

O efeito da discretização na classificação: um estudo comparativo de técnicas de aprendizagem supervisionada para caracterização de variáveis econômicas

Antônio M. T. de Araújo¹, Eduardo Palhares Júnior¹

¹Instituto Federal de Educação, Ciência e Tecnologia do Amazonas (IFAM)
Campus Manaus Zona Leste – Manaus, AM – Brasil

teixeira2gpt@gmail.com, eduardo.palharesjr@ifam.edu.br

Abstract. *This work proposes a comparative study between different machine learning techniques, applied to the analysis of the phases of the Brazilian economic cycle. Various macroeconomic indicators were used to build a model capable of identifying and predicting alternation points in the economic cycle, such as the beginning of a recession or a recovery. In previous works, it was found that the data discretization stage was decisive in the quality of the classification process, but with an asymmetry caused by the COVID-19 pandemic. In this work, a longer time interval is proposed that can introduce the effects of the pandemic into training, seeking to minimize overfitting problems.*

Resumo. *Este trabalho propõe um estudo comparativo entre diversas técnicas de aprendizado de máquina, aplicadas na análise das fases do ciclo econômico brasileiro. Foram utilizados diversos indicadores macroeconômicos para construir um modelo capaz de identificar e prever os pontos de alternância do ciclo econômico, como o início de uma recessão ou de uma recuperação. Em trabalhos anteriores, verificou-se que a etapa de discretização dos dados foi decisiva na qualidade do processo de classificação, mas com uma assimetria causada pela pandemia de COVID-19. Neste trabalho, é proposta um intervalo temporal maior que consiga introduzir os efeitos da pandemia no treinamento, buscando minimizar problemas de sobreajuste.*

1. Introdução

Os sistemas especialistas aplicados ao estudo de economia são utilizados para auxiliar os pesquisadores e profissionais na tomada de decisões em relação à previsão do mercado de ações, valores mobiliários e commodities. Compreender a dinâmica e interrelação das variáveis macroeconômicas com a evolução futura do ciclo econômico é um problema que desperta grande interesse tanto para os acadêmicos quanto para os participantes do mercado. Na literatura, podemos identificar duas correntes principais de pesquisa. A corrente teórica tenta explicar as características que determinam os pontos de virada do ciclo, através do comportamento da curva de juros como a inclinação, nível e curvatura, tomando como base em várias teorias financeiras [Hicks et al. 1975, Vasicek 1977, COX et al. 1981, BROWN and DYBVIG 1986, HO and LEE 1986, Nelson and Siegel 1987, Heath et al. 1992, Svensson 1994, Alexander et al. 2001]. Por outro lado, a corrente empírica busca implementar o uso de diferentes métodos

econométricos e técnicas de descoberta de conhecimento, com o objetivo de modelar a estrutura e a dinâmica dos ciclos econômicos.

Recentemente, o uso de técnicas de descoberta de conhecimento tem ficado cada vez mais comum, principalmente devido ao fato de serem capazes de capturar e lidar com não linearidades existentes entre as variáveis, bem como a complexidade envolvida na sazonalidade e nas rupturas estruturais [Ju et al. 1997, Kim and Noh 1997, Zimmermann et al. 2002, Jacovides 2008, Oh and Han 2000, Vela 2013, Gogas et al. 2014]. No entanto, existem muitas dificuldades envolvidas com o uso desse tipo de ferramenta, já que os padrões não linearidades consideradas em cada etapa da predição costumam ser extremamente complicados, tornando a interpretação inviável. Dessa forma, muitos desses métodos apesar de terem excelentes resultados, tem utilização limitada por serem considerados caixas-pretas.

A classificação de variáveis econômicas tem sido utilizada a muito tempo como instrumento para prever movimentos relevantes dentro dos estudos de economia, como [Burns and Mitchell 1946] que estudaram as desacelerações econômicas e recessões. Ao longo do tempo, várias variáveis diferentes foram propostas e avaliadas como indicadores econômicos, como discutido em [Estrella and Mishkin 1995], os quais são notadamente reconhecido como importantes na identificação de um ponto de inflexão nos movimentos de crescimento econômico.

Assim, este trabalho pretende desenvolver modelos baseados em aprendizado de máquina que analisem diversos indicadores econômicos, buscando uma metodologia que sinalize a possibilidade de um ponto de inversão do crescimento econômico, no caso, o início de uma fase de recessão em diversos horizontes temporais. Essa abordagem já foi aplicada e discutida anteriormente em [Palhares Júnior et al. 2024], no entanto como o estudo foi realizado durante a pandemia de Covid-19, levantou-se a hipótese que os eventos altamente não lineares nesse período estavam enviesando a qualidade dos resultados. Como esse trabalho considera um intervalo temporal maior, os dados referentes a da pandemia de Covid-19 fazem parte também do conjunto de treinamento, possibilitando assim isolar os efeitos caóticos desse período e focar a análise na arquitetura proposta, além de estimar o quanto as técnicas de aprendizagem de máquina propostos são robusto para aprender com esses eventos caóticos. Dessa forma o trabalho tem um enfoque maior nos efeitos da discretização utilizando diferentes quantidades de classe e intervalos.

2. Metodologia

A metodologia deste trabalho tem como objetivo estudar e caracterizar o comportamento de ciclos econômicos brasileiros utilizando técnicas de aprendizado estatístico. O foco é modelar o PIB como função de variáveis econômicas, considerando ciclos econômicos como flutuações entre estágios de expansão e contração. Além disso, diferentes teorias econômicas são abordadas, como os ciclos de dois, três ou quatro estágios [Bhaumin line].

O interesse pelo PIB como variável principal decorre do impacto da pandemia de Covid-19 e do risco de recessão global, mas o modelo é flexível, permitindo a análise de outras variáveis econômicas mediante ajustes nas hipóteses. Para facilitar o desenvolvimento e a análise, o projeto foi modularizado, permitindo depuração e compreensão dos estágios intermediários, conforme mostra o diagrama da figura 1.

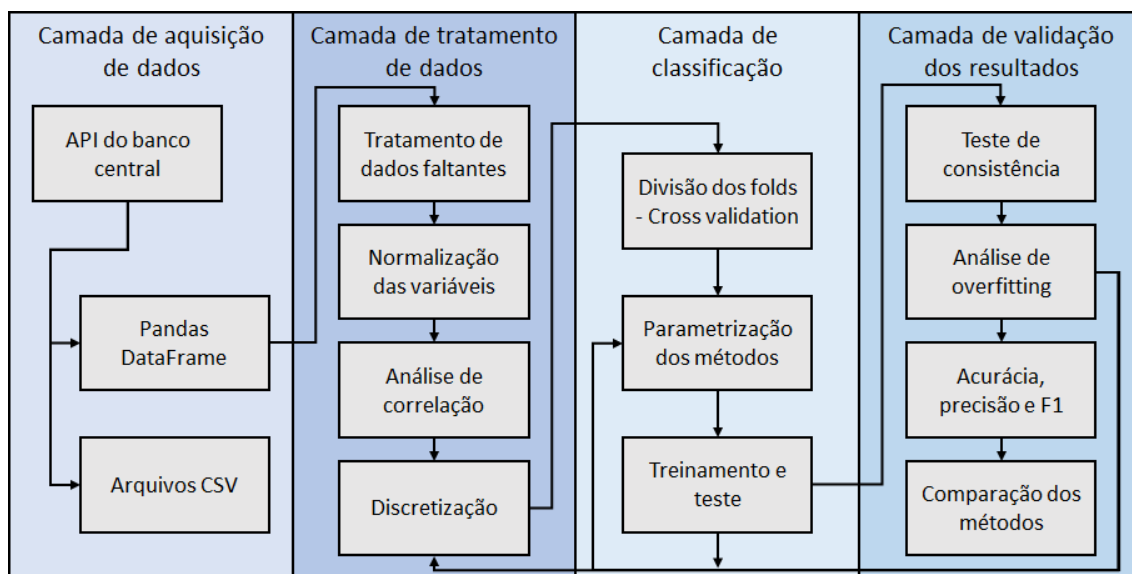


Figura 1. Arquitetura da metodologia

2.1. Aquisição e preparação de dados

Este trabalho refere-se a dados macroeconômicos da economia do Brasil conforme tabela 1, durante o período situado entre janeiro de 2002 e maio de 2024. Os dados foram obtidos diretamente do Banco Central do Brasil, através de uma API disponibilizada pela autoridade financeira do país.

Tabela 1. Lista e descrição das variáveis econômicas

Variável econômica	Descrição
PIB	Produto Interno Bruto mensal
IPA	Índice de preços ao produtor amplo
IPEM	Indicador da produção - extrativa mineral
IPIT	Indicadores da produção - indústria de transformação
IPBC	Indicadores da produção - bens de capital
IPBCD	Indicadores da produção - bens de consumo duráveis
IVVV	Índice volume de vendas no varejo - Automóveis, motocicletas, partes e peças - Brasil
VVCCL	Vendas de veículos pelas concessionárias - Comerciais leves
VVCC	Vendas de veículos pelas concessionárias - Caminhões
IEF	Índice de Expectativas Futuras
ICC	Índice de Confiança do Consumidor
Spub	Saldos das operações de crédito das instituições financeiras sob controle público
Spriv	Saldos das operações de crédito das instituições financeiras sob controle privado
M1	Meios de pagamento - M1 (média dos dias úteis do mês)
M2	Meios de pagamento - M2 (média dos dias úteis do mês)

As variáveis analisadas possuem granularidade mensal e não apresentam dados faltantes no período considerado. Para garantir comparabilidade, todas foram transformadas em variações percentuais mensais. Como o objetivo é prever a variação percentual do PIB, essas transformações foram usadas como entradas para o modelo classificador.

$$\Delta x_i = \frac{x_i - x_{i-1}}{x_{i-1}} \quad (1)$$

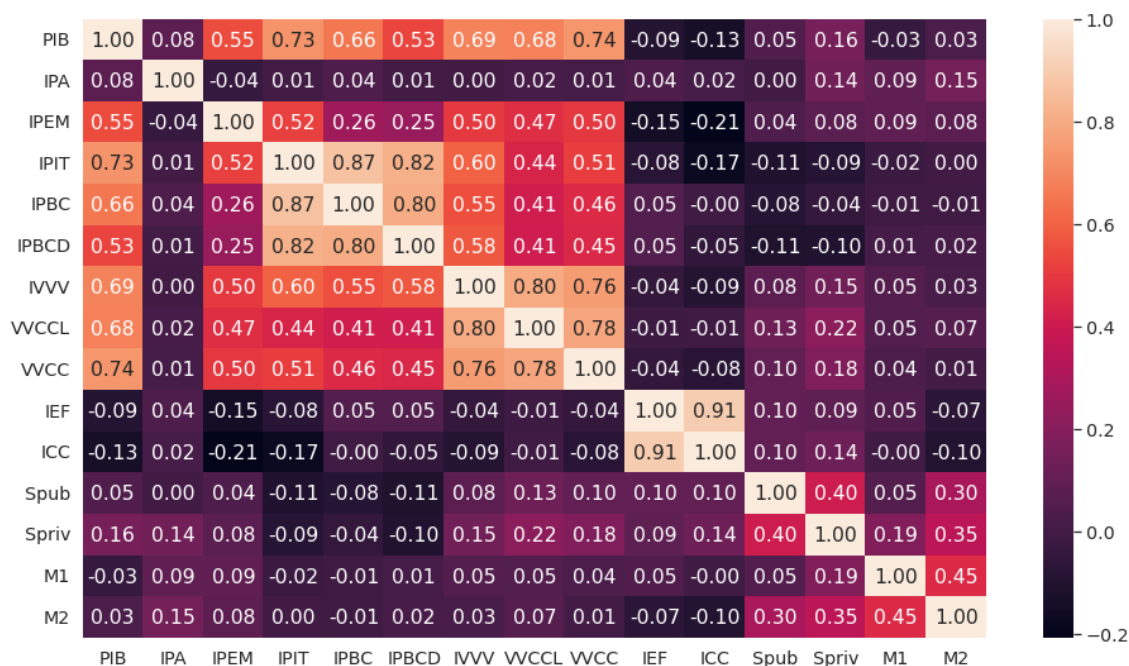


Figura 2. Correlação entre as variáveis econômicas - dataset completo.

Apesar da correlação de diversas variáveis serem aparentemente baixas, é preciso lembrar que os fenômenos analisados podem ser bastante não lineares. Buscando considerar essas sutilezas, foi proposta uma nova base de dados considerando somente as variáveis de correlação significativa, nomeada como base de dados restrita conforme figura 3.

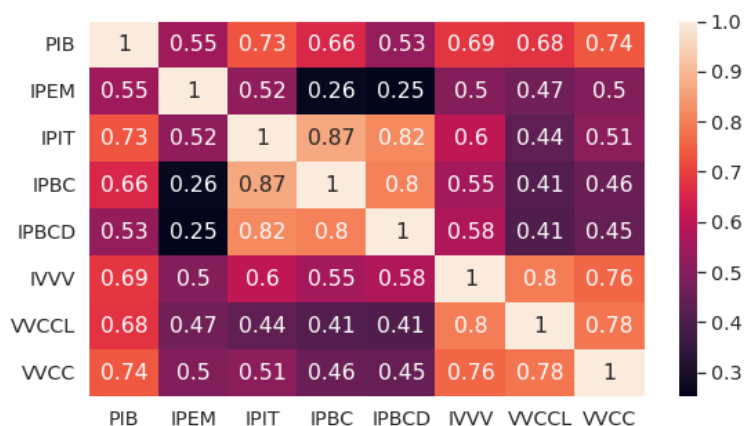


Figura 3. Correlação entre as variáveis econômicas - dataset restrito.

Todas as séries utilizadas foram testadas para estacionaridade através dos testes KPSS (Kwiatkowski–Phillips–Schmidt–Shin) e ADF (Augmented Dickey Fuller). Além disso, os estudos apresentados à seguir foram testados na base completa e restrita.

2.2. Discretização

A etapa de discretização é um estágio muito importante na construção do modelo, pois influencia na correta visualização do fenômeno. Neste trabalho foram propostas diversas discretizações, buscando uma representação suficientemente simples para evitar problemas numéricos, mas capaz de representar de maneira adequada as condições de tomada de decisão.

Intervalo com 3 categorias padrão: A discretização com 3 categorias separa os movimentos de alta e de queda significativa com um movimento de estagnação, em que a variação é inferior a um desvio padrão em relação à média, ou seja, aproximadamente 68% dos casos. Essa abordagem mostrou-se bastante promissora para uma aplicação prática, pois combina uma boa precisão na classificação com categorias efetivas na tomada de decisão.

$$\begin{cases} x \mapsto -1, & \text{if } \Delta x \leq \mu_x - \sigma_x; \\ x \mapsto 0, & \text{if } \mu_x - \sigma_x < \Delta x < \mu_x + \sigma_x; \\ x \mapsto 1, & \text{if } \Delta x \geq \mu_x + \sigma_x. \end{cases} \quad (2)$$

Intervalo com 5 categorias padrão: A classificação considerando 5 categorias é um refinamento em relação a 3 categorias, pois além de separar os casos de estagnação, também consegue separar casos muito raros, com 95% nas 3 categorias centrais.

$$\begin{cases} x \mapsto -2, & \text{if } \Delta x \leq \mu_x - 2 \cdot \sigma_x; \\ x \mapsto -1, & \text{if } \mu_x - 2 \cdot \sigma_x < \Delta x \leq \mu_x - \sigma_x; \\ x \mapsto 0, & \text{if } \mu_x - \sigma_x < \Delta x < \mu_x + \sigma_x; \\ x \mapsto 1, & \text{if } \mu_x + \sigma_x \leq \Delta x < \mu_x + 2 \cdot \sigma_x; \\ x \mapsto 2, & \text{if } \Delta x \geq \mu_x + 2 \cdot \sigma_x. \end{cases} \quad (3)$$

Apesar de ser uma descrição bastante interessante para o fenômeno, [Palhares Júnior et al. 2024] mostra que em casos com poucos dados, as categorias extremas podem ser bastante raras, fazendo com que alguns métodos apresentem dificuldade na aprendizagem dessas categorias, causando grande incidência falsos positivos e negativos, além de problemas de convergência para alguns classificadores.

Intervalo com 5 categorias modificado: Buscando minimizar o efeito dessas classes muito raras e seguindo o que foi proposto em [Palhares Júnior et al. 2024], um novo intervalo mais flexível representado na figura 4

$$\begin{cases} x \mapsto -2, & \text{if } \Delta x \leq \mu_x - 1,6745 \cdot \sigma_x; \\ x \mapsto -1, & \text{if } \mu_x - 1,6745 \cdot \sigma_x < \Delta x \leq \mu_x - 0,6745 \cdot \sigma_x; \\ x \mapsto 0, & \text{if } \mu_x - 0,6745 \cdot \sigma_x < \Delta x < \mu_x + 0,6745 \cdot \sigma_x; \\ x \mapsto 1, & \text{if } \mu_x + 0,6745 \cdot \sigma_x \leq \Delta x < \mu_x + 1,6745 \cdot \sigma_x; \\ x \mapsto 2, & \text{if } \Delta x \geq \mu_x + 1,6745 \cdot \sigma_x. \end{cases} \quad (4)$$

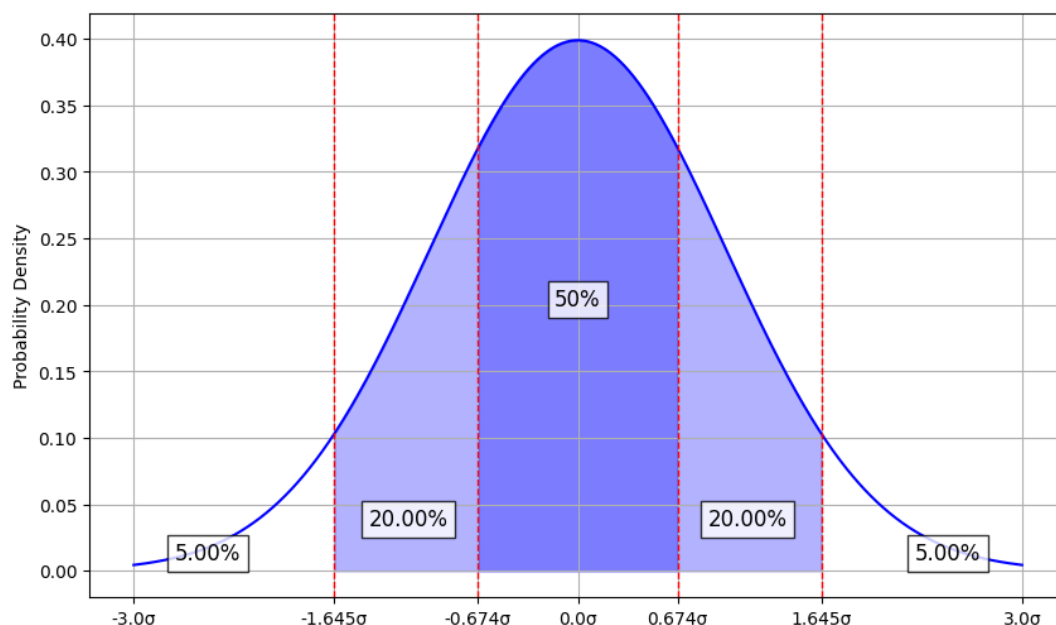


Figura 4. Distribuição dos dados em classes modificada

2.3. Classificação

2.3.1. Validação cruzada

Para evitar sobreajuste e viés, foi empregada a validação cruzada com divisão em subconjuntos de treinamento e teste, preservando a ordem cronológica das séries temporais. O método adotado foi o time series split cross validation, que expande iterativamente os conjuntos de treinamento e teste ao longo da série histórica. Após testar várias configurações, foi escolhida a divisão com 30% dos dados destinados ao teste e validação.

2.3.2. Métodos de classificação

Várias técnicas estatísticas e de mineração de dados foram aplicadas, as quais estão disponíveis como parte do conjunto de ferramentas da biblioteca Scikit-Learn do Python. Para todas essas técnicas, foram utilizadas as configurações padrão propostas pelo conjunto de ferramentas, uma vez que a otimização exaustiva dos parâmetros e arquiteturas do modelo estava fora do escopo do trabalho. No entanto, buscando solucionar alguns problemas de convergência, algumas dessas técnicas foram implementadas explicitamente de modo a testar com mais atenção alguns hiper parâmetros. Em detalhe, as técnicas aplicadas foram: k-vizinho mais próximo (KNN), gaussian Naive Bayes (NB), árvores de decisão (DT), florestas aleatórias (RF), logistic regression (LR), máquinas de vetores de suporte (SVC), redes neurais artificiais (NN).

3. Resultados

Para comparar os diversos cenários apresentados, é proposta uma análise comparativa entre todos os métodos, em relação à acurácia alcançada nas etapas de treinamento e de teste. Além disso, uma análise mais detalhada do comportamento do score F1 tanto

em relação as diferentes discretizações, quanto ao que se refere ao conjunto de variáveis explicativas utilizadas em cada cenário.

3.1. Acurácia

Para analisar a acurácia é preciso comparar todos os métodos propostos com todas as abordagens de discretização. A acurácia na etapa de treinamento apresentou um resultado comparável ao apresentando em [Palhares Júnior et al. 2024], mas focaremos nesse trabalho a discussão da etapa de teste, que possui mais relevância prática.

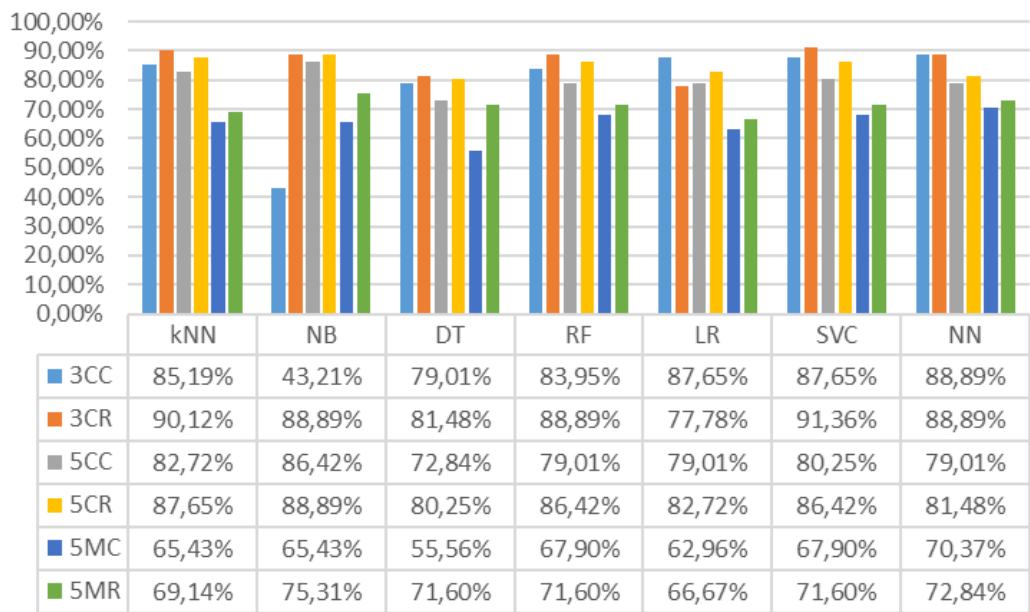


Figura 5. Comparação da acurácia no estágio de teste

Os rótulos das colunas referem-se a discretização utilizada e a quantidade de variáveis explicativas, por exemplo:

- 3CC → 3 categorias com a base completa (figura 2)
- 5MR → 5 categorias (modificada) com a base restrita (figura 3)

Quando comparamos os resultados de acurácia apresentados em [Palhares Júnior et al. 2024], vemos que os modelos tiveram uma melhora de performance na ordem de 10% em média. Além disso, o modelo com menos variáveis explicativas tem a tendência de uma acurácia maior, mas na maioria dos cenários apresentados essa diferença é marginal.

3.2. F-Score

A qualidade de cada método em relação a cada cenário é relativa, a depender do critério. No entanto, é importante comparar o efeito que cada cenário teve no desempenho de cada variação percentual no score F1 da base restrita (alguns parâmetros) em relação a base completa (todos os parâmetros). Como foram testados muitos cenários e em cada um deles há muitas categorias para se avaliar, a quantidade de resultados é significativamente grande. Dessa forma será apresentada uma análise qualitativa que demonstra qual o melhor cenário, o qual também será discutido de maneira quantitativa.

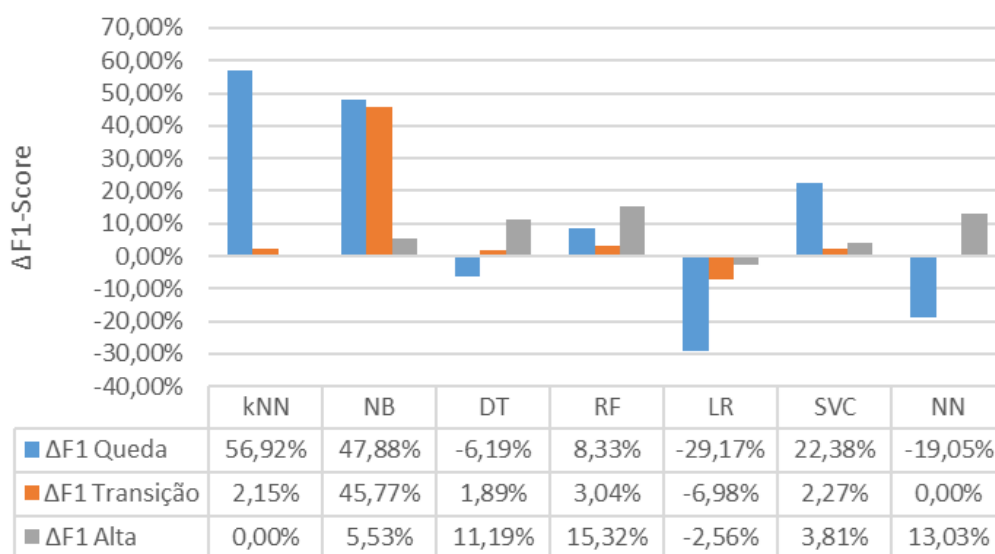


Figura 6. Variação do Score-f1 entre os métodos com 3 categorias padrão.

Considerando os cenários apresentados em [Palhares Júnior et al. 2024], podemos perceber um impacto mais significativo na restrição de variáveis. Observa-se uma maior variabilidade, especialmente em relação a prever movimentos de queda, mas de maneira geral é possível perceber que a redução de variáveis nesse cenário mostra-se vantajosa.

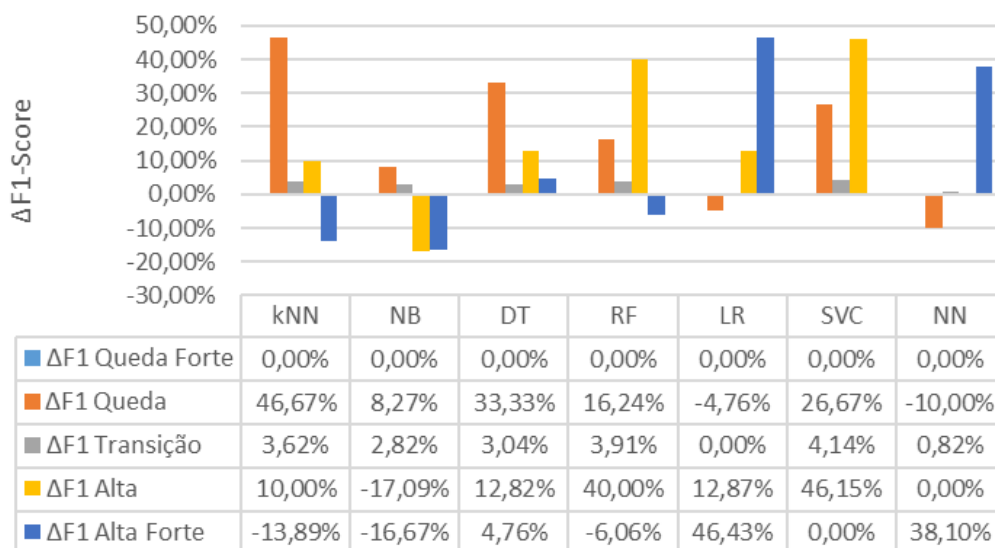


Figura 7. Variação do Score-f1 entre os métodos com 5 categorias padrão.

Avaliando o cenário com 5 categorias padrão, fica ainda mais evidenciado o impacto positivo da redução de variáveis. Não foi possível fazer uma comparação quantitativa, já vez que [Palhares Júnior et al. 2024] não traz nenhum resultado desse cenário. Qualitativamente, é possível perceber que o modelo de fato não consegue prever os movimentos de queda forte com nenhuma técnica, mas nesse trabalho foi possível verificar uma excelente capacidade de prever movimentos de alta forte.

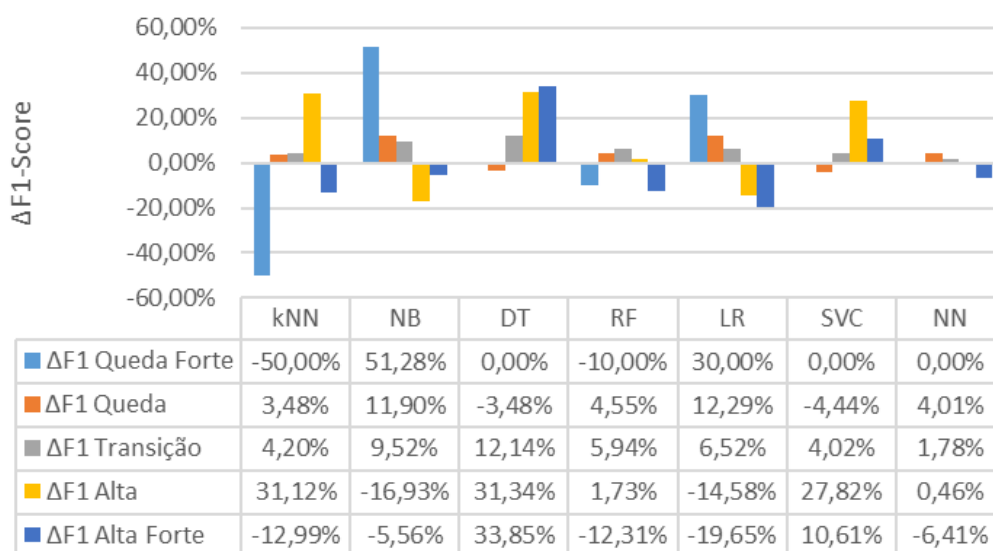


Figura 8. Variação do Score-f1 entre os métodos com 5 categorias modificada.

Novamente, observa-se uma melhora significativa ao restringir as variáveis explicativas, mas a variabilidade ainda é muito alta a depender do método, o que sugere que as limitações podem estar mais ligadas a erros de modelagem do que a capacidade preditiva dos métodos. Assim como no cenário anterior, há uma melhor capacidade preditiva para a categoria de alta forte em relação ao trabalho original.

Buscando resumir qual o ganho absoluto total alcançado em cada cenário quando se considera a base restrita, apesar da alta variância, a tabela 5 combina os ganhos líquidos absolutos de cada categoria em relação a cada cenário. Como são considerados 7 métodos, o ganho relativo é justamente o ganho líquido médio entre cada uma das 7 modelagens.

Tabela 2. Melhoria do score-f1 quando alguns parâmetros são removidos.

Categorias	3 categorias		5 categorias		5 modificada	
	absolute	relative	absolute	relative	absolute	relative
-2			0,00%	0,00%	21,28%	3,04%
-1	81,11%	11,59%	116,41%	16,63%	28,30%	4,04%
0	48,13%	6,88%	18,35%	2,62%	44,13%	6,30%
1	46,32%	6,62%	104,75%	14,96%	60,95%	8,71%
2			52,67%	7,52%	-12,46%	-1,78%
Total	175,56%	25,08%	292,18%	41,74%	167,12%	23,87%

De maneira geral, através da tabela 2 é possível verificar que um conjunto mais restrito de variáveis explicativas trouxe um ganho expressivo de performance em todos os casos analisados, o que não acontecia de forma tão significativa no trabalho original.

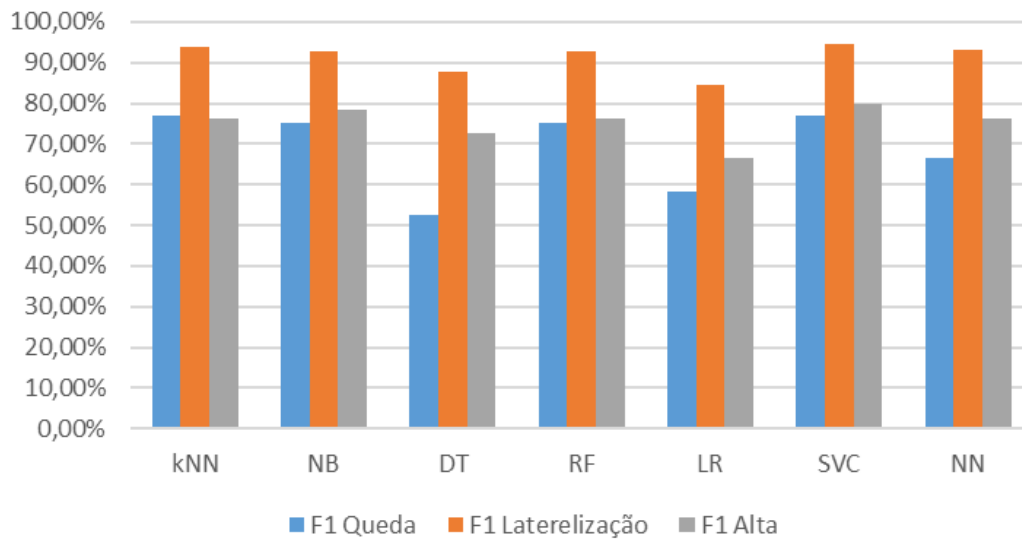


Figura 9. Resultados do Score-F1 da projeção do PIB brasileiro no modelo que considera 3 categorias padrão e conjunto restrito de variáveis

Conforme foi discutido em [Palhares Júnior et al. 2024], a figura 9 mostra discretização em 3 categorias apresenta um desempenho satisfatório e estável para todos os modelos, o que reforça que a modelagem é adequada. É notório que há ligeira dificuldade de prever os movimentos de queda, muito influenciado pelos efeitos do período da pandemia de COVID-19, que pode ser considerado um outlier.

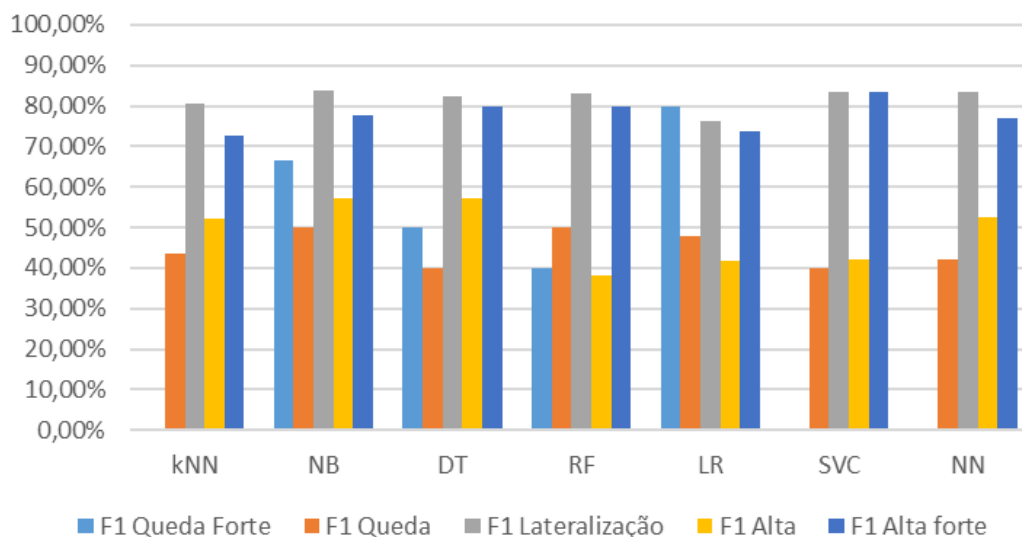


Figura 10. Resultados do Score-F1 da projeção do PIB brasileiro no modelo que considera 5 categorias modificada e conjunto restrito de variáveis

Ainda assim, analisando a figura 10, verificamos que o cenário com 5 categorias modificado é robusto para identificar os movimentos de alta forte, podendo ser aplicado

como um segundo estágio de previsão, ou seja, uma espécie de zoom nos cenários encontrados na discretização em 3 categorias. Além disso, a distorção causada pela pandemia de COVID-19 gera incertezas quanto a capacidade do método para prever quedas fortes, portanto, novos testes podem ser realizados excluindo esse período histórico com vista a mitigar esse possível problema de viés.

4. Conclusão

Este artigo propõe uma metodologia para modelar e prever o PIB do Brasil, utilizando indicadores macroeconômicos e técnicas de aprendizagem estatística para lidar com as não linearidades do modelo. A pesquisa foca na preparação e discretização dos dados, elementos decisivos para interpretar o comportamento do PIB. O estudo foi atualizado com maior intervalo temporal para que os efeitos dos outliers causados pela pandemia de Covid-19 fossem considerados no conjunto de treinamento. A metodologia é modular, permitindo melhorias contínuas, e mostrou que abordagens com menos categorias oferecem resultados práticos melhores, enquanto modelos mais sofisticados são úteis para uma análise mais detalhada dos dados. O cenário com 5 categorias modificado agora consegue prever altas fortes, mas ainda tem dificuldades nas classes de queda forte. A redução de variáveis aumentou a performance em todos os cenários, tanto em acurácia e especialmente em score-F1.

Diversos métodos de previsão, como máquinas de vetor de suporte e redes neurais, apresentaram dificuldades com muitas categorias, especialmente devido ao desbalanceamento nas categorias extremas. Em contraste, árvores de decisão, combinadas com florestas aleatórias, mostraram-se eficientes e interpretáveis. A quantidade de dados não foi limitante, mas mais dados poderiam melhorar os resultados. O ajuste de hiperparâmetros teve pouco impacto, indicando que os métodos empregados são robustos.

O estudo sugere várias direções para futuras investigações, como o uso de outras técnicas de aprendizagem de máquina (bagging, boosting), exploração de redes neurais profundas, análise de variáveis explicativas e modelos autoregressivos, além de técnicas com médias móveis para tendências sazonais. A metodologia pode ser aplicada a outras variáveis macroeconômicas e contextos internacionais.

Referências

- Alexander, G. J., Sharpe, W. F., and Bailey, J. V. (2001). *Fundamentals of investments*. Pearson Educación.
- Bhaumin, S. (2011 [Online].). Productivity and the economic cycle. BIS ECONOMICS PAPER NO. 12.
- BROWN, S. J. and DYBVIG, P. H. (1986). The empirical implications of the cox, ingersoll, ross theory of the term structure of interest rates. *The Journal of Finance*, 41(3):617–630.
- Burns, A. F. and Mitchell, W. C. (1946). *Measuring Business Cycles*. National Bureau of Economic Research, Inc.
- COX, J. C., INGERSOLL JR., J. E., and ROSS, S. A. (1981). A re-examination of traditional hypotheses about the term structure of interest rates. *The Journal of Finance*, 36(4):769–799.

- Estrella, A. and Mishkin, F. S. (1995). Predicting u.s. recessions: Financial variables as leading indicators. Working Paper 5379, National Bureau of Economic Research.
- Gogas, P., Papadimitriou, T., Matthaiou, M., and Chrysanthidou, E. (2014). Yield curve and recession forecasting in a machine learning framework. *Computational Economics*, 45(4):635–645.
- Heath, D., Jarrow, R., and Morton, A. (1992). Bond pricing and the term structure of interest rates: A new methodology for contingent claims valuation. *Econometrica*, 60(1):77–105.
- Hicks, J. R. et al. (1975). Value and capital: An inquiry into some fundamental principles of economic theory. *OUP Catalogue*.
- HO, T. S. Y. and LEE, S.-B. (1986). Term structure movements and pricing interest rate contingent claims. *The Journal of Finance*, 41(5):1011–1029.
- Jacovides, A. (2008). Forecasting interest rates from the term structure: Support vector machines vs neural networks. Master's thesis, University of Nottingham.
- Ju, Y., Kim, C., and Shim, J. (1997). Genetic-based fuzzy models: Interest rate forecasting problem. *Computers & Industrial Engineering*, 33(3):561–564. Selected Papers from the Proceedings of 1996 ICC&IC.
- Kim, S. H. and Noh, H. J. (1997). Predictability of interest rates using data mining tools: A comparative analysis of korea and the us. *Expert Systems with Applications*, 13(2):85–95.
- Nelson, C. R. and Siegel, A. F. (1987). Parsimonious modeling of yield curves. *The Journal of Business*, 60(4):473–489.
- Oh, K. J. and Han, I. (2000). Using change-point detection to support artificial neural networks for interest rates forecasting. *Expert Systems with Applications*, 19(2):105–115.
- Palhares Júnior, E., de Araujo, A. M. T., de Souza, A. H., da Silva, N. G., and da Silva Souza, W. (2024). Ensemble of machine learning applied to economic cycles analysis: a comparative study using antecedent macroeconomic indicators for brazilian gdp prediction classification. *Revista Brasileira de Planejamento e Desenvolvimento*. Submetido para publicação.
- Svensson, L. E. O. (1994). Estimating and interpreting forward interest rates: Sweden 1992-1994. *IMF Working Papers*, 94(114):1.
- Vasicek, O. (1977). An equilibrium characterization of the term structure. *Journal of Financial Economics*, 5(2):177–188.
- Vela, D. (2013). Forecasting Latin-American yield curves: An artificial neural network approach. BORRADORES DE ECONOMIA 010502, BANCO DE LA REPÚBLICA.
- Zimmermann, H., Tietz, C., and Grothmann, R. (2002). Yield curve forecasting by error correction neural networks and partial learning. In Verleysen, M., editor, *ESANN 2002, 10th Euroean Symposium on Artificial Neural Networks, Bruges, Belgium, April 24-26, 2002, Proceedings*, pages 407–412.