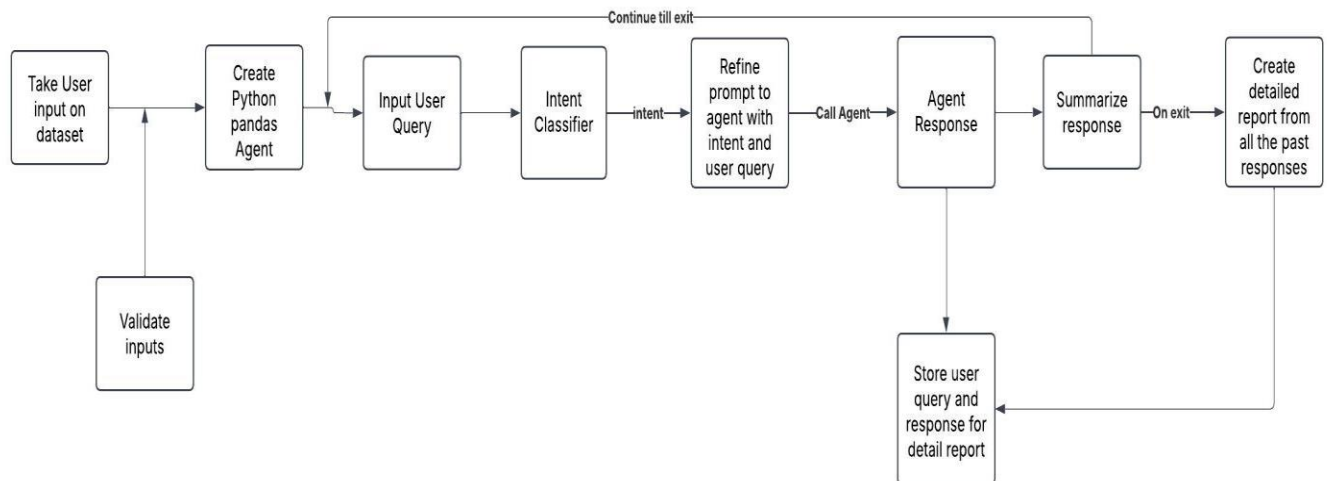


Semantic Spotter Project

Objective

Create a comprehensive data analysis agent on a given data set by converting natural language queries into detailed insights and reports.

Design



Major Components

- 1. Pandas Agent:** Used [Langchain Pandas Agent](#) for as agent for analyzing given dataset. This agent calls the Python agent under the hood, which executes LLM generated Python code.
- 2. Intent Classifier :** Analyses the input query from user and classifies the task in hand for the agent to execute. Rephrases the prompt to the agent based on the intent generated from LLM. Notify user if intent is unknown or not applicable as a data analysis task.

3. Report Builder : Generates a detailed report from a set of query and response pairs which the agent has generated

Model used

For evaluation of the project I have used Azure GPT 4O model. Following environment variables needs to be set to execute the model

AZURE_API_VERSION

AZURE_DEPLOYMENT_NAME

AZURE_OPENAI_API_KEY

AZURE_OPENAI_ENDPOINT

Code

The project is coded using Python 3 using langchain and AzureOpenAI modules. Following files(modules) are included in the package

dialog_flow_agent.py (Main dialog flow and entry point)

langchain_agent.py (Agent initialization and responses)

llm_helper.py (Additional llm helper methods)

Step to execute : python.exe dialog_flow_agent.py

Execution Flow

We have used loan dataset for testing the execution flow. Below is the execution flow snippet.

Welcome to your data analysis assistant!!

You can simulate a conversation with your data analysis assistant

You can ask questions about the dataset and get the analysis done

Enter the one word description of the dataset(eg loan, credit, insurance etc):Loan dataset

=====

Enter the filepath(csv) for the dataset:loan.csv

=====

Agent created for dataset: loan.csv

=====

You can ask questions about the dataset and get the analysis done

=====

Enter your query or type 'exit' to end the conversation:How many rows are present in the dataset?

Intent of the query: The intent of the query "How many rows are present in the dataset?" for the Loan dataset is "data summarization."

> Entering new AgentExecutor chain...

Invoking: `python_repl_ast` with `{ 'query': 'len(df)' }`

39717The dataset contains 39,717 rows.

> Finished chain.

Question: How many rows are present in the dataset?

=====

The dataset contains 39,717 rows.

=====

Brief summary of the response:

The dataset comprises 39,717 rows, indicating a substantial volume of data for analysis.

=====

Enter your query or type 'exit' to end the conversation:What the average annual income for all loans with loan status as fully paid?

Intent of the query: The intent of the query "What the average annual income for all loans with loan status as fully paid?" based on the dataset description "Loan dataset" is "data analysis".

> Entering new AgentExecutor chain...

Invoking: `python_repl_ast` with `{ 'query': "df[df['loan_status'] == 'Fully Paid']['annual_inc'].mean()}"`

69862.50332807285The average annual income for all loans with the loan status as "Fully Paid" is approximately \$69,862.50.

> Finished chain.

Question: What the average annual income for all loans with loan status as fully paid?

=====

The average annual income for all loans with the loan status as "Fully Paid" is approximately \$69,862.50.

=====

Brief summary of the response:

The average annual income for fully paid loans is \$69,862.50.

=====

Enter your query or type 'exit' to end the conversation:Whats the distribution of the loan amount?Also display the statistic of the field

Intent of the query: The intent of the query "Whats the distribution of the loan amount? Also display the statistic of the field." based on the dataset description "Loan dataset" is "data analysis".

> Entering new AgentExecutor chain...

Invoking: `python_repl_ast` with `{ 'query': "import matplotlib.pyplot as plt\nimport pandas as pd\n\n# Assuming df is already defined\nplt.hist(df['loan_amnt'], bins=20, edgecolor='black')\nplt.title('Distribution of Loan Amount')\nplt.xlabel('Loan Amount')\nplt.ylabel('Frequency')\nplt.show()"} `

Invoking: `python_repl_ast` with `{ 'query': "df['loan_amnt'].describe()"} `

count 39717.000000

mean 11219.443815

std 7456.670694

min 500.000000

25% 5500.000000

```
50%    10000.000000
75%    15000.000000
max     35000.000000
```

Name: loan_amnt, dtype: float64The distribution of the loan amount is visualized in the histogram above.

Here are the statistics for the loan amount field:

```
- **Count**: 39,717
- **Mean**: 11,219.44
- **Standard Deviation**: 7,456.67
- **Minimum**: 500
- **25th Percentile**: 5,500
- **Median (50th Percentile)**: 10,000
- **75th Percentile**: 15,000
- **Maximum**: 35,000
```

> Finished chain.

Question: Whats the distribution of the loan amount?Also display the statistic of the field

=====

The distribution of the loan amount is visualized in the histogram above.

Here are the statistics for the loan amount field:

```
- **Count**: 39,717
- **Mean**: 11,219.44
- **Standard Deviation**: 7,456.67
- **Minimum**: 500
- **25th Percentile**: 5,500
- **Median (50th Percentile)**: 10,000
```

- **75th Percentile**: 15,000

- **Maximum**: 35,000

=====

Brief summary of the response:

The loan amount distribution is shown in a histogram, with key statistics including a mean of 11,219.44, median of 10,000, and a range from 500 to 35,000, highlighting variability with a standard deviation of 7,456.67.

=====

Enter your query or type 'exit' to end the conversation: What percentage of rented owners has fully paid the loan?

Intent of the query: The intent of the query "What percentage of rented owners has fully paid the loan?" based on the dataset description "Loan dataset" is "data analysis".

> Entering new AgentExecutor chain...

```
Invoking: `python_repl_ast` with `{ 'query': "renters_fully_paid = df[(df['home_ownership'] == 'RENT') & (df['loan_status'] == 'Fully Paid')]\npercentage_fully_paid = (len(renters_fully_paid) / len(df[df['home_ownership'] == 'RENT'])) * 100\npercentage_fully_paid" }
```

82.76099264511349

```
Invoking: `python_repl_ast` with `{ 'query': "df['home_ownership'].unique()" }
```

['RENT' 'OWN' 'MORTGAGE' 'OTHER' 'NONE'] Approximately 82.76% of renters have fully paid their loans.

> Finished chain.

Question: What percentage of rented owners has fully paid the loan?

=====
Approximately 82.76% of renters have fully paid their loans.

=====
Brief summary of the response:

82.76% of renters have fully paid their loans.
=====

Enter your query or type 'exit' to end the conversation:exit

Thank you for using the Data Analysis Assistant!

Creating a detailed analysis report from your questions....

Detailed Data Analysis Report

1. **Dataset Overview**

- **Total Number of Rows**: The dataset comprises a total of 39,717 rows. This indicates a substantial amount of data, which can provide a comprehensive view of the loan-related information and trends.

2. **Average Annual Income for Fully Paid Loans**

- **Income Insight**: The average annual income for loans that have been fully paid is approximately \$69,862.50. This figure suggests that individuals with a higher income are more likely to successfully repay their loans, potentially indicating a correlation between income level and loan repayment success.

3. **Loan Amount Distribution and Statistics**

- **Distribution Visualization**: The loan amount distribution is represented through a histogram, which provides a visual understanding of how loan amounts are spread across the dataset.

- **Statistical Summary**:

- **Count**: 39,717 loans are recorded, matching the total number of rows in the dataset.

- **Mean Loan Amount**: \$11,219.44, indicating the average loan size.

- **Standard Deviation**: \$7,456.67, reflecting the variability in loan amounts.

- **Minimum Loan Amount**: \$500, showing the smallest loan issued.
- **25th Percentile**: \$5,500, suggesting that 25% of loans are below this amount.
- **Median (50th Percentile)**: \$10,000, indicating the middle value of the loan amounts.
- **75th Percentile**: \$15,000, showing that 75% of loans are below this amount.
- **Maximum Loan Amount**: \$35,000, representing the largest loan issued.
- **Insight**: The wide range and high standard deviation suggest diverse borrowing needs and financial capabilities among borrowers.

4. **Loan Repayment Among Renters**

- **Repayment Rate**: Approximately 82.76% of renters have fully paid their loans. This high percentage indicates a strong repayment performance among renters, which could be attributed to effective financial management or favorable loan terms for this demographic.

Conclusion

This report provides a comprehensive overview of the dataset, highlighting key insights into loan amounts, borrower income, and repayment behaviors. The data suggests that higher income levels are associated with successful loan repayment, and renters demonstrate a commendable repayment rate. The variability in loan amounts underscores the diverse financial needs and capabilities of borrowers. These insights can be valuable for financial institutions in tailoring loan products and assessing risk.