# SQL & Data Modeling Sprint

# Project

**Introduction**

In the era of big data and AI, structured data and relational databases remain the backbone of data science pipelines. Even as unstructured data (like text or images) gains prominence, most organizations still rely on structured, relational data for core analytics and decision-making. Mastering data modeling (designing how data is structured and stored) and SQL (Structured Query Language for querying databases) are therefore critical for data scientists. These skills ensure that raw data can be transformed into reliable, actionable insights that drive business value. Below, we explore why structured data, data modeling, and SQL are foundational to real-world data science, addressing key questions and illustrating with examples.

This research report will specifically address five guiding questions that explore key aspects of these skills:

1. Why structured data is crucial in data science pipelines?
2. What role does data modeling play in preparing data for analysis or machine learning?
3. How do relational databases support scalable and clean data practices in real-world data science projects?
4. Why is SQL still considered a foundational skill even with tools like Python and Pandas?
5. a practical example illustrating how SQL is used to extract insights before applying machine learning.

## Structured Data in Data Science Pipelines

What is Structured Data?

**Structured data** is clearly organized information (e.g., tables, rows, columns) stored in formats like SQL databases or Excel. It follows a fixed schema, making it easy to search, sort, and analyze.

Structured Data is Important for You in Data Science because of:

1. Ease of Analysis:

   o Structured data's organized nature makes it straightforward to perform statistical analysis and derive meaningful insights using tools like SQL, Python, or R.

2. Machine Learning and Modeling:

   o Structured data is essential for training machine learning models since algorithms typically require clear numerical and categorical inputs.

## Role of Data Modeling

Data modeling prepares data for analysis and ML by organizing it clearly. A good schema simplifies feature extraction. For example, churn prediction is easier with structured tables of customer activity.

## Relational Databases in Real-world Projects

Relational databases ensure clean, consistent, and scalable data using structured schemas and ACID rules. Netflix uses Redshift and Trino to manage large data volumes efficiently.

## SQL as a Foundational Skill

It handles large datasets better than Pandas and teaches relational thinking. SQL is widely used in hiring tests. Airbnb uses it daily for data insights.

## Practical Example of SQL Usage

Before machine learning models are deployed, SQL is extensively used to perform exploratory data analysis (EDA), extracting key insights from large datasets. For instance, SQL queries easily summarize customer engagement metrics such as purchase history and average usage. This preliminary analysis helps define effective ML features and guides model development. Airbnb frequently employs SQL-based EDA to optimize their recommendation algorithms and customer experience strategies.

# References

- ChatGpt

- Hevo Data. *What is Structured Data and Why is It Important? Explains structured vs unstructured data and emphasizes the popularity and benefits of structured data*. What is Structured Data and Why is It Important for You?

- KDnuggets. *Introduction to Databases in Data Science*. Data Modeling in Machine Learning Pipelines: Best Practices Using SQL and NoSQL Databases - DATAVERSITY

- Medium. *Why the Creator of Pandas Says Data Scientists Must Learn SQL*. Why the Creator of Pandas Says Data Scientists Must Learn SQL | by Dr. Ilyes SEDKA | Medium

- Milvus (Zilliz). (2025). *How is SQL Used in Data Analytics?* How is SQL used in data analytics?