

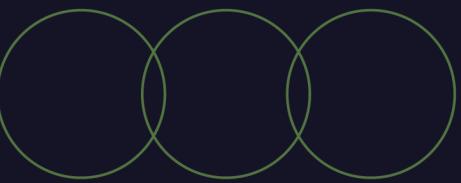
AI INNOVATIONS LAB

AI-Powered Data Cleaning System

Presented by Anant Arya, Ansh



Importance of Data Cleaning



Understanding the impact

Noisy Data

Raw data is often filled with **inaccuracies** and inconsistencies. These errors can lead to misleading insights and poor decision-making, ultimately affecting the quality of machine learning models.

Missing Values

When datasets contain **missing values**, they can distort analysis and lead to incorrect conclusions. Handling these gaps appropriately is crucial for maintaining the integrity of the data.

Duplicate Records

Duplicates can inflate metrics and skew results, creating an illusion of data reliability. Identifying and removing these records is essential to ensure precise analysis and reporting.

Project Objectives

Key goals for effective data cleaning and analysis

Detect Anomalies

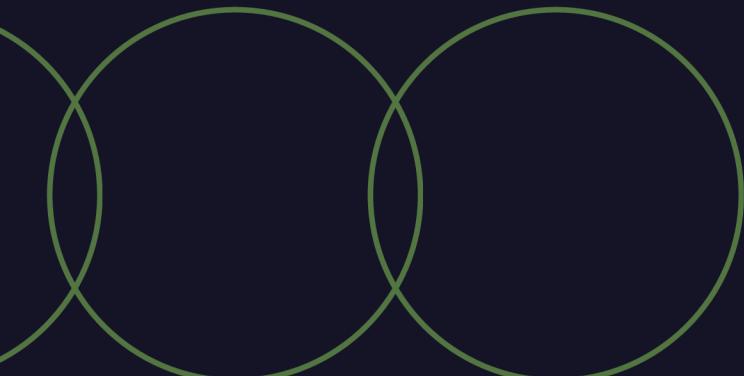
Our first objective is to **detect anomalies** in the dataset, ensuring that outliers are identified and flagged for further analysis.

Handle Missing Values

Addressing missing values is crucial; we implement strategies to fill these gaps, ensuring a complete and accurate dataset.

Identify Duplicates

We aim to identify and remove duplicate entries, which can skew metrics and lead to incorrect conclusions in data analysis.



Tools and Technologies



Essential components for success

Python Libraries

Python offers powerful libraries like **Pandas** for data manipulation and **sklearn** for machine learning, making it crucial for implementing effective data cleaning and preprocessing workflows.

Streamlit Framework

Streamlit enables the creation of **interactive web applications** effortlessly, allowing users to visualize data cleaning processes and results in real time, which enhances user experience and engagement.

MySQL Database

MySQL is employed for **storing and managing** raw datasets, providing a robust database solution that facilitates data retrieval and integration into the cleaning pipeline seamlessly.

System Architecture

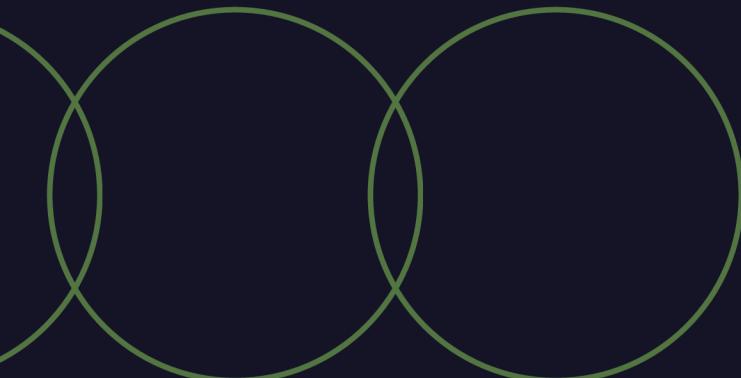
Overview of data flow and processing components

MySQL Database

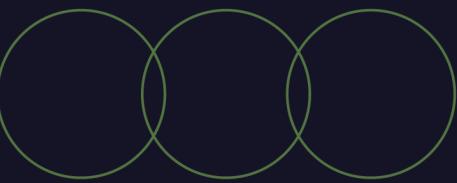
MySQL serves as the foundational storage for **raw datasets**, ensuring organized and efficient data retrieval for subsequent processing and analysis tasks within the system architecture.

Python Logic

Python orchestrates the data cleaning logic, utilizing libraries to manipulate data and integrate the **anomaly detection model**, thereby enhancing the overall data quality and reliability of outputs.



Database Design Overview



Structure of customer data

Customer Data Table

The **customer_data** table is crucial for storing user information. It includes essential fields such as **ID**, **name**, **email**, **age**, and **salary** to manage clients effectively.

Fields Overview

Each field serves a significant purpose: the **ID** uniquely identifies records, while **name**, **email**, **age**, and **salary** provide vital information necessary for analysis and customer interactions.

Data Quality Issues

This table often contains **messy** and **duplicate rows**, which complicate data analysis. Addressing these issues is necessary to ensure reliable results and maintain the integrity of the dataset.

Anomaly Detection

Harnessing IsolationForest for
intelligent outlier detection

Unsupervised Model

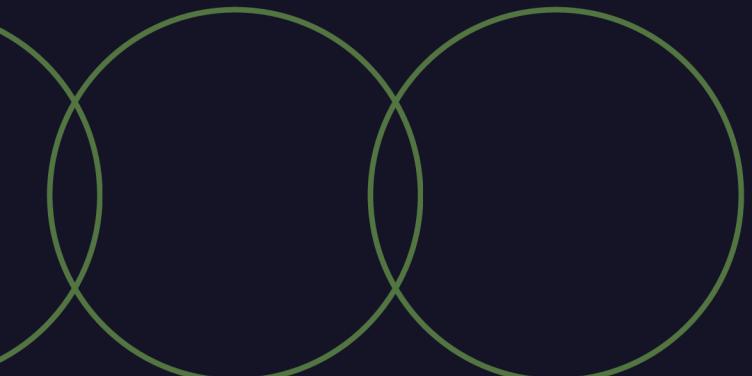
The IsolationForest model operates without labeled data, identifying anomalies based on the **structure of the dataset** rather than predefined categories.

Flags Outliers

This model effectively flags outliers by measuring their **isolation** within the data distribution, highlighting unusual instances that may distort analysis.

Salary Anomalies

IsolationForest is particularly useful for detecting salary anomalies, ensuring that extreme values are identified and addressed for accurate **data interpretation**.



Streamlit Application

Interactive interface for real-time
data cleaning tasks

Load Data

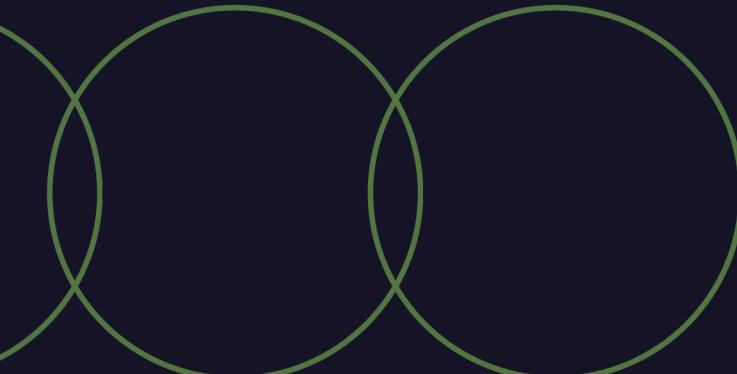
Users can effortlessly **load datasets** directly from the MySQL database, simplifying the data cleaning process significantly.

Display Issues

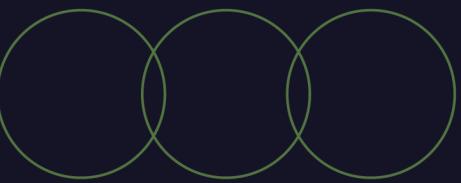
The application highlights **data quality issues** visually, allowing users to quickly identify missing values and anomalies for resolution.

Cleaning Options

Streamlit provides intuitive **cleaning options** for users, enabling them to fill missing values, remove duplicates, and flag outliers effortlessly.



Data Cleaning Pipeline



Essential steps for accuracy

Remove Duplicates

Removing duplicates is crucial to ensure that each record is unique, which prevents skewing analysis results and gives a clearer picture of the dataset's true characteristics.

Fill Missing Values

Filling missing values is essential for maintaining data integrity, as it allows for more accurate calculations and prevents gaps that could lead to misleading insights during analysis.

Remove Anomalies

Anomaly removal enhances the reliability of datasets by eliminating outliers that can distort statistical models, ensuring that machine learning algorithms operate on valid and representative data.

Results Summary

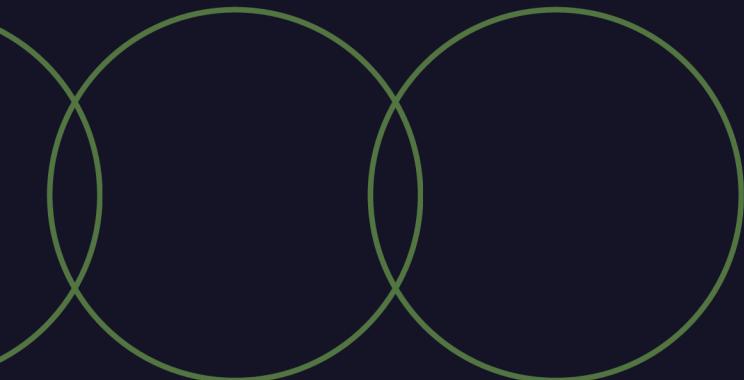
Demonstrating the impact of
effective data cleaning

Clean Dataset

A **clean dataset was generated**, showcasing significant improvements in data quality that enhance the accuracy of machine learning models and analytics.

Outlier Removal

All identified outliers were successfully **removed from the dataset**, ensuring that analysis outputs are reliable and reflective of real-world values without bias.



Conclusion

Summary of findings from the data cleaning project

AI Enhancements

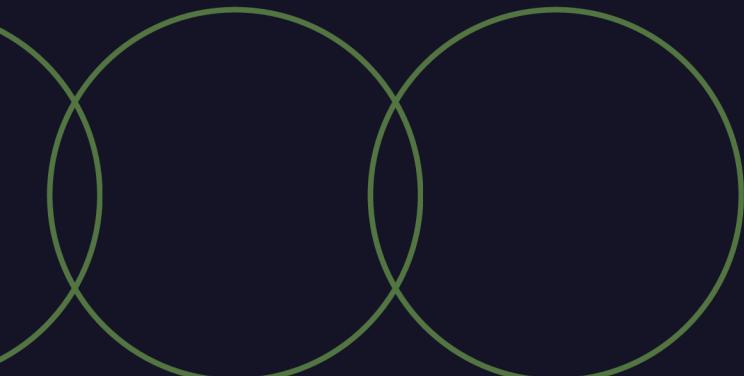
The integration of AI significantly enhances the **efficacy of data cleaning**, leading to higher accuracy in machine learning models.

MySQL Stability

Utilizing MySQL ensures a **robust data storage solution**, providing seamless access and management of raw and cleaned datasets.

Interactive Pipeline

The interactive cleaning pipeline created with Streamlit offers **real-time engagement**, allowing users to visualize and rectify data issues effectively.



Thank You

