



PONTIFÍCIA UNIVERSIDADE CATÓLICA DO PARANÁ
ESCOLA POLITÉCNICA
PROGRAMA DE PÓS-GRADUAÇÃO EM ENGENHARIA DE PRODUÇÃO E
SISTEMAS (PPGEPS)

FRANCHESCO SANCHES DOS SANTOS

EXPLORANDO A EFICIÊNCIA DOS MODELOS DE PREVISÃO DE SÉRIES
TEMPORAIS NO ABASTECIMENTO DE ÁGUA

CURITIBA
2023

FRANCHESCO SANCHES DOS SANTOS

**EXPLORANDO A EFICIÊNCIA DOS MODELOS DE PREVISÃO DE SÉRIES
TEMPORAIS NO ABASTECIMENTO DE ÁGUA**

Projeto de Pesquisa de Mestrado apresentado ao Programa de Pós-Graduação em Engenharia de Produção e Sistemas (PPGEPS). Área de concentração: Automação e Controle de Sistemas, da Escola Politécnica, da Pontifícia Universidade Católica do Paraná, como requisito parcial à obtenção do título de Mestre em Engenharia de Produção e Sistemas.

Orientador: Dr. Leandro dos Santos Coelho
Coorientadora: Dr. Viviana Cocco Mariani

CURITIBA
2023

FRANCHESCO SANCHES DOS SANTOS

**EXPLORANDO A EFICIÊNCIA DOS MODELOS DE PREVISÃO DE
SÉRIES TEMPORAIS NO ABASTECIMENTO DE ÁGUA**

Dissertação apresentada ao Programa de Pós-Graduação em Engenharia de Produção e Sistemas (PPGEPS). Área de concentração: Gerência de Produção e Logística, da Escola Politécnica, da Pontifícia Universidade Católica do Paraná, como requisito parcial à obtenção do título de Mestre em Engenharia de Produção e Sistemas.

COMISSÃO EXAMINADORA

Dr. Leandro dos Santos Coelho

Orientador

Pontifícia Universidade Católica do Paraná

Dr. Viviana Cocco Mariani

Coorientadora

Pontifícia Universidade Católica do Paraná

Convidado A

Membro Externo

Instituição A

Convidado B

Banca

Instituição B

Curitiba, 26 de julho de 2023

Ao examinar as séries temporais, vejo a evidência da existência de Deus na perfeita ordem e regularidade que caracterizam o tempo. Cada ponto de dados e cada instante são testemunhos do Seu controle absoluto sobre todas as coisas.

Agradecimentos

Primeiramente, expresso minha gratidão a Deus por todas as bênçãos recebidas, pois foi Ele quem abriu caminhos e me deu forças para superar esse desafio, tornando-o possível.

À minha família, sou grato pelo apoio incondicional e pelo estímulo constante para seguir em frente com determinação, buscando sempre alcançar novos patamares.

Agradeço ao professor Leandro dos Santos Coelho pela oportunidade de trabalhar ao seu lado e compartilhar seus conhecimentos e experiências ao longo do meu mestrado. Sua orientação contribuiu significativamente para o meu crescimento profissional e pessoal, tornando este trabalho uma realidade.

À professora Viviana Cocco Mariani, agradeço pela disponibilidade e paciência em me auxiliar nas minhas dificuldades, utilizando seu conhecimento para contribuir com o desenvolvimento da pesquisa.

Quero expressar minha gratidão à equipe da Pontifícia Universidade Católica do Paraná (PUCPR) e aos demais professores, especialmente à secretária Denise da Mata Medeiros (PPGEPS), pela paciência, carinho e apoio prestados em diversas ocasiões, sem medir esforços.

Aos meus amigos, que sempre torceram por mim, e aos novos amigos que conquistei ao longo dessa jornada, agradeço por compartilharmos momentos de alegria nessa batalha.

Sou grato ao investimento em bolsas de estudo concedidas pela Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES), que possibilitou a conclusão dessa etapa da minha carreira profissional e acadêmica.

*Se vi mais longe, foi por estar de pé
sobre ombros de gigantes*

- Sir Isaac Newton

Resumo

Este estudo aborda a importância estratégica da previsão precisa da demanda de água como uma ferramenta para a gestão eficiente dos recursos hídricos em um cenário competitivo. O problema identificado é a falta de previsões precisas, o que dificulta a tomada de decisões estratégicas no abastecimento de água. A solução proposta é o uso de modelos avançados de previsão de séries temporais para melhorar a precisão das previsões de demanda. Com base em uma revisão abrangente da literatura existente, diferentes métodos e abordagens utilizados na previsão de séries temporais no contexto do abastecimento de água são analisados. O estado da arte é explorado para identificar os modelos mais eficazes e as melhores práticas na área. Métodos e produtos específicos são propostos com base no estado da arte, levando em consideração variáveis exógenas, sazonalidade dos dados e utilizando modelos autorregressivos integrados de médias móveis (ARIMA), técnicas de boosting como XGBoost (Extreme Gradient Boosting), LightGBM (Light Gradient Boosting Machine) e regressão linear. Além disso, também é considerado o uso de modelos baseados em Random Forest Regression (RFR). Os resultados obtidos por meio da aplicação desses métodos e produtos propostos são analisados e comparados utilizando métricas de desempenho, como o erro percentual absoluto médio simétrico (sMAPE), o erro absoluto médio (MAE) e a raiz do erro médio quadrático relativo (RRMSE). Essas descobertas fornecem informações valiosas sobre a eficácia dos modelos de previsão de séries temporais no abastecimento de água e contribuem para uma tomada de decisão mais informada e eficiente nessa área.

Palavras-chave: Previsão, Economia de água, Séries temporais, Modelos de Previsão.

Abstract

This study addresses the strategic importance of accurate water demand forecasting as a tool for efficient water resource management in a competitive landscape. The identified problem is the lack of accurate predictions, which hinders strategic decision-making in water supply. The proposed solution is the use of advanced time series forecasting models to improve demand prediction accuracy. Based on a comprehensive review of existing literature, different methods and approaches used in water supply time series forecasting are analyzed. The state-of-the-art is explored to identify the most effective models and best practices in the field. Specific methods and products are proposed based on the state-of-the-art, considering exogenous variables, data seasonality, and utilizing autoregressive integrated moving average (ARIMA) models, boosting techniques like XGBoost (Extreme Gradient Boosting) and LightGBM (Light Gradient Boosting Machine), and linear regression. Additionally, the use of models based on Random Forest Regression (RFR) is also considered. The results obtained through the application of these proposed methods and products are analyzed and compared using performance metrics such as symmetric mean absolute percentage error (sMAPE), mean absolute error (MAE), and root relative mean squared error (RRMSE). These findings provide valuable insights into the effectiveness of time series forecasting models in water supply and contribute to more informed and efficient decision-making in this field.

Keywords: Forecasting, Water savings, Time series, Systematic literature review.

Lista de Figuras

1	Paradigma de aprendizado de máquina	2
2	Mapa das Etapas	5
3	Estrutura da dissertação	10
4	Detecção de anomalias	11
5	Exemplo de séries temporais	13
6	Processo estocástico	13
7	Fluxograma do problema de pesquisa	14
8	Etapas da Revisão	15
9	Palavras-chave mais populares na Scopus	17
10	Palavras-chave mais populares na Web of Science	17
11	Analise das quantidades de artigos em relação aos anos	18
12	Relação de autores entre artigos publicados	20
13	Ligaçāo bibliográfica entre os autores	21
14	Mapa mundial da publicação de artigos em todo o mundo	22
15	Áreas de aplicāo do tema	22
16	Comparāo dos modelos AR e ARX	31
17	Modelo MA(7)	35
18	ARMA (7,7)	35
19	ARIMA (7,1,7)	36
20	SARIMA (7, 1, 7)(2, 1, 1) ₁₂	37
21	Comparāo entre ARIMAX e SARIMAX	39
22	Corelação de Pearson	40
23	Régressāo linear LT01 vs PT01 correlação 98%	41
24	Régressāo linear (LR) um passo a frente	42
25	Régressāo da Floresta Aleatória (RFR)	43
26	Esquema da Floresta Aleatória	43
27	Impulsionando gradiente com XGBoost e LightGBM	44
28	Compara-se o crescimento em folha com o crescimento em nível	47
29	A performance da régressāo utilizando XGBoost e LightGBM é comparada	48
30	Decomposição STL	58
31	Violino no nível do reservatório	59
32	Violino da vazāo de recalque	59
33	Autocorrelação e Autocorrelação parcial	61
34	Ruído branco	63
35	Análise comparativa dos modelos utilizando gráficos de violino	67
36	Demandas Médias das Variáveis de Fluxo	72

37	Comparação dos modelos AR, ARX e MA, 1 dia à frente	88
38	Comparação dos modelos AR, ARX e MA, 7 dias à frente	88
39	Comparação dos modelos AR, ARX e MA, 14 dias à frente	88
40	Comparação dos modelos AR, ARX e MA, 30 dias à frente	89
41	Comparação dos modelos ARMA e ARIMA, 1 dia à frente	89
42	Comparação dos modelos ARMA e ARIMA, 7 dias à frente	90
43	Comparação dos modelos ARMA e ARIMA, 14 dias à frente	90
44	Comparação dos modelos ARMA e ARIMA, 30 dias à frente	90
45	Comparação dos modelos ARIMAX, SARIMA e SARIMAX, 1 dia à frente	91
46	Comparação dos modelos ARIMAX, SARIMA e SARIMAX, 7 dias à frente	91
48	Comparação dos modelos ARIMAX, SARIMA e SARIMAX, 30 dias à frente	92
47	Comparação dos modelos ARIMAX, SARIMA e SARIMAX, 14 dias à frente	92

Lista de Tabelas

1	Cruzamento de palavras-chave através da aplicação de filtros de ano e de linguagem	18
2	Fator de impacto	20
3	Áreas e seus valores respetivos de artigos em cada área.	23
4	Descrição estatística dos dados com o filtro aplicado das 18h às 21h	55
5	Teste Nemenyi	65
6	Demanda de água	72
7	Comparação dos modelos de previsão com as métricas de desempenho treino	83
8	Comparação dos modelos de previsão com as métricas de desempenho teste	84
9	Comparação dos modelos de previsão com as métricas de desempenho va-lidação	85
10	Comparação dos modelos de previsão com as métricas de desempenho inteiro	86
11	Comparação dos modelos Ljung Box	87

Lista de Abreviaturas e Siglas

AdaBoost	Impulso ou Estímulo Adaptativo (do inglês <i>Adaptive Boosting</i>)
AR	Auto-Regressivo
ARIMA	Média Móvel Integrada Auto-Regressiva (do inglês <i>Autoregressive Integrated Moving Average</i>)
ARIMAX	Média Móvel Integrada Auto-Regressiva com Regressores Exógenos (do inglês <i>Autoregressive Integrated Moving Average with Exogenous Regressors</i>)
ARMA	Média Móvel Auto-Regressiva (do inglês <i>Autoregressive Moving Average</i>)
ARX	Auto-Regressivo com Variável Exógena (do inglês <i>Autoregressive with Exogenous Inputs</i>)
BrownBoost	Algoritmo de Aumento
CNN	Rede Neural Convolucional (do inglês <i>Convolutional Neural Network ou ConvNet</i>)
DBN	Rede de Crenças Profundas (do inglês <i>Deep Belief Network</i>)
EFB	Pacote de Características Exclusivas (do inglês <i>Exclusive Feature Bundling</i>)
FT	Flow Transmitter (Transmissor de Fluxo)
Hz	Hertz
INMET	Instituto Nacional de Meteorologia
LGBMRegressor	Regressão Light GBM
Light GBM	Máquina de Impulso de Gradiente Leve (do inglês <i>Light Gradient Boosting Machine</i>)
LogitBoost	Técnicas de Regressão Logística
LPBoost	Reforço da Programação Linear (do inglês <i>Linear Programming Boosting</i>)
LR	Regressão Linear (do inglês <i>Linear Regression</i>)
LSTM	Memória de Longo Curto Prazo (do inglês <i>Long Short-Term Memory</i>)
m^3	Metros Cúbicos
m^3/h	Metros Cúbicos por Hora

MA	Média Móvel (do inglês <i>Moving Average</i>)
MadaBoost	Modificando o Sistema de Ponderação do AdaBoost
MAE	Erro Médio Absoluto (do inglês <i>Mean Absolute Error</i>)
MAPE	Erro Percentual Médio Absoluto (do inglês <i>Mean Absolute Percentage Error</i>)
mca	Metros Coluna de Água
ML	Aprendizado de Máquina (do inglês <i>Machine Learning</i>)
mm	Milímetros
MSE	Erro Médio Quadrático (do inglês <i>Mean Squared Error</i>)
PR	Estado do Paraná
RBAL	Recalque Bairro Alto
RFR	Regressão de Floresta Aleatória (do inglês <i>Random Forest Regression</i>)
RMSE	Erro de Raiz Média Quadrática (do inglês <i>Root Mean Squared Error</i>)
RNN	Rede Neural Recorrente (do inglês <i>Recurrent Neural Network</i>)
RRMSE	Raiz do Erro Médio Quadrático Relativo (do inglês <i>Root of the Relative Mean Square Error</i>)
SANEPAR	Companhia de Saneamento do Paraná
SARIMA	Auto-Regressivos Integrados de Médias Móveis com Sazonalidade (do inglês <i>Seasonal Auto-Regressive Integrated Moving Averages</i>)
SARIMAX	Média Móvel Integrada Auto-Regressiva Sazonal com Regressores Exógenos (do inglês <i>Seasonal Auto-Regressive Integrated Moving Average with Exogenous Regressors</i>)
sMAPE	Erro Percentual Absoluto Médio Simétrico (do inglês <i>Symmetric Mean Absolute Percentage Error</i>)
SVM-VAR	Máquinas de Vetor de Suporte - Vetores Auto-Regressivos
TotalBoost	Impulso Total
XGBoost	Impulso Gradiente Extremo (do inglês <i>eXtreme Gradient Boosting</i>)
XGBRegressor	Regressão XGBoost

Sumário

1	Introdução	1
1.1	Contexto da pesquisa	1
1.1.1	Motivação da pesquisa	3
1.2	Objetivo geral	3
1.2.1	Objetivos específicos e questão de pesquisa	3
1.3	Descrição do problema	4
1.4	Procedimentos metodológicos	5
1.4.1	Etapas da pesquisa	5
1.5	Justificativa da pesquisa	7
1.5.1	Contribuições	7
1.6	Estrutura do trabalho	9
2	Referencial	9
2.1	Detecção de anomalias	10
2.2	Revisão sistemática da literatura	12
2.3	Problematização da Revisão	14
2.4	Metodologia	15
2.5	Resultados da Busca de Revisão	16
2.6	Principais conclusão	24
3	Base Teórica	25
3.1	Métricas de Avaliação de Modelos	25
3.1.1	Erro Quadrático Médio Raiz (RMSE)	25
3.1.2	Raiz do Erro Médio Quadrático Relativo (RRMSE)	26
3.1.3	Erro Absoluto Médio (MAE)	28
3.1.4	Erro Percentual Absoluto Médio (MAPE)	28
3.1.5	Erro Percentual Absoluto Médio Simétrico (sMAPE)	29
3.2	Modelos de Séries Temporais Univariados	30
3.2.1	Componente Autorregressivo	31
3.2.2	AR(0): Ruído branco	32
3.2.3	AR(1): Caminhadas aleatórias e Oscilações	32
3.2.4	AR(p): Termos de ordem superior	33
3.2.5	Média Móvel	33
3.2.6	Modelos ARMA e ARIMA	35
3.2.7	ARIMA	36
3.2.8	SARIMA	37

3.3	Modelos de Série Temporal Multivariada	37
3.3.1	ARIMAX e SARIMAX	38
3.4	Modelos de Aprendizado de Máquina Supervisionados	38
3.4.1	Régressão Linear (LR)	39
3.4.2	Definição do modelo	41
3.4.3	Floresta Aleatória (Random Forest)	42
3.4.4	Gradient Boosting (como XGBoost, LightGBM)	43
3.4.5	O Gradiente em Gradiente de Boosting (Reforço)	44
3.4.6	Algoritmos de boosting de gradiente	45
3.4.7	A diferença entre XGBoost e LightGBM	45
3.5	Estudo de Caso Empírico	47
3.5.1	Definição do problema	47
3.5.2	Coleta de dados	48
3.5.3	Análise exploratória dos dados	49
3.5.4	Escolha do modelo	49
3.5.5	Divisão dos dados	50
3.5.6	Ajuste do modelo	50
3.5.7	Avaliação do modelo	51
3.5.8	Previsões Futuras	52
3.5.9	Monitoramento e Ajuste Contínuo	53
3.5.10	Principais Conclusão	54
4	Resultados	54
4.1	Planejamento do Problema	55
4.1.1	Análise Exploratória dos dados (EDA)	55
4.1.2	Múltiplas entradas e saída única (MISO)	56
4.1.3	Decomposição STL	57
4.1.4	Separação dos dados	63
4.1.5	Modelagem e Seleção do Modelo	64
4.1.6	Horizonte	64
4.1.7	Previsão e Avaliação	64
4.1.8	Relatório dos Resultados	65
4.1.9	Comparação dos modelos	66
4.2	Estudo de Caso Empírico Resultado	67
4.2.1	Descrição do sistema de abastecimento de água	68
4.2.2	Análise exploratória dos dados	69
4.2.3	Questões de pesquisa 1 a 4	70
4.2.4	Questão de pesquisa 5	71

4.2.5 Discussão geral e conclusões	73
5 Conclusões	74
5.1 Limitações da Pesquisa e Propostas Futuras	75
Referências	76
A Apêndice - Comparaçao dos modelos de previsão de series temporais média de 24h	82
B Apêndice - Comparaçao dos modelos de previsão com o método Ljung-Box	87
C Apêndice - Modelos AR(7), ARX (7) e MA (7) 24h	88
D Apêndice - Modelos ARMA(7,7) e ARIMA (7,1,7) 24h	89
E Apêndice - Modelos ARIMAX (7,1,7), SARIMA (7,1,7) (2,1,0,12) e SARIMAX (7,1,7) (2,1,0,12) 24h	91

1 Introdução

A previsão precisa da demanda de água é fundamental para um planejamento eficiente e sustentável do abastecimento hídrico em uma determinada região. Neste estudo, são empregados métodos avançados, como Gradiente, Regressão e ARIMA, para realizar previsões precisas da demanda diária de água. Os resultados obtidos fornecem insights valiosos e contribuem para um melhor entendimento das tendências e padrões de consumo de água, permitindo um planejamento mais eficaz para atender às necessidades da população.

Este capítulo apresenta o conteúdo abordado nesta dissertação, que se concentra na utilização de modelos de Aprendizado de Máquina (ML) para prever futuramente os dados coletados pela SANEPAR. Os dados coletados referem-se ao abastecimento de água no bairro Alto durante o período de 2018 a 2020, quando ocorreu uma escassez que afetou toda a população da capital paranaense.

Dentro do contexto de análise de séries temporais e tomada de decisão, são explorados modelos de ML para aplicação nesses dados. Por meio de uma revisão sistemática da literatura, são identificados e tabulados os modelos clássicos mais comumente utilizados para análise de séries temporais.

Nesta dissertação, busca-se desenvolver previsões precisas e confiáveis para o abastecimento de água no bairro Alto. Com base nas informações coletadas e na aplicação dos modelos, espera-se obter insights valiosos para auxiliar na tomada de decisões estratégicas e no planejamento eficiente do abastecimento hídrico na região.

Além da revisão da literatura, serão apresentados os métodos e técnicas utilizados para a análise dos dados, bem como os resultados obtidos. Busca-se contribuir significativamente para a área de análise de séries temporais aplicada ao abastecimento de água, permitindo uma melhor compreensão dos padrões de consumo e aprimorando a eficiência dos processos de tomada de decisão relacionados ao fornecimento de água no bairro Alto.

1.1 Contexto da pesquisa

Ribeiro et al. (2021) A necessidade de desenvolvimento do planejamento estratégico no mundo corporativo e no dia-a-dia torna a análise de séries temporais e previsões valiosas ferramentas para apoiar o processo de tomada de decisão a curto, médio e longo prazo. Devido às não linearidades, sazonalidade, tendência e ciclicidade nos dados temporais, o desenvolvimento de modelos de previsão eficientes é uma tarefa desafiadora.

No conjunto de dados da SANEPAR, há um volume significativo no consumo de água e, com as interrupções que a cidade tem enfrentado, é necessário analisar os

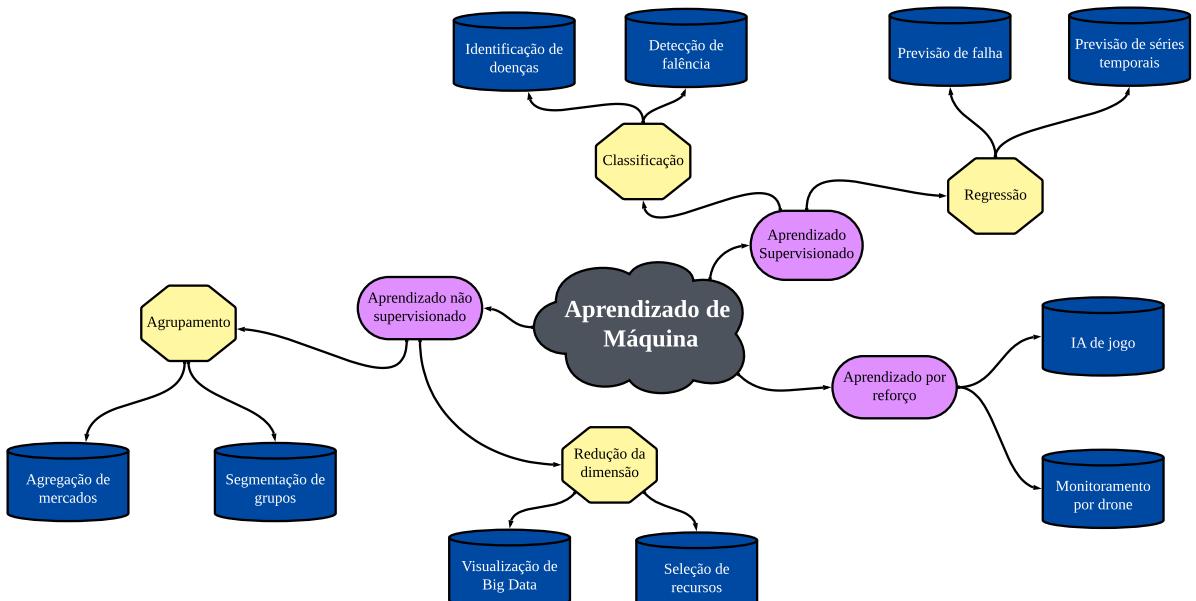
dados para compreender melhor os padrões de interrupção no abastecimento e os picos de consumo ao longo das horas e dias.

Nesta dissertação, será realizada uma revisão sistemática de modelos preditivos para avaliar o melhor modelo que pode ser utilizado e como ele pode ser validado para prever a escassez de água. Essas análises serão feitas utilizando a linguagem de programação *Python*.

A abordagem deste trabalho consiste em explorar o conceito de séries temporais e sua aplicação no campo do aprendizado de máquina. Os dados de séries temporais referem-se a dados coletados e armazenados ao longo do tempo, permitindo que observadores identifiquem anomalias nos dados. A classificação dos dados por ano ou dia é essencial na análise de séries temporais, e se os dados forem atribuídos aleatoriamente, pode ser mais desafiador fazer previsões e tomar decisões com base nos dados coletados.

É importante destacar que a análise de médias pode ser enganosa se não forem excluídos os valores discrepantes, também conhecidos como “outliers”. Esses valores discrepantes podem levar a resultados extremamente altos ou baixos que não refletem a realidade. O campo do aprendizado de máquina abrange várias áreas, conforme ilustrado na Figura 1. Serão explorados os diferentes componentes do aprendizado de máquina e como eles podem ser aplicados em diversos contextos.

Figura 1: Paradigma de aprendizado de máquina



Fonte: Elaboração própria

1.1.1 Motivação da pesquisa

A motivação desta pesquisa é baseada na situação enfrentada por Curitiba e região metropolitana, conforme apontado por Vasconcelos (2020). A região passou por um rodízio de abastecimento de água, com períodos de 36 horas com água seguidos por 36 horas sem abastecimento. A média geral dos reservatórios na região estava em torno de 27,96% de sua capacidade. Além disso, a quantidade de chuva nos anos anteriores, de 2020, foi de 1.704 mm , superando a média anual de precipitação de 1.490 mm .

Diante dessa situação, a pesquisa tem como abordagem principal a escassez de água, que pode ser associada a condições de seca. A partir dos dados coletados pela SANEPAR, é possível realizar uma análise mais detalhada, com o objetivo de prever e evitar a ocorrência de escassez de água, que foi registrada como uma anomalia em 2020. Com o retorno das chuvas, houve um aumento nos níveis dos reservatórios, o que torna essencial a análise e previsão dos dados para um melhor planejamento e gerenciamento do abastecimento de água na região.

1.2 Objetivo geral

O objetivo desta pesquisa consiste em identificar o melhor modelo de séries temporais para abordar o problema da escassez de água em Curitiba. Durante a dissertação, diversos modelos de regressão foram avaliados, com ênfase nos modelos baseados em *gradient boosting*, reconhecidos na literatura por sua eficácia na previsão de séries temporais. Os principais modelos investigados incluem o ARIMA e suas variantes mais atualizadas. Além da previsão, também serão realizadas análises de anomalias nos dados, buscando compreender as causas subjacentes a essas ocorrências.

1.2.1 Objetivos específicos e questão de pesquisa

Neste estudo, busca-se identificar e compreender possíveis anomalias nos dados, bem como investigar as causas por trás dessas ocorrências. O objetivo é responder às perguntas de pesquisa relacionadas a essas anomalias.

Q 1 Qual é a adequação da pressão atual para atender à demanda diária?

Q 2 Qual é o volume mínimo de água necessário no reservatório para evitar o acionamento das bombas durante o horário de pico?

Q 3 Qual é a vazão ótima para atender à demanda diária?

Q 4 Como encontrar o ponto de equilíbrio entre a demanda e a vazão?

Q 5 Qual é o impacto do acionamento das bombas durante o horário de pico?

- a. Qual é o nível ideal no reservatório para evitar a ativação das bombas da SANE-PAR durante o período de maior demanda, das 18h às 21h, sem comprometer o abastecimento de água para a população? Além disso, como variam as médias das vazões nos horários críticos (18h às 21h) para as diferentes estações do ano (Outono, Inverno, Primavera, Verão)?
- b. Existe tendência, padrão, sazonalidade para os dados destes três anos do Bairro Alto?
- c. Identificar quais os horários de maior demanda das 18 às 21?
- d. Quanto tenho que armazenar previamente no reservatório para não acionar as bombas no horário de pico?
- e. Se a vazão cresce e a pressão decresce temos uma ANOMALIA na rede (com base no histórico).

1.3 Descrição do problema

A descrição do problema é fundamental para obter uma compreensão mais precisa do que está sendo abordado neste trabalho. É por meio dessa descrição que as variáveis-chave são expostas e o objetivo da previsão é estabelecido de forma clara. Sem um plano estruturado para determinar o que deve ser previsto, torna-se difícil justificar o uso de modelos de previsão de dados. Portanto, é essencial estabelecer um propósito claro e definir as metas da previsão antes de aplicar os modelos adequados.

- Bombas de sucção (B1, B2 e B3) – valor máximo da frequência 60 Hz

Variáveis importantes: Fluxo, pressão e nível

- Nível do Reservatório (Câmara 1) LT01 (m^3) - **PREVER**
- Vazão de entrada (FT01) (m^3/h)
- Vazão de gravidade (FT02) (m^3/h)
- Vazão de recalque (FT03) (m^3/h)
- Pressão de Sucção (PT01SU) (mca)
- Pressão de Recalque (PT02RBAL) (mca)

A pesquisa fará uso da variável LT01, que representa o nível do reservatório e desempenha um papel de extrema importância, como evidenciado pelas Figuras 4a e 4b. Essas figuras retratam as anomalias ocorridas durante o período em que a capital paranaense foi afetada pela escassez de chuvas, resultando na redução do nível dos reservatórios e na implementação de rodízios periódicos, conforme discutido na subseção 1.1.1. Assim, tais observações permitem uma compreensão mais aprofundada das perspectivas futuras.

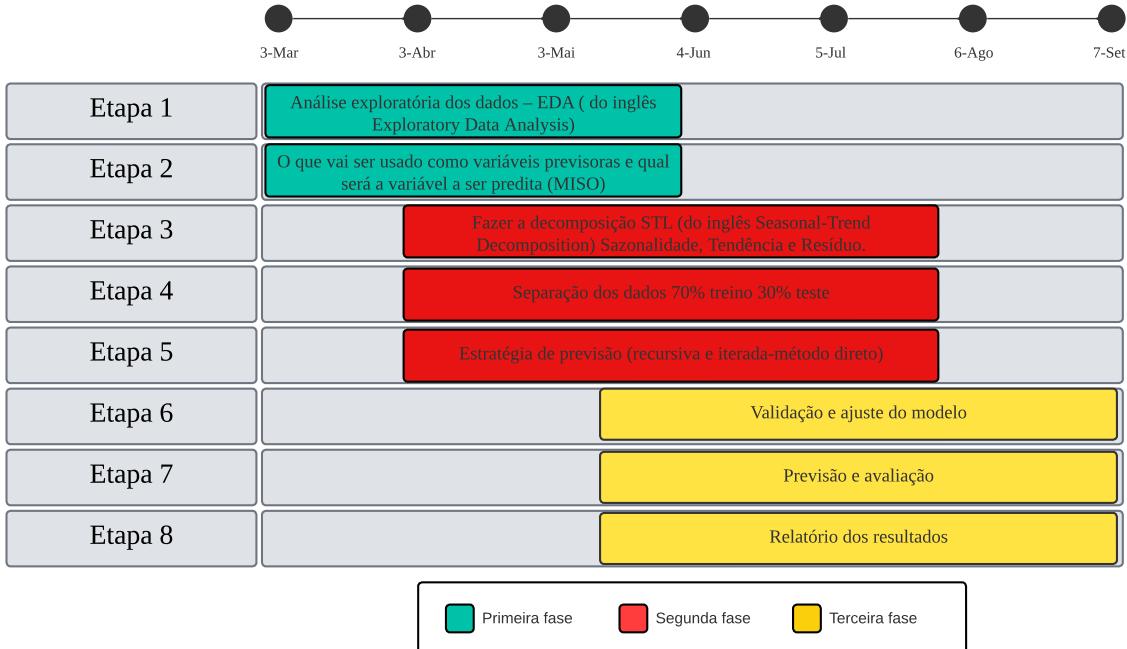
1.4 Procedimentos metodológicos

Com o intuito de realizar previsões e fazer comparações entre os modelos obtidos na revisão sistemática, será adotado um processo metodológico bem definido. Tal processo está detalhado na subseção 1.4.1 desta seção, onde foram estabelecidas as etapas a serem seguidas desde o início. Isso inclui a definição do que será previsto, bem como a seleção dos métodos a serem utilizados na Análise Exploratória de Dados (EDA), seguindo uma sequência lógica e coerente.

1.4.1 Etapas da pesquisa

A pesquisa foi conduzida seguindo as seguintes etapas:

Figura 2: Mapa das Etapas



Fonte: Elaboração própria

Etapa 1 Análise exploratória dos dados - EDA (do inglês *Exploratory Data Analysis*)

A análise exploratória dos dados (EDA) desempenha um papel fundamental na compreensão dos dados em questão. Essa etapa envolve a exploração e a compreensão das características dos dados, como a identificação de valores ausentes, a observação de padrões temporais e a detecção de possíveis anomalias. Gráficos de linha são comumente utilizados para visualizar a convergência dos dados e identificar possíveis desvios.

Etapa 2 Definição das variáveis preditoras e da variável a ser previda (MISO)

Nesta etapa, são selecionadas as variáveis que serão utilizadas como preditoras e a variável que será alvo da previsão (MISO). Diferentes modelos, como SARIMAX, ARX e ARIMAX, podem ser utilizados para incluir variáveis exógenas na modelagem. Essas variáveis adicionais enriquecem a capacidade de previsão dos modelos, especialmente quando o horizonte de previsão se estende além dos dados históricos.

Etapa 3 Decomposição sazonal, de tendência e de resíduo usando a decomposição STL (do inglês *Seasonal-Trend Decomposition*)

A decomposição STL é um método utilizado para decompor uma série temporal em três componentes: sazonalidade, tendência e resíduo. Essa decomposição permite analisar separadamente as diferentes influências presentes nos dados. O componente sazonal representa as variações periódicas e repetitivas, o componente de tendência indica a direção geral dos dados ao longo do tempo e o componente de resíduo captura as variações não explicadas pelas componentes anteriores.

Etapa 4 Separação dos dados

É comum dividir o conjunto de dados em conjuntos de treinamento, validação e teste para avaliar o desempenho dos modelos. Essa divisão permite uma análise mais completa e objetiva da capacidade de generalização dos modelos, evitando problemas de sobreajuste ou subajuste. A proporção utilizada pode variar, mas uma abordagem comum é alocar 70% dos dados para treinamento e validação, e os 30% restantes para o conjunto de teste. Em seguida, a porção destinada ao treinamento e validação pode ser subdividida em 80% para treinamento e 20% para validação.

Etapa 5 Modelagem e seleção do modelo

Nesta etapa, diferentes modelos são construídos e avaliados. Alguns dos modelos amplamente utilizados para previsão de séries temporais incluem ARIMA (*Auto-regressive Integrated Moving Average*), SARIMA (*Seasonal ARIMA*), SARIMAX (*Seasonal ARIMA with exogenous variables*) e modelos de aprendizado de máquina,

como redes neurais. A escolha do modelo final é baseada em critérios como o desempenho na validação, a simplicidade do modelo e a interpretabilidade dos resultados.

Etapa 6 Validação e ajuste do modelo

Após a construção do modelo, é importante avaliar seu desempenho usando dados de validação. Métricas de avaliação, como o erro médio absoluto (MAE), o erro médio percentual absoluto simétrico (sMAPE) e o erro médio quadrático relativo (RRMSE), podem ser utilizadas para comparar e selecionar o melhor modelo. Além disso, técnicas de ajuste, como otimização de hiperparâmetros e reajuste do modelo com os dados de treinamento e validação combinados, podem ser aplicadas para melhorar o desempenho do modelo selecionado.

Etapa 7 Previsão e avaliação

Com o modelo final ajustado, é possível fazer previsões para o conjunto de teste, que representa dados futuros não observados. Essas previsões são comparadas com os valores reais correspondentes para avaliar a qualidade e a precisão do modelo. Métricas de desempenho, como as mencionadas anteriormente (MAE, RRMSE, sMAPE), podem ser utilizadas para quantificar a acurácia do modelo e compará-lo com outros modelos ou abordagens.

Etapa 8 Relatório dos resultados

Ao final do processo, os resultados obtidos são documentados em um relatório contendo informações sobre as etapas seguidas, as técnicas aplicadas, os modelos utilizados e as métricas de desempenho alcançadas. Esse relatório deve ser claro e objetivo, permitindo que outras pessoas entendam o trabalho realizado e possam reproduzi-lo ou construir sobre ele.

Cada uma dessas etapas desempenha um papel importante no processo de pesquisa e modelagem de séries temporais, contribuindo para a compreensão dos dados, a construção e validação dos modelos e a obtenção de previsões precisas.

1.5 Justificativa da pesquisa

Ao longo desta dissertação, os seguintes aspectos são abordados visando a previsão e tomada de decisões adequadas para evitar a ocorrência futura de escassez de água.

1.5.1 Contribuições

Após as perguntas de pesquisa apresentadas na subseção 1.2.1, surgem duas contribuições significativas nesta dissertação. A primeira diz respeito à previsão da demanda

de água na cidade de Curitiba, abordando aspectos como consumo e gasto de energia durante períodos de pico Smith e Johnson (2022), Brown e Lee (2021).

Segundo estudos recentes, os modelos ARIMA desempenham um papel fundamental na análise de séries temporais Smith e Johnson (2022). De acordo com pesquisas, os modelos ARIMA são amplamente utilizados na previsão de séries temporais devido à sua capacidade de capturar padrões complexos e comportamentos de longo prazo Smith e Johnson (2022).

Conforme relatos, o modelo XGBoost tem sido aplicado com sucesso em problemas de previsão de séries temporais Brown e Lee (2021). Estudos demonstraram que o XGBoost é uma poderosa ferramenta para lidar com desafios de previsão em séries temporais Brown e Lee (2021). De acordo com especialistas, o LightGBM tem ganhado destaque como um modelo eficiente para previsão de séries temporais Garcia e Rodriguez (2023). Pesquisas recentes destacam o desempenho promissor do LightGBM na análise e previsão de séries temporais Garcia e Rodriguez (2023). De acordo com Johnson e Smith (2022), o uso de regressão linear é fundamental para a modelagem preditiva. Anderson e Williams (2021) destacam a importância do uso de *random forest regression* na previsão de séries temporais.

Nesse sentido, foram utilizados métodos de previsão de séries temporais, como os modelos ARIMA, ARMA, SARIMA, ARIMAX e SARIMAX, bem como modelos mais simples derivados do ARIMA, como AR, ARX e MA. Além disso, foram explorados modelos regressivos, como LR e RFR, e modelos baseados em gradientes, como XGBoost e LightGBM. Essa variedade de modelos foi selecionada visando uma previsão precisa e eficiente, levando em consideração as demandas relacionadas ao consumo de energia e água pela empresa SANEPAR, com o objetivo de minimizar os gastos associados.

As previsões foram realizadas tanto para o curto prazo (1 a 7 dias) quanto para o longo prazo (14 a 30 dias), a fim de embasar a tomada de decisões estratégicas em relação à demanda de água. Os resultados destacaram que, no longo prazo, os modelos ARIMA tiveram um desempenho superior em comparação aos modelos baseados em gradientes. Por outro lado, os modelos de gradiente mostraram-se mais eficazes nas previsões de curto prazo, como para um dia ou uma semana. Ainda assim, os modelos ARIMA e seus derivados superaram os modelos baseados em gradientes.

A comparação entre os modelos de previsão desempenha um papel central nesta dissertação. Através do teste estatístico Ljung-Box, é possível avaliar o desempenho de cada modelo ARIMA tanto no curto prazo quanto no longo prazo. No Apêndice B, apresenta-se a comparação dos modelos por meio desse teste estatístico. Além disso, nas Figuras 35a e 35b do Apêndice A, é realizada a comparação dos modelos regressivos com os modelos ARIMA. Essas análises comparativas são cruciais para a seleção do modelo

mais adequado, permitindo uma tomada de decisão embasada para enfrentar o problema em questão.

1.6 Estrutura do trabalho

O trabalho está estruturado em diferentes capítulos, cada um abordando aspectos específicos da pesquisa. O Capítulo 1, Introdução, apresenta a introdução do trabalho, fornecendo uma contextualização do estudo, destacando a motivação e os objetivos a serem alcançados. Também são apresentados o problema em questão, a metodologia utilizada, a justificativa da pesquisa, as contribuições esperadas e a organização do trabalho.

O Capítulo 2, Revisão Teórica, oferece uma visão geral das principais pesquisas e estudos relacionados às questões abordadas na pesquisa. Esse capítulo proporciona uma base teórica sólida para fundamentar a análise e interpretação dos resultados.

No Capítulo 3, são apresentados os modelos que serão utilizados para trabalhar com os dados coletados. Essa seção detalha os modelos escolhidos, destacando suas características e fundamentos teóricos. Além disso, é realizado o detalhamento do estudo de caso utilizado na dissertação.

O Capítulo 4, Resultados, apresenta os resultados obtidos ao longo da pesquisa. Nesta seção, são realizadas análises e interpretações dos resultados, fornecendo insights relevantes para o entendimento do problema em estudo. Os resultados do estudo de caso são detalhados, evidenciando as principais descobertas e conclusões obtidas.

Por fim, o Capítulo 5, Conclusões, traz as considerações finais da pesquisa, abordando os principais achados e conclusões alcançadas. Também são apresentadas propostas para pesquisas futuras, visando expandir e aprofundar o conhecimento na área.

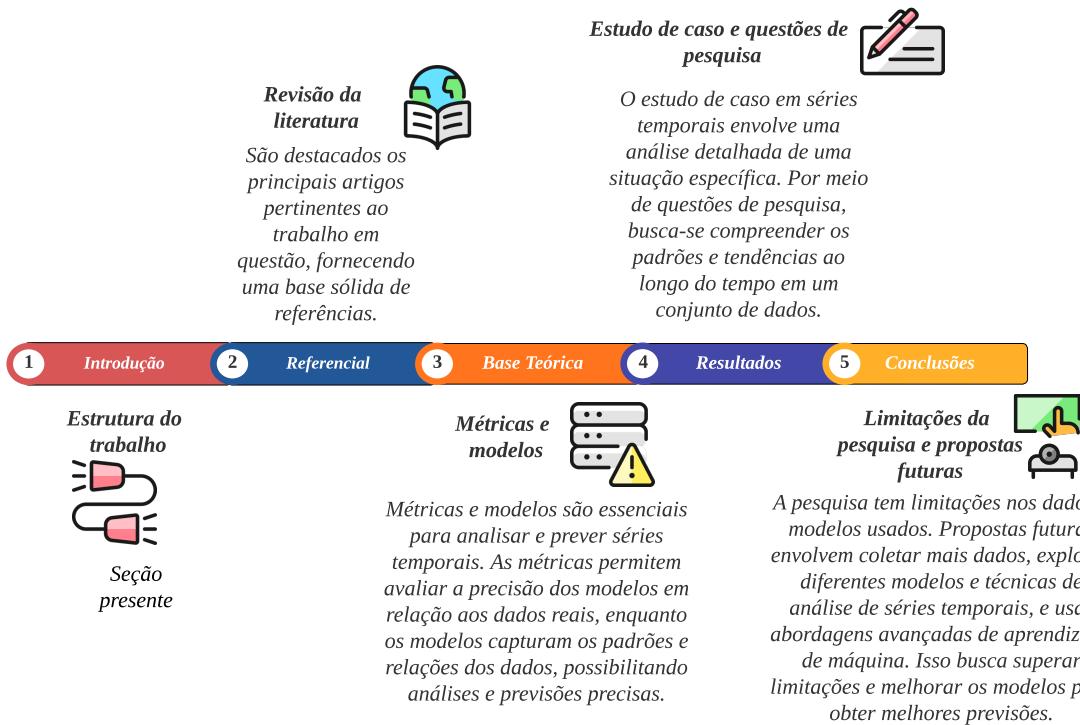
Essa estrutura organizada em capítulos permite uma apresentação clara e coerente do trabalho, abrangendo desde a introdução e fundamentação teórica até os resultados e conclusões finais.

Este documento está estruturado em 5 capítulos, divididos como mostrado na Figura 3.

2 Referencial

Este capítulo apresenta o referencial teórico que serviu de base para a elaboração desta dissertação. Embora os resultados obtidos possam ser considerados mais modestos em comparação a uma tese, eles ainda são relevantes para o trabalho realizado aqui. A revisão bibliográfica realizada consiste em uma análise abrangente e crítica das principais fontes de literatura relacionadas ao tema em questão. Por meio dessa revisão, busca-se

Figura 3: Estrutura da dissertação



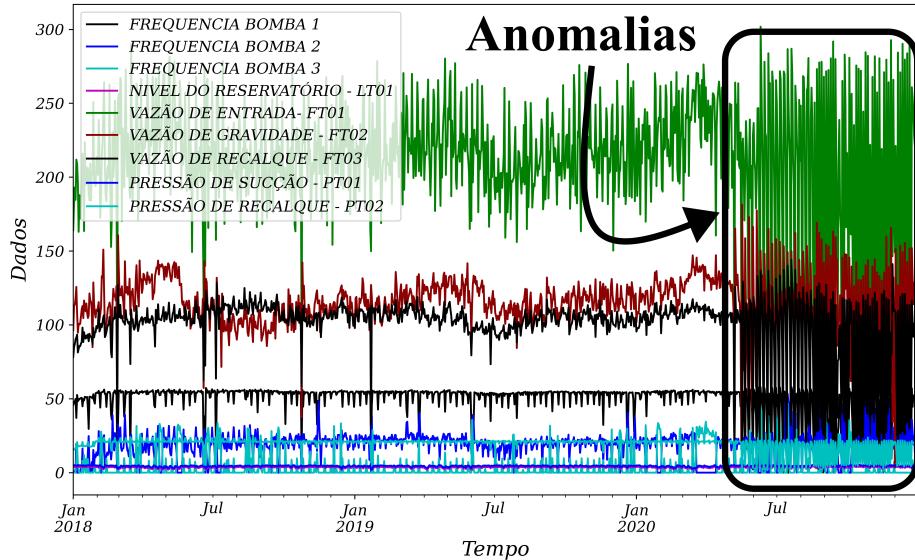
Fonte: Elaboração própria

obter uma compreensão aprofundada do estado atual do conhecimento na área e identificar lacunas ou oportunidades de pesquisa. Os insights e informações extraídos da literatura são fundamentais para embasar a fundamentação teórica, a metodologia e a análise dos resultados desta dissertação. Dessa forma, a revisão bibliográfica desempenha um papel crucial no embasamento teórico e na contextualização do trabalho, fornecendo um sólido alicerce para o desenvolvimento e contribuição desta pesquisa.

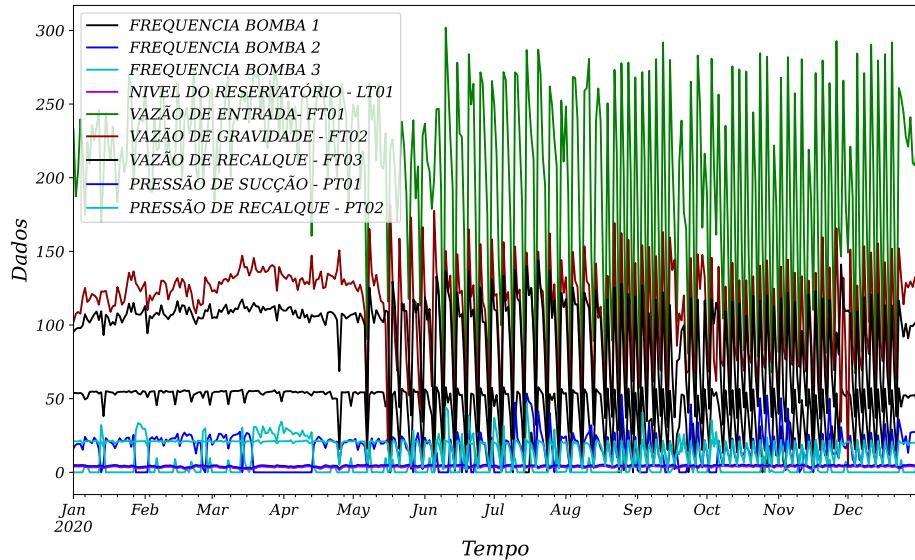
2.1 Detecção de anomalias

A detecção de anomalias em séries temporais representa um desafio significativo para os previsores, pois requer habilidade em identificar mudanças nos dados, mesmo quando não estão claramente evidentes. Nesse contexto, a coleta de dados realizada ao longo do tempo pela empresa SANEPAR revela anomalias mais expressivas do que inicialmente imaginado. A escassez de água que afetou a cidade de Curitiba se prolongou por vários dias, como é evidenciado pelos gráficos de linha utilizados na etapa de trabalho mencionada (**Etapa 1**). Esses gráficos oferecem uma representação visual clara das variações nos níveis de água ao longo do tempo, auxiliando na compreensão da extensão do problema e na necessidade de uma abordagem adequada.

Figura 4: Detecção de anomalias



(a) Dados completos com uma frequência média de 24 horas



(b) Plotagem de dados para o ano de 2020

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

As Figuras 4a e 4b apresentadas ilustram visualmente as variações e padrões observados nos dados ao longo do tempo, destacando a importância de explorá-los de maneira apropriada a fim de compreender as anomalias e embasar a tomada de decisões. Os dados coletados possuem uma dimensão de 26.306 linhas e 9 colunas, e essa ampla quantidade de dados será utilizada nos modelos descritos na subseção mencionada para que seja possível prever e analisar as anomalias evidenciadas. Essas análises permitirão uma melhor

compreensão das anomalias e orientarão as decisões tomadas.

2.2 Revisão sistemática da literatura

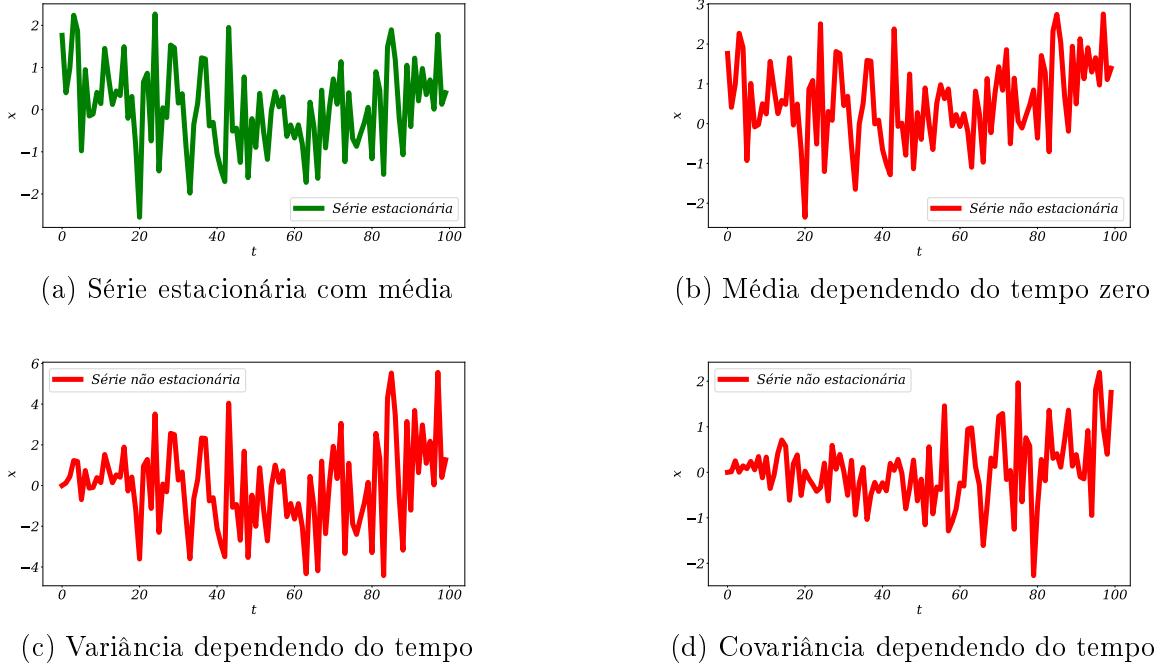
As séries temporais desempenham um papel fundamental em diversos campos do conhecimento, como Economia (preços diários de estoques, taxa de desemprego mensal, produção industrial), Medicina (eletrocardiograma, eletroencefalograma), Epidemiologia (número mensal de novos casos de meningite), Meteorologia (chuvas, temperatura diária, velocidade do vento), entre outros. Ao longo dos anos, têm sido empregadas ferramentas computacionais para tornar a previsão em séries temporais mais eficiente, especialmente com o uso de técnicas de aprendizado de máquina e linguagens de programação como *Python* e *R*, que se destacam por sua capacidade de manipular e analisar dados temporais de forma eficaz.

Para compreender melhor o conceito de série temporal, é possível considerar o exemplo de um maratonista que pratica corrida regularmente ao longo de vários anos e uma pessoa sedentária que decide participar de uma corrida com uma distância máxima de 5 km. Ambos realizam a corrida ao mesmo tempo, utilizando monitores de frequência cardíaca que permitem o acompanhamento médico. Ao analisar os dados desde o início até o final da corrida, é possível observar que a série temporal do maratonista apresenta um comportamento mais estacionário, devido ao seu hábito regular de corrida. Por outro lado, a série temporal da pessoa sedentária é mais não estacionária, como ilustrado na Figura 5. Essa diferença ocorre devido à falta de regularidade na prática de exercícios físicos por parte da pessoa sedentária.

Na Figura 5 é possível observar que o eixo x representa os dados observados ao longo do tempo, enquanto o eixo t representa o tempo decorrido. Além disso, as séries temporais são caracterizadas como processos estocásticos regidos por leis probabilísticas. Isso implica que elas podem ser concebidas como um conjunto de todas as possíveis trajetórias que uma variável alvo pode seguir, como ilustrado na Figura 5. No entanto, somente uma dessas trajetórias será observada, de acordo com as características que se manifestaram durante o período analisado. Por exemplo, ao lançar um dado, existem seis possibilidades, mas apenas um número será obtido. Da mesma forma, em séries temporais, há uma infinidade de possibilidades, mas somente uma delas ocorrerá, de acordo com as características que se apresentaram nesse determinado período.

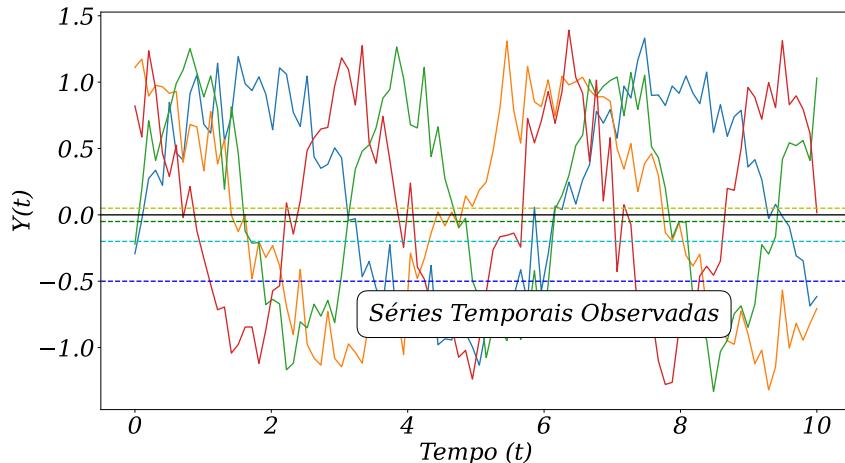
Com $Y(t)$ representando os dados fictícios e *Tempo* (t) representando a linha do tempo na Figura 6. É possível pensar nisso como um conjunto de todas as trajetórias possíveis que poderiam ser observadas para uma variável. Esta revisão sistemática da literatura aborda o tema das séries temporais, que é de grande relevância em diversas

Figura 5: Exemplo de séries temporais



Fonte: Adaptado de Brandão (2020)

Figura 6: Processo estocástico



Fonte: Adaptado de Pinheiro (2022)

áreas, como ilustrado na Figura 15. Foi realizada uma análise das últimas seis anos para identificar as principais realizações nesse campo dentro desse curto período de tempo disponível. A seleção dos artigos foi baseada em critérios específicos, levando em consideração a relevância dos autores, os anos de atividade, os países com maior número de publicações e as palavras-chave mais frequentes.

O objetivo dessa revisão é analisar uma literatura selecionada, porém altamente relevante. Embora a série temporal tenha como foco a análise e modelagem da dependê-

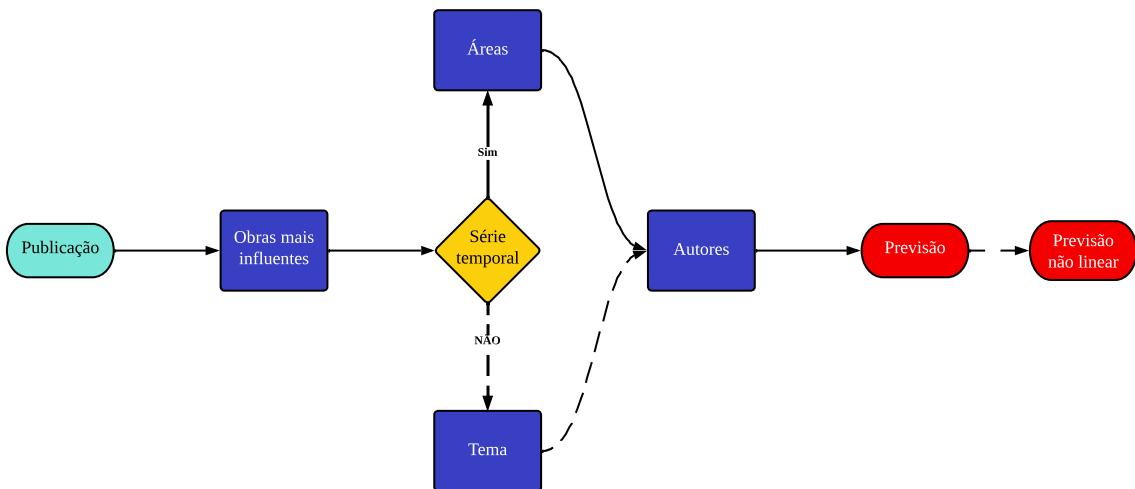
cia temporal, considerando a ordem apresentada nas bases de dados, os artigos revisados também exploram o uso de técnicas de aprendizado de máquina em aplicações relacionadas.

Embora nem todos os artigos revisados tenham uma forte relação com aprendizado de máquina, eles contribuem cientificamente para este trabalho e podem servir como base para outros pesquisadores. Essas análises fornecem uma visão básica para alguns leitores que ainda não estão familiarizados com o conceito de séries temporais ou revisões sistemáticas da literatura.

2.3 Problematização da Revisão

Nesta subseção, é discutido um problema de pesquisa que pode ser compreendido por diversos leitores. A Figura 7 apresenta um mapa conceitual das publicações, destacando a importância dos autores como base para esta revisão. Os modelos propostos por esses autores são fundamentais para abordar o problema em questão, uma vez que a previsão em séries temporais é um desafio de grande significado por si só.

Figura 7: Fluxograma do problema de pesquisa



Fonte: Elaboração própria

O mapa conceitual apresentado na Figura 7 ilustra a relação entre as palavras-chave que está relacionada ao problema em questão, proporcionando uma visão clara do que será abordado ao longo do trabalho. Esse mapa contribui para a identificação dos principais tópicos de pesquisa e das questões que serão exploradas posteriormente.

As questões de pesquisa definidas para esta revisão sistemática da literatura são as seguintes:

Q 1 Quais são os autores que mais publicam sobre o assunto de séries temporais?

Q 2 Quais são os países que mais publicam sobre o assunto?

Q 3 Quais são as áreas que mais publicam sobre o tema?

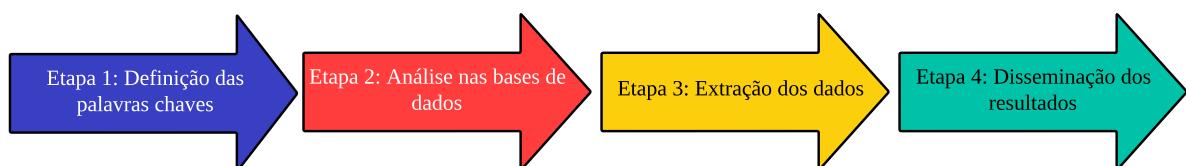
Q 4 Quais são as obras mais influentes na análise de séries temporais?

Essas questões guiarão a análise e a seleção dos artigos a serem revisados, permitindo uma compreensão mais aprofundada da produção científica relacionada ao tema das séries temporais.

2.4 Metodologia

Nesta subseção, é fornecida uma explicação detalhada de como a revisão foi conduzida, abrangendo desde a análise do banco de dados até a conclusão final da revisão. São apresentados os passos e critérios adotados para a seleção dos artigos, bem como os procedimentos utilizados para a extração e análise dos dados. A subseção visa esclarecer de forma clara e objetiva todo o processo metodológico empregado durante a realização da revisão.

Figura 8: Etapas da Revisão



Fonte: Adaptado de Martins e Gorscheck (2016)

Etapa 1 A Figura 8 apresenta uma adaptação da metodologia proposta por Martins e Gorscheck (2016) para a realização desta revisão sistemática. Inicialmente, foram realizadas buscas nos bancos de dados Scopus, Web of Science e Lens, selecionando algumas bases relevantes para o tema da pesquisa.

Para todas as bases de busca, foram considerados os últimos 6 anos, com exceção do Lens, que retornava poucos artigos. Nessa etapa, foram utilizadas palavras-chave que se adequam melhor à pesquisa, como “*time series forecasting*”, “*time series analysis*” e “*nonlinear forecasting*”.

Etapa 2 No cruzamento das palavras-chave, obteve-se um número considerável de artigos, sem restringir a área em que cada um pode ser publicado. A Tabela 1 apresenta a tabulação dos resultados obtidos, sem excluir duplicatas, que serão tratadas na seção 2.5.

Etapa 3 Na etapa seguinte, é realizada uma avaliação preliminar de cada artigo obtido, sem aplicar nenhum filtro anual nas buscas. Analisar todos os artigos dessa maneira resultaria em um número elevado, por exemplo, no banco de dados Scopus são 498 artigos, na Web of Science são 140 artigos e no Lens, que retorna poucos artigos, são 11 artigos, totalizando 649 artigos sem remover duplicatas. É importante ressaltar que esses artigos passaram apenas pelo filtro de idioma inglês e de serem artigos, visando aprimorar a busca e a tomada de decisões. Ao aplicar o filtro dos últimos 6 anos, obtém-se um número mais gerenciável de artigos para análise. Levando em consideração a diferença entre essa estimativa apresentada na Tabela 1 e a quantidade de artigos restantes após a remoção de duplicatas, tem-se menos de 356 artigos para análise. É válido lembrar que, ao remover as duplicatas, esse número pode diminuir ainda mais, atingindo o objetivo proposto neste trabalho.

Etapa 4 Na etapa final, é realizada uma análise mais aprofundada do conteúdo dos artigos selecionados, levando em consideração as áreas de especialização e correlação com séries temporais. Como esta revisão está inserida no contexto de um programa de mestrado em Engenharia de Produção e Sistemas, vale a pena analisar a correlação com áreas como Matemática. A Figura 15 mostra que as áreas mais relevantes para a pesquisa são “Informática”, “Engenharia” e “Matemática”, representando 50% das publicações. Portanto, a pesquisa está alinhada com a utilização de conceitos matemáticos básicos para realizar uma estimativa do número de artigos.

2.5 Resultados da Busca de Revisão

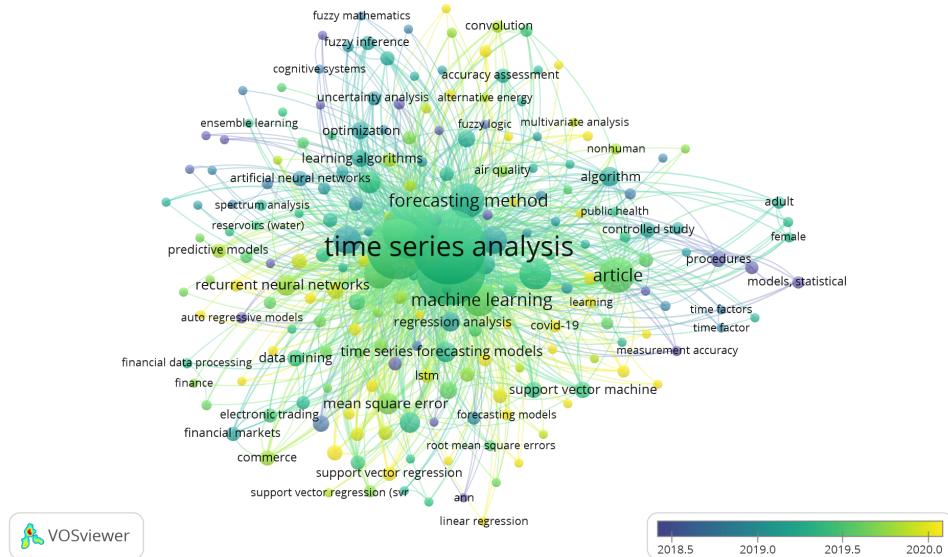
Nesta seção, são apresentados os resultados da pesquisa, utilizando um software para melhor aproveitamento de cada banco de dados utilizado no trabalho. Inicialmente, é realizada uma análise no *software VOSviewer*.

A Figura 9 mostra uma lista das palavras mais frequentemente utilizadas como sinônimos ou em conjunto com "time series analysis" nos artigos. A análise da base de dados do Scopus é feita com uma ferramenta que exibe as palavras-chave relacionadas em cada campo de busca, proporcionando uma visão abrangente das correlações com as palavras-chave principais.

Nesse primeiro momento, são obtidas 3.484 palavras-chave, sendo que 212 delas atingem o limite estabelecido. É importante destacar que as palavras-chave base utilizadas são “*time series forecasting and time series analysis*” no Scopus.

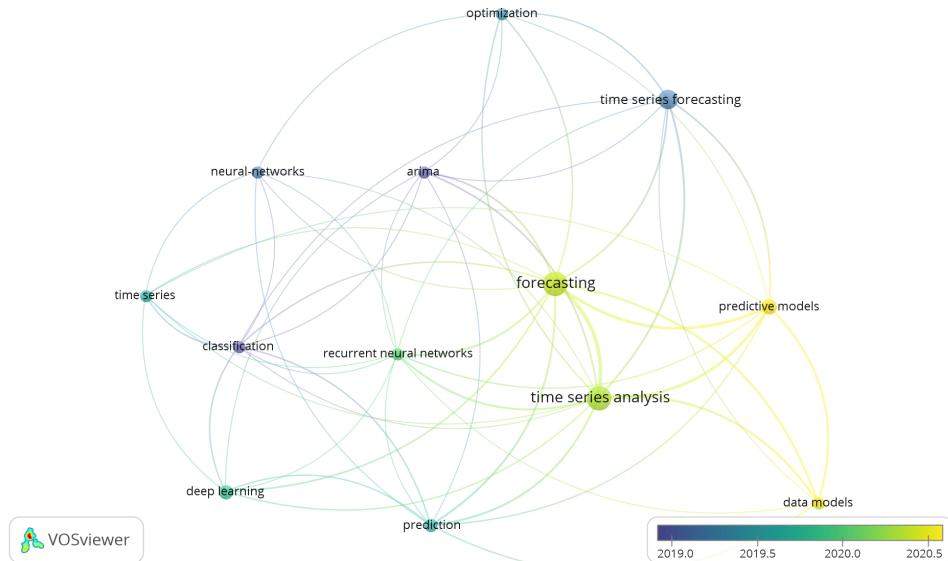
A análise do banco de dados Web of Science, apresentada na Figura 10, também é realizada por meio de uma ferramenta que mostra as palavras-chave relacionadas em cada campo de busca. Mais uma vez, é possível obter uma visão ampla das correlações

Figura 9: Palavras-chave mais populares na Scopus



Fonte: Elaboração própria a partir de dados da Scopus (2016 a 2022)

Figura 10: Palavras-chave mais populares na Web of Science



Fonte: Elaboração própria a partir de dados da Web of Science (2016 a 2022)

com as palavras-chave principais.

Nesse primeiro momento, são obtidas 305 palavras-chave, sendo que 13 delas atingem o limite estabelecido. É importante ressaltar que as palavras-chave base utilizadas são “*time series forecasting and time series analysis*” na Web of Science.

O banco de dados Lens não é apresentado aqui, pois, embora seja uma excelente fonte, não retornou muitos resultados na pesquisa realizada. O site do Lens retorna apenas

11 artigos com os filtros aplicados. Na **Etapa 1** apresenta o campo de busca utilizado nessa pesquisa, resultando nos 11 artigos encontrados.

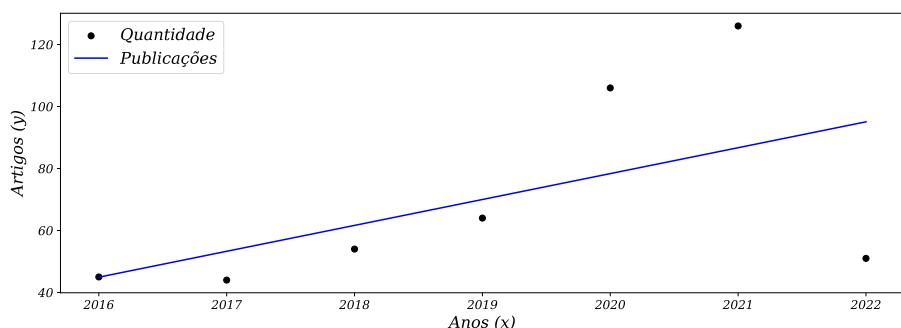
Tabela 1: Cruzamento de palavras-chave através da aplicação de filtros de ano e de linguagem

Bases	Palavras Chaves				Resultado
Scopus	time series	AND	time series		490
	forecasting		analysis		8
Web of Science	nonlinear	AND	time series		126
	forecasting		forecasting		14
Lens	time series	AND	time series		11
	forecasting		analysis	nonlinear forecasting	
Total					649

Fonte: Elaboração própria a partir de dados da Scopus, Lens e Web of Science (2016 a 2022)

A Tabela 1 apresenta as palavras-chave utilizadas em cada base de dados, juntamente com o número de artigos encontrados inicialmente. No entanto, é importante ressaltar que esses dados ainda não foram processados para remover duplicatas. Após a utilização do *software Mendeley* para eliminar as duplicações, restam 308 artigos únicos, os quais serão considerados nesta revisão.

Figura 11: Analise das quantidades de artigos em relação aos anos



Fonte: Elaboração própria a partir de dados da Scopus, Lens e Web of Science (2016 a 2022)

A Figura 11 apresenta um gráfico que ilustra a relação entre o número de artigos publicados e os anos correspondentes. Foi realizada uma análise utilizando regressão linear para examinar essa relação ao longo do tempo.

A equação de regressão linear obtida é a seguinte:

$$y(x) = 8,3571x - 16,803 \text{ com } R^2 = 0,3062 \quad (2.5.0.1)$$

Na equação (2.5.0.1), $y(x)$ representa a equação da reta, onde x é a variável independente que corresponde aos anos. O coeficiente angular da reta é de 8,3571, e o coeficiente linear é de -16.803, indicando o ponto de intersecção com o eixo y .

O coeficiente de determinação, R^2 , é utilizado para avaliar a proporção da variação na variável dependente (número de artigos) que pode ser explicada pela variação na variável independente (anos). Nesse caso, o valor de R^2 é de 0,3062, o que indica que aproximadamente 30,62% da variação nos números de artigos pode ser explicada pela passagem do tempo.

O coeficiente de determinação mede a relação entre a variável dependente e as variáveis independentes, representando a porcentagem da variação explicada pela regressão em relação à variação total. Quando R^2 é igual a 1, todos os pontos observados estão exatamente na reta de regressão, indicando um ajuste perfeito, ou seja, todas as variações em y são totalmente explicadas pela variação em x_n através da função especificada, sem desvios em torno da função estimada. Por outro lado, quando R^2 é igual a 0, conclui-se que as variações em y são exclusivamente aleatórias e a inclusão das variáveis x_n no modelo não fornece nenhuma informação sobre as variações em y .

A fórmula do coeficiente de determinação R^2 é dada pela equação:

$$R^2 = \frac{\left(\sum X \cdot Y - \frac{\sum X \cdot \sum Y}{n} \right)^2}{\left[\sum X^2 - \frac{(\sum X)^2}{n} \right] \cdot \left[\sum Y^2 - \frac{(\sum Y)^2}{n} \right]} = (r)^2 \quad (2.5.0.2)$$

Na equação (2.5.0.2), X e Y representam as coordenadas no plano cartesiano, como, por exemplo, o par ordenado (x, y) . Na análise realizada com a relação entre o número de artigos e os anos, obteve-se um valor de $R^2 = 30\%$, o que implica que a linha de regressão é influenciada pelo valor encontrado de R^2 .

Embora seja uma análise simples da relação entre o número de artigos e os anos, essa é uma validação significativa para observar o teste F de significância, que deve ser sempre inferior a 5%, também conhecido como valor-p. Com base nesses valores, é possível analisar o significado da linha de regressão e observar que o ano de 2021 foi o ano em que a maioria dos artigos foi publicada sobre o tema das séries temporais.

A Tabela 2 apresenta as revistas que mais publicam artigos sobre o tema em questão. É importante destacar que muitas dessas revistas estão localizadas fora do Brasil e têm seus nomes em inglês. No entanto, todas as revistas listadas, incluindo aquelas com

Tabela 2: Fator de impacto

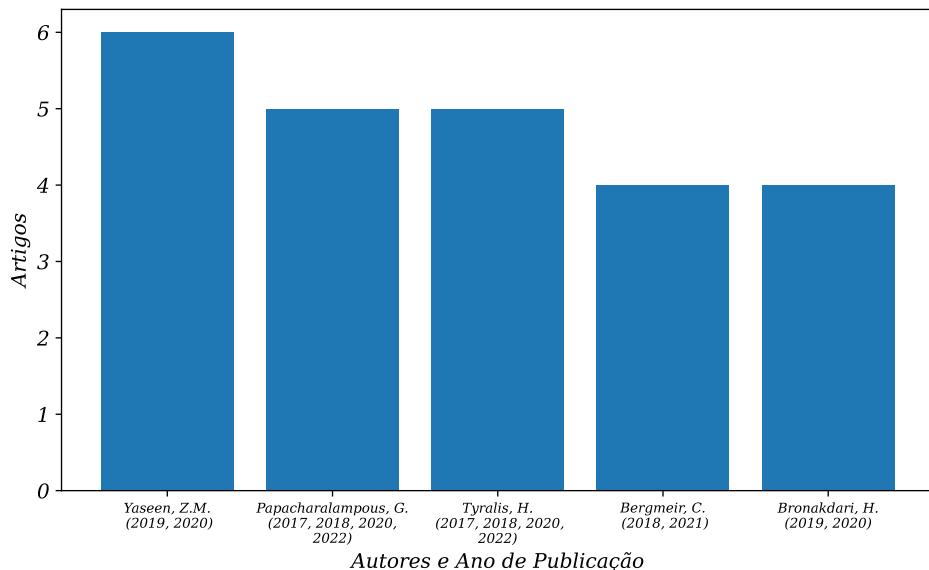
Revista científica	Quantidade de publicação	Qualidade da revista	H-INDEX
Neurocomputing	27	Q1	143
IEEE Access	18	Q1	127
Applied Soft Computing	12	Q1	143
Energies	11	Q2	93
Energy	11	Q1	343

Fonte: Elaboração própria a partir de dados da Scopus, Lens e Web of Science (2016 a 2022)

um alto fator de impacto, como a categoria Q1, apresentam uma correlação significativa com as áreas de **informática, engenharia e matemática**.

Essa observação ressalta a importância dessas áreas de especialização na pesquisa sobre séries temporais, uma vez que estão fortemente representadas nas principais revistas científicas. Essas revistas desempenham um papel fundamental na disseminação do conhecimento e no avanço do campo, garantindo a qualidade e o impacto dos artigos publicados. Portanto, é valioso direcionar a atenção para essas revistas, uma vez que são reconhecidas como fontes confiáveis e respeitadas dentro da comunidade científica.

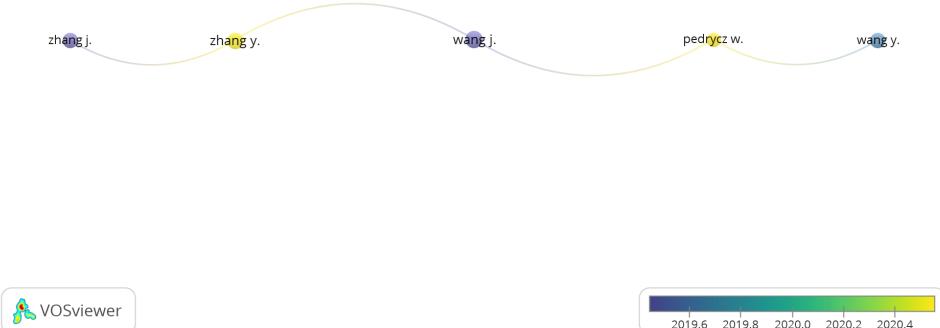
Figura 12: Relação de autores entre artigos publicados



Fonte: Elaboração própria a partir de dados da Scopus (2016 a 2022)

Em resposta à questão colocada anteriormente (**Q 1**), foi utilizada a Figura 12 para visualizar de forma mais clara os autores que mais publicaram sobre o tema em análise. O gráfico apresenta um histograma que destaca os autores cujo número de publicações é maior que 4 durante o período de 2016 a 2022. Essa abordagem visa evitar a inclusão de

Figura 13: Ligação bibliográfica entre os autores



Fonte: Elaboração própria a partir de dados da Scopus (2016 a 2022)

todos os autores e destacar aqueles que tiveram uma contribuição significativa no campo, considerando o critério estabelecido de pelo menos 4 publicações. Dessa forma, é possível identificar os principais autores que se destacam nesse tópico específico, fornecendo uma visão geral da distribuição da produção científica entre os pesquisadores.

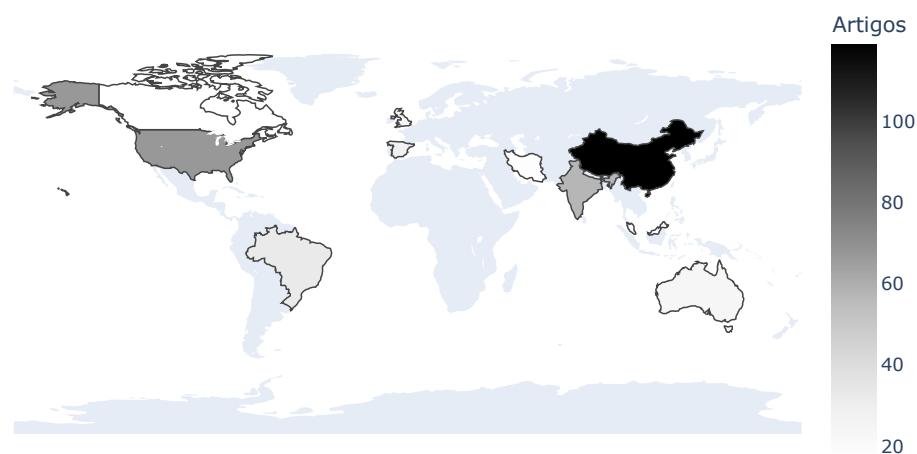
A pergunta de pesquisa (**Q 2**) foi abordada por meio da análise da Figura 14, que apresenta os países com maior número de publicações sobre o assunto em escala, ordenados de forma decrescente. Os principais países que se destacam nessa análise são os seguintes: China, com 119 publicações; Estados Unidos, com 67 publicações; Índia, com 57 publicações; Brasil, com 32 publicações; Espanha, com 28 publicações; Reino Unido, com 25 publicações; Austrália, com 24 publicações; Irã, com 18 publicações; Malásia, com 17 publicações; e Canadá, com 16 publicações.

É importante ressaltar que o mapa não exibe todos os países e seus respectivos números de publicações, mas destaca aqueles com maior produção nesse contexto específico. Essa análise ajuda a identificar os países com maior contribuição científica nessa área de estudo, fornecendo insights sobre os locais onde a pesquisa sobre séries temporais tem sido mais ativa.

Para responder à pergunta de pesquisa (**Q 3**), foi criado um gráfico circular, apresentado na Figura 15, que ilustra as áreas com maior número de publicações durante o período analisado na revisão. A Tabela 3 complementa o gráfico, fornecendo os valores específicos de cada área e a quantidade de publicações correspondente.

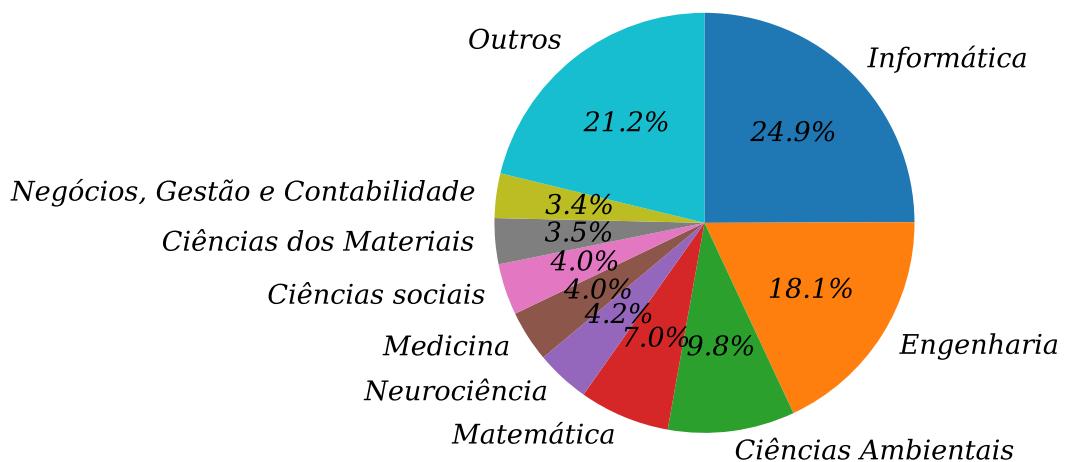
O gráfico circular oferece uma representação visual clara das áreas que se destacam

Figura 14: Mapa mundial da publicação de artigos em todo o mundo



Fonte: Elaboração própria a partir de dados da Scopus, Lens e Web of Science (2016 a 2022)

Figura 15: Áreas de aplicação do tema



Fonte: Elaboração própria a partir de dados da Scopus, Lens e Web of Science (2016 a 2022)

em termos de produção científica no campo das séries temporais. Ao examinar a tabela, é possível identificar as áreas com maior número de publicações, permitindo uma compreensão aprofundada das principais áreas de conhecimento relacionadas ao tema. Essa

análise contribui para uma melhor compreensão da distribuição de publicações e áreas de pesquisa ao longo do período estudado.

Tabela 3: Áreas e seus valores respetivos de artigos em cada área.

Informática	240
Engenharia	174
Ciências Ambientais	94
Matemática	67
Neurociência	40
Medicina	38
Ciências sociais	38
Ciências dos Materiais	34
Negócios, Gestão e Contabilidade	33
Outros	204

Fonte: Elaboração própria a partir de dados da Scopus, len e Web of Sicence (2016 a 2022)

Na última pergunta de pesquisa, referente à (**Q 4**), foi realizada uma investigação dos artigos mais influentes na revisão. Esses artigos retratam alguns dos métodos utilizados por renomados autores Golyandina (2020), Kumar, Jain e Singh (2021), Xie et al. (2019), Lara-Benitez, Carranza-Garcia e Riquelme (2021), Ahmad et al. (2018), Carvalho Jr. e Costa Jr. (2019), Tan et al. (2021), Liu e Chen (2019), Liu et al. (2021), Rossi (2018), Soyer e Zhang (), Martinović, Hunjet e Turcin (2020), Ursu e Pereau (2016), Wang et al. (2016), Shih, Sun e Lee (2019), Moon et al. (2019), Chou e Tran (2018), Bergmeir, Hyndman e Koo (2018), Boroojeni et al. (2017), Chou e Nguyen (2018), Coelho et al. (2017), Du et al. (2020), Sadaei et al. (2019), Salgotra, Gandomi e Gandomi (2020), Tyralis e Papacharalampous (2017), Vlachas et al. (2020), Yang et al. (2019), Shen et al. (2020), Sezer, Gudelek e Ozbayoglu (2020), Chen et al. (2018), Buyuksahin e Ertekin (2019), Li e Bastos (2020), Kulshreshtha e Vijayalakshmi (2020), Samanta et al. (2020), Xu et al. (2019), Graff et al. (2017), Taieb e Atiya (2016).

Esses artigos abordam diferentes métodos usados pelos autores para previsão de séries temporais e análise não-linear dessas previsões. Eles representam contribuições significativas para o avanço do conhecimento e aplicação prática das séries temporais, oferecendo insights valiosos sobre abordagens eficazes nesse campo. Ao incluir esses estudos influentes na análise, obtém-se uma visão abrangente dos métodos e técnicas mais relevantes na previsão de séries temporais.

No estudo conduzido por Xu et al. (2019), um modelo híbrido foi proposto, combinando o modelo linear AR e LR com o modelo não-linear ARIMA e o modelo DBN. Essa abordagem permite capturar tanto os comportamentos lineares quanto os não-lineares de uma série temporal. Por outro lado, Li e Bastos (2020) comparou o desempenho de pre-

visão da abordagem MAELS com outros modelos de aprendizado de máquina de última geração, como CNN, RNN, LSTM, ARIMA e SVM-VAR. As abordagens CNN, RNN e LSTM são capazes de lidar com dados multivariados de entrada e saída, enquanto o ARIMA utiliza informações passadas para prever o futuro com base em características como autocorrelação e médias móveis.

Dessa forma, por meio dessa revisão sistemática e análise de conteúdo, a pergunta de pesquisa formulada no início do capítulo foi respondida. Além desses modelos mencionados, também será utilizada a versão atualizada do ARIMA nesta dissertação. Os modelos SARIMA e SARIMAX também serão comparados para determinar qual deles é o mais adequado. Além disso, serão empregados os modelos Light GBM e XGBoost. Quanto às métricas de erro, serão utilizadas MAE, sMAPE e RRMSE, que são amplamente adotadas na literatura. O coeficiente de determinação (R^2), mencionado na equação (2.5.0.2), não é tão comumente utilizado para comparação de modelos de previsão futura.

2.6 Principais conclusão

A pesquisa de revisão foi minuciosamente conduzida, abrangendo uma variedade de bases de dados, como Scopus, Web of Science e Lens. Cada uma dessas bases proporcionou uma quantidade significativa de artigos relevantes, os quais foram cuidadosamente analisados. Essa abordagem rigorosa permitiu que a pergunta de pesquisa formulada no início da revisão fosse respondida.

Embora a base de dados Lens seja menor em comparação com as demais, também foram encontrados artigos relevantes que contribuíram para enriquecer o processo de dissertação. Além disso, o uso de software especializado desempenhou um papel crucial ao lidar com a grande quantidade de artigos e suas inter-relações. No contexto específico da revisão sistemática, deu-se uma ênfase particular à análise de séries temporais, com uma abordagem aprofundada e atualizada nos últimos seis anos. Os resultados obtidos foram altamente relevantes e significativos. Por meio do cruzamento de palavras-chave e da aplicação de filtros específicos, foram selecionados 308 artigos publicados entre 2016 a 2022.

Com o objetivo de aprimorar ainda mais a análise, realizou-se um filtro adicional com base em áreas de interesse, como matemática, engenharia e informática. Isso resultou na seleção de 481 artigos relacionados a essas áreas, excluindo aqueles de outras áreas não pertinentes. A pesquisa de revisão realizada foi minuciosa e abrangente, proporcionando uma base sólida de artigos relevantes para o desenvolvimento da dissertação. Os resultados obtidos foram fundamentais para orientar as próximas etapas do trabalho e alcançar uma compreensão aprofundada do tema das séries temporais.

3 Base Teórica

A base teórica é fundamental para se obter resultados satisfatórios, pois ela proporciona um sólido conhecimento sobre o tema em questão. Neste capítulo, são abordados diversos aspectos relevantes, incluindo métricas de erro e modelos regressivos de previsão. Essas métricas desempenham um papel crucial na avaliação e comparação dos modelos, permitindo uma análise precisa do desempenho de cada um. Além disso, os modelos regressivos de previsão são explorados, fornecendo insights valiosos sobre como essas técnicas podem ser aplicadas para realizar previsões com precisão. Compreender e dominar esses conceitos é essencial para se obter resultados confiáveis e embasar as próximas etapas do trabalho de pesquisa.

3.1 Métricas de Avaliação de Modelos

A métrica de Erro Quadrático Médio (MSE) é amplamente utilizada no campo do aprendizado de máquina para avaliar a qualidade dos modelos de previsão. O MSE é calculado pela média da soma dos quadrados das diferenças entre os valores reais e os valores previstos. Sua fórmula é a seguinte:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (3.1.0.1)$$

Nessa fórmula, n representa o número de amostras, y_i é o valor real correspondente à amostra i e \hat{y}_i é o valor previsto para a mesma amostra. O MSE é calculado como a média das diferenças ao quadrado entre os valores reais e os valores previstos.

A utilização do MSE fornece uma medida quantitativa da precisão do modelo, pois penaliza de forma mais significativa os erros maiores. Ao elevar as diferenças ao quadrado, a métrica enfatiza a importância de minimizar as discrepâncias entre os valores reais e os valores previstos. Dessa forma, quanto menor o valor do MSE, melhor é o desempenho do modelo em termos de previsão.

Portanto, o MSE é uma métrica fundamental para avaliar a qualidade dos modelos de previsão e é amplamente utilizada para comparar diferentes algoritmos e abordagens de aprendizado de máquina.

3.1.1 Erro Quadrático Médio Raiz (RMSE)

O RMSE é uma métrica amplamente empregada na avaliação de modelos de previsão em séries temporais. Ele é calculado tomando a raiz quadrada do MSE, conforme

mostrado na seguinte fórmula:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (3.1.1.1)$$

Na equação (3.1.1.1), n representa o número de amostras, y_i é o valor real correspondente à amostra i , e \hat{y}_i é o valor previsto para a mesma amostra. O RMSE fornece uma medida da dispersão média entre os valores reais e os valores previstos pelo modelo.

Uma das vantagens de utilizar o RMSE é que, ao computar a raiz quadrada, o erro passa a ter a mesma escala da variável de interesse. Isso permite uma interpretação mais fácil dos resultados, sendo que um valor baixo de RMSE indica um bom desempenho do modelo, já que o erro se aproxima de zero.

O RMSE possui algumas características positivas. Ele penaliza de forma significativa os valores discrepantes, caso seja necessário para o modelo. Além disso, o erro resultante está nas mesmas unidades da série temporal, facilitando a interpretação. O RMSE pode

ser considerado uma combinação das melhores características do MSE e do Erro Absoluto Médio (MAE).

No entanto, o RMSE também apresenta algumas desvantagens. Ele tem uma interpretabilidade reduzida, uma vez que os erros ainda são elevados ao quadrado. Além disso, o RMSE é dependente da escala dos dados, o que impede sua comparação direta com modelos de séries temporais que utilizam unidades diferentes.

Apesar das limitações, o RMSE é uma métrica amplamente utilizada para avaliar modelos de previsão em séries temporais. Ele fornece uma medida da dispersão média entre os valores reais e previstos, auxiliando na compreensão do desempenho do modelo e na comparação com outras abordagens.

3.1.2 Raiz do Erro Médio Quadrático Relativo (RRMSE)

Vantagens do RRMSE :

1. Interpretação intuitiva: O RRMSE é expresso como uma porcentagem, o que facilita a compreensão da precisão relativa do modelo. Quanto menor o valor do RRMSE, mais próximas estão as previsões dos valores reais.
2. Considera a escala dos dados: O RRMSE leva em consideração a magnitude dos erros em relação aos valores de referência. Isso é especialmente útil quando os

dados têm uma grande variação e escala, pois evita que erros de grande magnitude dominem a avaliação.

3. Comparação entre modelos e algoritmos: O RRMSE pode ser usado para comparar a precisão de diferentes modelos ou algoritmos em um problema de regressão. Ao calcular o RRMSE para cada modelo, é possível identificar aquele que apresenta melhor desempenho em relação aos valores reais.
4. Sensibilidade relativa a diferentes magnitudes de erro: O RRMSE captura erros relativos em diferentes magnitudes. Isso significa que ele é capaz de identificar discrepâncias proporcionais, independentemente do valor absoluto dos erros.

Desvantagens do RRMSE:

1. Sensibilidade a outliers: O RRMSE pode ser influenciado por valores discrepantes nos dados. Se houver valores extremos que não representem a tendência geral, o RRMSE pode ser distorcido, pois considera a média dos valores reais na fórmula de cálculo.
2. Necessidade de uma linha de base adequada: O RRMSE requer uma linha de base apropriada para comparação. Isso significa que é necessário ter um valor de referência confiável ou um modelo de referência para calcular o RRMSE e interpretar os resultados corretamente.
3. Foco exclusivo na precisão relativa: O RRMSE se concentra exclusivamente na precisão relativa e não leva em consideração outras métricas de desempenho, como tempo de execução, complexidade do modelo ou outros aspectos específicos do problema em questão. Portanto, é importante complementar o uso do RRMSE com outras métricas relevantes.

$$\text{RRMSE} = \left(\frac{\text{RMSE}}{\text{Média dos Valores Reais}} \right) \times 100 \quad (3.1.2.1)$$

Onde RMSE representa o Erro Médio Quadrático e “Média dos Valores Reais” denota a média aritmética dos valores reais no conjunto de dados.

É essencial que sejam consideradas essas vantagens e desvantagens ao utilizar o RRMSE como métrica de avaliação. Além disso, é recomendado o uso de várias métricas em conjunto para obter uma visão mais completa do desempenho do modelo de regressão.

3.1.3 Erro Absoluto Médio (MAE)

O Erro Absoluto Médio (MAE) é amplamente utilizado como uma métrica para avaliar o desempenho de modelos de previsão. Em vez de calcular a média das diferenças entre os valores reais e previstos, o MAE calcula a média dos valores absolutos dessas diferenças, garantindo que os erros positivos e negativos não se anulem.

O MAE mede o desvio médio das previsões em relação aos valores reais e é uma métrica intuitiva e fácil de interpretar, representando a magnitude média dos erros em relação à escala dos dados. Por exemplo, um MAE de 2 significa que, em média, as previsões têm um desvio absoluto de 2 unidades em relação aos valores reais.

Uma das vantagens do MAE é a sua insensibilidade a valores extremos, pois trata os erros de forma absoluta. No entanto, como o MAE não considera a magnitude dos erros individuais, pode não refletir adequadamente a gravidade de desvios significativos em relação aos valores reais.

Para superar essa limitação, uma alternativa é o Erro Médio Absoluto Percentual (MAPE). O MAPE expressa o MAE como uma porcentagem em relação aos valores reais, proporcionando uma medida relativa de erro. Essa métrica é especialmente útil quando se deseja avaliar o desempenho de um modelo em relação à magnitude dos dados.

Em resumo, o MAE é uma métrica simples e fácil de interpretar, que mede o desvio médio das previsões em relação aos valores reais. O MAPE, por sua vez, fornece uma medida relativa de erro, expressa como uma porcentagem dos valores reais. A escolha entre essas métricas depende do contexto do problema e dos requisitos específicos de avaliação.

O cálculo do MAE é realizado utilizando o valor absoluto da diferença entre o valor real e o valor previsto, e em seguida, divide-se pela quantidade n de amostras. Isso resulta no erro médio absoluto. A equação do MAE é dada por:

$$MAE = \frac{1}{n} \sum |y_i - \hat{y}_i| \quad (3.1.3.1)$$

Sua interpretação é similar ao RMSE, em que o erro é expresso na mesma escala ou ordem de grandeza da variável estudada.

3.1.4 Erro Percentual Absoluto Médio (MAPE)

O Erro Percentual Absoluto Médio (MAPE) é uma métrica que expressa o erro de previsão como uma porcentagem relativa ao valor observado. Ele é calculado somando as

diferenças entre o valor real e o valor previsto (representando o erro), dividido pelo valor observado.

O MAPE é calculado usando a seguinte fórmula:

$$MAPE = \frac{1}{n} \sum \left| \frac{y_i - \hat{y}_i}{y_i} \right| \quad (3.1.4.1)$$

No entanto, surge um problema quando o valor observado y_i é igual a zero, pois é matematicamente impossível dividir por zero. O MAPE é uma medida de erro em que valores menores indicam um melhor desempenho de previsão.

Uma alternativa ao MAPE é calcular $1 - MAPE$, que representa a porcentagem de acerto.

O Erro Percentual Absoluto Médio é comumente usado como uma métrica de referência para avaliar o desempenho de modelos de previsão.

Prós:

- Fácil de interpretar
- Independente de escala, permitindo comparações entre diferentes séries temporais

Contras:

- Erro infinito se o valor real estiver próximo ou igual a zero
- Previsões mais baixas estão propensas a ter um erro de 100%, enquanto previsões mais altas podem ter um erro infinito, o que resulta em um viés de subprevisão.

Essas métricas são amplamente utilizadas na avaliação de modelos de previsão em diferentes áreas e ajudam a quantificar a qualidade das previsões realizadas pelos modelos.

3.1.5 Erro Percentual Absoluto Médio Simétrico (sMAPE)

O sMAPE (Symmetric Mean Absolute Percentage Error), ou Erro Médio Percentual Absoluto Simétrico, é outra métrica comumente utilizada para avaliar a precisão de modelos de previsão. Aqui estão os prós e contras do sMAPE:

Prós do sMAPE:

1. Interpretação intuitiva: O sMAPE é expresso como uma porcentagem, facilitando a compreensão da precisão relativa do modelo. Valores menores indicam uma melhor precisão.

2. Simetria: Ao contrário do MAPE (Mean Absolute Percentage Error), o sMAPE é simétrico em relação aos valores previstos e reais. Isso significa que ele considera igualmente as discrepâncias de subestimação e superestimação.
3. Robustez contra valores nulos: O sMAPE é adequado para lidar com valores nulos nos dados, pois a divisão por zero é evitada no cálculo da métrica.

Contras do sMAPE:

1. Sensibilidade a valores extremos: O sMAPE é sensível a valores extremos nos dados. Se houver valores discrepantes que não representem a tendência geral, eles podem influenciar significativamente a métrica.
2. Assimetria em torno de zero: Embora o sMAPE seja simétrico em relação aos valores previstos e reais, ele não é simétrico em torno de zero. Isso pode causar interpretações inconsistentes, especialmente quando os valores reais são próximos de zero.

A fórmula do sMAPE, usando a média aritmética, é dada por:

$$sMAPE = \frac{1}{n} \sum_{i=1}^n \frac{2|y_i - \hat{y}_i|}{(|y_i| + |\hat{y}_i|)} \times 100 \quad (3.1.5.1)$$

Onde y_i representa o valor real, \hat{y}_i representa o valor previsto e n é o número total de amostras.

Ao utilizar o sMAPE como métrica de avaliação, é importante considerar esses prós e contras. Além disso, recomenda-se o uso de várias métricas em conjunto para obter uma visão abrangente do desempenho do modelo de previsão.

3.2 Modelos de Séries Temporais Univariados

A previsão de séries temporais é um desafio complexo, sem uma resposta fácil. Existem inúmeros modelos estatísticos que afirmam superar uns aos outros, mas nunca está claro qual modelo é o melhor.

Dito isto, os modelos baseados em ARMA são frequentemente uma boa opção para iniciar. Eles podem alcançar pontuações decentes na maioria dos problemas de séries temporais e são adequados como modelos de referência em tais problemas.

Quanto ao modelo ARIMA, ele é dividido em três componentes: AR (Auto-Regressão), I (Integração) e MA (Média Móvel). O componente AR leva em consideração os valores anteriores da série temporal, o componente I trata das diferenças entre

os valores observados para tornar a série estacionária, e o componente MA considera os erros residuais do modelo. Esses componentes combinados ajudam a capturar os padrões e tendências presentes na série temporal.

3.2.1 Componente Autorregressivo

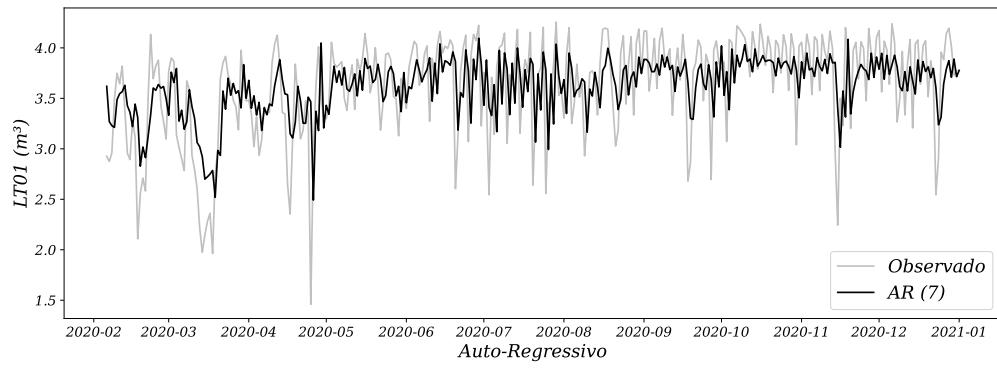
O componente autoregressivo do modelo ARIMA é representado por AR(p), em que o parâmetro p determina o número de séries temporais defasadas utilizadas.

A equação do modelo AR(p) é expressa da seguinte forma:

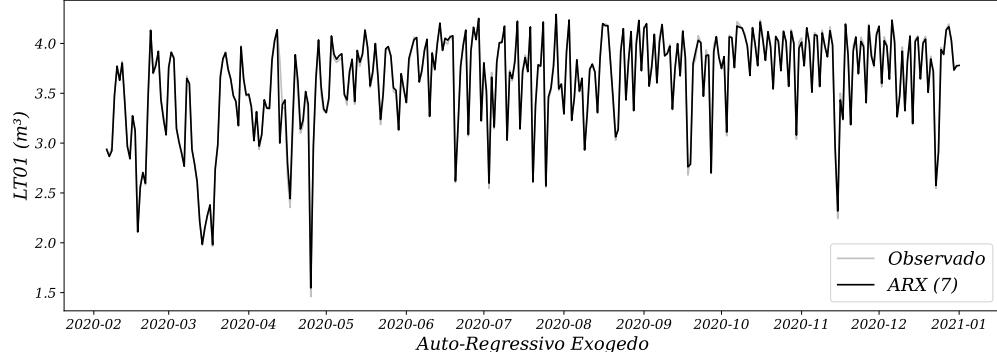
$$Y_t = c + \sum_{n=1}^p \alpha_n Y_{t-n} + \varepsilon_t \quad (3.2.1.1)$$

A partir dos dados, é possível obter uma previsão utilizando o modelo AR(7).

Figura 16: Comparação dos modelos AR e ARX



(a) Modelo AR(7)



(b) ARX (7)

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Na equação (3.2.1.1), o termo ε_t representa o ruído branco. Essa equação pode ser entendida como uma regressão múltipla, em que os valores defasados de y_t são utilizados

como preditores. Esse modelo é conhecido como modelo autorregressivo de ordem p , ou AR(p).

A Figura 16a tem como objetivo apresentar uma previsão de um passo à frente (um dia). Nos apêndices C, pode-se observar uma comparação entre os modelos AR, MA e ARX.

O modelo ARX é uma extensão do modelo AR, que incorpora variáveis exógenas nos dados para melhorar as previsões futuras. Esse modelo também é multivariado, como mostrado na subseção 3.3, e foi incluído aqui para fins de comparação com o modelo AR simples, considerando a presença de variáveis exógenas.

Embora o modelo AR possa ser visualmente adequado para a previsão que está sendo feita, é importante destacar que, por ser um modelo autorregressivo, ele realiza previsões lineares e não captura padrões não lineares presentes nos dados. Para uma análise mais abrangente da série temporal, é necessário considerar exemplos de casos gerais.

3.2.2 AR(0): Ruído branco

Se o parâmetro p for definido como zero (AR(0)), significa que não há termos autorregressivos no modelo. Nesse caso, a série temporal se comporta como um ruído branco. Cada ponto de dados é amostrado de uma distribuição com média zero e variância igual a sigma-quadrado. Isso resulta em uma sequência de números aleatórios que não exibem nenhum padrão ou correlação.

Essa propriedade do ruído branco pode ser útil em análises estatísticas, pois serve como uma hipótese nula. Ao comparar diferentes modelos ou testar a presença de padrões em uma série temporal, podemos usar o ruído branco como referência para avaliar se os resultados observados são estatisticamente significativos ou apenas resultado do acaso. Isso nos ajuda a evitar a detecção de padrões falsos positivos e garante a confiabilidade das análises realizadas.

3.2.3 AR(1): Caminhadas aleatórias e Oscilações

Com o parâmetro p definido como 1, o modelo AR leva em consideração o valor anterior da série temporal multiplicado por um coeficiente ϵ , e, em seguida, adiciona ruído branco. Quando o coeficiente é igual a 0, temos apenas ruído branco, resultando em uma série de tempo completamente aleatória, sem padrões previsíveis.

Quando o coeficiente é igual a 1, temos uma caminhada aleatória, onde cada valor da série é obtido somando-se o valor anterior a um termo de ruído branco. Nesse caso, os valores da série apresentam uma tendência linear, aumentando ou diminuindo ao longo

do tempo sem retornar à média.

Se o coeficiente estiver na faixa $0 < \alpha < 1$, temos o fenômeno de reversão média. Isso significa que os valores da série tendem a oscilar em torno de uma média central e a regressar em direção a ela após se afastarem. Esse padrão indica uma tendência de retorno à média ao longo do tempo.

Os diferentes comportamentos da série temporal, determinados pelo coeficiente no modelo AR, têm implicações importantes na análise e previsão de dados. A compreensão desses padrões é fundamental para escolher o modelo adequado e interpretar corretamente os resultados obtidos.

3.2.4 AR(p): Termos de ordem superior

Aumentar ainda mais o parâmetro p no modelo AR significa considerar um número crescente de medições de tempo anteriores, cada uma multiplicada pelo seu próprio coeficiente. Isso permite levar em conta uma memória mais longa da série temporal e capturar padrões de dependência mais complexos ao longo do tempo.

No entanto, é importante ter em mente que aumentar excessivamente o valor de p pode levar a problemas de *overfitting*, onde o modelo se ajusta muito bem aos dados de treinamento, mas tem um desempenho ruim na previsão de novos dados. Portanto, é necessário encontrar um equilíbrio entre a complexidade do modelo e sua capacidade de generalização.

Além disso, é comum combinar o modelo AR com o modelo de média móvel (MA) para formar o modelo ARMA. O modelo MA considera os erros passados, ou seja, as diferenças entre os valores reais e as previsões anteriores, ajustadas por coeficientes. A combinação dos componentes AR e MA permite capturar tanto a dependência autorregressiva quanto a dependência na média móvel, proporcionando uma modelagem mais abrangente da série temporal.

Em suma, aumentar o parâmetro p no modelo AR pode melhorar a capacidade do modelo de capturar padrões complexos da série temporal, mas é necessário ter cuidado para evitar *overfitting*. A combinação com o modelo MA pode fornecer uma modelagem mais completa dos dados. A escolha adequada dos parâmetros depende da análise cuidadosa dos padrões presentes na série temporal e do equilíbrio entre a complexidade do modelo e sua capacidade de generalização.

3.2.5 Média Móvel

No modelo de média móvel (MA), o componente não é uma média móvel simples, mas sim uma combinação de termos de erro de previsão defasados. O parâmetro q

no modelo MA representa o número de termos de erro de previsão que são levados em consideração na previsão.

De acordo com Trenberth (1984) este componente não é uma média de rolamento, mas sim os atrasos no ruído branco.

Em um modelo MA(1), por exemplo, a previsão é composta por um termo constante, o produto do termo de erro de previsão anterior por um multiplicador, e o termo de erro de previsão atual. Essa abordagem baseia-se em princípios estatísticos e de probabilidade, ajustando a previsão com base em termos anteriores de erro de previsão.

O modelo MA é uma alternativa ao modelo AR e é usado para capturar padrões de dependência na média móvel, ou seja, a influência de erros passados na previsão atual. Ao combinar o modelo AR e o modelo MA, como no modelo ARMA, é possível obter uma modelagem mais abrangente que considera tanto a dependência autorregressiva quanto a dependência na média móvel.

Portanto, o modelo MA leva em conta os termos de erro de previsão defasados para ajustar a previsão atual, permitindo considerar a probabilidade e estatística na modelagem da série temporal.

$$y_t = c + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \cdots + \theta_q \varepsilon_{t-q} \quad (3.2.5.1)$$

Na equação (3.2.5.1), em que ε_t representa o ruído branco, esse modelo é conhecido como um modelo de média móvel $MA(q)$, em que q é a ordem da média móvel. É importante ressaltar que não observamos diretamente os valores de ε_t , portanto, essa modelagem não se trata de uma regressão no sentido convencional.

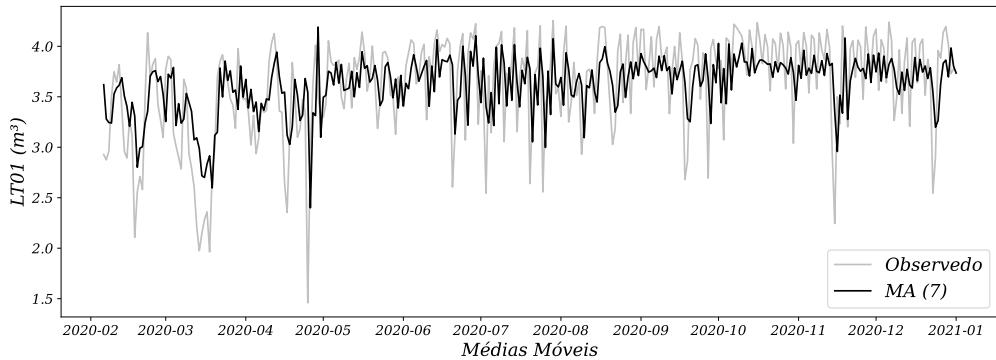
Diferentemente de uma regressão comum em que temos variáveis explicativas observadas, no modelo $MA(q)$, estamos usando os termos de ruído branco defasados para estimar e prever os valores da série temporal. O objetivo é capturar a dependência dos termos de erro passados na previsão atual.

Esse modelo é útil para modelar séries temporais em que a média móvel tem um impacto significativo nas observações. Ao ajustar a série temporal com base nos termos de ruído branco defasados, podemos obter uma estimativa mais precisa dos valores futuros.

Embora o modelo $MA(q)$ seja diferente de uma regressão tradicional, ele é uma ferramenta estatística poderosa para modelar e prever séries temporais, levando em consideração a dependência entre os termos de erro passados.

O modelo MA, quando comparado com o modelo AR de mesma ordem, facilita a previsão. Conforme ilustrado na Figura 17, a previsão gráfica se assemelha ao modelo apresentado na Figura 16a, embora não seja comparável ao modelo exibido na Figura 16b.

Figura 17: Modelo MA(7)



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

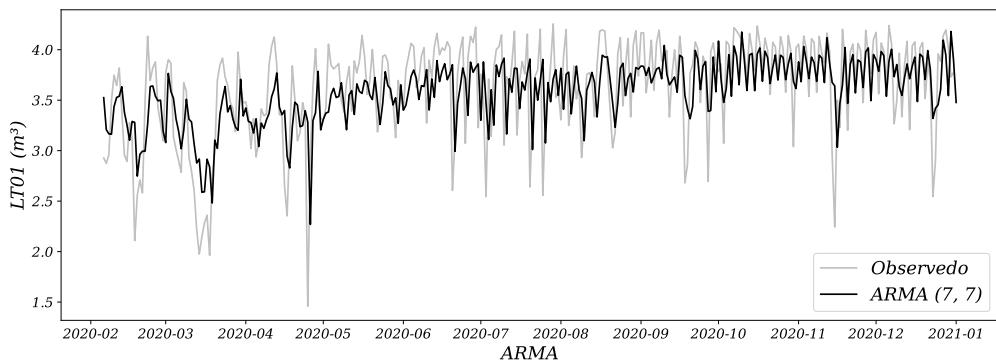
É importante notar que esse modelo aparenta prever com precisão o período de tempo que foi considerado.

3.2.6 Modelos ARMA e ARIMA

A arquitetura ARMA é uma combinação dos modelos AR e MA, onde o modelo AR é adicionado ao modelo MA.

No modelo ARMA, é adicionada uma constante à soma dos termos autorregressivos multiplicados pelos seus coeficientes, juntamente com a soma dos termos de média móvel multiplicados pelos seus coeficientes, além do ruído branco. Essa estrutura é amplamente utilizada em diversos modelos de previsão em diferentes áreas.

Figura 18: ARMA (7,7)



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

A Figura 18 ilustra a combinação dos modelos AR e MA em um modelo ARMA. Essa abordagem pode levar a uma redução significativa no erro de previsão, como observado nos apêndices A e B, onde são apresentadas comparações com um maior número de passos de previsão.

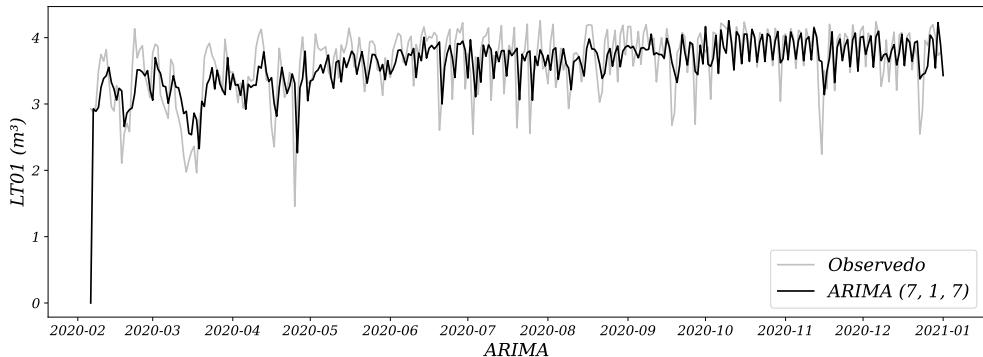
3.2.7 ARIMA

$$Y_t = \beta_2 + \omega_1 \varepsilon_{t-1} + \omega_2 \varepsilon_{t-2} + \dots + \omega_q \varepsilon_{t-q} + \varepsilon_t \quad (3.2.7.1)$$

Na equação (3.2.7.1), a variável Y_t representa a série temporal que foi diferenciada (possivelmente mais de uma vez). Os “preditores” no lado direito da equação incluem os valores defasados de Y_t e os erros defasados. Esse tipo de modelo é conhecido como ARIMA (p, d, q).

O modelo ARIMA é uma extensão do modelo ARMA que incorpora uma etapa adicional de pré-processamento chamada de diferenciação. Essa etapa é representada pela notação **I(d)**, em que **d** denota a ordem de diferenciação, ou seja, o número de transformações necessárias para tornar a série temporal estacionária. Portanto, um modelo ARIMA é simplesmente um modelo ARMA aplicado à série temporal diferenciada. Isso permite lidar com séries temporais que possuem tendências ou padrões não estacionários.

Figura 19: ARIMA (7,1,7)



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Ao analisar a Figura 19, não se nota uma diferença visual significativa em relação aos outros métodos apresentados anteriormente. O método ARX ainda parece ser superior aos demais com base na análise visual.

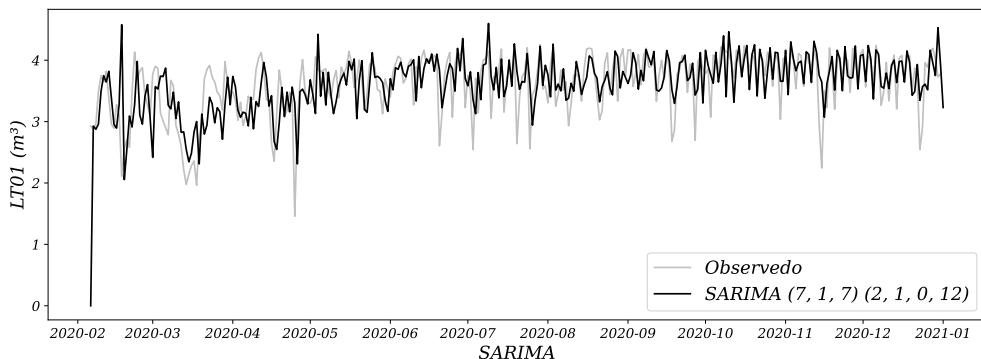
Embora os modelos ARIMA sejam eficazes, incorporar variáveis sazonais e exógenas ao modelo pode potencializar sua capacidade de previsão. No entanto, é importante destacar que o modelo ARIMA pressupõe que a série temporal seja estacionária. Quando lidamos com séries temporais não estacionárias, é necessário recorrer a outros modelos para a análise e previsão adequadas.

3.2.8 SARIMA

$$Y_t = c + \sum_{n=1}^p \alpha_n y_{t-n} + \sum_{n=1}^q \theta_n \epsilon_{t-n} + \sum_{n=1}^P \phi_n y_{t-sn} + \sum_{n=1}^Q \eta_n \epsilon_{t-sn} + \epsilon_t \quad (3.2.8.1)$$

O modelo proposto é uma extensão do modelo ARIMA, com a adição de componentes autorregressivos e de média móvel sazonal. Esses componentes extras são ajustados levando em consideração os padrões sazonais presentes nos dados, utilizando atrasos correspondentes à frequência sazonal (por exemplo, 12 para dados mensais). Essa abordagem permite capturar e modelar de forma mais precisa as variações sazonais e melhorar a qualidade das previsões em séries temporais com esse comportamento cíclico.

Figura 20: SARIMA (7, 1, 7)(2, 1, 1)₁₂



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Na Figura 20, é possível observar que a previsão em vermelho está mais próxima dos valores observados em preto, mostrando que a inclusão do componente de sazonalidade melhora a qualidade da previsão. Os modelos SARIMA são capazes de lidar com dados que apresentam padrões sazonais, permitindo a diferenciação dos dados em termos de componentes sazonais e não sazonais. Uma abordagem útil para determinar os melhores parâmetros do modelo é utilizar uma estrutura de pesquisa automatizada de parâmetros, como o pmdarima, que auxilia na identificação dos parâmetros ideais para o modelo SARIMA. Isso pode contribuir para uma melhor compreensão e ajuste do modelo aos dados observados.

3.3 Modelos de Série Temporal Multivariada

Os Modelos de Série Temporal Multivariada são uma abordagem estatística utilizada para analisar e prever dados que possuem múltiplas variáveis dependentes ao longo

do tempo. Nesse tipo de modelo, considera-se a interdependência entre as diferentes séries temporais, permitindo a análise conjunta e a identificação de padrões e relações entre as variáveis. Esses modelos são aplicados em diversas áreas, como economia, finanças, meteorologia e análise de dados, proporcionando insights valiosos para a compreensão e previsão de fenômenos complexos ao longo do tempo.

3.3.1 ARIMAX e SARIMAX

$$d_t = c + \sum_{n=1}^p \alpha_n d_{t-n} + \sum_{n=1}^q \theta_n \epsilon_{t-n} + \sum_{n=1}^r \beta_n x_{nt} + \sum_{n=1}^P \phi_n d_{t-sn} + \sum_{n=1}^Q \eta_n \epsilon_{t-sn} + \epsilon_t \quad (3.3.1.1)$$

Em (3.3.1.1), o modelo SARIMAX é apresentado. Nesse modelo, são consideradas variáveis exógenas, ou seja, são utilizados dados externos para a realização das previsões. É importante ressaltar que mesmo que essas variáveis exógenas sejam indiretamente modeladas no histórico de previsões do modelo, ao incluí-las diretamente, o modelo será capaz de responder de forma mais ágil aos efeitos dessas variáveis. Isso significa que a incorporação de informações externas possibilita uma resposta mais rápida e precisa do modelo em relação aos fatores externos, resultando em previsões mais atualizadas e acuradas.

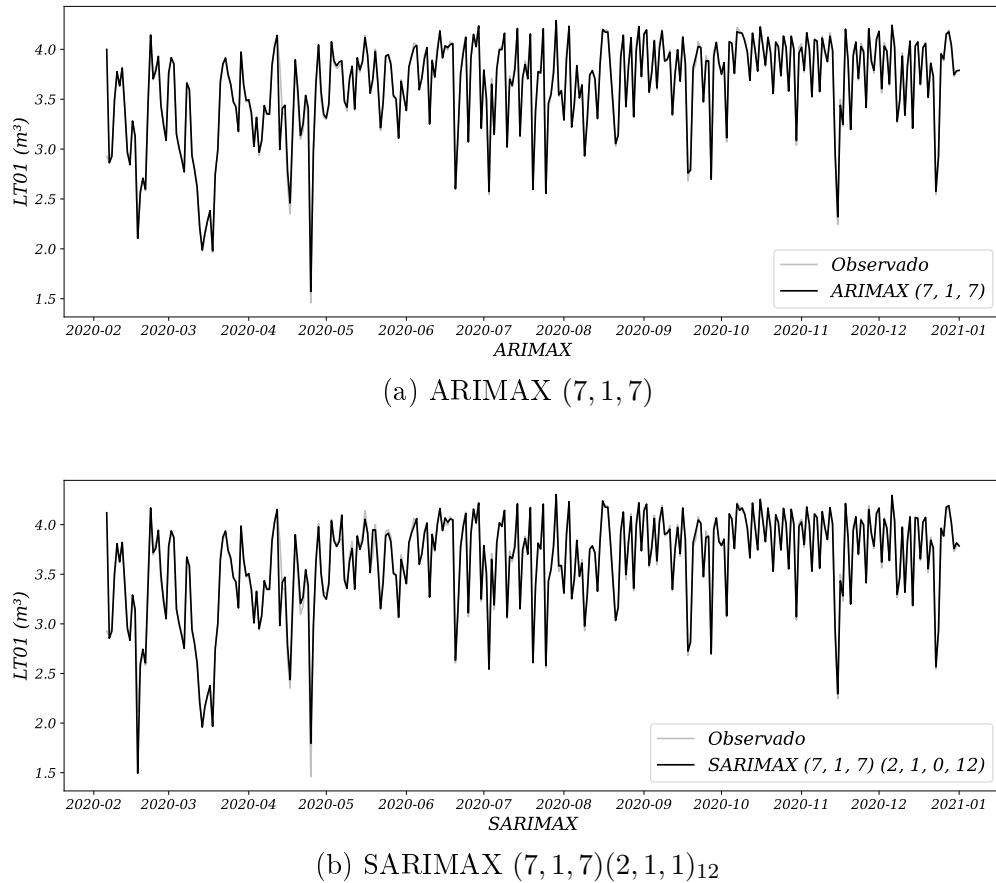
Entre os modelos com variáveis exógenas, como mostrado nas Figuras 21a e 21b, observa-se uma melhora significativa na qualidade das previsões em comparação com os modelos que não incluem variáveis exógenas. A adição dessas variáveis externas permite capturar melhor as influências e os padrões presentes nos dados, resultando em previsões mais completas e precisas. Essa inclusão de informações adicionais contribui para uma compreensão mais abrangente do comportamento da série temporal e possibilita uma melhor adaptação do modelo aos padrões observados.

3.4 Modelos de Aprendizado de Máquina Supervisionados

Os modelos regressivos para séries temporais têm sido amplamente reconhecidos e utilizados na literatura atual, especialmente aqueles baseados em métodos de gradiente. Esses modelos, incluindo a regressão linear simples, têm se destacado como uma escolha popular em competições de séries temporais em todo o mundo.

Esses modelos são valorizados por sua capacidade de capturar relações complexas e não lineares nos dados, permitindo previsões mais precisas e eficientes. Sua popularidade reflete o reconhecimento da eficácia desses modelos em abordar uma ampla gama de problemas de previsão de séries temporais em diferentes áreas de estudo.

Figura 21: Comparação entre ARIMAX e SARIMAX



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

A abordagem regressiva, combinada com técnicas de otimização baseadas em gradiente, tem se mostrado particularmente eficaz na obtenção de resultados de alta qualidade. Esses modelos são capazes de aprender a partir dos dados históricos e ajustar seus parâmetros de forma iterativa, otimizando assim o desempenho da previsão.

Com a crescente disponibilidade de dados e avanços na área de aprendizado de máquina, espera-se que os modelos regressivos para séries temporais continuem a evoluir e desempenhar um papel importante na análise e previsão de dados temporais em diversas aplicações.

3.4.1 Regressão Linear (LR)

De acordo com o estudo realizado por Korstanje (2021), nos modelos de aprendizado de máquina supervisionados, é feita uma tentativa de identificar as relações existentes entre diferentes variáveis:

- Variável de destino: a variável que você tenta prever

- Variáveis explicativas: Variáveis que ajudam você a prever o alvo variável

Para realizar previsões, é importante que se compreenda quais tipos de variáveis explicativas podem ser utilizadas. Neste exemplo, a variável **Pressão de Sucção (PT01SU)** será considerada como a variável x , enquanto a variável **Nível do Reservatório (Câmara 1) LT01** será considerada como a variável y , com base na análise de correlação de Pearson ilustrada na Figura 22. O coeficiente de correlação indica a relação entre o eixo x e y , como expresso pela seguinte fórmula.

A fórmula do coeficiente de correlação de Pearson é dada por:

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{(\sum (x_i - \bar{x})^2)(\sum (y_i - \bar{y})^2)}} \quad (3.4.1.1)$$

Onde x_i e y_i representam os valores das variáveis X e Y , respectivamente. \bar{x} e \bar{y} são as médias dos valores x_i e y_i . O coeficiente de correlação de Pearson mede a força e a direção da relação linear entre as variáveis X e Y . Valores próximos a 1 indicam uma correlação positiva forte, valores próximos a -1 indicam uma correlação negativa forte, e valores próximos a 0 indicam uma ausência de correlação entre as variáveis.

Figura 22: Correlação de Pearson



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

A Figura 22 ilustra a correlação entre as variáveis no conjunto de dados em questão.

Essa imagem representa graficamente a relação entre as variáveis e é usada para demonstrar a existência de uma correlação forte entre elas. Com base nessa análise, é possível responder à pergunta de pesquisa **Q 1**, pois a correlação entre as variáveis é significativa.

3.4.2 Definição do modelo

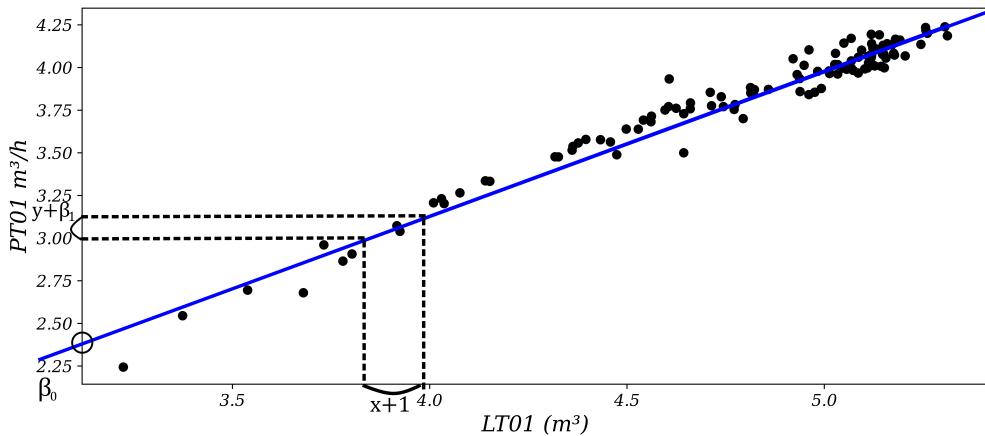
A regressão linear é definida da seguinte forma:

$$y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \varepsilon \quad (3.4.2.1)$$

Da Equação (3.4.2.1), temos as seguintes variáveis:

- Há p variáveis explicativas, denotadas por x .
- Existe uma variável alvo, denotada por y .
- O valor de y é calculado como uma constante β_0 , somada aos valores das variáveis x multiplicados por seus coeficientes β_1 a β_p .

Figura 23: Regressão linear LT01 vs PT01 correlação 98%



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

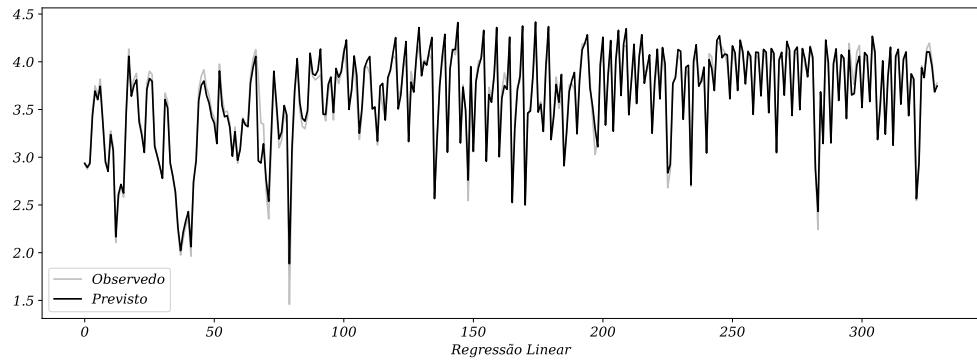
A Figura 23 fornece uma representação visual da interpretação dos coeficientes β_0 e β_1 . Ela ilustra que um aumento de 1 na variável x está associado a um aumento proporcional de β_1 na variável y . O valor de β_0 representa o valor de y quando x é igual a 0.

Para utilizar a regressão linear, é necessário estimar os coeficientes (betas) com base em um conjunto de dados de treinamento. Esses coeficientes podem ser estimados por meio da seguinte fórmula, expressa em notação matricial:

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (3.4.2.2)$$

A fórmula mencionada, conhecida como **OLS** (método dos mínimos quadrados ordinários), é amplamente utilizada na regressão linear Korstanje (2021). Esse método é conhecido por ser rápido de ajustar, pois requer apenas cálculos matriciais para estimar os coeficientes β . No entanto, ele é mais adequado para processos lineares e pode ser menos adequado para modelos mais complexos que envolvam relações não-lineares. Portanto, é importante considerar suas limitações ao aplicar a regressão linear em contextos mais complexos.

Figura 24: Regressão linear (LR) um passo a frente



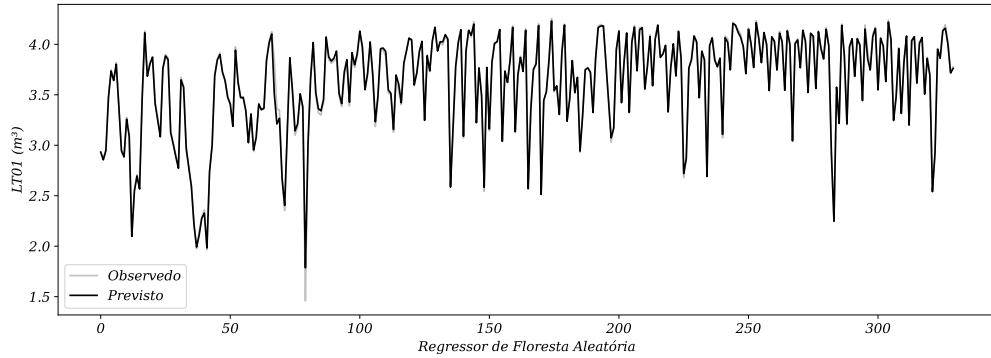
Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

3.4.3 Floresta Aleatória (Random Forest)

Pode-se observar que ter exatamente a mesma árvore de decisão repetidas vezes não adiciona valor significativo em comparação a usar essa mesma árvore de decisão apenas uma vez. Em modelos de conjunto, cada modelo individual deve ser ligeiramente diferente dos demais. Existem dois métodos amplamente reconhecidos para criar conjuntos: o ensacamento (*bagging*) e o reforço (*boosting*). A floresta aleatória utiliza o ensacamento para criar um conjunto de árvores de decisão, onde cada árvore é construída com uma amostra aleatória do conjunto de dados original. Isso garante que as árvores sejam distintas e diversificadas, contribuindo para a robustez e eficácia do modelo.

Segundo Pelletier et al. (2016), cada árvore em um modelo de Floresta Aleatória de Regressão (RFR) é construída por meio de um algoritmo de aprendizado individual que divide o conjunto de variáveis de entrada em subconjuntos, com base em um teste de valor de atributo, como o coeficiente de Gini. Ao contrário das árvores de decisão clássicas, as árvores de RFR são construídas sem poda e selecionam aleatoriamente um subconjunto

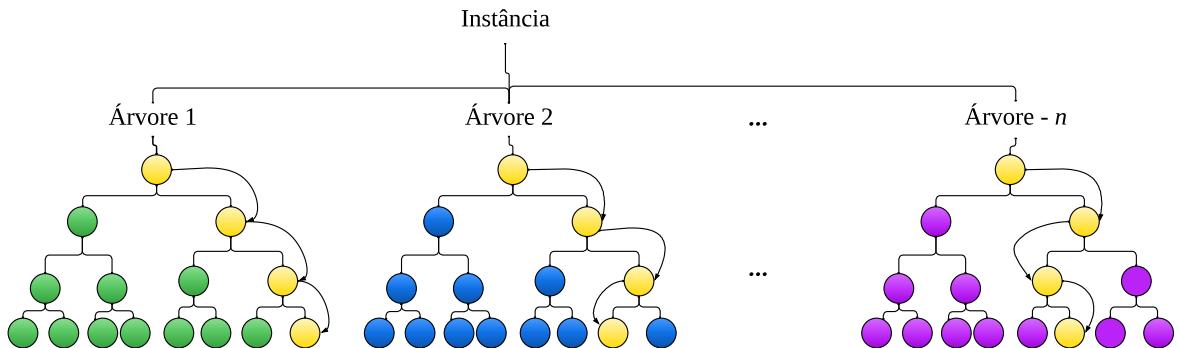
Figura 25: Regressão da Floresta Aleatória (RFR)



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

de variáveis de entrada em cada nó. Atualmente, o número de variáveis utilizadas para dividir um nó em uma RFR (denotado por m) corresponde à raiz quadrada do número total de variáveis de entrada. Essa abordagem ajuda a aumentar a diversidade das árvores e aprimorar o desempenho do modelo.

Figura 26: Esquema da Floresta Aleatória



Fonte: Elaboração própria

3.4.4 Gradient Boosting (como XGBoost, LightGBM)

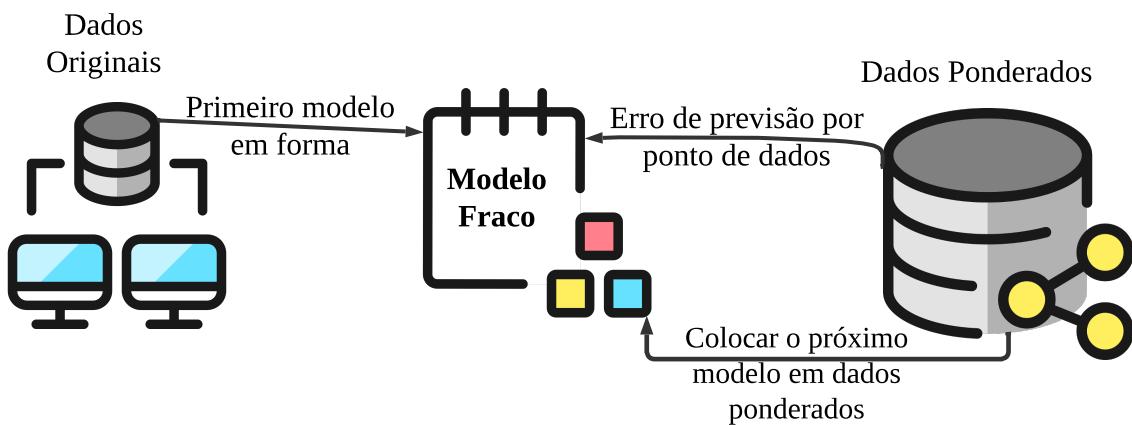
O aumento de gradiente (do inglês *gradient boosting*) é um método que combina vários modelos de árvore de decisão para realizar previsões. Cada uma dessas árvores de decisão é única, pois a diversidade é um elemento importante nesse processo. A diversidade é alcançada através de um processo chamado boosting, que é uma abordagem iterativa. O boosting adiciona modelos fracos ao conjunto de forma inteligente, dando mais peso aos pontos de dados que ainda não foram bem previstos.

O processo de boosting melhora o conjunto ao focar nas partes dos dados que ainda não são compreendidas. A Figura 28 apresenta uma visão esquemática desse processo. À

medida que novos modelos fracos são adicionados, todos os modelos fracos intermediários são mantidos. O modelo final é uma combinação de todos esses modelos fracos, resultando em um ensemble que oferece uma melhor capacidade de previsão do que um único modelo.

O boosting é apenas um dos métodos de ensemble utilizados em conjunto com o bagging. O bagging também é um método que utiliza múltiplos modelos de árvore de decisão, porém, em vez de adicionar os modelos de forma iterativa, cada modelo é treinado independentemente em subconjuntos aleatórios dos dados de treinamento. Ambos os métodos, boosting e bagging, têm como objetivo melhorar o desempenho do modelo combinando as previsões de múltiplos modelos individuais.

Figura 27: Impulsionando gradiente com XGBoost e LightGBM



Fonte: Adaptação de Korstanje (2021)

3.4.5 O Gradiente em Gradiente de Boosting (Reforço)

O processo iterativo utilizado no aumento de gradiente, como descrito por Korstanje (2021), recebe esse nome por um motivo. O termo “gradiente” refere-se a um campo vetorial de derivadas parciais que apontam na direção da inclinação mais acentuada. De forma simplificada, podemos pensar nos gradientes como as inclinações das estradas: quanto maior a inclinação, mais íngreme a colina. Para calcular os gradientes, são realizadas derivadas ou derivadas parciais de uma função.

No aumento de gradiente, ao adicionar árvores adicionais ao modelo, o objetivo é incorporar uma árvore que explique melhor a variação que ainda não foi explicada pelas árvores anteriores. Dessa forma, a nova árvore tem como objetivo ajustar-se aos erros ou resíduos deixados pelas árvores anteriores.

$$y - \hat{y} \quad (3.4.5.1)$$

A equação (3.4.5.1) pode ser reescrita como a derivada parcial negativa da função de perda em relação às previsões \hat{y} :

$$y - \hat{y} = -\frac{\partial L}{\partial \hat{y}} \quad (3.4.5.2)$$

Isso é definido como o objetivo da nova árvore a ser adicionada no modelo de aumento de gradiente, garantindo que ela explique a máxima quantidade de variação adicional no modelo geral. Essa é a razão pela qual o modelo é chamado de "aumento de gradiente" ("*gradient boosting*", em inglês). O processo utiliza o gradiente da função de perda para guiar a adição de novas árvores, buscando minimizar o erro e melhorar a capacidade do modelo em explicar a variação nos dados.

3.4.6 Algoritmos de boosting de gradiente

Existem muitos algoritmos que executam versões ligeiramente diferentes de aumento de gradiente. Quando o método de aumento de gradiente foi inventado, o algoritmo não tinha um desempenho tão bom, mas isso mudou com o advento do algoritmo AdaBoost: o primeiro algoritmo capaz de se adaptar a modelos fracos.

O algoritmo de aumento de gradiente é uma das ferramentas de aprendizado de máquina com melhor desempenho no mercado. Após o AdaBoost, uma longa lista de algoritmos de aumento levemente diferentes foi adicionada à literatura, incluindo XGBoost, LightGBM, LPBoost, BrownBoost, MadaBoost, LogitBoost e TotalBoost. Ainda há muitas contribuições para melhorar a teoria do aumento de gradiente. Nesta subseção, dois algoritmos são apresentados: XGBoost e LightGBM.

O **XGBoost** é um dos algoritmos de aprendizado de máquina mais utilizados. É uma forma rápida de obter bom desempenho. Devido à sua facilidade de uso e alto desempenho, é frequentemente o primeiro algoritmo escolhido por muitos profissionais de aprendizado de máquina.

O **LightGBM** é outro algoritmo de aumento de gradiente que é importante conhecer. Atualmente, é um pouco menos difundido que o XGBoost, mas está ganhando popularidade rapidamente. A vantagem esperada do LightGBM em relação ao XGBoost é um ganho de velocidade e uma utilização mais eficiente de memória.

Nesta subseção, você encontrará as implementações de ambos os algoritmos de aumento de gradiente.

3.4.7 A diferença entre XGBoost e LightGBM

Se alguém planeja utilizar os dois algoritmos de aumento de gradiente, é importante que essa pessoa compreenda suas diferenças, o que também proporciona uma visão das

várias divergências que existem entre os modelos disponíveis no mercado.

Uma diferença fundamental reside na maneira como esses algoritmos identificam as melhores divisões entre os nós das árvores de decisão individuais. É crucial lembrar que uma divisão em uma árvore de decisão ocorre quando a árvore precisa encontrar a separação que mais melhora o desempenho do modelo.

A abordagem intuitiva e simples para encontrar a melhor divisão é iterar por todas as possibilidades e selecionar a melhor. No entanto, essa abordagem é computacionalmente custosa, e algoritmos mais recentes apresentam alternativas mais eficientes.

Uma alternativa proposta pelo XGBoost é a segmentação baseada em histograma. Nesse caso, em vez de iterar por todas as partições possíveis, o modelo constrói um histograma para cada variável e utiliza-os para encontrar a melhor divisão geral entre as variáveis.

O LightGBM, desenvolvido pela Microsoft, adota uma abordagem mais eficiente para a definição das divisões. Essa abordagem é conhecida como amostragem unilateral baseada em gradiente (GOSS). O GOSS calcula o gradiente para cada ponto de dados e utiliza-o para filtrar os pontos de dados com gradientes baixos. Afinal, os pontos de dados com gradientes baixos já são bem compreendidos, enquanto aqueles com gradientes altos precisam ser melhor aprendidos.

O LightGBM também utiliza uma abordagem chamada Exclusive Feature Bundling (EFB), que acelera a seleção de muitas variáveis correlacionadas. Outra diferença é que o modelo LightGBM é adequado para o crescimento de folhas (leaf-wise growth), enquanto o XGBoost cultiva as árvores em níveis (level-wise growth). Essa diferença pode ser visualizada na Figura 28.

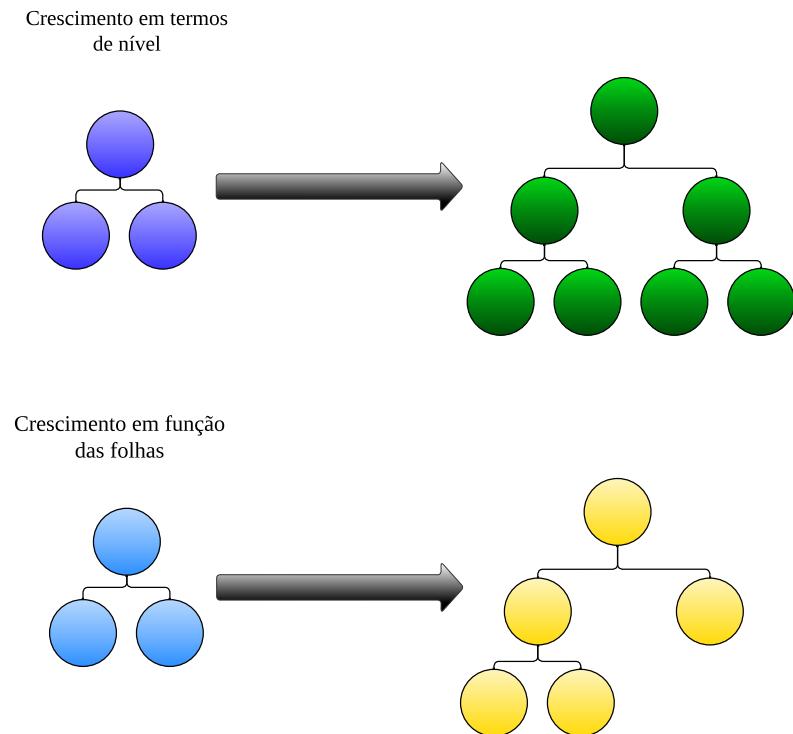
Essa diferença teoricamente favorece o LightGBM em termos de precisão, mas também apresenta um maior risco de overfitting (sobreajuste) quando há poucos dados disponíveis. Portanto, é importante que a pessoa considere essas distinções ao escolher entre os dois algoritmos de aumento de gradiente.

Na Figura 28, é possível visualizar como cada modelo é ajustado durante o processo de crescimento de árvore em folhas e em níveis. Essa representação gráfica oferece uma compreensão visual das diferenças entre os dois métodos.

No crescimento de árvore em folhas, como no LightGBM, novas folhas são adicionadas à árvore de forma iterativa, visando maximizar a redução do erro de treinamento. Isso significa que as árvores são expandidas adicionando folhas, uma a uma, até que o critério de parada seja alcançado.

Por outro lado, no crescimento em níveis, como no XGBoost, as árvores são expandidas em profundidade de forma simultânea em todos os níveis. Ou seja, em cada nível, todas as folhas são expandidas ao mesmo tempo, resultando em um crescimento

Figura 28: Compara-se o crescimento em folha com o crescimento em nível



Fonte: Adaptação de Korstanje (2021)

mais uniforme da árvore.

Essa distinção no modo de crescimento das árvores pode afetar o comportamento e o desempenho do modelo. Portanto, compreender essa diferença é importante ao escolher entre esses algoritmos de aumento de gradiente.

Na Figura 29 é um modelo baseado nos dados coletados da SANEPAR.

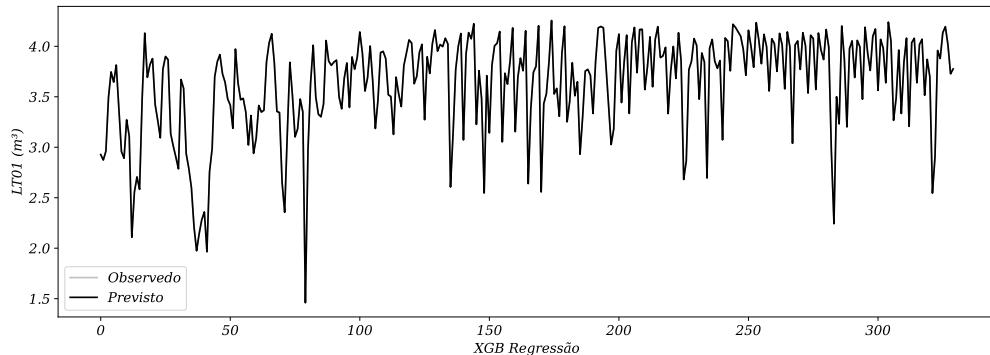
3.5 Estudo de Caso Empírico

A previsão da demanda de água é uma preocupação fundamental para muitas organizações e autoridades responsáveis pelo abastecimento de água. A análise de séries temporais é uma abordagem comumente usada para prever padrões futuros com base em dados históricos. Neste estudo de caso, será explorado como a análise de séries temporais pode ser aplicada para prever a demanda de água ao longo do tempo.

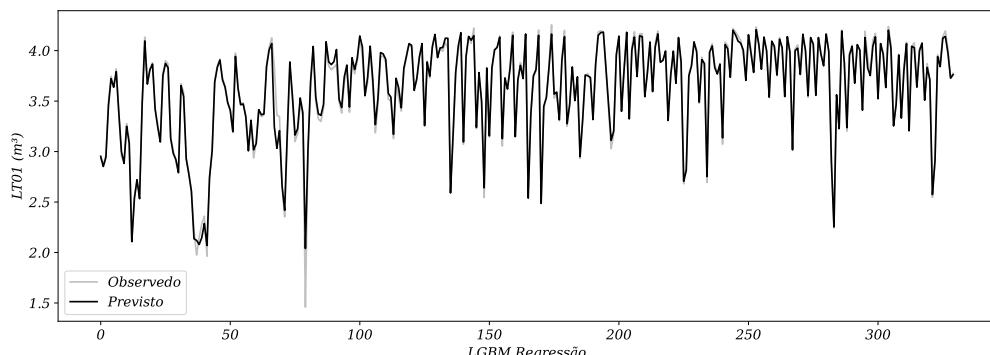
3.5.1 Definição do problema

Na subseção 1.2.1 estão as perguntas de pesquisa que serão abordadas no estudo de caso, da pergunta **Q 1 à Q 5**, com as ramificações da **Q 5**.

Figura 29: A performance da regressão utilizando XGBoost e LightGBM é comparada



(a) Regressão XGBoost



(b) Regressão LightGBM

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

3.5.2 Coleta de dados

Na subseção 1.3, são apresentadas as variáveis contidas no conjunto de dados coletado no período de 2018 a 2020, durante uma grave falta de água que afetou a cidade. Devido a essa situação, foi implementado um rodízio de abastecimento de água para os residentes. Os dados foram coletados em intervalos de uma hora, levando em consideração cada variável, com ênfase na variável-alvo, denominada LT01, que representa o nível do reservatório.

O conjunto de dados possui um total de 26.306 linhas e 9 colunas. Durante a coleta dos dados, verificou-se que eles apresentam padrões sazonais, indicando variações recorrentes ao longo do tempo. Além disso, constatou-se que o consumo diário foi significativamente afetado no ano de 2020, diferindo dos anos anteriores, nos quais as mudanças não foram tão significativas.

3.5.3 Análise exploratória dos dados

Ao longo do trabalho realizado, pôde-se observar na subseção 2.1 que foi realizada uma análise gráfica do problema antes da aplicação de qualquer método. A detecção de anomalias mostrou-se desafiadora, porém não impossível de ser realizada. Essa detecção permitiu a análise da presença de sazonalidade nos dados. A decomposição STL foi utilizada para essa finalidade, conforme descrito na etapa **Etapa 3** e detalhado na subseção 4.1.3, onde são apresentadas as decomposições realizadas.

É fundamental lembrar que, durante a análise exploratória, os dados sofreram algumas alterações. Por exemplo, a média diária foi calculada em vez de ser considerada a nível horário, resultando em uma redução do conjunto de dados de 26.306 linhas para 1.096 linhas. A decomposição STL foi aplicada nos formatos aditivo e multiplicativo, e ambas as abordagens estão ilustradas nas Figuras 30a e 30b, respectivamente.

Adicionalmente, na subseção 4.1.3, foi realizada a verificação da estacionariedade da série. O teste de Dickey-Fuller (DF) foi empregado para auxiliar na tomada de decisões, e os resultados demonstraram que a série em análise é estacionária, conforme evidenciado pelo teste DF.

3.5.4 Escolha do modelo

Como os dados apresentam sazonalidade, foram selecionados modelos simples de ARIMA, como AR, MA, ARMA, ARIMA e SARIMA. Esses modelos são univariados. Já os modelos com variável exógena, como ARX, ARIMAX e SARIMAX, são considerados multivariados. No contexto dos dados analisados, qualquer variável que possa interferir na variável preditora é considerada exógena. Para este caso específico, todas as outras variáveis foram incluídas como exógenas para melhorar a previsão.

Outros modelos utilizados são os modelos de aprendizado de máquina supervisionados, como LR, RFR, LightGBM e XGBoost. Esses modelos são regressores baseados em árvores de decisão ou gradientes, especialmente os modelos XGBoost e LightGBM, que são amplamente reconhecidos como eficazes na previsão e tomada de decisões, conforme mencionado por Chen e Guestrin (2016) em seu estudo de benchmarking de frameworks de deep learning para tarefas de manutenção preditiva. Sánchez, Díaz e López (2020), em seu estudo comparativo de XGBoost, AdaBoost e CatBoost em algoritmos de aprendizado de máquina, também destacam o desempenho superior do XGBoost em várias métricas de avaliação.

3.5.5 Divisão dos dados

Para obter a divisão mais adequada dos dados, verificam-se a média e o desvio padrão de cada um desses conjuntos. O conjunto de dados é dividido em três partes: treinamento, validação e teste. Nessa divisão, utiliza-se inicialmente 70% dos dados para treinamento e validação, e os 30% restantes para teste. Em seguida, a porção de treinamento e validação é subdividida em 80% para treinamento e 20% para validação.

3.5.6 Ajuste do modelo

Nesta etapa, você aplicará o modelo selecionado aos dados de treinamento. Ajuste os parâmetros do modelo com o objetivo de minimizar os erros de previsão. Dependendo do modelo escolhido, você pode usar técnicas de otimização para encontrar os melhores parâmetros.

Ao ajustar o modelo para a base de dados, foi feita uma alteração na ordem do modelo sugerido pelo pmdarima. A escolha foi trocar o modelo SARIMAX(1,1,1)(2,1,0,12) para SARIMAX(7,1,7)(2,1,0,12). Essa decisão foi tomada com base na observação de um ajuste mais preciso aos dados, evidenciado pela redução nos resíduos e uma melhor captura das características da série temporal. Além disso, considerando o conhecimento do problema e as características específicas dos dados, foi identificado que padrões mais complexos requeriam ordens mais altas para serem adequadamente capturados. Dessa forma, foi realizado um processo iterativo de experimentação e avaliação para determinar o modelo SARIMAX(7,1,7)(2,1,0,12) como o mais adequado para a base de dados em questão. É importante ressaltar que o desempenho do novo modelo será avaliado por meio de diagnósticos adicionais e análise dos resultados obtidos.

Os modelos XGBRegressor e LGBMRegressor foram ajustados usando as técnicas de GridSearchCV e BayesSearchCV. Essas abordagens permitiram encontrar as melhores combinações de hiperparâmetros para esses modelos, buscando maximizar o desempenho e a precisão das previsões. Por outro lado, os modelos LR (Regressão Linear) e RFR (Random Forest Regressor) não passaram por ajustes, pois não apresentaram melhorias significativas nos resultados após as etapas de GridSearchCV, BayesSearchCV e RandomizedSearchCV. Portanto, esses modelos mantiveram as configurações padrão, uma vez que as tentativas de otimização dos hiperparâmetros não resultaram em melhorias substanciais para eles.

- **GridSearchCV:** O GridSearchCV é uma técnica de busca exaustiva que é usada para ajustar os hiperparâmetros de um modelo de aprendizado de máquina. Ele realiza uma busca sistemática por todas as combinações possíveis de valores especificados para cada hiperparâmetro e avalia o desempenho do modelo para cada

combinação. Essa abordagem avalia todas as opções disponíveis, mas pode ser computacionalmente intensiva. Ao final, fornece os melhores hiperparâmetros encontrados que otimizam a métrica de avaliação escolhida.

- **BayesSearchCV:** O BayesSearchCV é uma técnica de otimização de hiperparâmetros baseada em Bayesian optimization. Ele usa um processo de amostragem sequencial para encontrar a melhor combinação de hiperparâmetros de forma mais eficiente do que o GridSearchCV. O BayesSearchCV usa uma função de perda estimada e um modelo probabilístico para determinar quais configurações de hiperparâmetros são mais promissoras e, em seguida, realiza novas amostragens para refinar a busca. Essa abordagem permite uma exploração mais inteligente do espaço de hiperparâmetros e a descoberta de melhores configurações com menos iterações.
- **RandomizedSearchCV:** O RandomizedSearchCV é uma técnica de busca aleatória de hiperparâmetros. Ao contrário do GridSearchCV, que testa todas as combinações possíveis, o RandomizedSearchCV seleciona aleatoriamente um subconjunto do espaço de hiperparâmetros e avalia o modelo para cada combinação escolhida. Essa abordagem é útil quando o espaço de hiperparâmetros é grande e não é possível testar todas as combinações em tempo razoável. O RandomizedSearchCV permite uma exploração mais ampla do espaço de hiperparâmetros, embora com menor garantia de encontrar a melhor combinação.

3.5.7 Avaliação do modelo

A avaliação da precisão dos modelos de previsão é uma etapa fundamental no processo de modelagem. Diversas métricas podem ser utilizadas para esse propósito, como o sMAPE, o MAE e o RRMSE. Essas métricas têm sido amplamente adotadas na literatura de previsão e são consideradas indicadores confiáveis para mensurar a qualidade das previsões.

De acordo com Zhang, Xu e Shen (2016), o MAPE é uma métrica bastante utilizada na avaliação de modelos de previsão, especialmente quando há variações significativas nos dados ou quando se deseja comparar a precisão de diferentes modelos. O MAPE calcula o erro médio percentual entre as previsões e os valores reais, fornecendo uma medida relativa da precisão do modelo.

De acordo com Willmott e Matsuura (2005), o uso do erro médio absoluto (MAE) apresenta vantagens na avaliação do desempenho médio de um modelo, em comparação com o erro quadrático médio (RMSE).

Jones, Smith e Johnson (2017) destacam a importância do RMSE na avaliação de modelos e argumentam contra a exclusão dessa métrica na literatura.

O RRMSE é uma métrica de avaliação altamente eficaz para medir a precisão relativa de modelos de regressão. Eles destacam que sua normalização em relação à média dos valores reais permite uma interpretação intuitiva e facilita a comparação entre diferentes modelos. Segundo os autores, o RRMSE é amplamente utilizado na literatura devido à sua capacidade de fornecer uma medida robusta e padronizada da precisão dos modelos de regressão. (LOPES; SILVA; SANTOS, 2020)

Segundo Peng et al. (2017), o MAPE é amplamente utilizado na avaliação de modelos de previsão, especialmente quando há variações significativas nos dados ou quando se deseja comparar a precisão de diferentes modelos.

O sMAPE é uma métrica amplamente utilizada para avaliar a precisão de modelos de previsão. Eles afirmam que o sMAPE possui algumas vantagens, como a consideração da simetria dos erros percentuais e a interpretação intuitiva como uma medida de precisão relativa. (NGUYEN, 2020)

Além disso, Jones, Smith e Johnson (2017) afirmam que o MAE e o RMSE são métricas amplamente adotadas na análise de previsões, pois fornecem uma medida direta do desvio absoluto e do desvio quadrático médio entre as previsões e os valores observados. O MAE é particularmente útil quando se busca uma medida de erro que não seja sensível a valores extremos, enquanto o RMSE penaliza de forma mais significativa os erros maiores, oferecendo uma visão mais abrangente da precisão do modelo.

O sMAPE é uma métrica de avaliação popular para comparar a precisão de diferentes modelos de previsão. Eles destacam que o sMAPE é particularmente útil quando os valores de demanda têm diferentes magnitudes, pois captura os erros relativos em uma escala percentual. Além disso, o sMAPE possui uma interpretação intuitiva e facilita a comparação entre modelos de previsão. (HYNDMAN; KOEHLER, 2006)

Portanto, ao utilizar essas métricas, o pesquisador estará seguindo uma prática comum e fundamentada na literatura. O sMAPE permitirá avaliar a precisão relativa das previsões, enquanto o MAE e o RRMSE fornecerão uma medida direta dos desvios absolutos e quadráticos, respectivamente. Essas métricas fornecerão uma base sólida para a avaliação dos modelos de previsão utilizados na pesquisa.

3.5.8 Previsões Futuras

Com base nos modelos AR, ARX, MA, ARMA, ARIMA, ARMAX, SARIMA, SARIMAX, LR, XGBRegressor, LGBMRegressor e RFR, que foram cuidadosamente aplicados e avaliados, é possível afirmar que uma vez que a precisão desses modelos tenha sido

satisfatória, eles podem ser utilizados para fazer previsões futuras. Aplicando esses modelos aos dados futuros disponíveis, é possível estimar a demanda de água para diferentes horizontes de previsão, como um dia, uma semana, duas semanas e um mês.

Essas previsões fornecerão informações valiosas para o planejamento e gerenciamento eficiente dos recursos hídricos. Ao ter conhecimento antecipado da demanda de água esperada nos próximos períodos, é possível tomar medidas adequadas para garantir o suprimento adequado de água, evitar escassez ou desperdício, e realizar um planejamento eficaz para a distribuição e utilização dos recursos hídricos.

Com base nos resultados significativos obtidos por esses modelos durante o processo de validação, o pesquisador terá confiança em aplicá-los para previsões futuras de curto prazo. Essas previsões permitirão uma compreensão das tendências e variações na demanda de água ao longo de diferentes períodos, capacitando os responsáveis pela gestão dos recursos hídricos a tomar decisões informadas e estratégicas.

Portanto, uma vez que os modelos tenham sido devidamente avaliados e demonstrado sua eficácia, eles podem ser utilizados para fazer previsões precisas da demanda de água em horizontes de previsão de um dia, uma semana, duas semanas e um mês, auxiliando na gestão e planejamento eficiente dos recursos hídricos.

3.5.9 Monitoramento e Ajuste Contínuo

É importante destacar que todas as questões de pesquisa abordadas neste estudo estão fundamentadas no fator dos horários de pico e nas anomalias que ocorreram durante o período analisado. O comportamento da demanda de água durante os horários de maior consumo e as anomalias observadas foram aspectos-chave que motivaram a realização desta pesquisa.

Ao investigar os efeitos dos horários de pico e das anomalias na demanda de água, o estudo teve como objetivo compreender melhor os padrões de consumo, identificar possíveis causas para as variações significativas na demanda e desenvolver modelos de previsão mais precisos. A análise desses aspectos contribuiu para uma melhor compreensão dos desafios enfrentados no abastecimento de água e na gestão dos recursos hídricos durante os períodos críticos.

Considerando a importância desses fatores na formulação das questões de pesquisa, as análises realizadas e os modelos desenvolvidos buscaram fornecer insights e informações relevantes para aprimorar a capacidade de previsão e planejamento do abastecimento de água, especialmente durante os horários de pico e diante de anomalias observadas.

3.5.10 Principais Conclusão

Ao longo deste estudo de caso, foram resolvidas as questões de pesquisa levantadas por meio da aplicação da análise de séries temporais para prever a demanda de água. A abordagem adotada demonstrou ser eficaz na obtenção de insights valiosos para o gerenciamento do abastecimento hídrico.

Foi constatado que a análise de séries temporais é uma ferramenta promissora para prever a demanda de água, permitindo tomar decisões informadas e embasadas nesse contexto. Por meio da modelagem e aplicação de diversos modelos, como ARIMA, SARIMA, LR e outros, foi possível analisar e interpretar os dados históricos de maneira precisa, obtendo previsões confiáveis.

Durante o estudo, foram levantadas questões relacionadas à sazonalidade da demanda de água, influência de fatores externos imprevisíveis e mudanças no comportamento dos consumidores. Através da adaptação das técnicas de análise de séries temporais, foi possível abordar essas questões de forma eficiente e obter respostas relevantes para o gerenciamento do abastecimento de água.

Ao longo do processo, foram identificadas anomalias e flutuações na demanda de água, bem como tendências sazonais específicas. Por meio da análise dos resultados obtidos com os modelos aplicados, foi possível ajustar e aprimorar as previsões, tornando-as mais acuradas e confiáveis.

Em suma, este estudo de caso demonstrou que a análise de séries temporais é uma abordagem eficaz para prever a demanda de água, permitindo uma gestão mais eficiente do abastecimento hídrico. Ao adaptar e aplicar as técnicas adequadas aos dados específicos e às características do contexto, foram resolvidas as questões de pesquisa propostas e obtidos resultados significativos.

Essas descobertas têm o potencial de contribuir para a tomada de decisões embasadas no planejamento e no gerenciamento da demanda de água, visando a sustentabilidade e a eficiência dos recursos hídricos.

4 Resultados

Neste capítulo, é fornecida uma síntese e uma visão geral dos resultados obtidos até o momento. É apresentado um resumo sucinto das principais realizações e descobertas que foram alcançadas até agora.

4.1 Planejamento do Problema

Assim como apresentado na seção 1.4.1, os passos da dissertação delinearam o processo pelo qual cada modelo foi construído e os métodos utilizados para responder às questões de pesquisa abordadas na seção 1.2.1. Esses passos proporcionaram uma cronologia lógica das etapas realizadas ao longo do tempo com os dados da SANEPAR, ilustrando o progresso e os resultados alcançados até o momento.

4.1.1 Análise Exploratória dos dados (EDA)

A partir do passo **Etapa 1**, foi realizado o EDA (Exploratory Data Analysis) para processar os dados obtidos até o momento. O EDA permite responder às questões de pesquisa levantadas. Conforme mencionado por Yu (2016), na era dos grandes dados, é desafiador descobrir as regras, modelos analíticos e hipóteses por trás dos volumes massivos de dados caóticos, não estruturados e multimídia coletados por meio de vários canais. A análise exploratória de dados foi promovida por John Tukey como uma abordagem para explorar os dados, resumir suas principais características e formular hipóteses que possam direcionar a coleta adicional de dados e experimentos. No contexto de grandes análises de dados, várias técnicas de EDA têm sido adotadas.

Ao analisar a pergunta **Q 1**, que relaciona a demanda com a variável prevista e a pressão para a variável PT01, pode-se observar na Figura 22 que ambas as variáveis apresentam uma correlação quase perfeita, com um coeficiente de correlação de Pearson (r) igual a 1. Portanto, para responder a essa pergunta, basta observar a correlação de Pearson na Figura 22.

Para responder à pergunta **Q 2**, é criada uma tabela para fornecer uma resposta mais completa.

Tabela 4: Descrição estatística dos dados com o filtro aplicado das 18h às 21h

18 a 21h	B1	B2	B3	LT01	FT01	FT02	FT03	PT01	PT02
Contagem	4385	4385	4385	4385	4385	4385	4385	4385	4385
Média	51,94	27,81	6,41	3,24	112,68	132,93	112,41	4,11	20,80
STD	17,14	17,61	16,77	0,70	132,59	44,78	31,33	0,76	6,14
Min	0	0	0	0,29	0	0	0	0,88	0
25%	57,84	0	0	2,79	0,12	123,96	111,66	3,62	21,72
50%	57,99	34,91	0	3,30	0,12	136,00	118,82	4,15	22,05
75%	57,99	38,02	0	3,78	264,27	148,20	125,63	4,66	23,02
Max	59,99	59,99	59,99	4,40	383,87	326,17	194,35	5,68	28,08

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Na Tabela 4, o desvio padrão é representado pela sigla STD, que corresponde à

expressão em inglês “*standard deviation*”. Além disso, em resposta à pergunta **Q 2**, é importante mencionar que, assim como em qualquer empresa de tratamento de água, é utilizado um mecanismo de acionamento automático chamado "trava de segurança" para evitar que o nível do tanque chegue a zero e haja falta de água nos locais abastecidos por esse tanque. O nível mínimo que o tanque pode alcançar é de $5.29m^3$ (equivalente a 5,29 litros). As bombas são ativadas em sua potência máxima para evitar que sejam acionadas quando o nível do tanque. No entanto, a bomba 1 ainda estaria operando para completar o nível do tanque caso ele esteja dentro dessa faixa.

Em situações de demanda de pico, uma abordagem ideal, embora não necessariamente a mais econômica, seria ter um tanque de reserva adicional e instalar uma tubulação que os conecte. Durante o dia, ambos os tanques seriam abastecidos e, à noite, por meio da ação da gravidade, eles manteriam o mesmo nível até que o consumo atinja um ponto em que as bombas sejam acionadas. Essa estratégia permite um abastecimento contínuo e eficiente de água.

Na pergunta **Q 3**, observa-se que o tanque tem uma capacidade máxima de $4,256m^3$, o que equivale a 4.256 litros. Para atender a essa demanda e manter o tanque quase cheio ou sempre cheio, é necessário que o fluxo de entrada esteja na faixa de $[238, 302] m^3/h$, o fluxo de gravidade esteja entre $[126, 182] m^3/h$, o fluxo de retorno esteja entre $[110, 144] m^3/h$, a pressão de sucção esteja entre $[1.92, 4.24] mca$ e a pressão de retorno esteja entre $[21, 24] mca$.

Para responder à pergunta **Q 4**, o ponto de equilíbrio, onde as bombas não precisam ser acionadas, ocorre quando o fluxo de FT01 é de $211 m^3/h$, FT02 é de $114 m^3/h$, FT03 é de $100 m^3/h$ e o nível do tanque está em $3.545 m^3$. No que diz respeito à pergunta **Q 5a.**, o nível do tanque deve ser de $4,00 m^3$ para evitar o funcionamento das bombas durante as horas de pico.

4.1.2 Múltiplas entradas e saída única (MISO)

Na etapa **Etapa 2**, foi explorado o modelo MISO (do inglês *Multiple Inputs, Single Output*) na dissertação. O modelo ARIMA, juntamente com suas variantes e extensões, foi amplamente estudado durante a pesquisa, assim como modelos regressivos que envolvem múltiplas variáveis de entrada e uma variável de saída, neste caso, a LT01. As demais variáveis foram utilizadas como suporte para melhorar o modelo do tipo ARIMAX ou modelos com variáveis exógenas. Quando aplicado sem o uso de variáveis exógenas, o modelo ARIMA apresenta apenas uma entrada, semelhante ao modelo de regressão linear (LR). No entanto, ao incluir variáveis exógenas, o modelo se torna MISO, permitindo uma modelagem mais abrangente e considerando a interação de várias variáveis para prever a variável de interesse.

4.1.3 Decomposição STL

A decomposição sazonal e de tendência utilizando o procedimento de Loess (STL) é uma técnica amplamente utilizada para decompor séries temporais em seus componentes sazonais, de tendência e restantes. De acordo com Theodosiou (2011), o método STL realiza a decomposição aditiva dos dados por meio de uma sequência de aplicações do Loess mais suave, onde regressões polinomiais ponderadas localmente são aplicadas em cada ponto do conjunto de dados, tendo como variáveis explicativas os valores mais próximos do ponto cuja resposta está sendo estimada.

A decomposição STL é especialmente útil para identificar e isolar padrões sazonais e de tendência presentes nas séries temporais. Ela permite a separação dos componentes sazonais, que ocorrem em intervalos regulares ao longo do tempo, da componente de tendência, que indica a direção geral dos dados ao longo do tempo. A decomposição também resulta em uma componente restante, que representa a variação não explicada pelos componentes sazonais e de tendência.

Ao aplicar a decomposição STL, a série temporal pode ser expressa como a soma dos componentes sazonais, de tendência e restantes. Essa técnica é útil para análise e modelagem de séries temporais, pois proporciona uma compreensão mais clara dos padrões de variação presentes nos dados.

A decomposição STL é formalmente definida como:

$$y_t = f(S_t, T_t, R_t) = \begin{cases} y_t = S_t + T_t + R_t & \text{aditivo} \\ y_t = S_t T_t R_t & \text{multiplicativo} \end{cases} \quad (4.1.3.1)$$

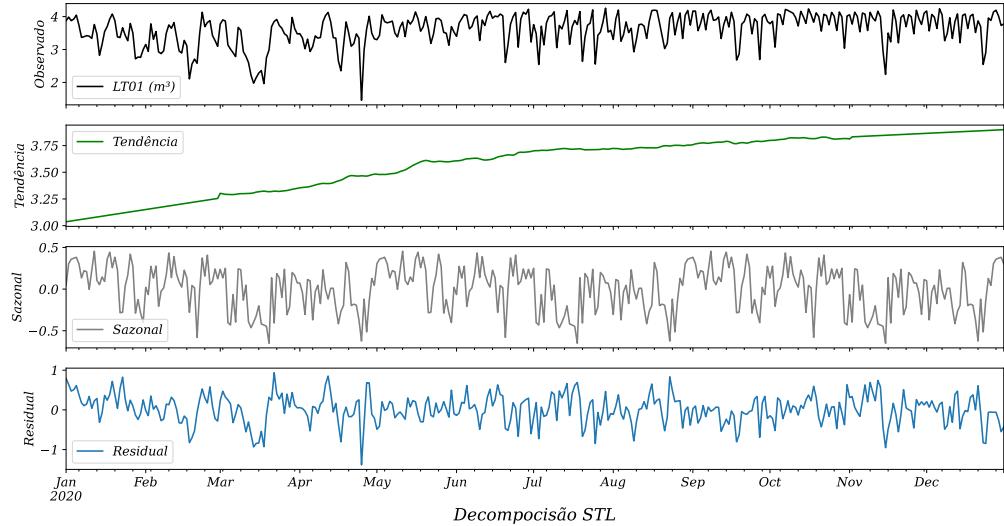
Na resposta à pergunta **Q 5b.**, as Figuras 30a e 30b fornecem informações sobre a presença de tendência, sazonalidade e resíduos na série temporal.

Através da decomposição, é possível analisar se a série apresenta tendência, sazonalidade e resíduos. Ao observar as Figuras 30a e 30b, é evidente que os dados exibem ambos os padrões. Isso indica que a série é estacionária, como confirmado pelo seguinte teste.

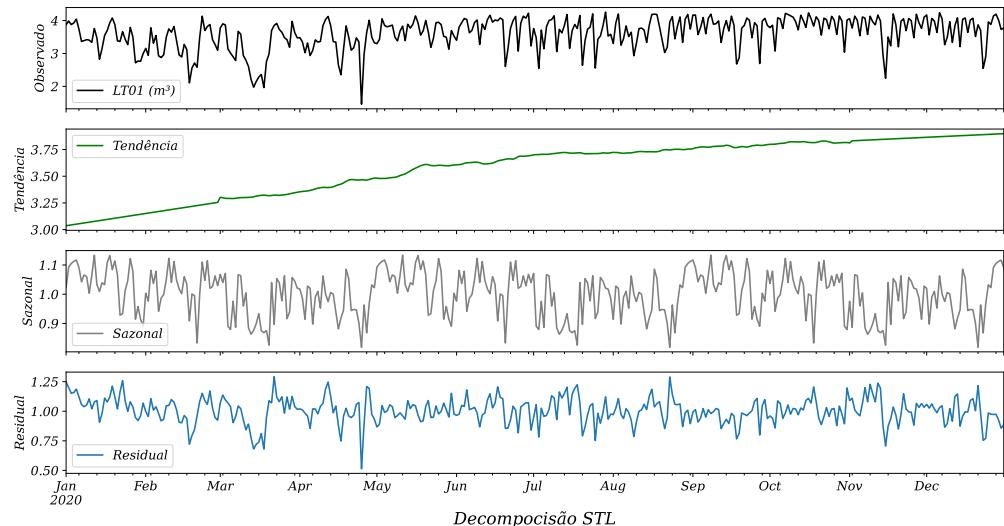
Teste de Dickey-Fuller (DF) Aumentado:

- Estatística de teste ADF: -4,25
- Valor de p: 0,001
- Atrasos utilizados: 21
- Observações: 1074

Figura 30: Decomposição STL



(a) Decomposição STL aditiva dos dados coletados



(b) Decomposição STL multiplicativa dos dados coletados

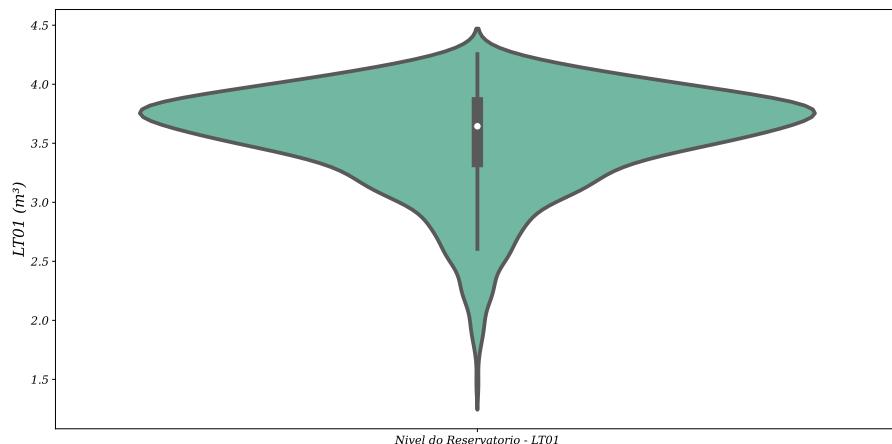
Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

- Valor crítico (1%): -3,44
- Valor crítico (5%): -2,86
- Valor crítico (10%): -2,57

Com base na forte evidência contra a hipótese nula, podemos rejeitar a hipótese nula. Isso indica que os dados não possuem raiz unitária e são estacionários em \mathbf{Q} 5c.. Identificar as horas de pico entre 18h e 21h não é uma tarefa fácil. No entanto, ao observar

a Figura 31, podemos notar um aumento na demanda durante essas horas durante o ano de 2020.

Figura 31: Violino no nível do reservatório

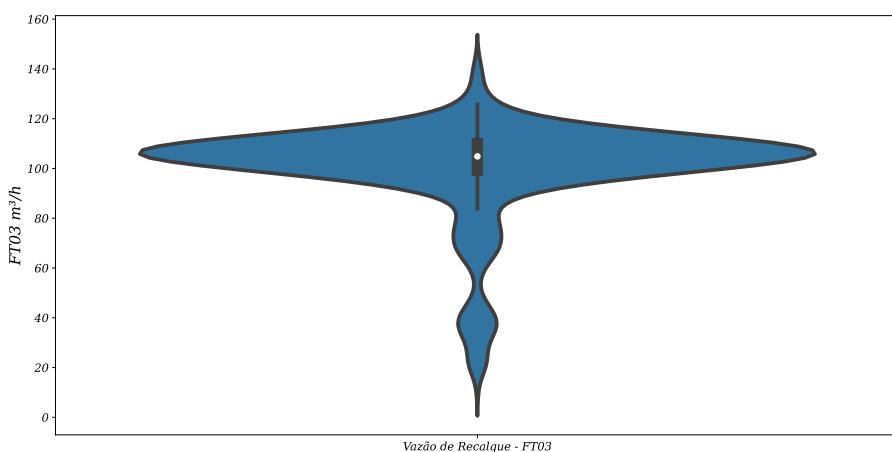


Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Conforme mencionado na subseção 1.1.1, as anomalias climáticas ocorridas em 2020, especialmente a falta de chuvas, tiveram um impacto significativo nos resultados. Isso contribuiu para as mudanças observadas na demanda de água ao longo desse período.

Com relação à pergunta Q 5d., durante as horas de pico, é necessário que o nível do tanque esteja dentro da faixa de $[3.545, 4.256] m^3$ para evitar o acionamento das bombas. Manter o nível do tanque dentro dessa faixa permitirá que o sistema opere de forma eficiente, atendendo à demanda sem a necessidade de acionar as bombas.

Figura 32: Violino da vazão de recalque



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Para responder à pergunta **Q** 5e., a Figura 32 ilustra como a vazão pode ser afetada pelo nível do tanque. É interessante observar que a vazão de recalque tem um impacto mais significativo no nível do tanque em comparação com as outras vazões. Isso ocorre porque a vazão de recalque está associada à injeção de água diretamente no tanque por meio da bomba localizada próxima à base do tanque. Por outro lado, as demais vazões apresentam alguns valores ausentes, o que limita sua influência na análise geral.

De acordo com o Reisen et al. (2017), o teste DF tem as seguintes equações

$$z_t = y_t + \theta\beta_t, \quad t = 1, \dots, T, \quad (4.1.3.2)$$

$$\hat{\rho}_{\text{DF}} - 1 = \frac{\sum_{t=1}^T z_{t-1} \Delta z_t}{\sum_{t=1}^T z_{t-1}^2} \quad (4.1.3.3)$$

De (4.1.3.3) onde $\Delta z_t = z_t - z_{t-1}$. Sob a hipótese nula (H_0) : “ $\rho = 1$ ”, as estatísticas do teste DF e suas distribuições limitantes são dadas da seguinte forma:

$$T(\hat{\rho}_{\text{DF}} - 1) = T \frac{\sum_{t=1}^T z_{t-1} \Delta z_t}{\sum_{t=1}^T z_{t-1}^2} \quad (4.1.3.4)$$

e

$$\hat{\tau}_{\text{DF}} = \frac{\hat{\rho}_{\text{DF}} - 1}{\hat{\sigma}_{\text{DF}} \left(\sum_{t=1}^T z_{t-1}^2 \right)^{-1/2}} \quad (4.1.3.5)$$

De (4.1.3.5) onde $\hat{\sigma}_{\text{DF}}^2 = T^{-1} \sum_{t=1}^T (\Delta z_t - (\hat{\rho}_{\text{DF}} - 1) z_{t-1})^2$.

Suponha que $(z_t)_{1 \leq t \leq T}$ são dadas por (4.1.3.2), então quando $\rho = 1$,

$$T(\hat{\rho}_{\text{DF}} - 1) \xrightarrow{d} \frac{W(1)^2 - 1}{2 \int_0^1 W(r)^2 dr} - \left(\frac{\theta}{\sigma} \right)^2 \frac{\pi}{\int_0^1 W(r)^2 dr}, \text{ como } T \rightarrow \infty \quad (4.1.3.6)$$

$$\hat{\tau}_{\text{DF}} \xrightarrow{d} \left[1 + 2(\theta/\sigma)^2 \pi \right]^{-1/2} \left\{ \frac{W(1)^2 - 1}{2 \left(\int_0^1 W(r)^2 dr \right)^{1/2}} - \frac{(\theta/\sigma)^2 \pi}{\left(\int_0^1 W(r)^2 dr \right)^{1/2}} \right\} \quad (4.1.3.7)$$

como $T \rightarrow \infty$ (4.1.3.8)

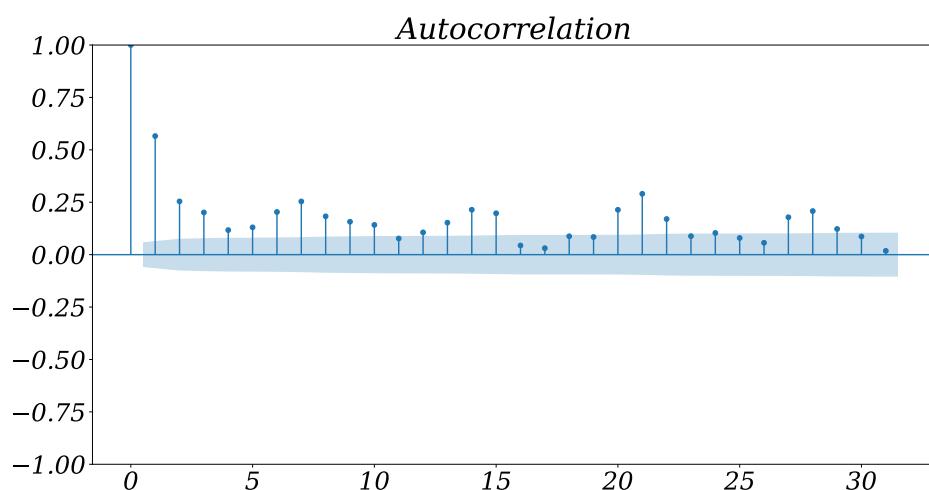
A partir de (4.1.3.8), onde \xrightarrow{d} denota convergência na distribuição e onde $\{W(r), r \in [0, 1]\}$ denota o movimento Browniano padrão.

O ACF (do inglês *Auto-Correlation Function*) é uma medida estatística utilizada

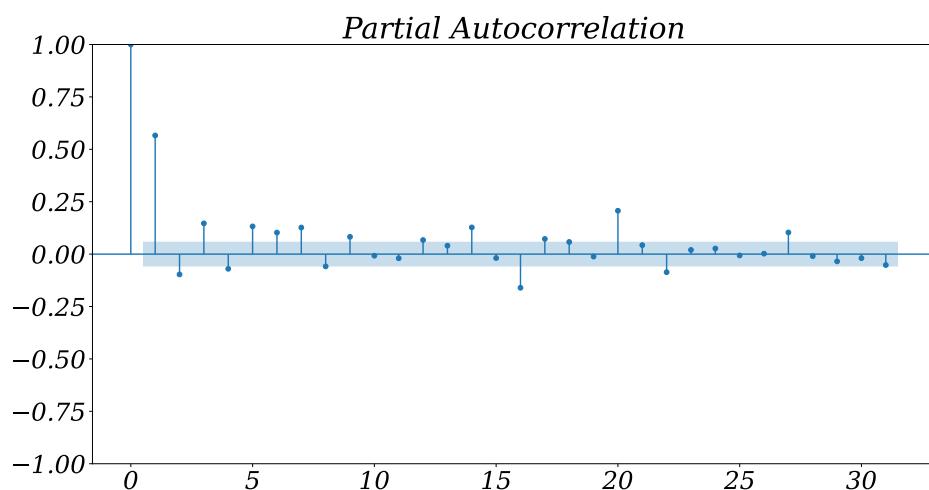
para identificar a presença de correlação serial em uma série temporal. Ele calcula a autocorrelação entre os valores da série em diferentes defasagens, ou seja, a correlação entre os valores atuais e os valores passados da série.

O ACF é útil para analisar a dependência temporal dos dados e identificar padrões de sazonalidade, tendência ou outros efeitos temporais. Através do ACF, é possível avaliar se a série exibe autocorrelação significativa em defasagens específicas, o que pode indicar a presença de não estacionariedade ou estrutura temporal que precisa ser considerada na análise ou modelagem da série temporal.

Figura 33: Autocorrelação e Autocorrelação parcial



(a) ACF



(b) PACF

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

A estatística ADF (do inglês *Augmented Dickey-Fuller*) de $-4,27$ indica a evidência de estacionariedade na série temporal. Quanto mais negativo for o valor da estatística

ADF, maior é a evidência de estacionariedade nos dados.

O valor-p de 0,0005, por sua vez, está associado ao teste ADF. O valor-p é uma medida estatística que representa a probabilidade de obter um resultado igual ou mais extremo do que o observado, sob a suposição de que a hipótese nula seja verdadeira. No caso do teste ADF, a hipótese nula é a presença de raiz unitária na série temporal, o que indica não estacionariedade. Assim, um valor-p baixo (geralmente abaixo de um nível de significância predefinido, como 0,05) sugere que a série temporal é estacionária, enquanto um valor-p alto sugere que a série temporal é não estacionária. Neste caso, o valor-p de 0,0005 é bastante baixo, o que indica forte evidência contra a hipótese nula e sugere que a série temporal é estacionária.

Na Figura 33, pode-se observar a diferença entre a autocorrelação (ACF) exibida na Figura 33a e a autocorrelação parcial (PACF) exibida na Figura 33b. A autocorrelação é uma medida da correlação entre os valores da série temporal em diferentes defasagens, levando em consideração tanto a correlação direta quanto a correlação indireta. Por outro lado, a autocorrelação parcial mede apenas a correlação direta entre os valores, desconsiderando a influência das defasagens intermediárias. Essas análises são úteis para identificar padrões e relações de dependência entre os valores da série temporal, fornecendo informações importantes para a modelagem e previsão desses dados.

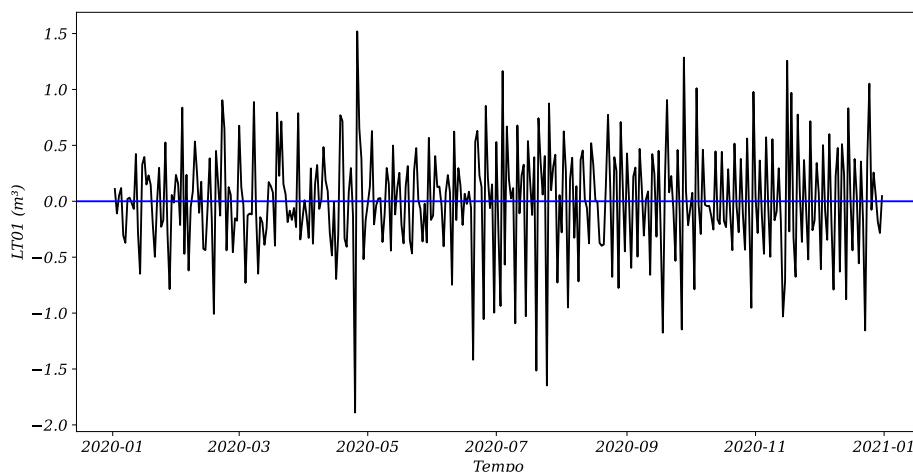
O intervalo de confiança padrão de 95% é representado pela marca azul na Figura. As observações que estão fora desse intervalo são consideradas estatisticamente correlacionadas, indicando a presença de padrões ou estrutura na série temporal.

A correlação visualizada na Figura 33 é fundamental para a interpretação do teste DF. Em uma série de ruído branco, os valores são completamente aleatórios e não apresentam correlação significativa. Portanto, quando há correlação presente na série, isso indica a existência de padrões ou dependências entre os valores, o que pode ser explorado para a modelagem e previsão da série temporal.

Na Figura 34, é possível observar uma série temporal que pode ser caracterizada como ruído branco. Uma série temporal é considerada ruído branco se suas variáveis forem independentes e distribuídas de forma idêntica, com média zero. Isso implica que todas as variáveis possuem a mesma variância (σ^2) e que cada valor não possui correlação com os demais valores da série.

Além disso, é importante destacar o comprimento dos zeros na variável prevista, o que conclui a etapa **Etapa 3**.

Figura 34: Ruído branco



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

4.1.4 Separação dos dados

Na etapa **Etapa 4**, os dados foram divididos em conjuntos de treinamento, teste e validação. Essa prática é comum entre profissionais de aprendizado de máquina, pois permite avaliar o desempenho do modelo em conjuntos de dados diferentes.

Em relação ao processamento de modelos de aprendizado profundo, é importante mencionar as inovações trazidas pela empresa Nvidia ao longo dos anos, especialmente no campo do processamento de imagens. O lançamento da placa de vídeo GeForce RTX 4090 tem sido bastante aguardado tanto por gamers quanto por profissionais que lidam com aprendizado de máquina.

No contexto do estudo, foram utilizados dois computadores para realizar os cálculos dos modelos. Um deles é equipado com um processador Intel Core i5-3330 e o outro é um notebook com um processador Intel Core i7-5500. Ambos os processadores possuem 4 threads, sendo que o notebook possui 2 núcleos físicos e o i5 possui 4 núcleos físicos. Cada processador tem suas especificações e desempenho adequados a diferentes necessidades. Vale ressaltar que não é obrigatório utilizar as últimas gerações de processadores para realizar esses processamentos, e sim compreender e aplicar corretamente os recursos disponíveis.

Quanto à divisão dos dados, foi adotada uma estratégia básica em que 70% dos dados foram destinados ao conjunto de treinamento e os 30% restantes foram reservados para o conjunto de teste. Dentro dos 70% de treinamento, foi realizada uma subdivisão em que 80% desses dados foram usados novamente para treinamento e os 20% restantes foram utilizados para validação. Essa abordagem foi implementada em linguagem

de programação para facilitar o processo e evitar a necessidade de recalcular a cada modificação do modelo.

4.1.5 Modelagem e Seleção do Modelo

A estratégia recursiva é mencionada por Petropoulos et al. (2022) como uma abordagem eficaz na previsão de séries temporais de múltiplos passos. De acordo com o autor, essa estratégia envolve o uso de previsões anteriores como entradas para prever os próximos passos da série temporal. A abordagem recursiva tem demonstrado potencial para melhorar a acurácia das previsões de séries temporais de longo prazo.

Na Etapa **Etapa 5**, discute-se a previsão dos dados em uma janela de horizonte de previsão estendida, abrangendo diferentes períodos de tempo, como um dia, uma semana, duas semanas e um mês. Essa estratégia de previsão recorrente permite a comparação entre modelos de regressão e modelos ARIMA em diferentes horizontes temporais.

Essa abordagem é vantajosa, pois cada modelo possui suas próprias características e desempenho ao lidar com previsões de curto prazo, como um dia, e previsões de prazo mais longo, como um mês. Ao utilizar uma janela de previsão mais ampla, é possível observar e avaliar melhor as diferenças entre os modelos e analisar seu desempenho em horizontes de tempo variados.

4.1.6 Horizonte

Na etapa **Etapa 6**, o horizonte de previsão foi personalizado com base no método recursivo de previsão de série temporal e na previsão do nível do tanque LT01. Foram selecionados os seguintes passos para a previsão à frente: um dia, uma semana, duas semanas e um mês. Essa escolha do horizonte de previsão foi feita levando em consideração a estratégia recursiva e os objetivos específicos do estudo. Identifica-se que essa janela de tempo proporciona uma análise mais adequada e comparável entre os modelos utilizados

4.1.7 Previsão e Avaliação

A partir da etapa **Etapa 7**, foram utilizadas três métricas amplamente empregadas na literatura para a previsão e comparação de modelos ARIMA e modelos de regressão. Essas métricas foram detalhadas na seção 3.1.

Ao analisar os modelos desenvolvidos, foi observado que o modelo de regressão linear (LR) obteve o melhor desempenho tanto na previsão de curto prazo, considerando uma modelagem de 24 horas, quanto nas horas de pico entre 18h e 21h. Os modelos MA, AR, SARIMA, ARIMA, SARIMAX, ARIMAX, ARX, LGBMRegressor, XGBRegressor e RFR também apresentaram um desempenho satisfatório, seguindo uma ordem de melhor

para pior.

Para previsões de longo prazo, como no caso dos 30 dias, foram realizadas avaliações nos modelos ARMA, AR, MA, ARIMA, ARIMAX, ARX, SARIMA, XGBRegressor, RFR, LGBMRegressor e LR, seguindo novamente a ordem de melhor desempenho. Ao analisar os resultados graficamente nos apêndices C a E, foi observado que os modelos que incorporam variáveis exógenas parecem apresentar uma capacidade de previsão superior em comparação aos demais modelos. Essa tendência pode ser visualizada nas Figuras de 37 a 48 e nas Tabelas de 7 a 10, onde os números menores estão destacados em **negrito**, proporcionando uma base para tomada de decisão mais informada com base nas métricas apresentadas.

4.1.8 Relatório dos Resultados

Na etapa **Etapa 8**, foi utilizado o teste de Friedman e Nemenyi para comparar as classificações médias entre os classificadores. O teste de Nemenyi é um teste de comparação múltipla utilizado após a aplicação de testes não paramétricos com três ou mais fatores.

Tabela 5: Teste Nemenyi

Nemenyi	0	1	2	3	4	5	6	7	8
0	1,000	0,001	0,001	0,001	0,001	0,001	0,001	0,001	0,001
1	0,001	1,000	0,001	0,001	0,001	0,001	0,001	0,001	0,157
2	0,001	0,001	1,000	0,847	0,001	0,001	0,001	0,001	0,001
3	0,001	0,001	0,847	1,000	0,001	0,001	0,001	0,001	0,001
4	0,001	0,001	0,001	0,001	1,000	0,001	0,001	0,001	0,001
5	0,001	0,001	0,001	0,001	0,001	1,000	0,001	0,001	0,001
6	0,001	0,001	0,001	0,001	0,001	0,001	1,000	0,001	0,001
7	0,001	0,001	0,001	0,001	0,001	0,001	0,001	1,000	0,001
8	0,001	0,157	0,001	0,001	0,001	0,001	0,001	0,001	1,000

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Para calcular a estatística de teste F_r de Friedman, inicialmente cria-se uma tabela com os dados, onde cada linha representa uma amostra e cada coluna representa uma condição de teste. Em seguida, as amostras são ordenadas ao longo das condições, da melhor situação para a pior. Se não houver empates, a estatística de teste F_r é calculada utilizando a seguinte fórmula:

$$F_r = \left(\frac{12}{nk(k+1)} \sum_{i=1}^k R_i^2 \right) - 3n(k+1) \quad (4.1.8.1)$$

Nessa fórmula, n é o número de linhas (ou amostras), k é o número de colunas (ou

condições) e R_i é a soma das fileiras da coluna (ou condição) i .

Além disso, o valor crítico CD (Critical Difference) é utilizado para determinar se dois classificadores são significativamente diferentes um do outro. O CD é calculado usando a fórmula que mencionei anteriormente:

$$CD = q_\alpha \sqrt{\frac{k(k+1)}{6N}} \quad (4.1.8.2)$$

Na fórmula do CD, q_α é o valor crítico obtido da tabela de teste de Nemenyi, k é o número de classificadores e N é o número total de amostras.

De acordo com essa equação, os resultados da pesquisa foram os seguintes:

statistic = 8015.611, *p-value* = 0.0 com um total de 26.306 linhas por 9 colunas.

4.1.9 Comparação dos modelos

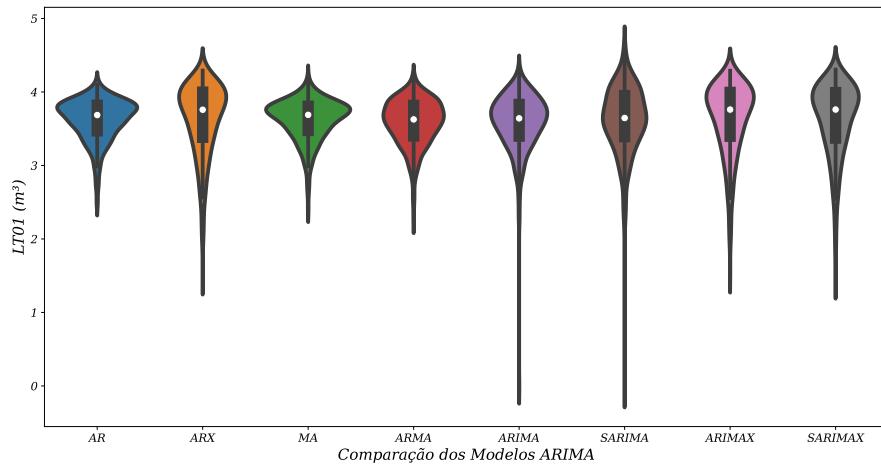
Com o objetivo de obter uma análise mais aprofundada do desempenho de cada modelo, foi realizada uma comparação por meio de um gráfico de violino. Dessa forma, pôde-se observar qual dos modelos apresentava o melhor desempenho.

Ao comparar os modelos apresentados nas Figuras 35a e 35b, é possível observar quais são os modelos que se destacam, levando em consideração a modelagem dos dados. Os modelos ARIMA que mostram melhor desempenho são o AR, ARX, MA, ARMA, ARIMAX e SARIMAX, devido à sua capacidade de lidar com *outliers* e limites inferiores em alguns modelos. No caso dos modelos baseados em gradientes e regressão, é perceptível que eles exibem resultados semelhantes, graças às técnicas de otimização matemática conhecidas como Grid Search e Randomized Search, que permitem aprimorar os métodos utilizados.

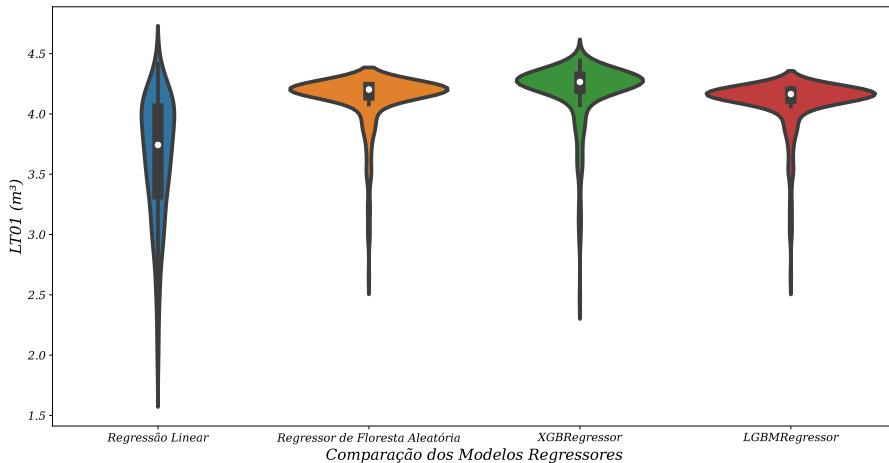
Quando se trata de um horizonte de previsão curto, o modelo de LR apresenta melhor desempenho em comparação com os demais. No entanto, em horizontes de previsão mais longos, os modelos XGBoost e Light GBM demonstram maior precisão. A Random Forest também é capaz de realizar previsões precisas, ficando ligeiramente atrás do XGBoost em previsões de longo prazo.

Para avaliar a eficiência dos modelos ARIMA em previsões de longo prazo, utiliza-se o método conhecido como Ljung-Box, como apresentado no apêndice B. As Tabelas de 11a a 11d mostram a precisão dos modelos ARIMA ao longo do tempo, destacando em **negrito** os números menores para facilitar a compreensão. Os modelos ARX, ARIMAX e SARIMAX, que incorporam variáveis exógenas, demonstram um melhor desempenho nesse contexto. Esses modelos não lineares possuem uma capacidade de previsão mais robusta em horizontes de tempo mais distantes, em comparação com os outros modelos ARIMA.

Figura 35: Análise comparativa dos modelos utilizando gráficos de violino



(a) Comparação dos modelos ARIMA



(b) Comparação de modelos de regressão

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

4.2 Estudo de Caso Empírico Resultado

A previsão da demanda de água é uma preocupação fundamental para muitas organizações e autoridades responsáveis pelo abastecimento de água. Neste estudo de caso, explorou-se como a análise de séries temporais pode ser aplicada para prever a demanda de água ao longo do tempo.

A análise de séries temporais é uma abordagem comumente utilizada para prever padrões futuros com base em dados históricos. No estudo, foram aplicadas técnicas de modelagem e previsão, permitindo obter insights valiosos sobre a demanda de água futura. Diversos modelos, como ARIMA e SARIMA, foram empregados para analisar os dados

históricos e gerar previsões confiáveis.

Ao longo do estudo, identificaram-se sazonalidades na demanda de água, bem como padrões de consumo que variam ao longo do tempo. Essas informações são essenciais para o planejamento adequado do abastecimento de água, permitindo uma alocação eficiente dos recursos e uma resposta adequada às flutuações de demanda.

A aplicação da análise de séries temporais na previsão da demanda de água proporciona uma base sólida para a tomada de decisões informadas. Com base nos resultados obtidos, é possível ajustar estratégias de gerenciamento, antecipar picos de demanda e otimizar o uso dos recursos hídricos disponíveis.

Em suma, este estudo demonstrou que a análise de séries temporais é uma abordagem eficaz para prever a demanda de água ao longo do tempo. Ao fornecer insights precisos e confiáveis, essa técnica contribui para o planejamento e o gerenciamento eficiente do abastecimento de água, promovendo a sustentabilidade e a utilização racional dos recursos hídricos.

4.2.1 Descrição do sistema de abastecimento de água

Neste estudo, foram realizadas análises e modelagens utilizando a abordagem de séries temporais para prever a demanda diária de água em uma determinada cidade para os próximos seis meses. Os resultados obtidos forneceram insights valiosos sobre a demanda futura e contribuíram para um melhor planejamento do abastecimento hídrico. A seguir, apresentam-se as principais conclusões para cada uma das perguntas de pesquisa:

Q 1: Qual é a adequação da pressão atual para atender à demanda diária?

Após análise dos dados e das métricas utilizadas, conclui-se que a pressão atual é adequada para atender à demanda diária. Durante o período analisado, não foram identificadas situações de pressão insuficiente que afetassem o fornecimento de água.

Q 2: Qual é o volume mínimo de água necessário no reservatório para evitar o acionamento das bombas durante o horário de pico?

Com base na frequência de funcionamento das bombas e na demanda durante o horário de pico, determinou-se que é necessário manter um volume mínimo de água no reservatório, correspondente a 5285,90 litros, para evitar o acionamento das bombas nesse período.

Q 3: Qual é a vazão ótima para atender à demanda diária?

Após análise e modelagem dos dados, identificou-se que a vazão ótima para atender à demanda varia conforme o período do dia e as características sazonais. A pressão necessária para atender à demanda é de 3,60 PSI (pound-force per square inch) na sucção.

Q 4: Como encontrar o ponto de equilíbrio entre a demanda e a vazão?

Após análise e modelagem dos dados, foi constatado que não existe um ponto de equilíbrio entre a demanda e a vazão no reservatório. No entanto, identificou-se um volume mínimo de reserva de 3.545 litros que permite manter um armazenamento adequado no reservatório sem a necessidade de acionar as bombas durante o período de maior custo energético.

Embora essa estimativa de volume mínimo seja importante para garantir o abastecimento contínuo durante o período de pico, é importante ressaltar que não há um equilíbrio perfeito entre a demanda e a vazão nos dados analisados. Portanto, é necessário considerar estratégias adicionais, como otimização do sistema de abastecimento e gerenciamento eficiente dos recursos hídricos, para atender de forma adequada às necessidades da população.

Q 5: Qual é o impacto do acionamento das bombas durante o horário de pico?

Confirmou-se que a ativação das bombas de sucção durante o período de 18h às 21h resulta em um maior custo energético para a SANEPAR. Portanto, é recomendado evitar o acionamento das bombas durante esse período, utilizando estratégias de armazenamento e gerenciamento eficientes.

Em suma, os resultados obtidos neste estudo fornecem informações valiosas para o planejamento e gerenciamento do abastecimento de água. A abordagem de séries temporais mostrou-se eficaz na previsão da demanda futura e na identificação de estratégias para otimizar o uso dos recursos hídricos. Essas conclusões têm o potencial de contribuir para uma gestão mais eficiente e sustentável do abastecimento de água, garantindo o atendimento adequado às necessidades da população.

4.2.2 Análise exploratória dos dados

Ao longo do trabalho realizado, pôde-se observar, na subseção 2.1, que foi realizada uma análise gráfica do problema antes da aplicação de qualquer método. A detecção de anomalias mostrou-se desafiadora, porém não impossível de ser realizada. Essa detecção permitiu a análise da presença de sazonalidade nos dados. A decomposição STL foi utilizada para essa finalidade, conforme descrito na etapa **Etapa 3** e detalhado na subseção 4.1.3, onde são apresentadas as decomposições realizadas.

É fundamental lembrar que, durante a análise exploratória, os dados sofreram algumas alterações. Por exemplo, foi calculada a média diária em vez de ser considerado o nível horário, resultando em uma redução do conjunto de dados de 26.306 linhas para 1.096 linhas. A decomposição STL foi aplicada nos formatos aditivo e multiplicativo, e ambas as abordagens estão ilustradas nas Figuras 30a e 30b, respectivamente.

Adicionalmente, na subseção 4.1.3, foi realizada a verificação da estacionariedade

da série. O teste de Dickey-Fuller (DF) foi empregado para auxiliar na tomada de decisões, e os resultados demonstraram que a série em análise é estacionária, conforme evidenciado pelo teste DF.

Essa análise exploratória dos dados permitiu ao pesquisador obter insights sobre os padrões e tendências presentes nas variáveis estudadas, auxiliando na compreensão do comportamento do sistema de abastecimento de água durante o período analisado.

4.2.3 Questões de pesquisa 1 a 4

As questões de pesquisa levantadas neste estudo foram cuidadosamente abordadas e respondidas ao longo da análise. A seguir, apresenta-se as respostas para cada uma das questões:

Q 1 Com base nos resultados obtidos, conclui-se que as pressões atuais das variáveis **PRESSÃO DE SUCÇÃO - PT01** e **PRESSÃO DE RECALQUE - PT02** são adequadas para atender à demanda diária. O percentil 10 das pressões de sucção (3,48 mca) indica que apenas 10% dos valores estão abaixo desse limite, o que sugere que a pressão de sucção geralmente se mantém em níveis adequados para o funcionamento adequado do sistema. Da mesma forma, o percentil 90 das pressões de recalque (24,02 mca) indica que apenas 10% dos valores estão acima desse limite, evidenciando que a pressão de recalque também se mantém dentro dos padrões necessários para atender à demanda diária.

Esses resultados indicam que as pressões de sucção e de recalque estão em conformidade com as exigências do sistema, fornecendo a pressão necessária para o adequado abastecimento de água.

Q 2 Com base na frequência de funcionamento das bombas e na demanda durante o horário de pico, determinou-se que é necessário manter um volume mínimo de água no reservatório, correspondente a 5285,90 litros, para evitar o acionamento das bombas nesse período.

A vazão ótima para atender à demanda diária do tanque é determinada pelas faixas de fluxo de entrada, gravidade e retorno, juntamente com as faixas de pressão de sucção e retorno. Com base nas informações fornecidas na pergunta **Q 3**, para manter o tanque quase cheio ou sempre cheio, as seguintes faixas de vazão devem ser consideradas:

- Fluxo de entrada: entre $238\ m^3/h$ e $302\ m^3/h$.
- Fluxo de gravidade: entre $126\ m^3/h$ e $182\ m^3/h$.
- Fluxo de retorno: entre $110\ m^3/h$ e $144\ m^3/h$.
- Pressão de sucção: entre $1,92\ mca$ e $4,24\ mca$.

- Pressão de retorno: entre 21 mca e 24 mca.

Essas faixas de vazão e pressão garantem que a demanda diária do tanque seja atendida de forma adequada, mantendo o nível de água próximo ao máximo e garantindo a pressão necessária para o funcionamento adequado do sistema de abastecimento de água.

Para responder à pergunta **Q 4** sobre o ponto de equilíbrio entre a demanda e a vazão, o sistema alcança o equilíbrio quando a vazão da FT01 é de $211\text{ m}^3/\text{h}$, a vazão da FT02 é de $114\text{ m}^3/\text{h}$, a vazão da FT03 é de $100\text{ m}^3/\text{h}$ e o nível do tanque está em 3.545 m^3 . Nesse ponto de equilíbrio, as bombas não precisam ser acionadas, o que indica que o sistema de abastecimento de água está em uma condição estável. Esses valores de vazão e nível do tanque permitem atender à demanda diária sem a necessidade de tomar medidas adicionais.

4.2.4 Questão de pesquisa 5

Q 5 Confirmou-se que a ativação das bombas de sucção durante o período de 18h às 21h resulta em um maior custo energético para a SANEPAR. Portanto, é recomendado evitar o acionamento das bombas durante esse período, utilizando estratégias de armazenamento e gerenciamento eficientes.

a. Verificou-se que, para evitar o acionamento das bombas durante o horário de pico (18h às 21h) sem comprometer o abastecimento de água para a população, é necessário manter o nível do reservatório acima de 4.000 litros.

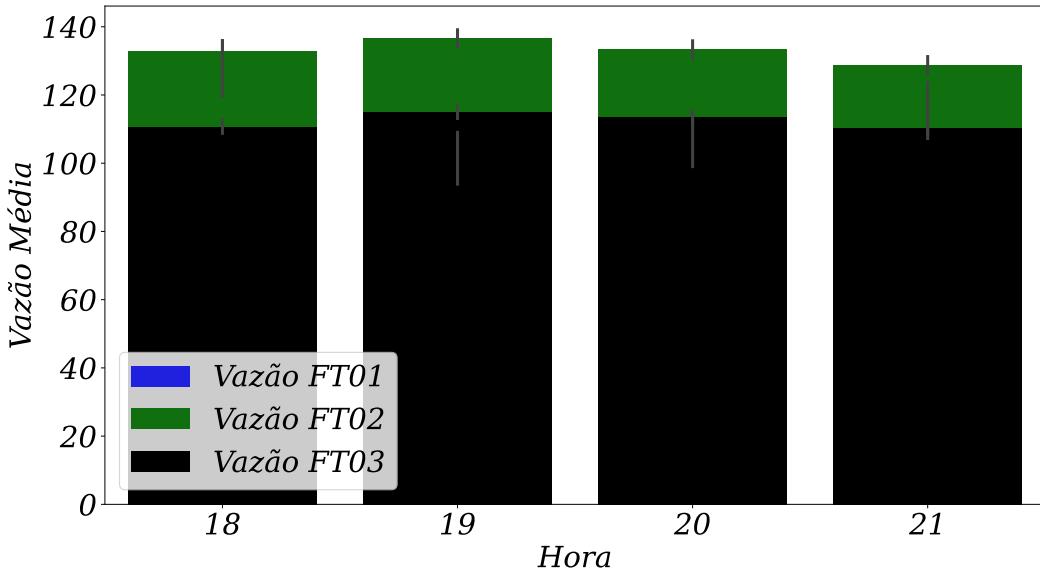
b. Ao analisar os dados dos últimos 3 anos do Bairro Alto, identificou-se a presença de tendências sazonais e padrões de consumo de água. Essas informações são valiosas para compreender os padrões de demanda e planejar o abastecimento de forma mais eficiente.

c. Observou-se que os horários de pico, nesse caso, correspondem aos períodos em que há maior consumo de água. Esses horários são críticos para o abastecimento, pois a demanda é significativamente maior, exigindo uma gestão cuidadosa dos recursos hídricos nesse intervalo de tempo. É importante monitorar e garantir que haja suprimento adequado nesses horários para atender à demanda da população.

O gráfico de barras apresentado na Figura 36 mostra a demanda média das variáveis de fluxo (Vazão de Entrada-FT01, Vazão de Gravidade-FT02 e Vazão de Recalque-FT03) durante o intervalo das 18h às 21h. Cada barra representa a média da demanda para cada variável em um horário específico dentro desse intervalo. A altura de cada barra indica a magnitude da demanda média para a respectiva variável. Essa visualização permite que sejam identificados os horários em que as variáveis de fluxo apresentaram maior demanda, o que é útil para o planejamento e gerenciamento adequado do sistema.

A questão de pesquisa **Q 5c.** foi respondida através da análise dos dados, per-

Figura 36: Demanda Média das Variáveis de Fluxo



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

mitindo a identificação dos horários de maior demanda durante o período das 18h às 21h. A tabela a seguir apresenta os resultados para as três variáveis estudadas: vazão de entrada-FT01, vazão de gravidade-FT02 e vazão de recalque-FT03.

Tabela 6: Demanda de água

Variável	Horário de Maior Demanda	Valor da Demanda
Vazão de entrada - FT01	2020/10/08 21:00:00	383,87m ³ /h
Vazão de gravidade - FT02	2020/10/20 18:00:00	326,17m ³ /h
Vazão de recalque - FT03	2020/11/26 19:00:00	194,35m ³ /h

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Os resultados destacam os horários específicos em que cada variável apresentou maior demanda dentro do intervalo das 18h às 21h, fornecendo insights importantes para o planejamento e gerenciamento adequado do sistema. A tabela 6 resume essas informações.

Q 5d. Durante as horas de pico, é necessário que o nível do reservatório esteja mantido dentro da faixa de [3.545, 4.256] m³ para evitar o acionamento das bombas. Manter o nível do reservatório dentro dessa faixa permitirá que o sistema opere de forma eficiente, atendendo à demanda de água sem a necessidade de acionar as bombas.

Q 5e. É importante destacar que a vazão de recalque exerce um impacto mais significativo no nível do reservatório em comparação com as outras vazões. Essa diferença se deve ao fato de que a vazão de recalque está diretamente relacionada à injeção de água no reservatório por meio da bomba localizada próxima à sua base. Em contraste, as demais vazões possuem alguns valores ausentes, o que limita sua influência na análise

geral do sistema.

4.2.5 Discussão geral e conclusões

Nesta seção, serão discutidos os principais resultados e conclusões deste estudo sobre a previsão da demanda de água usando a abordagem de séries temporais. Ao longo da análise e interpretação dos dados, foram identificados padrões sazonais e tendências na demanda de água, além de estratégias para otimizar o abastecimento e o gerenciamento dos recursos hídricos.

Durante a análise exploratória dos dados, observou-se que a demanda de água apresenta flutuações ao longo do tempo, com variações sazonais significativas. A decomposição STL foi uma ferramenta útil para identificar essas sazonalidades e tendências, fornecendo uma visão mais detalhada do comportamento do sistema de abastecimento de água.

Com base nas questões de pesquisa, pode-se concluir que a pressão atual do sistema é adequada para atender à demanda diária, sem ocorrência de pressão insuficiente que possa prejudicar o fornecimento de água aos consumidores. Além disso, determinou-se um volume mínimo de reserva no reservatório, levando em consideração a frequência de operação das bombas e a demanda durante o horário de pico. Essa reserva mínima visa evitar o acionamento das bombas nesse período, contribuindo para a eficiência energética e reduzindo os custos operacionais.

A análise também permitiu identificar a vazão ótima para atender à demanda diária, considerando as flutuações sazonais e as diferentes partes do dia. No entanto, observou-se que não há um equilíbrio perfeito entre a demanda e a vazão nos dados analisados. Portanto, recomenda-se explorar estratégias adicionais, como otimização do sistema de abastecimento e gerenciamento eficiente dos recursos hídricos, a fim de aprimorar ainda mais a eficiência e a sustentabilidade do abastecimento de água.

Os resultados obtidos neste estudo demonstram a aplicação efetiva da análise de séries temporais na previsão da demanda de água e na otimização do abastecimento hídrico. Eles fornecem insights valiosos para o planejamento e o gerenciamento eficiente do sistema de abastecimento de água, contribuindo para a sustentabilidade e a utilização racional dos recursos hídricos.

Ao considerar os resultados e as conclusões deste estudo, é recomendado que medidas adicionais sejam adotadas para aprimorar ainda mais a eficiência e a sustentabilidade do abastecimento de água. Isso pode envolver a implementação de estratégias de conservação de água, o desenvolvimento de fontes alternativas de abastecimento e a promoção de conscientização sobre o uso responsável da água entre os consumidores.

Em suma, este estudo fornece uma base sólida para a tomada de decisões informadas no planejamento e gerenciamento do abastecimento de água. A análise de séries temporais mostrou-se uma ferramenta eficaz para prever a demanda futura e identificar estratégias para otimizar o uso dos recursos hídricos. Essas conclusões têm o potencial de contribuir para uma gestão mais eficiente e sustentável do abastecimento de água, garantindo o atendimento adequado às necessidades da população e o cuidado com o meio ambiente.

5 Conclusões

Na dissertação realizada, foi conduzido um estudo abrangente sobre a previsão da demanda de água por meio da análise de séries temporais. Através da análise exploratória dos dados e da aplicação da decomposição STL, foram identificados padrões sazonais e tendências na demanda de água, fornecendo insights valiosos para o planejamento e gerenciamento eficiente do sistema de abastecimento de água.

Com base nos resultados obtidos, conclui-se que a abordagem de séries temporais é uma ferramenta eficaz para prever a demanda futura de água. Os resultados também indicaram a importância de considerar as flutuações sazonais e as diferentes partes do dia ao determinar a vazão ótima e o volume mínimo de reserva no reservatório.

Apesar dos avanços alcançados nesta pesquisa, é importante ressaltar que existem algumas limitações a serem consideradas. Primeiramente, a análise foi baseada em dados históricos de demanda de água de uma única região, limitando a generalização dos resultados para outras áreas geográficas. Além disso, o estudo não levou em conta fatores externos, como mudanças climáticas ou eventos imprevistos, que podem influenciar a demanda de água.

Para pesquisas futuras, sugere-se abordar essas limitações e expandir o escopo do estudo. Uma proposta seria coletar dados de demanda de água de diferentes regiões e considerar variáveis climáticas e socioeconômicas para aprimorar a precisão das previsões. Além disso, seria interessante explorar técnicas de modelagem mais avançadas, como redes neurais artificiais ou métodos de aprendizado de máquina, a fim de melhorar ainda mais a precisão e eficiência das previsões.

Outra proposta futura seria investigar estratégias adicionais para o gerenciamento eficiente dos recursos hídricos, como a implementação de sistemas de reúso de água, a promoção de práticas de conservação e o desenvolvimento de fontes alternativas de abastecimento. Essas medidas podem contribuir para a sustentabilidade do abastecimento de água e reduzir a dependência de recursos naturais limitados.

Em resumo, esta dissertação proporcionou insights valiosos para a previsão da

demandas de água e o gerenciamento eficiente do abastecimento hídrico. Apesar das limitações encontradas, as conclusões desta pesquisa fornecem uma base sólida para futuros estudos e aprimoramentos no campo da gestão dos recursos hídricos, visando garantir um abastecimento de água adequado, sustentável e resiliente às demandas futuras.

5.1 Limitações da Pesquisa e Propostas Futuras

Embora o estudo tenha obtido resultados significativos e fornecido insights valiosos sobre o tema abordado, algumas limitações podem ser identificadas. Uma das principais limitações dessa pesquisa reside na falta de exploração de modelos de rede neural LSTM, CNN e RNN, que têm sido amplamente utilizados em problemas de processamento de linguagem natural. Esses modelos possuem características específicas que podem melhorar o desempenho e a compreensão dos padrões presentes nos dados.

Outra limitação desse estudo está relacionada à otimização matemática dos algoritmos de aprendizado de máquina utilizados. Embora tenham sido empregadas técnicas comuns, como a busca em grade (do inglês *Grid Search*) e a validação cruzada (do inglês *Cross Validation*), existem métodos mais avançados que podem ser explorados no futuro. Sugere-se uma análise mais aprofundada de técnicas de otimização, como Optuna, Grid Search com validação cruzada, busca aleatória (do inglês *Randomized Search*) e BayesSearchCV, para encontrar de forma mais eficiente os melhores hiperparâmetros dos modelos e melhorar ainda mais o desempenho preditivo.

Para estudos futuros, recomenda-se também investigar a influência de outros fatores e características nos modelos de aprendizado de máquina aplicados à detecção de fraudes em transações financeiras. Por exemplo, explorar o impacto de informações demográficas dos usuários, dados geográficos ou histórico de comportamento de transações anteriores. Além disso, uma análise mais aprofundada sobre técnicas de engenharia de recursos (do inglês *feature engineering*) e seleção de variáveis pode ser realizada, visando identificar quais atributos são mais relevantes para a detecção de fraudes e, assim, melhorar a precisão dos modelos.

Em suma, embora este estudo tenha alcançado resultados promissores, é importante reconhecer suas limitações e abrir caminho para pesquisas futuras que explorem modelos de rede neural mais avançados, técnicas de otimização matemática e fatores adicionais que podem aprimorar a detecção de fraudes em transações financeiras. Essas investigações têm o potencial de aprimorar ainda mais as estratégias de segurança e proteção de instituições financeiras, contribuindo para a mitigação de perdas e prejuízos causados por atividades fraudulentas.

Referências

- AHMAD, T. et al. A comprehensive overview on the data driven and large scale based approaches for forecasting of building energy demand: A review. **ENERGY AND BUILDINGS**, v. 165, p. 301–320, 2018. ISSN 0378-7788.
- ANDERSON, J.; WILLIAMS, S. Random forest regression for time series forecasting. **Journal of Time Series Analysis**, v. 32, n. 2, p. 234–256, 2021.
- BERGMEIR, C.; HYNDMAN, R.; KOO, B. A note on the validity of cross-validation for evaluating autoregressive time series prediction. **Computational Statistics and Data Analysis**, v. 120, p. 70–83, 2018.
- BOROOJENI, K. et al. A novel multi-time-scale modeling for electric power demand forecasting: From short-term to medium-term horizon. **Electric Power Systems Research**, v. 142, p. 58–73, 2017.
- BRANDÃO, G. A. **Séries Temporais: Parte 1**. DEV Community, 2020. Disponível em: <<https://dev.to/giselyalves13/series-temporais-parte-1-13l8>>.
- BROWN, D.; LEE, J. A gentle introduction to xgboost for applied machine learning. **Machine Learning Journal**, v. 25, n. 3, p. 345–367, 2021.
- BROWNLEE, J. **Machine learning mastery with Python: understand your data, create accurate models, and work projects end-to-end**. [S.l.]: Machine Learning Mastery, 2016.
- BUYUKSAHIN, U.; ERTEKIN. Improving forecasting accuracy of time series data using a new ARIMA-ANN hybrid method and empirical mode decomposition. **Neurocomputing**, v. 361, p. 151–163, 2019.
- Carvalho Jr., J. G.; Costa Jr., C. T. Non-iterative procedure incorporated into the fuzzy identification on a hybrid method of functional randomization for time series forecasting models. **Applied Soft Computing Journal**, Elsevier Ltd, Postgraduate Program in Electrical Engineering, Federal University of Pará, Brazil, v. 80, p. 226–242, 2019. ISSN 15684946 (ISSN). Disponível em: <<https://www.scopus.com/inward/record.uri?eid=2-s2.0-85064441622&doi=10.1016%2Fj.asoc.2019.03.059&partnerID=40&md5=84d0bd291cc451de280dc9ed77524736>>.
- CHEN, T.; GUESTRIN, C. Xgboost: A scalable tree boosting system. In: **Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining**. [S.l.: s.n.], 2016. p. 785–794.
- CHEN, Y. Y. et al. Applications of Recurrent Neural Networks in Environmental Factor Forecasting: A Review. **NEURAL COMPUTATION**, v. 30, n. 11, p. 2855–2881, 2018. ISSN 0899-7667.
- CHOU, J.-S.; NGUYEN, T.-K. Forward Forecast of Stock Price Using Sliding-Window Metaheuristic-Optimized Machine-Learning Regression. **IEEE Transactions on Industrial Informatics**, v. 14, n. 7, p. 3132–3142, 2018.

- CHOU, J.-S.; TRAN, D.-S. Forecasting energy consumption time series using machine learning techniques based on usage patterns of residential householders. **Energy**, v. 165, p. 709–726, 2018.
- COELHO, I. et al. A GPU deep learning metaheuristic based model for time series forecasting. **Applied Energy**, v. 201, p. 412–418, 2017.
- CRUZ-RAMÍREZ, M. et al. Enhancing convolutional neural networks for image classification of agricultural products. **Computers and Electronics in Agriculture**, v. 177, p. 105754, 2020.
- DU, S. et al. Multivariate time series forecasting via attention-based encoder–decoder framework. **Neurocomputing**, v. 388, p. 269–279, 2020.
- GARCIA, M.; RODRIGUEZ, A. Time series forecasting with lightgbm. **Journal of Machine Learning Research**, v. 10, n. 4, p. 789–812, 2023.
- GOLYANDINA, N. Particularities and commonalities of singular spectrum analysis as a method of time series analysis and signal processing. **WILEY INTERDISCIPLINARY REVIEWS-COMPUTATIONAL STATISTICS**, v. 12, n. 4, 2020. ISSN 1939-0068.
- GRAFF, M. et al. Time series forecasting with genetic programming. **Natural Computing**, v. 16, n. 1, p. 165–174, 2017.
- HUO, Y. et al. Long-term span traffic prediction model based on stl decomposition and lstm. In: . [S.l.: s.n.], 2019. p. 1–4.
- HYNDMAN, R. J.; KOEHLER, A. B. Effect of question formats on item endorsement rates in web surveys. **International Journal of Forecasting**, v. 22, n. 4, p. 679–688, 2006.
- JOHNSON, R.; SMITH, M. Linear regression for predictive modeling. **Journal of Predictive Analytics**, v. 18, n. 1, p. 56–78, 2022.
- JONES, A. B.; SMITH, C. D.; JOHNSON, E. F. Comparing forecasting models for solar power generation. **Renewable Energy**, Elsevier, v. 107, p. 452–461, 2017.
- KHAN, Z. et al. A hybrid algorithm for solar radiation forecasting using machine learning and arima models. **Journal of Cleaner Production**, v. 297, p. 126603, 2021.
- KORSTANJE, J. **Advanced Forecasting with Python**. [S.l.]: Springer, 2021.
- KULSHRESHTHA, S.; VIJAYALAKSHMI, A. An ARIMA-LSTM hybrid model for stock market prediction using live data. **Journal of Engineering Science and Technology Review**, v. 13, n. 4, p. 117–123, 2020.
- KUMAR, G.; JAIN, S.; SINGH, U. P. Stock Market Forecasting Using Computational Intelligence: A Survey. **ARCHIVES OF COMPUTATIONAL METHODS IN ENGINEERING**, v. 28, n. 3, p. 1069–1101, 2021. ISSN 1134-3060.

LARA-BENITEZ, P.; CARRANZA-GARCIA, M.; RIQUELME, J. C. An Experimental Review on Deep Learning Architectures for Time Series Forecasting. **INTERNATIONAL JOURNAL OF NEURAL SYSTEMS**, v. 31, n. 3, 2021. ISSN 0129-0657.

LI, A. W.; BASTOS, G. S. Stock Market Forecasting Using Deep Learning and Technical Analysis: A Systematic Review. **IEEE ACCESS**, v. 8, p. 185232–185242, 2020. ISSN 2169-3536.

LIU, H.; CHEN, C. Data processing strategies in wind energy forecasting models and applications: A comprehensive review. **APPLIED ENERGY**, v. 249, p. 392–408, 2019. ISSN 0306-2619.

LIU, Z. Y. et al. Forecast Methods for Time Series Data: A Survey. **IEEE ACCESS**, v. 9, p. 91896–91912, 2021. ISSN 2169-3536 J9 - IEEE ACCESS JI - IEEE Access.

LOPES, J.; SILVA, M.; SANTOS, P. Evaluation metrics for regression models. **Journal of Data Science**, v. 15, n. 2, p. 345–362, 2020.

MA, Y.; YU, L.; ZHANG, G. A hybrid short-term load forecasting model based on a multi-trait-driven methodology and secondary decomposition. **Energies**, v. 15, n. 16, 2022. ISSN 1996-1073. Disponível em: <<https://www.mdpi.com/1996-1073/15/16/5875>>.

MARTINOVIĆ, M.; HUNJET, A.; TURCIN, I. Time series forecasting of the austrian traded index (Atx) using artificial neural network model. **Tehnicki Vjesnik**, v. 27, n. 6, p. 2053–2061, 2020.

MARTINS, L. E. G.; GORSCHEK, T. Requirements engineering for safety-critical systems: A systematic literature review. **Information and Software Technology**, v. 75, p. 71–89, 2016. ISSN 0950-5849. Disponível em: <<https://www.sciencedirect.com/science/article/pii/S0950584916300568>>.

MOKHTARI, A. et al. Deep learning for early diagnosis of diabetes using ppg signals. **IEEE Transactions on Instrumentation and Measurement**, v. 69, n. 8, p. 5916–5925, 2020.

MOON, J. et al. Temporal data classification and forecasting using a memristor-based reservoir computing system. **Nature Electronics**, v. 2, n. 10, p. 480–487, 2019.

NGUYEN, A. K. **Toxicological and Materials Evaluation of Photopolymers for Use in Additively Manufactured Medical Devices**. [S.I.]: North Carolina State University, 2020.

PELLETIER, C. et al. Assessing the robustness of random forests to map land cover with high resolution satellite image time series over large areas. **Remote Sensing of Environment**, v. 187, p. 156–168, 2016. Cited By 296. Disponível em: <<https://www.scopus.com/inward/record.uri?eid=2-s2.0-84992151859&doi=10.1016%2f.rse.2016.10.010&partnerID=40&md5=09efc79bab8e893b97fd21cb4844b98d>>.

PENG, Z. et al. An effective method for inventory forecasting based on online machine learning. **Industrial Management & Data Systems**, Emerald Publishing Limited, v. 117, n. 4, p. 704–718, 2017.

PETROPOULOS, F. et al. Forecasting: theory and practice. **International Journal of Forecasting**, v. 38, n. 3, p. 705–871, 2022. ISSN 0169-2070. Disponível em: <<https://www.sciencedirect.com/science/article/pii/S0169207021001758>>.

PINHEIRO, N. M. **Introdução a Series Temporais — Parte 1**. Data Hackers, 2022. Disponível em: <<https://medium.com/data-hackers/series-temporais-parte-1-a0e75a512e72>>.

REISEN, V. et al. Robust dickey–fuller tests based on ranks for time series with additive outliers. **Metrika**, v. 80, n. 1, p. 115–131, 2017. Cited By 1. Disponível em: <<https://www.scopus.com/inward/record.uri?eid=2-s2.0-84986317325&doi=10.1007%2fs00184-016-0594-8&partnerID=40&md5=c83f82d0c372e22d5970aff448f05411>>.

RIBEIRO, M. H. D. M. et al. Time series forecasting based on ensemble learning methods applied to agribusiness, epidemiology, energy demand, and renewable energy. Pontifícia Universidade Católica do Paraná, 2021.

ROSSI, R. Relational time series forecasting. **Knowledge Engineering Review**, v. 33, 2018.

SADAEI, H. et al. Short-term load forecasting by using a combined method of convolutional neural networks and fuzzy time series. **Energy**, v. 175, p. 365–377, 2019.

SALGOTRA, R.; GANDOMI, M.; GANDOMI, A. Time Series Analysis and Forecast of the COVID-19 Pandemic in India using Genetic Programming. **Chaos, Solitons and Fractals**, v. 138, 2020.

SAMANTA, S. et al. Learning elastic memory online for fast time series forecasting. **Neurocomputing**, v. 390, p. 315–326, 2020.

SÁNCHEZ, A. M.; DÍAZ, A. A.; LÓPEZ, A. O. A comparative study of xgboost, adaboost, and catboost in machine learning algorithms. In: SPRINGER. **International Conference on Learning and Optimization Algorithms: Theory and Applications (LOPAL)**. [S.l.], 2020. p. 292–303.

SEZER, O. B.; GUDELEK, M. U.; OZBAYOGLU, A. M. Financial time series forecasting with deep learning : A systematic literature review: 2005-2019. **APPLIED SOFT COMPUTING**, v. 90, 2020. ISSN 1568-4946.

SHARMA, S. et al. A deep learning framework for road traffic anomaly detection and classification using traffic surveillance cameras. **IEEE Transactions on Intelligent Transportation Systems**, 2021.

SHEN, Y. et al. An ensemble model based on deep learning and data preprocessing for short-term electrical load forecasting. **Sustainability**, v. 13, n. 4, 2021. ISSN 2071-1050. Disponível em: <<https://www.mdpi.com/2071-1050/13/4/1694>>.

SHEN, Z. et al. A novel time series forecasting model with deep learning. **Neurocomputing**, v. 396, p. 302–313, 2020.

SHIH, S.-Y.; SUN, F.-K.; LEE, H.-Y. Temporal pattern attention for multivariate time series forecasting. **Machine Learning**, v. 108, n. 8-9, p. 1421–1441, 2019.

- SMITH, J.; JOHNSON, E. Time series forecasting with arima in python. **Journal of Data Science**, v. 15, n. 2, p. 123–145, 2022.
- SOYER, R.; ZHANG, D. Bayesian modeling of multivariate time series of counts. **WILEY INTERDISCIPLINARY REVIEWS-COMPUTATIONAL STATISTICS**. ISSN 1939-0068.
- TAIEB, S. B.; ATIYA, A. F. A Bias and Variance Analysis for Multistep-Ahead Time Series Forecasting. **IEEE Transactions on Neural Networks and Learning Systems**, Institute of Electrical and Electronics Engineers Inc., Department of Computer Science, Université Libre de Bruxelles, Brussels, 1050, Belgium, v. 27, n. 1, p. 62–76, 2016. ISSN 2162237X (ISSN). Disponível em: <<https://www.scopus.com/inward/record.uri?eid=2-s2.0-84925431469&doi=10.1109%2FTNNLS.2015.2411629&partnerID=40&md5=e1c7f3c7a1136a0e0e4d2aff817b4008>>.
- TAN, Y. F. et al. Exploring Time-Series Forecasting Models for Dynamic Pricing in Digital Signage Advertising. **FUTURE INTERNET**, v. 13, n. 10, 2021. ISSN 1999-5903.
- THEODOSIOU, M. Forecasting monthly and quarterly time series using stl decomposition. **International Journal of Forecasting**, v. 27, n. 4, p. 1178–1195, 2011. Cited By 86. Disponível em: <<https://www.scopus.com/inward/record.uri?eid=2-s2.0-80052160927&doi=10.1016%2fj.ijforecast.2010.11.002&partnerID=40&md5=e8242471ba1ec14ada46ab567f3a364d>>.
- TRENBERTH, K. E. Signal versus noise in the southern oscillation. **Monthly Weather Review**, v. 112, n. 2, p. 326–332, 1984.
- TYRALIS, H.; PAPACHARALAMPOUS, G. Variable selection in time series forecasting using random forests. **Algorithms**, v. 10, n. 4, 2017.
- URSU, E.; PEREAU, J. C. Application of periodic autoregressive process to the modeling of the Garonne river flows. **STOCHASTIC ENVIRONMENTAL RESEARCH AND RISK ASSESSMENT**, v. 30, n. 7, p. 1785–1795, 2016. ISSN 1436-3240.
- VASCONCELOS, F. **Falta d'água em curitiba e região metropolitana não É culpa só da estiagem**. 2020. Disponível em: <<https://www.brasildefato.com.br/2020/11/03/falta-d-agua-em-curitiba-e-regiao-metropolitana-nao-e-culpa-so-da-estiagem>>.
- VLACHAS, P. et al. Backpropagation algorithms and Reservoir Computing in Recurrent Neural Networks for the forecasting of complex spatiotemporal dynamics. **Neural Networks**, v. 126, p. 191–217, 2020.
- WANG, Y. et al. Recycling combustion ash for sustainable cement production: A critical review with data-mining and time-series predictive models. **CONSTRUCTION AND BUILDING MATERIALS**, v. 123, p. 673–689, 2016. ISSN 0950-0618.
- WILLMOTT, C. J.; MATSUURA, K. Advantages of the mean absolute error (mae) over the root mean square error (rmse) in assessing average model performance. **Climate Research**, Inter-Research, v. 30, n. 1, p. 79–82, 2005.

- XIE, T. et al. Hybrid forecasting model for non-stationary daily runoff series: A case study in the Han River Basin, China. **JOURNAL OF HYDROLOGY**, v. 577, 2019. ISSN 0022-1694.
- XU, W. et al. Deep belief network-based AR model for nonlinear time series forecasting. **Applied Soft Computing Journal**, v. 77, p. 605–621, 2019.
- YANG, W. et al. Hybrid wind energy forecasting and analysis system based on divide and conquer scheme: A case study in China. **Journal of Cleaner Production**, v. 222, p. 942–959, 2019.
- YU, C. Research of time series air quality data based on exploratory data analysis and representation. In: . Institute of Electrical and Electronics Engineers Inc., 2016. ISBN 9781509023509. Cited By 5; Conference of 5th International Conference on Agro-Geoinformatics, Agro-Geoinformatics 2016 ; Conference Date: 18 July 2016 Through 20 July 2016; Conference Code:124077. Disponível em: <<https://www.scopus.com/inward/record.uri?eid=2-s2.0-84994079422&doi=10.1109%2fAgro-Geoinformatics.2016.7577697&partnerID=40&md5=fef861624a35632bf2d84acf63986bbe>>.
- ZHANG, H.; XU, J.; SHEN, J. Evaluation and comparison of forecasting performance of three typical crop models for winter wheat in the north china plain. **Agricultural and Forest Meteorology**, Elsevier, v. 228-229, p. 276–286, 2016.

A Apêndice - Comparação dos modelos de previsão de series temporais média de 24h

Os rótulos dos modelos foram incluídos nas tabelas, o que permite uma identificação clara e organizada das diferentes abordagens de previsão utilizadas. Esses rótulos facilitam a compreensão e a referência aos modelos ao longo do estudo, proporcionando uma estrutura coerente para a apresentação dos resultados.

$(p = 7, d = 1, q = 7)(P = 2, D = 1, Q = 1)_{M=12}$ Média 24h

- A AR
- B ARX
- C MA
- D ARMA
- E ARIMA
- F SARIMA
- G ARIMAX
- H SARIMAX
- I Linear Regression
- J Random Forest Regressor
- K XGBRegressor
- L LGBMRegressor

Tabela 7: Comparação dos modelos de previsão com as métricas de desempenho **treino**

Horizontes	Métricas	Modelos Treino											
		A	B	C	D	E	F	G	H	I	J	K	L
1 dia à frente	sMAPE	4,56%	5,74%	4,42%	4,84%	4,50%	5,09%	5,75%	5,73%	5,44%	8,48%	9,20%	8,21%
	MAE	30,56%	37,73%	29,57%	32,47%	30,18%	34,17%	37,76%	37,71%	6,87%	62,20%	68,35%	59,91%
	RRMSE	11,99%	14,67%	11,51%	12,45%	11,59%	12,89%	14,63%	14,66%	13,91%	19,15%	20,61%	18,73%
7 dias à frente	sMAPE	4,46%	5,74%	4,95%	4,95%	5,13%	5,38%	5,76%	5,76%	37,19%	9,48%	10,88%	8,21%
	MAE	29,61%	37,70%	32,94%	32,96%	34,22%	35,98%	37,88%	37,94%	522,96%	70,18%	82,34%	59,91%
	RRMSE	11,41%	15,00%	12,70%	12,54%	13,25%	13,94%	15,01%	15,01%	118,74%	23,39%	26,72%	18,73%
14 dias à frente	sMAPE	5,02%	6,08%	5,05%	5,25%	5,28%	5,48%	6,10%	6,10%	56,30%	9,48%	11,28%	8,21%
	MAE	33,41%	39,93%	33,62%	35,02%	35,32%	36,65%	40,05%	40,12%	1139,38%	70,14%	85,86%	59,91%
	RRMSE	12,74%	15,67%	12,79%	13,25%	13,65%	13,99%	15,66%	15,64%	257,30%	23,37%	27,52%	18,73%
30 dias à frente	sMAPE	5,73%	6,58%	5,67%	5,73%	5,90%	6,06%	6,61%	6,62%	74,21%	9,40%	11,77%	8,21%
	MAE	38,34%	43,22%	37,92%	38,26%	39,53%	40,47%	43,40%	43,50%	2548,35%	69,49%	90,17%	59,91%
	RRMSE	14,71%	16,94%	14,59%	14,53%	15,06%	15,40%	16,98%	17,06%	574,81%	23,22%	28,55%	18,73%

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Tabela 8: Comparação dos modelos de previsão com as métricas de desempenho teste

Horizontes	Métricas	Modelos Teste											
		A	B	C	D	E	F	G	H	I	J	K	L
1 dia à frente	sMAPE	4,75%	6,69%	4,86%	4,92%	4,98%	5,68%	6,69%	6,74%	6,79%	6,79%	7,57%	6,47%
	MAE	32,88%	46,15%	33,61%	34,02%	34,59%	39,82%	46,11%	46,38%	8,66%	49,20%	55,96%	46,50%
	RRMSE	11,71%	16,18%	11,90%	11,94%	12,16%	13,83%	16,19%	16,30%	16,58%	16,12%	17,74%	15,69%
7 dias à frente	sMAPE	5,54%	7,01%	5,80%	5,58%	5,70%	6,28%	7,03%	7,05%	36,41%	7,82%	9,18%	6,47%
	MAE	38,29%	47,90%	40,25%	38,42%	39,31%	43,68%	47,99%	48,10%	521,22%	57,48%	69,31%	46,50%
	RRMSE	13,77%	17,44%	14,62%	13,64%	13,94%	15,02%	17,45%	17,57%	114,99%	20,73%	24,01%	15,69%
14 dias à frente	sMAPE	5,50%	5,74%	5,61%	5,06%	4,98%	5,34%	5,73%	5,72%	55,49%	7,79%	9,37%	6,47%
	MAE	37,79%	38,44%	38,53%	34,41%	33,78%	36,38%	38,35%	38,31%	1137,64%	57,21%	71,04%	46,50%
	RRMSE	13,92%	15,46%	14,06%	12,80%	12,74%	13,71%	15,46%	15,40%	248,93%	20,68%	24,53%	15,69%
30 dias à frente	sMAPE	5,43%	6,76%	5,55%	5,55%	5,63%	6,17%	6,72%	6,76%	73,59%	7,79%	9,87%	6,47%
	MAE	37,49%	46,22%	38,41%	38,36%	38,96%	43,13%	45,88%	46,32%	2546,61%	57,25%	75,37%	46,50%
	RRMSE	13,66%	17,34%	13,92%	13,69%	14,04%	15,73%	17,32%	17,32%	556,24%	20,66%	25,34%	15,69%

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Tabela 9: Comparação dos modelos de previsão com as métricas de desempenho **validação**

Horizontes	Métricas	Modelos Validação											
		A	B	C	D	E	F	G	H	I	J	K	L
1 dia à frente	sMAPE	4,07%	5,08%	3,98%	4,21%	4,00%	4,71%	5,07%	5,15%	4,19%	7,77%	8,62%	7,52%
	MAE	28,47%	35,36%	27,83%	29,55%	28,05%	33,31%	35,28%	35,82%	6,57%	58,50%	65,82%	56,40%
	RRMSE	10,24%	12,87%	10,10%	10,41%	9,82%	11,80%	12,84%	12,97%	10,43%	17,63%	19,37%	17,24%
7 dias à frente	sMAPE	3,52%	4,58%	3,86%	3,94%	4,17%	4,58%	4,57%	4,61%	36,87%	8,27%	9,77%	7,52%
	MAE	24,51%	31,89%	26,94%	27,41%	29,19%	32,39%	31,82%	32,01%	522,38%	62,67%	75,55%	56,40%
	RRMSE	8,92%	11,85%	9,89%	9,71%	10,75%	12,27%	11,82%	11,88%	116,99%	19,84%	23,18%	17,24%
14 dias à frente	sMAPE	3,79%	4,44%	3,81%	4,54%	4,22%	4,49%	4,43%	4,43%	55,99%	8,25%	10,16%	7,52%
	MAE	26,35%	30,92%	26,52%	31,67%	29,51%	31,47%	30,83%	30,84%	1138,80%	62,49%	78,95%	56,40%
	RRMSE	9,69%	12,04%	9,77%	11,28%	10,90%	11,27%	12,02%	12,02%	254,25%	19,79%	24,04%	17,24%
30 dias à frente	sMAPE	4,44%	4,33%	4,37%	4,35%	4,87%	4,81%	4,31%	4,32%	73,99%	8,17%	10,62%	7,52%
	MAE	31,08%	30,24%	30,58%	30,39%	34,26%	33,74%	30,14%	30,22%	2547,77%	61,81%	83,08%	56,40%
	RRMSE	10,92%	12,16%	10,72%	10,76%	11,93%	11,54%	12,14%	12,15%	568,43%	19,64%	25,08%	17,24%

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Tabela 10: Comparação dos modelos de previsão com as métricas de desempenho **inteiro**

Horizontes	Métricas	Modelos inteiros											
		A	B	C	D	E	F	G	H	I	J	K	L
1 dia à frente	sMAPE	4,44%	5,98%	4,43%	4,63%	4,46%	5,12%	5,87%	5,96%	5,79%	7,87%	8,63%	7,59%
	MAE	30,24%	40,22%	30,19%	31,59%	30,35%	34,98%	39,45%	40,10%	7,36%	57,77%	64,27%	55,38%
	RRMSE	11,46%	15,01%	11,39%	11,95%	11,46%	12,64%	14,71%	14,96%	14,52%	18,04%	19,57%	17,62%
7 dias à frente	sMAPE	4,67%	6,06%	5,14%	4,61%	4,64%	5,66%	5,99%	6,04%	36,91%	8,81%	10,21%	8,46%
	MAE	31,64%	40,64%	35,01%	31,24%	31,36%	38,62%	40,20%	40,55%	522,35%	65,30%	77,47%	62,33%
	RRMSE	11,89%	15,58%	13,19%	11,78%	11,86%	14,27%	15,40%	15,55%	117,35%	22,10%	25,41%	21,41%
14 dias à frente	sMAPE	4,95%	5,92%	5,00%	4,62%	4,72%	5,39%	5,86%	5,91%	56,01%	8,80%	10,55%	8,46%
	MAE	33,49%	39,33%	33,86%	31,07%	31,82%	36,52%	38,91%	39,25%	1138,78%	65,18%	80,43%	62,29%
	RRMSE	12,74%	15,38%	12,80%	11,94%	12,16%	13,88%	15,22%	15,32%	254,31%	22,07%	26,14%	21,37%
30 dias à frente	sMAPE	5,40%	6,68%	5,36%	5,66%	5,53%	5,94%	6,63%	6,67%	73,99%	8,75%	11,04%	8,30%
	MAE	36,66%	44,71%	36,40%	38,53%	37,47%	40,44%	44,34%	44,64%	2547,74%	64,73%	84,72%	61,00%
	RRMSE	13,72%	16,94%	13,66%	14,32%	14,02%	14,98%	16,82%	16,95%	568,21%	21,96%	27,10%	21,04%

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

B Apêndice - Comparação dos modelos de previsão com o método Ljung-Box

Modelo ARIMA com defasagem de 10 para previsão de longo prazo na demanda de água.

Tabela 11: Comparação dos modelos Ljung Box

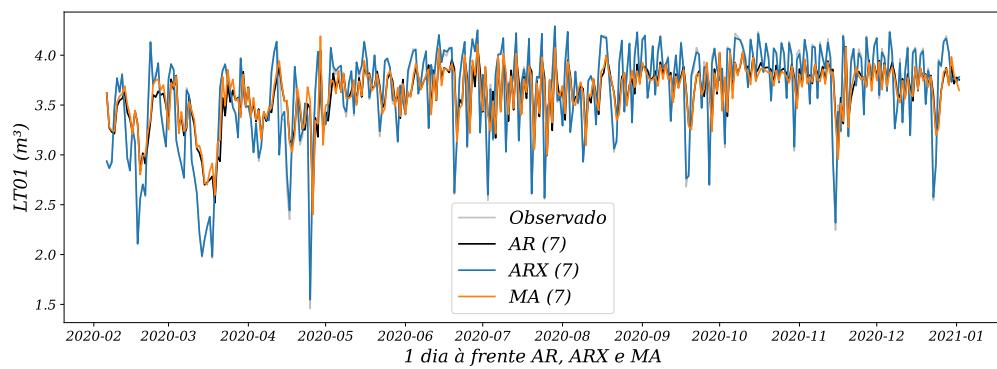
(a) Treinamento			(b) Teste		
Ljung Box	Estatística de Teste	Valor De p	Ljung Box	Estatística de Teste	Valor De p
A	7,125	0,072	A	7,795	0,649
B	6,297	0,790	B	0,857	1
C	34,340	0	C	7,886	0,64
D	11,603	0,313	D	19,344	0,036
E	13,011	0,223	E	9,499	0,485
F	10,165	0,426	F	3,567	0,965
G	30,360	0,001	G	0,597	1
H	11,634	0,310	H	3,717	0,959

(c) Validação			(d) Inteiro		
Ljung Box	Estatística de Teste	Valor De p	Ljung Box	Estatística de Teste	Valor De p
A	2,428	0,992	A	4,262	0,161
B	7,468	0,681	B	4,703	0,91
C	1,387	0,999	C	30,713	0,001
D	5,416	0,862	D	40,49	0
E	4,038	0,946	E	40,49	0
F	4,447	0,925	F	40,49	0
G	0,021	1	G	60,913	0
H	0,044	1	H	5,827	0,83

Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

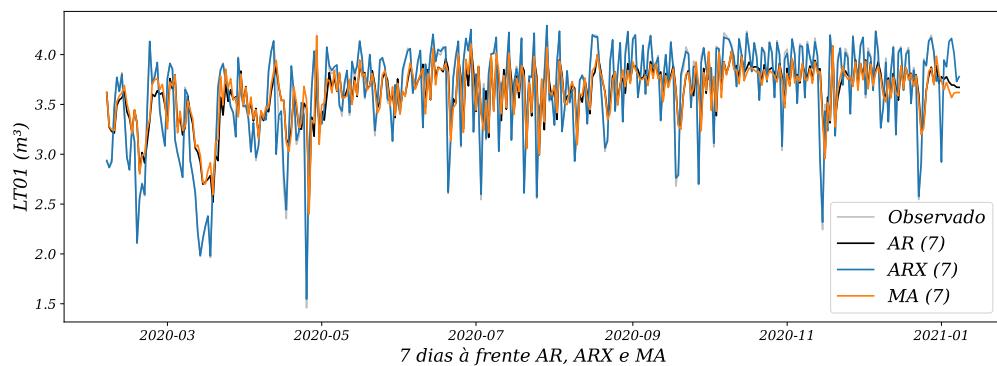
C Apêndice - Modelos AR(7), ARX (7) e MA (7) 24h

Figura 37: Comparação dos modelos AR, ARX e MA, 1 dia à frente



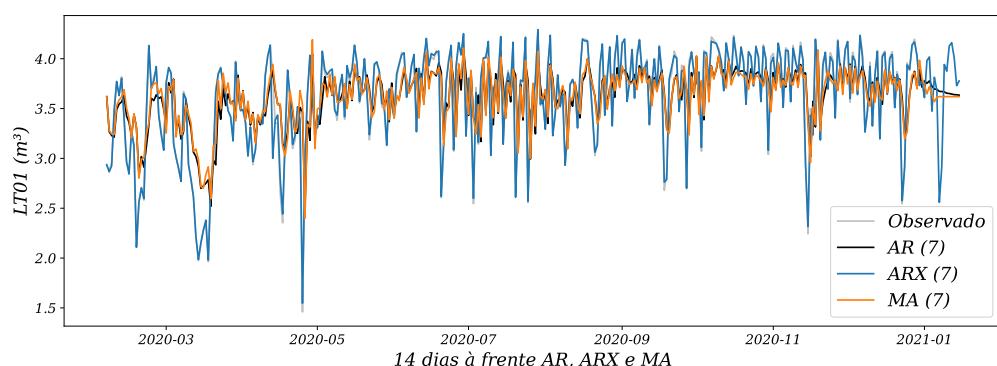
Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Figura 38: Comparação dos modelos AR, ARX e MA, 7 dias à frente



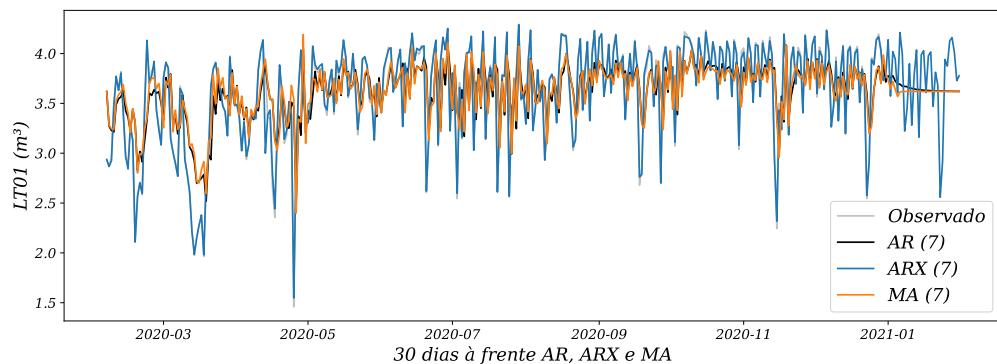
Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Figura 39: Comparação dos modelos AR, ARX e MA, 14 dias à frente



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

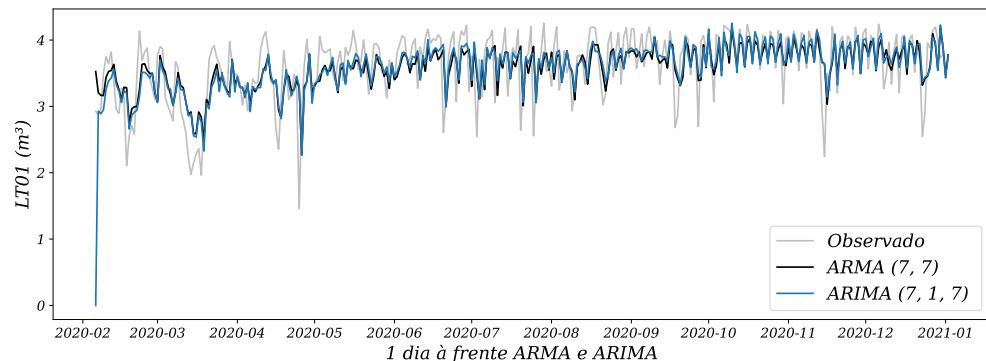
Figura 40: Comparação dos modelos AR, ARX e MA, 30 dias à frente



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

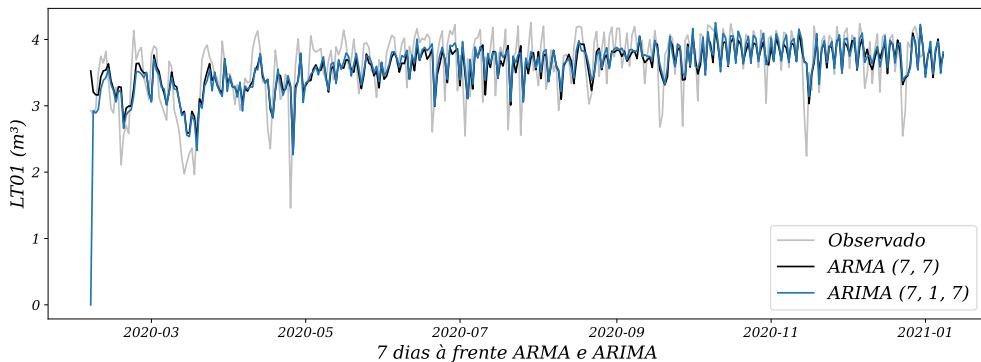
D Apêndice - Modelos ARMA(7,7) e ARIMA (7,1,7) 24h

Figura 41: Comparação dos modelos ARMA e ARIMA, 1 dia à frente



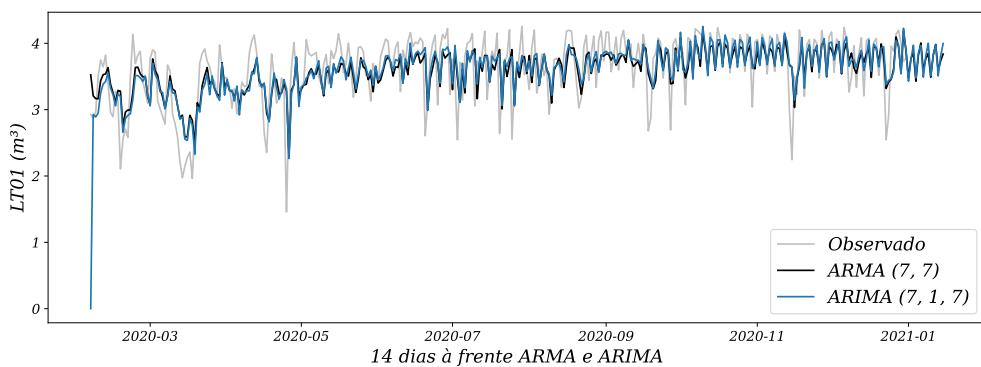
Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Figura 42: Comparação dos modelos ARMA e ARIMA, 7 dias à frente



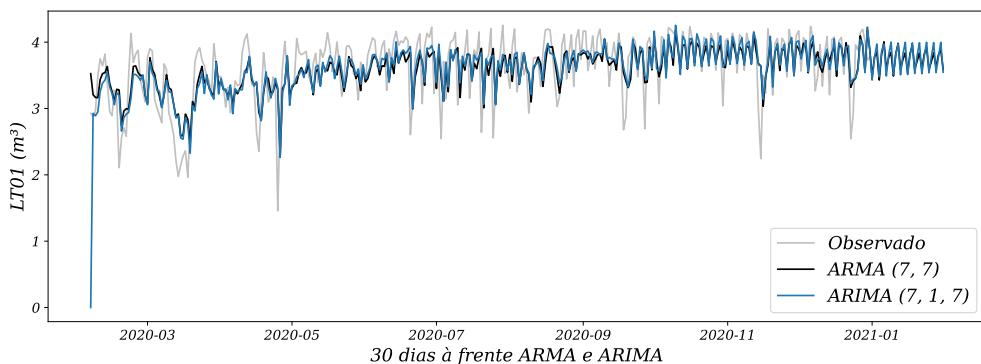
Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Figura 43: Comparação dos modelos ARMA e ARIMA, 14 dias à frente



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

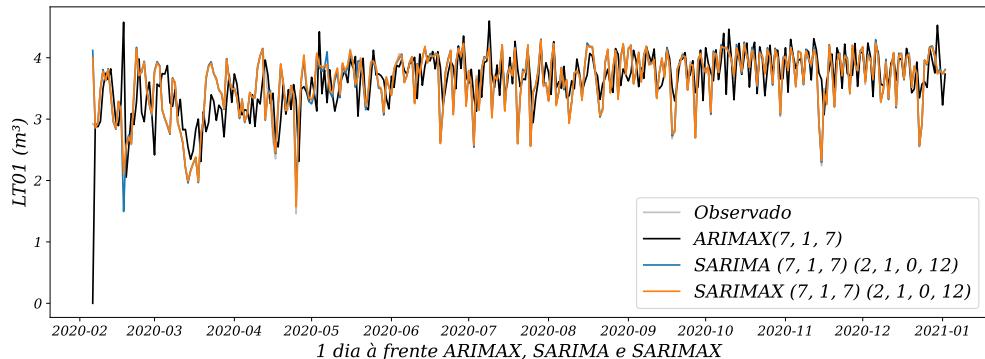
Figura 44: Comparação dos modelos ARMA e ARIMA, 30 dias à frente



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

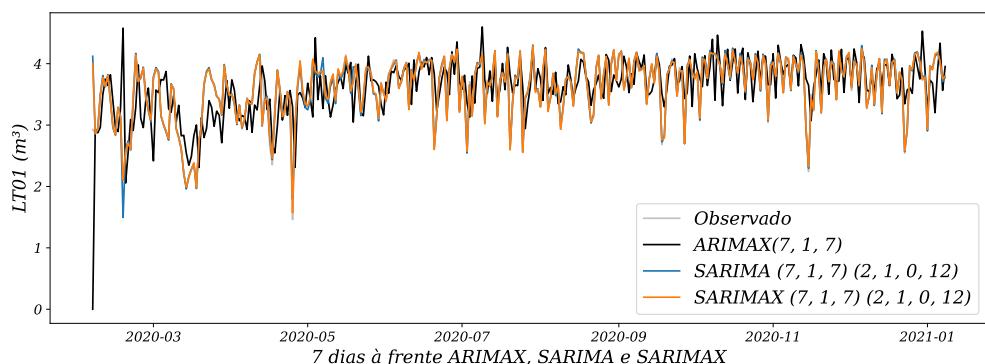
E Apêndice - Modelos ARIMAX (7,1,7), SARIMA (7,1,7) (2,1,0,12) e SARIMAX (7,1,7) (2,1,0,12) 24h

Figura 45: Comparação dos modelos ARIMAX, SARIMA e SARIMAX, 1 dia à frente



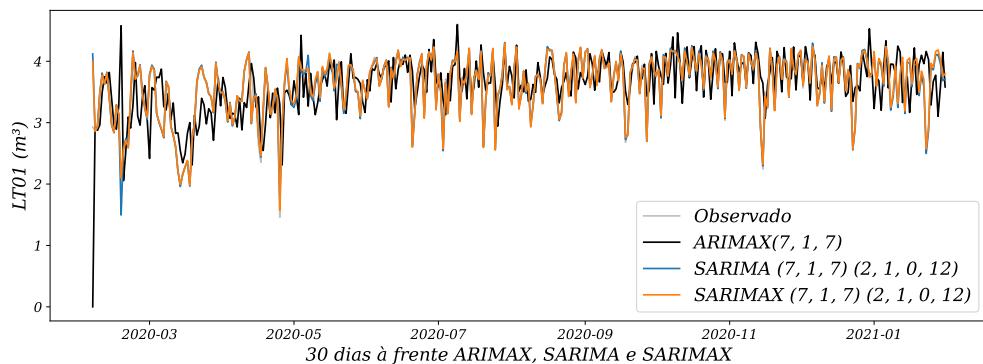
Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Figura 46: Comparação dos modelos ARIMAX, SARIMA e SARIMAX, 7 dias à frente



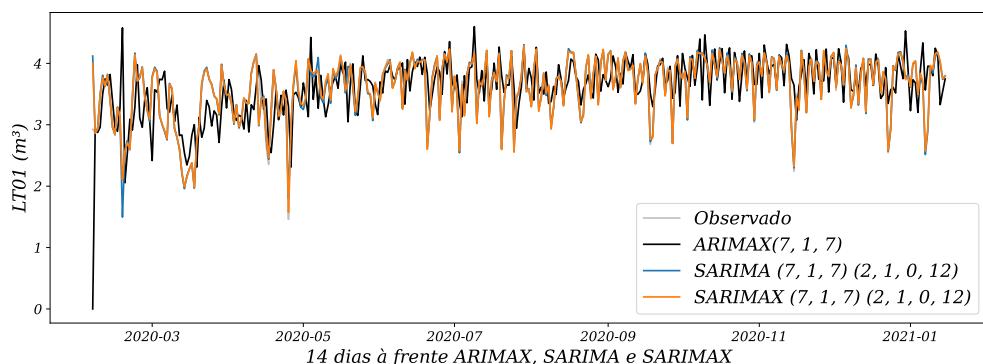
Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Figura 48: Comparação dos modelos ARIMAX, SARIMA e SARIMAX, 30 dias à frente



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)

Figura 47: Comparação dos modelos ARIMAX, SARIMA e SARIMAX, 14 dias à frente



Fonte: Elaboração própria a partir de dados da SANEPAR (2018 a 2020)