

Construção de um Modelo Orientado a Dados para Detecção de Fraudes em Cartões de Crédito utilizando Dados Sintéticos

Alexandre dos Santos, Roger Passos, Luis Tarrataca, Douglas Cardoso,
Diego Haddad, Felipe Henriques

PPCIC, Cefet/RJ - Rio de Janeiro - RJ - Brasil
Center of Linguistics, University of Porto - Porto - Portugal

Setembro de 2024

Sumário

1. Introdução
2. Simulador de Dados
3. Modelo Orientado a Dados
4. Avaliação Experimental
5. Conclusão

Sumário

1. Introdução
2. Simulador de Dados
3. Modelo Orientado a Dados
4. Avaliação Experimental
5. Conclusão

Introdução

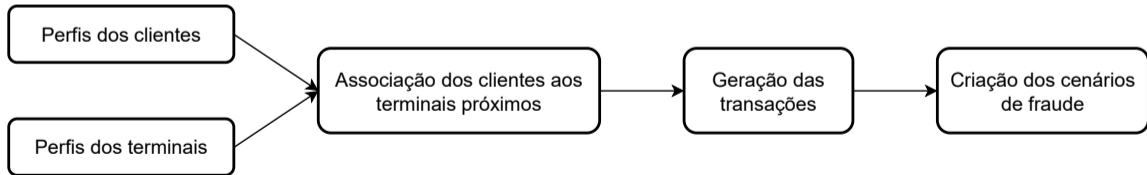
- Justificativa
 - Fraudes em cartões de crédito acarretam grandes prejuízos às empresas e à sociedade em geral
 - É um problema complexo que exige uma abordagem sistemática
 - A escassez de bases de dados públicas prejudica a reprodutibilidade de experimentos nessa área
- Objetivos
 - Construção de um simulador de dados sintéticos de transações utilizando como base um simulador inicial
 - Construção de modelos orientados a dados para identificação das fraudes
 - Avaliação experimental utilizando diferentes métricas de desempenho

Sumário

1. Introdução
2. Simulador de Dados
3. Modelo Orientado a Dados
4. Avaliação Experimental
5. Conclusão

Descrição do Simulador

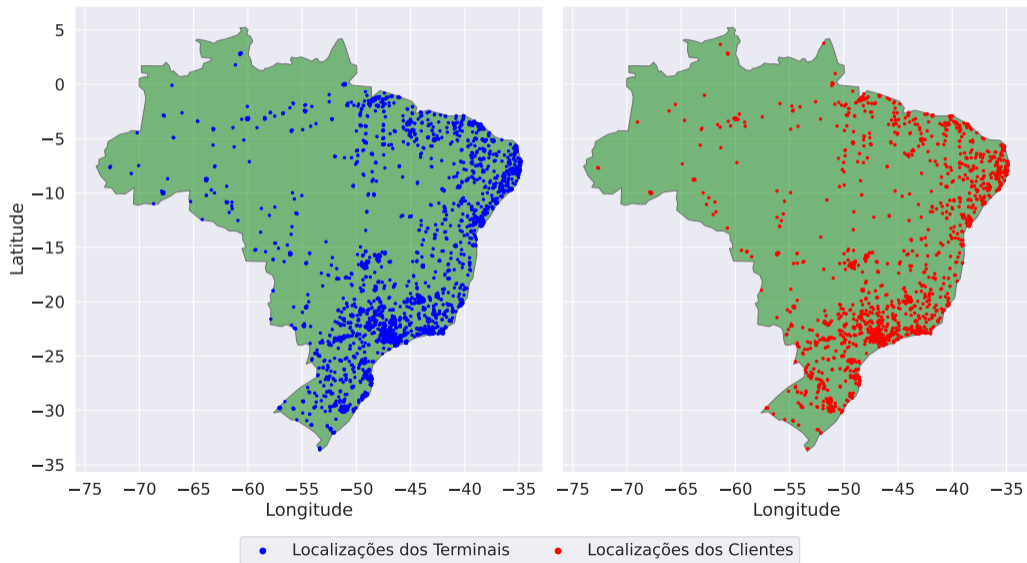
1. Tabela com propriedades dos clientes
2. Tabela com propriedades dos terminais
3. Obtenção dos terminais disponíveis para transações CP de cada cliente
4. Tabela de transações
5. Adição dos rótulos a partir de 4 cenários de fraude



Contribuições – Localização Geográfica

- Coordenadas geográficas reais do Brasil
- Municípios mais populosos → maior número de clientes e terminais
- Escolhido o município → Localização dentro de um raio de 12 km do centro geográfico do município. Valor definido utilizando uma aproximação da área mediana dos municípios brasileiros
- Tipos de Localizações
 - Localização do cliente (ou de cobrança)
 - Localização do terminal
 - Localização de entrega

Contribuições – Localização Geográfica



Contribuições – Tipo da Transação

- Diferenciação entre transações CP e CNP
- Afeta:
 - Localização do terminal
 - Localização de entrega
 - Cenários de fraude
 - Horário das transações

Contribuições – Cenários de Fraude

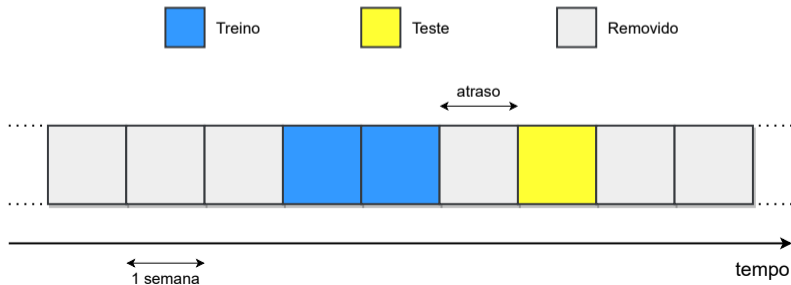
- Adição de um novo cenário de fraude
 - Cartão perdido ou roubado
 - Fraudador gasta o máximo possível o mais rápido possível, antes do cartão ser bloqueado
- Modificações dos atributos das transações marcadas como fraude
 - Quantia
 - Tipo da transação
 - Data
 - Localizações

Sumário

1. Introdução
2. Simulador de Dados
3. Modelo Orientado a Dados
4. Avaliação Experimental
5. Conclusão

Engenharia de Atributos e Divisão dos Dados

- Engenharia de Atributos
 - Derivados diretamente
 - Agregação
 - Codificação de risco
- Divisão dos Dados
 - Levou em consideração: Mudança de Contexto e o *Feedback* atrasado
 - Divisão Prequencial: Treino/Validação e Treino/Teste



Treino e Validação dos Modelos

- Algoritmos de Aprendizado de Máquina
 - **Classificação:** *Random Forest, Logistic Regression, K-Nearest Neighbors*
 - **Detecção de Anomalias:** *Isolation Forest, Elliptic Envelope*
- Estratégias de pré-processamento
- Otimização de Hiperparâmetros
- Melhor modelo escolhido a partir do desempenho na métrica AP de Validação

Sumário

1. Introdução
2. Simulador de Dados
3. Modelo Orientado a Dados
- 4. Avaliação Experimental**
5. Conclusão

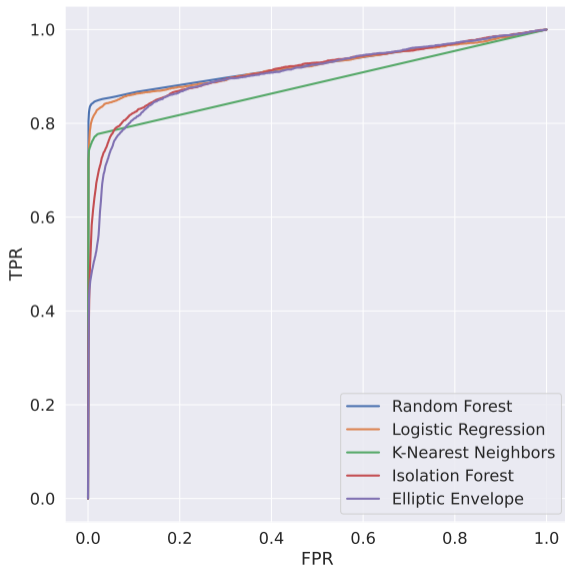
Dependentes de Limiar

- TPR (*Recall*)
- Precisão
- FPR
- G-Mean
- F1-Score

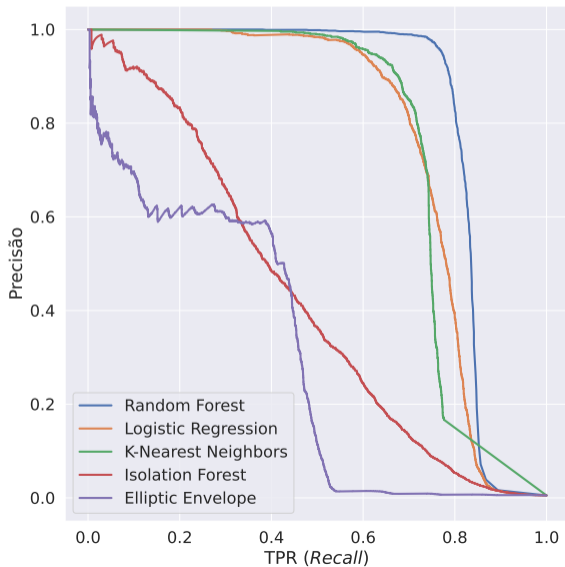
Independentes de Limiar

- Curva ROC e AUC-ROC
- Curva *Precision Recall* e AP

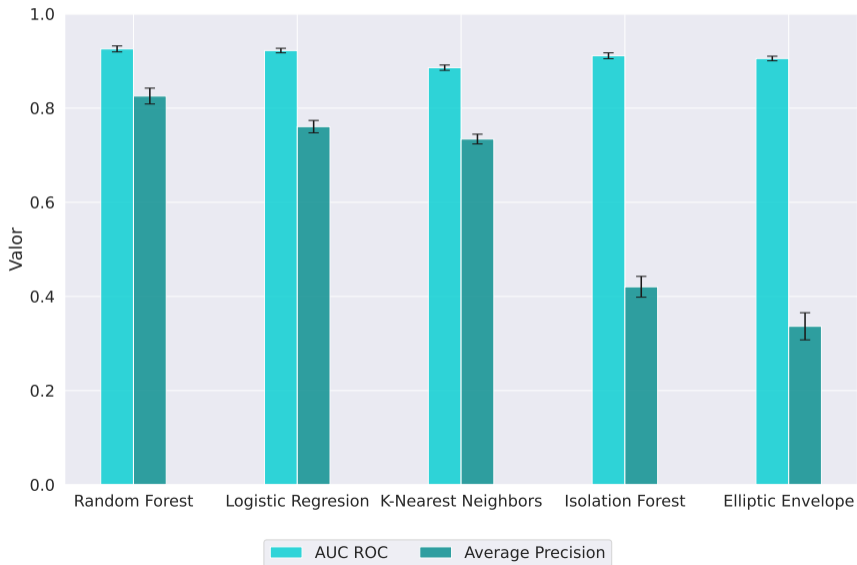
Comparação de Desempenho – Curva ROC



Comparação de Desempenho – Curva PR



Comparação de Desempenho – Áreas sob as Curvas



Comparação de Desempenho – Max F1-Score

- Limiar que maximiza a métrica F1-Score

Algoritmo	Métricas				
	FPR	TPR	Precisão	G-Mean	F1-Score
<i>Random Forest</i>	0.0 ± 0.0	0.767 ± 0.017	0.964 ± 0.013	0.876 ± 0.01	0.854 ± 0.014
<i>Logistic Regression</i>	0.001 ± 0.0	0.685 ± 0.031	0.864 ± 0.043	0.827 ± 0.019	0.763 ± 0.012
<i>K-Nearest Neighbors</i>	0.0 ± 0.0	0.685 ± 0.02	0.9 ± 0.024	0.828 ± 0.012	0.778 ± 0.012
<i>Isolation Forest</i>	0.002 ± 0.001	0.425 ± 0.033	0.474 ± 0.042	0.65 ± 0.025	0.446 ± 0.018
<i>Elliptic Envelope</i>	0.002 ± 0.0	0.413 ± 0.038	0.587 ± 0.019	0.641 ± 0.03	0.484 ± 0.028

Sumário

1. Introdução
2. Simulador de Dados
3. Modelo Orientado a Dados
4. Avaliação Experimental
5. Conclusão

Considerações Finais

- Todas as etapas para a elaboração de modelos orientados a dados foram desenvolvidas
- Reprodutibilidade dos experimentos
- Desempenho inferior dos detectores de anomalias
- Trabalhos futuros
 - Simulador de dados
 - Modelos orientados a dados

Construção de um Modelo Orientado a Dados para Detecção de Fraudes em Cartões de Crédito utilizando Dados Sintéticos

Alexandre dos Santos, Roger Passos, Luis Tarrataca, Douglas Cardoso,
Diego Haddad, Felipe Henriques

PPCIC, Cefet/RJ - Rio de Janeiro - RJ - Brasil
Center of Linguistics, University of Porto - Porto - Portugal

Setembro de 2024