

Luan Ferreira de Souza

**Modelo *Fuzzy* Evolutivo Interpretável para  
Predição de Séries Temporais no Mercado  
Financeiro**

Brasil

2024

Luan Ferreira de Souza

## **Modelo *Fuzzy* Evolutivo Interpretável para Predição de Séries Temporais no Mercado Financeiro**

Projeto de Pesquisa desenvolvido durante a disciplina de Trabalho de Conclusão de Curso II e apresentado à banca avaliadora do Curso de Engenharia de Controle e Automação da Escola Superior de Tecnologia da Universidade do Estado do Amazonas, como pré-requisito para a obtenção do título de Engenheiro de Controle e Automação.

Universidade do Estado do Amazonas – UEA

Escola Superior de Tecnologia – EST

Engenharia de Controle e Automação

Orientador: Prof. Dr. Luiz Alberto Queiroz Cordovil Júnior

Coorientador: Prof. Dr. Rodrigo Farias Araújo

Brasil


2024

Luan Ferreira de Souza

## Modelo *Fuzzy* Evolutivo Interpretável para Predição de Séries Temporais no Mercado Financeiro

Projeto de Pesquisa desenvolvido durante a disciplina de Trabalho de Conclusão de Curso II e apresentado à banca avaliadora do Curso de Engenharia de Controle e Automação da Escola Superior de Tecnologia da Universidade do Estado do Amazonas, como pré-requisito para a obtenção do título de Engenheiro de Controle e Automação.

Assinado por:



3B34FC4E04B249D...

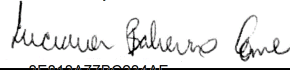
**Prof. Dr. Luiz Alberto Queiroz  
Cordovil Júnior**

Universidade do Estado do Amazonas (UEA)  
- Orientador



**Prof. Dr. Almir Kimura Júnior**  
Universidade do Estado do Amazonas (UEA)

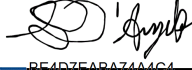
Assinado por:



3E013A77B0234AE...

**Profa. Dra. Luciana Balieiro Cosme**  
Instituto Federal de Educação, Ciência e  
Tecnologia do Norte de Minas Gerais  
(IFNMG)

Signed by:



BE4D7EAD74A4C4...

**Prof. Dr. Marcos Flávio Silveira  
Vasconcelos D'Angelo**  
Universidade Estadual de Montes Claros  
(UNIMONTES)

Brasil  
2024

*À minha mãe, cujo apoio inabalável e encorajamento foram fundamentais para a conclusão deste trabalho. Quando pensei em desistir, foi a sua força e a sua crença em mim que me mantiveram firme. Sou eternamente grato.*

# Agradecimentos

Agradeço, primeiramente, à minha mãe, que me deu apoio e a motivação para perdurar e concluir esta jornada, mesmo quando eu já não sentia que conseguiria.

Ao meu Coorientador, Prof. Dr. Rodrigo Farias Araújo por ter sido o professor que aceitou me orientar inicialmente com um tema inusitado para o curso, além de ter sido crucial nos ensinamentos de suas matérias ministradas.

Ao meu Orientador, Prof. Dr. Luiz Alberto Queiroz Cordovil Júnior por ter me aceito como seu primeiro orientando em sua posição como docente, e por todo o tempo disponibilizado para me guiar com recomendações de artigos, fornecendo feedbacks valiosos durante a construção do trabalho e pela paciência com todas as minhas dúvidas.

Aos meus amigos, que estavam perto ou longe, mas que me acompanharam durante este período turbulento e tanto me ajudaram a manter a sanidade em cheque.

Por fim, agradeço a todos os que estiveram presentes em algum momento desta jornada, seja me ajudando diretamente com o trabalho ou indiretamente.

*“Um dos grandes segredos da sabedoria econômica é saber aquilo que não se sabe.  
(Galbraith, John)”*

# Resumo

Sistemas financeiros contêm muitas variáveis, estas sendo determinísticas e heurísticas, assim tornando-os matematicamente complexos de representar. Desenvolveu-se uma estrutura capaz de lidar com as incertezas presentes no mercado financeiro, fazendo uso da lógica *fuzzy* e algoritmo de agrupamento de dados, com foco na semântica dos dados. Em vista disso, este projeto teve como objetivo a construção de um sistema capaz de efetuar previsões assertivas quanto ao valor futuro de ativos financeiros, de forma a trazer credibilidade para as previsões utilizando técnicas de explicabilidade. Para alcançar tal objetivo, escolheu-se utilizar um modelo *fuzzy* evolutivo para a representação não-linear dos dados e adaptou-se o algoritmo de clusterização, *Online Elliptical Clustering* (OEC), para o agrupamento de dados e detecção de pontos de mudança. O modelo *fuzzy* evolutivo faz uso da combinação de vários sistemas lineares em faixas distintas, também chamados de modelos locais, em que tais faixas são selecionadas através da base de conhecimento. A adaptação do algoritmo OEC fornece uma abordagem evolutiva quanto a representação semântica dos dados, utilizando hiper-elipsóides para o agrupamento dos dados e detecção de pontos de mudança para atualização dos antecedentes. Optou-se pela aplicação de modelos autoregressivos, para representação de cada modelo local. Sua implementação foi realizada em *Python*, aproveitando-se de bibliotecas existentes *open-source* para o tratamento dos dados, construção dos *clusters* e representação dos modelos. Por fim, o modelo proposto obteve resultados adequados quanto às suas previsões, através de métricas de avaliação como o MAE, RMSE, MAPE e sMAPE, conseguindo atribuir suas devidas explicações quanto a previsão efetuada e as nuances do comportamento da série.

**Palavras-chave:** Finanças Quantitativas, Mercado Financeiro, Modelos de Predição Interpretáveis, Predição de Séries Temporais, Modelagem *fuzzy*, *Online Clustering*, *XAI*.

# Abstract

Financial systems encompass numerous variables, both deterministic and heuristic, rendering them intricately complex to model mathematically. A framework was developed to address uncertainties inherent in financial markets, employing fuzzy logic and a data clustering algorithm with a focus on semantic data interpretation. Accordingly, this project aimed to construct a system capable of making precise predictions regarding the future values of financial assets, thereby enhancing prediction credibility through explainability techniques. To achieve this objective, an evolving fuzzy model was employed for non-linear data representation, alongside the adaptation of the Online Elliptical Clustering (OEC) algorithm for data clustering and change point detection. The evolving fuzzy model integrates multiple linear systems within distinct ranges, termed local models, which are selected through knowledge-based methods. The adaptation of the OEC algorithm facilitates an evolving approach to semantic data representation, utilizing hyper-ellipsoids for data clustering and change point detection to update antecedents. Autoregressive models were employed to represent each local model. Implementation was conducted in Python, utilizing existing open-source libraries for data handling, cluster construction, and model representation. Ultimately, the proposed model yielded suitable predictive outcomes, assessed through metrics including MAE, RMSE, MAPE, and sMAPE, effectively providing elucidations for predictions and nuances in the series behavior.

**Keywords:** Quantitative Finance, Financial Market, Interpretable Prediction Models, Time-series forecasting, Fuzzy modelling, Online Clustering, XAI.

# Lista de ilustrações

Figura 1 – Representação esquemática do modelo geral. . . . .	22
Figura 2 – Representação esquemática do modelo ARX. . . . .	23
Figura 3 – Representação esquemática do modelo ARMAX. . . . .	23
Figura 4 – Estrutura de atualização EEFIG. . . . .	30
Figura 5 – <i>Cluster</i> Inicial com sua zona de guarda e o Rastreador de Estado no início da mudança de estado do sistema. . . . .	34
Figura 6 – Relação da Interpretabilidade e do Desempenho preditivo de modelos de IA. . . . .	36
Figura 7 – Fluxograma de desenvolvimento do projeto de pesquisa. . . . .	38
Figura 8 – Fluxograma de desenvolvimento do OEC. . . . .	47
Figura 9 – Divisão de modelos locais por regras. . . . .	49
Figura 10 – Janela de Observação móvel baseado no horizonte de memória. . . . .	51
Figura 11 – Série do preço de fechamento do índice IBOV. . . . .	52
Figura 12 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 1. . . . .	53
Figura 13 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 2. . . . .	54
Figura 14 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 3. . . . .	55
Figura 15 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 3. . . . .	55
Figura 16 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 4. . . . .	56
Figura 17 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 4. . . . .	57
Figura 18 – Simulação I - Gráfico da série original e a predição para a região pertencente a Regra 1. . . . .	58
Figura 19 – Simulação I - Gráfico da série original e a predição para a região pertencente a Regra 2. . . . .	58
Figura 20 – Simulação I - Gráfico da série original e a predição para a região pertencente a Regra 3. . . . .	59
Figura 21 – Simulação I - Gráfico da série original e a predição para a região pertencente a Regra 4. . . . .	60
Figura 22 – Série Original destacando os setores locais identificados por cada regra, referente ao experimento com 2 variáveis. . . . .	61

Figura 23 – Métricas de Avaliação referentes ao experimento <i>out-of-sample</i> com 2 variáveis, resultando em 4 modelos locais, dispostas em <i>boxplots</i> . . . . .	63
Figura 24 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 1. . . . .	65
Figura 25 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 2. . . . .	66
Figura 26 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 2. . . . .	66
Figura 27 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 3. . . . .	67
Figura 28 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 3. . . . .	68
Figura 29 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 4. . . . .	69
Figura 30 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 4. . . . .	69
Figura 31 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 5. . . . .	70
Figura 32 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 5. . . . .	71
Figura 33 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 6. . . . .	71
Figura 34 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 6. . . . .	72
Figura 35 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 1. . . . .	73
Figura 36 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 2. . . . .	74
Figura 37 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 3. . . . .	74
Figura 38 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 4. . . . .	75
Figura 39 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 5. . . . .	76
Figura 40 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 6. . . . .	76
Figura 41 – Série Original destacando os setores locais identificados por cada regra, referente ao experimento com 4 variáveis. . . . .	77

Figura 42 – Métricas de Avaliação referentes ao experimento *out-of-sample* com 4 variáveis, resultando em 6 modelos locais, dispostas em *boxplots*. . . . 79

# Lista de Algoritmos

Algoritmo 1 – Classe Regra - Função Rule e InitRule . . . . .	41
Algoritmo 2 – Função EvaluateIDCAD . . . . .	42
Algoritmo 3 – Classe Regra - Função Step . . . . .	43
Algoritmo 4 – Classe Regra - Função WeightedStep . . . . .	44
Algoritmo 5 – Função StateTrackerIDCAD . . . . .	45
Algoritmo 6 – Função CSeparation . . . . .	46

# Lista de tabelas

Tabela 1 – Métricas de Avaliação de Performance do Modelo. . . . .	27
Tabela 2 – Variáveis utilizadas para Simulação I. . . . .	52
Tabela 3 – Métricas de Avaliação utilizadas na Simulação I. . . . .	61
Tabela 4 – Importância Normalizada de cada atributo, por modelo, destacando os com maior pertinência para a Simulação I. . . . .	62
Tabela 5 – Variáveis utilizadas para Simulação II. . . . .	64
Tabela 6 – Métricas de Avaliação utilizadas na Simulação II. . . . .	77
Tabela 7 – Importância Normalizada de cada atributo, por modelo, destacando os com maior pertinência para a Simulação II. . . . .	78

# Lista de abreviaturas e siglas

EEFIG	Evolving Ellipsoidal Fuzzy Information Granules
T-S	Takagi-Sugeno
ARX	<i>Autoregressive with exogenous inputs</i>
ARMAX	<i>Autogressive moving average with exogenous inputs</i>
MA	Média móvel
EMQ	Estimador de Mínimos Quadrados
ERMQ	Estimador Recursivo de Mínimos Quadrados
PJG	<i>Principle of Justifiable Granularity</i>
IG	<i>Information Granules</i>
OEC	<i>Online Elliptical Clustering</i>
wIDCAD	<i>Weighted Incremental Data Capture Anomaly Detection</i>
MAE	<i>Mean Absolute Error</i>
MDAE	<i>Median Absolute Error</i>
MSE	<i>Mean Squared Error</i>
RMSE	<i>Root Mean Squared Error</i>
MAPE	<i>Mean Absolute Prediction Error</i>
sMAPE	<i>Symmetric Mean Absolute Prediction Error</i>
IA	Inteligência Artificial
XAI	<i>Explainable Artificial Intelligence</i>
RNP	Redes Neurais Profundas
IBOV	Índice Bovespa
SELIC	Sistema Especial de Liquidação e Custódia
IPCA	Índice Nacional de Preços ao Consumidor Amplo

ADF	<i>Dickey-Fuller Aumentado</i>
API	<i>Application Programming Interface</i>
IDE	<i>Integrated Development Environment</i>
SGS	Sistema Gerenciador de Séries Temporais
FAC	Função de Autocorrelação
FACP	Função de Autocorrelação Parcial

# Sumário

<b>1</b>	<b>INTRODUÇÃO</b>	<b>17</b>
<b>1.1</b>	<b>Objetivos</b>	<b>18</b>
1.1.1	Objetivo geral	18
1.1.2	Objetivos específicos	18
<b>1.2</b>	<b>Contribuições</b>	<b>19</b>
<b>1.3</b>	<b>Estrutura do Documento</b>	<b>19</b>
<b>2</b>	<b>REFERENCIAL TEÓRICO</b>	<b>21</b>
<b>2.1</b>	<b>Identificação de Sistemas Dinâmicos</b>	<b>21</b>
2.1.1	Representação em Tempo Discreto	21
2.1.2	Modelo ARX	22
2.1.3	Modelo ARMAX	22
2.1.4	Estimador de Mínimos Quadrados	23
2.1.5	Estimador Recursivo de Mínimos Quadrados	24
<b>2.2</b>	<b>Métricas de Avaliação</b>	<b>25</b>
<b>2.3</b>	<b>Modelagem fuzzy</b>	<b>27</b>
2.3.1	Modelo Takagi-Sugeno (T-S)	27
2.3.2	Modelo T-S baseado em grânulos evolutivos elipsoidais <i>fuzzy</i> (EEFIG)	29
<b>2.4</b>	<b><i>Online Elliptical Clustering</i></b>	<b>31</b>
2.4.1	Parametrização dos <i>clusters</i> elipsoidais	32
2.4.2	Condição de atualização (wIDCAD)	32
2.4.3	Novos Clusters baseados em Rastreador de Estado do Sistema	33
<b>2.5</b>	<b><i>Explainable Artificial Intelligence (XAI)</i></b>	<b>35</b>
<b>3</b>	<b>METODOLOGIA</b>	<b>38</b>
<b>3.1</b>	<b>Coleta de dados</b>	<b>39</b>
<b>3.2</b>	<b>Análise de dados</b>	<b>39</b>
<b>3.3</b>	<b>Regras <i>Fuzzy</i> via Cluster</b>	<b>40</b>
3.3.1	Parametrização Inicial das regras	40
3.3.2	Atualização Recursiva das regras	42
3.3.3	Rastreador de Estados	45
3.3.4	Criação de Nova Regra	46
<b>3.4</b>	<b>Modelos ARX e ARMAX</b>	<b>48</b>
<b>3.5</b>	<b>Interpretabilidade e Explicabilidade do Modelo</b>	<b>49</b>
<b>4</b>	<b>EXPERIMENTOS COMPUTACIONAIS</b>	<b>50</b>

<b>4.1</b>	<b>Simulação I</b> . . . . .	<b>52</b>
4.1.1	Caracterização dos Dados . . . . .	52
4.1.2	Resultados . . . . .	57
<b>4.2</b>	<b>Simulação II</b> . . . . .	<b>64</b>
4.2.1	Caracterização dos Dados . . . . .	64
4.2.2	Resultados . . . . .	72
<b>5</b>	<b>CONSIDERAÇÕES FINAIS</b> . . . . .	<b>80</b>
<b>5.1</b>	<b>Trabalhos Futuros</b> . . . . .	<b>81</b>
	<b>REFERÊNCIAS</b> . . . . .	<b>82</b>

# 1 Introdução

A teoria dos mercados eficientes trata sobre a perspectiva de que existe uma racionalidade por trás das mudanças e do dinamismo do mercado (JONES; NETTER, 2008). Pode-se citar, por exemplo, o mercado de ações brasileiro, em que os agentes tentam explicar e especular valores futuros de ativos através de análises fundamentalistas baseadas em indicadores financeiros. Uma abordagem mais racional e objetiva, retirando o viés humano, é a análise quantitativa através de modelos baseados em dados (KADIMA GESTÃO DE INVESTIMENTOS LTDA, 2023).

Ativos presentes em mercados financeiros são um conjunto de dados que ocorrem em sequência no tempo, assim sendo um sistema dinâmico, habitualmente chamado de série temporal. Comumente, a identificação da dinâmica de sistemas é feita utilizando métodos estocásticos, tais como os descritos em Aguirre (2014). Estes consideram a evolução da dinâmica de um sistema de forma probabilística, para estimar os parâmetros de um sistema a partir de uma base de dados, incluindo na dinâmica ruídos e distúrbios aproximando-se do sistema real.

De acordo com Taylor e Karlin (1994), modelos estocásticos tratam-se de aproximações e logo contêm erros de estimação e previsão, causando a perda de dinâmica do modelo. Para que possa haver uma compensação pelos erros acarretados do modelo, técnicas de modelagem não-linear são utilizadas.

A lógica *fuzzy* tornou-se uma abordagem muito comum, para sistemas não-lineares, desde sua primeira menção nos anos 60 (ZADEH, 1965), permitindo lidar com imprecisões e incertezas em dados e informações, mantendo ainda sua interpretabilidade e explicabilidade quanto a saída do modelo. Ao contrário de algoritmos de aprendizado de máquina supervisionados, que são muito utilizados por sua versatilidade de uso, os mesmos não possuem uma boa explicabilidade quanto aos valores previstos, e conforme a complexidade dos dados e/ou modelo aumenta, mais poder computacional era necessário (GERON, 2017; LUGER, 2013; DU; LIU; HU, 2019).

Sistemas baseados em lógica *fuzzy* apresentam uma explicabilidade inerente devido a utilização de regras que conduzem as decisões tomadas. Por outro lado, quando não há conhecimento prévio sobre os dados estudados, é necessário utilizar técnicas capazes de extrair informações e, a partir disso, criar as regras. O algoritmo *Online Elliptical Clustering*, apresentado por Moshtaghi, Leckie e Bezdek (2016), demonstra uma capacidade evolutiva para agrupamentos de dados, de forma que os agrupamentos podem ser utilizados como regras em um sistema *fuzzy* para captura de dinâmicas não-lineares.

Nesse contexto, a teoria de controle moderna *fuzzy* foi aplicada na área de finanças

por Hachicha, Jarboui e Siarry (2011), propondo um controlador *fuzzy* otimizado por um algoritmo evolutivo com intuito de modelar a dinâmica de mercados financeiros. Brož e Dostál (2013) propõem um modelo *fuzzy* que oferece suporte à tomada de decisão para investimentos a longo prazo, avaliando o status atual de mercados financeiros a partir dos dados históricos.

Para manter a interpretabilidade e explicabilidade de uma predição de valor futuro em uma série temporal, a combinação da lógica *fuzzy* e modelos estocásticos se torna um método que é eficiente na predição de séries temporais ao não demandar uma complexidade computacional alta ao mesmo tempo que provê a explicabilidade necessária para manter sua coerência semântica.

Portanto, parte-se da hipótese de que é possível utilizar técnicas modernas de controle em combinação com métodos de identificação de sistemas dinâmicos para criar um modelo *fuzzy* capaz de prever o valor futuro de ativos financeiros corretamente, além de ter a capacidade de ser interpretável.

## 1.1 Objetivos

### 1.1.1 Objetivo geral

Desenvolver um modelo *fuzzy* evolutivo baseado em dados, que incorpore técnicas de explicabilidade de modelos, com o propósito de modelar e prever valores futuros do Índice Bovespa B3.

### 1.1.2 Objetivos específicos

1. Validar as predições dos modelos locais *out-of-sample* utilizando as métricas de desempenho baseadas em estimativa de erro.
2. Verificar a capacidade de identificação do OEC de diferentes regiões de operação em séries temporais, abrangendo adequadamente as dinâmicas em cada uma.
3. Avaliar a explicabilidade dos modelos através das importâncias de cada atributo, inferindo suas representatividades em cada região de operação.
4. Testar possibilidades de otimização do modelo a partir de variáveis promissoras no quesito explicabilidade.

## 1.2 Contribuições

A principal contribuição deste trabalho abrange a utilização de técnicas de clusteração de dados temporais para identificação de diferentes dinâmicas, com o intuito de aplicação em um modelo de predição *fuzzy*. Assim, evidenciando um método para captar variabilidades distintas do mercado de ações brasileiro, além de proporcionar uma visão mais explicativa quanto às predições realizadas.

De forma estruturada, as seguintes adaptações foram feitas:

1. O algoritmo se baseia na identificação de anomalias para atualização e criação de novos *clusters*. Nenhum dado pode ser desconsiderado em sistemas financeiros e têm-se o pressuposto que observações próximas no tempo são relacionadas, então dados considerados anômalos que não pertencem a nenhuma zona de guarda foram alocados para o *cluster* em que se apresentou o maior grau de pertinência apresentado. Assim, limitou-se a perda de informações disponíveis para o modelo.
2. O Rastreador de Estados, utilizado para acompanhar mudanças rápidas dos dados, é baseado em um fator de esquecimento. A partir deste fator, o conceito de horizonte de memória é aplicado com o intuito de limitar a quantidade de dados passados utilizados. Da mesma forma que o rastreador utiliza uma quantidade limitada de dados passados, escolheu-se aplicar o horizonte de memória para o treinamento de cada modelo local. Com esta adaptação foi possível verificar a qualidade semântica de cada região identificada, tendo que uma boa modelagem de regras acarreta em uma boa representatividade semântica dos dados pertencentes a cada regra.

## 1.3 Estrutura do Documento

No primeiro capítulo foi apresentado o contexto inicial do tema deste trabalho, abordando a teoria do mercado eficiente e a motivação para a representação e predição de sistemas financeiros. Este capítulo estabelece a base teórica e justifica a importância do estudo.

O segundo capítulo é dedicado à revisão dos referenciais teóricos relacionados à identificação de sistemas, também conhecida como econometria. Aqui, são exploradas as métricas de acurácia, o conceito de lógica *fuzzy* e sua aplicação em modelos de predição, bem como o campo da inteligência artificial explicável. Este capítulo fornece o embasamento teórico necessário para a compreensão das metodologias aplicadas posteriormente.

No terceiro capítulo são descritas as metodologias utilizadas para a aplicação dos conceitos apresentados no capítulo anterior. Isso inclui o uso de bibliotecas pré-existentes

para a implementação de modelos auto-regressivos e o foco na construção do algoritmo de clusterização. Este capítulo detalha as técnicas e ferramentas empregadas na pesquisa.

O quarto capítulo apresenta as análises e os resultados obtidos a partir das metodologias aplicadas. As discussões são realizadas com base nos dados analisados, permitindo o entendimento da representatividade da dinâmica do mercado de ações brasileiro. Este capítulo finaliza o estudo com a interpretação dos achados e suas implicações.

Por fim, o quinto capítulo apresenta uma síntese das principais descobertas, destacando que a proposta inicial do projeto foi atingida, e explorando as nuances da pesquisa. Também são sugeridas direções para trabalhos futuros, com diferentes objetivos e possíveis aplicações, ampliando o impacto e a aplicabilidade dos resultados obtidos neste estudo.

## 2 Referencial Teórico

Este capítulo tem como objetivo apresentar a fundamentação teórica da pesquisa. Por se tratar de um foco multidisciplinar, escolheu-se abordar primeiramente os principais modelos estocásticos para identificação de sistemas dinâmicos, muito utilizados para predição de séries temporais. Na sequência, são apresentados os conceitos de lógica *fuzzy*, seguido de um modelo de identificação *fuzzy* Takagi-Sugeno, tornando possível incorporar a noção de incerteza no modelo. Na terceira Seção, é apresentado o algoritmo de clusterização OEC a ser adaptado para o modelo. Por último, será destacada a necessidade da explicabilidade e interpretabilidade quanto a modelos de predição de séries temporais.

### 2.1 Identificação de Sistemas Dinâmicos

#### 2.1.1 Representação em Tempo Discreto

De forma a se trabalhar com sistemas dinâmicos, em especial séries temporais, funções de transferência e modelos em espaço de estados não são representações adequadas para o controle desejado dos dados de saída. Assim, representações discretas se tornam as mais adequadas.

Tem-se um modelo de estimação geral dado em [Aguirre \(2014\)](#) como (ver Fig. 1)

$$A(q)y(k) = \frac{B(q)}{F(q)}u(k) + \frac{C(q)}{D(q)}\nu(k) \quad (2.1)$$

sendo  $q^{-1}$  como um operador de atraso, assim  $y(k)q^{-1} = y(k-1)$ ,  $\nu(k)$  ruído branco e  $A(q)$ ,  $B(q)$ ,  $C(q)$ ,  $D(q)$ ,  $F(q)$  representam os polinômios do tipo:

$$A(q) = 1 - a_1q^{-1} - \dots - a_{n_y}q^{-n_y},$$

$$B(q) = b_1q^{-1} + \dots + b_{n_y}q^{-n_y},$$

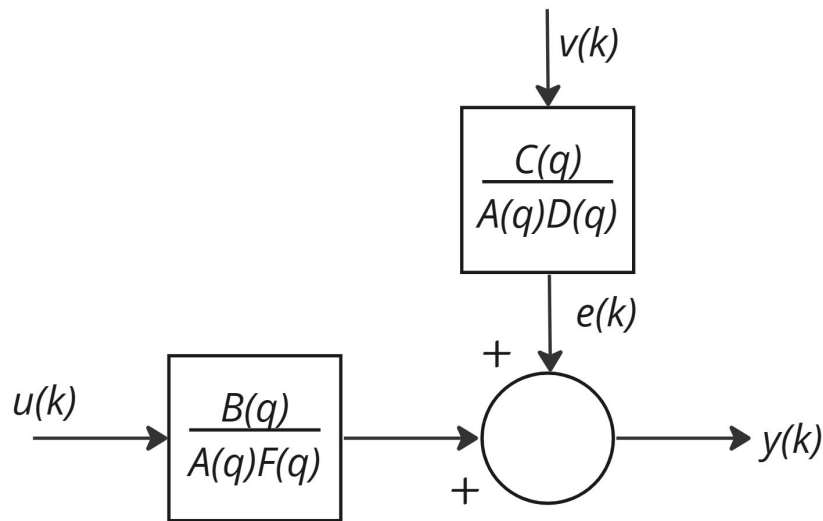
$$C(q) = 1 + c_1q^{-1} + \dots + c_{n_y}q^{-n_y},$$

$$D(q) = 1 + d_1q^{-1} + \dots + d_{n_y}q^{-n_y},$$

$$F(q) = 1 + f_1q^{-1} + \dots + f_{n_y}q^{-n_y}.$$

Especialmente, dois modelos de identificação de sistemas são relevantes para esta pesquisa, em que, o primeiro contém um ruído não branco mas é modelado como um processo branco e no segundo, o erro é modelado como ruído branco.

Figura 1 – Representação esquemática do modelo geral.



Fonte: Aguirre (2014)

### 2.1.2 Modelo ARX

O modelo autoregressivo com entradas exógenas (ARX do inglês *autoregressive with exogenous inputs*) pertence a classe de modelos de *erro na equação*, ou seja, o ruído  $\nu(k)$  afeta diretamente a equação.

A partir do modelo geral (2.1), o ARX é obtido modificando-se  $C(q) = D(q) = F(q) = 1$ , assim

$$y(k) = \frac{B(q)}{A(q)}u(k) + \frac{1}{A(q)}\nu(k)$$

em que coloca em destaque as “funções de transferência” do sistema e de ruído, como pode ser visto na Fig. 2.

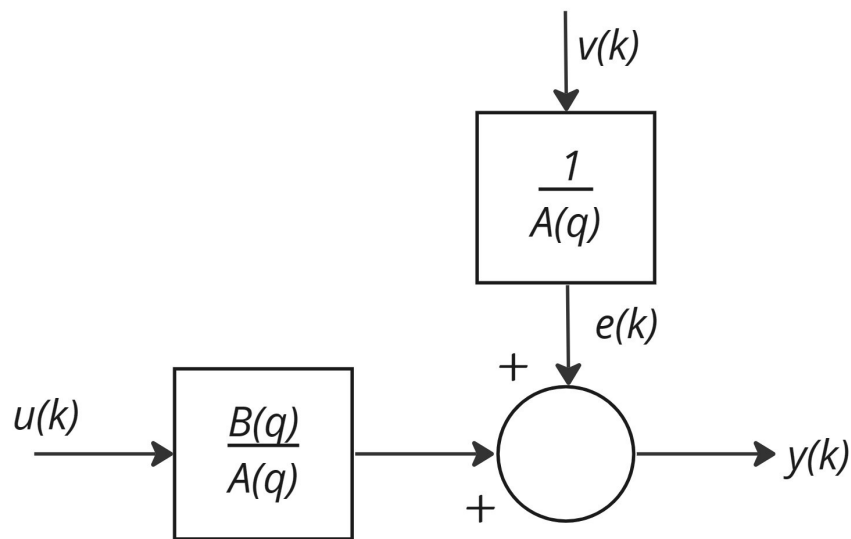
### 2.1.3 Modelo ARMAX

De forma similar ao modelo ARX, o modelo autoregressivo com média móvel e entradas exógenas (ARMAX do inglês *autoregressive moving average with exogenous inputs*) também faz parte da classe de modelos de *erro na equação*. Neste caso o erro na equação é representado como um operador de média móvel (MA), e o polinômio  $C(q)/A(q)$  se torna um filtro de ruídos.

A partir do modelo geral (2.1), o ARMAX é obtido modificando-se  $D(q) = F(q) = 1$ , assim (ver Fig. 3)

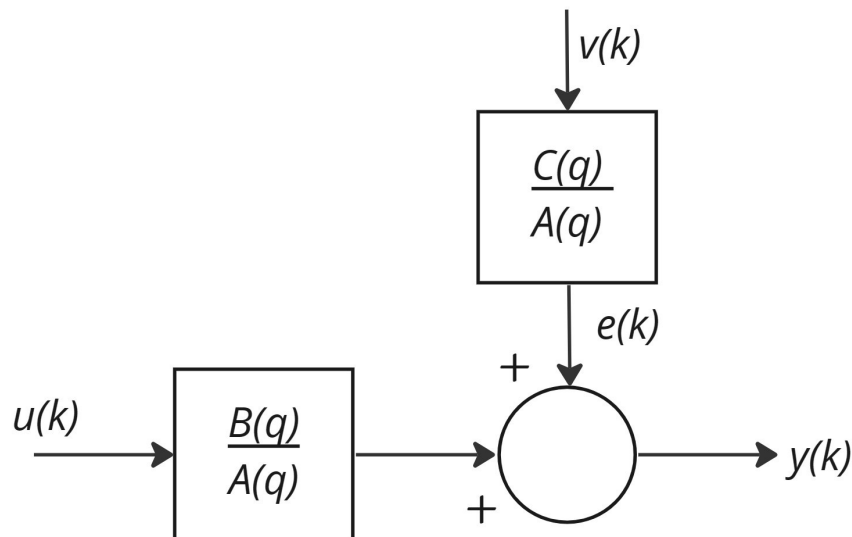
$$y(k) = \frac{B(q)}{A(q)}u(k) + \frac{C(q)}{A(q)}\nu(k)$$

Figura 2 – Representação esquemática do modelo ARX.



Fonte: Aguirre (2014)

Figura 3 – Representação esquemática do modelo ARMAX.



Fonte: Aguirre (2014)

#### 2.1.4 Estimador de Mínimos Quadrados

Tendo modelos com representações discretas desenvolvidas, faz-se necessário estimar os parâmetros que explicam a dinâmica do sistema. O método utilizado, o Estimador de Mínimos Quadrados (EMQ), considera que o valor estimado do vetor de parâmetros é conhecido, e que haja um erro em referência a saída observado com base em um vetor

de regressores medido. Portanto, o EMQ trata de um problema de minimização de uma função custo, esta sendo o erro em referência a saída observada.

- $\hat{\theta}$  - vetor de parâmetros;
- $\xi$  - erro em referência a saída observada;
- $\Psi$  - vetor de regressores;
- $y$  - saída observada.

Ao tratarmos as variáveis de forma matricial, o erro  $\xi$  torna-se o vetor de resíduos  $\xi$ , sendo definido em [Aguirre \(2014\)](#) como:

$$\xi = y - \Psi \hat{\theta}$$

Conseqüentemente, temos o EMQ estabelecido como

$$\hat{\theta}_{MQ} = [\Psi^T \Psi]^{-1} \Psi^T y$$

em que as variáveis regressoras são armazenadas no vetor  $\Psi$ .

Para os modelos *ARX* e *ARMAX*, o vetor de regressores se dá por:

$$\Psi = [y(k-1) \quad \dots \quad y(k-n_y) \quad u(k-1) \quad \dots \quad u(k-n_u)]^T$$

Como solução, o sinal de entrada deve ser aleatório, na verdade pseudo-aleatório. Dadas essas condições, o método dos mínimos quadrados pode identificar um sistema dinâmico e lidar com o ruído presente no sinal de saída por meio de um modelo parametrizado.

A polarização pode ocorrer em modelos *ARX* e *ARMAX*, onde o ruído é colorido, quando os erros de regressão são auto correlacionados e quando o modelo inclui regressores que empacotam saídas. Para resolver este problema, é preciso tornar os sinais de entrada aleatórios, embora, na prática, sejam apenas pseudo-aleatórios. Considerando as circunstâncias apresentadas, o método dos mínimos quadrados apresenta-se como uma técnica eficaz para a identificação de sistemas dinâmicos por meio de um modelo parametrizado, bem como para o tratamento de ruídos que possam estar presentes no sinal de saída.

### 2.1.5 Estimador Recursivo de Mínimos Quadrados

O EMQ considera a disponibilidade total dos dados *a priori*, ou seja, antes da estimação dos parâmetros. Em situações práticas, o caso citado acontece quando os dados de entrada e saída são previamente coletados com intuito *a posteriori*. Em particular, para esta pesquisa é relevante a utilização da estimação recursiva, onde os dados são coletados de forma sequencial.

A recursividade de estimação depende de um período de amostragem para cada instante de coleta, corrigindo a estimativa anterior com os novos dados, definido em [Aguirre \(2014\)](#) por

$$\hat{\theta}_k = \hat{\theta}_{k-1} - \eta_k \mathbf{d}_k, \quad (2.3)$$

onde o termo na extremidade direita representa a correção da estimativa anterior em que  $\eta_k$  representa o grau da correção e  $\mathbf{d}_k$  determina como a correção é distribuída entre cada componente do vetor de regressores.

O estimador de mínimos quadrados recursivo (ERMQ) estabelece dois novos parâmetros para introduzir a forma sequencial de estimação. O primeiro se trata de uma matriz de ganho  $K_k$ , de forma análoga ao grau de correção em (2.3), em que para encontrar  $K_k$  impõe-se uma minimização da covariância do vetor de parâmetros estimados no instante  $k$ ,  $\hat{\theta}_k$ ; o segundo parâmetro é  $P_k$ , chamado de matriz de covariância na  $k$ -ésima iteração.

De forma geral, o algoritmo ERMQ é calculado na forma e ordem

$$\begin{cases} K_k = P_{k-1} \psi_k [\psi_k^T P_{k-1} \psi_k + 1]^{-1}; \\ \hat{\theta}_k = \hat{\theta}_{k-1} + K_k [y(k) - \psi_k^T \hat{\theta}_{k-1}]; \\ P_k = P_{k-1} - K_k \psi_k^T P_{k-1}. \end{cases}$$

em que  $\psi_k = \psi(k-1)$ ,  $P_k$  multiplica o componente de *inovação* ( $y(k) - \psi_k^T \hat{\theta}_{k-1}$ ), e o vetor de parâmetros estimados  $\hat{\theta}$  é definido a cada iteração, à medida que novos dados são incorporados.

[Aguirre \(2014\)](#) traz uma atenção especial a matriz de covariância, em que os elementos da diagonal principal são as variâncias dos respectivos elementos no vetor de parâmetros, ou seja, estas indicam um nível de confiança em relação a estimação dos elementos de  $\hat{\theta}_k$ . Porém, se os valores da matriz de covariância se tornarem muito pequenos e mais próximo de ser singular, pior condicionado será o sistema ([ZHU; STEC, 2006](#)).

## 2.2 Métricas de Avaliação

Para avaliação de modelos de identificação de sistemas dinâmicos, diversas métricas são utilizadas de acordo com a literatura. [Steurer, Hill e Pfeifer \(2021\)](#) demonstram que as métricas podem ser caracterizadas em 7 classes, das quais 3 são relevantes: Diferença Absoluta, Diferença quadrática, Proporção Absoluta.

Com o foco de predição de valores de ativos financeiros, os valores reais serão chamados de  $p_n$  e as predições de  $\hat{p}_n$ , onde  $n = 1, \dots, N$  são os índices dos dados. Para facilitar a interpretabilidade e comparação, as métricas estão dispostas de forma que quanto menor os valores absolutos, melhor o desempenho.

## Métricas de Diferença Absoluta

Com o foco na dispersão média do erro, métricas de diferença absoluta limitam o impacto de *outliers* individuais na performance do modelo, sendo muito utilizados em situações onde há erro dos dados de entrada ou problemas com a qualidade dos dados.

- Erro Médio Absoluto (MAE do inglês *Mean Absolute Error*)

$$MAE = \frac{1}{N} \sum_{n=1}^N |p_n - \hat{p}_n| \quad (2.4)$$

- Erro Mediano Absoluto (MDAE do inglês *Median Absolute Error*)

$$MDAE = med|p_n - \hat{p}_n| \quad (2.5)$$

## Métricas de Diferença Quadrática

Em contraste a média de erros absolutos, diferenças quadráticas são mais sensíveis a *outliers*, sendo úteis em que previsões com erros muito discrepantes tenham que ser minimizados.

- Erro Quadrático Médio (MSE do inglês *Mean Squared Error*)

$$MSE = \frac{1}{N} \sum_{n=1}^N (p_n - \hat{p}_n)^2 \quad (2.6)$$

- Raiz do Erro Quadrático Médio (RMSE do inglês *Root Mean Squared Error*)

$$RMSE = \sqrt{\frac{\sum_{n=1}^N (p_n - \hat{p}_n)^2}{N}} \quad (2.7)$$

## Métricas de Proporção Absoluta

Métricas de proporção têm como foco destacar a grandeza do erro em relação ao dados tratados, levando em consideração sua escala.

- Erro Médio Absoluto de Predição (MAPE do inglês *Mean Absolute Prediction Error*)

$$MAPE = \frac{1}{N} \sum_{n=1}^N \left| \left( \frac{p_n}{\hat{p}_n} \right) - 1 \right| \quad (2.8)$$

Este não é simétrico em dispersão, ou seja, quando  $\hat{p} \leq p$ ,  $p/\hat{p}$  não tem limite, podendo variar entre 1 e infinito. Por outro lado, quando  $\hat{p} \geq p$ ,  $p/\hat{p}$  fica limitado entre 0 e 1. Devido a essa assimetria, o MAPE simétrico surge.

- Erro Médio Absoluto de Predição Simétrico (sMAPE do inglês *Symmetric Mean Absolute Prediction Error*)

$$sMAPE = \frac{1}{N} \sum_{n=1}^N \left( \frac{|p_n - \hat{p}_n|}{p_n + \hat{p}_n} \right) \quad (2.9)$$

As métricas utilizadas neste trabalho dispostas anteriormente podem ser agrupadas como na Tabela 1.

Sigla	Nome	Fórmula
MAE	Erro Médio Absoluto	$\frac{1}{N} \sum_{n=1}^N  p_n - \hat{p}_n $
RMSE	Raiz do Erro Quadrático Médio	$\sqrt{\frac{\sum_{n=1}^N (p_n - \hat{p}_n)^2}{N}}$
MAPE	Erro Médio Absoluto de Predição	$\frac{1}{N} \sum_{n=1}^N \left  \frac{p_n}{\hat{p}_n} - 1 \right $
sMAPE	Erro Médio Absoluto de Predição Simétrico	$\frac{1}{N} \sum_{n=1}^N \left( \frac{ p_n - \hat{p}_n }{p_n + \hat{p}_n} \right)$

Tabela 1 – Métricas de Avaliação de Performance do Modelo.

## 2.3 Modelagem *fuzzy*

### 2.3.1 Modelo Takagi-Sugeno (T-S)

Um dos principais ganhos da utilização de técnicas baseadas na lógica *fuzzy*, é que esta se trata de uma extensão da lógica booleana, em que só existem duas condições: 0 ou 1, verdadeiro ou falso, pertence ou não pertence. A lógica *fuzzy* considera que existam várias condições intermediárias entre os estados de pertencer e não pertencer, atribuindo variáveis linguísticas como “mais”, “menos” e “mais ou menos” (ZADEH, 1994). Assim, é possível que uma variável possa coabitar dois setores lineares distintos, em que cada variável linguística é representada por uma função de pertinência, atribuindo pesos a essa variável referente a cada região. Assim, pode-se representar de forma mais assertiva a condição de um sistema, que possui não-linearidades, em determinado instante.

A técnica de modelagem *fuzzy* Takagi-Sugeno (T-S) consiste em uma abordagem para modelar sistemas não-lineares, utilizando combinações de sistemas lineares em faixas distintas. A modelagem é realizada por meio da implementação de regras que indicam as relações de entrada e saída lineares em determinadas regiões de operação de um sistema dinâmico não-linear (TANAKA; WANG, 2001). O modelo completo do sistema é obtido pela combinação dos modelos locais criados pelas regras *SE-ENTÃO fuzzy*, aferidos seus pesos através das funções de pertinência. A principal diferença entre a modelagem *fuzzy*

Takagi-Sugeno e outras técnicas de modelagem *fuzzy* é que as regras *SE-ENTÃO* geram equações como resultado (ANGELOV; FILEV, 2004; PAN; WANG; HUANG, 2019).

De acordo com Tanaka e Wang (2001), a  $i$ -ésima regra do modelo *fuzzy* T-S é definida, para sistemas discretos com variáveis premissas  $z_1(t) \dots z_p(t)$ , como:

**Regra  $i$ :**

**SE**  $z_1(t)$  é  $M_{i1}$  e ... e  $z_p(t)$  é  $M_{ip}$ , **ENTÃO**

$$\begin{cases} x(t+1) = A_i x(t) + B_i u(t), \\ y(t) = C_i x(t). \end{cases}$$

Em que:

- $M_{ip}$  é o conjunto *fuzzy* da  $p$ -ésima variável premissa na  $i$ -ésima regra;
- $i = 1, 2, \dots, r$  denota a  $i$ -ésima regra do modelo;
- $x(t) \in \mathbb{R}^n$  é o vetor de estados;
- $u(t) \in \mathbb{R}^m$ ;
- $y(t) \in \mathbb{R}^q$  o vetor de saída;
- $z_1(t) \dots z_p(t)$  são as variáveis premissa, onde o vetor discreto de premissas é  $z(t)$ .

$A_i \in \mathbb{R}^{n \times n}$ ,  $B_i \in \mathbb{R}^{n \times m}$  e  $C_i \in \mathbb{R}^{q \times n}$  são as matrizes do  $i$ -ésimo modelo linear local. Para tornar o processo de *defuzzificação* mais factível é tomada como verdadeira a premissa de que as variáveis  $z_1(t) \dots z_p(t)$  não são funções de  $u(t)$ .

Logo, dado um conjunto  $(x(t), u(t))$  de dados de variáveis de estado e entrada, a saída do modelo nebuloso é dada para sistemas discretos como:

$$\mathbf{x}(\mathbf{t} + \mathbf{1}) = \sum_{i=1}^r h_i(z(t)) \{A_i x(t) + B_i u(t)\} \quad (2.10)$$

$$\mathbf{y}(\mathbf{t}) = \sum_{i=1}^r h_i(z(t)) C_i x(t) \quad (2.11)$$

Em que:

$$\begin{cases} z(t) = [z_1(t) \quad z_2(t) \quad \dots \quad z_p(t)], \\ w_i(z(t)) = \prod_{j=1}^p M_{ij}(z_j(t)), \\ h_i(z(t)) = \frac{w_i(z(t))}{\sum_{i=1}^r w_i(z(t))}, \quad \forall t. \end{cases}$$

O termo  $M_{ij}(z(t))$  é o nível de pertinência de  $z_j(t)$  em  $M_{ij}$ . Dado que para  $i = 1, 2, \dots, r$

$$\begin{cases} \sum_{i=1}^r w_i(z(t)) > 0, \\ w_i(z(t)) \geq 0. \end{cases}$$

temos

$$\begin{cases} \sum_{i=1}^r z_i(z(t)) = 1, \\ h_i(z(t)) \geq 0, \end{cases} \quad \text{para todo } t.$$

As Eq. (2.10) e (2.11) representam o processo da *defuzzyficação*, em que se inicia a combinação do modelos lineares para formação do modelo não-linear.

### 2.3.2 Modelo T-S baseado em grânulos evolutivos elipsoidais *fuzzy* (EEFIG)

No contexto de obter representações baseadas em dados, o Princípio de Granularidade Justificada (PJJ do inglês *Principle of Justifiable Granularity*) é uma referência importante por sua capacidade de produzir representações abstratas a partir de dados mantendo sua semântica. O conceito por trás do PJJ se baseia na formação de grânulos de informação (IG do inglês *Information Granules*), que maximizam as evidências experimentais e a semântica, sendo retratados pelas variáveis *coverage* e *specificity*. Em destaque, o PJJ é indicado como um problema multiobjetivo para calcular os limites ótimos dos grânulos, utilizando de algoritmos como o algoritmo genético.

Em Cordovil et al. (2020) é proposto um modelo Takagi-Sugeno baseado em grânulos evolutivos elipsoidais para reconstrução de dados e predições um passo a frente, em que, é desenvolvido um processo de granulação *online*, baseados nos grânulos evolutivos elipsoidais *fuzzy* (EEFIG do inglês *Evolving Ellipsoidal Fuzzy Information Granules*), e um processo de parametrização de fluxo de dados, baseado no PJJ.

#### Grânulos evolutivos elipsoidais

Dado uma hiper-elipsoide  $\varepsilon$  com centro em  $\mu$  e uma matriz de covariância  $C$ , Cordovil et al. (2020) define

$$\varepsilon = \{x \in \mathbf{X} \mid \mathcal{M} \leq \kappa\},$$

em que

$$\begin{cases} \mathcal{M} = (x - \mu)^\top C^{-1}(x - \mu), \text{ é a distância Mahalabonis entre um dado } x \text{ e } \mu; \\ \kappa = (\chi)_p^{-1}(\gamma), \text{ é um limiar embasado no inverso de um distribuição chi ao quadrado.} \end{cases}$$

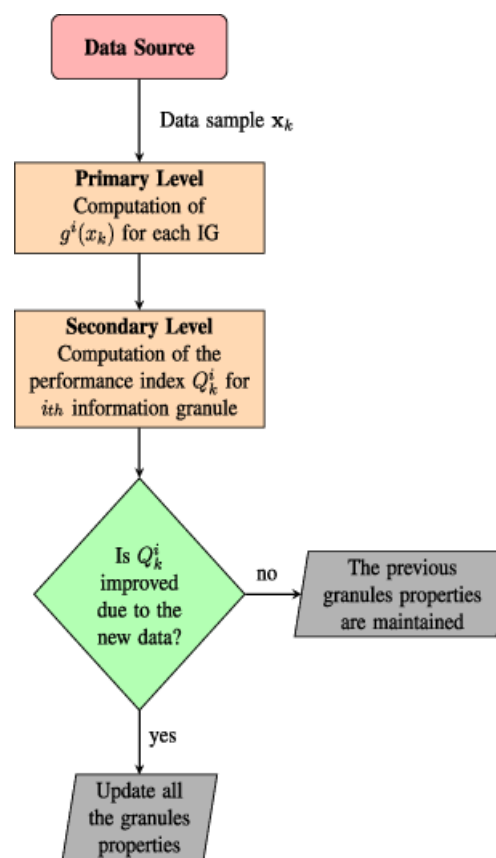
O limiar  $\kappa$  tem a função de reconhecer um dado considerado como anomalia em relação a um dado normal que foi inserido no elipsoide. A formulação do conjunto de dados é representado em termos da média ( $\mu_j$ ), limiar esquerdo da média ( $L_j$ ) e limiar direito da média ( $R_j$ ).

Com o objetivo de formular um processo *online* de granulação para um fluxo dados, tendo em vista que todo  $x_k$  está disponível a todo instante  $k$ , uma estrutura de atualização de grânulos, ver Fig. 4, é apresentada baseada em duas etapas.

Etapa 1. Referente a aferição das funções de pertinência, embasado nos IGs atuais, assim inserindo um novo dado em um IG.

Etapa 2. Quando um novo dado é inserido, é acarretado uma condição de atualização das propriedades de cada IG, se o índice de performance,  $Q_j$ , tiver melhora.

Figura 4 – Estrutura de atualização EFIG.



Fonte: Cordovil et al. (2020).

## Modelo *fuzzy* T-S granular

Como visto na subseção 2.3.1, a modelagem *fuzzy* T-S consiste de uma combinação de sistemas locais para formação da não-linearidade. A representação de modelos locais de médias EEFIG permite a identificação da dinâmica de tal setor, justificado pelo critério de desempenho. O índice de desempenho determina a qualidade do modelo feito, em que, ao atualizar os grânulos, os limiares dos parâmetros das funções de pertinência também serão atualizados resultando na modificação do conseqüente das regras *fuzzy* do modelo.

Com intuito de desenvolver a estrutura para modelagem e estimação, tem-se como regra padrão:

$$\text{Regra } i : \text{ SE } \mathbf{x}_k \text{ é } \mathcal{G}^i \text{ ENTÃO } \mathbf{x}_{k+1} = \mathbf{x}_k \Theta^{i^\top},$$

em que  $\mathbf{x}_{k+1}$  é o conseqüente da regra, definido em Cordovil et al. (2020) por

$$\mathbf{x}_{k+1} = \sum_{i=1}^{N_k} g_k^i f(\mathbf{x}_k),$$

com  $g_k^i$  sendo o grau de ativação normalizado da  $i$ -ésima regra.

Em consequência da atualização dos grânulos, instigado pelo fluxo de novos dados, as conseqüentes das regras são atualizadas utilizando o Estimador Recursivo de Mínimos Quadrados (ERMQ) estipulando uma relação entrada-saída.

## 2.4 Online Elliptical Clustering

Algoritmos de clusterização são utilizados para explicar a estrutura de um fluxo de dados, de forma a ser um aprendizado não-supervisionado. Habitualmente são divididos em 2 classes, a primeira sendo uma clusterização generalista de dados, onde não há uma seqüência a se seguir e se faz necessário saber a quantidade de *clusters* previamente (XU; TIAN, 2015).

A segunda classe adota uma ordem sequencial para os dados (séries temporais), em que se tem a prerrogativa de que observações próximas, em passos de tempo, são relacionadas, estas comumente utilizando técnicas de detecção de pontos de mudança no fluxo de dados para determinar a quantidade de *clusters* (XU; TIAN, 2015).

Moshtaghi, Leckie e Bezdek (2016) propõe um algoritmo de clusterização baseado em *clusters* evolutivos hiper-elipsoidais para estruturação do fluxo de dados, baseado na ideia de detecção de anomalias por captura incremental de dados ponderada (wIDCAD do inglês *Weighted Incremental Data Capture Anomaly detection*), e um rastreador de estados rápido capaz de detectar novos *clusters* através do conceito de c-separação.

### 2.4.1 Parametrização dos *clusters* elipsoidais

Dessa forma, dado um fluxo de dados  $\mathbf{x}_k = [x_0, \dots, x_k]$  sendo as primeiras  $k$  amostras, onde cada uma é um vetor  $p \times 1 \in \mathbb{R}^p$ . O hiper-elipsoide com raio  $t$  centrado na média das primeiras  $k$  amostras,  $m_k$ , com matriz de covariância  $S_k$  é definida como

$$e_k(m_k, S_k^{-1}; t) = \{x \in \mathbb{R}^p \mid (x - m_k)^\top S_k^{-1} (x - m_k) \leq t^2\}, \quad (2.12)$$

em que  $(x - m_k)^\top S_k^{-1} (x - m_k)$  é a distância Mahalanobis do dado  $x$  para o centroide do *cluster*  $m_k$  e  $S_k^{-1}$  é a matriz de  $e_k$ . Assim, pode-se definir o limite do *cluster* sendo

$$\delta_{e_k}(m_k, S_k^{-1}; t) = \{x \in \mathbb{R}^p \mid (x - m_k)^\top S_k^{-1} (x - m_k) = t^2\}. \quad (2.13)$$

Considerando os parâmetros do *cluster*, determina-se como anomalia para com  $e_k$  qualquer vetor de dados  $\mathbf{x} \in \mathbb{R}^p$  que está fora de  $e_k$ , de acordo com:

$$\mathbf{x} \text{ é anômalo para } e_k \Leftrightarrow (x - m_k)^\top S_k^{-1} (x - m_k) > t^2. \quad (2.14)$$

A medida  $t^2$  é utilizada como parâmetro para descrever o que  $e_k$  contém, seus limites e o que não contém. A mesma é calculada através de uma medida estatística utilizada para testes de hipótese em dados distribuídos em categorias, se tratando de variáveis qualitativas, chamada Qui-Quadrado,  $\chi^2$ . Dessa forma, temos

$$t^2 = (\chi^2_p)^{-1}(\gamma) \quad (\text{i.e., o inverso do qui-quadrado com } p \text{ graus de liberdade}).$$

### 2.4.2 Condição de atualização (wIDCAD)

O desenvolvimento de um *cluster* hiper-elipsoidal demonstrado anteriormente considera somente um conjunto de dados de tamanho  $k$  inicial e os parâmetros  $e_k$  e  $\delta_{e_k}$ , que representam a dinâmica do *cluster*. Como o objetivo é formular um processo *on-line* de clusterização para um fluxo de dados, utiliza-se o algoritmo de Detecção de Anomalias por Captura Incremental de Dados Ponderada (wIDCAD do inglês *Weighted Incremental Data Capture Anomaly Detection*) para atualização dos *clusters* em  $k + 1$ .

O wIDCAD atualiza o centroide, Eq. (2.15), e a matriz inversa de covariância, Eq. (2.16), onde há um peso  $w_k$  respectivo para cada observação. Os pesos são os graus de pertinência normalizados de  $x_{k+1}$  em relação a cada *cluster*, designado como  $w_{k+1,i}$ , onde  $i$  representa o conjunto de *clusters*  $|C|$ , calculados pelas Eq. (2.12) e (2.13), no tempo  $k$  ( $\sum w_{k+1,i} = 1$ ).

$$m_{k+1,i} = m_k + \frac{w_{k+1,i}}{\sum_{j=1}^{k+1} w_{ij}} (x_{k+1} - m_{t,i}); 1 \leq i \leq |C| \quad (2.15)$$

$$S_{k+1,i}^{-1} = \chi_{ki} \times \left[ S_{ki}^{-1} - \frac{S_{ki}^{-1}(x_{k+1} - m_{ki})(x_{k+1} - m_{ki})^\top S_{ki}^{-1}}{\delta_{ki}(x_{k+1} - m_t)^\top S_{ki}^{-1}(x_{k+1} - m_{k,i})} \right] \quad (2.16)$$

em que,

$$\beta_{k,i} = \sum_{j=1}^k w_{i,j}, \quad (\text{Somatório dos pesos atribuídos});$$

$$\alpha_{k,i} = \sum_{j=1}^k w_{i,j}^2, \quad (\text{Somatório dos pesos atribuídos ao quadrado});$$

$$\chi_{k,i} = \frac{\beta_{k,i}(\beta_{k+1,i}^2 - \alpha_{k+1,i})}{\beta_{k+1,i}(\beta_{k,i}^2 - \alpha_{k,i})}, \quad (\text{Variável direta de } t^2);$$

$$\delta_{k,i} = \frac{\beta_{k+1,i}(\beta_{k,i}^2 - \alpha_{k,i})}{\beta_{k,i}w_{k+1,i}(\beta_{k+1,i} + w_{k+1,i} - 2)}, \quad (\text{Limite do cluster}).$$

Com o intuito de resguardar o algoritmo de ruídos ou *outliers* de valores expressivos, foram considerados dois valores para  $\gamma$  em (2.14),  $\gamma_1 = 0,99$  para o limite normal de cada *cluster* e  $\gamma_2 = 0,999$  para fornecer uma “zona de guarda” ao redor do limite normal. Dados novos que caem na zona de guarda são considerados anomalias que resultam na atualização do *cluster*. Para garantir poder estatístico aos limites, é considerado um período de estabilização  $n_s$  em que a zona de guarda é ativada somente após.

Dada as fórmulas e condição para atualização dos clusters, os passos abaixo se seguem:

1. Com a chegada de  $x_{k+1}$ , calcula-se o conjunto de  $|C|$  distâncias Mahalonobis, onde  $M_{k+1,i} : 1 \leq i \leq |C|$ , de  $x_{k+1}$  aos  $|C|$  centroides dos *clusters* atuais  $m_{k,i}$ .
2. Testar a observação em relação a Eq. (2.14) e os limites elipsoidais definidos por  $\gamma_2$ .
3. Seja  $C' \subseteq C$  um conjunto de *clusters* em que a nova amostra encontra-se na sua zona de guarda. Os pesos associados com a nova observação para cada *cluster*  $C_i$  em  $C'$ , são calculados utilizando a Eq. (2.17).

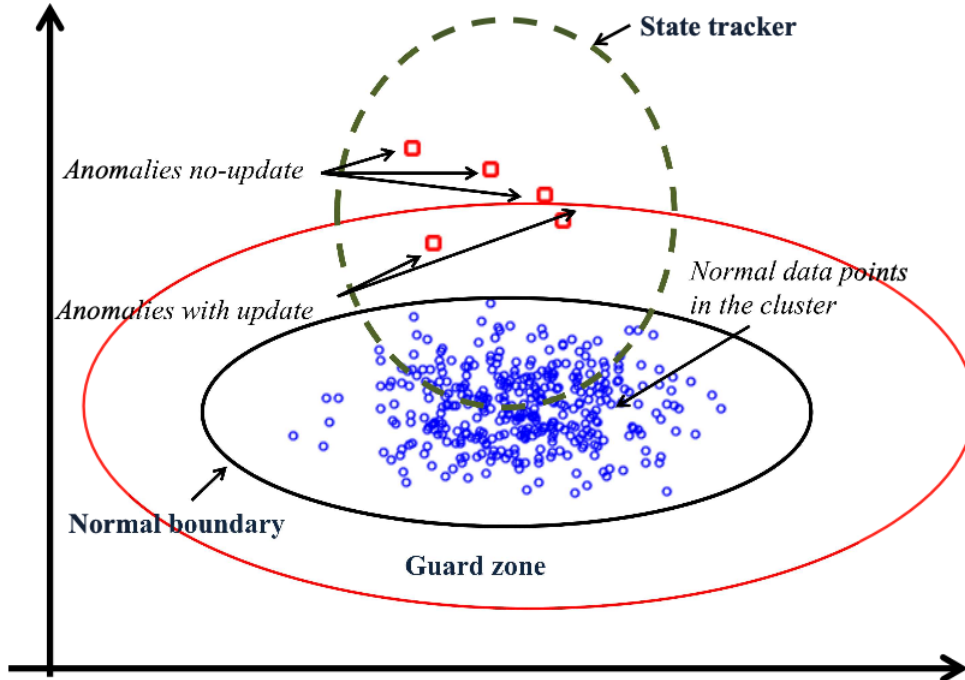
$$w_{k+1,i} = \frac{e^{-\frac{1}{2}M_{k+1,i}}}{\sum_{j=1}^{|C'|} e^{-\frac{1}{2}M_{k+1,j}}}. \quad (2.17)$$

### 2.4.3 Novos Clusters baseados em Rastreador de Estado do Sistema

A identificação de novos *clusters* em séries temporais, comumente são embasados em algoritmos de detecção de pontos de mudança baseados em limiares definidos previamente por um usuário. Com o intuito de identificar novos *clusters* sem limiares pré-definidos, Moshtaghi, Leckie e Bezdek (2016) apresentam um rastreador de estado do sistema utilizando um modelo de *cluster* único com fator de esquecimento.

Na Fig. 5 pode-se observar o funcionamento do modelo proposto por Moshtaghi, Leckie e Bezdek (2016), onde temos os seguintes fatores destacados:

Figura 5 – *Cluster* Inicial com sua zona de guarda e o Rastreador de Estado no início da mudança de estado do sistema.



Fonte: Moshtaghi, Leckie e Bezdek (2016)

1. **Normal boundary:** Limite normal do *cluster* representado pela medida  $t^2$  utilizando  $\gamma_1 = 0.99$ .
2. **Guard zone:** Zona de guarda ao redor do limite normal representado pela medida  $t^2$ , utilizando  $\gamma_2 = 0.999$ .
3. **Normal data points in the cluster:** Observações pertencentes ao *cluster*.
4. **Anomalies with update:** Observações fora dos limites normais do *cluster*, mas pertencentes a zona de guarda, resultante em atualização dos limites do *cluster*.
5. **Anomalies no-update:** Observações fora dos limites normais e fora da zona de guarda, sem atualizações.
6. **State tracker:** Rastreador de Estados.

Para atualização dos parâmetros do rastreador de estado,  $m_{k+1,\lambda}$  e  $S_{k+1,\lambda}^{-1}$ , no tempo  $k + 1$  com fator de esquecimento  $\lambda \in [0,9; 1)$ , as Eq. (2.18) e (2.19) são utilizadas de forma evolutiva.

$$m_{k+1,\lambda} = \lambda m_{k,\lambda} + (1 - \lambda)x_{k+1} \quad (2.18)$$

$$S_{k+1,\lambda}^{-1} = \frac{kS_{k\lambda}^{-1}}{\lambda(k-1)} \times \left[ I - \frac{(x_{k+1} - m_{k\lambda})(x_{k+1} - m_{k\lambda})^\top S_{k\lambda}^{-1}}{\frac{(k-1)}{\lambda} + (x_{k+1} - m_{k\lambda})^\top S_{k\lambda}^{-1}(x_{k+1} - m_{k\lambda})} \right] \quad (2.19)$$

Para que não haja um crescimento descontrolado do rastreador, o valor de  $k$  é restringido ao seu valor efetivo baseado no horizonte de memória do rastreador. Após  $N_{eff} = 3\tau$  observações, então  $k = N_{eff}$ , onde  $\tau = \frac{1}{1-\lambda}$  é o horizonte de memória com fator de esquecimento exponencial  $\lambda$ .

A partir do rastreador, é definida a regra para criação de um novo *cluster* fundamentado no conceito de  $c$ -separação no contexto de misturas gaussianas. Duas densidades gaussianas  $p$ -variadas  $N_p(\mu_i, \Sigma_i)$  e  $N_p(\mu_j, \Sigma_j)$  são  $c$ -separadas se

$$\|\mu_i - \mu_j\| \geq c \sqrt{p \times \max(d_{max}(\Sigma_i), d_{max}(\Sigma_j))}, \quad (2.20)$$

em que  $d_{max}(\Sigma)$  é o maior autovalor de  $\Sigma$ , e o valor de  $c \in [0, \infty)$  representa o grau de separação de duas distribuições, onde uma mistura 2-separada corresponde a duas distribuições gaussianas quase completamente separadas.

Assim, um *cluster* novo é criado a partir do momento em que o rastreador se tornar 2-separado dos *clusters* existentes, normalmente acontecendo com um salto de anomalias consequentes. Para fornecer informações sobre o início de mudança em um fluxo de dados, é utilizado um *buffer* de anomalias consequentes, baseados no limiar  $\mu_i = 0,99$ .

Com a condição de criação de um novo *cluster* alcançada, a observação inicial do *buffer* é sinalizada como o ponto de mudança e os dados contidos nele são usados para criar o novo *cluster*, e o mesmo é zerado. Para que não seja computacionalmente custoso executar a condição de  $c$ -separação para cada nova observação, Moshtaghi, Leckie e Bezdek (2016) determinam que o *buffer* necessita de um mínimo de  $p + 1$  anomalias consequentes no *buffer*, em consideração do cálculo da matriz inversa de covariância em  $p$ -dimensões.

## 2.5 Explainable Artificial Intelligence (XAI)

Os primeiros sistemas de inteligência artificial (IA) eram notáveis por sua interpretabilidade, contrastando com o surgimento posterior de sistemas de decisão opacos, exemplificados pelas Redes Neurais Profundas (RNP). Esse avanço tecnológico tem suscitado uma demanda crescente por transparência, especialmente entre os diversos intervenientes no campo da IA. O êxito empírico de modelos de Aprendizado Profundo, como as RNP, tem sua origem na intersecção entre algoritmos de aprendizado eficazes e espaços paramétricos vastos. A complexidade desses modelos, caracterizados como caixas-pretas, apresenta desafios para a compreensão direta de seus mecanismos internos, levantando preocupações

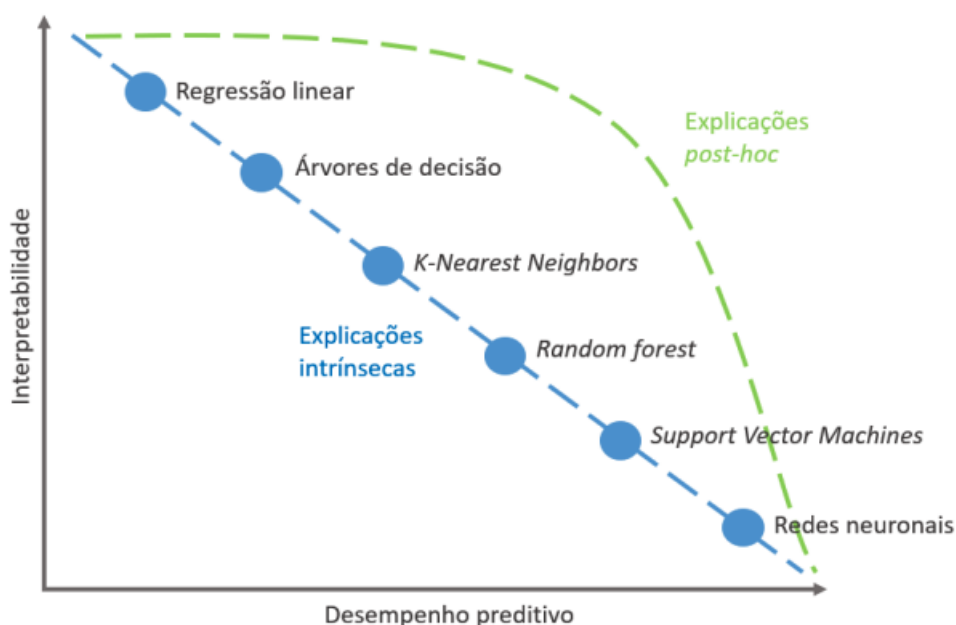
em relação à justificabilidade e legitimidade das decisões que produzem (ARRIETA et al., 2019).

Embora o objetivo de sistemas de IA seja primordialmente obter resultados com erros mínimos, seja por otimização de uma função de perda ou outra metodologia, o propósito do mundo real é disponibilizar informações úteis. Trindade (2020) detalha a importância da interpretabilidade e da explicabilidade para dois casos:

1. Campos de atuação onde há consequências graves: Medicina, Justiça criminal e Mercados Financeiros.
2. Problemática nova, sem comprovação ou pesquisa científica de forma a creditar na escolha do sistema.

Modelos de IA caixa-preta apresentam alta complexidade, muitas vezes impossibilitando o entendimento do seu funcionamento e o “porquê” por trás dos resultados, porém comumente a acurácia dos resultados do modelo é inversa ao desempenho preditivo do mesmo, como exemplificado por Trindade (2020) na Fig. 6.

Figura 6 – Relação da Interpretabilidade e do Desempenho preditivo de modelos de IA.



Fonte: Trindade (2020).

Em meio a esta circunstância, o campo de pesquisa de Inteligência Artificial Explicável (XAI do inglês *Explainable Artificial Intelligence*) surgiu com intuito de que tais sistemas sejam mais fáceis de serem compreendidos por seres humanos. Ao mesmo tempo

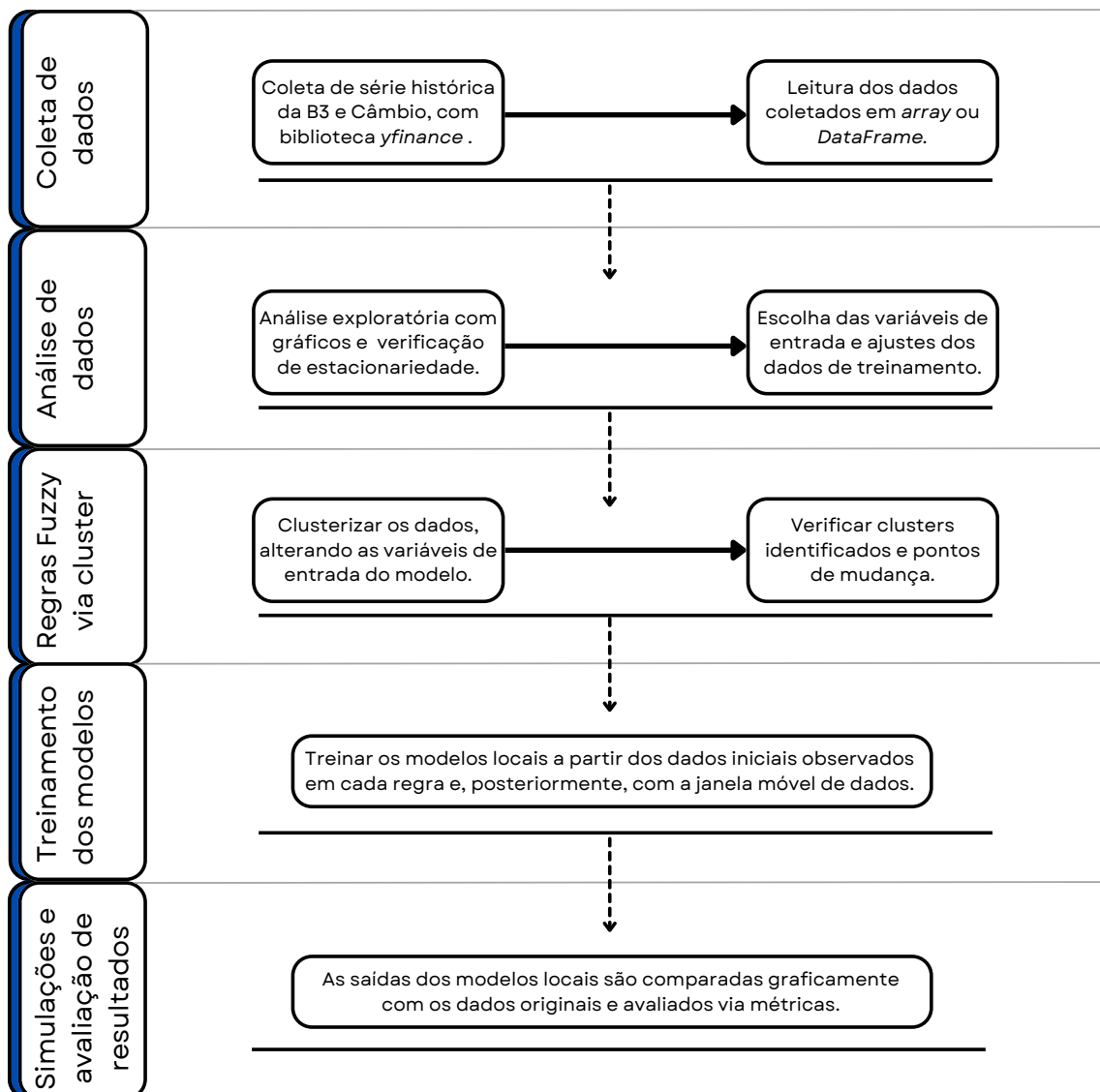
que ainda mantenham alta eficiência, os humanos precisam entender o modelo para que haja legitimidade do mesmo (GKATZIA; LEMON; RIESER, 2016; YE; JOHNSON, 1995; RANJBAR; MOMTAZI; HOMAYOONPOUR, 2024). No contexto de XAI, Biran e Cotton (2017) denotam de forma mais direta os conceitos de Interpretabilidade e Explicabilidade:

- Interpretabilidade refere-se a aptidão de um modelo ter seu funcionamento compreendido por um ser humano. Alguns modelos mais simples, como um sistema baseado em regras, tem essa capacidade de forma inerente, porém sistemas mais complexos necessitam de abordagens complementares.
- Explicabilidade refere-se ao “porquê” de um resultado específico, estabelecendo uma forte relação quanto a entrada e a saída do modelo, podendo, como exemplo, apontar alguma observação ou evento determinado que teve maior impacto para decisão do sistema.

### 3 Metodologia

Com o objetivo de apresentar a abordagem utilizada no projeto, este capítulo esquematiza os passos tomados durante o desenvolvimento do mesmo e suas metodologias, divididas por funcionalidade de acordo com a Figura 7.

Figura 7 – Fluxograma de desenvolvimento do projeto de pesquisa.



Fonte: O Autor.

## 3.1 Coleta de dados

De forma inicial, realiza-se a coleta dos dados a partir de algum banco de dados ou arquivo em texto. Como o projeto trata de sistemas financeiros, o objeto da pesquisa é o do índice IBOV (Ibovespa B3), representativo da dinâmica da bolsa de valores brasileira. Destaca-se que outras variáveis de entrada são necessárias para que o modelo tenha informações externas necessárias, assim [Jacques, Borges e Miranda \(2020\)](#) evidencia a forte relação entre indicadores macroeconômicos e empresas listadas na B3, estes indicadores sendo o Câmbio(dólar), Juros (SELIC) e inflação (IPCA), em que a relação do câmbio é evidenciada também em [GUIDE INVESTIMENTOS \(2023\)](#).

Para tal coleta, utilizou-se a linguagem *Python*<sup>1</sup>, sendo possível obter os dados de ativos financeiros presentes na B3 utilizando a biblioteca *yfinance*<sup>2</sup>, esta sendo uma biblioteca *open-source*. Nesta fase inicial, os dados coletados a partir da biblioteca *yfinance* tem o formato de *Data Frame* (biblioteca *Pandas*<sup>3</sup>, mas também podem ser trabalhados em formatos de listas (estrutura inerente a linguagem *Python*) ou de *Array* (biblioteca *NumPy*<sup>4</sup>).

Para obtenção dos índices macroeconômicos, a biblioteca *Pandas* apresenta funções capazes de ler resultados de *sites* em formatos específicos. Para isso, utilizando a API<sup>5</sup> (do inglês *Application Programming interface* do Banco Central conectado ao Sistema Gerenciador de Séries Temporais (SGS) foi possível obter os dados históricos referentes ao Juros e Inflação, em formato *JSON*<sup>6</sup>.

## 3.2 Análise de dados

Após coleta e leitura dos dados, é importante verificar se os formatos e estruturas dos dados são coerentes para funcionamento adequado dos algoritmos. Além disso, por se tratar de séries temporais, há de ser feita uma análise exploratória quanto a sua estacionariedade, uma hipótese comumente falsa em sistemas financeiros, em que as propriedades de uma série não dependam do tempo em que está sendo observada. [Kwiatkowski et al. \(1992\)](#) apresentam uma forma de transformar séries não estacionárias para estacionárias, através da diferenciação, em que é feita a diferença entre observações consecutivas.

Para verificação da condição de não estacionariedade de uma série, utiliza-se o teste de *Dickey-Fuller* Aumentado (ADF), primeiramente mencionado por [Said e Dickey \(1984\)](#). Também chamado de teste de raiz unitária, o mesmo verifica se o critério do teste de hipótese

---

<sup>1</sup> <https://www.python.org/>

<sup>2</sup> <https://pypi.org/project/yfinance/>

<sup>3</sup> <https://pandas.pydata.org/>

<sup>4</sup> <https://numpy.org/>

<sup>5</sup> <http://api.bcb.gov.br/dados/serie/bcdata.sgs./dados?formato=json>

<sup>6</sup> <https://www.json.org/json-en.html>

nula é alcançado ou não, verificando o *P-value* do teste sendo:  $P\text{-value} > 0,05$ , Hipótese nula aceita (não-estacionária);  $P\text{-value} < 0,05$ , Hipótese nula rejeitada(estacionária) (BUENO, 2011). Para efetuar o teste ADF em *Python* utiliza-se a biblioteca *statsmodels*<sup>7</sup>, a mesma provendo um vasto conjunto de funções para estimação de modelos estatísticos, teste de hipótese e exploração estatística de dados.

Além disso, análises gráficas quanto ao comportamento dos dados no tempo são necessárias, como um aprofundamento quanto a visualização das funções de autocorrelação (FAC) e autocorrelação parcial (FACP) para determinar os parâmetros iniciais, determinadas por Bueno (2011), dos modelos, posteriormente. De forma a ajustar os parâmetros determinados inicialmente, alguns critérios de informação podem ser utilizados, de forma que minimize uma função baseada nos resíduos, aplicando uma penalidade com o aumento do número de parâmetros.

### 3.3 Regras *Fuzzy* via Cluster

O processo de geração das regras *fuzzy* para um modelo *fuzzy* é responsável por indicar a relação de entrada(s) e saída de forma linear, com o intuito de encontrar regiões com dinâmicas lineares em um sistema com dinâmica não-linear. Tratando-se de dados sequenciais, o fator tempo se torna crucial para lidarmos com a criação destas regras, partindo do pressuposto que dados consecutivos no tempo possuem algum tipo de similaridade entre si.

Assim, para estruturar o fluxo de dados sequenciais e gerar as regras, os dados serão agrupados em hiper-elipsoides evolutivas embasadas na detecção de anomalias de forma incremental. De acordo com Moshtaghi, Leckie e Bezdek (2016) e o processo de parametrização dos *clusters* elipsoidais (veja Subseção 2.4.1), inicialmente tem-se um vetor de dados de tamanho  $k$ , com raio  $t$  centrado na média das amostras iniciais.

#### 3.3.1 Parametrização Inicial das regras

Para uma execução hábil, foi criada uma Classe chamada *Regra*, em que inicialmente pode-se inicializar um *cluster* vazio ou um *cluster* com parâmetros baseados nas primeiras  $k$  amostras, de acordo com o Algoritmo 1. Os parâmetros do *cluster* serão utilizados durante todo o processo, estes sendo:

- **k**: Quantidade de dados pertencentes ao *cluster*.
- **Matriz**: Matriz de covariância  $S_k^{-1}$ .
- **Centros**: Centróide do cluster  $m_k$ .

<sup>7</sup> <https://www.statsmodels.org/stable/index.html>

- **PesoSoma**: Somatório dos pesos atribuídos (utilizado na Eq. (2.16)).
- **PesoSomaExp**: Somatório dos pesos atribuídos ao quadrado (utilizado na Eq. (2.16)).

---

**Algoritmo 1:** Classe Regra - Função Rule e InitRule
 

---

```

Propriedades k, Matriz, Centros, PesoSoma, PesoSomaExp;
/
/ Função para inicialização dos parâmetros um Cluster
/
function Rule(D);
Cluster = Cluster.InitRule(D);
return [Cluster];
/
/ Função para inicialização de um Cluster
/
function InitRule(D, Dados=None);
Regra = Rule;
/
if Dados == None then
| Inicializar as propriedades de Regra com valor nulo;
else
| Inicializar as propriedades k, PesoSoma, PesoSomaExp da Regra com valor
| nulo;
|
| Regra.Centros ← média das variáveis das duas primeiras observações;
|
| Atualizar os parâmetros recursivamente com os dados restantes
|
| for i = 2 a tamanho(Dados) do
| | Limiar1 ←  $(\chi^2)_D^{-1} * (0,99)$ ;
| | [Regra, _] = Regra.StepWeighted(Regra,1,Dados[i,:D], Limiar1, 1)
| end
end
return [Regra];

```

---

A função *InitRule* consegue receber duas variáveis de entrada, necessárias para essa parametrização inicial do primeiro *cluster*, assim:

- **D**: A quantidade de variáveis a serem agrupadas, ou seja, se temos uma saída e uma entrada, então teria o valor igual a 2. No caso de trabalhar com mais de uma entrada, por exemplo utilizando uma média móvel, com uma janela fixa, com intuito de agregar uma característica de comportamento em médio prazo, então a variável teria o valor de 3.
- **Dados**: O vetor de amostras iniciais com  $D$  colunas e  $k$  linhas. Se essa variável não for recebida pela função, então a mesma inicia com parâmetros nulos. Caso contrário, ela

inicializa o centroide como uma média das duas primeiras amostras e, após, atualiza recursivamente com os dados restantes, utilizando a função *StepWeighted*.

### 3.3.2 Atualização Recursiva das regras

Após essa inicialização utilizando as primeiras  $k$  observações, as condições de atualização são executadas de forma recursiva a cada nova observação. Assim, a função *EvaluateIDCAD*, exemplificada com o Algoritmo 2, calcula as distâncias da nova observação ao centroide do(s)  $cluster(s)$  e assim associa um peso da amostra referente a pertinência àquela regra. A partir dos pesos calculados, é verificada a condição de atualização e anomalia. A função *EvaluateIDCAD* recebe 3 variáveis, a saber:

- **Clusters:** A lista que contém as Regras, criadas a partir da Classe *Regra*, contendo todos seus parâmetros.
- **Dado:** Observação no tempo  $k + 1$ .
- **Limiar1:** Limite normal do *Cluster*, definido na Subseção 2.4.2.

---

#### Algoritmo 2: Função EvaluateIDCAD

---

```

Data: [Clusters, Dado, Limiar1]
Result: [Clusters, Anomalia]
/
x ← Dado;
mk ← Clusters.Centros;
Sk ← Clusters.Matriz;
Inicializar ns (período de estabilização);
/
/ Calcular as distâncias Mahalonobis como nas Eqs. (2.12) e (2.13)
/
for j = 0 a tamanho(Dado) do
  for i = 0 a tamanho(Clusters) do
    mahaldist ← (xj - mi)⊤ Si-1 (xj - mi);
    distance[i, j] ← mahaldist; / Variável q armazena as distâncias
    calculadas/
  end
end
/ Testar a nova observação contra Eq. (2.13) e o limite γ2 (Subseção
  2.4.1)
/
Atividade ← Eq. (2.17); / Calcular os pesos associados/
for i = 0 a tamanho(Clusters) do
  Clusters[i], Anomalia ← Clusters[i].Step(x, ns, Limiar1, Atividade[i]);
end

```

---

Para verificar se a nova amostra é uma anomalia e/ou condiz com as condições de atualização dos parâmetros de algum *cluster*, a função *Step*, Algoritmo 3, pertencente a Classe *Regra*, é utilizada de forma a percorrer todas os *clusters* existentes, retornando o *cluster* atualizado ou não, e se a amostra é uma anomalia ou não.

A função *Step* recebe 5 variáveis, onde:

- **Cluster**: Uma regra pertencente a lista de regras, caso não especificada, considera como a regra em questão.
- **Dado**: Observação no tempo  $k + 1$ .
- **$n_s$** : Período de estabilização estatística.
- **Limiar1**: Limite normal do Cluster, definido na Subseção 2.4.2.
- **Peso**: Peso da observação associado a regra. Caso não especificada, considera como *None*.

---

**Algoritmo 3:** Classe Regra - Função Step
 

---

```

/                                                                 /
/ Função para efetuar um passo com verificação do Peso e Período de
  estabilização                                                                 /
/                                                                 /
function Step(Cluster=self, Dado,  $n_s$ , Limiar1, Peso=None);
modelo ← Cluster;
if Peso = None then
  | Peso ← 1;
else
  | Peso ← Peso;
end
/                                                                 /
/ Condição do período de estabilização                                                                 /
/                                                                 /
if modelo.k <  $n_s$  then
  | [novoModelo, Anomalia] = modelo.StepWeighted(Peso, Dado, Limiar1,
  |   Atualizar = 1);
else
  | [novoModelo, Anomalia] = modelo.StepWeighted(Peso, Dado, Limiar1,
  |   Atualizar = 2);
end
return [novoModelo, Anomalia];

```

---

A função *Step*, em suma, realiza a verificação da existência de um peso associado e se a regra em questão já atingiu o período de estabilização, ou seja, se o *cluster* já atingiu

a quantidade de dados mínimo necessário. Nesta condição, a função *StepWeighted* é a responsável pela verificação de anomalia e atualização dos parâmetros do *cluster*.

A função *StepWeighted*, Algoritmo 4, inicialmente calcula a distância Mahalonobis do dado em referência ao centróide da regra. Em seguida, efetua a verificação se essa distância é maior que o limite normal do *cluster*, em que no caso de ser verdadeiro identifica a variável *Anomalia* como verdadeira. Ainda na condição verdadeira para anomalia, também é verificado se o dado está fora da zona de guarda, assim retornando o *cluster* sem atualização.

---

**Algoritmo 4:** Classe Regra - Função WeightedStep
 

---

```

/
/ Função para atualização de um passo dos parâmetros do cluster de
  forma ponderada
/
function WeightedStep(Cluster=self, Peso, Dado, Limiar1, Atualizar);
Inicializa Anomalia ← 0 ;           / 0: Não é anomalia/
x ← Dado;
mk ← Cluster.Centros;
Sk ← Cluster.Matriz;
/
mahaldist ← (x - mk)⊤ Sk-1 (x - mk) ;
if mahaldist > Limiar1 then
  Anomalia ← 1;
  if mahaldist > limite γ2 & Atualizar > 1 then
    Atualizar ← 0;
    return;
  end
else
  Atualizar as propriedades do cluster de acordo com as Eqs. (2.15) e (2.16);
  Cluster.k = Cluster.k+1;
end
return [Cluster, Anomalia];

```

---

Caso a distância *mahaldist* for menor que o limite normal do *cluster*, então as propriedades da regra são atualizadas de acordo com a Subseção 2.4.2. Dentre as propriedades da regra, duas são utilizadas diretamente para o cálculo de atualização:

- **SomaPeso**: Denotado como  $w_{i,j}$  no calculo do somatório dos pesos atribuídos,  $\beta_{k,i}$ .
- **SomaPesoExp**: Denotado como  $w_{i,j}^2$  no calculo do somatório dos pesos atribuídos,  $\alpha_{k,i}$ .

A função *StepWeighted* recebe 5 variáveis, onde:

- **Cluster:** Uma regra pertencente a lista de regras, caso não especificada, considera como a regra em questão.
- **Peso:** Peso da observação associado a regra. Caso não especificada, considera como *None*. Denotado por  $w_{k+1,i}$ .
- **Dado:** Observação no tempo  $k + 1$ .
- **Limiar1:** Limite normal do *Cluster*, definido na Subseção 2.4.2.
- **Atualizar:** Condição de Atualização baseado no período de estabilização, no caso do dado ser uma anomalia que pertence a zona de guarda, determinada na função *Step*.

### 3.3.3 Rastreador de Estados

De forma paralela a inicialização do primeiro *cluster*, também é necessário fazer o mesmo com o rastreador de estados, de modo que sua matriz covariância se inicialize como uma matriz identidade com tamanho da variável  $D$ . O seu centroide é inicializado de forma igual ao centroide do primeiro *cluster*, com as primeiras  $k$  amostras.

---

#### Algoritmo 5: Função StateTrackerIDCAD

---

```

Data: [RastreadorMatriz, RastreadorCentro, TempoDados, Lambda, Dado]
Result: [RastreadorMatriz, RastreadorCentro]
/
 $x_{k+1} \leftarrow \text{Dado};$ 
 $m_{k\lambda} \leftarrow \text{Clusters.Centros};$ 
 $S_{k\lambda}^{-1} \leftarrow \text{Clusters.Matriz};$ 
 $\lambda \leftarrow \text{Lambda};$ 
 $\tau \leftarrow \frac{1}{1-\lambda};$  / Horizonte de Memória/
 $N_{eff} \leftarrow 3 * \tau;$  / Valor efetivo de  $k$ /
 $k \leftarrow \min(\text{TempoDados}, N_{eff});$  / Limita o  $k$  no seu valor efetivo/
/
/ Calcular um passo de atualização recursiva do Rastreador de
Estados (Subseção 2.4.3)
/
denMatriz  $\leftarrow \frac{(k-1)}{\lambda} + (x_{k+1} - m_{k\lambda})^\top S_{k\lambda}^{-1} (x_{k+1} - m_{k\lambda});$ 
numMatriz  $\leftarrow (x_{k+1} - m_{k\lambda})(x_{k+1} - m_{k\lambda})^\top S_{k\lambda}^{-1};$ 
/
Mult  $\leftarrow \frac{k}{\lambda * (k-1)};$ 
RastreadorMatriz  $\leftarrow \text{Mult} * (S_{k\lambda}^{-1} - \frac{S_{k\lambda}^{-1} * \text{numMatriz}}{\text{denMatriz}});$  / Eq. (2.19)/
/
RastreadorCentro  $\leftarrow \lambda * m_{k\lambda} + (1 - \lambda) * x_{k+1};$  / Eq. (2.18)/

```

---

Com os parâmetros do rastreador inicializadas, a função *StateTrackerIDCAD*, Algoritmo 5, atualiza os parâmetros do rastreador de forma recursiva no tempo  $k + 1$  com fator de esquecimento  $\lambda = 0,9$ , baseado no horizonte de memória e o valor efetivo para  $k$ .

A função *StateTrackerIDCAD* recebe 5 variáveis, onde:

- **RastreadorMatriz:** Matriz de covariância do rastreador de estados.
- **RastreadorCentro:** Centroe do rastreador de estados.
- **TempoDados:** Passo de tempo atual, indicando a quantidade de observações novas recebidas.
- **Lambda:** Fator de esquecimento, necessário para o cálculo do horizonte de memória e o valor efetivo de  $k$ .
- **Dado:** Observação no tempo  $k + 1$ .

### 3.3.4 Criação de Nova Regra

Com o Rastreador de Estados e a(s) regra(s) atualizados, inicia-se a verificação da condição de c-separação que calcula o grau de separação de duas distribuições gaussianas, que caso verdadeiro indica duas distribuições quase completamente separadas. A função *CSeparation*, Algoritmo 6, é responsável por verificar a condição imposta pela Eq. (2.20), retornando 1 para verdadeiro e 0 para falso.

---

#### Algoritmo 6: Função CSeparation

---

```

Data: [Clusters, RastreadorMatriz, RastreadorCentro, Separação]
Result: [CSeparado]
/
p ← tamanho(RastreadorCentro);
dmax( $\sum_{cluster}$ ) ← max(autovalor(Clusters[j].Centros));
dmax( $\sum_{rastreador}$ ) ← max(autovalor(RastreadorCentro));
Inicializar CSeparado ← 1 ; / Iniciar como separado/
/
/ Calcular a condição de CSeparação (Subseção 2.4.3) /
/
for j = 0 a tamanho(Clusters) do
    termo1 = norma(Clusters[i].Centros - RastreadorCentro);
    termo2 = Separação *  $\sqrt{p} \times \max(d_{max}(\sum_{cluster}), d_{max}(\sum_{rastreador}))$ ;
    / Condição da Eq. (2.20) /
    if termo1 < termo2 then
        | CSeparado ← 0;
    end
end

```

---

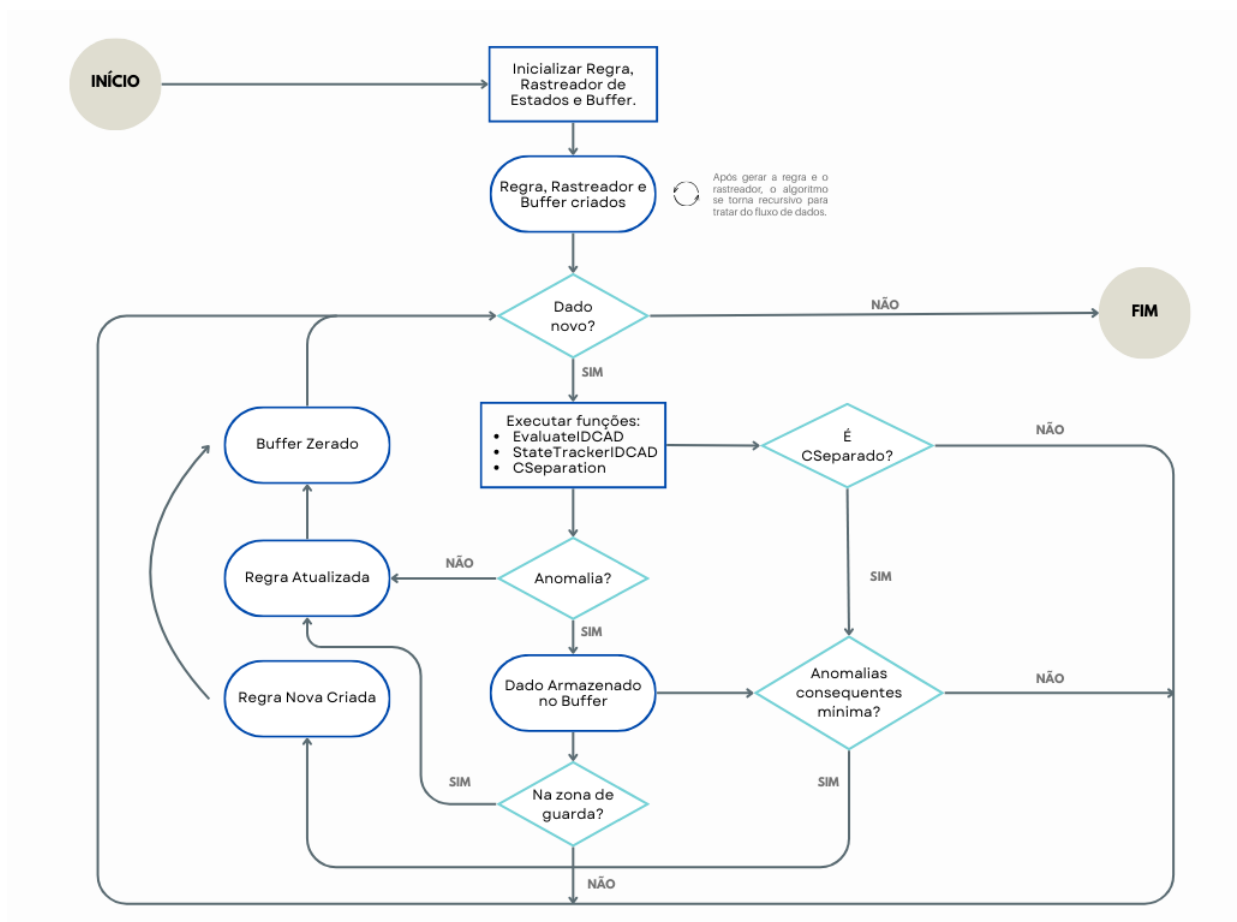
A função *CSeparation* recebe 4 variáveis, onde:

- **Clusters:** A lista que contém as Regras, criadas a partir da Classe Regra, contendo todos seus parâmetros.
- **RastreadorMatriz:** Matriz de covariância do rastreador de estados.
- **RastreadorCentro:** Centróide do rastreador de estados.
- **Separação:** Fator de esquecimento, necessário para o cálculo do horizonte de memória e o valor efetivo de  $k$ .
- **Dado:** Observação no tempo  $k + 1$ .

Comumente, essa condição é verdadeira quando ocorre uma explosão de anomalias consecuentes, assim um *buffer* é utilizado para armazenar estas anomalias consecuentes e para criar a nova regra, sendo resetado após.

Em suma, o algoritmo de clusterização funciona de acordo com a Fig. 8.

Figura 8 – Fluxograma de desenvolvimento do OEC.



Fonte: O Autor.

## 3.4 Modelos ARX e ARMAX

A utilização de análises gráficas quanto a Função de Autocorrelação e Autocorrelação Parcial para o estudo de modelos autoregressivos, no que tange a estimação inicial das ordens do modelo. Camilo (2012) demonstra que a FAC permite identificar a ordem  $q$  de um processo de Médias Móveis (MA), e que a FACP permite identificar a ordem  $p$  de um processo Auto-Regressivo (AR).

Utilizando a biblioteca *statsmodels*<sup>8</sup>, é possível criar os modelos necessários com funções específicas para a modelagem do tipo autoregressiva. As funções abaixo foram utilizadas:

- `ARIMA(endog, exog, order=(p, d, q))`: Essa Classe possibilita o estudo de modelos ARIMA (do inglês *AutoRegressive Integrated Moving Average*), inicializando o modelo de forma a receber 3 parâmetros:
  - `endog`: A variável endógena representa a série temporal observada.
  - `exog`: A variável exógena representa fatores independentes e externas ao modelo, mas com forte relação a variável estudada.
  - `order=(p, d, q)`: Parâmetro que descreve as ordens do modelo, onde:
    - \*  $p$ : Ordem da parte autoregressiva;
    - \*  $d$ : Ordem de diferenciação;
    - \*  $q$ : Ordem da parte de média móvel.
- `fit()`: Função que estima os parâmetros do modelo.
- `predict()`: Função que faz previsões in-sample e out-of-sample, recebendo um início e fim.
- `forecast()`: Função que faz previsões out-of-sample, recebendo o número de passos a frente.
- `get_forecast()`: Função que faz previsões out-of-sample, incluindo os intervalos de confiança.

**Note:** A Classe ARIMA pertence ao caminho `statsmodels.tsa.arima.model`, e as funções restantes pertencem a Classe ARIMA.

<sup>8</sup> <https://www.statsmodels.org/stable/index.html>

### 3.5 Interpretabilidade e Explicabilidade do Modelo

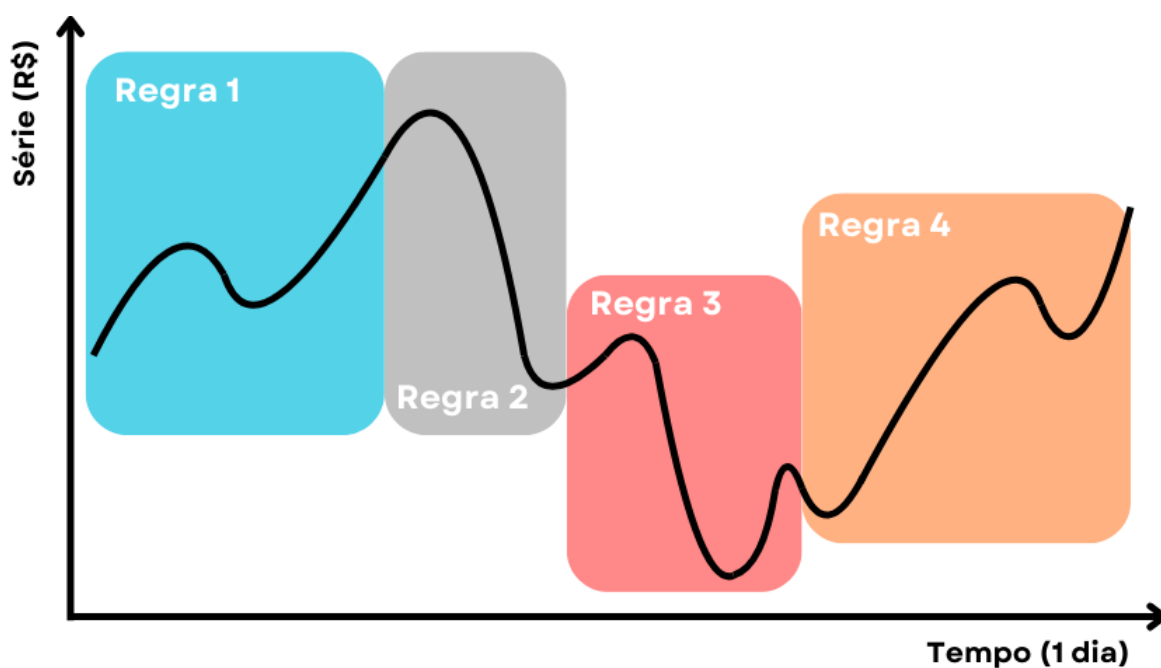
Os modelos autoregressivos comumente são os primeiros a serem estudados na linha de modelos de identificação de sistemas dinâmicos, por sua estrutura simples de entender e sua natureza polinomial, fazendo com que não necessite de um extenso conhecimento matemático e estatístico para seu entendimento, além de serem extensivamente aplicados em sistemas financeiros, chamando-os de modelos econométricos (BUENO, 2011).

Com uma alta capacidade de interpretabilidade, esses modelos conseguem estabelecer uma relação clara quanto a entrada e saída utilizados, devido sua estrutura parametrizada, podendo apontar que a observação de um dia específico teve um peso maior que de outros dado um coeficiente de importância do modelo para algum parâmetro.

Além disso, como na Figura 9, regiões na série temporal podem ser identificadas através de aplicação de algoritmos de clusterização, como o *Online Elliptical Clustering*, com intuito de manter dados com semântica similar agrupados, ou seja, é possível compreender eventos associados com a dinâmica da série em cada região. Assim, podemos também verificar quando uma mudança de dinâmica na série está ocorrendo.

Tratando de sistemas baseados em regras, sua capacidade de interpretabilidade se torna inerente, fazendo com que seu funcionamento seja melhor compreendido. Com isso, além de serem utilizadas para prever a dinâmica da série, também auxiliam na investigação quanto ao impacto de variáveis exógenas no comportamento da série estudada.

Figura 9 – Divisão de modelos locais por regras.



Fonte: O Autor.

## 4 Experimentos Computacionais

Neste capítulo serão apresentados 4 conjuntos de dados referentes à série temporal financeira do índice IBOV. O primeiro conjunto trata do índice IBOV e o restante das entradas exógenas, que serão utilizadas para criação das regras e identificação das diferentes regiões dentro da série. Para cada experimento serão discriminados as variáveis utilizadas, a análise exploratória inicial dos dados, as regiões locais identificadas e suas implicações, assim como os resultados e seus pontos principais. É importante destacar que o objetivo não é explorar minuciosamente as propriedades da série e dos modelos locais gerados, mas demonstrar, de maneira objetiva, a aplicação dos métodos propostos provendo a explicabilidade necessária.

Dois cenários serão apresentados para validação do modelo proposto, dado um cenário inicial mais simples utilizando somente a variável de saída e uma variável exógena; e um cenário mais complexo em que outras variáveis exógenas são acrescentadas para auxiliar na investigação do comportamento do índice IBOV. As simulações serão executadas em um ambiente local para *Python* utilizando a IDE (do inglês *Integrated Development Environment*) *DataSpell*<sup>1</sup> desenvolvido pela empresa *JetBrains*©, específico para análise exploratória de dados e prototipagem de modelos de aprendizado de máquina.

A primeira simulação será executada de forma *out-of-sample*, utilizando os dados da B3 e os dados do câmbio (dólar) para criação das regras e assim geração dos modelos locais. Somente a variável exógena do câmbio está sendo utilizada por ter frequência diária, em contraste com as outras variáveis macroeconômicas indicadas por [Jacques, Borges e Miranda \(2020\)](#). A análise exploratória, no que tange os gráficos de FAC e FACP, será feita de forma local para cada região identificada. A fim de verificar a qualidade semântica de criação das regras, será feita uma análise gráfica *in-sample* quanto a representatividade de cada regra identificada. Além disso, pretende-se verificar a acurácia dos modelos locais através das métricas de avaliação citadas na Subseção 2.2.

A segunda simulação, de forma similar a primeira, será executada de forma *out-of-sample* acrescentando duas médias móveis de curto prazo da série IBOV ([VARELLA, 2012](#)), com intuito de melhor inferir a dinâmica do sistema estudado no que tange a criação das regras e logo a representação dos modelos locais. Será verificado se o uso de mais variáveis para criação das regras foi eficaz na investigação das diferentes dinâmicas da série, demonstrando uma maior explicabilidade quanto às possíveis mudanças de comportamento.

Ressalta-se que o intuito com a aplicação do algoritmo *OEC* é o agrupamento de dados por similaridade ao longo de um período de tempo, não sua capacidade em

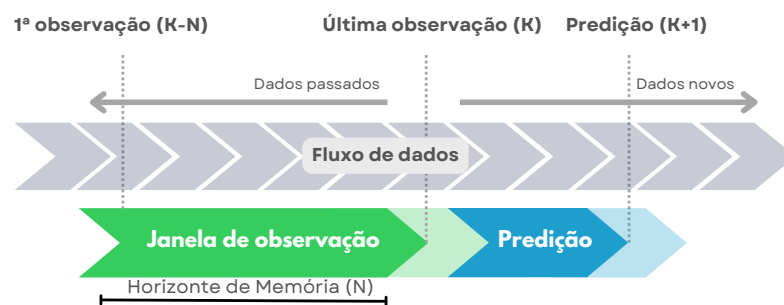
<sup>1</sup> <https://www.jetbrains.com/dataspell/>

identificar pontos *outliers*.

Ainda, o número de anomalias consequentes usadas para criar novos *clusters* pode ser insuficiente para estimar os parâmetros do modelo inicial, prejudicando as previsões iniciais. Para solucionar este problema de conjuntos amostrais iniciais muito pequenos em modelos autoregressivos com quebra na série, [Pesaran e Timmermann \(2005\)](#) demonstraram que a utilização de 10 a 20 dados pré-quebra melhoram a performance inicial do modelo.

Em ambas as simulações, os modelos serão testados de forma evolutiva, fazendo previsões um passo a frente e estimando os parâmetros dos modelos com os dados disponíveis até o tempo de cada iteração. Baseado no horizonte de memória, oriundo do fator de esquecimento utilizado no *OEC*, os modelos locais terão uma janela móvel de observação, como visto na Fig. 10.

Figura 10 – Janela de Observação móvel baseado no horizonte de memória.



Fonte: O Autor.

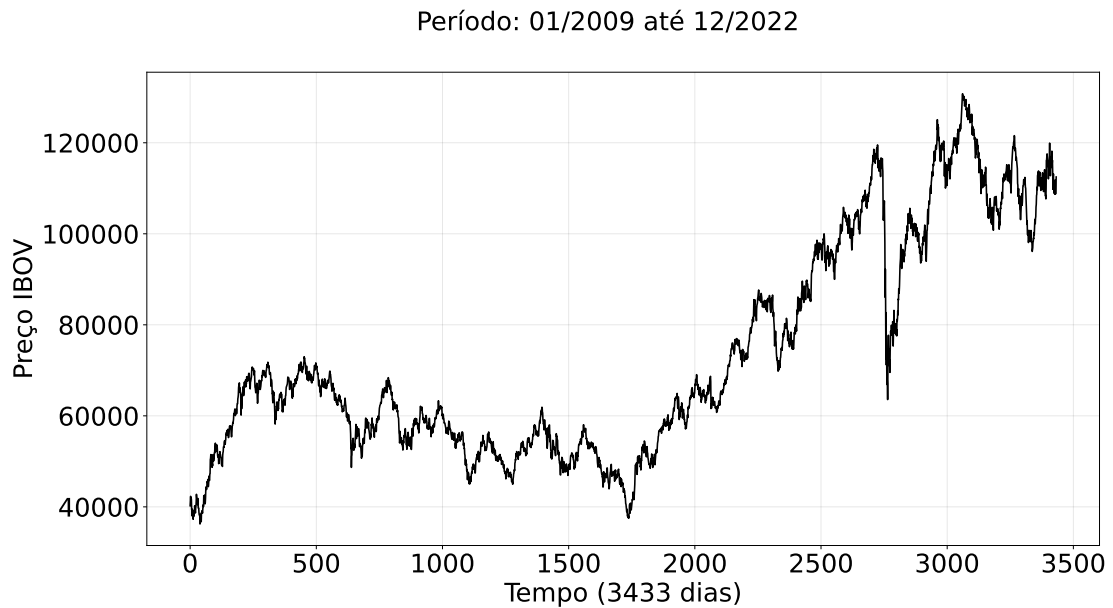
## Índice Ibovespa B3

Apresentamos a série do preço de fechamento diário do índice IBOV na Fig. 11, onde temos um total de 3433 observações coletadas, entre as datas de 01/2009 a 12/2022 com frequência diária. Inicialmente, pode-se reconhecer a natureza não-linear e caótica da série de forma visual. Assim, como ambas as simulações tratam da mesma série, o teste de *Dickey-Fuller aumentado* é realizado de forma prévia a fim de verificar a condição de não estacionariedade da mesma.

Executando o teste ADF inicialmente com a série original, resultando em um *P-value* de 0,534 para significância de 5%, a hipótese nula é aceita, logo, a série não é estacionária. Utilizando do método de diferenciação, apresentado por [Kwiatkowski et al. \(1992\)](#), para tornar a série estacionária, um *P-value* de  $2,25^{-28}$  para significância de 5% é alcançado, assim, a série diferenciada de um termo pode ser considerada estacionária.

Em termos práticos, a série diferenciada de um termo representa a variação entre o preço de fechamento de uma observação no tempo em  $k + 1$  e  $k$ .

Figura 11 – Série do preço de fechamento do índice IBOV.



Fonte: O Autor.

## 4.1 Simulação I

As variáveis utilizadas nesta simulação são descritas na Tabela 2.

<b>Descrição das Variáveis Endógenas e Exógenas</b>	
Variáveis	Informações
B3	Série do preço de fechamento do índice IBOV.
USD	Série do preço de fechamento do câmbio.

Tabela 2 – Variáveis utilizadas para Simulação I.

### 4.1.1 Caracterização dos Dados

A partir dos dados das variáveis citadas, pode-se gerar as regras e identificar as regiões locais na série, com aplicação do algoritmo OEC, assim viabilizando a análise para cada regra identificada, com auxílio dos gráficos FAC e FACP. Foram identificadas 4 regiões ao todo durante a série, em que cada regra será caracterizada quanto ao tamanho dos conjuntos de dados, a representatividade de cada regra e os possíveis parâmetros a serem utilizados nos modelos locais, através das análises gráficas.

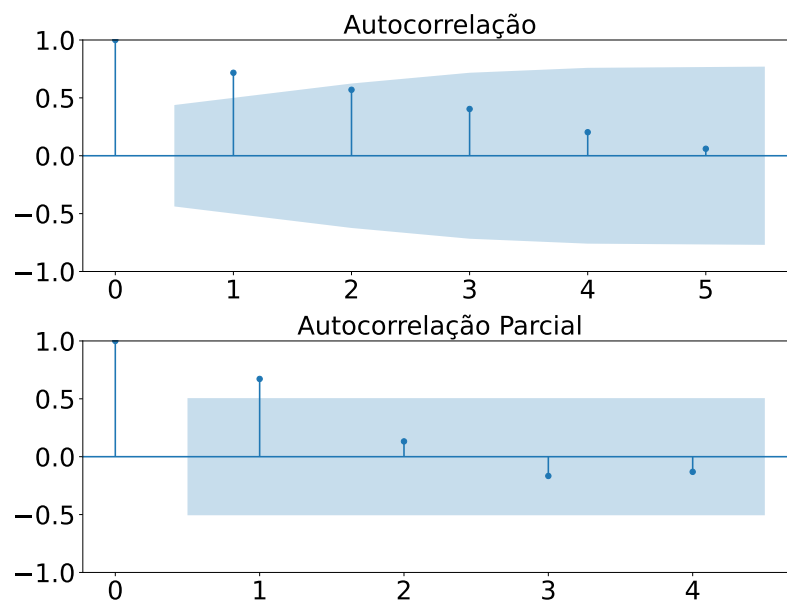
Regra 1:

Ao trabalhar com a primeira regra identificada, consideramos a janela de observação inicial de dados sendo equivalente ao horizonte de memória, 100 observações, tendo em

vista que não há conhecimento prévio em relação aos dados. Foram identificados 2225 observações atribuídas a esta regra, entretanto, somente os dados pertencentes a janela de observação inicial são utilizados para a análise inicial.

Pode-se observar, na Fig. 12, os gráficos de autocorrelação da primeira regra identificada, utilizando 25 atrasos.

Figura 12 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 1.



Fonte: O Autor.

A partir disto, pode-se inferir os parâmetros do modelo a serem utilizados:

- $p$  (Ordem da parte autoregressiva): 8;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

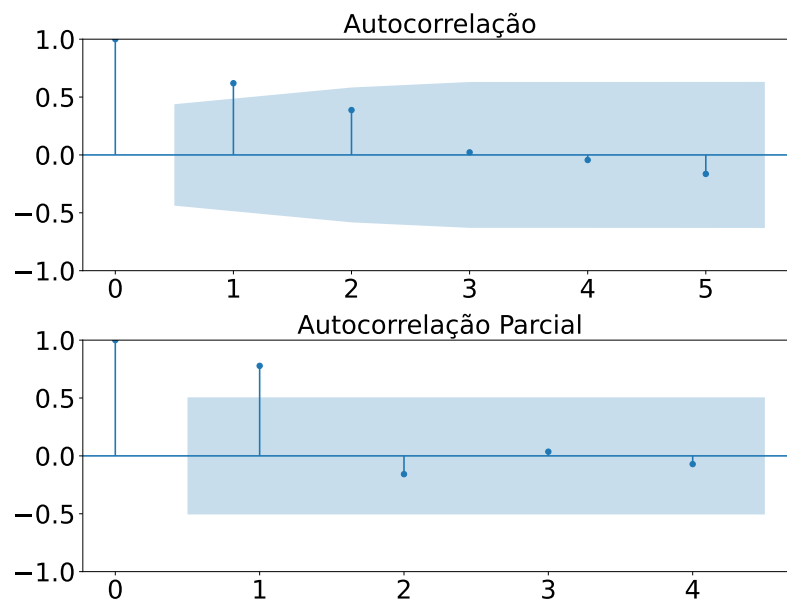
Regra 2:

A partir da segunda regra identificada, a janela de observação inicial é representada pela quantidade de anomalias consequentes observadas que atenderam a condição de criação de uma nova regra.

Esta janela inicial contém 15 observações, e ao total 100 pontos foram atribuídos a esta regra. Assim somente é necessário a análise dos dados iniciais.

Pode-se observar, na Fig. 13, os gráficos de autocorrelação da segunda regra identificada, utilizando 5 e 4 atrasos.

Figura 13 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 2.



Fonte: O Autor.

Assim, pode-se inferir os parâmetros do modelo a serem utilizados:

- $p$  (Ordem da parte autoregressiva): 2;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

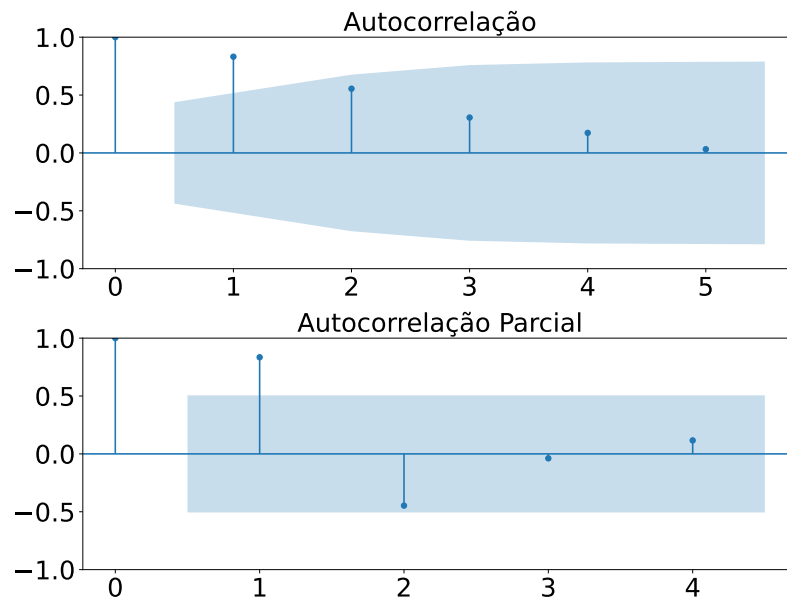
Regra 3:

A janela de observação inicial da terceira regra identificada contém 10 observações. Ao total, 736 dados foram atribuídos a essa regra, assim a análise é necessária para a janela inicial de dados e também para a janela de dados móvel. A partir da Fig. 14 pode-se inferir os parâmetros do modelo inicial, temos:

- $p$  (Ordem da parte autoregressiva): 1;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

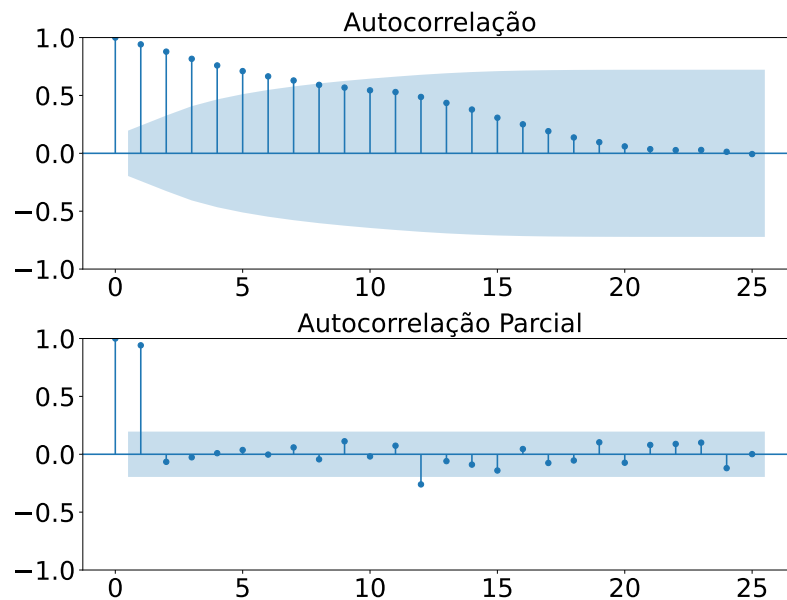
Em sequência, com a janela móvel de dados sendo atingida, infere-se os parâmetros do modelo novamente, com auxílio dos gráficos de autocorrelação, ver Fig. 15:

Figura 14 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 3.



Fonte: O Autor.

Figura 15 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 3.



Fonte: O Autor.

- $p$  (Ordem da parte autoregressiva): 7;
- $d$  (Ordem de diferenciação): 1;

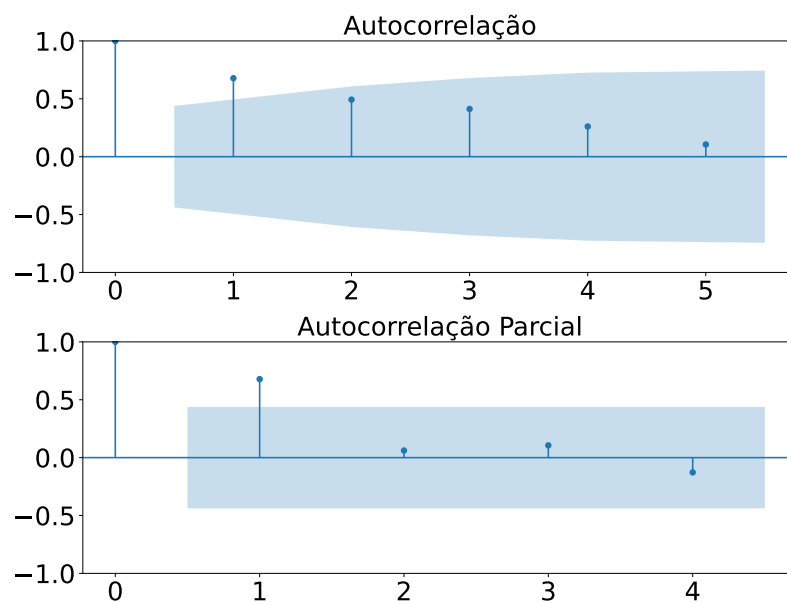
- $q$  (Ordem da parte de média móvel): 1.

#### Regra 4:

A janela de observação inicial da quarta regra identificada, oriunda de anomalias consequentes, contém 10 observações. Ao total, 372 dados foram atribuídos a essa regra, assim a análise é necessária para a janela inicial de dados e também para a janela de dados móvel. Observa-se os gráficos de autocorrelação da janela inicial pertencente a segunda regra identificada, ver Fig. 16. Inferindo os parâmetros do modelo inicial, temos:

- $p$  (Ordem da parte autoregressiva): 1;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Figura 16 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 4.

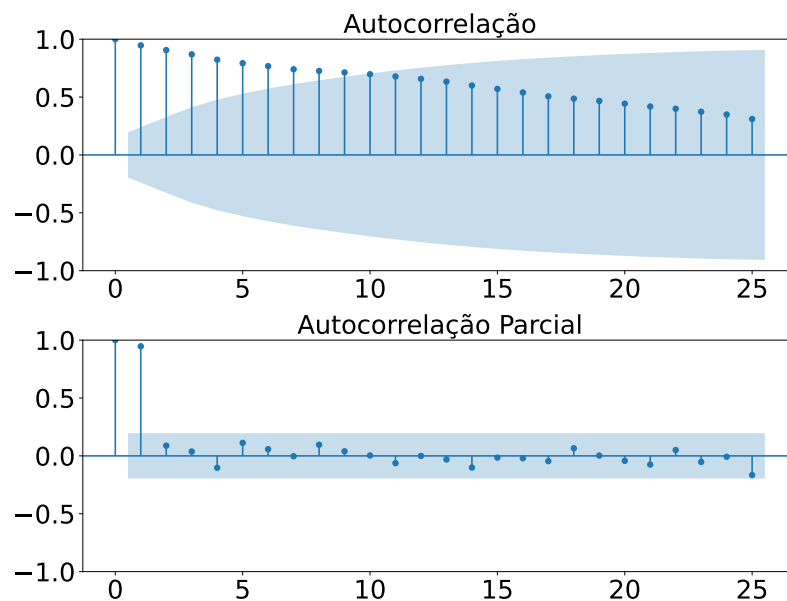


Fonte: O Autor.

Em sequência, com a janela móvel de dados sendo atingida, infere-se os parâmetros do modelo novamente, com auxílio dos gráficos de autocorrelação, ver Fig. 17:

- $p$  (Ordem da parte autoregressiva): 7;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Figura 17 – Simulação I - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 4.



Fonte: O Autor.

#### 4.1.2 Resultados

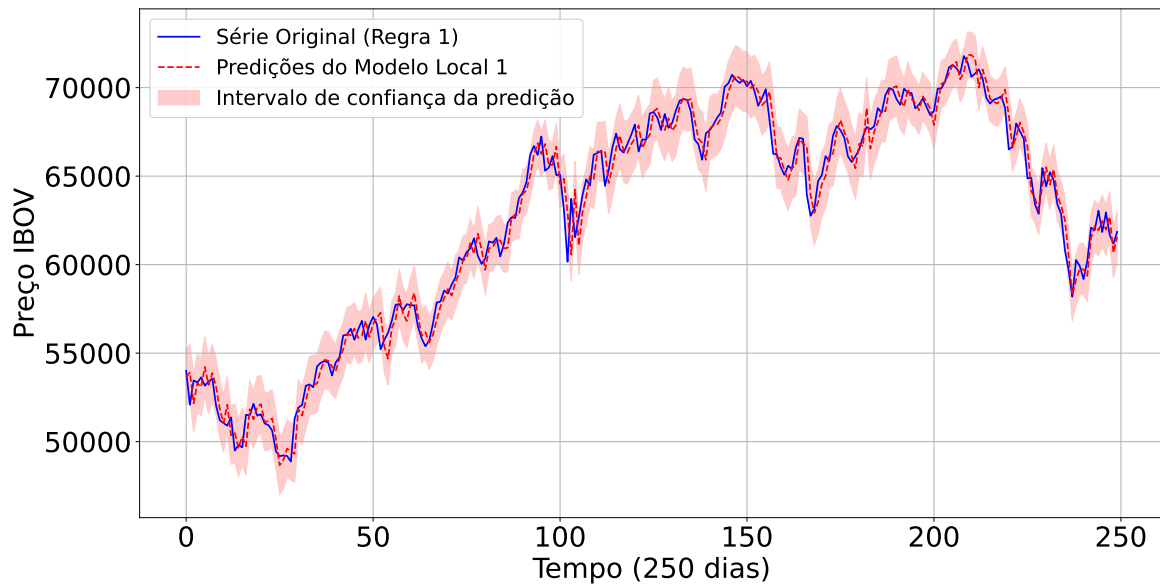
Efetuando-se a análise exploratória inicial quanto a condição de estacionariedade, foi possível verificar que, após uma diferenciação, a mesma tornou-se estacionária de acordo com o teste de *Dickey-Fuller aumentado*. Assim, foram utilizados de gráficos de autocorrelação e autocorrelação parcial para escolha inicial dos parâmetros de cada modelo, podendo testá-los de forma *out-of-sample*.

Com o intuito inicial de identificar regiões com dinâmicas similares, a utilização de uma janela de dados de tamanho fixo não ocasionará em perda de representatividade, assumindo que o algoritmo recebeu informações suficientes, evitando assim a geração de regras incompletas ou incorretas.

Traz-se atenção a utilização do conceito de memória de horizonte para a escolha inicial dos parâmetros, fazendo com que os modelos que representem regiões com um conjunto de dados maior que o horizonte de memória tenham a necessidade de uma estipulação inicial dos parâmetros com os dados iniciais e uma segunda estipulação de parâmetros, posteriormente. O modelo inicial contém poucos parâmetros e poucos dados para treiná-lo, tornando-o mais simples e possivelmente com menor acurácia em seus resultados.

Dessa forma, o primeiro modelo local estipulado teve como parâmetros  $(8,1,1)$ , onde é possível observar as predições efetuadas contra os dados reais na Fig. 18.

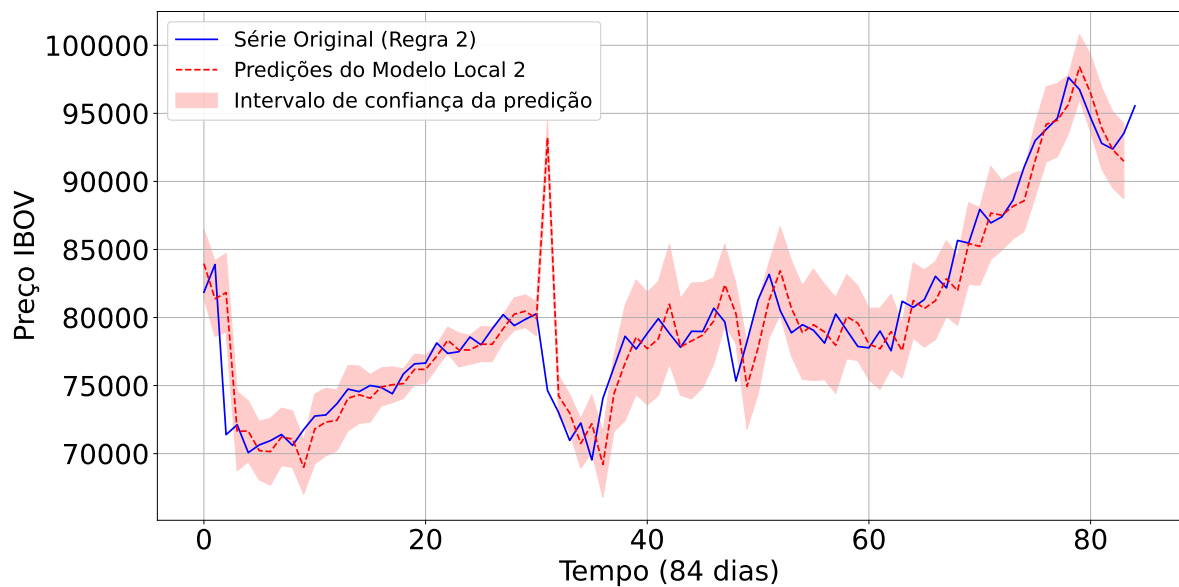
Figura 18 – Simulação I - Gráfico da série original e a predição para a região pertencente a Regra 1.



Fonte: O Autor.

Em caráter de exceção, o segundo modelo local que não atingiu a quantidade de dados suficiente para utilização da janela de observação móvel de dados. Com isso, somente um modelo inicial com parâmetros  $(2,1,1)$  foi utilizado.

Figura 19 – Simulação I - Gráfico da série original e a predição para a região pertencente a Regra 2.



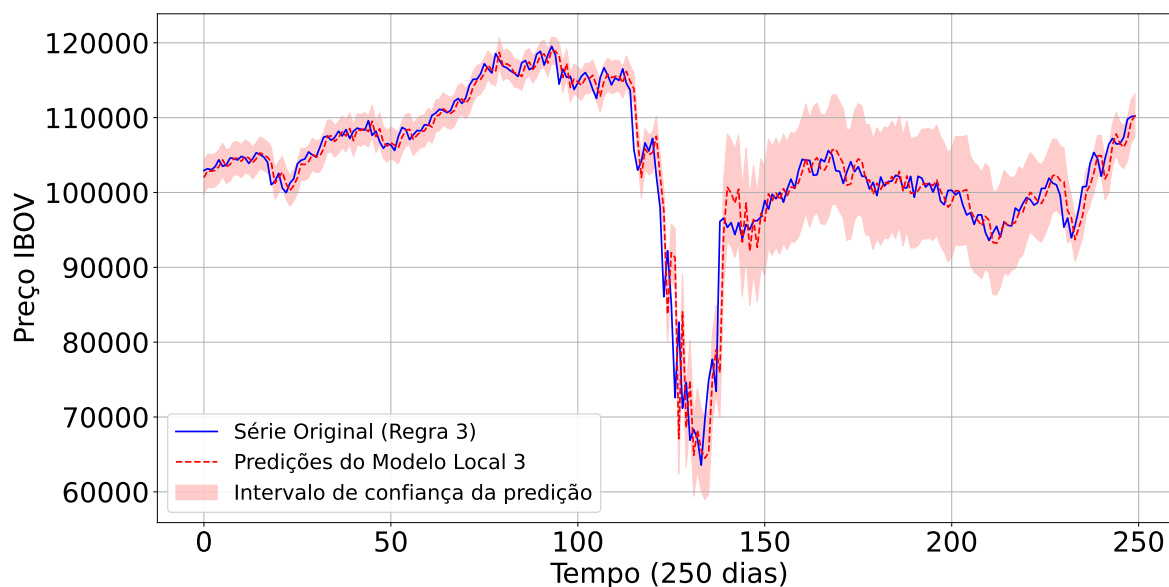
Fonte: O Autor.

Sendo um modelo de poucos parâmetros e baixa complexidade, acarretando resultados inadequados, apresentou alta incerteza em suas previsões, podendo ser visualizada com intervalos de confiança relativamente grandes em relação a média, além de uma predição significativamente desconectada da dinâmica, ver Fig. 19.

Quanto aos gráficos de previsões dos modelos locais contra os dados reais, destaca-se os intervalos de confiança de cada predição plotado juntamente aos gráficos trazendo uma melhor visualização quanto a incerteza do resultado, onde quanto maior o nível de confiança maior tendem a ser os intervalos de confiança (SILVA; NAGHETTINI; PORTELA, 2011). O nível de confiança utilizado para a elaboração da pesquisa é de 95%, assim, pode-se interpretar que a probabilidade do intervalo de confiança conter o valor real da variável é 95%.

Igualmente, o terceiro modelo local traz uma alta variabilidade, de acordo com seus intervalos de confiança, onde em um certo período apresenta intervalos de confiança com amplitudes baixas e posteriormente aumenta drasticamente. Essa variabilidade pode ser vista na Fig. 20, onde há um evento de mudança brusca de amplitude da série em que o modelo não consegue acompanhar a dinâmica de modo apropriado.

Figura 20 – Simulação I - Gráfico da série original e a predição para a região pertencente a Regra 3.



Fonte: O Autor.

No terceiro modelo, inicialmente foram utilizados os parâmetros estimados (2,1,1) e posteriormente, com a janela móvel de dados, os parâmetros (7,1,1). Verificou-se que dos 8 termos autoregressivos utilizados no modelo, somente os 5 primeiros tinham significância estatística suficiente, assim um modelo com os parâmetros (5,1,1) se torna mais adequado.

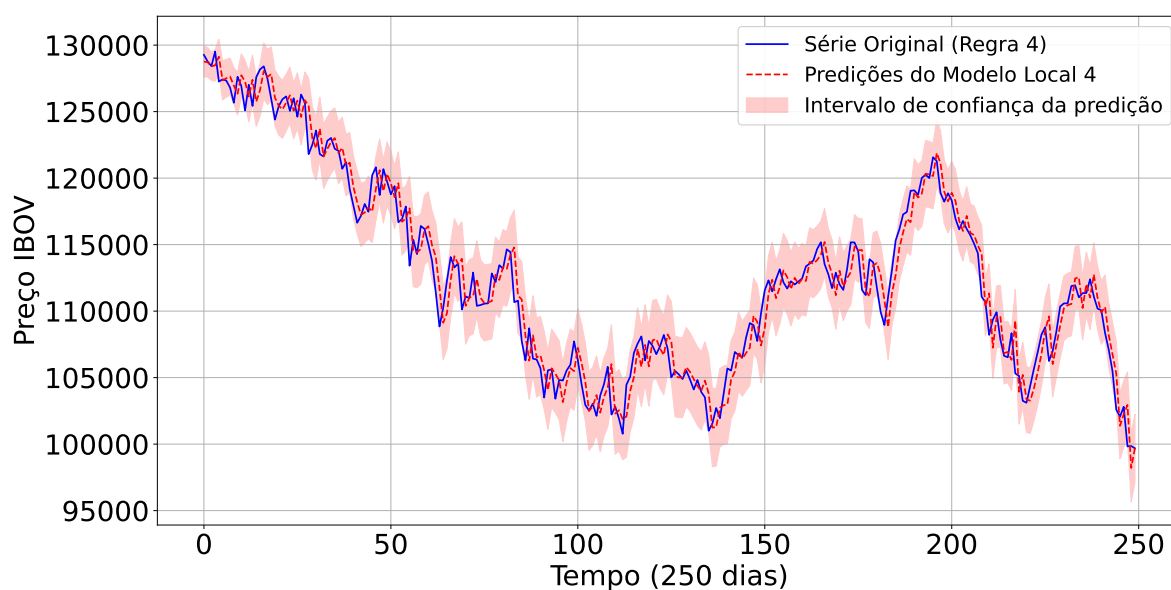
Com o advento da globalização econômica, ou seja, os entrelaços financeiros que ocorrem entre os diferentes países, suas economias locais são afetadas por eventos externos fazendo com que os mercados financeiros de cada país apresentem trajetórias similares.

O evento ocorrido deu início em dezembro de 2019 com o surgimento de casos de pneumonia em Wuhan, província chinesa de Hubei. Posteriormente, o vírus *SARS-CoV-2* foi identificado como a causa dos casos de pneumonia, que logo se espalharam por outros países. Assim, o governo brasileiro impôs o fechamento de estabelecimentos públicos e privados, além do isolamento em domicílio da população para contenção do vírus (FEYISA, 2020).

Com a crise sanitária da COVID-19 e o fechamento dos comércios locais, o mercado financeiro reagiu de forma extrema, ocasionando o acionamento do *Circuit breaker* diversas vezes. O *Circuit breaker* é uma ferramenta de proteção contra quedas abruptas da bolsa de valores, em que as negociações de ativos são interrompidas quando há uma queda maior que 10%. Ao todo, no mês de março de 2020, onde o índice Ibovespa retraiu mais de 30% e o câmbio teve um aumento aproximado de 16% (FREITAS; SANTIAGO; CARVALHO, 2023).

Dado o evento causado por fatores externos não diretamente relacionados com a bolsa de valores, não foi possível identificar a mudança de dinâmica que ocorreu, assim sugerindo que não houve informações suficientes para que a mudança de momento do mercado fosse identificada, resultando em um dimensionamento insatisfatório das regras.

Figura 21 – Simulação I - Gráfico da série original e a predição para a região pertencente a Regra 4.



Fonte: O Autor.

Assim, no quarto modelo foram estipulados os parâmetros iniciais e os posteriores a janela móvel de dados, (1,1,0) e (9,1,1), respectivamente. É possível observar as previsões efetuadas contra os dados reais na Fig. 21.

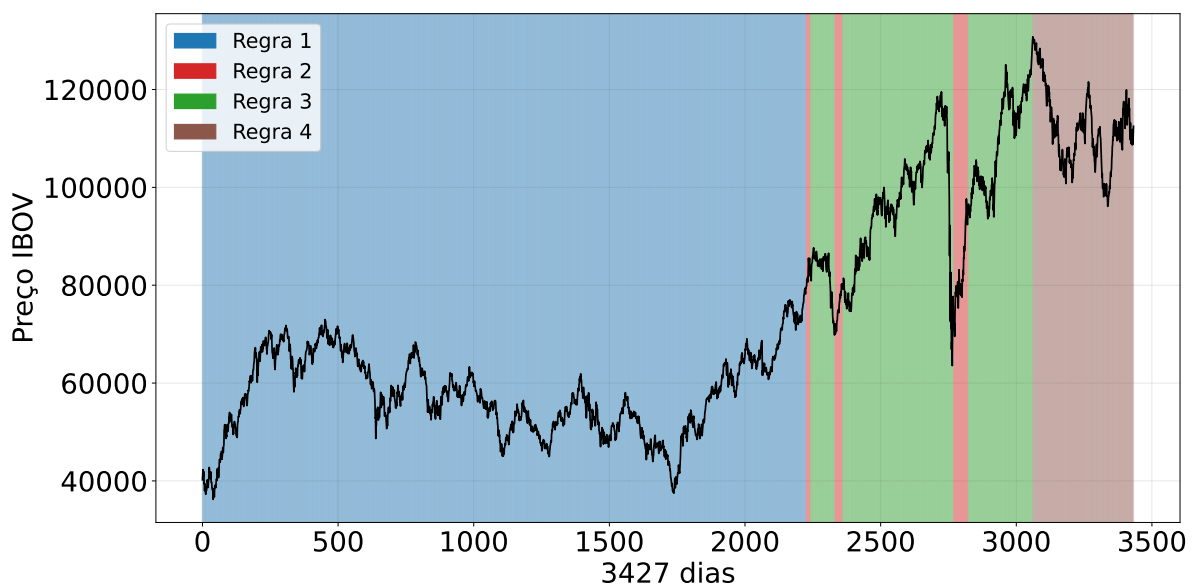
Além disso, a Tab. 3 apresenta todas as métricas de avaliação para cada modelo local, auxiliadas pela visualização das mesmas em *boxplots* apresentados na Fig. 23. Dessa forma, consolidando os apontamentos realizados quanto a alta incerteza nas previsões do segundo modelo sendo confirmados pela métrica RMSE e MAE de valores altos.

Métricas para a Simulação I (Média $\pm$ Desvio Padrão)				
Modelos	MAE	RMSE	MAPE (%)	SMAPE (%)
1°	627.2 $\pm$ 96.3	809.9 $\pm$ 133.5	1.1 $\pm$ 0.2	0.6 $\pm$ 0.1
2°	1823.2 $\pm$ 587.3	3285.3 $\pm$ 725.5	2.4 $\pm$ 0.8	1.2 $\pm$ 0.4
3°	1340.7 $\pm$ 593.5	1930.4 $\pm$ 1140.1	1.4 $\pm$ 0.7	0.7 $\pm$ 0.3
4°	1149.4 $\pm$ 80.0	1464.3 $\pm$ 121.0	1.1 $\pm$ 0.1	0.5 $\pm$ 0.0

Tabela 3 – Métricas de Avaliação utilizadas na Simulação I.

Evidenciando a separação dos setores identificados, ver Fig. 22, nota-se que apesar das métricas com níveis adequados para o primeiro modelo, a primeira região contém 65% de todas as observações, levantando um questionamento quanto a qualidade da representatividade de informações e de uma mesma dinâmica se prevalecendo por um tempo tão prolongado.

Figura 22 – Série Original destacando os setores locais identificados por cada regra, referente ao experimento com 2 variáveis.



Fonte: O Autor.

Sabendo que as informações fornecidas não foram suficientes, devido a não identificação da mudança de dinâmica causada pela crise sanitária da COVID-19, pode-se supor que possam conter outras regiões insuficientemente coerentes também. Ademais, observa-se que apesar da não identificação da mudança de dinâmica citada anteriormente, a Regra 3 conseguiu identificar, de forma gráfica, setores com tendências significativas de alta.

A Tab. 4 exibe a importância normalizada de cada atributo, calculada a partir dos coeficientes dos parâmetros, entre entradas e regressores do modelo, de forma percentual, destacando os 3 mais importantes por modelo.

Atributos		Importância dos Atributos Normalizada por Modelo Local			
Tipo	Atraso/Grau	1	2	3	4
Exógena	1°	<b>46.65</b> %	<b>21.85</b> %	<b>41.64</b> %	0.26 %
AR	1°	0.50 %	<b>31.73</b> %	4.24 %	<b>34.27</b> %
AR	2°	9.24 %	3.34 %	8.05 %	0.40 %
AR	3°	9.10 %	-	<b>21.13</b> %	3.50 %
AR	4°	2.63 %	-	<b>20.07</b> %	<b>5.66</b> %
AR	5°	<b>15.70</b> %	-	2.29 %	5.29 %
AR	6°	1.40 %	-	-	5.21 %
AR	7°	<b>13.01</b> %	-	-	4.45 %
AR	8°	1.24 %	-	-	2.28 %
AR	9°	-	-	-	0.50 %
MA	1°	0.55 %	<b>43.10</b> %	2.58 %	<b>38.17</b> %

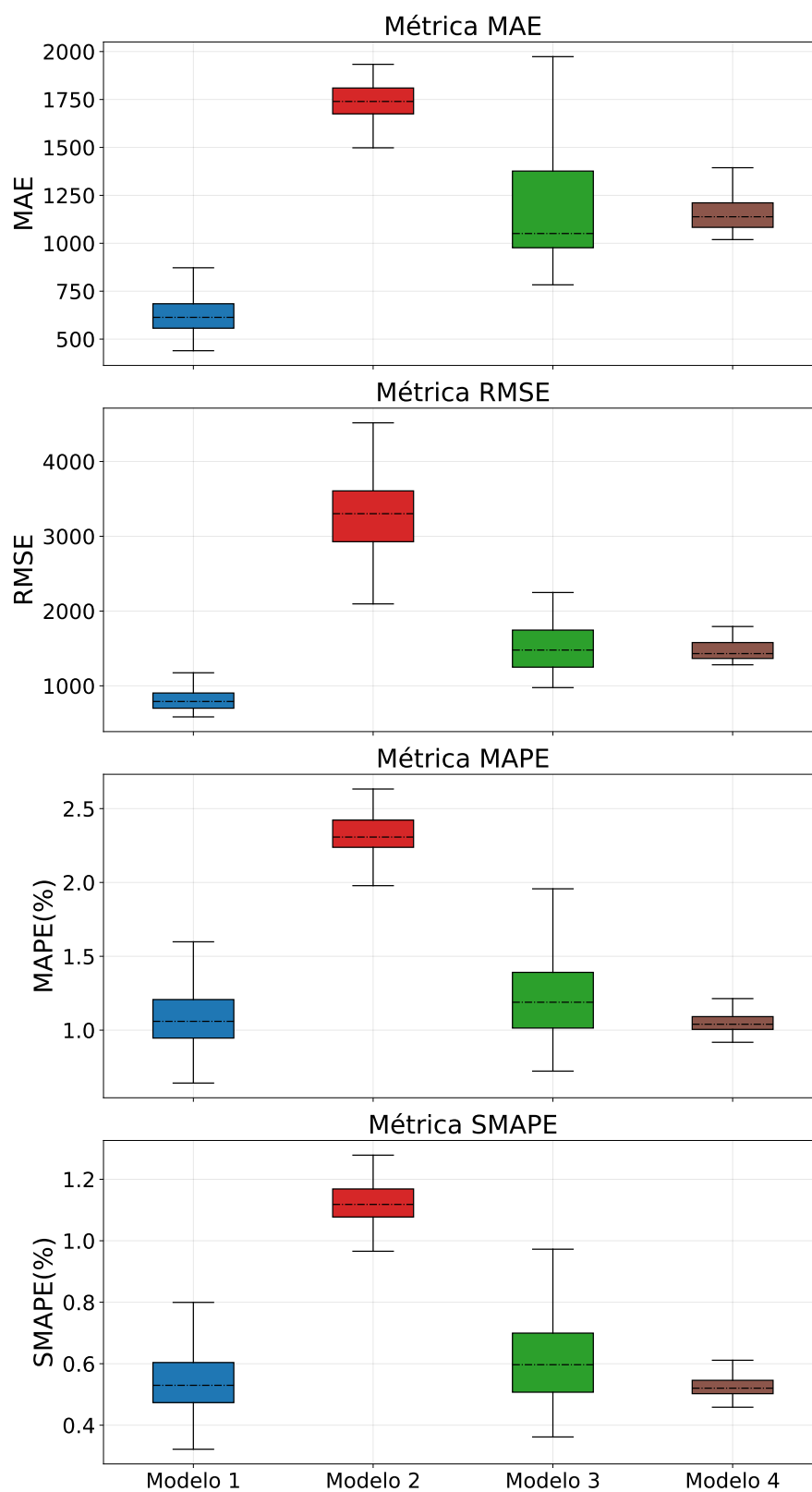
Tabela 4 – Importância Normalizada de cada atributo, por modelo, destacando os com maior pertinência para a Simulação I.

Cumprir destacar que, o câmbio (entrada exógena) apresentou alta relevância nas predições realizadas pelos 3 primeiros modelos e que a média móvel de 1 atraso apresentou alta relevância no segundo e quarto modelos, indicando o forte entrelaçamento do câmbio com as predições feitas e o vínculo da média móvel de curto prazo com as predições feitas nas regiões com quantidades menores de dados.

De outra forma, destaca-se o baixo vínculo do parâmetro de média móvel com o primeiro modelo, reforçando o questionamento quanto a uma única dinâmica prevalecer por um período de tempo extenso como visto na primeira região. Porém, os autoregressores de atraso 5 e 7 apresentaram relevância significativa para a predição, indicando que as observações de 5 e 7 dias atrás, em referência ao momento atual da predição, contribuíram fortemente para o valor predito.

Visto o vínculo da média móvel com regiões identificadas menos extensas, esse parâmetro se torna um possível catalisador para otimizar a identificação das regiões e melhorar a representatividade dos modelos propostos, dado sua natureza de comportamento de tendências.

Figura 23 – Métricas de Avaliação referentes ao experimento *out-of-sample* com 2 variáveis, resultando em 4 modelos locais, dispostas em *boxplots*.



Fonte: O Autor.

## 4.2 Simulação II

Tendo em vista os resultados da primeira simulação, a relação de médias móveis com regiões menos extensas se apresentou como um bom catalisador para otimizar a identificação das regiões com diferentes dinâmicas. Assim, duas séries de média móvel foram acrescentadas como variáveis exógenas, com o objetivo de identificar de forma mais adequada as diferentes dinâmicas do IBOV.

As variáveis utilizadas nesta simulação são descritas na Tabela 5.

Descrição das Variáveis Endógenas e Exógenas	
Variáveis	Informações
B3	Série do preço de fechamento do índice IBOV.
USD	Série do preço de fechamento do câmbio.
MA <sub>3</sub>	Série da média móvel de 3 dias do preço de fechamento.
MA <sub>7</sub>	Série da média móvel de 7 dias do preço de fechamento.

Tabela 5 – Variáveis utilizadas para Simulação II.

### 4.2.1 Caracterização dos Dados

A partir dos dados das variáveis citadas, pode-se gerar as regras e identificar as regiões locais na série, com aplicação do algoritmo OEC, assim viabilizando a análise para cada regra identificada, com auxílio dos gráficos FAC e FACP. Foram identificadas 6 regiões ao todo durante a série, em que cada regra será caracterizada quanto ao tamanho dos conjuntos de dados, a representatividade de cada regra e os possíveis parâmetros a serem utilizados nos modelos locais, através das análises gráficas.

Regra 1:

Ao trabalhar com a primeira regra identificada, consideramos a janela de observação inicial de dados sendo equivalente ao horizonte de memória, 100 observações, tendo em vista que não há conhecimento prévio em relação aos dados.

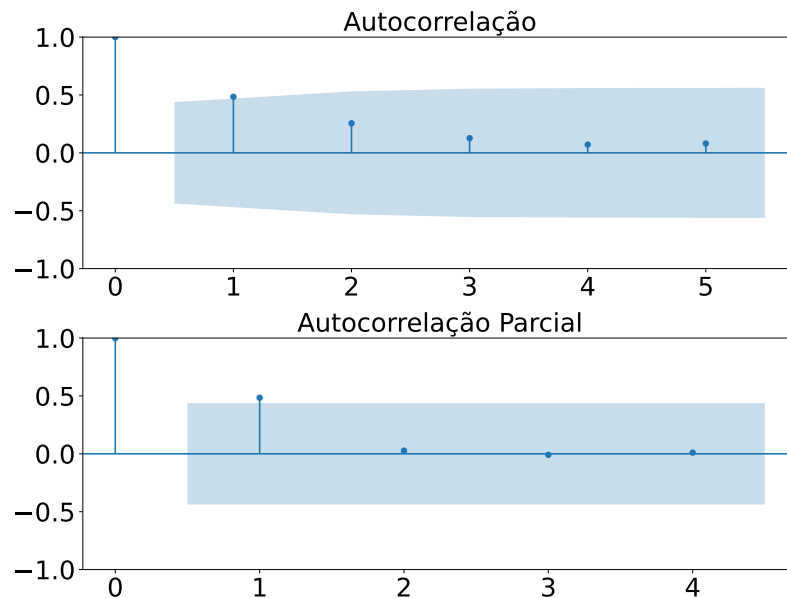
Foram identificados 1225 observações atribuídas a esta regra, entretanto, somente os dados pertencentes a janela de observação inicial são utilizados para a análise inicial. Pode-se observar, na Fig. 24, os gráficos de autocorrelação da primeira regra identificada, utilizando 25 atrasos.

Dessa forma, infere-se os seguintes parâmetros:

- $p$  (Ordem da parte autoregressiva): 8;
- $d$  (Ordem de diferenciação): 1;

- $q$  (Ordem da parte de média móvel): 1.

Figura 24 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 1.



Fonte: O Autor.

#### Regra 2:

A partir da segunda regra identificada, a janela de observação inicial é representada pela quantidade de anomalias consequentes utilizada para criação da nova regra. A janela de inicial contém 10 observações. Ao total, 700 dados foram atribuídos a essa regra, assim a análise é necessária para a janela inicial e a janela de dados móvel.

A partir da Fig. 25 pode-se inferir os parâmetros do modelo inicial, temos:

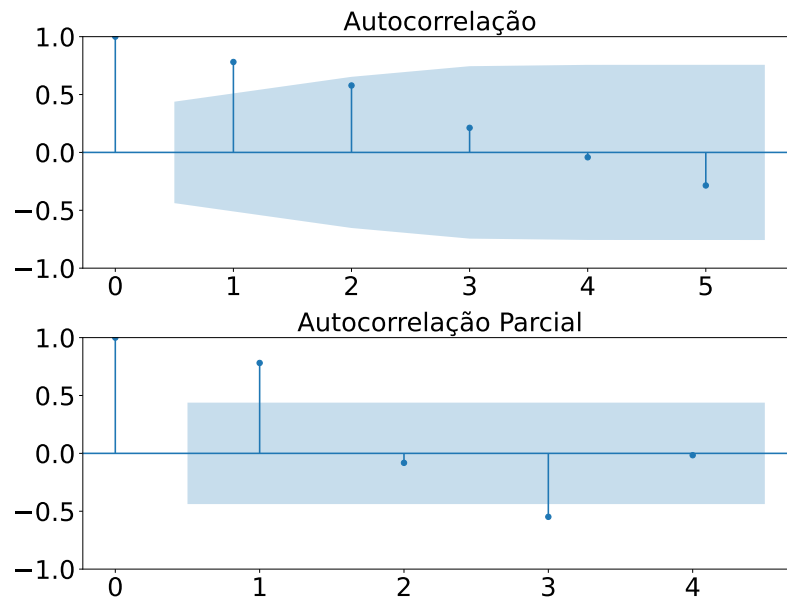
- $p$  (Ordem da parte autoregressiva): 1;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Em sequência, com a janela móvel de dados preenchida, infere-se os parâmetros do modelo novamente, ver Fig. 26:

- $p$  (Ordem da parte autoregressiva): 7;
- $d$  (Ordem de diferenciação): 1;

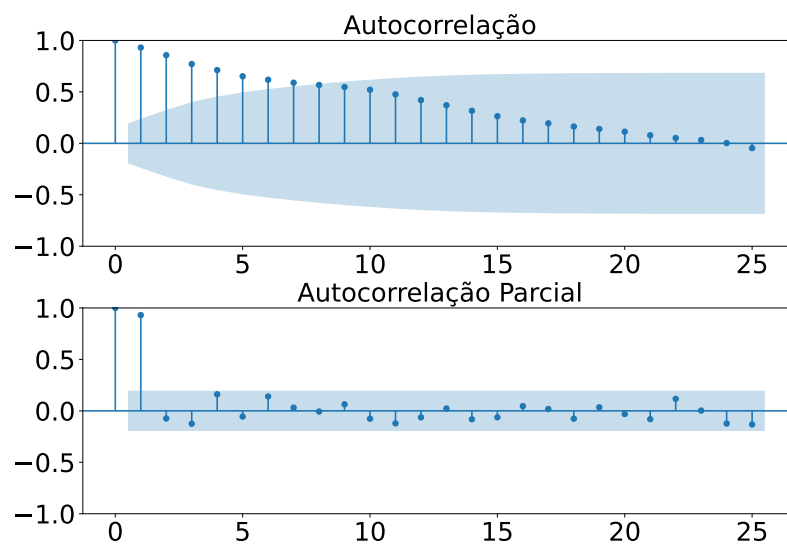
- $q$  (Ordem da parte de média móvel): 1.

Figura 25 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 2.



Fonte: O Autor.

Figura 26 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 2.



Fonte: O Autor.

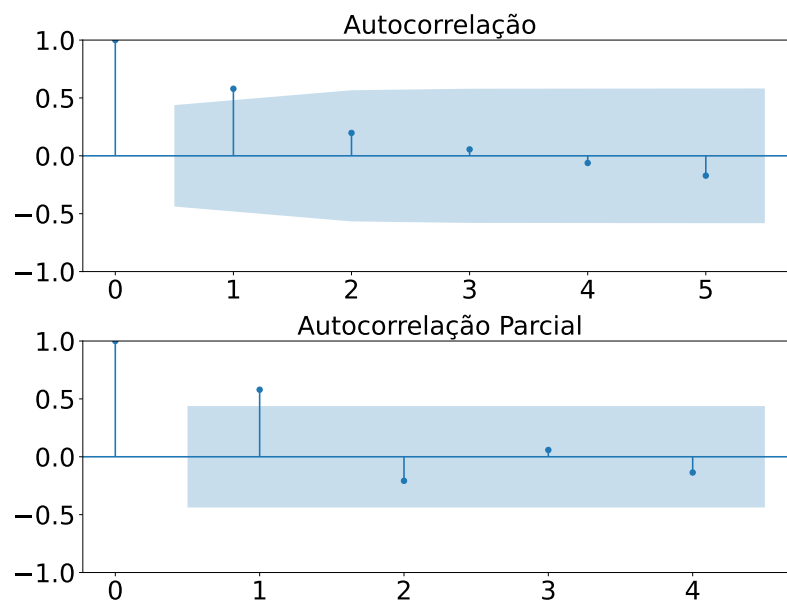
## Regra 3:

A janela de observação inicial da terceira regra identificada contém 10 observações. Ao total, 177 dados foram atribuídos a essa regra, assim a análise é necessária para a janela inicial de dados e também para a janela de dados móvel.

A partir da Fig. 27 pode-se inferir os parâmetros do modelo inicial, temos:

- $p$  (Ordem da parte autoregressiva): 1;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Figura 27 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 3.

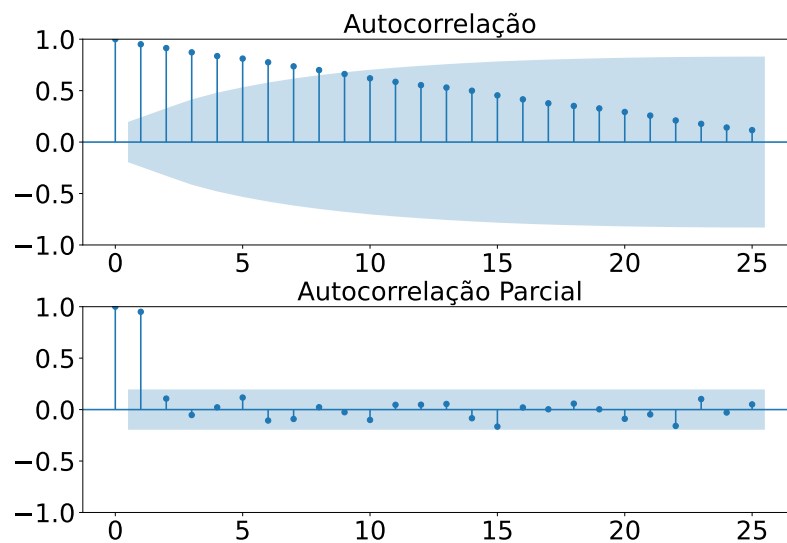


Fonte: O Autor.

Em sequência, com a janela móvel de dados sendo atingida, infere-se os parâmetros do modelo novamente, com auxílio dos gráficos de autocorrelação, ver Fig. 28:

- $p$  (Ordem da parte autoregressiva): 8;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Figura 28 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 3.



Fonte: O Autor.

#### Regra 4:

A janela de observação inicial da terceira regra identificada contém 10 observações. Ao total, 189 dados foram atribuídos a essa regra, assim a análise é necessária para a janela inicial de dados e também para a janela de dados móvel.

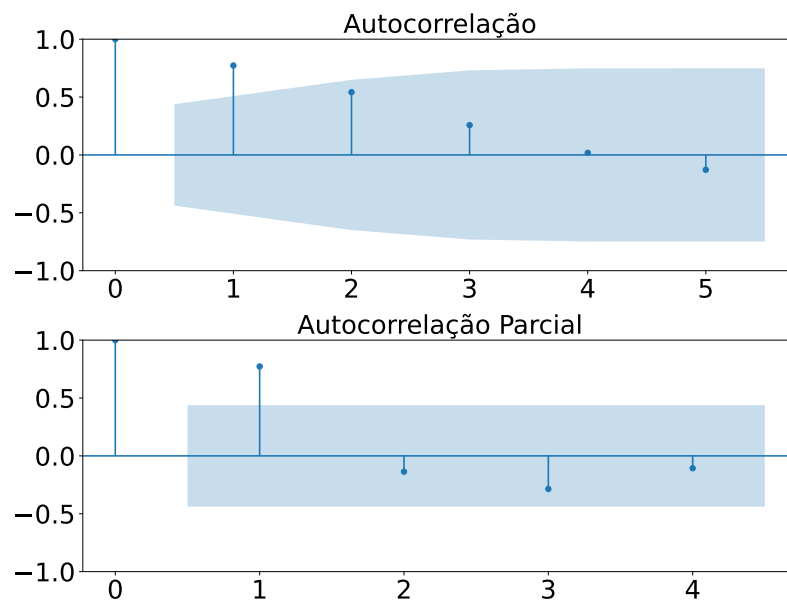
A partir da Fig. 29 pode-se inferir os parâmetros do modelo inicial, temos:

- $p$  (Ordem da parte autoregressiva): 1;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Em sequência, com a janela móvel de dados sendo atingida, infere-se os parâmetros do modelo novamente, com auxílio dos gráficos de autocorrelação, ver Fig. 30:

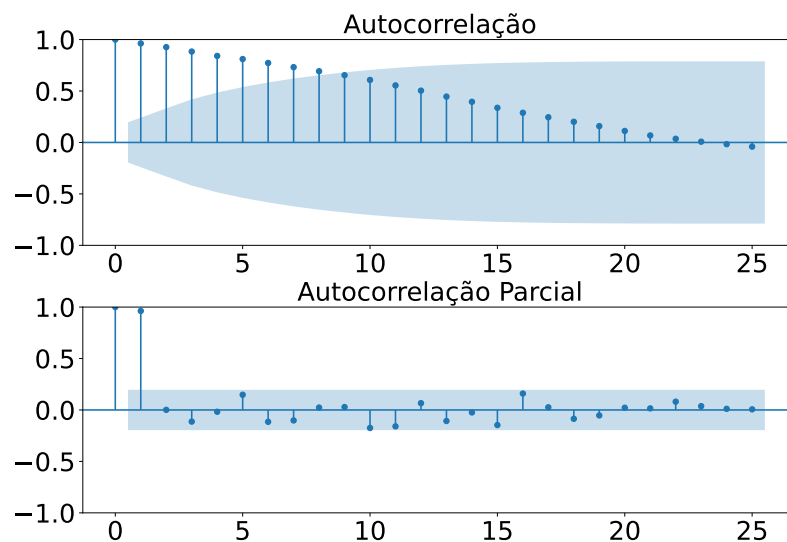
- $p$  (Ordem da parte autoregressiva): 7;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Figura 29 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 4.



Fonte: O Autor.

Figura 30 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 4.



Fonte: O Autor.

Regra 5:

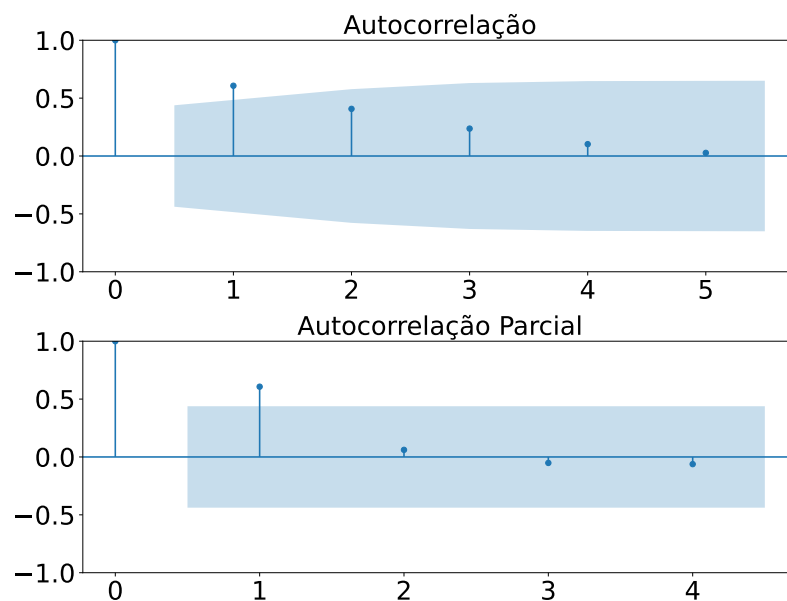
A janela de observação inicial da terceira regra identificada contém 13 observações. Ao total, 443 dados foram atribuídos a essa regra, assim a análise é necessária para a

janela inicial de dados e também para a janela de dados móvel.

A partir da Fig. 31 pode-se inferir os parâmetros do modelo inicial, temos:

- $p$  (Ordem da parte autoregressiva): 1;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Figura 31 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 5.



Fonte: O Autor.

Em sequência, com a janela móvel de dados sendo atingida, infere-se os parâmetros do modelo novamente, com auxílio dos gráficos de autocorrelação, ver Fig. 32:

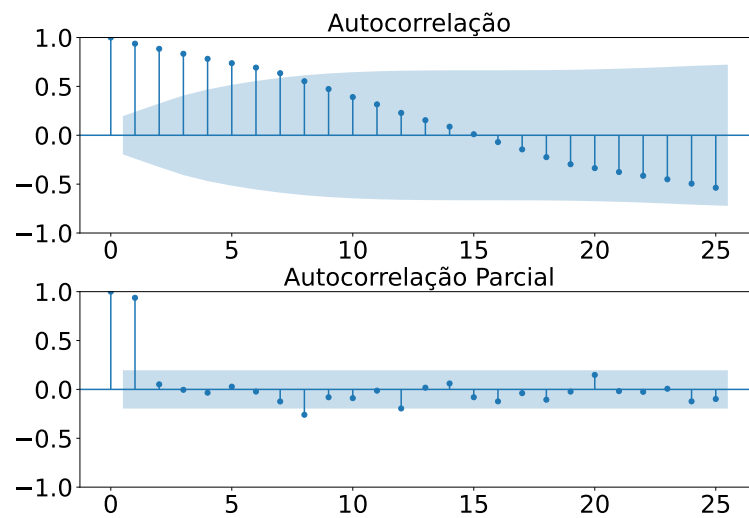
- $p$  (Ordem da parte autoregressiva): 7;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Regra 6:

A janela de observação inicial da terceira regra identificada contém 10 observações. Ao total, 693 dados foram atribuídos a essa regra, assim duas análises são necessárias.

A partir da Fig. 33 pode-se inferir os parâmetros do modelo inicial, temos:

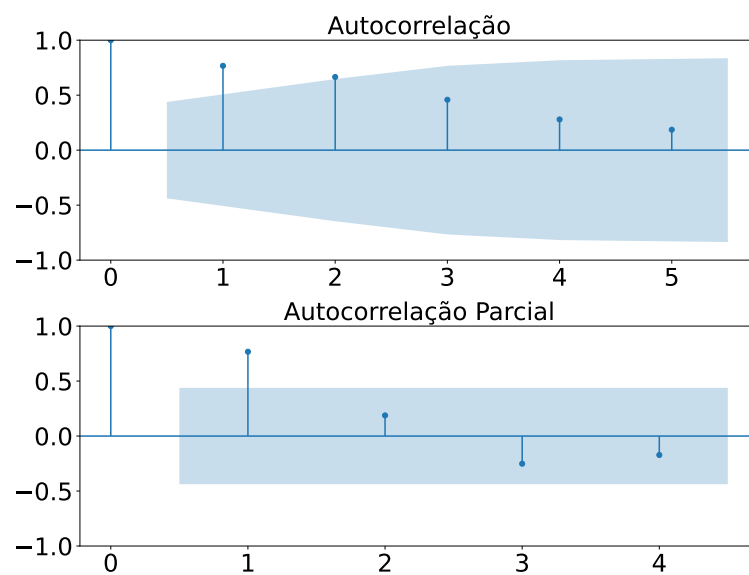
Figura 32 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 5.



Fonte: O Autor.

- $p$  (Ordem da parte autoregressiva): 2;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Figura 33 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados iniciais pertencentes ao Cluster 6.

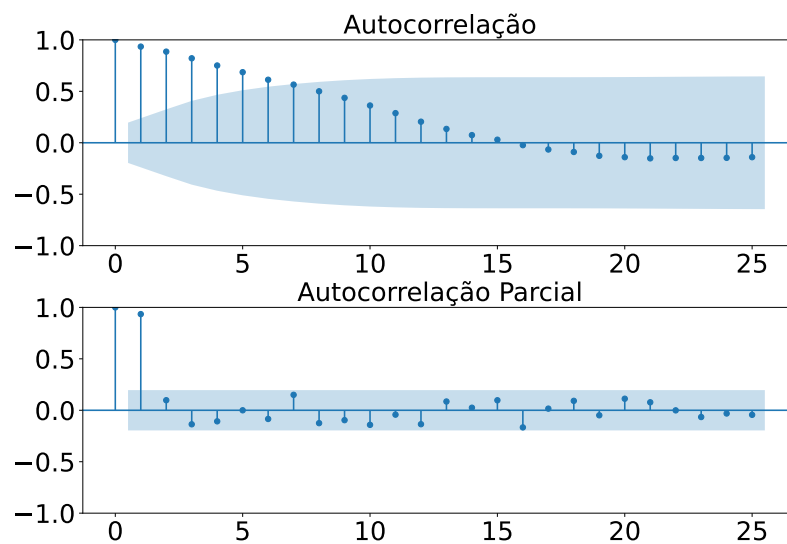


Fonte: O Autor.

Em sequência, com a janela móvel de dados sendo atingida, infere-se os parâmetros do modelo novamente, com auxílio dos gráficos de autocorrelação, ver Fig. 34:

- $p$  (Ordem da parte autoregressiva): 7;
- $d$  (Ordem de diferenciação): 1;
- $q$  (Ordem da parte de média móvel): 1.

Figura 34 – Simulação II - Autocorrelação e Autocorrelação Parcial dos dados no horizonte de memória pertencentes ao Cluster 6.



Fonte: O Autor.

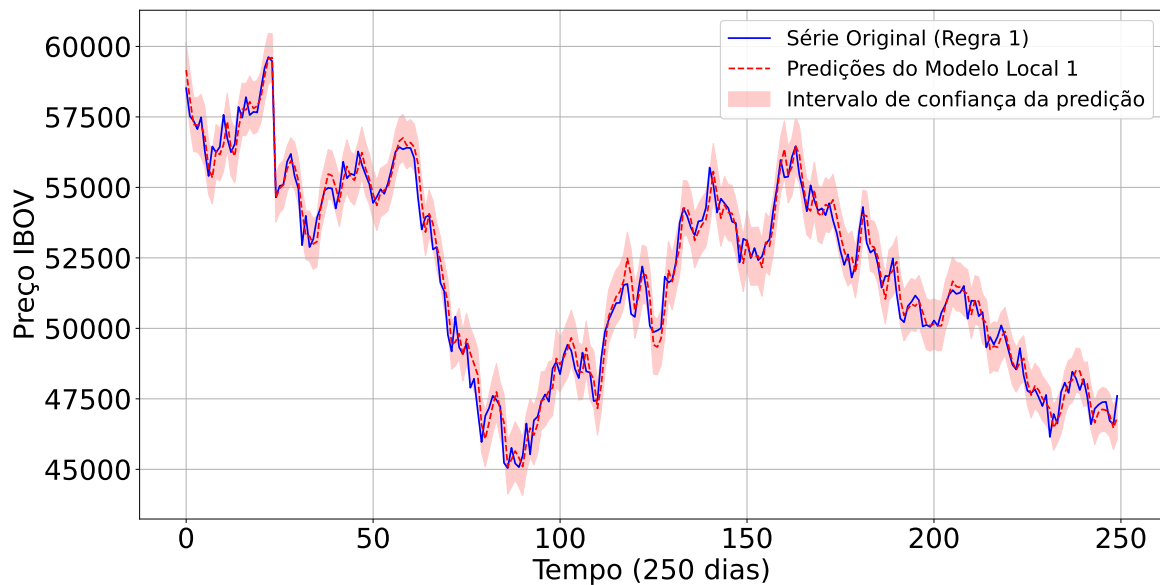
## 4.2.2 Resultados

Salienta-se que a série foi considerada estacionária de acordo com o teste de *Dickey-fuller aumentado*, após uma diferenciação. Assim, os gráficos de autocorrelação e autocorrelação parcial foram utilizados para escolha inicial dos parâmetros de cada modelo, assim viabilizando as simulações *out-of-sample*.

Com o objetivo inicial de identificar regiões com dinâmicas similares, a utilização de uma janela de dados de tamanho fixo não comprometerá a representatividade, desde que o algoritmo receba informações adequadas, evitando, assim, a geração de regras incompletas ou incorretas. Inicialmente, observa-se que, ao utilizar as médias móveis como variáveis exógenas, duas dinâmicas novas foram identificadas, totalizando 6 dinâmicas totais. Pode-se inferir, de forma prévia, que ocorreu uma melhora na representação dos dados.

Dessa forma, o primeiro modelo local estimado a partir das análises exploratórias iniciais teve como parâmetros (9,1,1). Verificou-se que dos 9 termos autoregressivos utilizados neste modelo, somente os 4 primeiros tinham significância estatística suficiente, assim um modelo de parâmetros (4,1,1) se tornou mais adequado, ver Fig. 35.

Figura 35 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 1.



Fonte: O Autor.

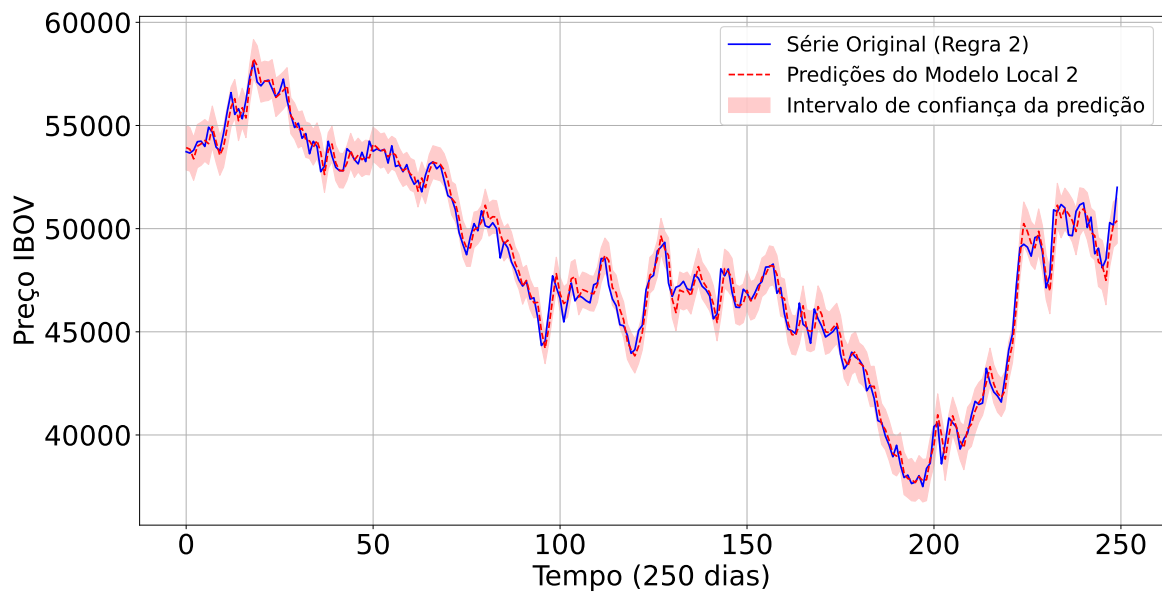
Torna-se visível no primeiro modelo local, em comparação com a simulação anterior, a melhora no desempenho das predições. A diminuição de amplitude dos intervalos de confiança, atribuindo uma menor incerteza quanto as pedições efetuadas, além de um acompanhamento da dinâmica de forma que não se apresenta defasada.

Apesar das observações preliminares, é imperativo prosseguir com as análises detalhadas das próximas regiões identificadas, a fim de garantir a precisão e a abrangência dos resultados obtidos.

No segundo modelo, inicialmente foram utilizados os parâmetros estimados (1,1,1). Posteriormente, com a aplicação de uma janela móvel de dados, foram adotados os parâmetros (8,1,1). As predições realizadas por este segundo modelo local podem ser comparadas aos dados reais, conforme ilustrado na Fig. 36.

Visualmente, observa-se uma baixa incerteza nas predições, como indicado pelos estreitos intervalos de confiança. Além disso, o modelo demonstra um excelente acompanhamento da dinâmica dos dados, sem apresentar defasagem aparente. Estas características indicam uma robustez e precisão significativas do modelo na captura das variações temporais dos dados analisados.

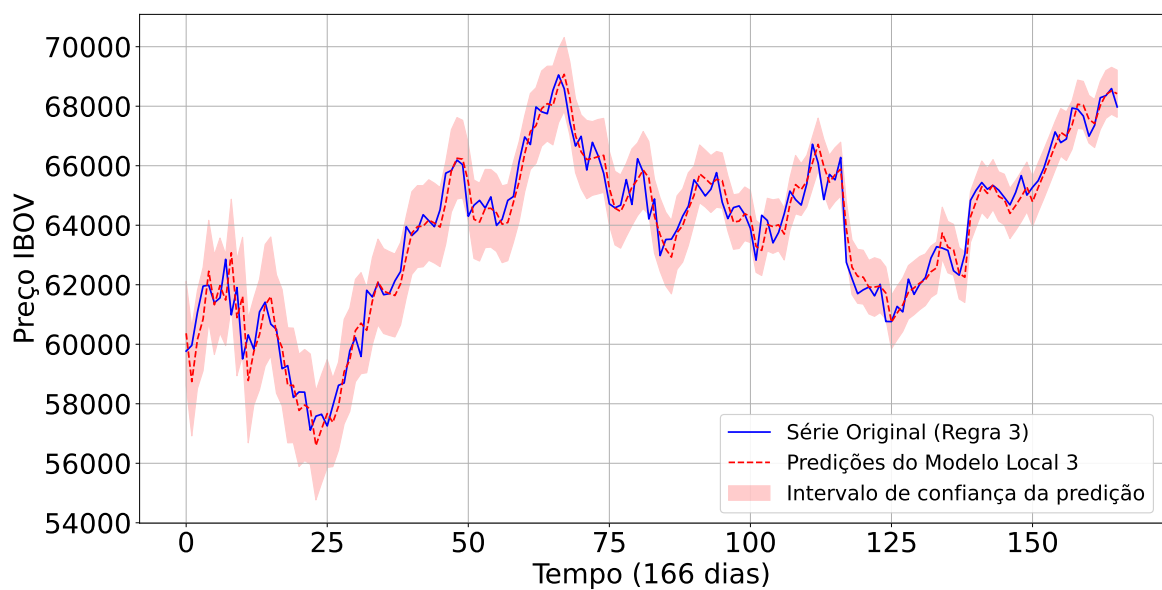
Figura 36 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 2.



Fonte: O Autor.

Já no terceiro modelo, constatou-se uma amplitude significativamente maior dos intervalos de confiança, ver Fig. 37.

Figura 37 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 3.



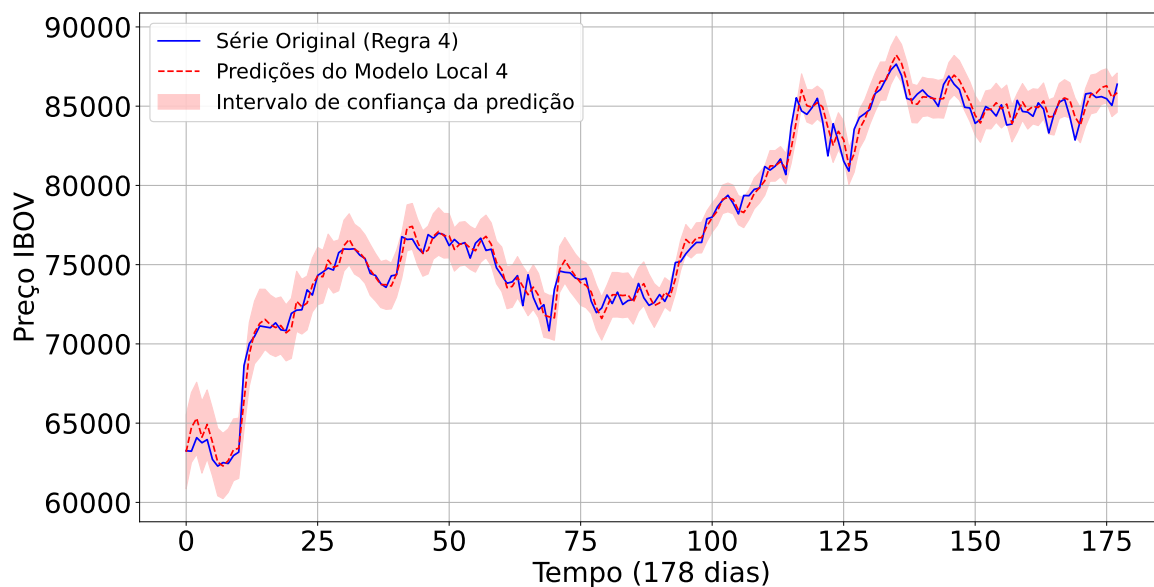
Fonte: O Autor.

Este comportamento sugere uma sensibilidade aumentada do modelo, que pode ser atribuída a quantidade de dados limitadas identificadas para esta região. Os parâmetros utilizados inicialmente e posteriormente a janela de dados móveis foram de  $(1,1,1)$  e  $(8,1,1)$ , respectivamente.

De forma semelhante, no quarto modelo, verificou-se uma incerteza relativamente maior nos resultados gerados pelo modelo inicial antes de atingir a janela móvel de dados, ver Fig. 38. Esse aumento na incerteza sugere uma sensibilidade acentuada do modelo, possivelmente devido à quantidade limitada de dados atribuída à quarta região identificada, resultando em um modelo com parâmetros mais simples. Os parâmetros utilizados inicialmente e posteriormente a janela de dados móveis foram de  $(1,1,1)$  e  $(7,1,1)$ , respectivamente.

Para tentar mitigar essa incerteza, pode-se acrescentar novas características à região, garantindo que o modelo possa capturar adequadamente as dinâmicas presentes, assim como feito nesta segunda simulação, dado os resultados da primeira simulação executada. Entretanto, o modelo ainda foi capaz de acompanhar as variações temporais dos dados.

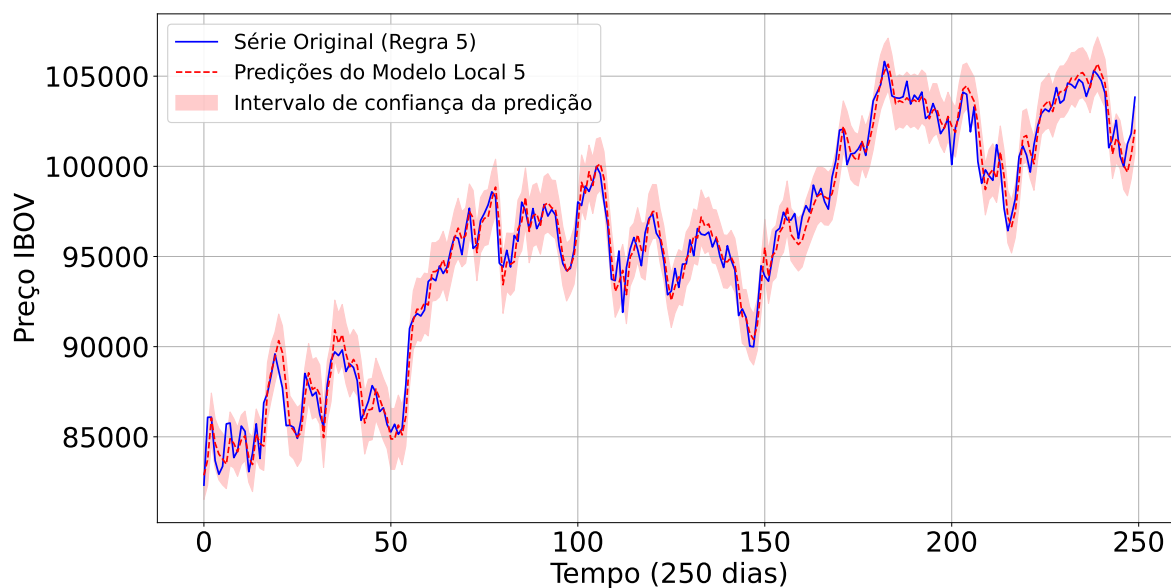
Figura 38 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 4.



Fonte: O Autor.

Pode-se observar que a quinta região identificada apresenta movimentações extremamente voláteis, em que mudanças bruscas de amplitude ocorrem há um aumento de amplitude nos intervalos de confiança do modelo. Contudo, o modelo conseguiu acompanhar tais variações de modo adequado, ver Fig. 39. Os parâmetros utilizados inicialmente e posteriormente, a janela de dados móveis, foram de  $(1,1,1)$  e  $(7,1,1)$ , respectivamente.

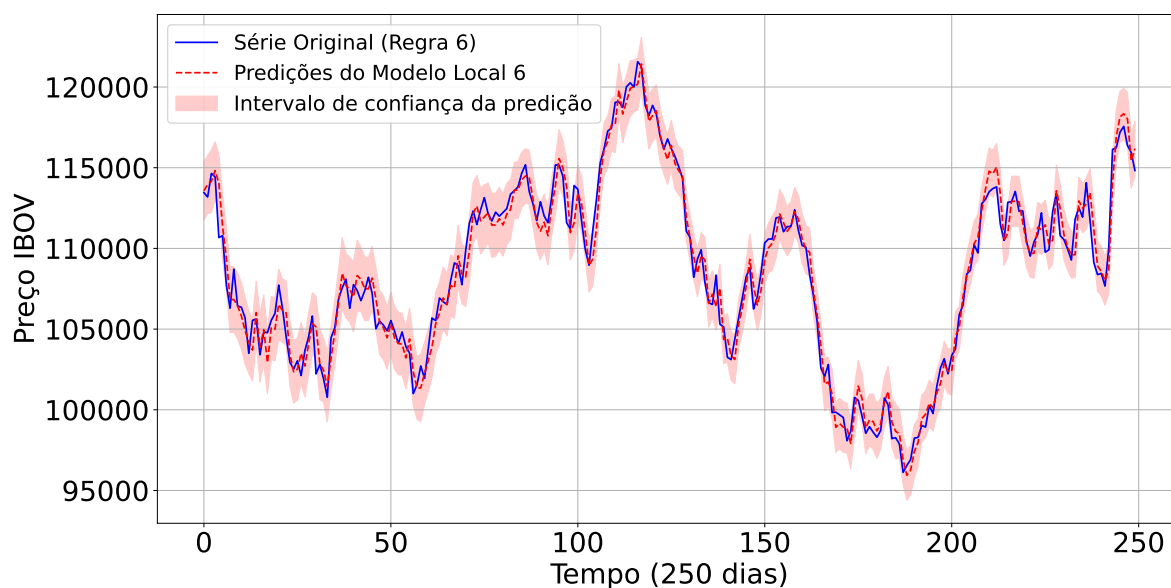
Figura 39 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 5.



Fonte: O Autor.

Igualmente, a sexta região identificada apresenta alta volatilidade. Os parâmetros utilizados para o sexto modelo local inicialmente e posteriormente, a janela de dados móveis, foram de (2,1,1) e (7,1,1), respectivamente.

Figura 40 – Simulação II - Gráfico da série original e a predição para a região pertencente a Regra 6.



Fonte: O Autor.

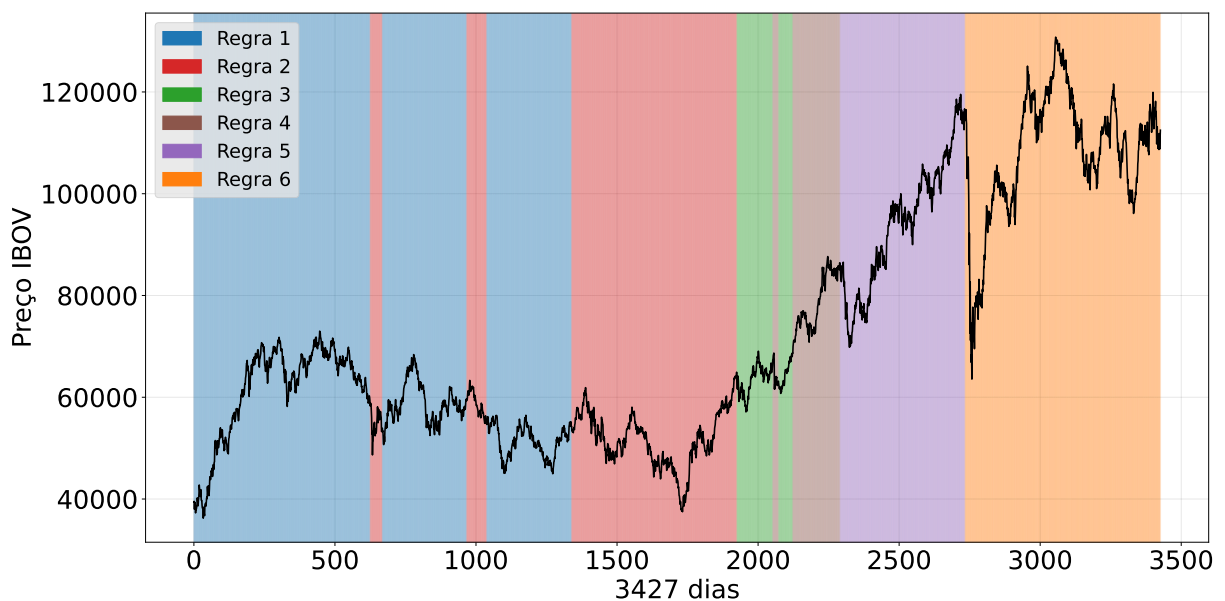
A Tab. 6 apresenta todas as métricas de avaliação para cada modelo local, auxiliadas pela visualização das mesmas em boxplots, ver Fig. 42. Dessa forma, reafirmando os apontamentos realizados quanto a melhora na representação dos dados, através das métricas alcançadas.

Métricas para a Simulação II (Média $\pm$ Desvio Padrão)				
Modelos	MAE	RMSE	MAPE (%)	SMAPE (%)
1°	453.5 $\pm$ 24.5	591.1 $\pm$ 34.9	0.74 $\pm$ 0.04	0.37 $\pm$ 0.02
2°	449.5 $\pm$ 77.3	578.7 $\pm$ 102.1	0.89 $\pm$ 0.17	0.44 $\pm$ 0.09
3°	457.4 $\pm$ 25.9	567.4 $\pm$ 40.1	0.71 $\pm$ 0.05	0.36 $\pm$ 0.02
4°	498.9 $\pm$ 35.4	642.6 $\pm$ 39.9	0.65 $\pm$ 0.03	0.32 $\pm$ 0.02
5°	672.2 $\pm$ 50.8	841.6 $\pm$ 60.2	0.73 $\pm$ 0.12	0.36 $\pm$ 0.06
6°	800.2 $\pm$ 166.3	1042.4 $\pm$ 313.7	0.75 $\pm$ 0.25	0.38 $\pm$ 0.12

Tabela 6 – Métricas de Avaliação utilizadas na Simulação II.

Observa-se valores baixos da métrica RMSE para os quatro primeiros modelos, com um aumento sutil nos últimos dois modelos. Sem altas discrepâncias entre as métricas, sugere-se que as predições tiveram alta acurácia ao longo do tempo sem apresentar pontos outliers de alta magnitude, destacando somente o segundo modelo que apresentou um desvio padrão aumentado, sendo melhor visualizado na Fig. 42.

Figura 41 – Série Original destacando os setores locais identificados por cada regra, referente ao experimento com 4 variáveis.



Fonte: O Autor.

Evidenciando as regiões identificadas, ver Fig. 41, nota-se que houve uma melhor identificação das diferentes dinâmicas da série. Vale destacar que houve uma clara iden-

tificação da mudança de dinâmica na série causada pela crise sanitária da COVID-19, transição entre a regra 5 e 6, que apesar de ocasionada por eventos externos ainda assim foi possível capturar a variação temporal dos dados.

Além disso, salienta-se que as regiões não necessariamente são contínuas no tempo. Vemos que a Regra 2 é criada próximo de 600 observações e logo após a dinâmica da série retorna para a Regra 1. Assim, conseguimos inferir que o algoritmo não necessariamente indica que está ocorrendo uma mudança de dinâmica imediata, mas podendo sugerir que uma mudança futura está por vir.

A Tab. 7 exibe a importância normalizada de cada atributo, entre entradas e regressores do modelo, de forma percentual, destacando os 3 mais importantes por modelo.

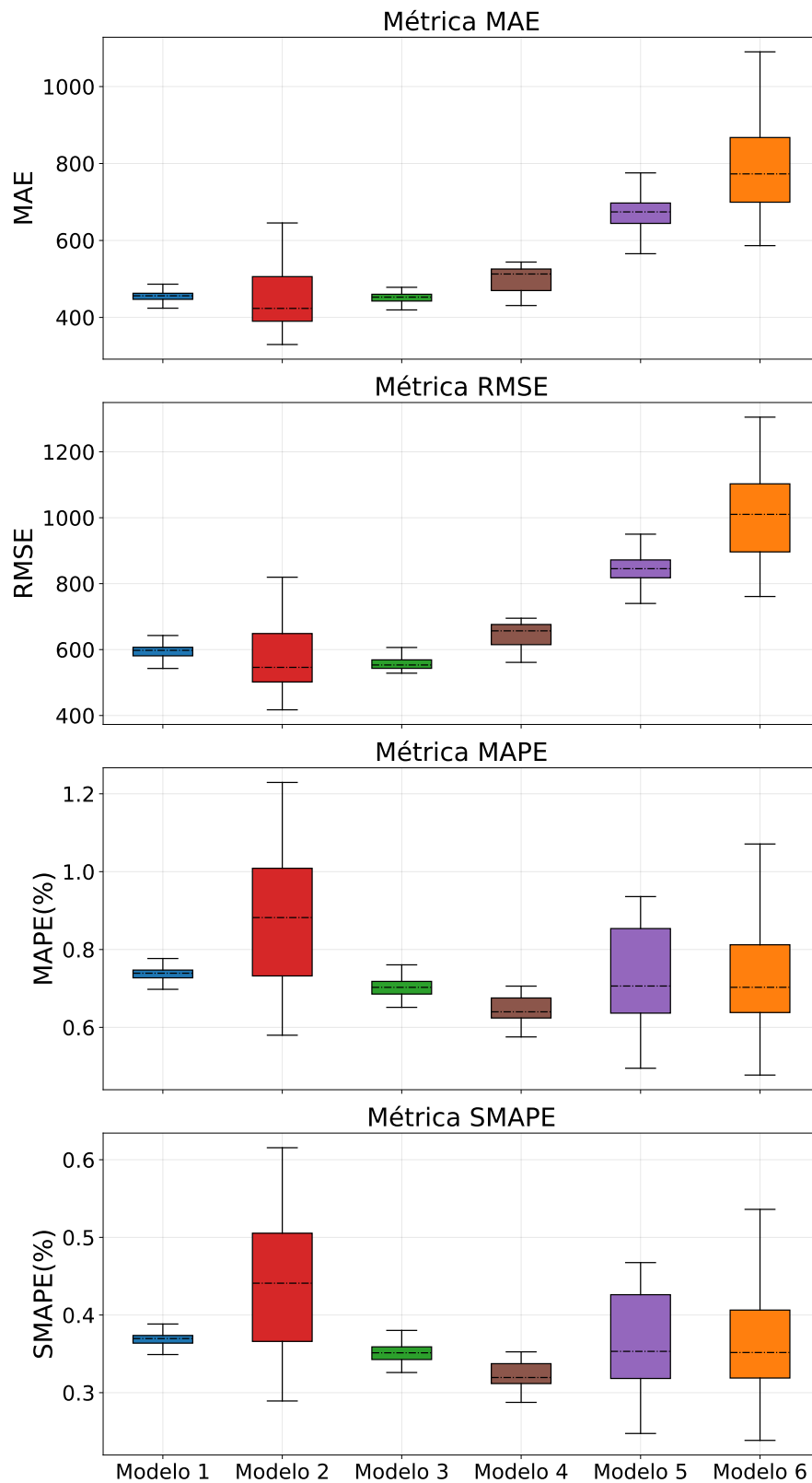
Atributo		Importância dos Atributos Normalizada por Modelo Local					
Tipo	Atraso/Grau	1	2	3	4	5	6
Exógena	1°	4.36 %	13.28 %	6.41 %	<b>10.86 %</b>	8.80 %	<b>12.72 %</b>
Exógena	2°	<b>41.39 %</b>	<b>40.98 %</b>	<b>24.37 %</b>	<b>54.76 %</b>	<b>40.14 %</b>	<b>34.54 %</b>
Exógena	3°	<b>13.20 %</b>	<b>12.21 %</b>	7.86 %	<b>20.95 %</b>	<b>15.67 %</b>	9.35 %
AR	1°	4.27 %	6.76 %	<b>10.15 %</b>	0.95 %	1.75 %	1.37 %
AR	2°	9.76 %	7.41 %	<b>10.74 %</b>	4.30 %	<b>15.04 %</b>	<b>16.62 %</b>
AR	3°	2.47 %	0.09 %	8.14 %	0.44 %	0.17 %	1.19 %
AR	4°	5.13 %	0.60 %	7.55 %	0.64 %	3.45 %	11.04 %
AR	5°	-	1.43 %	5.54 %	0.68 %	2.44 %	4.99 %
AR	6°	-	0.82 %	5.00 %	0.31 %	1.64 %	6.22 %
AR	7°	-	1.64 %	4.85 %	2.33 %	1.46 %	-
AR	8°	-	2.91 %	2.87 %	1.53 %	-	-
MA	1°	<b>19.42 %</b>	<b>11.89 %</b>	6.51 %	2.25 %	9.46 %	1.96 %

Tabela 7 – Importância Normalizada de cada atributo, por modelo, destacando os com maior pertinência para a Simulação II.

É relevante ressaltar que a partir da importância de cada atributo, ao contrário de dias específicos apresentarem uma maior causalidade com as predições, as entradas exógenas de médias móveis se tornaram os maiores pivôs para as predições alcançadas. Indicando que a movimentação dos dados agrupados tem um efeito maior sobre os resultados alcançados, exceto pelas ressalvas dos modelos 3, 5 e 6, em que termos autoregressivos tiveram maior importância.

Ademais, nos modelos locais em que as médias móveis foram mais relevantes, sugere-se uma dinâmica mais previsível, onde o comportamento agregado dos dados proporciona uma representação mais precisa. Por outro lado, nos modelos em que os termos autoregressivos se mostraram mais influentes, observa-se que eventos individuais impactam significativamente a dinâmica, resultando em uma maior volatilidade e, conseqüentemente, maior incerteza.

Figura 42 – Métricas de Avaliação referentes ao experimento *out-of-sample* com 4 variáveis, resultando em 6 modelos locais, dispostas em *boxplots*.



Fonte: O Autor.

## 5 Considerações Finais

O presente estudo buscou explorar as dinâmicas do Índice Bovespa (IBOV) com o intuito de trazer uma abordagem mais objetiva e interpretável no que tange a especulação de valores futuros de ativos financeiros, agregando a justificabilidade e legitimidade para as decisões produzidas. O mesmo foi realizado através da aplicação de modelos autoregressivos e de conceitos da lógica *fuzzy* via métodos de clusterização de dados, além da utilização de diferentes janelas de dados.

Durante a primeira simulação, a falta de informações adequadas nos dados utilizados impediu a identificação devida das dinâmicas do IBOV, resultando em uma representação inadequada dos modelos e métricas de acurácia pouco congruentes. Contudo, possibilitou compreender o comportamento do modelo e inferir melhorias ao utilizar novos dados de entrada na segunda simulação, obtendo métricas pertinentes. Além disso, esta simulação conseguiu identificar de forma adequada as diferentes dinâmicas do IBOV, ajustando-se às mudanças de momento do mercado, destacando o reconhecimento da mudança de dinâmica causada pela COVID-19.

Os resultados indicaram que, nos modelos onde as médias móveis apresentaram maior relevância, observou-se uma dinâmica mais previsível, com o comportamento agregado dos dados oferecendo uma representação mais precisa. Por outro lado, nos modelos em que os termos autoregressivos mostraram maior importância, constatou-se que eventos individuais impactam significativamente a dinâmica, resultando em maior volatilidade e incerteza.

Ademais, os modelos analisados destacaram a importância de ajustar as características das regiões identificadas para mitigar a incerteza observada, especialmente nos casos em que a quantidade de dados era limitada. Essa abordagem assegura que o algoritmo receba informações suficientes para evitar a geração de regras incompletas ou incorretas, garantindo assim a robustez e a precisão das predições.

Em síntese, este trabalho contribuiu para uma compreensão mais aprofundada das dinâmicas do IBOV, evidenciando que diferentes técnicas de modelagem podem captar variabilidades distintas do mercado. As análises sugerem que a combinação de médias móveis e termos autoregressivos pode proporcionar uma visão mais abrangente e precisa, adaptando-se às nuances dos dados financeiros provendo a legitimidade necessária para os resultados produzidos.

## 5.1 Trabalhos Futuros

Conforme destacado previamente, os modelos alcançados durante os experimentos computacionais indicaram a necessidade de ajustar as características das regiões identificadas com o intuito de mitigar a incerteza observada, com destaque para regiões amostrais menores. Evidências empíricas sugerem que representações de modelos autoregressivos de séries temporais financeiros e macroeconômicas estão sujeitas a quebras estruturais (DUFAYS; ROMBOUTS, 2020; PESARAN; TIMMERMANN, 2002).

Quebras estruturais no contexto de representação de séries temporais financeiras por modelos autogressivos podem ser caracterizados por uma mudança na série ou nos parâmetros do modelo. No primeiro cenário, a quebra pode ser ocasionada devido a uma mudança brusca da média dos dados, seja aumentando ou diminuindo o valor, sugerindo uma micro mudança na série. No segundo cenário, a quantidade de regressores utilizados acarreta na degradação da matriz de covariância do modelo, indicando um dimensionamento ruim dos parâmetros (PESARAN; TIMMERMANN, 2005).

À luz do contexto, a incorporação de um método capaz de lidar com o reconhecimento de quebras estruturais na série, dentro de cada região identificada, é um caminho evidente a ser seguido, agregando na exploração das dinâmicas do IBOV de forma a identificar micro regiões, dentro das regiões identificadas pelo OEC.

# Referências

- AGUIRRE, L. A. Introdução à identificação de sistemas - técnicas lineares e não lineares aplicadas a sistemas: Teoria e aplicação. *Editora UFMG*, v. 4, 2014. Citado 6 vezes nas páginas 17, 21, 22, 23, 24 e 25.
- ANGELOV, P. P.; FILEV, D. P. An approach to online identification of takagi-sugeno fuzzy models. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, v. 34, n. 1, p. 484 – 498, 2004. Citado na página 28.
- ARRIETA, A. B. et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. 10 2019. Citado na página 36.
- BIRAN, O.; COTTON, C. V. Explanation and justification in machine learning : A survey or. In: . [S.l.: s.n.], 2017. Citado na página 37.
- BROŽ, Z.; DOSTÁL, P. Fuzzy logic decision support for long-term investing in the financial markets. In: ZELINKA, I. et al. (Ed.). *Nostradamus: Modern Methods of Prediction, Modeling and Analysis of Nonlinear Systems*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013. p. 113–121. Citado na página 18.
- BUENO, R. d. L. d. S. *Econometria de séries temporais*. 2. ed. São Paulo: Cengage Learning, 2011. Citado 2 vezes nas páginas 40 e 49.
- CAMILO, E. V. Modelos de previsão utilizando séries temporais. *Universidade Estadual da Paraíba*, 2012. Citado na página 48.
- CORDOVIL, L. A. Q. et al. Uncertain data modeling based on evolving ellipsoidal fuzzy information granules. *IEEE Transactions on Fuzzy Systems*, v. 28, n. 10, p. 2427–2436, Oct 2020. ISSN 1941-0034. Citado 3 vezes nas páginas 29, 30 e 31.
- DU, M.; LIU, N.; HU, X. Techniques for interpretable machine learning. *Communications of the ACM*, v. 63, n. 1, p. 68–77, Dec. 2019. ISSN 0001-0782. Citado na página 17.
- DUFAYS, A.; ROMBOUTS, J. V. Relevant parameter changes in structural break models. *Journal of Econometrics*, v. 217, n. 1, p. 46 – 78, 2020. Citado na página 81.
- FEYISA, H. The world economy at covid-19 quarantine: contemporary review. *International Journal of Economics, Finance and Management Sciences*, v. 8, n. 2, p. 63–74, 2020. Citado na página 60.
- FREITAS, C. M.; SANTIAGO, Y. V.; CARVALHO, S. M. S. Downside risk aplicado a carteiras de ações brasileiras durante período pandêmico da covid-19. *Trends in Computational and Applied Mathematics*, Sociedade Brasileira de Matemática Aplicada e Computacional - SBMAC, v. 24, n. 3, p. 557–574, Jul 2023. ISSN 2676-0029. Citado na página 60.
- GERON, A. *Hands-on Machine Learning with Scikit-Learn and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. Sebastopol, CA: O'Reilly Media, 2017. ISBN 978-1491962299. Citado na página 17.

- GKATZIA, D.; LEMON, O.; RIESER, V. Natural language generation enhances human decision-making with uncertain information. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Berlin, Germany: Association for Computational Linguistics, 2016. p. 264–268. Citado na página 37.
- GUIDE INVESTIMENTOS. *A relação inversa entre o Ibovespa e a taxa de câmbio*. São Paulo, SP: [s.n.], 2023. <<https://conteudos.guide.com.br/textos/a-relacao-inversa-entre-o-ibovespa-e-a-taxa-de-cambio/>>. Acesso em: 27 jun. 2024. Citado na página 39.
- HACHICHA, N.; JARBOUI, B.; SIARRY, P. A fuzzy logic control using a differential evolution algorithm aimed at modelling the financial market dynamics. *Information Sciences*, v. 181, n. 1, p. 79–91, 2011. Citado na página 18.
- JACQUES, K.; BORGES, S.; MIRANDA, G. Relações entre os indicadores econômico-financeiros e as variáveis macroeconômicas dos segmentos empresariais da b3. *Revista de Administração, Contabilidade e Economia da Fundace*, v. 11, 01 2020. Citado 2 vezes nas páginas 39 e 50.
- JONES, S. L.; NETTER, J. M. Efficient capital markets. *The Library of Economics and Liberty*, 2008. Coleção: Corporations and Financial Markets. Citado na página 17.
- KADIMA GESTÃO DE INVESTIMENTOS LTDA. *Carta Trimestral - Dezembro 2023*. Rio de Janeiro: [s.n.], 2023. Cartas de Gestão. Disponível em: <<https://www.kadimaasset.com.br/cartas/>>. Acesso em: 22 jun. 2024. Citado na página 17.
- KWIATKOWSKI, D. et al. Testing the null hypothesis of stationarity against the alternative of a unit root: How sure are we that economic time series have a unit root? *Journal of Econometrics*, v. 54, n. 1, p. 159–178, 1992. ISSN 0304-4076. Citado 2 vezes nas páginas 39 e 51.
- LUGER, G. *Inteligência Artificial*. [S.l.]: PEARSON BRASIL, 2013. ISBN 9788581435503. Citado na página 17.
- MOSHTAGHI, M.; LECKIE, C.; BEZDEK, J. Online clustering of multivariate time-series. In: . [S.l.: s.n.], 2016. p. 360–368. Citado 6 vezes nas páginas 17, 31, 33, 34, 35 e 40.
- PAN, M.; WANG, H.; HUANG, J. T-s fuzzy modeling for aircraft engines: The clustering and identification approach. *Energies*, v. 12, n. 17, 2019. Citado na página 28.
- PESARAN, M.; TIMMERMANN, A. Market timing and return prediction under model instability. *Journal of Empirical Finance*, v. 9, n. 5, p. 495 – 510, 2002. Citado na página 81.
- PESARAN, M. H.; TIMMERMANN, A. Small sample properties of forecasts from autoregressive models under structural breaks. *Journal of Econometrics*, v. 129, n. 1-2, p. 183 – 217, 2005. Citado 2 vezes nas páginas 51 e 81.
- RANJBAR, N.; MOMTAZI, S.; HOMAYOONPOUR, M. Explaining recommendation system using counterfactual textual explanations. *Machine Learning*, v. 113, n. 4, p. 1989 – 2012, 2024. Citado na página 37.
- SAID, S. E.; DICKEY, D. A. Testing for unit roots in autoregressive-moving average models of unknown order. *Biometrika*, v. 71, n. 3, p. 599, 1984. Citado na página 39.

- SILVA, A.; NAGHETTINI, M.; PORTELA, M. Sobre a estimação de intervalos de confiança para os quantis de variáveis aleatórias hidrológicas. *Recursos Hídricos*, v. 32, p. 63–75, 11 2011. Citado na página 59.
- STEURER, M.; HILL, R. J.; PFEIFER, N. Metrics for evaluating the performance of machine learning based automated valuation models. *Journal of Property Research*, v. 38, n. 2, p. 99–129, 2021. Citado na página 25.
- TANAKA, K.; WANG, H. O. *Fuzzy Control Systems Design and Analysis*. [S.l.: s.n.], 2001. Citado 2 vezes nas páginas 27 e 28.
- TAYLOR, H.; KARLIN, S. *An introduction to stochastic modeling*. Rev. ed. San Diego [u.a.]: Academic Press, 1994. Citado na página 17.
- TRINDADE, J. F. V. *Interpretabilidade em Modelos de Sistemas de Recomendação*. Dissertação (Mestrado) — Departamento de Ciências de Computadores, Mestrado em Ciências de Dados (Data Science), 2020. Citado na página 36.
- VARELLA, T. F. Construção e avaliação de estratégias de investimento com o uso de médias móveis como único indicador técnico. *Universidade Federal do Rio Grande do Sul*, 2012. Citado na página 50.
- XU, D.; TIAN, Y. A comprehensive survey of clustering algorithms. *Annals of Data Science*, Springer, v. 2, n. 2, p. 165–193, 2015. Citado na página 31.
- YE, L.; JOHNSON, P. The impact of explanation facilities on user acceptance of expert systems advice. *MIS Quarterly: Management Information Systems*, v. 19, n. 2, p. 157–172, June 1995. ISSN 0276-7783. Citado na página 37.
- ZADEH, L. Soft computing and fuzzy logic. *IEEE Software*, v. 11, n. 6, p. 48–56, 1994. Citado na página 27.
- ZADEH, L. A. Fuzzy sets. *Information and Control*, v. 8, n. 3, p. 338–353, 1965. ISSN 0019-9958. Citado na página 17.
- ZHU, Y.; STEC, P. Simple control-relevant identification test methods for a class of ill-conditioned processes. *Journal of Process Control*, v. 16, n. 10, p. 1113–1120, 2006. ISSN 0959-1524. Citado na página 25.