

Adaptação do algoritmo LAX-Reg – *Lean Adaptive XGBoost Regressor* para suportar a regressão de fluxo de dados com mudança de conceito recorrente

Felipe Kenji Nishioka Jojima, Fabiano Baldo

INTRODUÇÃO

O crescimento de dispositivos conectados às redes impulsionou a demanda por algoritmos de aprendizado de máquina capazes de lidar com fluxos de dados ininterruptos e não estacionários (Grando, 2023). Portanto, o algoritmo deve aprender de forma incremental e processar grandes volumes de dados rapidamente (Ikononovska, Gama e Džeroski, 2015). Além disso, o aprendizado deve refletir as mudanças de conceito na distribuição dos dados, que ocorrem quando as propriedades estatísticas da variável alvo mudam de maneira imprevista (Liao e Wang, 2018). Um algoritmo que se destaca em precisão e velocidade de resposta para esse tipo de tarefa é o LAX-Reg, proposto por Grando (2023). Entretanto, em sua implementação foram utilizadas bibliotecas que tiveram sua atualização descontinuada, como, por exemplo, *scikit-multiflow* (Montiel *et al.*, 2018). Portanto, este trabalho deve como objetivo adaptar o LAX-Reg, que é um algoritmo baseado em XGBoost (Chen *et al.*, 2015), para suportar mudanças de conceito recorrentes utilizando novas bibliotecas, como as bibliotecas *River* (Montiel *et al.*, 2021) e *Maki Nage* (Picard, 2020).

DESENVOLVIMENTO

Foi feita uma análise e substituição dos métodos que utilizavam a biblioteca *scikit-multiflow* (Montiel *et al.*, 2018) do trabalho original pelos métodos da biblioteca *River* (Montiel *et al.*, 2021), buscando encontrar métodos equivalentes aos utilizados pela biblioteca antiga. A maior mudança foi a do método “*Evaluate Prequential*”, o qual testa o modelo de análise do fluxo de dados e, depois, treina-o com os dados usados no teste. Esse método foi implementado de duas maneiras, sendo a primeira delas utilizando apenas os métodos da biblioteca *River* e a segunda maneira utilizando o framework *Maki Nage* (Picard, 2020), visando melhorar o tempo de execução. Para lidar com a mudança de conceito, foi utilizada o ADWIN – *Adaptive Sliding Window* –, que monitora uma sequência de dados do fluxo usando janelas deslizantes (Barddal, 2019). Essa estratégia foi implementada nos algoritmos LAX-A, que é uma das versões do LAX-Reg (Grando, 2023) e AFXGB (Souza, Grando e Baldo, 2022). A partir da refatoração foram gerados os algoritmos LAX, LAX-T, LAX-A, AFXGB com extensão *_R* e com extensão *_M*, representando a implementação *River* e a *Maki Nage*, respectivamente.

RESULTADOS

Para os experimentos, foram considerados os mesmos cenários do trabalho de Grando (2023). Para fins comparativos, como a utilização do framework *Maki Nage* ou da biblioteca *River* não interferem nos resultados de Erro Médio Quadrático (MSE), apenas os resultados do *River* foram considerados. Sobre o MSE, os valores de erro aumentaram em comparação aos obtidos pelo algoritmo original de Grando (2023), com exceção do algoritmo AFXGB_R, como observado na Tabela 1, o que indica um possível gargalo na refatoração dos métodos de avaliação. O teste de Friedman indicou diferenças significativas entre os algoritmos com 95% de confiança, o que permitiu fazer o teste de Nemenyi, apresentado na Figura 1, que mostrou que nenhum algoritmo apresentou um desempenho significativamente diferente dos demais. Sobre o tempo de execução e a memória, houve um aumento de 2 a 3 vezes o valor obtido na implementação original nas duas aplicações, destacando que a implementação com o *Maki*

Nage – que visava um ganho no critério de tempo – apresentou os maiores valores nessa métrica, sendo de 10 a 50 milissegundos mais lento que a adaptação do *River*.

CONSIDERAÇÕES FINAIS

As adaptações dos algoritmos para bibliotecas mais recentes, com exceção do AFXGB, apresentaram uma queda no desempenho, enquanto a implementação do ADWIN nos algoritmos LAX mostrou uma melhora em relação aos demais algoritmos. Esses resultados mostram a necessidade de um aprimoramento nos métodos substituídos para torná-los mais competitivos, verificando possíveis gargalos presentes nos métodos de avaliação.

Palavras-chave: refatoração, fluxo de dados, regressão, mudança de conceito, xgboost.

ILUSTRAÇÕES

Tabela 1. Comparativo entre os valores de MSE dos algoritmos LAX, LAX-T, LAX-A e AFXGB nas execuções pelo Scikit-Multiflow e pelo River.

Dataset	LAX	LAX-T	LAX-A	AFXGB	LAX_R	LAX-T_R	LAX-A_R	AFXGB_R
metro	1738735 ⁽²⁾	1536842⁽¹⁾	2885753 ⁽⁷⁾	2009869 ⁽⁵⁾	2600338 ⁽⁶⁾	1750625 ⁽³⁾	3984337 ⁽⁸⁾	1813253 ⁽⁴⁾
waze	2,810 ⁽⁴⁾	2,890 ⁽⁵⁾	2,790 ⁽³⁾	3,350 ⁽⁸⁾	3,034 ⁽⁶⁾	3,086 ⁽⁷⁾	2,314⁽¹⁾	2,393 ⁽²⁾
bikes	589 ⁽⁶⁾	732 ⁽⁷⁾	2683 ⁽⁸⁾	360 ⁽⁴⁾	273 ⁽³⁾	523 ⁽⁵⁾	151⁽¹⁾	243 ⁽²⁾
barra	351⁽¹⁾	355 ⁽²⁾	439 ⁽³⁾	594 ⁽⁵⁾	922 ⁽⁸⁾	484 ⁽⁴⁾	625 ⁽⁷⁾	606 ⁽⁶⁾
costa	163⁽¹⁾	172 ⁽²⁾	207 ⁽³⁾	316 ⁽⁶⁾	426 ⁽⁷⁾	260 ⁽⁴⁾	567 ⁽⁸⁾	305 ⁽⁵⁾
fried1(A)	78 ⁽⁵⁾	33 ⁽²⁾	175 ⁽⁸⁾	49 ⁽⁴⁾	172 ⁽⁶⁾	34 ⁽³⁾	174 ⁽⁷⁾	32⁽¹⁾
fried1(G)	82 ⁽⁵⁾	36⁽¹⁾	135 ⁽⁶⁾	57 ⁽⁴⁾	182 ⁽⁷⁾	41 ⁽²⁾	189 ⁽⁸⁾	44 ⁽³⁾
reg(A)	3987⁽¹⁾	5022 ⁽⁴⁾	4728 ⁽³⁾	5208 ⁽⁶⁾	7316 ⁽⁷⁾	5028 ⁽⁵⁾	8113 ⁽⁸⁾	4322 ⁽²⁾
reg(G)	4173⁽¹⁾	5194 ⁽³⁾	4906 ⁽²⁾	5520 ⁽⁵⁾	7824 ⁽⁷⁾	6038 ⁽⁶⁾	8544 ⁽⁸⁾	5280 ⁽⁴⁾

Fonte: O autor.

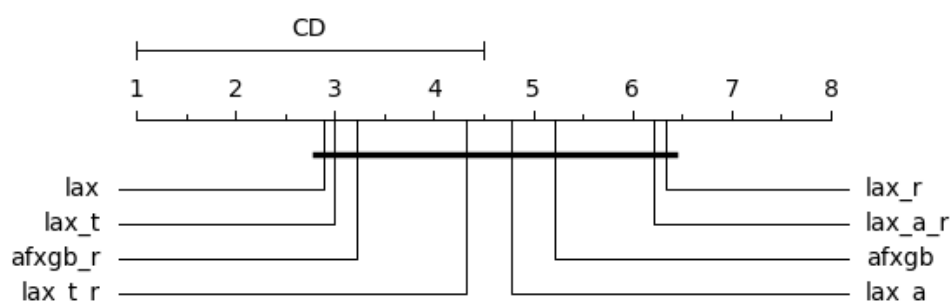


Figura 1. Teste de Nemenyi das classificações médias dos algoritmos para o MSE. Fonte: O autor.

REFERÊNCIAS BIBLIOGRÁFICAS

BARDDAL, J. P. Vertical and horizontal partitioning in data stream regression ensembles. In: IEEE. **2019 International Joint Conference on Neural Networks (IJCNN)**. Curitiba, BR, 2019. p. 1–8.

CHEN, T. et al. Xgboost: extreme gradient boosting. **R package version 0.4-2**, v. 1, n. 4, p. 1–4, 2015.

GRANDO, Julia. **Adaptação ao XGBoost para suporte à regressão de fluxo de dados não-estacionários**. 60 p. Joinville, 2023. Mestrado em computação aplicada – Universidade do Estado de Santa Catarina, Centro de Ciências Tecnológicas. 04 jul. 2023.

HALFORD, M.; BOLMIER, G.; SOURTY, R.; VAYSSE, R.; ZOUITINE, A. **Crepe, a Python library for online machine learning**. 2019. Disponível em: <<https://github.com/MaxHalford/crepe>>. Acesso em: 16 ago. 2025.

IKONOMOVSKA, E.; GAMA, J.; DŽEROSKI, S. Online tree-based ensembles and option trees for regression on evolving data streams. **Neurocomputing**, v. 150, p. 458–470, 2015. ISSN 0925-2312. Special Issue on Information Processing and Machine Learning for Applications of Engineering Solving Complex Machine Learning Problems with Ensemble Methods Visual Analytics using Multidimensional Projections. Disponível em: <<https://www.sciencedirect.com/science/article/abs/pii/S0925231214012338>>.

LIAO, Z.; WANG, Y. Rival learner algorithm with drift adaptation for online data stream regression. In: **Proceedings of the 2018 International Conference on Algorithms, Computing and Artificial Intelligence**. New York, NY, USA: Association for Computing Machinery, 2018. (ACAI 2018). ISBN 9781450366250. Disponível em: <<https://dl.acm.org/doi/10.1145/3302425.3302475>>.

MONTIEL, J.; READ, J.; BIFET, A.; ABDESSALEM, T. Scikit-multiflow: A multi-output streaming framework. **The Journal of Machine Learning Research**, 19, 72, p. 1-5, 2018. Disponível em: <<http://jmlr.org/papers/v19/18-251.html>>. Acesso em: 16 ago. 2025.

MONTIEL, J.; HALFORD, M.; MASTELINI, S. M.; BOLMIER, G.; SOUTRY, R.; VAYSSE, R.; ZOUITINE, A.; GOMES, H. M.; READ, J.; ABDESSALEM, T.; BIFET, A. River: machine learning for streaming data in Python. **The Journal of Machine Learning Research**, 22, 110, p. 1-8, 2021. Disponível em: <<http://jmlr.org/papers/v22/20-1380.html>>. Acesso em: 16 ago. 2025.

PICARD, R. **Maki Nage**: Stream Processing made Easy. 2021. Disponível em: <<https://www.makinage.org/>>. Acesso em: 16 ago. 2025.

SOUZA, F. M. de; GRANDO, J.; BALDO, F. Adaptive fast XGBoost for regression. In: XAVIER-JUNIOR, J. C.; RIOS, R. A. (Ed.). **Intelligent Systems**. Cham: Springer International Publishing, 2022. p. 92–106. ISBN 978-3-031-21686-2

DADOS CADASTRAIS

BOLSISTA: Felipe Kenji Nishioka Jojima

MODALIDADE DE BOLSA: PROBIC/UDESC (IC)

VIGÊNCIA: 09/2024 a 08/2024 – Total: 12 meses

ORIENTADOR(A): Fabiano Baldo

CENTRO DE ENSINO: CCT

DEPARTAMENTO: DEPARTAMENTO DE CIENCIA DA COMPUTACAO CCT

ÁREAS DE CONHECIMENTO: Ciências Exatas e da Terra/ Ciência da Computação

TÍTULO DO PROJETO DE PESQUISA: StreamMining v.3.0 Novas Abordagens para Algoritmos de Aprendizagem em Fluxos de Dados Não Estacionários

Nº PROTOCOLO DO PROJETO DE PESQUISA: NPP3216-2023