4$ e $n = 100$



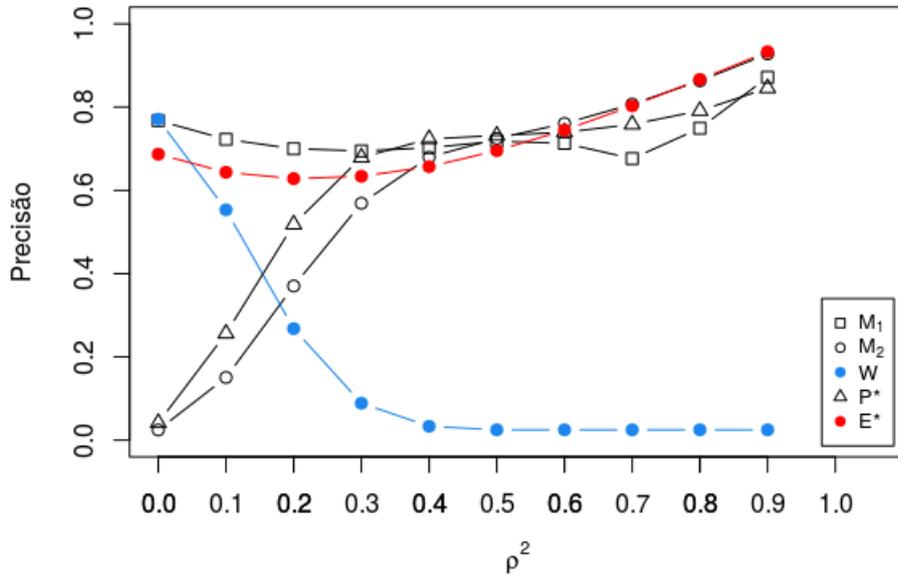
Fonte: Do autor.

Figura 62 – Precisão dos intervalos referente ao modelo de regressão onde $k = 5$ e $n = 15$



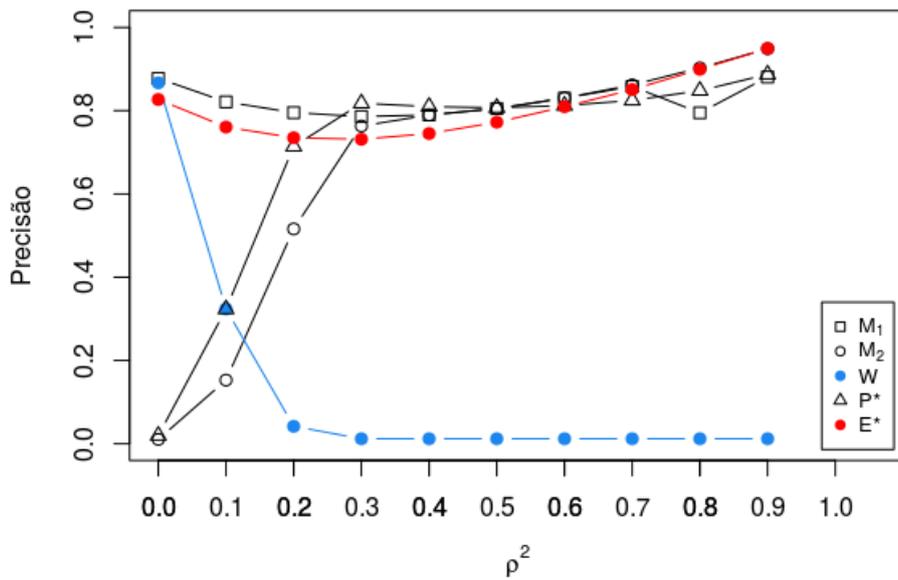
Fonte: Do autor.

Figura 63 – Precisão dos intervalos referente ao modelo de regressão onde $k = 5$ e $n = 50$



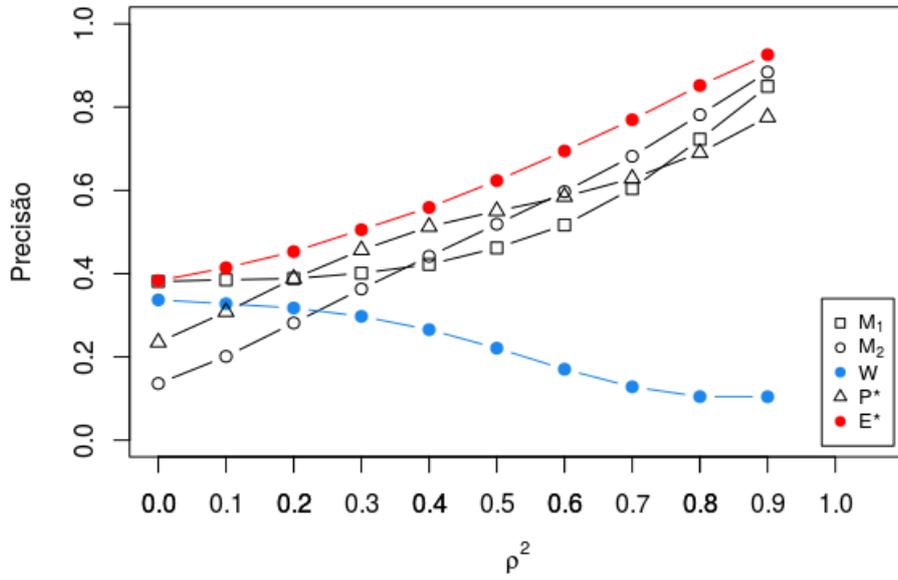
Fonte: Do autor.

Figura 64 – Precisão dos intervalos referente ao modelo de regressão onde $k = 5$ e $n = 100$



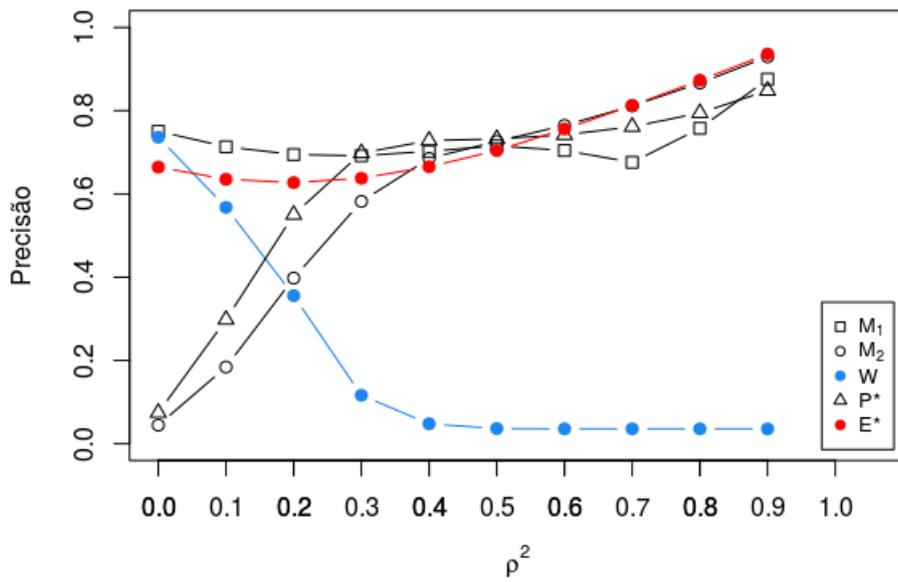
Fonte: Do autor.

Figura 65 – Precisão dos intervalos referente ao modelo de regressão onde $k = 6$ e $n = 15$



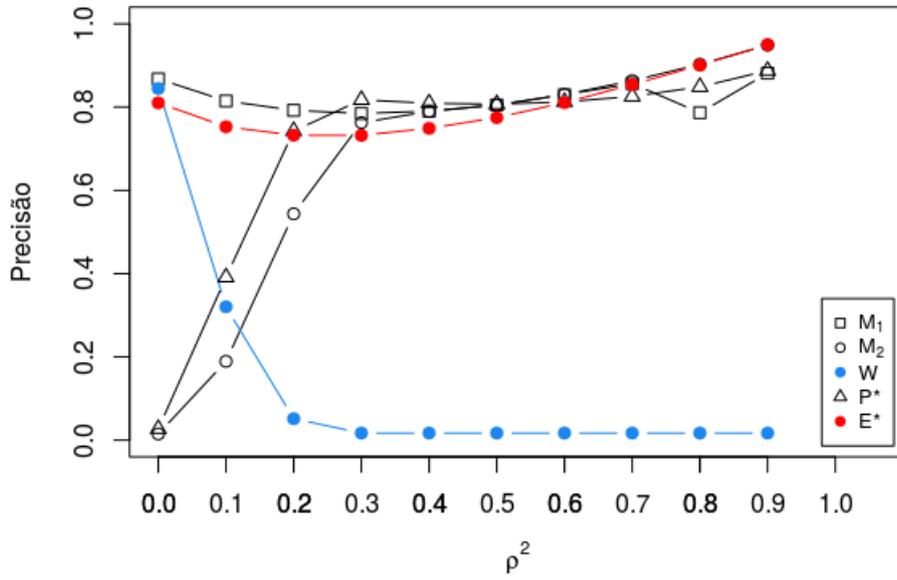
Fonte: Do autor.

Figura 66 – Precisão dos intervalos referente ao modelo de regressão onde $k = 6$ e $n = 50$



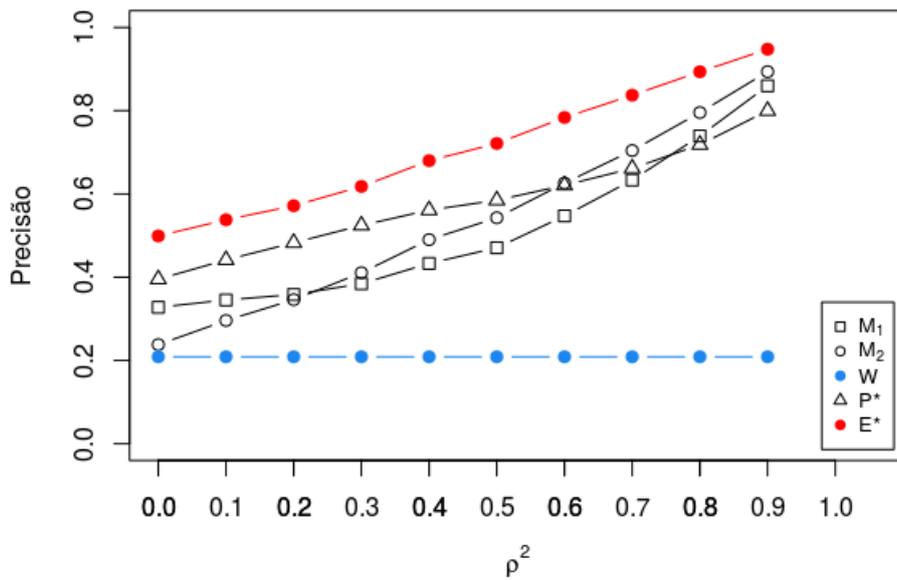
Fonte: Do autor.

Figura 67 – Precisão dos intervalos referente ao modelo de regressão onde $k = 6$ e $n = 100$



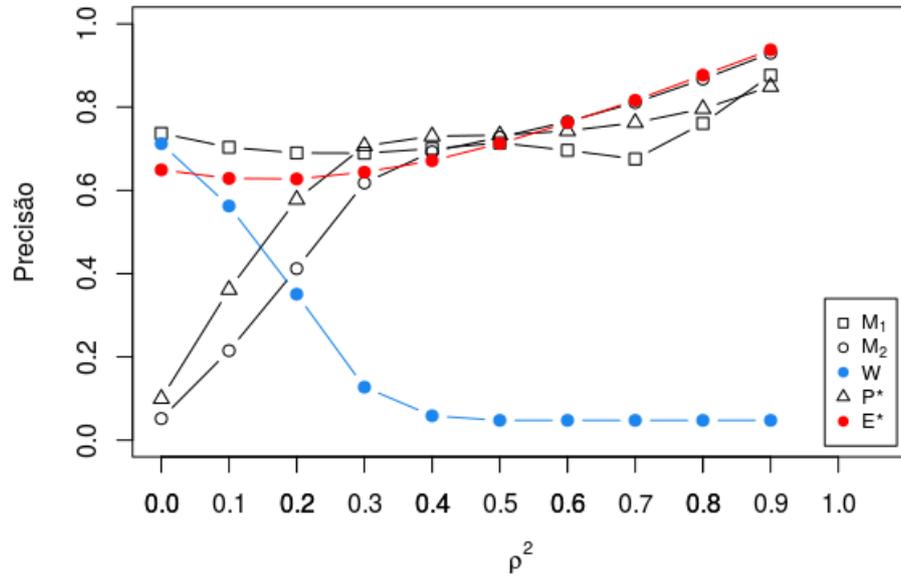
Fonte: Do autor.

Figura 68 – Precisão dos intervalos referente ao modelo de regressão onde $k = 7$ e $n = 15$



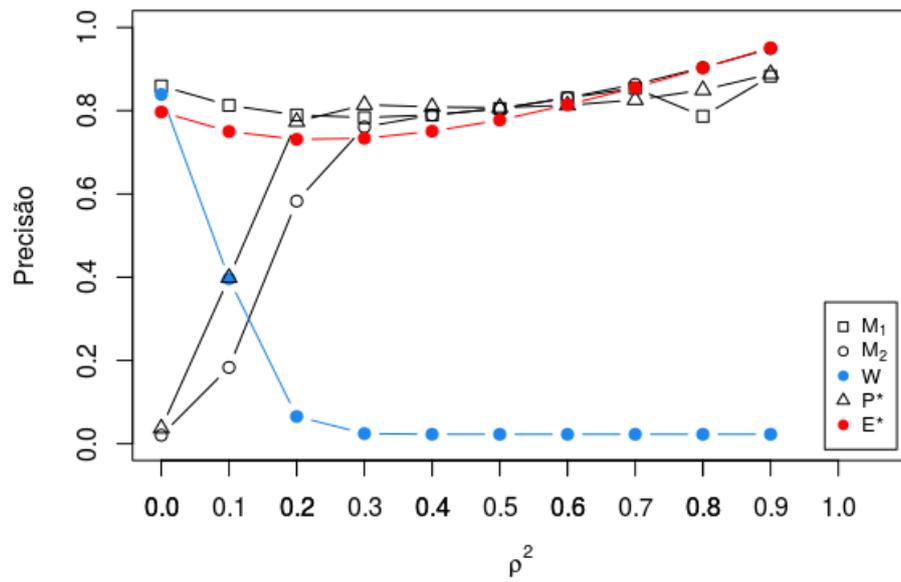
Fonte: Do autor.

Figura 69 – Precisão dos intervalos referente ao modelo de regressão onde $k = 7$ e $n = 50$



Fonte: Do autor.

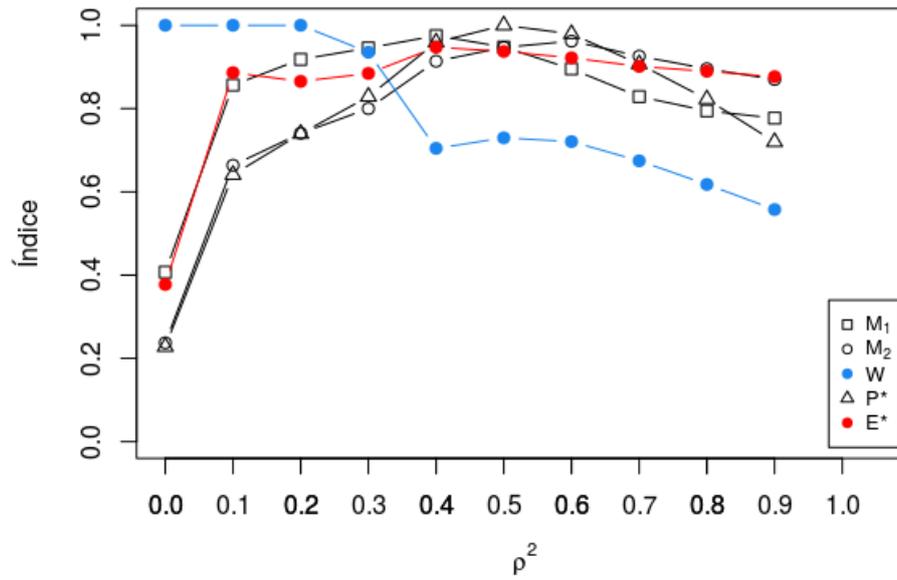
Figura 70 – Precisão dos intervalos referente ao modelo de regressão onde $k = 7$ e $n = 100$



Fonte: Do autor.

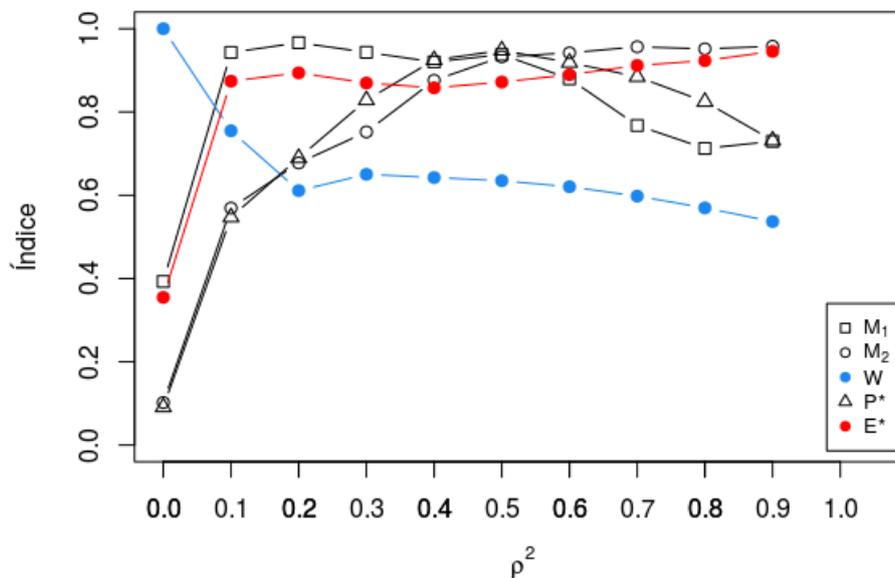
APÊNDICE C

Figura 71 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 2$ e $n = 15$



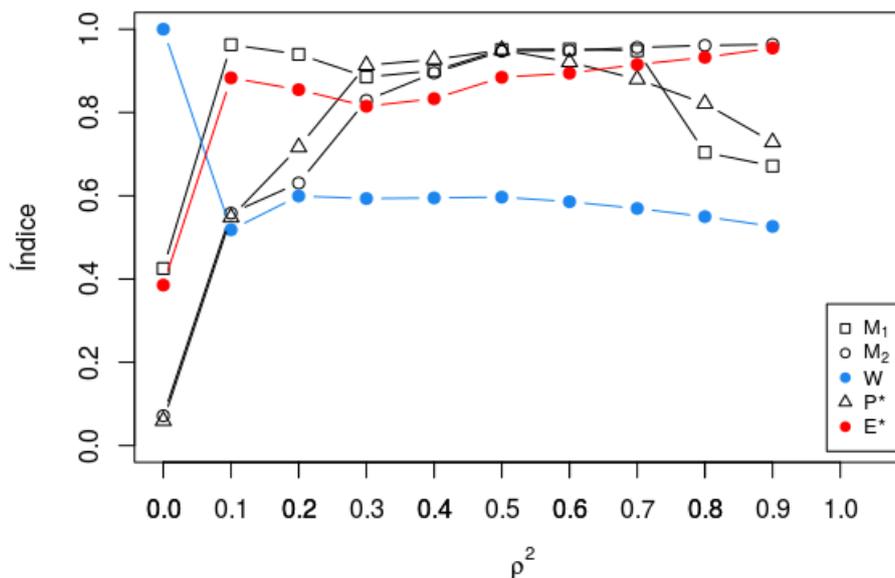
Fonte: Do autor.

Figura 72 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 2$ e $n = 50$



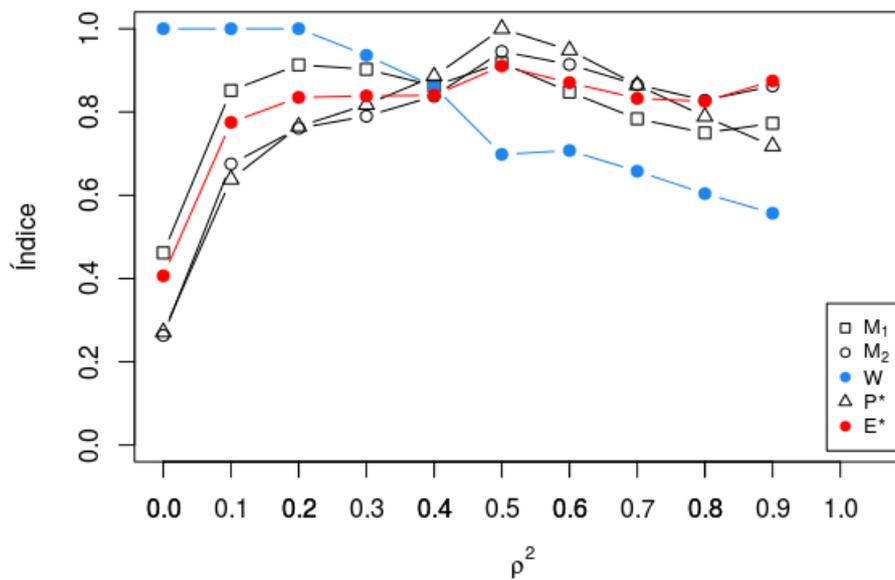
Fonte: Do autor.

Figura 73 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 2$ e $n = 100$



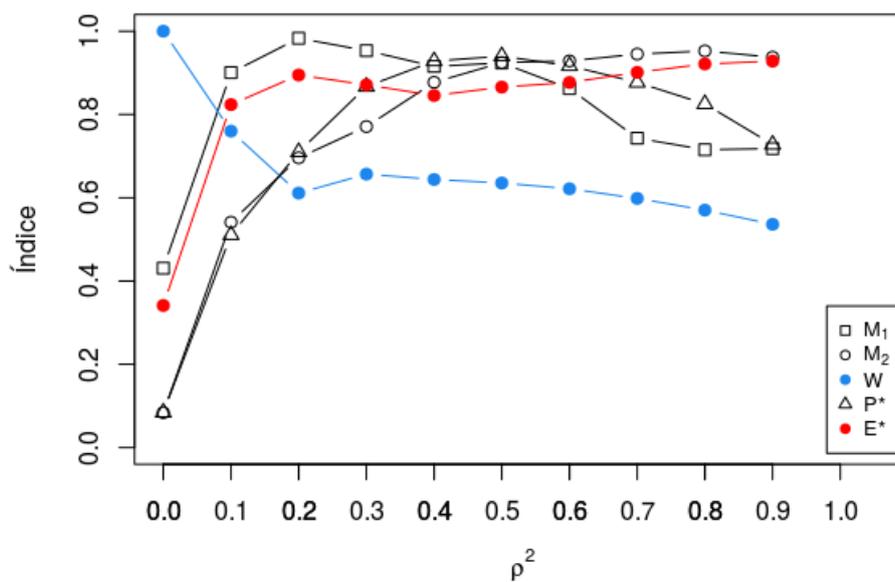
Fonte: Do autor.

Figura 74 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 3$ e $n = 15$



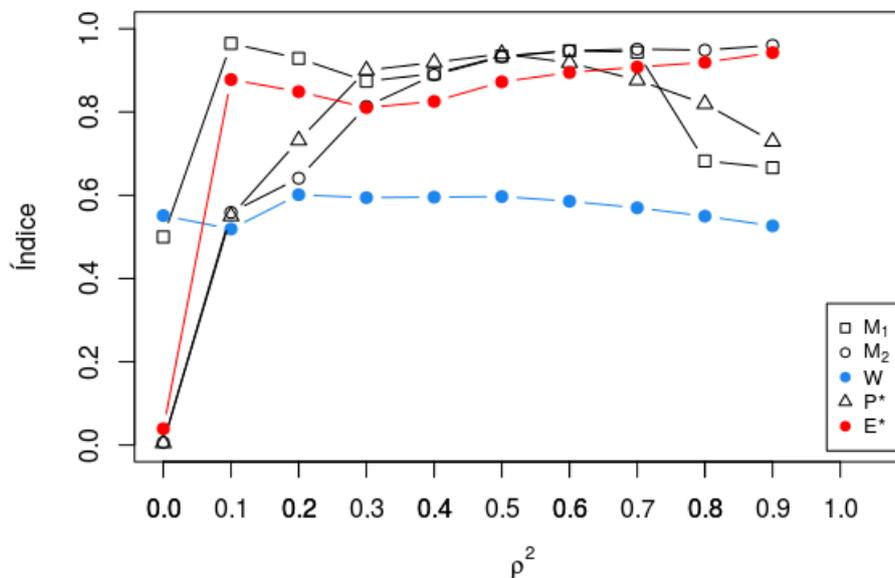
Fonte: Do autor.

Figura 75 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 3$ e $n = 50$



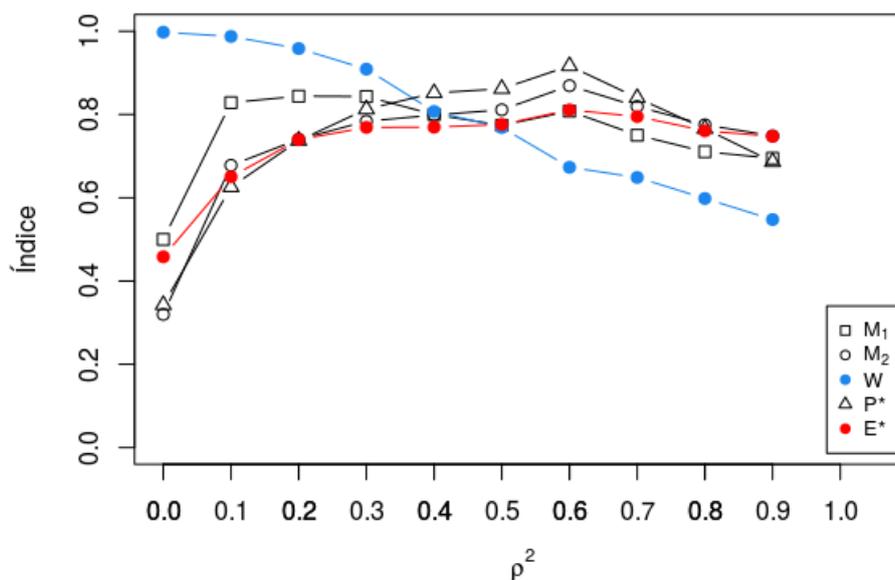
Fonte: Do autor.

Figura 76 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 3$ e $n = 100$



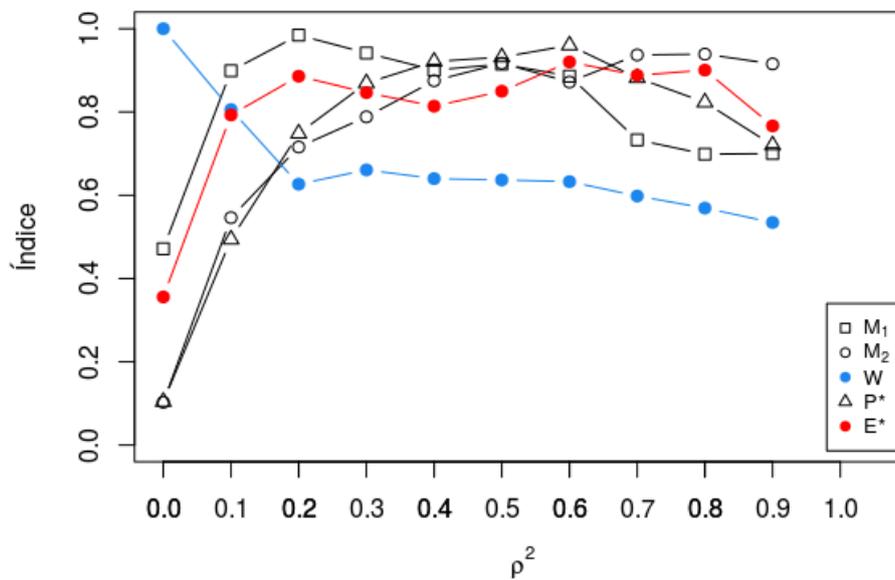
Fonte: Do autor.

Figura 77 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 4$ e $n = 15$



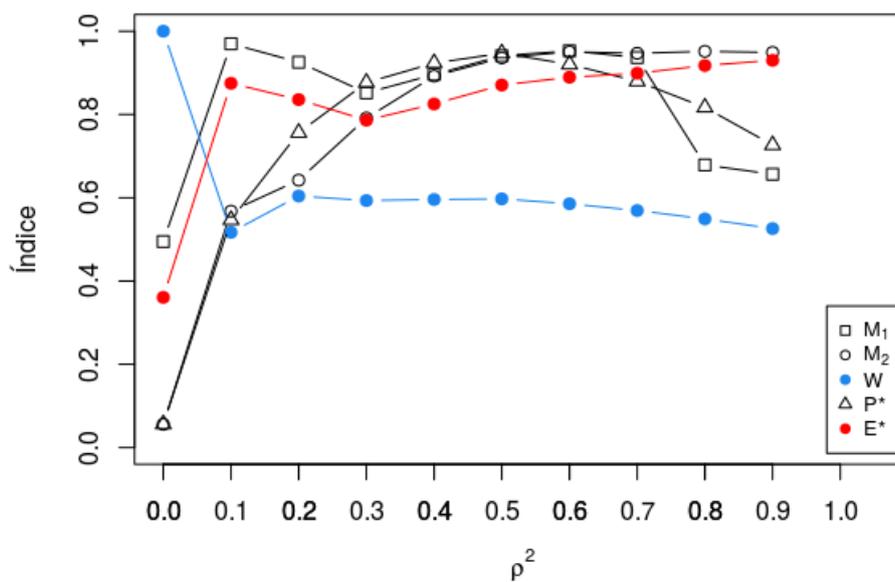
Fonte: Do autor.

Figura 78 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 4$ e $n = 50$



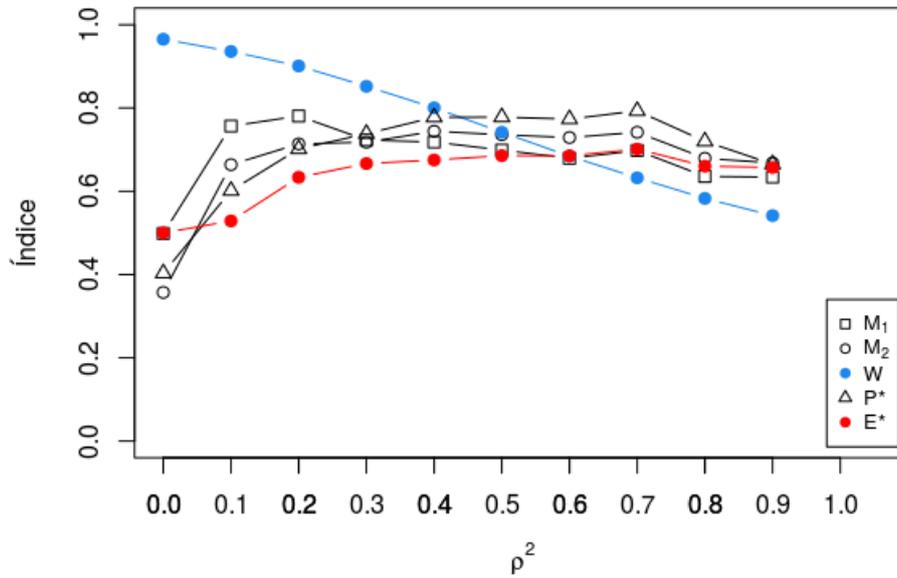
Fonte: Do autor.

Figura 79 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 4$ e $n = 100$



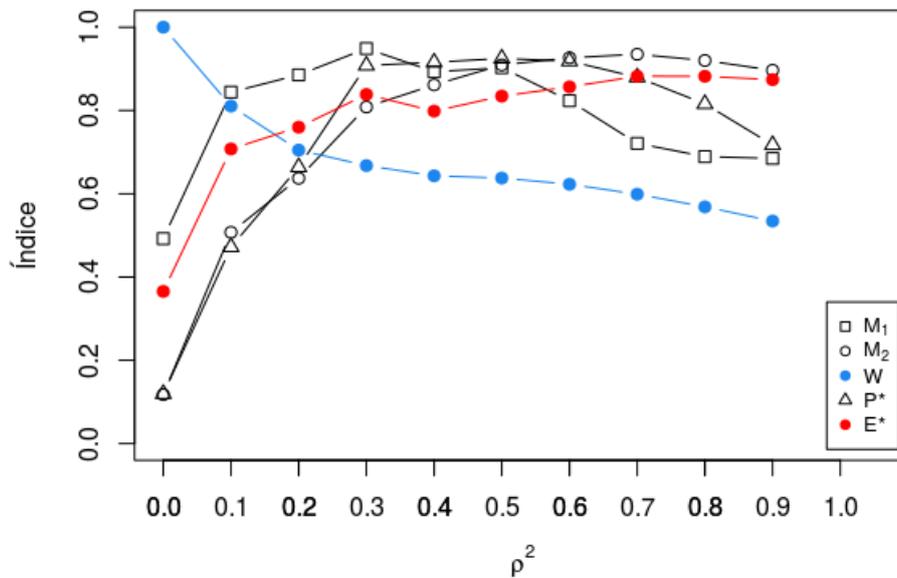
Fonte: Do autor.

Figura 80 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 5$ e $n = 15$



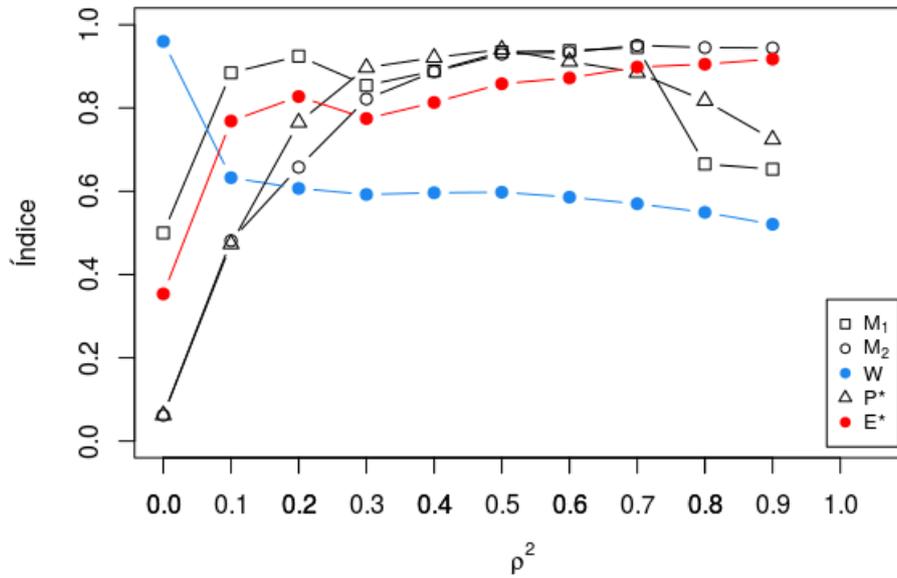
Fonte: Do autor.

Figura 81 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 5$ e $n = 50$



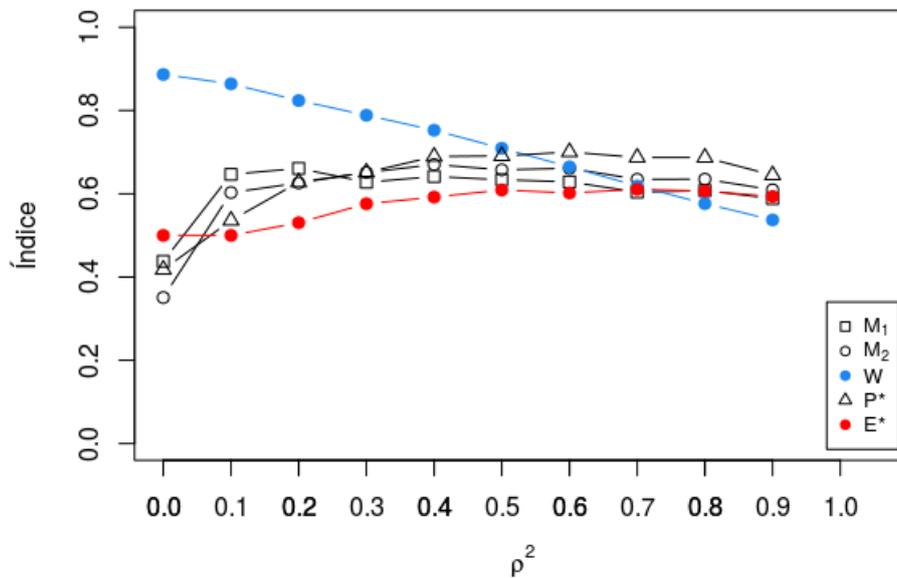
Fonte: Do autor.

Figura 82 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 5$ e $n = 100$



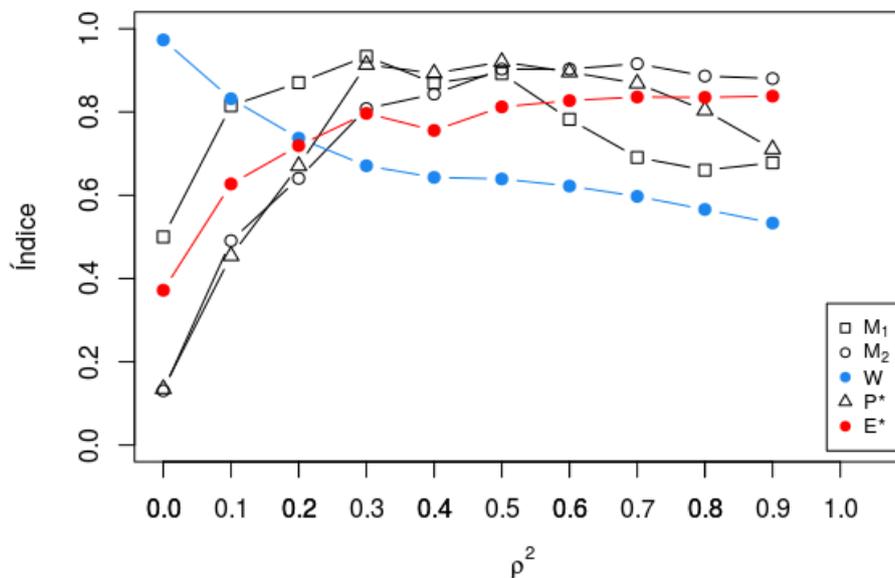
Fonte: Do autor.

Figura 83 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 6$ e $n = 15$



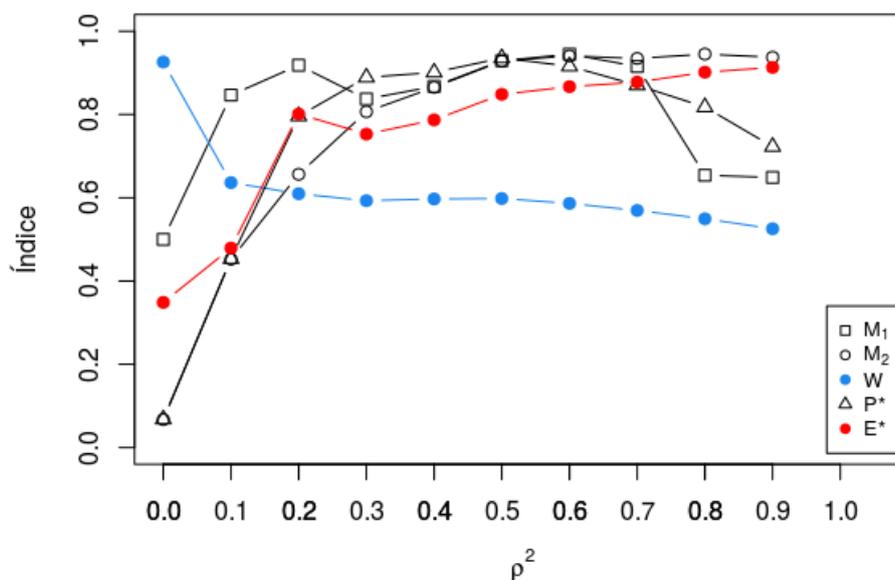
Fonte: Do autor.

Figura 84 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 6$ e $n = 50$



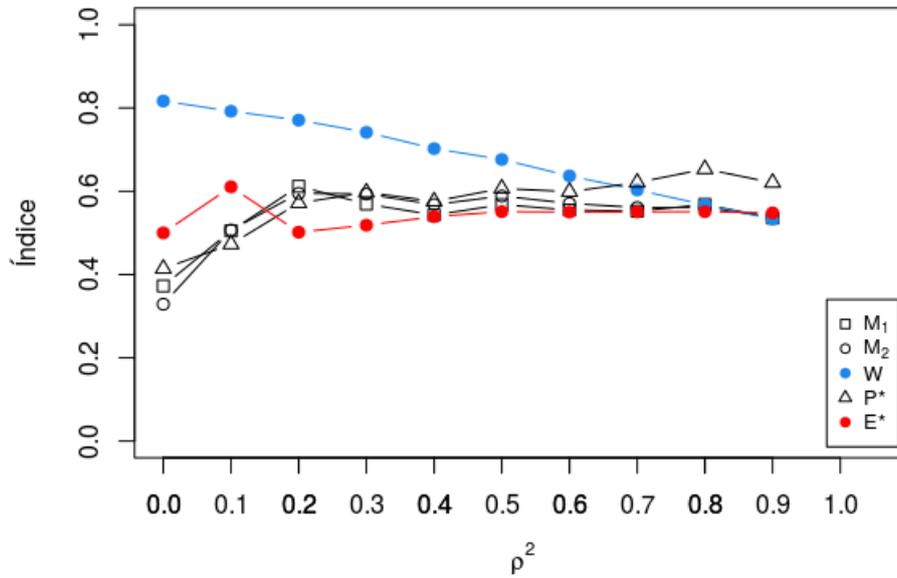
Fonte: Do autor.

Figura 85 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 6$ e $n = 100$



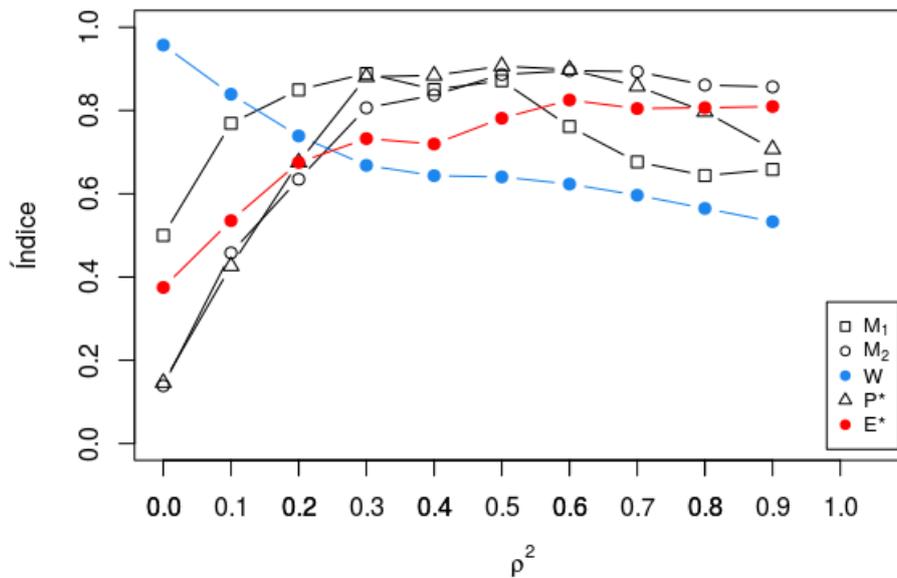
Fonte: Do autor.

Figura 86 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 7$ e $n = 15$



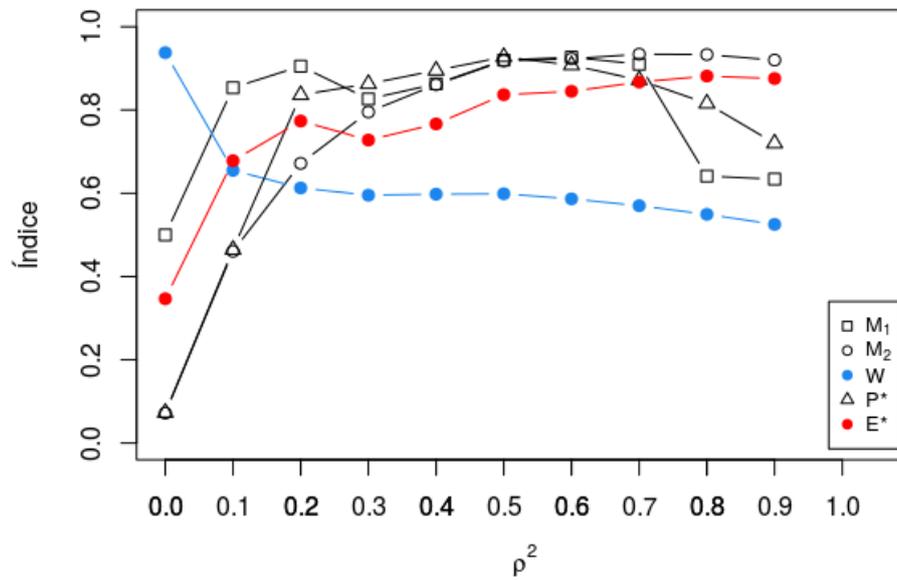
Fonte: Do autor.

Figura 87 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 7$ e $n = 50$



Fonte: Do autor.

Figura 88 – Índices dos estimadores respectivos ao cenário ao modelo de regressão onde $k = 7$ e $n = 100$



Fonte: Do autor.