

# Estimação dos principais direcionadores dos custos operacionais das empresas brasileiras de transmissão de energia elétrica utilizando modelos de regressão e programação linear

Igor Mazzeto Resende Soares

Orientador: Prof. Dr. Marcelo Azevedo Costa

Universidade Federal de Minas Gerais

Departamento de Estatística - ICEx

29 de agosto de 2023

## 1 Introdução

- Justificativa
- Objetivos

## 2 Metodologia

## 3 Estudo de caso - Resultados

- Estatísticas descritivas
- Regressão linear múltipla
- Regressão não linear - modelo gama
- Programação linear
- Modelo de programação linear com aplicação de bootstrap

## 4 Conclusões

- Comparação de resultados
- Conclusões

- Importância da estimação do PMSO para determinação da RAP e consequentemente.
- Aplicação de técnicas estatísticas considerando restrições físicas e técnicas intrínsecas à operação de transmissão de energia elétrica .

## São objetivos deste trabalho:

- Analisar as correlações existentes entre as variáveis e sua representatividade e importância para a realização de previsões.
- Definir um modelo de regressão linear múltipla ou de programação linear para o custo operacional que respeite as restrições impostas pela natureza da operação.
- Discutir os resultados encontrados e apresentar uma alternativa ao regulador para o modelo utilizado na definição dos custos operacionais.

- 1 Pesquisa descritiva quantitativa e que dispõe de um estudo de caso.
- 2 Serão respondidas questões relacionadas ao funcionamento de fenômenos que ocorrem no âmbito do SEB.
- 3 Foi utilizado o software R na versão 4.3.1 “*Beagle Scouts*” com execução em computador pessoal do modelo *MacBook Pro* da marca *Apple*, contendo processador de 2,4 GHz *Intel Core i5 Quad-Core* e memória *RAM* da especificação 8 GB 2133 MHz LPDDR3 no sistema operacional *macOS Ventura* 13.4.1

Tabela: Pacotes do R utilizados no trabalho

Pacote	Uso no trabalho
<i>boot</i>	Aplicação na técnica de <i>bootstrap</i>
<i>dplyr</i>	Manipulação de tabelas
<i>exploreR</i>	Manipulação rápida de dados
<i>forecast</i>	Criação de modelos de previsão
<i>ggplot2</i>	Criação de gráficos
<i>ggalt</i>	Geometrias adicionais para criar gráficos
<i>janitor</i>	Limpeza e preparação de dados
<i>lpsolve</i>	Interface para programação linear
<i>MASS</i>	Ajuste de modelo estatístico de regressão
<i>mvshapiro</i>	Teste de Shapiro
<i>openxlsx</i>	Leitura de planilhas eletrônicas da extensão <i>.xlsx</i>
<i>RcolorBrewer</i>	Implementação de paleta de cores nos gráficos
<i>scales</i>	Formatação das escalas dos gráficos
<i>tidyverse</i>	Ecossistema necessário para funcionamento de pacotes

Fonte: Autor.

- Base de dados fornecidas pela NT nº 097/2022–SRM/ANEEL por meio de uma planilha eletrônica
- Base composta por 125 observações e 20 colunas
- Variável de interesse - custos operacionais compostos pelas contas de pessoal, materiais, serviços de terceiros e outros (*PMSO*)

## Descrição dos direcionadores:

- ①  $X_1$ : Extensão de equipamentos de rede com tensionamento maior de 230 kV
- ②  $X_2$ : Módulos de manobra com tensão igual ou superior a 230 kV
- ③  $X_3$ : Equipamentos de rede com tensionamento maior de 230 kV
- ④  $X_4$ : Potência aparente total, em MVA, de equipamentos de subestação
- ⑤  $X_5$ : Potência reativa total, em Mvar, de equipamentos de subestação
- ⑥  $X_6$ : Equipamentos de subestação com tensão inferior a 230 kV
- ⑦  $X_7$ : Módulos de manobra com tensão inferior a 230 kV
- ⑧  $X_8$ : Extensão de equipamentos de rede com tensionamento maior de 230 kV

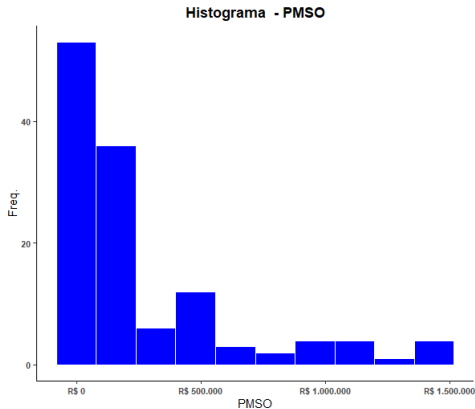


**Tabela:** Estatísticas descritivas das variáveis dependente e independentes

	<i>PMSO</i>	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$	$X_7$	$X_8$
<b>Mínimo</b>	1.273,0	0,0	4,0	7,0	0,0	0,0	0,0	0,0	0,0
<b>1º Quartil</b>	30.523,0	773,6	68,0	56,0	2.348,0	977,0	2,0	17,0	0,0
<b>Mediana</b>	103.467,0	3.517,7	138,0	116,0	7.500,0	4.227,0	5,0	43,0	45,2
<b>Média</b>	270.875,0	4.693,6	270,4	209,0	18.441,0	6.467,0	64,8	313,4	787,9
<b>3º Quartil</b>	262.355,0	6.848,9	345,0	275,0	19.527,0	9.485,0	80,0	375,0	500,7
<b>Máximo</b>	1.439.704,0	18.376,7	1.218,0	763,0	98.256,0	41.208,0	346,0	1.841,0	7.297,6

Fonte: Autor.

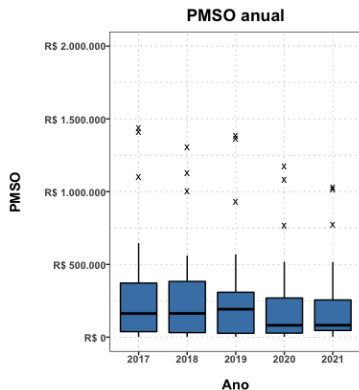
Figura: Histograma PMSO



Fonte: Autor.

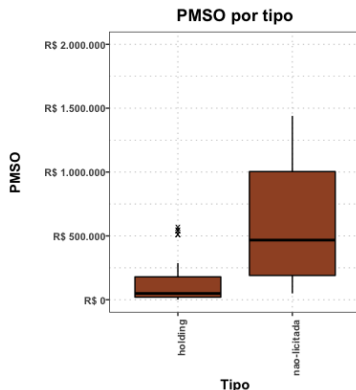
# Boxplots: anual e por tipo de empresa

Figura: Boxplot anual



Fonte: Autor.

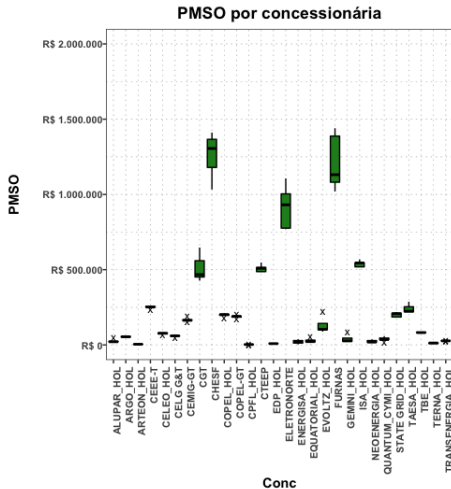
Figura: Boxplot por tipo



Fonte: Autor.

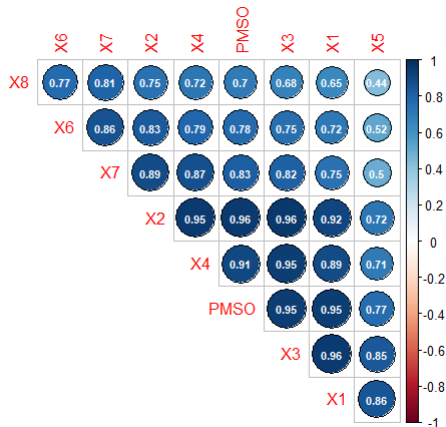
# Boxplot por concessionária

Figura: Boxplot por concessionária



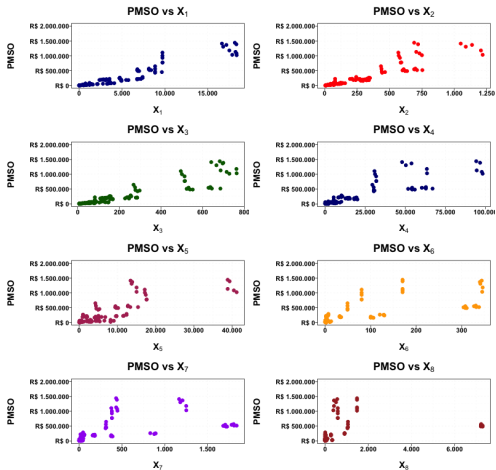
Fonte: Autor.

Figura: Matriz de correlação



Fonte: Autor.

Figura: Gráficos de dispersão



Fonte: Autor.

# Ajuste modelo linear múltiplo

Tabela: Coeficientes do ajuste do modelo

<i><b>Regressão linear</b></i>	
	<i><b>Coeficientes</b></i>
$\beta_0$	-64.039,21
$\beta_1$	13,32
$\beta_2$	794,1
$\beta_3$	169,74
$\beta_4$	<b>-0,52</b>
$\beta_5$	8,86
$\beta_6$	621,45
$\beta_7$	<b>-189,72</b>
$\beta_8$	<b>-8,01</b>

## Características do ajuste:

- $R^2 = 0.92$
- Modelo ajustado:  
$$\hat{y} = -64.03 + 13.32x_1 + 794.1x_2 + 169.7x_3 - 0.5x_4 + 8.8x_5 + 621.4x_6 - 189.7x_7 - 8x_8$$

Fonte: Autor.

# Testes dos pressupostos do modelo de regressão

Tabela: Pressupostos do modelo

Testes dos pressupostos do modelo						
Teste	Pressuposto	Estatística	Hipótese nula	valor-p	Significância	Veredito
Teste-F	Significância	$F = 0,92$	$\beta_1 = \dots \beta_k = 0$	$2,20 \cdot 10^{-16}$	0,05	Rejeitado
Breusch-Pagan	Homoscedasticidade	$LM = 2,25$	$\delta_1 = \dots \delta_k = 0$	$6,50 \cdot 10^{-8}$	0,05	Rejeitado
Durbin-Watson	Autocorrelação	$d = 1,07$	Correlação = 0	$2,46 \cdot 10^{-10}$	0,05	Rejeitado

Fonte: Autor.



# Testes de normalidade para o modelo de regressão

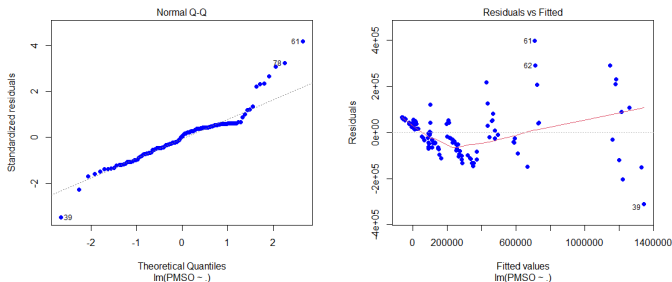
Tabela: Testes de normalidade

Testes de normalidade para os resíduos					
Teste	Estatística	Hipótese nula	valor-p	Significância	Veredicto
Shapiro-Wilk	$W = 0,92$	$H_0 : X \sim N$	$5,54 \cdot 10^{-6}$	0,05	Rejeitado
Anderon-Darling	$A = 2,25$	$H_0 : X \sim N$	$9,76 \cdot 10^{-6}$	0,05	Rejeitado
Kolmogorov-Smirnov	$D = 0,15$	$H_0 : X \sim N$	$2,03 \cdot 10^{-7}$	0,05	Rejeitado

Fonte: Autor.

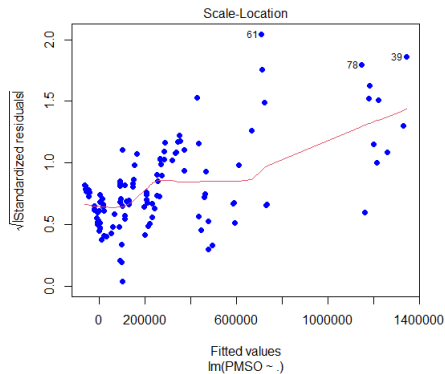
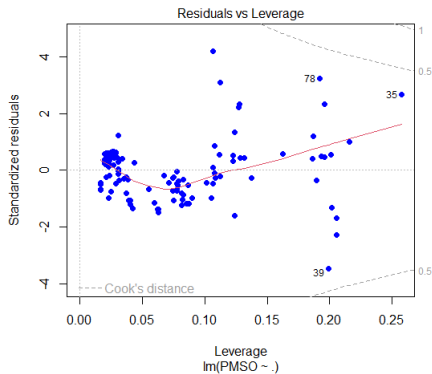
Gráfico de probabilidades  $Q - Q$  e gráfico de resíduos *versus* ajuste

Figura: Gráficos de resíduos



Fonte: Autor.

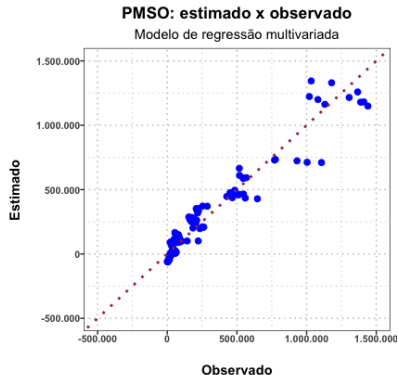
Figura: Gráficos de resíduos



Fonte: Autor.

# Ajustes para modelo de regressão multivariado

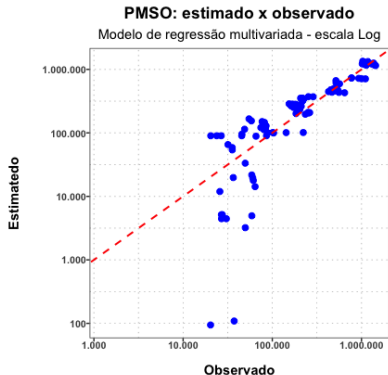
Figura: Modelo multivariado



Fonte: Autor.

$$R^2 : 0.92$$

Figura: Modelo multivariado - Log



Fonte: Autor.

$$R^2 : 0.97$$

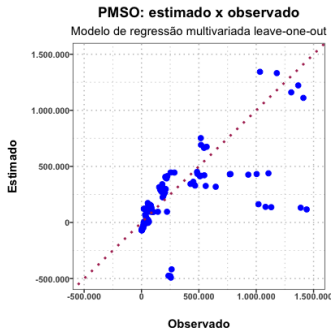
# Ajuste linear múltiplo leave-one-out com validação cruzada

Implementação de validação cruzadas em conjunto com a técnica leave-one-out:

- 1 Criar um novo banco de dados com informações de uma empresa escolhida: base de dados de validação
- 2 Criar outro banco de dados com informações das empresas remanescentes: base de treinamento
- 3 Implementar modelo de regressão linear múltiplo utilizando o banco de dados de treino
- 4 Implementar a técnica de  $AIC$  para escolher o melhor modelo e conjunto de variáveis
- 5 Aplicar o modelo de regressão linear múltiplo na base de dados de validação para estimar os valores de PMSO do modelo
- 6 Armazenar os dados do PMSO estimado para a empresa do banco de dados de validação
- 7 Realizar esse procedimento para todas as 28 empresas

# Ajustes para modelo de regressão *leave-one-out*

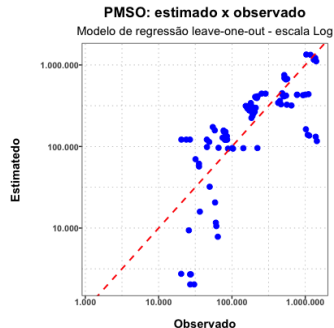
Figura: Modelo multivariado



Fonte: Autor.

$$R^2 : 0.32$$

Figura: Modelo multivariado - Log



Fonte: Autor.

$$R^2 : 0.93$$

# Ajuste - modelo gama

Implementação do modelo gama para lidar com os problemas encontrados anteriormente.

Foi escolhido o modelo gama devido a evidência de proporcionalidade quadrática na variância dos resíduos.

Implementação do modelo com as seguintes abordagens:

- 1 Ajuste com implementação do critério de informação de Akaike (AIC)  
– Regressão gama (1)
- 2 Ajuste com implementação do AIC e partição da base de dados em dois, um conjunto para treino o modelo e outro para aplicação do modelo preditivo (análogo ao realizado na seção anterior) – Regressão gama (2)

# Ajuste - Modelo Gama

Tabela: Ajuste - modelo gama

Resultados - ajuste modelo Gama		
Coefficientes	(1) Step AIC	(2) Step AIC - Partição
$\beta_0$	-1.754,9	-1.273,8
$\beta_1$	19,6	17,9
$\beta_2$	376,0	431,0
$\beta_3$	-	-
$\beta_4$	-	-
$\beta_5$	-	-
$\beta_6$	1.006,3	585,7
$\beta_7$	-	-
$\beta_8$	-25,5	-

Fonte: Autor.

Tabela: Qualidade do ajuste

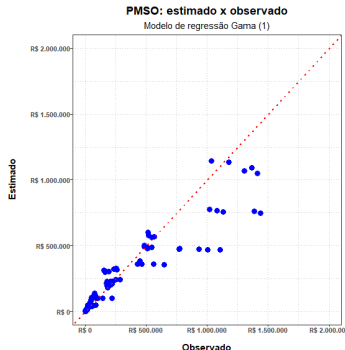
Modelo	$R^2$	Deviance
Regressão Gama (1)	0,82	1,00
Regressão Gama (2)	0,72	1,00

Fonte: Autor.



# Resultados para modelo de regressão gama

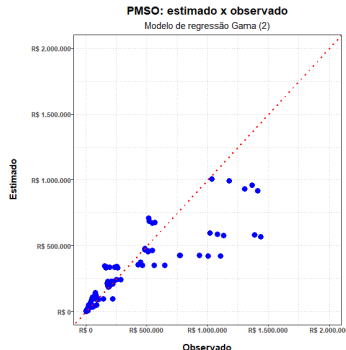
Figura: Modelo (1)



Fonte: Autor.

$R^2 : 0.82$   
p-valor: 1,00

Figura: Modelo (2)

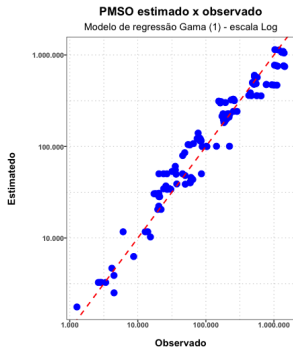


Fonte: Autor.

$R^2 : 0.72$   
p-valor: 1,00

# Resultados para modelo de regressão gama - logaritmo

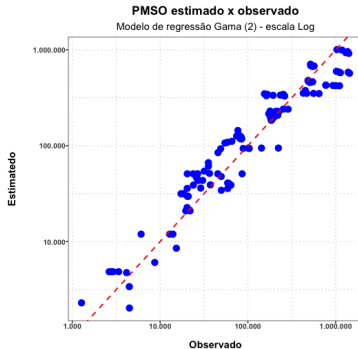
Figura: Modelo (1)



Fonte: Autor.

$$R^2 : 0.94$$

Figura: Modelo (2)



Fonte: Autor.

$$R^2 : 0.92$$

## Abaixo, algumas ponderações para implementação do modelo de programação linear:

- Modelo oriundo da regressão quantílica a  $y = \beta_0 + \beta x + \varepsilon$  com minimização de erros
- Função objetivo: minizar erros com fator  $\tau$
- Restrições:
  - $\beta \geq 0$
  - $\beta_0 = \alpha$
- Particionamento do erro em dois termos
- Implementação de modelo clássico e abordagem *leave-one-out*

**A abordagem escolhida foi a da mediana, portanto  $\tau = 0.5$**

$$\text{minimizar: } \sum_{j=1}^n \tau e_{1j} + (1 - \tau) e_{2j}$$

$$\text{sujeito a: } y_j = \alpha_j + \beta_{kj} x_{kj} + e_{1j} - e_{2j}, \quad \forall j \in \{1, \dots, n\}$$

$$\beta_k \geq 0, \quad \forall k \in \{1, \dots, 8\}$$

$$e_j = e_{1j} + e_{2j}, \quad \forall j \in \{1, \dots, n\}$$

$$N \in \{0, \dots, 128\}$$

# Algoritmo *leave-one-out*

- 1 Remover as observações para a concessionária na qual queremos estimar o PMSO.
- 2 Resolver o modelo de programação linear utilizando o conjunto de dados com as concessionárias que restaram.
- 3 Com o output do modelo contendo os coeficientes de interesse, estimar o PMSO para a concessionária de interesse.
- 4 Armazenar resultado.
- 5 Realizar o procedimento 28 vezes, número de concessionárias presentes no estudo.

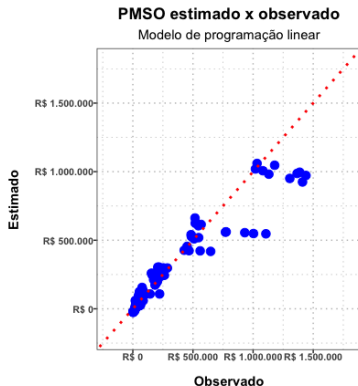
Tabela: Solução - programação linear

Programação Linear		
Coeficientes	Modelo completo	Leave-one-out
$\beta_0$	-29.345,2	27.379,5
$\beta_1$	20,7	20,9
$\beta_2$	450,0	446,3
$\beta_3$	-	-
$\beta_4$	1,8	1,7
$\beta_5$	3,7	3,7
$\beta_6$	-	11,9
$\beta_7$	-	-
$\beta_8$	-	-
$R^2$	0,88	0,79

Fonte: Autor.

# Soluções dos modelos: PMSO estimado x observado

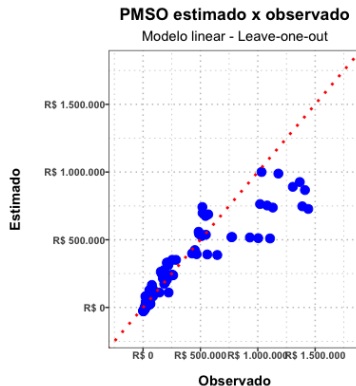
Figura: Modelo linear



Fonte: Autor.

$$R^2 : 0.88$$

Figura: Modelo linear - *leave-one-out*

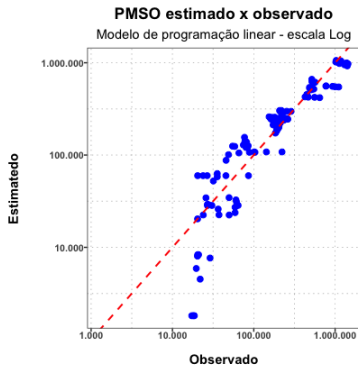


Fonte: Autor.

$$R^2 : 0.79x$$

# PMSO estimado x observado - logaritmo

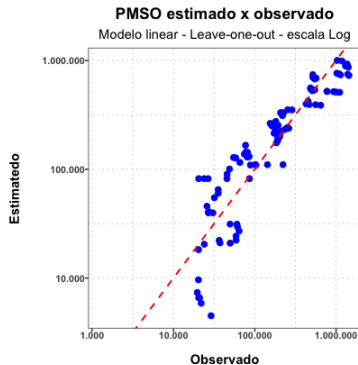
Figura: Modelo linear - log



Fonte: Autor.

$$R^2 : 0.80$$

Figura: *Leave-one-out* - log



Fonte: Autor.

$$R^2 : 0.80$$



# Modelo linear - bootstrap

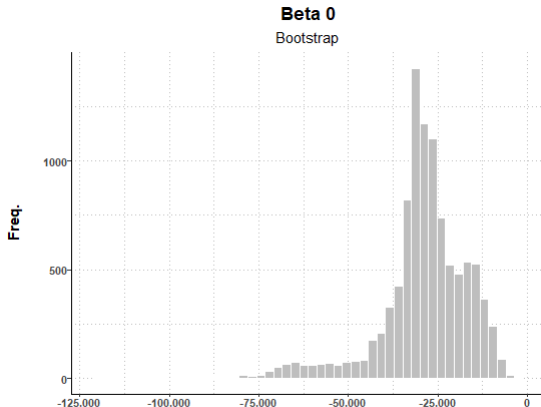
Para a verificar a robustez das restrições do modelo de programação linear foi implementada a técnica de bootstrap e para intervalo de confiança percentílico. O modelo completo foi escolhido por apresentar maior coeficiente de determinação.

Nesse procedimento foram realizadas 10.000 simulações com observações aleatoriamente no software *R*. Abaixo os passos do procedimento:

- 1 É gerada uma nova base de dados utilizando a técnica de reamostragem com reposição
- 2 É implementado o algoritmo SIMPLEX para resolver o modelo com a nova base de dados oriunda da reamostragem
- 3 A solução contendo os coeficientes de regressão é computada
- 4 Os passos de 1 a 3 são realizados 10.000 vezes

# Resultados para procedimento bootstrap - $\beta_0$

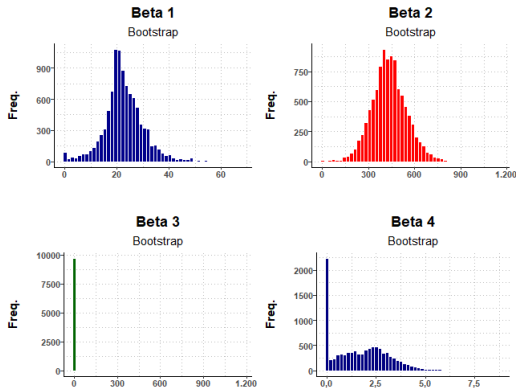
Figura: Histograma de soluções para  $\beta_0$



Fonte: Autor.

# Resultados para procedimento bootstrap - coeficientes

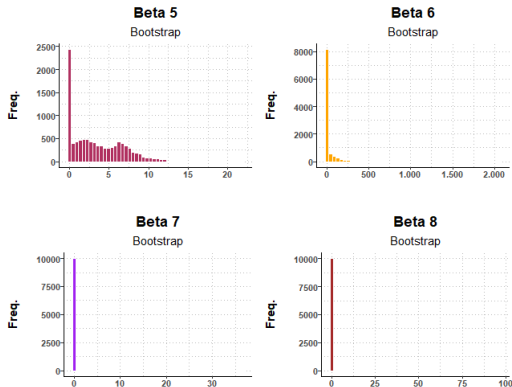
Figura: Histograma de soluções para  $\beta_1, \beta_2, \beta_3, \beta_4$



Fonte: Autor.

# Resultados para procedimento bootstrap - coeficientes

Figura: Histograma de soluções para  $\beta_5, \beta_6, \beta_7, \beta_8$



Fonte: Autor.

# Intervalos de confiança bootstrap

O nível de confiança escolhido para construir o intervalo foi de 95%. Dessa forma, obteve-se os seguintes intervalos com seus limites definidos:

**Tabela:** Intervalo de Confiança Percentílico - Bootstrap

<b>Intervalo de Confiança Percentílico</b>		
<i>Bootstrap</i>		
Coeficientes	2,5%	97,5%
$\beta_0$	-64.426,7	-10.021,2
$\beta_1$	6,3	40,0
$\beta_2$	206,3	672,2
$\beta_3$	0,0	51,3
$\beta_4$	0,0	4,6
$\beta_5$	0,0	11,1
$\beta_6$	0,0	281,0
$\beta_7$	0,0	0,0
$\beta_8$	0,0	0,0

Fonte: Autor.

# Resultados agrupados - finais

- Para comparação dos resultados, foram considerados os modelos implementados com as técnicas de validação cruzada e leave-one-out na escala logarítmica.
- Tal abordagem foi utilizada para atenuação das distorções presentes no banco de dados aqui descritas.
- Dessa forma, no estudo, foram considerados os resultados dos ajustes de modelos de **regressão linear**, **regressão não-linear (modelo gama)** e de **programação linear**.

Tabela: Resultados Finais

Comparativo de $R^2$	
Modelo	$R^2$
Regressão linear - Leave-one-out - Log	<b>0,93</b>
Regressão Gama (2) - Log	<b>0,92</b>
Modelo Linear - Leave-one-out - Log	<b>0,80</b>

Fonte: Autor.

# Conclusões

- O ajuste produzido pelo modelo linear múltiplo se mostrou inadequado por violar os pressupostos de normalidade e as restrições operacionais. Já o ajuste gama se mostrou razoável em relação a capacidade preditiva, porém com problemas de subdispersão.
- O modelo de programação linear se mostrou **o mais efetivo** para estimar o *PMSO* pois respeita as restrições operacionais e apresenta boa capacidade preditiva.
- Cinco entre as oito variáveis mostraram-se relevantes quando aplicado o modelo linear.
- Atualmente a ANEEL implementa uma metodologia que **utiliza direcionadores de custo redundantes podendo comprometer ou enviesar as estimativas de eficiência**.
- **Contribuição:** A abordagem inovadora de integrar restrições de uma regressão em um modelo de programação linear mostra a versatilidade da técnica, permitindo utilizá-la uma ampla gama de cenários do mundo real, destacando seu potencial para otimização e melhorias na eficiência em várias áreas.