# GOVERNMENT COLLEGE OF TECHNOLOGY

PROJECT NAME: CSV QUERY BOT SYSTEM

PROJECT MEMBERS:

- SIVARANJANI S(7177211740)
- PRIYANKA S (71772117133)

DEPT: BE CSE

COURSE: ARTIFICIAL INTELLIGENCE

COURSE CODE:18SPC702

# CSV Query Bot System

## Project Overview:

**Talk with CSV** is an innovative tool designed to simplify data interaction and analysis by enabling users to query CSV data using natural language inputs. The solution aims to bridge the gap between non-technical users and complex data analysis, transforming how data is interpreted, analyzed, and presented.

## Objectives:

The primary objectives of this project include:
- Enabling users to query CSV data using natural language.
- Providing meaningful responses, including textual data, tables, and visualizations.
- Ensuring a user-friendly interface that makes data analysis intuitive and accessible.

## Technical Stack:

The project leverages a robust technical stack to deliver its functionalities:
- Streamlit: For building an interactive and user-friendly web interface.
- Python (Pandas): For efficient data manipulation and analysis.
- Google Generative AI (Gemini API): For processing and responding to natural language queries.
- dotenv & json Libraries: Used for environment configuration and data serialization.

## Project Architecture and Workflow:

- **CSV Upload**: Users begin by uploading their CSV file through the Streamlit web interface.
- **Natural Language Query**:Users can input a natural language query related to their data.
- **Query Processing**: The system constructs a prompt and sends it to the Gemini API for natural language processing.
- **Response Handling**: The API returns a response, which is then processed to generate the desired output (e.g., text, tables, visualizations).
- **Result Display**: The output is displayed on the Streamlit interface, making data interaction easy and intuitive for users.
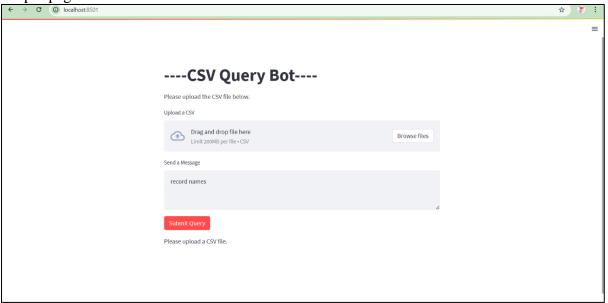
## Key Features:

- **Natural Language Querying**: Users can interact with data using everyday language without the need for complex code or SQL queries.
- **User-Friendly Interface**: Built with Streamlit, the interface is designed to make data interaction seamless and accessible.
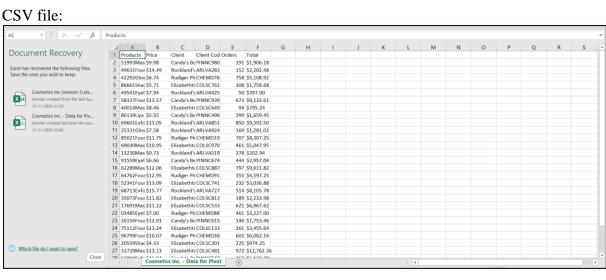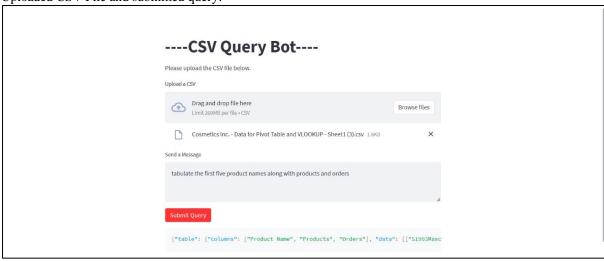
# Code Implementation:



```python
import os
import json
import pandas as pd
import streamlit as st
from dotenv import load_dotenv
import google.generativeai as genai

# Load environment variables
load_dotenv()

# Configure Gemini API key
genai.configure(api_key=os.getenv("GEMINI_API_KEY"))

def call_gemini_api(prompt):
    """
    Calls the Gemini API with a specified prompt.
    """
    try:
        model = genai.GenerativeModel('gemini-1.5-flash')
        response = model.generate_content(prompt)

        # Print the full response object for inspection
        print("Full Response:", response)

        # Check if 'text' is in response; return error if missing
        if hasattr(response, 'text'):
            return response.text
        else:
            return '{"answer": "Unexpected response structure from the Gemini API."}'
    except Exception as e:
        print(f"Request failed: {e}")
        return '{"answer": "An error occurred while contacting the Gemini API."}'

def csv_tool(file) -> pd.DataFrame:
    """
    Reads a CSV file and returns it as a DataFrame.
    """
```



```python
def csv_tool(file) -> pd.DataFrame:
    df = pd.read_csv(file)
    return df

def ask_agent(df, query):
    """
    Creates a prompt, sends it to the Gemini API, and decodes the response.
    """
    prompt = (
        """
        Let's decode the way to respond to the queries. The responses depend on the type of information requested in the query.

        1. If the query requires a table, format your answer like this:
            {"table": {"columns": ["column1", "column2", ...], "data": [[value1, value2, ...], [value1, value2, ...], ...]}}

        2. For a bar chart, respond like this:
            {"bar": {"columns": ["A", "B", "C", ...], "data": [25, 24, 10, ...]}}

        3. If a line chart is more appropriate, your reply should look like this:
            {"line": {"columns": ["A", "B", "C", ...], "data": [25, 24, 10, ...]}}

        4. For a plain question that doesn't need a chart or table, your response should be:
            {"answer": "Your answer goes here"}

        5. If the answer is not known or available, respond with:
            {"answer": "I do not know."}
        Example output: {"columns": ["Products", "Orders"], "data": [["51993Masc", 191], ["49631Foun", 152]]}

        Here's the query to work on:
        """ + query
    )

    response_text = call_gemini_api(prompt)
    return str(response_text)
```



```python
def write_answer(response_dict: dict):
            st.bar_chart(df)
        except Exception as e:
            st.write(f"Error creating bar chart: {e}")
            print(f"Error creating bar chart: {e}")

    # Check if the response is a line chart
    if "line" in response_dict:
        try:
            data = response_dict["line"]
            df_data = {
                col: [x[i] for x in data['data']] for i, col in enumerate(data['columns'])
            }
            df = pd.DataFrame(df_data)
            df.set_index(df.columns[0], inplace=True)  # Automatically set the first column as index
            st.write("Generated Line Chart:")
            st.line_chart(df)
        except Exception as e:
            st.write(f"Error creating line chart: {e}")
            print(f"Error creating line chart: {e}")

# Streamlit App Layout
st.set_page_config(page_title="📊 Talk with your CSV")
st.title("----CSV Query Bot----")
st.write("Please upload the CSV file below.")

data = st.file_uploader("Upload a CSV", type="csv")
query = st.text_area("Send a Message")

if st.button("Submit Query", type="primary"):
    if data:
        df = csv_tool(data)
        response = ask_agent(df, query)
        # Decode the response.
        decoded_response = decode_response(response)

        # Write the response to the Streamlit app.
        write_answer(decoded_response)
```

# Output Snapshots:

Output page:



CSV file:



Uploaded CSV File and submitted query:

Desired output:

```
*Untitled - Notepad                          —    □    ×
File  Edit  Format  View  Help
{"table":
{"columns": ["Product Name", "Products", "Orders"],
 "data": [["51993Masc", 191, "Masc"],
          ["49631Foun", 152, "Foun"],
          ["14072Cushi", 117, "Cushi"],
          ["14072Cushi", 117, "Cushi"],
          ["47364Wom", 105, "Wom"]]
}
}
```

## Challenges Faced:

During the development of **Talk with CSV**, we encountered several challenges:
- Handling Complex Queries: Ensuring accurate responses for diverse and complex queries posed a significant challenge. This required fine-tuning of prompts sent to the Gemini API and error handling mechanisms.
- Data Processing Efficiency: Processing large CSV files efficiently while maintaining a smooth user experience demanded optimization efforts.
- User Interface Design: Balancing a simple, user-friendly interface with the provision of comprehensive functionality required iterative design improvements.

## Solutions and Optimizations:

To overcome these challenges, we implemented the following solutions:
- Query Optimization: We improved the accuracy of responses by enhancing the prompt construction and integrating validation mechanisms.
- Efficient Data Handling:Using Python's Pandas library enabled us to process large datasets efficiently, reducing response times.
- User-Centric Design: Continuous feedback and usability testing allowed us to enhance the interface for improved user experience.

## Real-World Applications:

The versatility of **Talk with CSV** makes it applicable in various scenarios, including:
- **Education**: Students can analyze and interpret data without extensive technical knowledge.
- **Business Analysis**: Professionals can gain insights from sales data, customer data, and more, using simple queries.
- **Data Exploration**: Researchers and data analysts can quickly explore and visualize datasets.

## Future Enhancements:

We plan to further enhance **Talk with CSV** by:

- Adding Support for Additional Data Formats: Enabling interaction with other data formats beyond CSV files.
- Improving Natural Language Understanding: Refining query processing capabilities to handle more complex queries accurately.
- Expanding Visualization Options: Introducing more chart types and customizable data visualizations.
- Integrating Additional AI Capabilities: Leveraging other AI models to improve response generation and user engagement.