



Análise da Detecção de Fraudes em Transações Financeiras com Modelos de Aprendizado de Máquina

Frederico dos Santos, Gabriel de
Souza, Gabriel Lima de Souza, Nikolas
Louret

Sumário

01

introdução

02

metodologia

03

resultados

04

conclusão

01

introdução

Contexto

Crescimento das fraudes financeiras com o aumento das transações online

Impactos negativos na estabilidade econômica e confiança dos clientes

Desafio crescente devido à complexidade das operações financeiras

Métodos de Detecção de Fraudes Bancárias

Detecção de Uso Indevido:

- classifica transações com base em padrões conhecidos de fraude

Detecção de Anomalias:

- identifica desvios em relação ao comportamento normal com base em dados históricos.

Objetivo

Desenvolver um modelo de Machine Learning para detecção de fraudes com base em dados reais (Kaggle)

02

Metodologia

matplotlib



seaborn



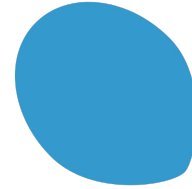
NumPy



pandas

kaggle

Google colab



scikit

learn

FERRAMENTAS

Configuração do Ambiente e Carregamento dos Dados

Configuração do Ambiente:

- Google Colab

Carregamento do Dataset:

- Fonte de dados: Kaggle – **Credit Card Fraud Detection**
- Utilização da biblioteca Pandas para importar os dados em um DataFrame

Análise Inicial do Dataset:

- Transações anonimizadas
- Dados reais de empresas europeias (Setembro de 2013)

Distribuição das Classes:

- Não fraudulentas: 99,828%
- Fraudulentas: 0,172%

Pré-processamento e Balanceamento dos Dados

Análise inicial:

- Verificação de valores ausentes ou duplicados.

Distribuição inicial

- Desequilíbrio entre classes identificada

Balanceamento

- Técnica aplicada: **Undersampling**
- **Motivação:** Equilibrar as classes e evitar vieses do modelo

Resultado (conjunto baseado em 984 transações):

- **Não fraudulentas:** 50%
- **Fraudulentas:** 50%

Modelos de Aprendizado de Máquina

- **Random Forest**
- **Support Vector Machine (SVM)**
- **Logistic Regression**
- **Naive Bayes**
- **K-Nearest Neighbors (KNN)**

objetivo: comparar o desempenho na detecção de fraudes

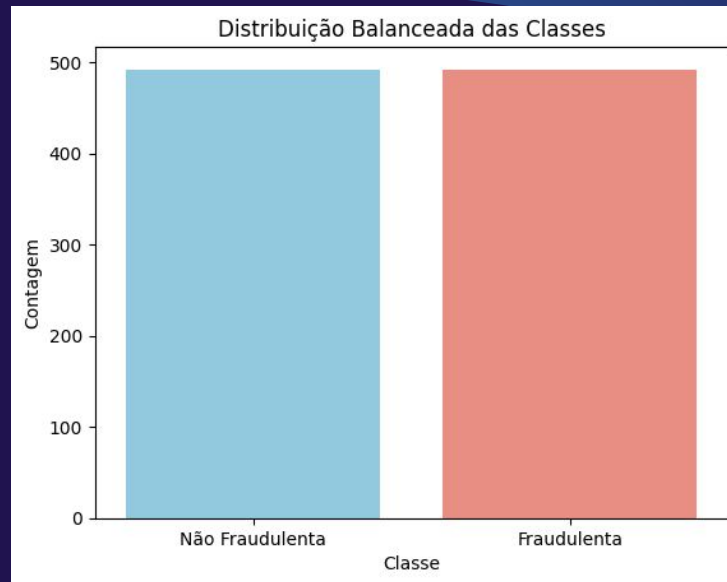
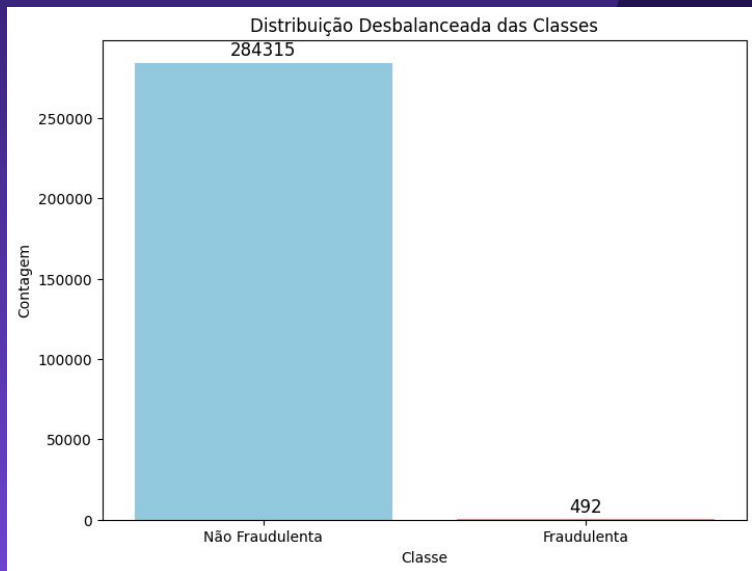
Avaliação dos Modelos

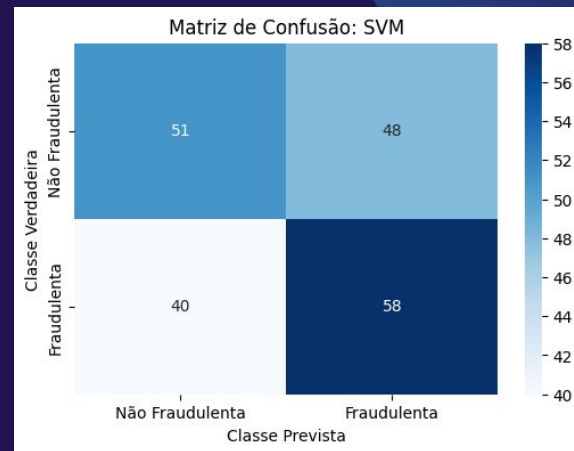
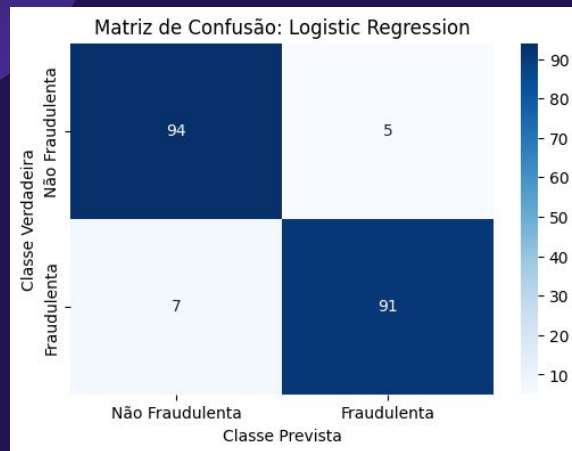
- **Acurácia**
- **Matriz de confusão**
- **Precisão, revocação e F1-score**
- **AUC da curva ROC**

justificativa: análise ampla e detalhada do desempenho

03

Resultados

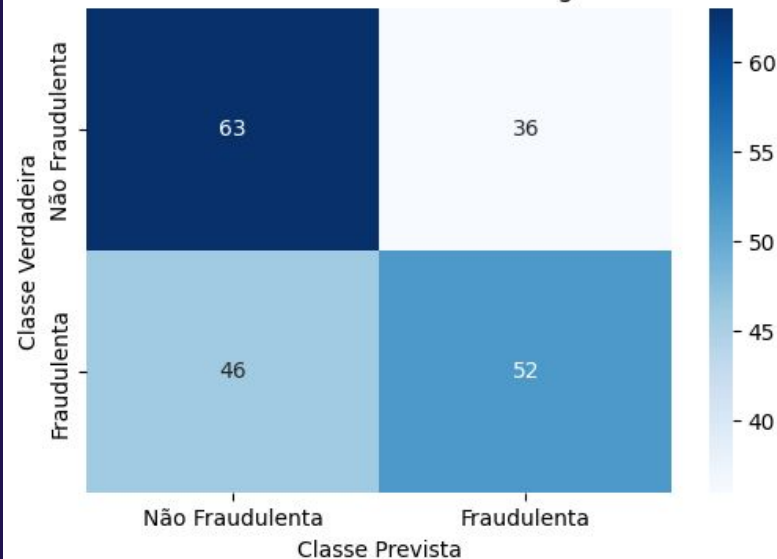




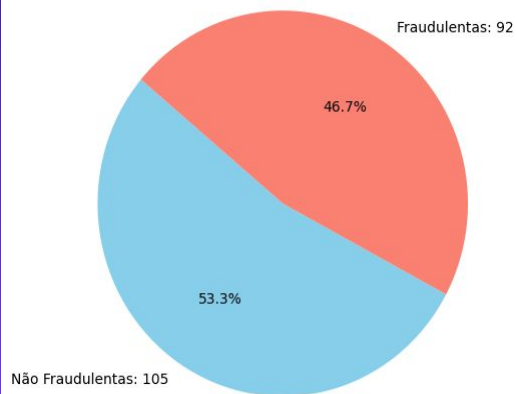
Matriz de Confusão: Naive Bayes



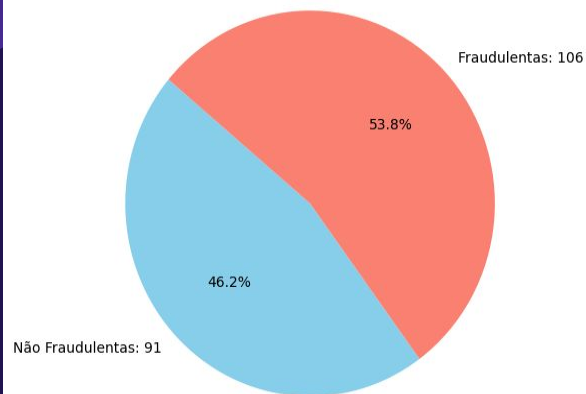
Matriz de Confusão: K-Nearest Neighbors



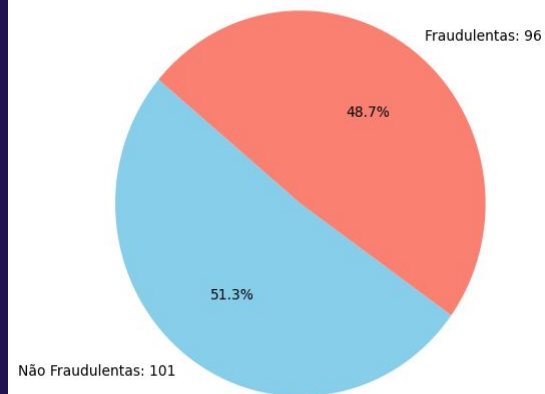
Distribuição de Previsões: Random Forest



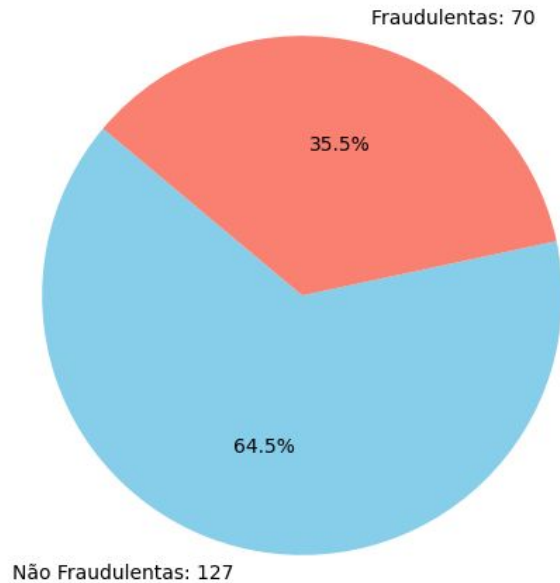
Distribuição de Previsões: SVM



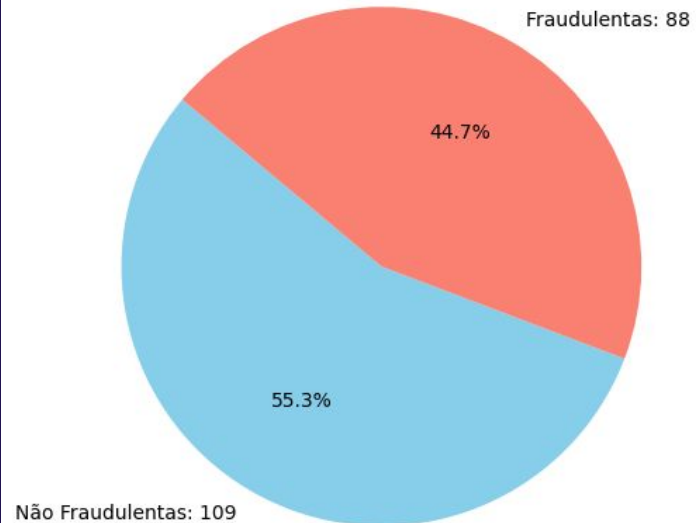
Distribuição de Previsões: Logistic Regression



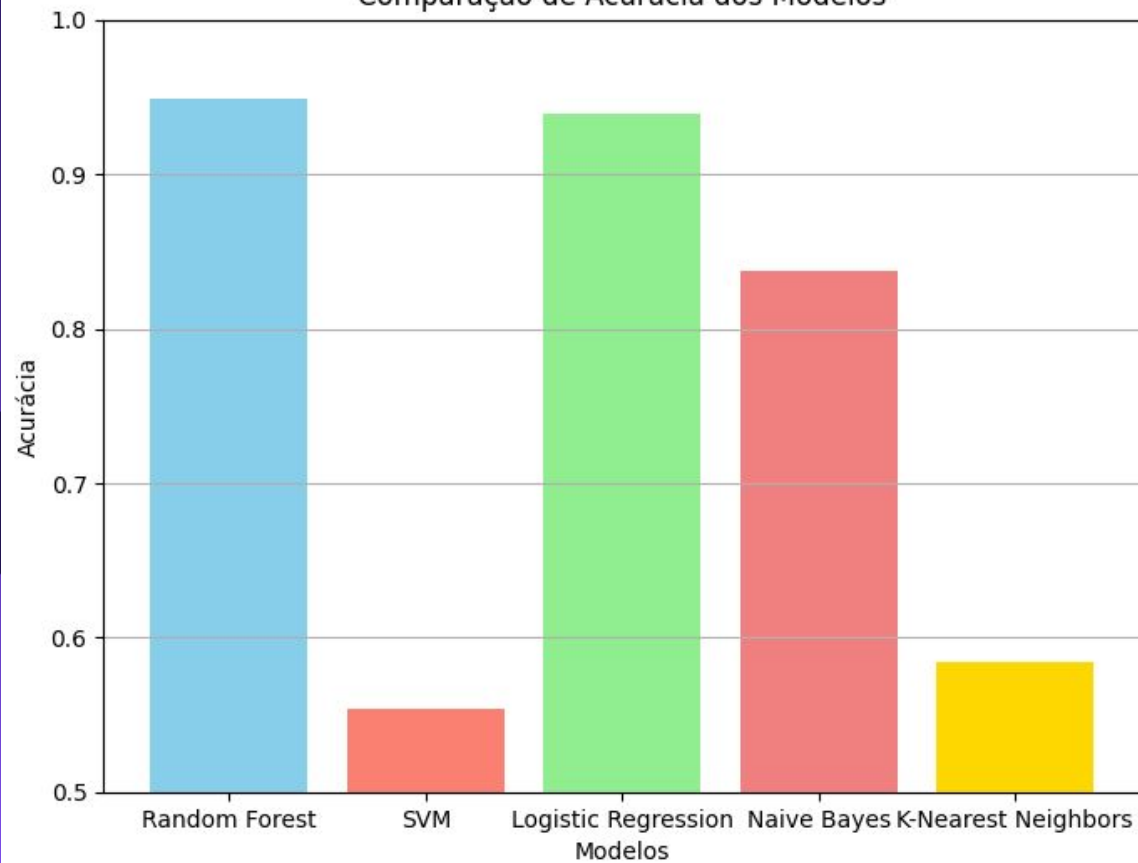
Distribuição de Previsões: Naive Bayes



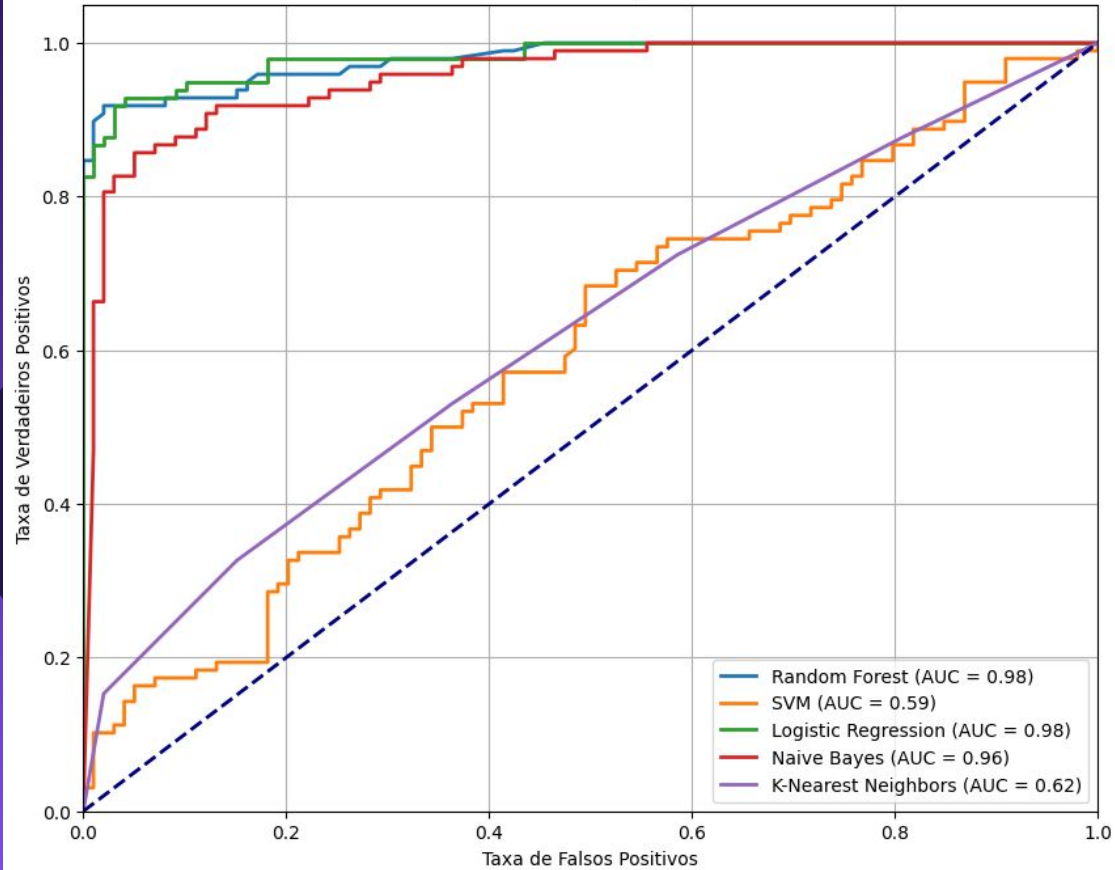
Distribuição de Previsões: K-Nearest Neighbors



Comparação de Acurácia dos Modelos



Curva ROC Comparativa



04

Conclusão

Desempenho dos Modelos e Impacto do Balanceamento

Desempenho dos Modelos

- **Random Forest (RF)**: Melhor desempenho geral (AUC > 0,97, acurácia de 94,92%, F1-score de 93,69%).
- **Logistic Regression (LR)**: Alta acurácia e equilíbrio entre precisão e revocação.
- **SVM e KNN**: Desempenho insatisfatório (acurácia < 60%).
- **Naive Bayes (NB)**: Boa revocação, mas baixa precisão (muitos falsos positivos).

Impacto do Balanceamento (Undersampling)

- **RF e LR**: Menos sensíveis à redução da classe majoritária, mantendo boa generalização.
- **SVM e KNN**: Desempenho comprometido pela alteração na distribuição das classes.

Conclusões Gerais e Trabalhos Futuros

Conclusões Gerais

- RF e LR são adequados para detectar fraudes em bases desbalanceadas.
- Técnicas de balanceamento e escolha do modelo impactam diretamente os resultados.

Trabalhos Futuros

- Explorar técnicas avançadas de balanceamento (como o SMOTE).
- Avaliar modelos mais complexos, como redes neurais, para maior precisão e robustez.

Obrigado