# Big Data – Hadoop - Hive
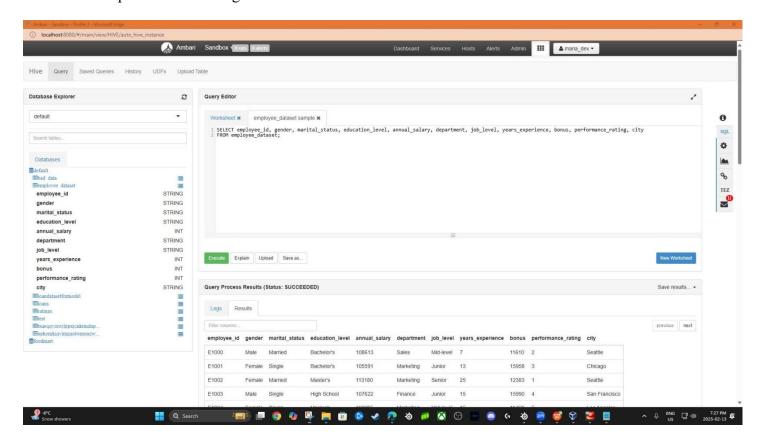
Date: 02/13/2025

**Scenario Question:**

You are a data analyst working for a company that manages a large workforce. Your manager has asked you to analyze the employee data to gather insights on employee performance, salary distribution, and department-wise headcount. You are given an employee dataset with the following columns: Employee_ID, Gender, Marital_Status, Education_Level, Annual_Salary, Department, Job_Level, Years_Experience, Bonus, Performance_Rating, and City that needs to be uploaded to **Hive** (Hadoop) for further analysis.
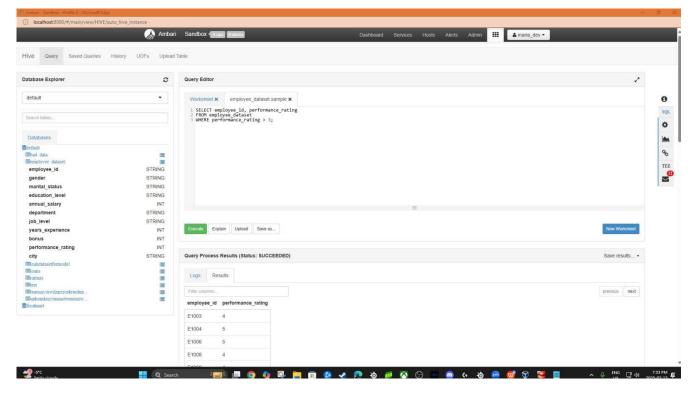
**Task 1: Select Basic Data**

- Your first task is to retrieve a list of all employees along with their key details. This includes their ID, gender, marital status, education level, annual salary, bonus, and performance rating.
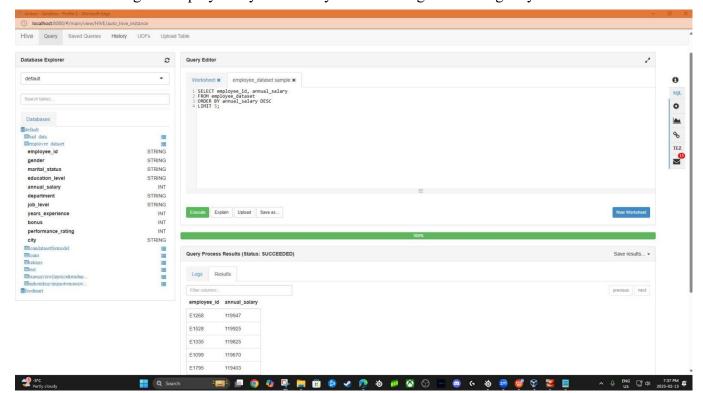
## Task 2: Filtering Data

- Next, you need to filter the data to find employees who have a **performance rating of 4 or higher**. This information will help you identify top performers in the organization.
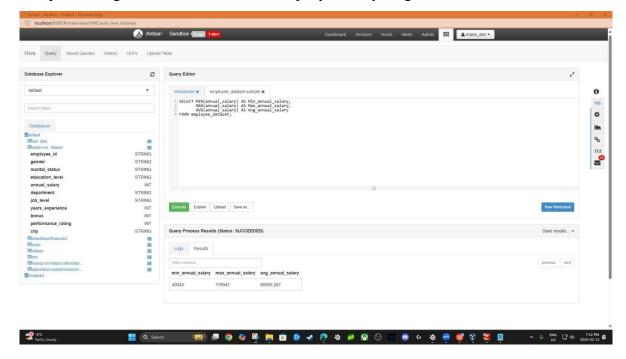


## Task 3: Sorting Data

- You are then tasked with identifying the **top 5 highest-paid employees** in the company.

- Sorting the employees by their salary in descending order will give you this information.
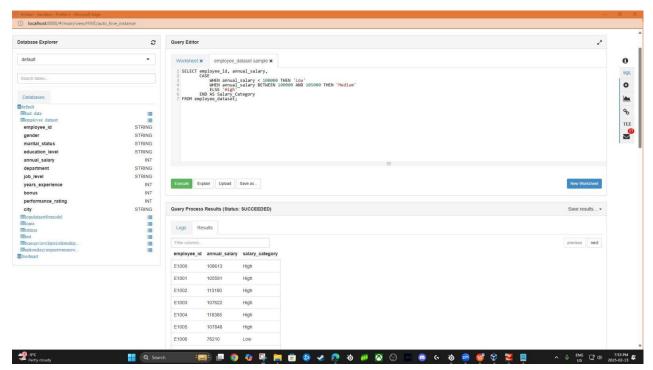
## Task 4: Salary Statistics

- Your manager is interested in understanding the distribution of employee salaries. You are required to retrieve the **minimum**, **maximum**, and **average salary** from the dataset to provide a general overview of the company's salary range.



## Task 5: Salary Category

- To make the data more insightful, you are asked to categorize employees into **"Low,"** **"Medium,"** and **"High"** salary groups based on their annual salary. This categorization will be used in further analysis.

**Task 6: Grouping by Department**

- Finally, your manager wants to see how employees are distributed across different **departments**. You will need to group the data by department and count the number of employees in each department.