# Project 8: Excel Q&A Assistant

# Key Concepts

# 1. Core Technologies and Libraries

## Pandas

- DataFrame Manipulation: At the heart of this project is Pandas, the primary tool for data manipulation in Python. You will use it to read, clean, transform, and analyze the data from Excel files.

- Multi-Sheet Handling: Learn to work with Excel files that contain multiple sheets, selecting, and combining data as needed.
- Data Aggregation: Use groupby() and aggregation functions (sum, mean, count, etc.) to summarize data based on different categories.

## LangChain

- Agents and Tools: This project will introduce you to the concept of LangChain agents. You will create a "tool" that the agent can use, which in this case will be the ability to execute Python (Pandis) code.
- LLM Integration: LangChain provides a standardized interface to interact with Large Language Models (LLMs) like GPT-3.5 or GPT-4.
- Prompt Engineering: You will design prompts that instruct the LLM on how to convert a natural language question into a Pandas command.

## Plotly

- Dynamic Visualizations: Plotly is a powerful library for creating interactive charts and graphs. Unlike static images, Plotly charts can be zoomed, panned, and hovered over to reveal more data.
- Chart Types: You will learn to generate various chart types like bar charts, line charts, scatter plots, and pie charts based on the data and the user's query.

# 2. Key Concepts in AI and Data Analysis

## Natural Language to Code

- Semantic Parsing: This is the process of converting natural language into a machine-readable format. In this project, you're converting a user's question into executable Python code.

- **Few-Shot Learning:** You might provide the LLM with a few examples of questions and their corresponding Pandas code to "teach" it how to perform the conversion. This is a form of in-context learning.

## Data Profiling

- Automated Data Analysis: Before you can answer questions about a dataset, you need to understand its structure. Data profiling is the process of automatically summarizing a dataset, including:
  - Data Types: Identifying which columns are numbers, dates, text, etc.
  - Descriptive Statistics: Calculating mean, median, standard deviation, etc., for numerical columns.
  - Value Distributions: Counting the occurrences of different values in categorical columns.

## Sandboxing

- Secure Code Execution: Since the LLM will be generating Python code that gets executed on your server, it's a major security risk.
- A "sandbox" is a restricted environment where you can run this code without it being able to access the filesystem, network, or other sensitive resources. This is crucial to prevent malicious code from being executed.

# 3. Architectural Concepts

## Milestone-Based Development

- Iterative Approach: This project is broken down into milestones, allowing you to build and test the application in stages. This is a common practice in software development to manage complexity and ensure quality.
- Quality Gates: Each milestone has a review process. This ensures that each part of the application is well-designed, secure, and performs well before you build on top of it.

# API-First Design

- Decoupled Frontend and Backend: While you can build this as a single Streamlit application, a more scalable approach is to have a separate backend (e.g., using FastAPI) that handles the data processing and AI logic.
- The frontend (e.g., a React app) would then communicate with the backend through an API. This separation of concerns makes the application easier to maintain and scale.