

Data Analysis

Exploratory Data Analysis (EDA) - LLM



Problem

▼ Problem description

Context

The original dataset contains 1000 records with 20 categorical attributes prepared by Prof. Hofmann.

In this dataset, each record represents a person taking out a loan from a bank.

Each person is classified as **good** or **bad** credit risk based on the set of attributes.

Link to the original dataset at

[UCI Machine Learning](#)

Content

The selected attributes:

- **Age** (numeric)
- **Sex** (text: male, female)
- **Job** (numeric: 0 - unskilled and non-resident, 1 - unskilled and resident, 2 - skilled, 3 - highly skilled)
- **Housing** (text: own, rent, or free)
- **Saving accounts** (text - little, moderate, quite rich, rich)
- **Checking account** (numeric, in DM - Deutsch Mark)
- **Credit amount** (numeric, in DM)
- **Duration** (numeric, in month)
- **Purpose**(text: car, furniture/equipment, radio/TV, domestic appliances, repairs, education, business, vacation/others)
- **Risk** (Value target - Good or Bad Risk)

Objective

Train a model to predict from new data whether a person applying for a loan represents a good or bad risk

PROMPT

Develop a Python program that performs the following exploratory data analysis (EDA) tasks on a CSV file provided by user `german_credit_data.csv`:

Reading CSV file: Load the CSV file into a DataFrame using the pandas library.

Initial display: Displays the first 5 rows of the dataset.

Dataset overview: Provides an informative summary of the dataset, including the number of rows and columns, the data types of each column, and the number of non-null values.

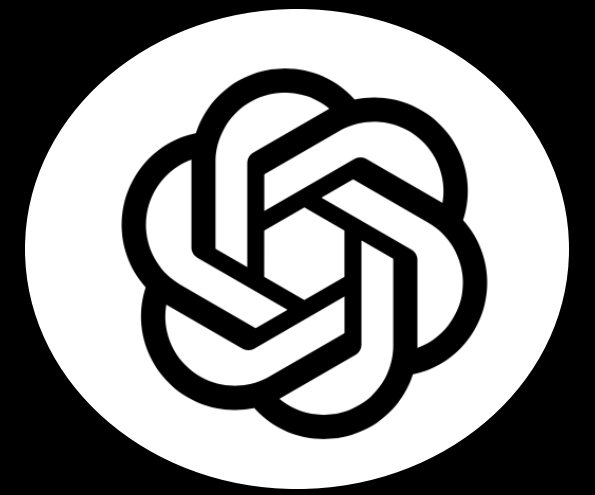
Statistical description: Generates a statistical summary of the dataset, including metrics such as mean, median, standard deviation, minimum, maximum, and quartiles for numeric columns.

Correlation matrix: Calculates the correlation matrix for numeric variables in the dataset.

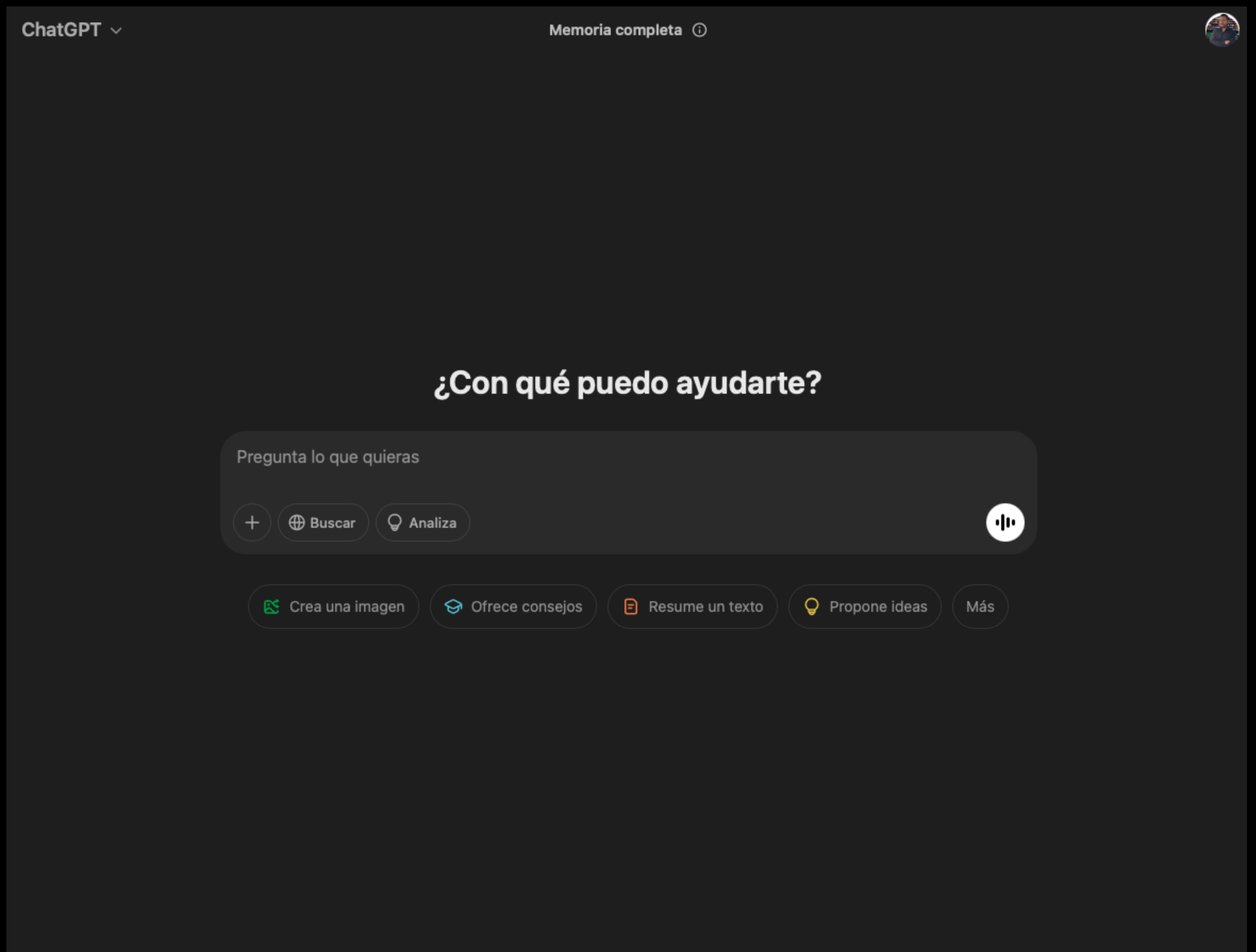
Heatmap: Visualize the correlation matrix using a heatmap with the seaborn or matplotlib library, making sure that the correlation values are readable and the plot is well labeled.

The program should be modular, easy to understand, and well-commented. It should also include error handling for cases such as a missing CSV file, non-numeric columns in the correlation matrix, or formatting issues in the file.

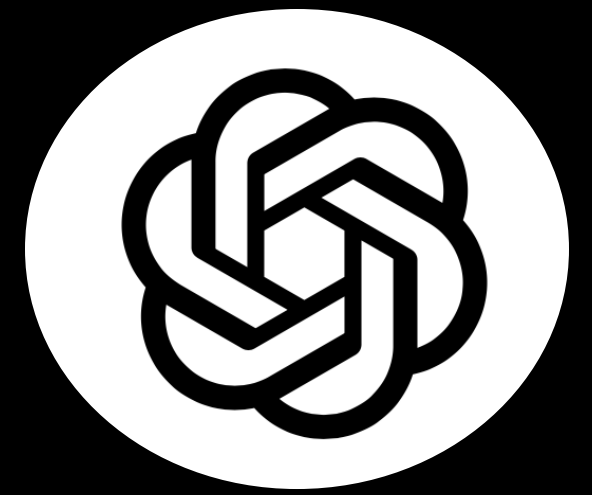
ChatGPT



Go to ChatGPT website



ChatGPT



Introduce a prompt to create EDA in python

ChatGPT

Memoria completa

¿Con qué puedo ayudarte?

german_credit_data.csv

Hoja de cálculo

Develop a Python program that performs the following exploratory data analysis (EDA) tasks on a CSV file provided by user german_credit_data.csv:
Reading CSV file: Load the CSV file into a DataFrame using the pandas library.
Initial display: Displays the first 5 rows of the dataset.
Dataset overview: Provides an informative summary of the dataset, including the number of rows and columns, the data types of each column, and the number of non-null values.
Statistical description: Generates a statistical summary of the dataset, including metrics such as mean, median, standard deviation, minimum, maximum, and quartiles for numeric columns.
Correlation matrix: Calculates the correlation matrix for numeric variables in the dataset.
Heatmap: Visualize the correlation matrix using a heatmap with the seaborn or matplotlib library

+

Buscar

Analiza

Crea una imagen

Ofrece consejos

Resume un texto

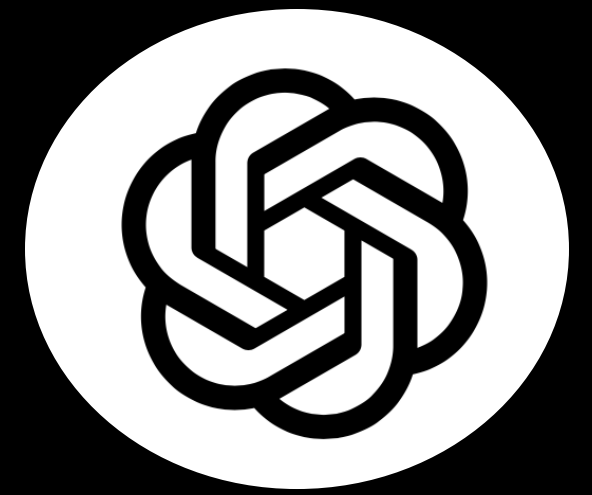
Propone ideas

Más

ChatGPT puede cometer errores. Comprueba la información importante.

?

ChatGPT



See result of LLM

ChatGPT

Memoria completa

Compartir

columns.

Correlation matrix: Calculates the correlation matrix for numeric variables in the dataset.

Heatmap: Visualize the correlation matrix using a heatmap with the seaborn or matplotlib library, making sure that the correlation values are readable and the plot is well labeled.

The program should be modular, easy to understand, and well-commented. It should also include error handling for cases such as a missing CSV file, non-numeric columns in the correlation matrix, or formatting issues in the file.

I'll develop a Python program that performs the requested Exploratory Data Analysis (EDA) tasks on the uploaded `german_credit_data.csv` file. The program will include modular functions, error handling, and clear comments. Here's the plan:

1. **Read the CSV file** and load it into a Pandas DataFrame.
2. **Display the first 5 rows** of the dataset.
3. **Provide a dataset overview** including shape, data types, and non-null values.
4. **Generate a statistical summary** for numeric columns.
5. **Calculate the correlation matrix** for numeric columns.
6. **Visualize the correlation matrix** using a heatmap.

I'll now write and execute the Python script.

It looks like there was an issue accessing the file. I'll retry loading the CSV file to see if we can proceed with the analysis. [↻]

Pregunta lo que quieras

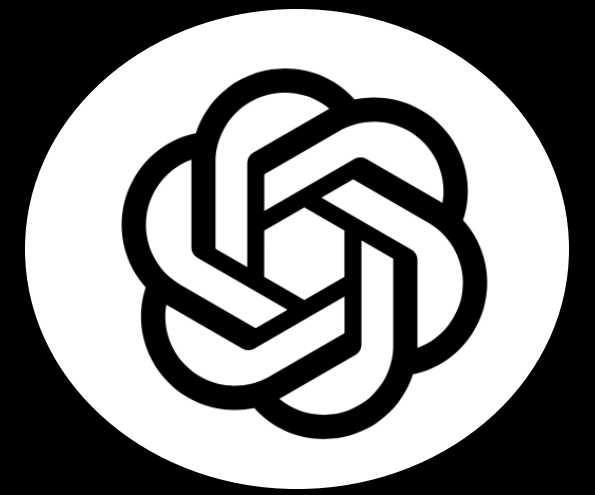
+

Buscar

Analiza

ChatGPT puede cometer errores. Comprueba la información importante.

ChatGPT



Copy Python code

ChatGPT

Memoria completa

Compartir

Análisis

python

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Define file path
file_path = "/mnt/data/german_credit_data.csv"

def load_data(file_path):
    """Load the CSV file into a Pandas DataFrame with error handling"""
    try:
        df = pd.read_csv(file_path)
        print("File successfully loaded!")
        return df
    except FileNotFoundError:
        print("Error: File not found.")
        return None
    except pd.errors.EmptyDataError:
        print("Error: File is empty.")
        return None
    except pd.errors.ParserError:
        print("Error: File could not be parsed.")
        return None

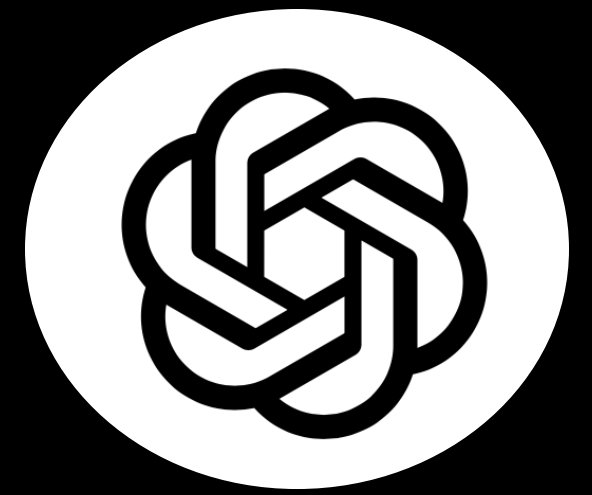
def display_head(df, n=5):
```

Mostrar siempre los detalles Copiar

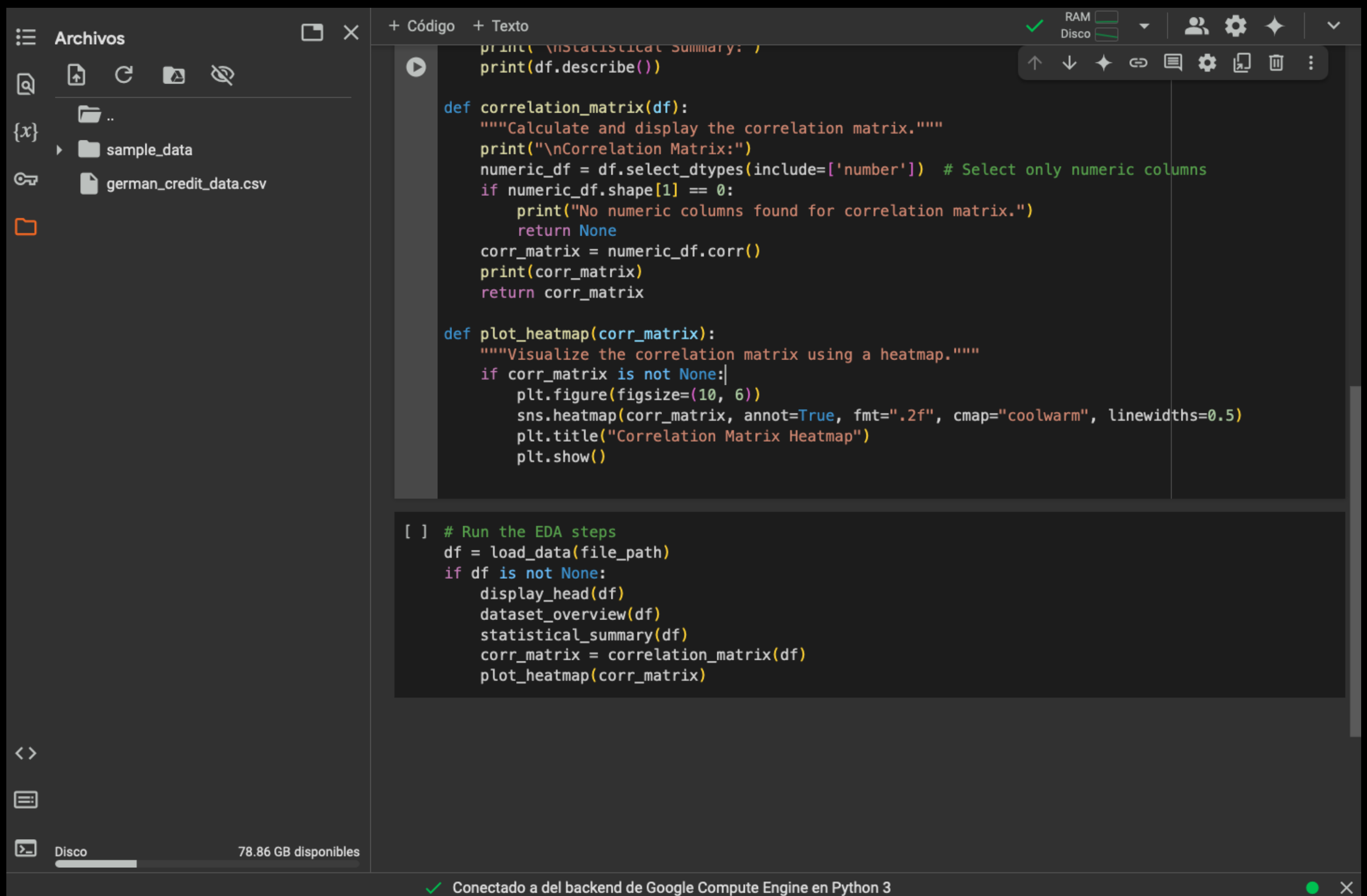
+ Buscar Analiza

ChatGPT puede cometer errores. Comprueba la información importante.

ChatGPT+ Colab



Go to Colab, paste code & check results

A screenshot of the Google Colab web interface. The left sidebar shows a file explorer with a folder named 'sample_data' and a file named 'german_credit_data.csv'. The main area contains two code cells. The top cell defines two functions: 'correlation_matrix(df)' which calculates the correlation matrix for numeric columns, and 'plot_heatmap(corr_matrix)' which visualizes the correlation matrix as a heatmap. The bottom cell is a code block that runs the EDA steps: loading the data, displaying the head, dataset overview, statistical summary, and then calling the two functions defined above. The interface includes a top bar with RAM and Disco usage, and a bottom status bar indicating connection to Google Compute Engine.

```
+ Código + Texto
print("\nStatistical Summary: ")
print(df.describe())

def correlation_matrix(df):
    """Calculate and display the correlation matrix."""
    print("\nCorrelation Matrix:")
    numeric_df = df.select_dtypes(include=['number']) # Select only numeric columns
    if numeric_df.shape[1] == 0:
        print("No numeric columns found for correlation matrix.")
        return None
    corr_matrix = numeric_df.corr()
    print(corr_matrix)
    return corr_matrix

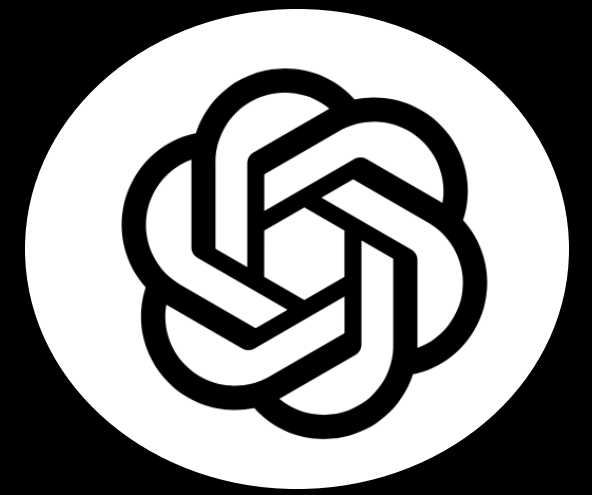
def plot_heatmap(corr_matrix):
    """Visualize the correlation matrix using a heatmap."""
    if corr_matrix is not None:
        plt.figure(figsize=(10, 6))
        sns.heatmap(corr_matrix, annot=True, fmt=".2f", cmap="coolwarm", linewidths=0.5)
        plt.title("Correlation Matrix Heatmap")
        plt.show()

[ ] # Run the EDA steps
df = load_data(file_path)
if df is not None:
    display_head(df)
    dataset_overview(df)
    statistical_summary(df)
    corr_matrix = correlation_matrix(df)
    plot_heatmap(corr_matrix)
```

Disco 78.86 GB disponibles

✓ Conectado a del backend de Google Compute Engine en Python 3

ChatGPT+ Colab



Go to Colab, paste code & check results

EDA_ChatGPT.ipynb ☆

Archivo Editar Ver Insertar Entorno de ejecución Herramientas Ayuda

Archivos

sample_data

german_credit_data.csv

+ Código + Texto

```
# Run the EDA steps
df = load_data(file_path)
if df is not None:
    display_head(df)
    dataset_overview(df)
    statistical_summary(df)
    corr_matrix = correlation_matrix(df)
    plot_heatmap(corr_matrix)
```

File successfully loaded!

First 5 rows of the dataset:

Unnamed: 0	age	sex	job	housing	saving_accounts	checking account
0	0	67	male	2	own	NaN
1	1	22	female	2	own	little
2	2	49	male	1	own	little
3	3	45	male	2	free	little
4	4	53	male	2	free	little

credit amount duration purpose risk

0	1169	6	radio/TV	good
1	5951	48	radio/TV	bad
2	2096	12	education	good
3	7882	42	furniture/equipment	good
4	4870	24	car	bad

Dataset Overview:

Number of Rows: 1000, Number of Columns: 11

Column Data Types and Non-null Counts:

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 1000 entries, 0 to 999

Data columns (total 11 columns):

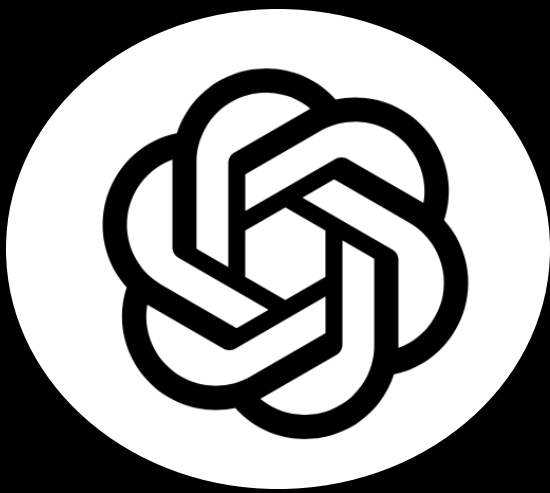
#	Column	Non-Null Count	Dtype
0	Unnamed: 0	1000 non-null	int64
1	age	1000 non-null	int64
2	sex	1000 non-null	object

Disco

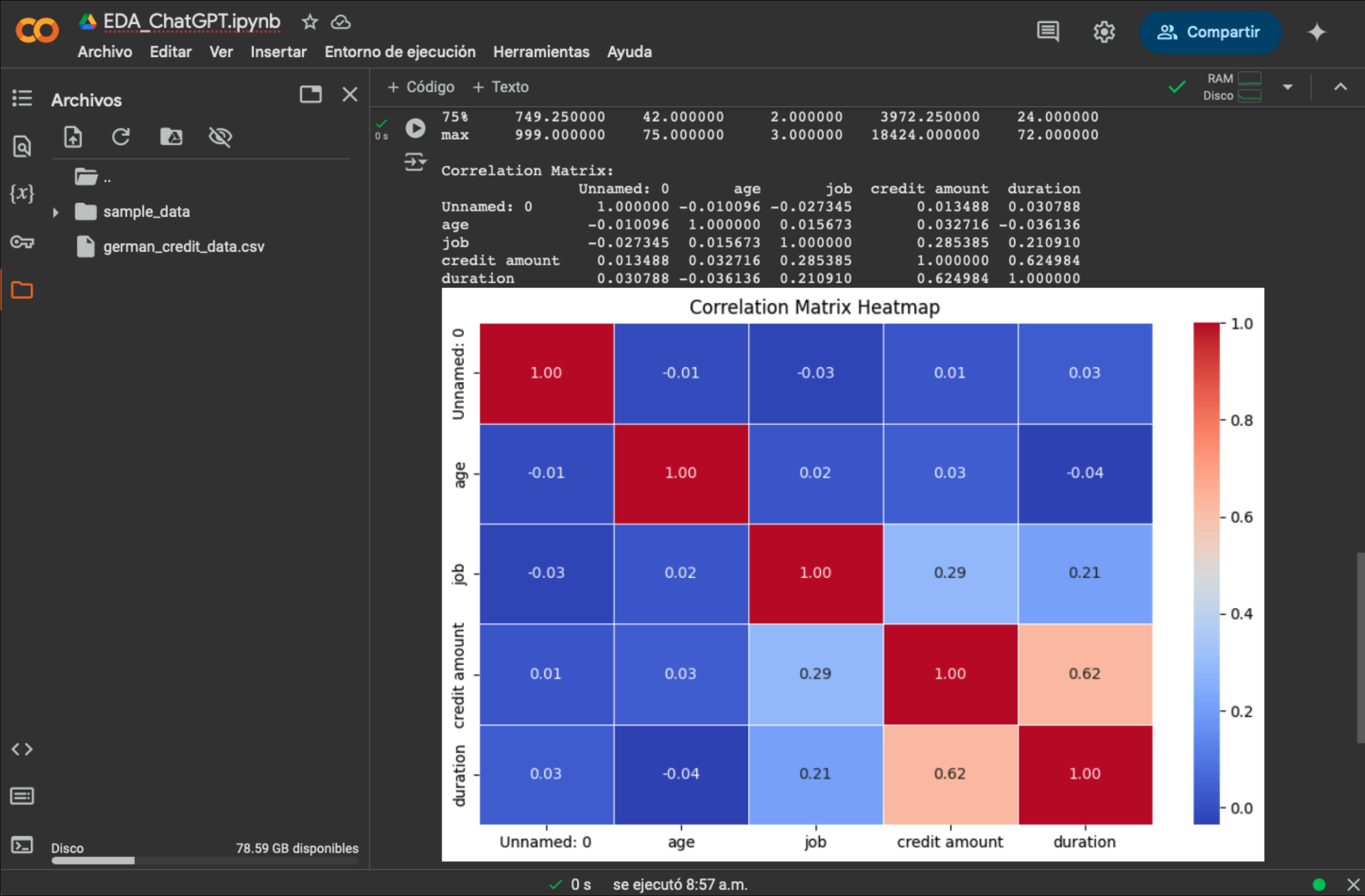
78.59 GB disponibles

0 s se ejecutó 8:57 a.m.

ChatGPT+ Colab



Go to Colab, paste code & check results



Qwen



Go to Qwen website

Qwen2.5-Max

Set as default

E

Good morning, Edgar Rios Linares

How can I help you today?

Thinking (QwQ)Web SearchToolsCodeMake a planNewsMore

Create image

Code

Make a plan

News

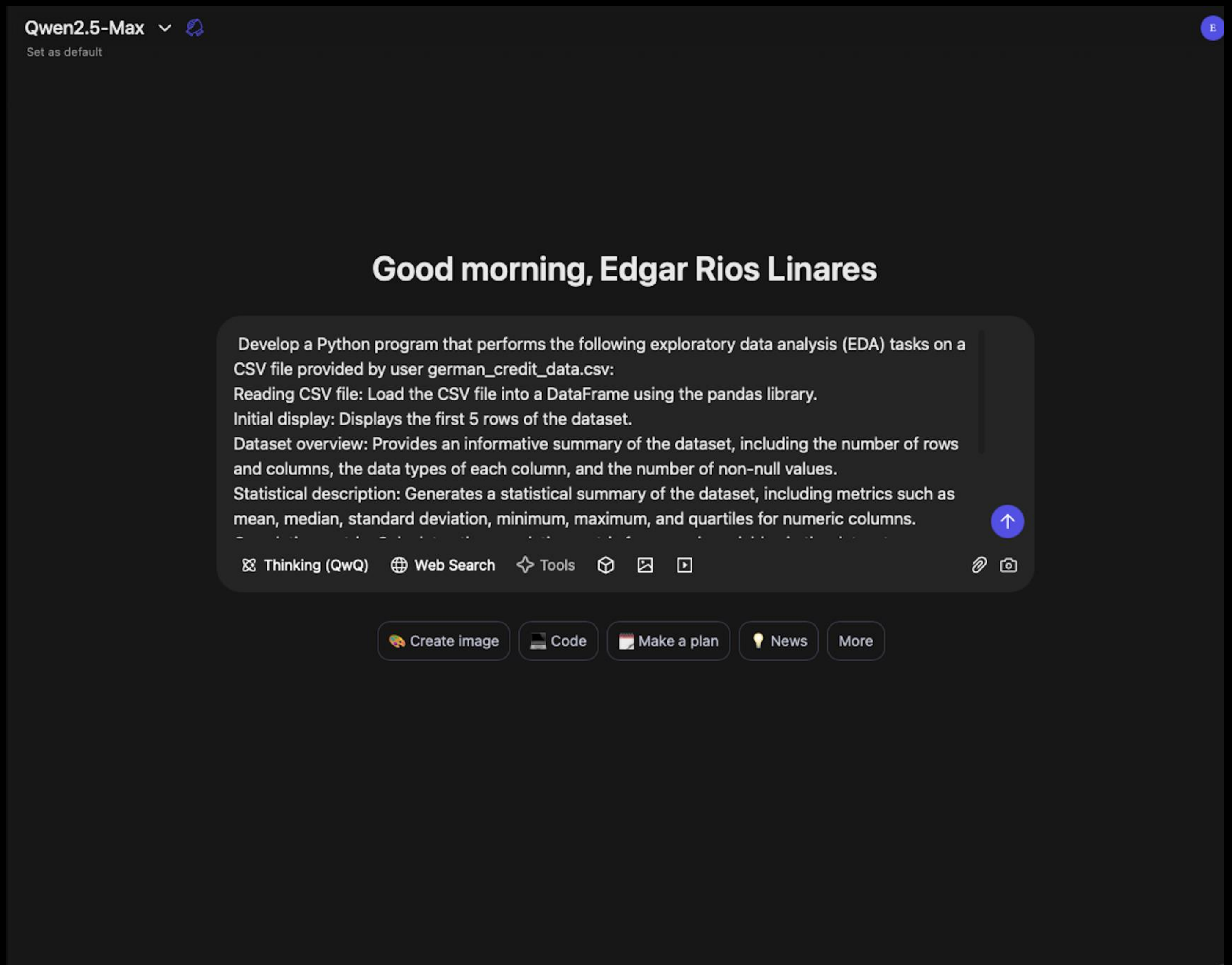
More

?

Qwen



Introduce a prompt to create EDA in python



Qwen



See result of LLM

Qwen2.5-Max

Qwen2.5-Max 8:59 AM

Below is a Python program that performs the requested exploratory data analysis (EDA) tasks on the `german_credit_data.csv` file. The program is modular, well-commented, and includes error handling for common issues.

python

Copy

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 def load_csv(file_path):
6     """
7     Load a CSV file into a pandas DataFrame.
8
9     Parameters:
10         file_path (str): Path to the CSV file.
11
12     Returns:
13         pd.DataFrame: Loaded DataFrame.
14     """
15     try:
16         df = pd.read_csv(file_path)
17         print("CSV file successfully loaded.")
18         return df
19     except FileNotFoundError:
20         print(f"Error: The file '{file_path}' was not found. Please check the file path.")
21         return None
22     except pd.errors.EmptyDataError:
23         print("Error: The file is empty or contains no data.")
24         return None
25     except pd.errors.ParserError:
26         print("Error: The file could not be parsed. Check the file format.")
27         return None
28
29 def display_initial_stats(df):
```

Send a Message

Thinking (QwQ) Web Search Tools

Thinking (QwQ) is a preview model and is still being updated.

Qwen



See result of LLM

Qwen2.5-Max

Explanation of the Code:

- Modular Design :**
 - Each task is encapsulated in its own function (`load_csv` , `display_initial_rows` , `dataset_overview` , etc.), making the code reusable and easy to understand.
- Error Handling :**
 - The `load_csv` function handles errors such as missing files, empty files, and parsing issues.
 - The `calculate_correlation_matrix` function checks if there are numeric columns before attempting to calculate the correlation matrix.
- Dataset Overview :**
 - The `dataset_overview` function provides a detailed summary of the dataset, including the number of rows, columns, data types, and descriptive statistics.
- Correlation Matrix and Heatmap :**
 - The correlation matrix is calculated only for numeric columns using `select_dtypes` .
 - The heatmap is visualized using `seaborn.heatmap` , with annotations for better readability.
- User Interaction :**
 - The program prompts the user to input the file path, making it flexible for different datasets.
- Visualization :**
 - The heatmap is styled with a color map (`coolwarm`) and annotated with correlation values rounded to two decimal places.

How to Run the Program:

- Save the code in a Python file, e.g., `eda_script.py` .
- Place the `german_credit_data.csv` file in the same directory or provide the full path to the file.

Send a Message

Thinking (QwQ)

Web Search

Tools

Thinking (QwQ) is a preview model and is still being updated.

Qwen



Copy Python code

Qwen2.5-Max

Qwen2.5-Max 8:59 AM

Below is a Python program that performs the requested exploratory data analysis (EDA) tasks on the `german_credit_data.csv` file. The program is modular, well-commented, and includes error handling for common issues.

python

Copy

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 def load_csv(file_path):
6     """
7     Load a CSV file into a pandas DataFrame.
8
9     Parameters:
10         file_path (str): Path to the CSV file.
11
12     Returns:
13         pd.DataFrame: Loaded DataFrame.
14     """
15     try:
16         df = pd.read_csv(file_path)
17         print("CSV file successfully loaded.")
18         return df
19     except FileNotFoundError:
20         print(f"Error: The file '{file_path}' was not found. Please check the file path.")
21         return None
22     except pd.errors.EmptyDataError:
23         print("Error: The file is empty or contains no data.")
24         return None
25     except pd.errors.ParserError:
26         print("Error: The file could not be parsed. Check the file format.")
27         return None
28
29 def display_initial_stats(df):
```

Send a Message

Thinking (QwQ) Web Search Tools

Thinking (QwQ) is a preview model and is still being updated.

Qwen + Colab



Go to Colab, paste code & check results

CO

EDA_QWen.ipynb

☆

Archivo Editar Ver Insertar Entorno de ejecución Herramientas Ayuda

Archivos

..

sample_data

german_credit_data.csv

+ Código + Texto

def main():
 """
 Main function to perform EDA on the provided CSV file.
 """
 # Prompt user for the file path
 file_path = "german_credit_data.csv"

 # Step 1: Load the CSV file
 df = load_csv(file_path)
 if df is None:
 return

 # Step 2: Display the first 5 rows
 display_initial_rows(df)

 # Step 3: Provide an overview of the dataset
 dataset_overview(df)

 # Step 4: Calculate the correlation matrix
 corr_matrix = calculate_correlation_matrix(df)

 # Step 5: Plot the heatmap
 plot_heatmap(corr_matrix)

if __name__ == "__main__":
 main()

CSV file successfully loaded.

First 5 rows of the dataset:
 Unnamed: 0 age sex job housing saving_accounts checking account \
0 0 67 male 2 own NaN little
1 1 22 female 2 own little moderate
2 2 49 male 1 own little NaN
3 3 45 male 2 free little little
4 4 53 male 2 free little little

credit amount duration purpose risk

Disco

78.59 GB disponibles

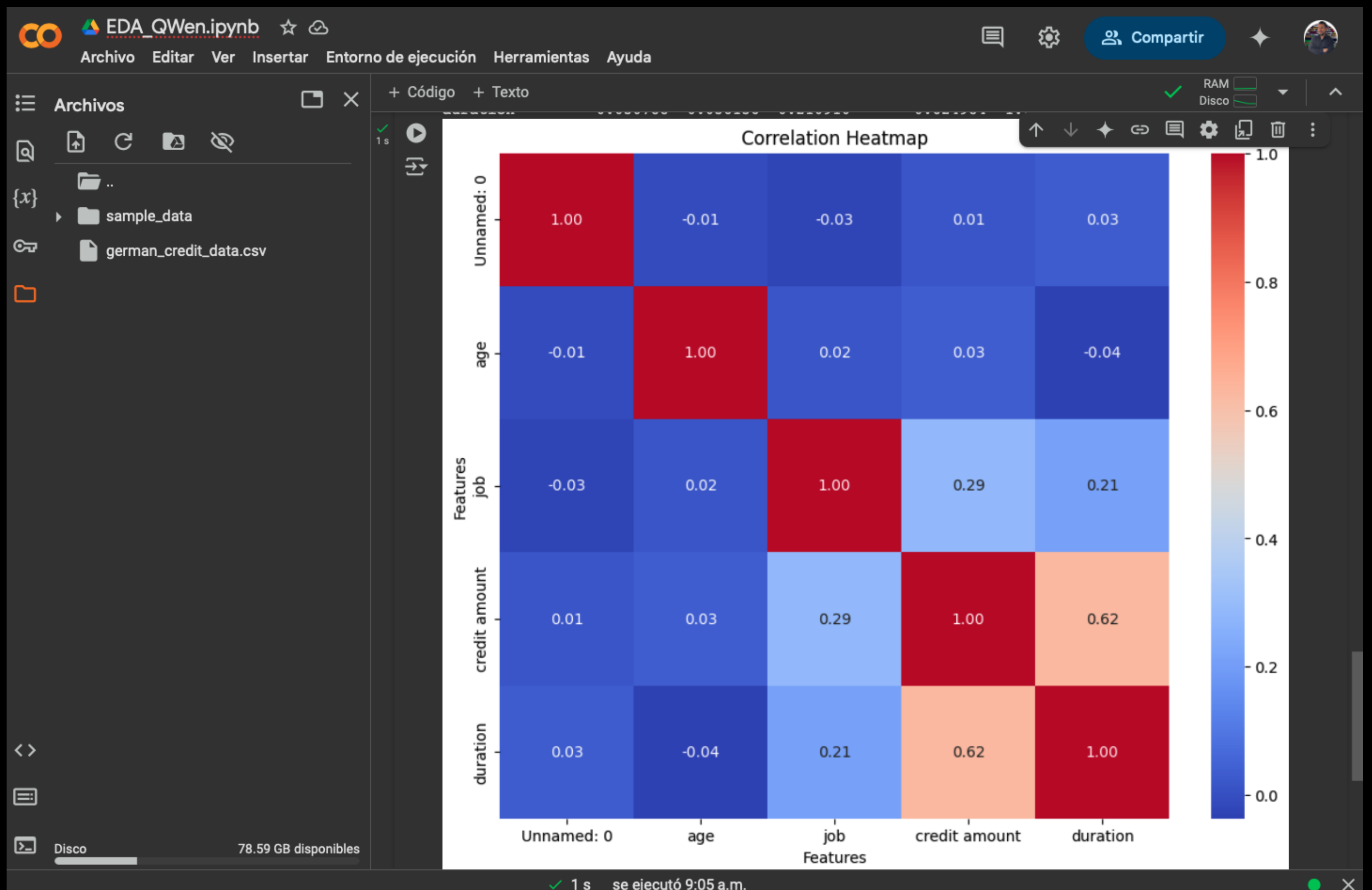
1 s

se ejecutó 9:05 a.m.

Qwen + Colab



Go to Colab, paste code & check results



The work is done, great job!"

You have a basic EDA quickly



Educator in AI

**Artificial
Intelligence**

Data Engineering



Machine Learning

Data Science

📌 **Linkedin** —> <https://www.linkedin.com/in/erlinares/>

👋 **Follow us on X**: <https://x.com/erlinares>^[SEP]

💻 **GitHub**: https://github.com/erlinares/365_AI_Journey/

💬 **Discord**: <https://discord.gg/5fFM2zh8>



Edgar Rios Linares