# Project 7: Data Warehousing with IBM Cloud Db2 Warehouse

**Phase 5: Project Documentation & Submission.**

## PROJECT OBJECTIVE

Our project main ob to tackle the challenge of modernizing our data management by designing and implementing a robust data warehouse using IBM Cloud Db2 Warehouse. This project is driven by the need to harness the potential of the organization's data, coming from diverse sources, and to empower data architects with the tools needed for insightful data analysis and informed decision-making.

**The project encompasses the following key objectives:**

i. Data Warehouse Structure: We will define a flexible and scalable schema and structure for our data warehouse that can efficiently handle data from various sources.
ii. Data Integration: The project will involve identifying key data sources within our organization and devising a strategy for integrating this data seamlessly into the data warehouse.
iii. ETL (Extract, Transform, Load) Processes: To ensure data accuracy and relevance, we will develop robust ETL processes for extracting data from source systems, transforming it into a usable format, and loading it efficiently into the data warehouse.
iv. Data Exploration: We aim to create an intuitive and interactive environment that allows data architects to explore and analyze the data effectively.
v. Actionable Insights: Ultimately, our goal is to deliver actionable insights from the data analysis, enabling our organization to make data-driven decisions.

## DESIGN THINKING

➤ **Data Warehouse Structure**

To define the data warehouse structure, we will take the following practical steps:

- Data Profiling: We will conduct thorough data profiling to gain insights into the data's characteristics, including data types, relationships, and quality.
- Entity-Relationship Diagram (ERD): We will create an ERD to visualize the data model, making it easier to design schemas and tables.

➤ **Data Integration**

For effective data integration, we will follow a realistic approach:

- Source Identification: We will identify all potential data sources within the organization, including databases, third-party APIs, and legacy systems.
- Data Extraction Strategy: We will determine data extraction methods and frequency, aligning them with our organization's data needs.
- Transformation Rules: We will establish clear data transformation rules and procedures to maintain data consistency and quality.

➤ **ETL Processes**

Pragmatic ETL processes will be designed and implemented as follows:

- ETL Tool Selection: We will select ETL tools or custom scripting languages based on our organization's resources and expertise.

➢ **Data Exploration**

To promote practical data exploration, we will adopt these strategies:

- User-Centric Design: Our user interface using flask will prioritize user-friendliness, making it easier for data architects to interact with the data warehouse.
- Query Building: We will provide user-friendly query-building capabilities to empower data architects to create custom queries.
- Data Visualization: Implementing data visualization components, such as interactive charts and graphs, will enhance the data exploration experience.

➢ **Actionable Insights**

Our approach to delivering actionable insights will be grounded in reality:

➢ Analysis Templates: We will create predefined analysis templates for common use cases to expedite insights generation.

# TOOLS AND TECHNOLOGIES PLANNED TO USE

## Data Modelling

Creating a data model that represents the structure of our data warehouse by involving tables, relationships, and attributes.

We are going to use the IBM Data Architect tool which can provide assistance in data modeling.

## Setting Up IBM Cloud Db2 Warehouse

Configuring our database instance, including setting up security, scalability, and performance options.

## Data Loading

Loading the transformed data into your **IBM Db2** Warehouse instance. We will use IBM DataStage for this purpose.

Ensuring data quality and consistency during the loading process

## Data Access and Analysis

Analysis using **IBM Watson Studio** tool. By building dashboards, reports, and data analysis pipelines to empower your data architects and analysts.
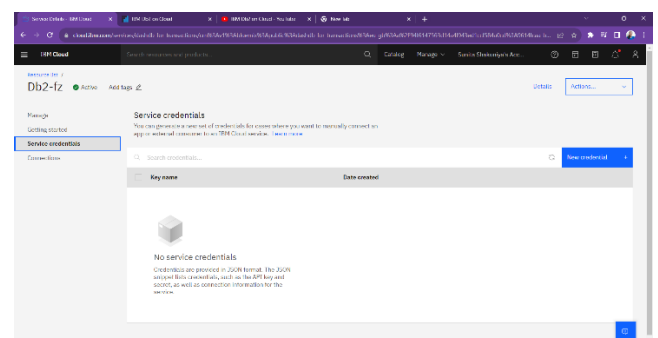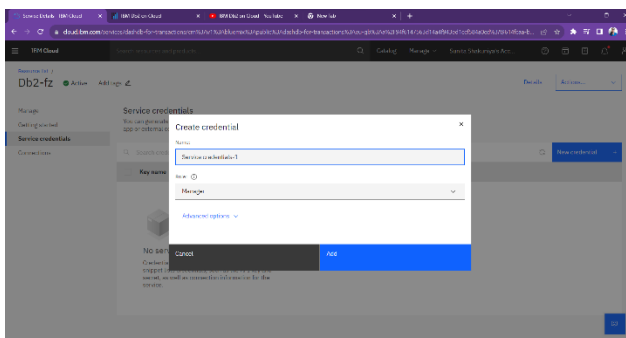
Utilize the capabilities of Db2 Warehouse, such as its in-database analytics and machine learning features, to perform advanced analyses.
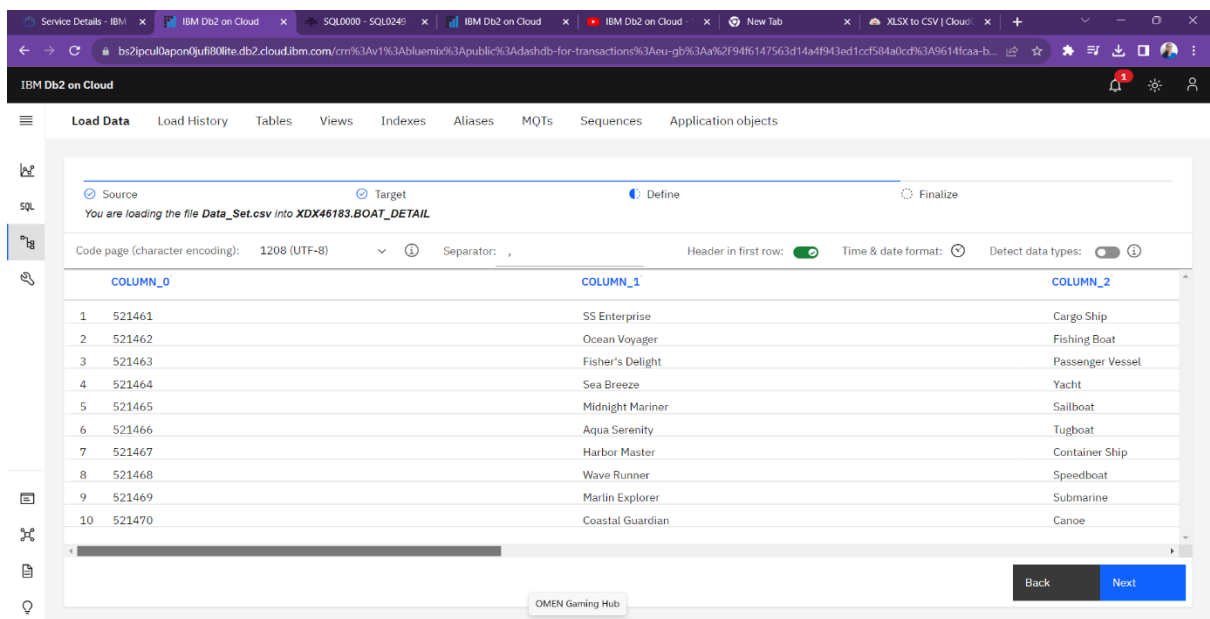
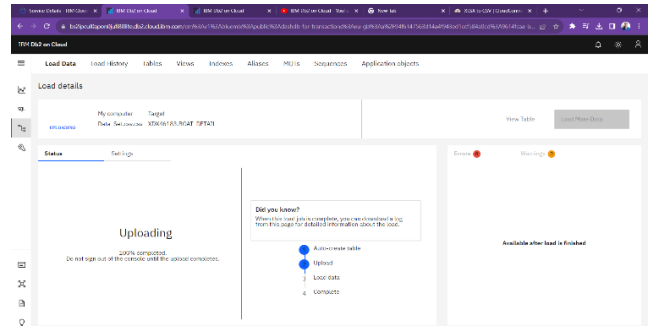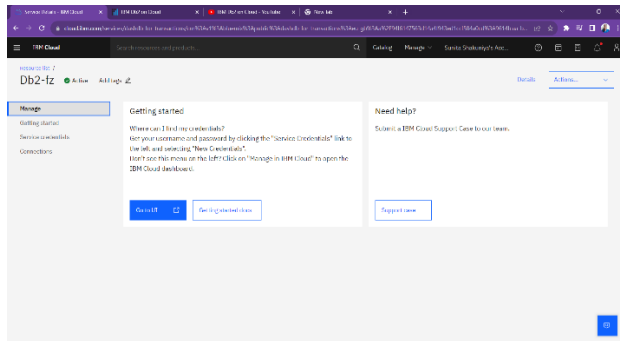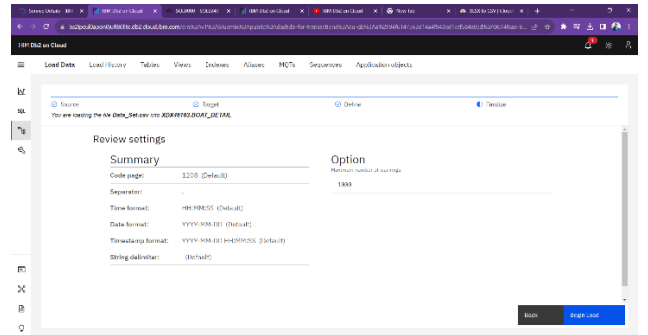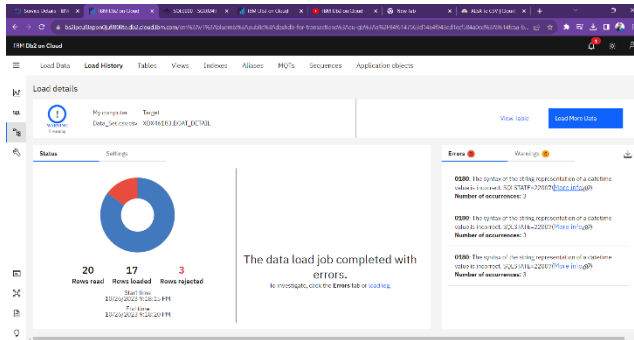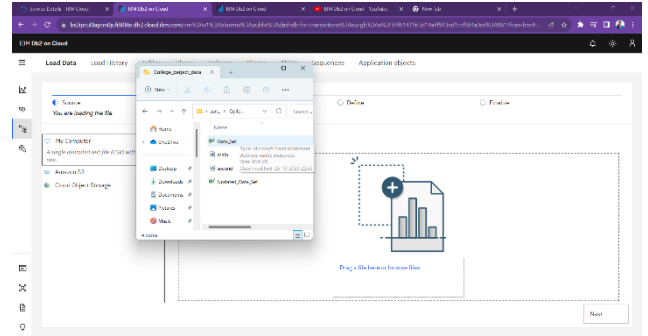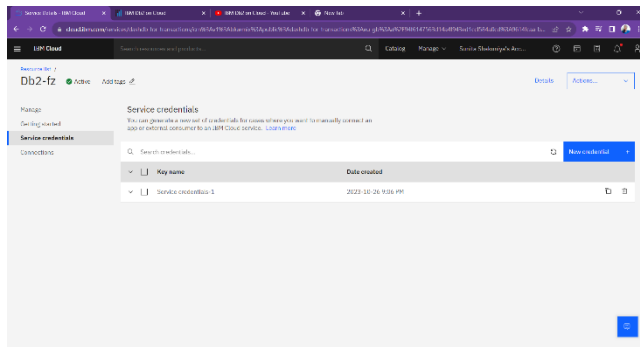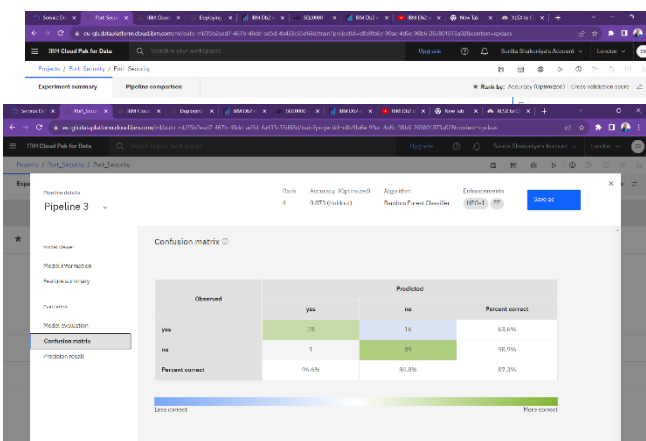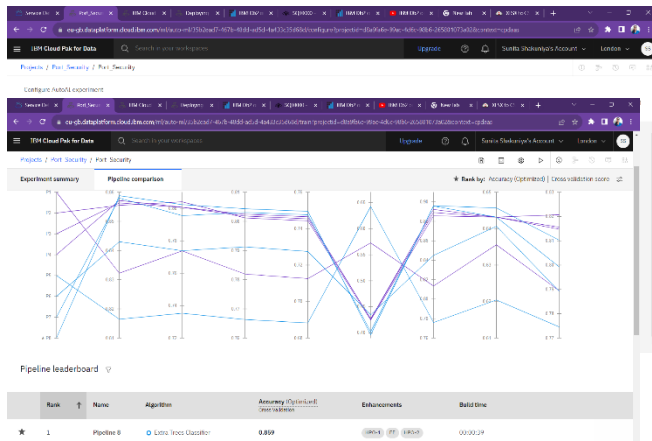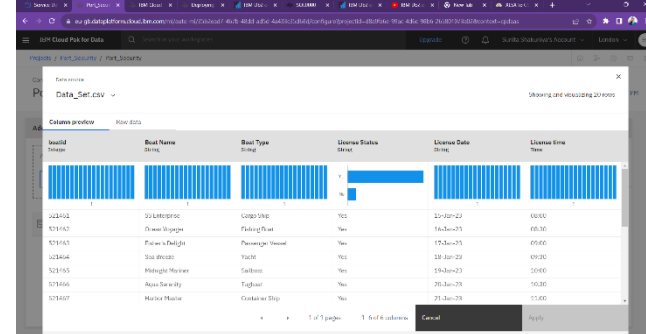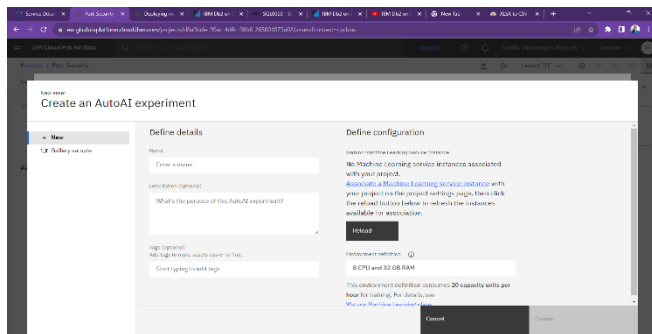# TECHNICAL IMPLEMENTATION

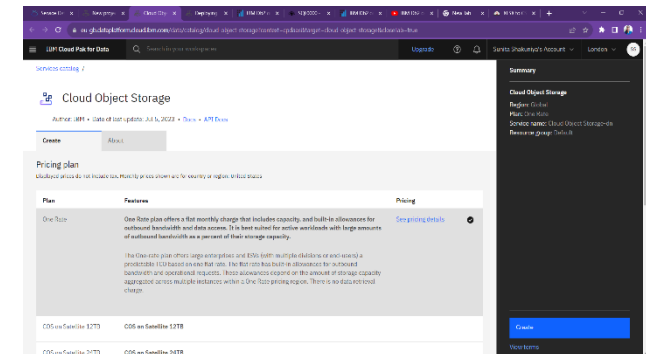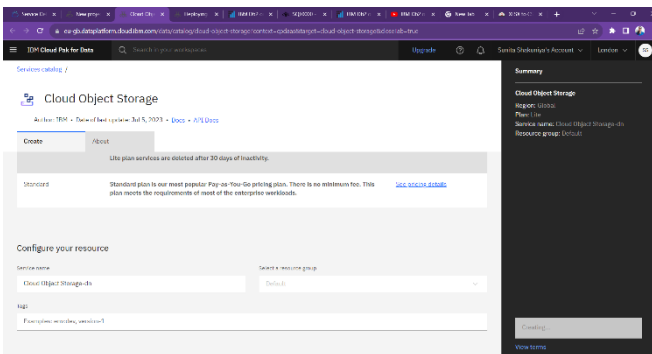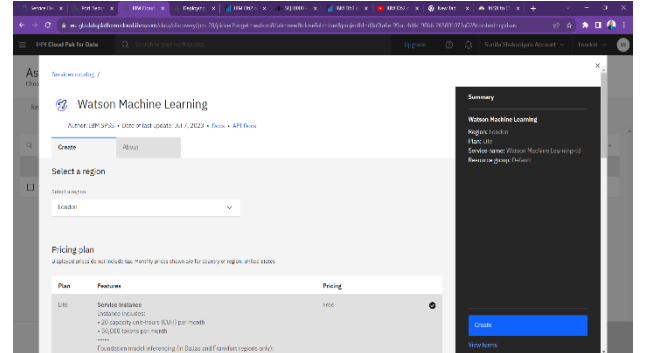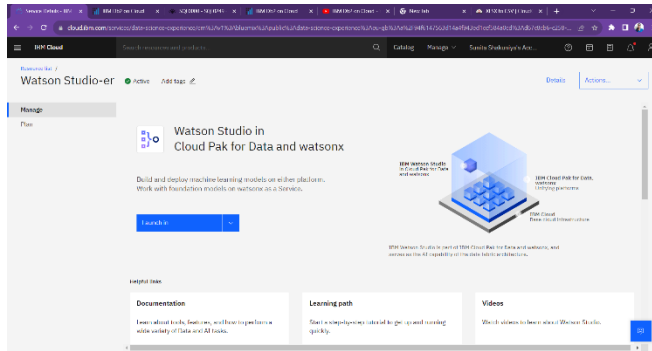➢ **Step 1:** Creating IBM cloud account





➢ **Step 2:** Setting up the db2 instance and experimenting.

> **Step 3:** Setting up the Watson Studio instance and experimenting.

➢ **Step 4:** As we can perform ETL process over Watson studio so we further proceed with panadas   library and MySql data base.

- Creating Excel data set



- Coding implementation over Vscode

```python
27
28      # Convert letters in 'Boat Name' and 'Boat Type' to uppercase
29      df['Boat Name'] = df['Boat Name'].str.upper()
30      df['Boat Type'] = df['Boat Type'].str.upper()
31
32      # The name of the table where you want to insert the data
33      table_name = 'boat_security'
34
35      # Write the cleaned and transformed data from the DataFrame to MySQL
36      df.to_sql(table_name, engine, if_exists='replace', index=False)
37
38
39      # Execute an SQL query to retrieve the data from the MySQL table
40      sql_query = f'SELECT * FROM {table_name}'
41      df = pd.read_sql_query(sql_query, engine)
42
43      # Specify the path for the output Excel file
44      output_excel_file = 'Updated_Data_Set.xlsx'
45
46      # Save the data from the SQL query to an Excel file
47      df.to_excel(output_excel_file, index=False)
```

```
PS C:\Users\sunit\College_parject_data> & C:/Users/sunit/AppData/Local/Programs/Python/Python310/python.exe c:/Users/sunit/College_parject_data/main.py
PS C:\Users\sunit\College_parject_data> & C:/Users/sunit/AppData/Local/Programs/Python/Python310/python.exe c:/Users/sunit/College_parject_data/main.py
PS C:\Users\sunit\College_parject_data> & C:/Users/sunit/AppData/Local/Programs/Python/Python310/python.exe c:/Users/sunit/College_parject_data/main.py
PS C:\Users\sunit\College_parject_data>
```
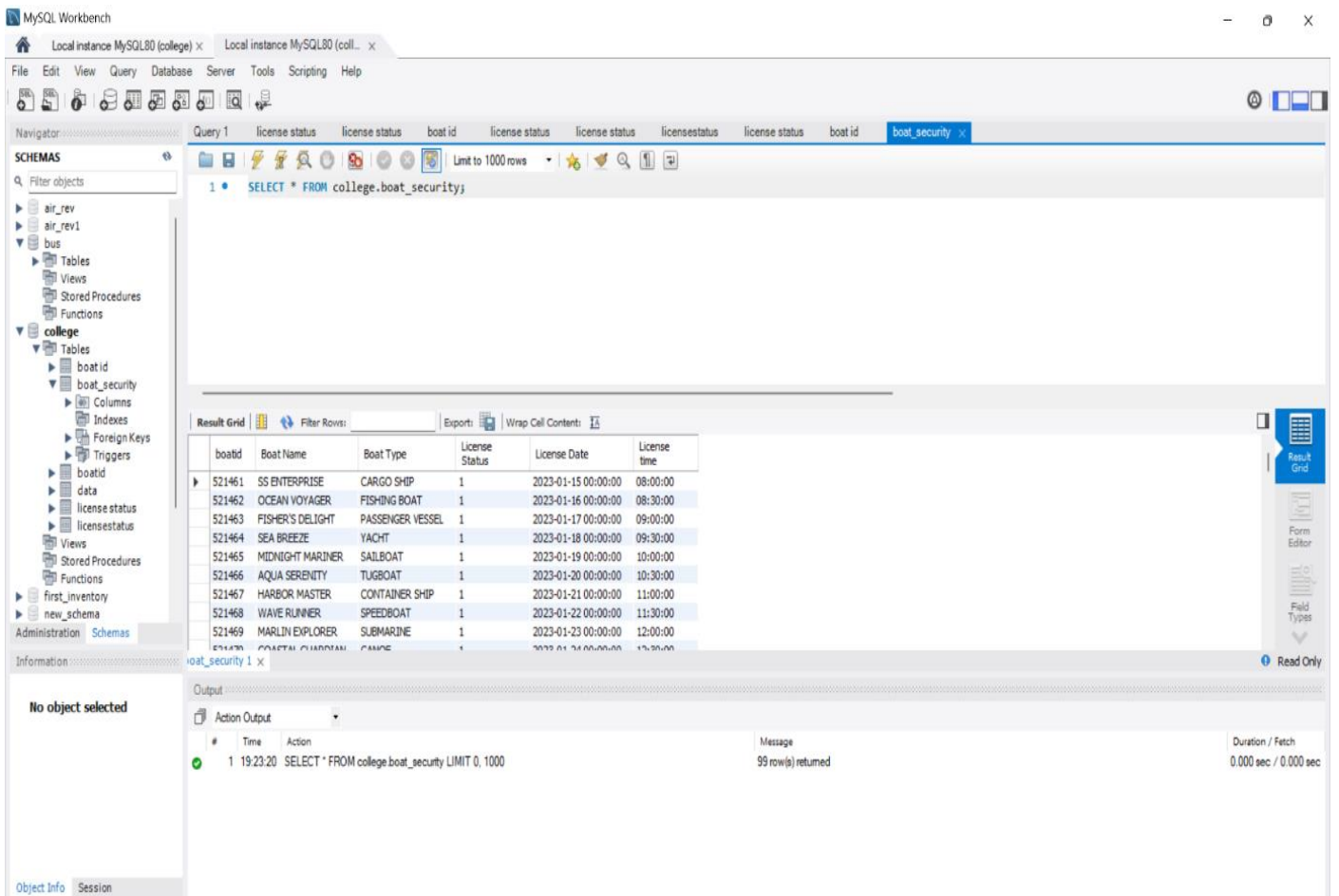
- MySQL Workbench

```sql
1 • SELECT * FROM college.boat_security;
```

| boatid | Boat Name | Boat Type | License Status | License Date | License time |
|---|---|---|---|---|---|
| 521461 | SS ENTERPRISE | CARGO SHIP | 1 | 2023-01-15 00:00:00 | 08:00:00 |
| 521462 | OCEAN VOYAGER | FISHING BOAT | 1 | 2023-01-16 00:00:00 | 08:30:00 |
| 521463 | FISHER'S DELIGHT | PASSENGER VESSEL | 1 | 2023-01-17 00:00:00 | 09:00:00 |
| 521464 | SEA BREEZE | YACHT | 1 | 2023-01-18 00:00:00 | 09:30:00 |
| 521465 | MIDNIGHT MARINER | SAILBOAT | 1 | 2023-01-19 00:00:00 | 10:00:00 |
| 521466 | AQUA SERENITY | TUGBOAT | 1 | 2023-01-20 00:00:00 | 10:30:00 |
| 521467 | HARBOR MASTER | CONTAINER SHIP | 1 | 2023-01-21 00:00:00 | 11:00:00 |
| 521468 | WAVE RUNNER | SPEEDBOAT | 1 | 2023-01-22 00:00:00 | 11:30:00 |
| 521469 | MARLIN EXPLORER | SUBMARINE | 1 | 2023-01-23 00:00:00 | 12:00:00 |
| 521470 | COASTAL GUARDIAN | CANOE | 1 | 2023-01-24 00:00:00 | 12:30:00 |

Action Output

| # | Time | Action | Message | Duration / Fetch |
|---|---|---|---|---|
| 1 | 19:23:20 | SELECT * FROM college.boat_security LIMIT 0, 1000 | 99 row(s) returned | 0.000 sec / 0.000 sec |

- Updated excel sheet



➢ **Step 5**: Conclusion, we finally performed the ETL process using the pandas' library over vscode and Excel data and after running the code for applying ETL queries we got the updated data back in the same format.

➢ Final Outcome

1. DB2
2. Watson studio
3. ETL process
4. Pandas
5. MySQL