

Fluxo de Trabalho Detalhado para Projeto de Detecção de Fraudes e Previsão de Tendências

Este documento descreve de forma detalhada e explicativa cada etapa do pipeline de Data Science e Inteligência Artificial, desde a ingestão dos dados até a apresentação dos resultados finais.

Nota de Escopo

Importante: O treinamento de um modelo de linguagem do zero não fará parte da fase inicial do projeto, pois exige alto consumo de recursos e grande volume de dados. Na etapa inicial, iremos focar em **fine-tuning** de modelos pré-treinados e em técnicas de quantização e grounding em fases posteriores.

1. Carregar Dados

Objetivo: Disponibilizar os dados brutos em memória para análise e processamento.

Ações:

- Importar bibliotecas essenciais (`pandas`, `json`, `glob`).
- Ler o arquivo `cards_data.csv` para um DataFrame `df_cards`:

```
df_cards = pd.read_csv("/mnt/data/cards_data.csv")
```

- Ler o arquivo `users_data.csv` para um DataFrame `df_users`:

```
df_users = pd.read_csv("/mnt/data/users_data.csv")
```

- Carregar o JSON de mapeamento MCC para um dicionário `mcc_map`:

```
with open("/mnt/data/mcc_codes.json", "r") as f:  
    mcc_map = json.load(f)
```

- Verificar integridade básica: exibir dimensões e primeiras linhas de cada DataFrame.
-

2. Explorar Dados (EDA)

Objetivo: Compreender a estrutura, qualidade e relacionamentos dos conjuntos de dados.

Ações:

- **Visão Geral:** usar `df.info()` e `df.describe()` para cada DataFrame.
- **Distribuição de Variáveis:** histogramas de valores de transação, contagem de transações por usuário e por categoria MCC.

- **Valores Ausentes:** calcular porcentagem de NaN em cada coluna.
 - **Relações Básicas:** cruzar df_cards e df_users por user_id para verificar correspondência de usuários.
 - **Mapeamento de MCC:** aplicar mcc_map para converter códigos em descrições legíveis e contar frequência por categoria.
-

3. Pré-processar Dados

Objetivo: Garantir qualidade e consistência para modelagem.

Ações:

1. **Tratamento de Valores Ausentes:**
 2. Remoção ou imputação (média/mediana/moda) conforme relevância da coluna.
 3. **Correção de Tipos:**
 4. Converter colunas de data (transaction_date) para datetime .
 5. Garantir formato numérico em valores monetários.
 6. **Encoder de Categorias:**
 7. Transformar colunas categóricas (ex.: merchant_category) usando OneHotEncoder ou LabelEncoder .
 8. **Feature Engineering Inicial:**
 9. Extrair ano , mês , dia_semana , hora a partir de transaction_date .
 10. Criar indicador is_weekend e is_night_transaction .
-

4. Mesclar Dados

Objetivo: Consolidar informações de cartões, usuários e categorias MCC.

Ações:

- Unir df_cards e df_users por user_id (inner join).
 - Incluir coluna merchant_category_desc usando mapeamento de MCC.
 - Resultado: DataFrame df_full contendo atributos transacionais, demográficos e descritivos.
-

5. Engenharia de Recursos para Detecção de Fraudes

Objetivo: Criar variáveis que capturem padrões suspeitos.

Ações:

- **Velocidade de Transação:** tempo entre transações consecutivas do mesmo usuário.
 - **Recursos Baseados em Valor:** razão entre valor da transação atual e média histórica do usuário.
 - **Recursos Temporais:** número de transações nas últimas 1h, 24h e 7 dias.
 - **Geolocalização (se disponível):** distância entre localizações de transações consecutivas.
-

6. Dividir Dados

Objetivo: Separar dados para treinamento e avaliação imparcial.

Ações:

- Definir variável alvo `is_fraud` (se disponível) ou usar rótulo fornecido.
 - Usar `train_test_split` do Scikit-learn (ex.: 70% treino / 30% teste), garantindo estratificação por `is_fraud`.
-

7. Treinar Modelo de Detecção de Fraudes

Objetivo: Aprender padrões que distinguem transações legítimas de fraudulentas.

Ações:

- Selecionar algoritmos iniciais: `RandomForestClassifier`, `XGBoostClassifier`.
 - Ajustar hiperparâmetros com `GridSearchCV` ou `RandomizedSearchCV`.
 - Treinar modelo em `X_train`, `y_train`.
-

8. Avaliar Modelo de Detecção de Fraudes

Objetivo: Medir desempenho e identificar ajustes.

Ações:

- Prever em `X_test` e calcular métricas: `precision`, `recall`, `f1-score`, `ROC AUC`.
 - Analisar curva ROC e curva de precisão-recall.
 - Matriz de confusão para visualizar tipos de erro.
-

9. Identificar Padrões e Insights

Objetivo: Entender características de fraudes e recomendações.

Ações:

- Avaliar importância de features (`feature_importances_`).
 - Investigar transações mal classificadas para descobrir novos padrões.
 - Relacionar perfis de usuário com maior risco.
-

10. Engenharia de Recursos para Previsão de Séries Temporais

Objetivo: Preparar dados para modelagem de tendências ao longo do tempo.

Ações:

- Agregar contagem e valor total de transações por período (mensal/trimestral).
 - Criar variáveis defasadas (`lag_1` , `lag_3` , `lag_12`).
 - Incluir indicadores sazonais (`mês` , `trimestre`).
-

11. Treinar Modelo de Previsão de Séries Temporais

Objetivo: Prever volume e frequência de fraudes até 2040.

Ações:

- Selecionar modelos: `Prophet` , `ARIMA` , `LSTM` .
 - Ajustar parâmetros sazonais e de tendência.
 - Treinar com dados agregados até o presente.
-

12. Avaliar Modelo de Séries Temporais

Objetivo: Garantir precisão e confiabilidade.

Ações:

- Métricas: `MAE` , `RMSE` , `MAPE` .
 - Cross-validation temporal (backtesting) para validar robustez.
-

13. Prever Tendências Futuras

Objetivo: Gerar projeções até o ano de 2040.

Ações:

- Usar modelo treinado para produzir previsões em horizonte de 15 anos.
 - Armazenar resultados em DataFrame `df_forecast` .
-

14. Visualizar Resultados

Objetivo: Facilitar interpretação dos stakeholders.

Ações:

- Gráficos de dispersão e linha para detecção de fraudes (análise histórica vs. predita).
 - Séries temporais com tendências e intervalos de confiança.
 - Dashboards interativos (ex.: Plotly, Tableau).
-

15. Resumir Descobertas

Objetivo: Destilar insights principais para tomada de decisão.

Ações:

- Síntese dos padrões de fraude (horários críticos, categorias de risco).
 - Principais variáveis preditoras e recomendações de monitoramento.
 - Tendências futuras de fraude e implicações de longo prazo.
-

16. Conclusão e Entrega

Objetivo: Apresentar resultados e plano de ação.

Ações:

- Preparar relatório executivo com principais métricas e gráficos.
- Documentar recomendações de negócio (alertas em tempo real, políticas de risco).
- Definir próximos passos: manutenção do modelo, monitoramento e ciclo de retraining.