# PROJECT 8

## AUTOMATED DATA QUALITY MONITORING IN CLOUD DATA WAREHOUSES

## INTRODUCTION

In today's data-driven world, organizations rely heavily on cloud data warehouses such as Snowflake, Google BigQuery, and Amazon Redshift to store and manage large volumes of critical business data. However, as data grows in volume and complexity, maintaining high data quality becomes increasingly challenging. Poor data quality—such as missing values, duplicates, anomalies, or schema mismatches can lead to inaccurate analytics, faulty business decisions, and regulatory risks. Manual monitoring is not scalable or reliable for large datasets, and often fails to catch issues in real-time.

## OBJECTIVES OF DATA QUALITY MONITORING

The goal of this project is to automate the monitoring of data quality within cloud-based data warehouses using artificial intelligence techniques. The system is designed to:

- Continuously assess key data quality metrics (e.g., completeness, consistency, accuracy).

- Automatically detect anomalies that could indicate data corruption or process failures.

- Generate real-time alerts via email when data integrity issues are found.

- Support strong data governance practices by ensuring reliable, clean, and trustworthy data.

## TOOLS AND PACKAGES REQUIRED

- **Jupyter Notebook:** The development environment used to create, test, and present the system step by step.

- **pandas:** For loading, processing, and validating structured data (typically from cloud exports in CSV or database queries).

- **scikit-learn (sklearn):** For building AI models like Isolation Forest to detect outliers and data anomalies.

- **smtplib:** For sending automated email alerts when issues are detected in the data.

# WORKFLOW OVERVIEW

## 1. Data Ingestion

- Load the data exported from a **cloud data warehouse** (e.g., Snowflake, BigQuery) using pandas.read_csv() or via a database connector.

- This can be scheduled or triggered regularly to align with ETL jobs.

## 2. Data Quality Checks

- **Completeness Check**: Use pandas to check for null or missing values in critical columns.

- **Uniqueness Check**: Detect duplicate records using df.duplicated().

- **Data Type Consistency**: Ensure each column matches expected data types (e.g., dates, integers).

## 3. AI-Based Anomaly Detection

- Apply **Isolation Forest** from sklearn on numeric features to automatically detect **statistical anomalies**.

- These anomalies may indicate unexpected shifts in data (e.g., unusual spikes, missing records, invalid entries).

## 4. Generate a Quality Report

- Summarize all issues found (nulls, duplicates, outliers) into a readable report.

## 5. Send Email Alerts

- Use smtplib to send the data quality report via email to data engineers or governance officers.

- This ensures that issues are addressed proactively and transparently.

## KEY FEATURES & BENEFITS

- **AI Integration**: Machine learning models identify subtle and complex patterns in data anomalies.

- **Automation**: Reduces manual inspection and makes monitoring scalable.

- **Data Governance**: Ensures high standards of trust, compliance, and data reliability.

- **Real-Time Alerts**: Stakeholders are notified immediately when issues occur.

### ⭕ Example Use Case

Imagine a business uploads sales data daily to BigQuery. This notebook can:

- Detect if daily sales suddenly drop to zero.

- Flag missing customer records.

- Identify data entry errors (e.g., negative prices).

- Email the data team immediately so they can fix upstream pipeline issues.

### ⭕ Output

- **Data Quality Summary Report** (printed and emailed)

- **Anomaly Detection Scores** per row

- **Automated Alert System** triggered by integrity issues

## CONCLUSION

This project successfully demonstrates how **AI can be used to automate data quality monitoring** within cloud-based data warehouses using a simple yet effective toolset: **Jupyter Notebook**, **pandas**, **scikit-learn**, and **smtplib**. By combining traditional data validation techniques with machine learning algorithms like **Isolation Forest**, the system is capable of continuously detecting anomalies and integrity issues in real-time.

The integration of **automated alerting via email** ensures that data stakeholders are promptly notified of any issues, allowing for quicker resolution and maintaining the **high standards of data governance** required in modern data-driven organizations. Furthermore, the use of open-source tools makes the solution cost-effective, customizable, and easily scalable.