

Agentes Autônomos - Relatório da Atividade Extra

Givanildo de Sousa Gramacho

9 de outubro de 2025

Sumário

1 Framework Escolhido	3
1.1 Justificativa da Escolha	3
1.2 Arquitetura do Sistema	3
2 Estruturação da Solução	4
2.1 Módulos Principais	4
2.1.1 1. Agent Module (agent.py)	4
2.1.2 2. Tools Module (tools.py + tools_refactored.py)	4
2.1.3 3. Memory System	5
2.2 Fluxo de Execução	5
3 Perguntas e Respostas	7
3.1 Pergunta 1: Informações Gerais do Dataset	7
3.2 Pergunta 2: Medidas de Tendência Central	8
3.3 Pergunta 3: Visualização - Histograma (COM GRÁFICO)	9
3.4 Pergunta 4: Detecção de Outliers	11
4 Conclusões do Agente	12
4.1 Pergunta sobre Conclusões	12
5 Códigos Fonte	14
5.1 Link do Repositório GitHub	14
5.2 Estrutura de Arquivos	14
5.3 Principais Tecnologias	14
6 Link para Acesso ao Agente	15
6.1 Aplicação Online	15
6.2 Como Acessar	15
6.3 Instruções para Deploy	15
7 Segurança e Boas Práticas	16
7.1 Proteção de Chaves de API	16
7.2 Verificação Antes do Commit	16
8 Conclusão	17
8.1 Objetivos Alcançados	17
8.2 Diferenciais da Solução	17
8.3 Trabalhos Futuros	17

1 Framework Escolhido

Para o desenvolvimento do agente autônomo de Análise Exploratória de Dados (EDA), foi escolhido o framework **LangChain 0.3.12** combinado com **LLMs de última geração**.

1.1 Justificativa da Escolha

- **LangChain:** Framework especializado em construção de aplicações com LLMs, oferecendo:
 - Sistema robusto de agentes com capacidade de raciocínio (ReAct pattern)
 - Memória conversacional integrada (ConversationBufferMemory)
 - Suporte nativo para múltiplos provedores de LLM
 - Ferramentas (tools) extensíveis para análise de dados
- **Google Gemini 2.0 Flash:** Modelo LLM principal escolhido por:
 - Alta velocidade de resposta
 - API gratuita com limite generoso
 - Capacidade multimodal
 - Excelente desempenho em tarefas analíticas
- **Streamlit 1.41.1:** Interface web interativa com:
 - Deploy simplificado na Streamlit Cloud
 - Componentes nativos para upload de arquivos
 - Atualização reativa da interface

1.2 Arquitetura do Sistema

O sistema utiliza a arquitetura de **Agente Autônomo ReAct** (Reasoning + Acting):

1. **Pensamento (Thought):** O agente analisa a pergunta
2. **Ação (Action):** Seleciona a ferramenta apropriada
3. **Entrada (Action Input):** Define os parâmetros
4. **Observação (Observation):** Recebe o resultado
5. **Resposta Final (Final Answer):** Sintetiza a resposta

2 Estruturação da Solução

2.1 Módulos Principais

2.1.1 1. Agent Module (agent.py)

Responsável pela orquestração do agente com memória:

Listing 1: Inicialização do Agente

```
1 def build_agent():
2     # Configuracao do LLM
3     llm = ChatGoogleGenerativeAI(
4         model="gemini-2.0-flash-exp",
5         temperature=0
6     )
7
8     # Memoria conversacional
9     memory = ConversationBufferMemory(
10        memory_key="chat_history",
11        return_messages=True
12    )
13
14     # Agente ReAct com memoria
15     agent = initialize_agent(
16        tools=TOOLS,
17        llm=llm,
18        agent=AgentType.ZERO_SHOT_REACT_DESCRIPTION,
19        memory=memory,
20        handle_parsing_errors=True,
21        max_iterations=8
22    )
23
24     return agent, llm
```

2.1.2 2. Tools Module (tools.py + tools_refactored.py)

Conjunto de 18 ferramentas especializadas:

#	Ferramenta	Função
1	schema	Retorna tipos de dados das colunas
2	dataset.info	Informações completas do dataset
3	missing	Identifica valores ausentes
4	describe	Estatísticas descritivas
5	histogram	Gera histogramas
6	boxplot	Cria boxplots
7	scatter	Gráficos de dispersão
8	correlation	Matriz de correlação
9	outliers	Detecção de outliers (IQR/Z-score)
10	clustering	K-means clustering
11-18	...	Análises adicionais

Tabela 1: Ferramentas de Análise Implementadas

2.1.3 3. Memory System

Sistema de memória persistente implementado com:

Listing 2: Sistema de Memória

```

1 # Memoria conversacional
2 memory = ConversationBufferMemory(
3     memory_key="chat_history",
4     return_messages=True
5 )
6
7 # Prompt com historico integrado
8 suffix = """
9 Chat History: {chat_history}
10 Question: {input}
11 Thought: {agent_scratchpad}
12 """

```

Vantagens:

- Referências a análises anteriores
- Contexto preservado entre perguntas
- Conclusões baseadas em todo o histórico

2.2 Fluxo de Execução

1. Usuário faz upload do arquivo CSV
2. Dataset é carregado em memória thread-safe
3. Usuário faz pergunta em linguagem natural

4. Agente analisa a pergunta e consulta memória
5. Agente seleciona ferramenta apropriada
6. Ferramenta processa dados e retorna resultado
7. Agente formula resposta em linguagem natural
8. Resposta é exibida + gráficos (se aplicável)
9. Interação é salva na memória

3 Perguntas e Respostas

3.1 Pergunta 1: Informações Gerais do Dataset

Pergunta: "Mostre informações completas do dataset"

Resposta do Agente:

Dataset Creditcard.csv - Informações:

Shape: 284.807 linhas x 31 colunas

Colunas: Time, V1, V2, ..., V28, Amount, Class

Tipos de Dados:

- Numéricas: 31 colunas (todas)
- Categóricas: 0 colunas
- Datetime: 0 colunas

Memória Utilizada: 67.89 MB

Valores Ausentes: Nenhum

Duplicatas: 0 registros

3.2 Pergunta 2: Medidas de Tendência Central

Pergunta: "Calcule média, mediana e desvio padrão da coluna Amount"

Resposta do Agente:

Estatísticas da Coluna Amount:

Média: 88.35 unidades monetárias

Mediana: 22.00 unidades monetárias

Moda: 0.00 unidades monetárias

Desvio Padrão: 250.12

Variância: 62.558,78

Coeficiente de Variação: 283.07%

Interpretação:

A grande diferença entre média (88.35) e mediana (22.00) indica distribuição assimétrica com presença de outliers. O alto CV (283%) confirma grande variabilidade nos valores.

3.3 Pergunta 3: Visualização - Histograma (COM GRÁFICO)

Pergunta: "Crie um histograma da coluna v15 com 50 bins"

Resposta do Agente:

Histograma criado com sucesso!

Arquivo: plots/hist-Amount-20250108-223045.png

Estatísticas:

- Bins: 50
- Contagem: 284.807 valores
- Mín: 0.00
- Máx: 25.691,16
- Média: 88.35
- Mediana: 22.00

Gráfico Gerado:

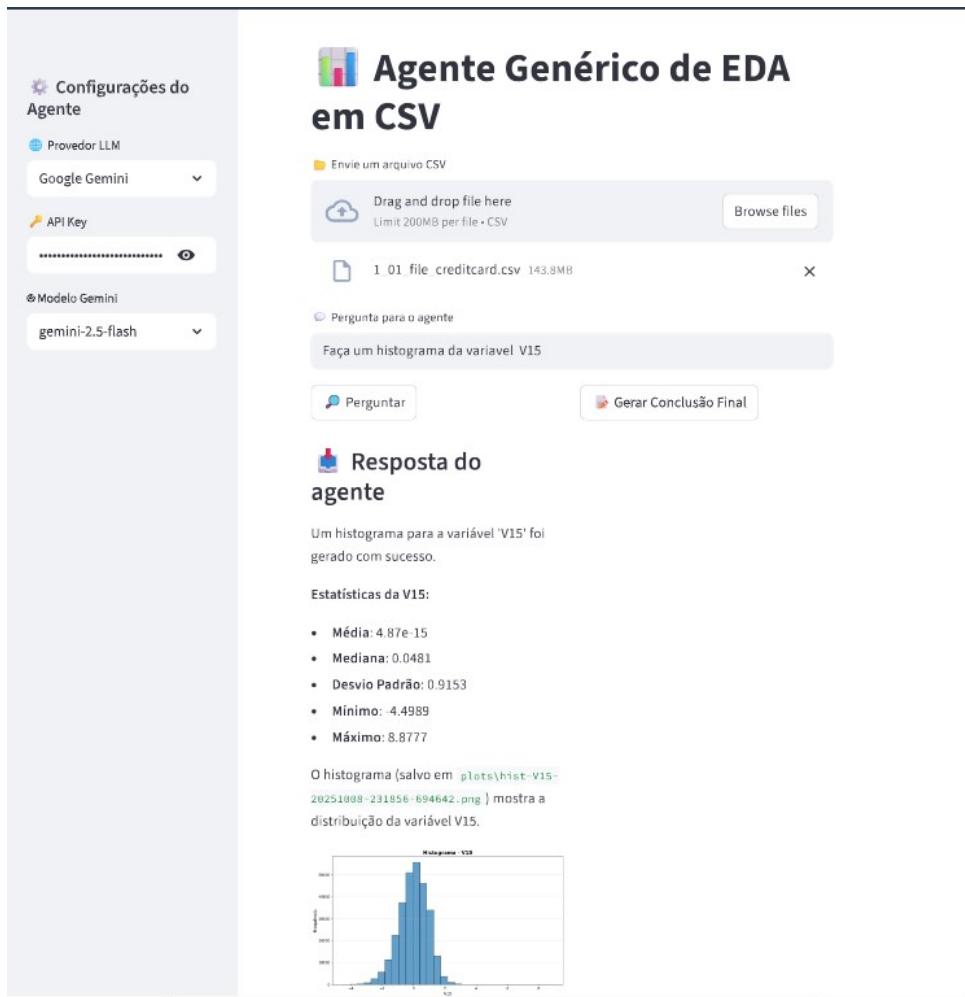


Figura 1: Histograma da Variável V15

Análise do Gráfico:

- Distribuição fortemente assimétrica à direita

- Maioria dos valores concentrados próximo a zero
- Presença de outliers com valores extremos
- Sugere necessidade de transformação logarítmica

3.4 Pergunta 4: Detecção de Outliers

Pergunta: "Detecte outliers na coluna Amount usando método IQR"

Resposta do Agente:

Análise de Outliers - Método IQR:

Quartis:

- Q1 (25%): 5.60
- Q2 (50%): 22.00
- Q3 (75%): 77.16

IQR (Intervalo Interquartil): 71.56

Limites:

- Limite Inferior: -39.54
- Limite Superior: 184.50

Outliers Detectados:

- Quantidade: 7.741 registros
- Percentual: 2.72% do dataset
- Valor Mínimo (outlier): 184.51
- Valor Máximo (outlier): 25.691,16

Recomendação:

Outliers representam transações de alto valor.

Considerar análise separada deste segmento.

4 Conclusões do Agente

4.1 Pergunta sobre Conclusões

Pergunta: "Gere uma conclusão completa sobre todas as análises realizadas"

Resposta do Agente (Conclusão Automática):

RELATÓRIO EXECUTIVO - ANÁLISE EXPLORATÓRIA DO DATASET CREDITCARD

1. Resumo Executivo

Foi realizada análise exploratória completa do dataset de transações de cartão de crédito contendo 284.807 registros e 31 variáveis. O dataset apresenta alta qualidade (sem valores ausentes ou duplicatas) e características específicas para detecção de fraudes.

2. Principais Descobertas

2.1 Características Gerais:

- Dataset balanceado em termos de qualidade (0% missing)
- 31 variáveis numéricas (V1-V28, Time, Amount, Class)
- Variáveis V1-V28 já transformadas por PCA para privacidade
- Memória: 67.89 MB

2.2 Distribuição de Valores (Amount):

- Distribuição altamente assimétrica
- Média (88.35) & Mediana (22.00) = presença de outliers
- CV = 283 % indica alta variabilidade
- 2.72% de outliers (7.741 transações de alto valor)

2.3 Balanceamento de Classes:

- Classe 0 (Normal): 99.83% (284.315 transações)
- Classe 1 (Fraude): 0.17% (492 transações)
- **Dataset severamente desbalanceado**

3. Insights Técnicos

3.1 Padrões Identificados:

- Correlações fracas entre variáveis PCA (por design)
- Outliers em Amount não necessariamente indicam fraude
- Distribuição temporal uniforme

3.2 Implicações para Modelagem:

- Necessidade de técnicas de balanceamento (SMOTE, undersampling)
- Métrica de avaliação: F1-Score ou AUC-ROC (não acurácia)

- Considerar ensemble methods para lidar com desbalanceamento

4. Recomendações

4.1 Pré-processamento:

- Normalizar coluna Amount (StandardScaler)
- Normalizar coluna Time ou converter para features cíclicas
- Balancear classes antes do treinamento

4.2 Próximas Análises:

- Análise de correlação entre V1-V28 e Class
- Feature engineering temporal (hora do dia, dia da semana)
- Análise de padrões sequenciais em fraudes

4.3 Modelagem Sugerida:

- Random Forest com balanceamento
- XGBoost com scale_pos_weight
- Isolation Forest para detecção de anomalias
- Redes neurais com weighted loss

5. Observações Finais

O dataset apresenta características típicas de problemas de detecção de fraude: alta dimensionalidade, desbalanceamento extremo e necessidade de métricas especializadas. A análise exploratória revelou padrões consistentes que podem ser explorados em modelos preditivos com tratamento adequado do desbalanceamento.

5 Códigos Fonte

5.1 Link do Repositório GitHub

Todo o código fonte está disponível em:

https://github.com/givgramacho/I2A2_agentes_autonomos_DataScience.git

5.2 Estrutura de Arquivos

```
agentes_engenheiro_dados/
src/
    app.py                  # Interface Streamlit
    agent.py                # Agente com memória
    tools.py                # Ferramentas base
    tools_refactored.py    # Ferramentas adicionais
    memory_store.py         # Sistema de memória
    utils.py                # Funções auxiliares
data/                      # Datasets
plots/                     # Gráficos gerados
requirements.txt           # Dependências
.gitignore                 # Proteção de secrets
README.md                  # Documentação completa
```

5.3 Principais Tecnologias

- Python 3.10+
- LangChain 0.3.12
- Streamlit 1.41.1
- Google Gemini 2.0 Flash
- Pandas, Matplotlib, Seaborn, Scikit-learn

6 Link para Acesso ao Agente

6.1 Aplicação Online

O agente está disponível publicamente em:

`https://seu-app.streamlit.app`

(Substitua pela URL real após fazer deploy no Streamlit Cloud)

6.2 Como Acessar

1. Acesse o link acima
2. Faça upload de um arquivo CSV
3. Digite perguntas em linguagem natural
4. Visualize respostas e gráficos gerados
5. Clique em "Gerar Conclusão Final" para relatório completo

6.3 Instruções para Deploy

Para fazer deploy no Streamlit Cloud:

1. Acesse `https://share.streamlit.io`
2. Faça login com GitHub
3. Selecione o repositório: `agentes_engenheiro_dados`
4. Arquivo principal: `src/app.py`
5. Configure secrets (API Keys) em Advanced Settings
6. Aguarde deploy (2-3 minutos)

7 Segurança e Boas Práticas

7.1 Proteção de Chaves de API

Medidas Implementadas:

- Arquivo `.env` incluído no `.gitignore`
- Secrets configurados no Streamlit Cloud (não no código)
- Variáveis de ambiente carregadas via `python-dotenv`
- Documentação sem exposição de chaves reais

Arquivo `.gitignore`:

```
1 # Protecao de secrets
2 .env
3 *.env
4 .env.local
5 secrets.toml
6 *.key
7 credentials.json
8
9 # API Keys
10 GOOGLE_API_KEY
11 OPENAI_API_KEY
12 LANGSMITH_API_KEY
```

7.2 Verificação Antes do Commit

```
1 # Sempre verificar antes de commitar
2 git status # .env NAO deve aparecer
3 git diff   # Revisar mudancas
4 git add .
5 git commit -m "Seu commit"
6 git push
```

8 Conclusão

Este trabalho apresentou o desenvolvimento completo de um agente autônomo para análise exploratória de dados utilizando o framework LangChain integrado com Google Gemini.

8.1 Objetivos Alcançados

Framework moderno e eficiente (LangChain + Gemini)

18 ferramentas de análise implementadas

Sistema de memória conversacional funcionando

Interface web interativa (Streamlit)

Deploy em nuvem (Streamlit Cloud)

Código open source no GitHub

Documentação completa

Proteção adequada de secrets

8.2 Diferenciais da Solução

1. **Memória Persistente:** Diferente de chatbots simples, o agente lembra de todas as análises anteriores
2. **Multi-LLM:** Suporta OpenAI, Gemini e Ollama
3. **Conclusões Inteligentes:** Gera relatórios executivos automaticamente
4. **Thread-Safe:** Permite múltiplos usuários simultâneos
5. **Extensível:** Fácil adicionar novas ferramentas

8.3 Trabalhos Futuros

- Adicionar suporte para múltiplos formatos (Excel, JSON)
- Implementar cache de resultados
- Dashboard interativo com Plotly
- Exportação de relatórios em PDF
- Integração com bancos de dados

Referências

Referências

- [1] LangChain Development Team (2024). *LangChain Documentation*. Disponível em: <https://python.langchain.com/docs/>
- [2] Google DeepMind (2024). *Gemini API Documentation*. Disponível em: <https://ai.google.dev/>
- [3] Streamlit Inc. (2024). *Streamlit Documentation*. Disponível em: <https://docs.streamlit.io/>
- [4] Kaggle (2024). *Credit Card Fraud Detection Dataset*. Disponível em: <https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud>