

Identificação de Anomalias em Transações Financeiras e Avaliação de Ações Preventivas no Ecosistema de Pagamentos

Exame de Qualificação de Mestrado

Candidato: Marcos Vinícius Piaia
Orientador: Prof. Dr. Zanoni Dias
Coorientador: Prof. Dr. Hélio Pedrini

Instituto de Computação
Unicamp

29 de Fevereiro de 2024

Agenda

- 1 Introdução
- 2 Objetivos
- 3 Fundamentação Teórica
- 4 Metodologia
- 5 Métricas
- 6 Resultados Preliminares
- 7 Plano de Trabalho e Cronograma de Execução
- 8 Referências

- O sistema financeiro, especialmente o brasileiro, possui uma construção muito flexível quando consideramos os elementos que fazem parte do mesmo e sua organização.
- Cada vez mais, o sistema vem se tornando mais acessível a diversas camadas da população, e essa integração se dá, principalmente, através do universo digital.
- As formas pelas quais o dinheiro transita dentro deste ecossistema tornaram-se mais rápidas e dinâmicas. Pode-se considerar aqui as diferenças e a evolução dos meios de pagamento entre TED [2] e Pix [1] como exemplos.

- Ecossistema complexo e altamente dinâmico.
- Impacto econômico e social.
- Necessidade de segurança e confiabilidade.

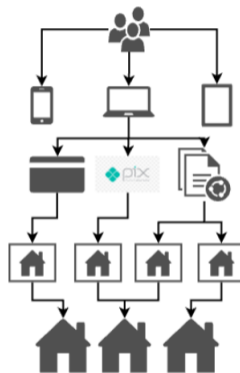


Figura: Ecossistema de pagamentos.

Objetivo Geral

Desenvolvimento de um fluxo de operações para garantir a rápida identificação de transações suspeitas e a automatização do tratamento destas.

Objetivos Específicos

- Avaliar a eficácia da aplicação de diferentes técnicas de aprendizado de máquina na identificação de fraudes, tratando-as como anomalias.
- Desenvolver um sistema de alerta de transações suspeitas.
- Avaliar a eficácia do sistema de alerta proposto.

Hipóteses do Trabalho

- Um dos modelos propostos será eficaz na identificação de anomalias nas várias bases consideradas.
- Este modelo será capaz de generalizar a identificação de transações fraudulentas sem a necessidade da introdução de dados sintéticos para o balanceamento da base.
- Na busca por evitar o uso de métodos de geração de dados sintéticos uma alternativa é a utilização de algoritmos focados na identificação do comportamento característico da normalidade nas transações.

Identificação de Anomalias

- Em geral, o evento que caracteriza a anomalia se mantém desconhecido até o seu acontecimento.
- Essas anomalias podem se apresentar agrupadas em classes bastante heterogêneas, dificultando sua generalização.
- Anomalias são eventos naturalmente raros, criando bases de dados desbalanceadas.

Identificação de Anomalias

- Considerando-se a existência de tipos muito diversos de anomalias, ressaltam-se os seguintes:
 - Anomalias pontuais, onde instâncias individuais dos dados são consideradas anômalas com respeito ao restante da base.
 - Anomalias condicionais, nas quais ainda se consideram instâncias individuais dos dados como sendo anômalos não com relação a uma característica de toda base, mas sim, como consequência gerada a partir de um contexto, caracterizado por uma combinação de outros eventos da mesma base.
 - Anomalias de grupo, nos quais se tratam grupos de dados como anômalos com relação a outros da mesma base. Por exemplo, um conjunto de contas falsas em uma plataforma de mídia social. A análise individual destas contas pode mostrar perfis válidos com ações válidas, enquanto estas, agrupadas, revelam uma ação fraudulenta.

Desafios

Considerando as características do problema, Pang et al. [6] resumiram os seguintes desafios:

- Baixa taxa de recuperação de detecção de anomalias.
- Detecção de anomalias em bases de dados de grande dimensionalidade e/ou dados não-independentes.
- Aprendizado eficiente da normalidade/anormalidade dos dados.
- Resiliência ao ruído dos dados.
- Detecção de anomalias complexas.
- Explicabilidade de uma identificação.

Pang et al. [6] avaliaram um conjunto de soluções para o tratamento de anomalias resumindo-as da seguinte forma:

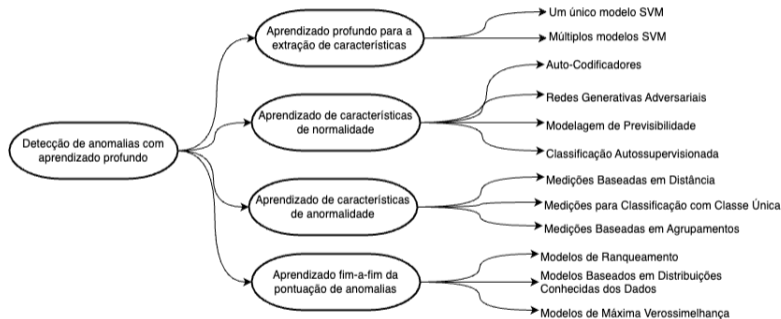


Figura: Distribuição de algoritmos proposta por Pang et al. [6].

Identificação de Fraudes

Tratando fraudes no sistema financeiro como anomalias:

- Dados desbalanceados, com poucas referências de anomalias em comparação ao montante de dados.
- Dados de baixa dimensionalidade visto que na maioria das bases disponíveis não encontramos grandes vetores de características identificando cada transação.
- Dados tabulares, derivados de transações de usuários.

Identificação de Fraudes

Considerando os itens anteriores, as seguintes abordagens se destacaram na Identificação de Fraudes:

- Regressão Logística [4].
- Rede de Desvios [7].
- Aumento de Gradiente Extremo [3].
- *Transformers* [8].

Identificação de Fraudes Utilizando Regressão Logística

Comparativo do desempenho entre os três algoritmos, demonstrando a melhor eficácia da Regressão Logística, em relação ao KNN e ao Naïve Bayes, como proposto por Itoo et al. [4].

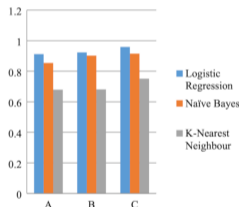


Figura: Acurácia apresentada na base proposta por Itoo et al. [5].

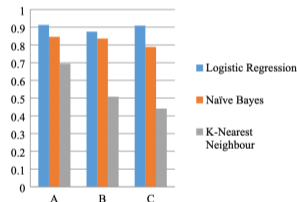


Figura: F1 apresentado na base proposta por Itoo et al. [5].

Identificação de Fraudes com Rede de Desvios

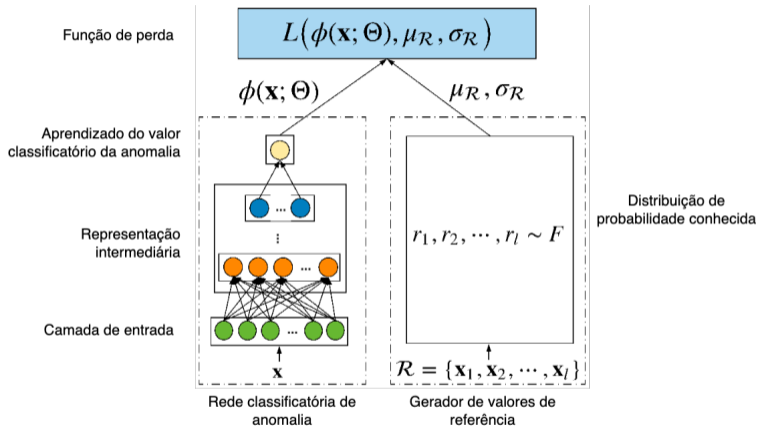


Figura: Modelo reproduzido a partir do proposto por Pang et al. [7].

Identificação de Fraudes com Aumento de Gradiente Extremo

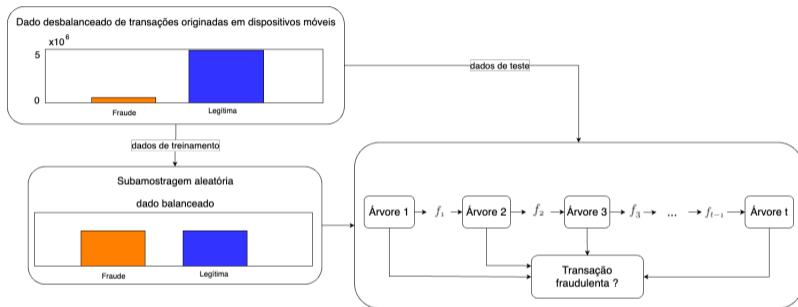


Figura: Fluxo do Aumento de Gradiente Extremo (XGBoost) com Subamostragem Aleatória (RUS) para a identificação de fraude [3].

Identificação de Fraudes com Transformers

Alteração do modelo tradicional de *Transformers* para a utilização nos grafos de transações propostos por Zhang et al. [8] pode ser vista na figura abaixo.

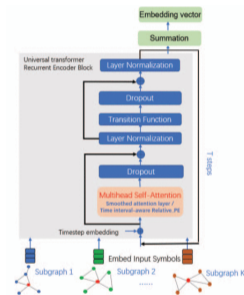


Figura: Arquitetura do DynGraphTrans [8] com duas modificações nos módulos de auto-atenção em relação à abordagem tradicional de *Transformers*: 1) uma camada extra de atenção suavizada 2) codificação posicional levando em consideração a abordagem de intervalos temporais nos dados.

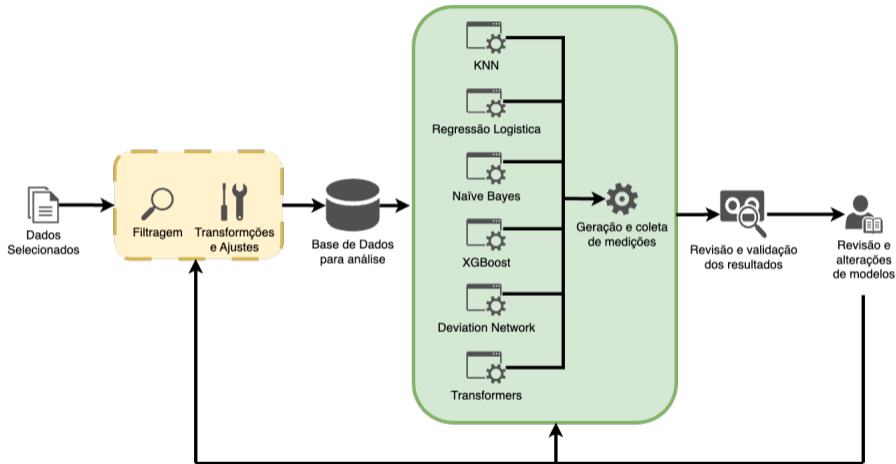


Figura: Fluxo proposto para o estudo.

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}} \quad (1)$$

$$\text{Precisão} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (2)$$

$$\text{Revocação} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (3)$$

$$\text{AUC} = \int_0^1 \text{Revocação}(T) \times \frac{d}{dT} \text{FPR}(T) dT \quad (4)$$

$$\text{F1} = 2 \times \frac{\text{Precisão} \times \text{Revocação}}{\text{Precisão} + \text{Revocação}} \quad (5)$$

- Experimentos iniciais com a base de dados proposta por Itoo et al. [5].

Algoritmo	Classe	Precisão	Revocação	F1
KNN	Legítima	1.00	1.00	1.00
	Fraude	0.94	0.74	0.83
Regressão Logística	Legítima	1.00	1.00	1.00
	Fraude	0.88	0.63	0.74

Tabela: Precisão, Revocação e F1 para a base de dados *Credit Card Fraud Detection* utilizando KNN e Regressão Logística.

Plano de Trabalho e Cronograma de Execução

Atividades	1 ^o ano						2 ^o ano					
	1	2	3	4	5	6	1	2	3	4	5	6
Etapa 1 - Preparação												
Disciplinas de Pós-Graduação	•	•	•	•	•	•						
Pesquisa bibliográfica	•	•	•	•	•	•	•	•	•			
Preparação da base de dados							•	•	•			
Etapa 2 - Implementações												
Implementação dos algoritmos							•	•	•			
Treinamento e validação dos modelos							•	•	•			
Verificação dos resultados								•	•	•		
Etapa 3 - Testes												
Aplicação dos modelos em múltiplas bases								•	•	•		
Revisão dos resultados								•	•	•		
Comparação com outros trabalhos								•	•	•		
Etapa 4 - Conclusão												
Publicação dos resultados										•	•	
Escrita da dissertação										•	•	•
Defesa da dissertação												•

Tabela: Cronograma de atividades dividido em bimestres.

Referências I

- [1] B. C. do Brasil.
O que é Pix?
<https://www.bcb.gov.br/estabilidadefinanceira/pix>.
Accessed at 2023-06-04.
- [2] B. C. do Brasil.
TED, DOC e book transfer: entenda como funcionam os tipos de transferências entre contas.
<https://www.bcb.gov.br/detalhenoticia/327/noticia>.
Accessed: 2023-06-04.
- [3] P. Hajek, M. Z. Abedin, and U. Sivarajah.
Fraud detection in mobile payment systems using an XGBoost-based framework.
Information Systems Frontiers, 2022.
- [4] F. Itoo, Meenakshi, and S. Singh.
Comparison and analysis of logistic regression, Naïve Bayes and KNN machine learning algorithms for credit card fraud detection.
International Journal of Information Technology, 13:1503–1511, 2021.
- [5] Kaggle.
Credit Card Fraud Detection.
<https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud>.
Accessed at 2023-07-16.
- [6] G. Pang, C. Shen, L. Cao, and A. V. D. Hengel.
Deep learning for anomaly detection: A review.
ACM Computing Surveys, 54(2):1–38, 2021.

Referências II

- [7] G. Pang, C. Shen, and A. van den Hengel.
Deep anomaly detection with deviation networks.
In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 353–362, 2019.
- [8] S. Zhang, T. Suzumura, and L. Zhang.
Dyngraphtrans: Dynamic Graph Embedding via Modified Universal Transformer Networks for Financial Transaction Data.
In *IEEE International Conference on Smart Data Services (SMDS)*, pages 184–191, 2021.

Obrigado.