Data Analytics

# CleanStream Framework ✦

**BI Data Cleaning**
**By Rakesh Mahakur**          8637291540          rakeshmahakur555@gmail.com

# CleanStream Framework: Automating Data Cleaning for BI

CleanStream revolutionizes data preparation by automating essential tasks like **duplicate removal**, handling missing values, and ensuring schema consistency, leading to reliable BI and ML insights.



DATA CLEANING

| 01 | 01 | 2. | 104 | 034 | 38 | 13 |
|----|----|----|-----|-----|-----|-----|
| CSV | Data Ingersion | Data Pipeline | Output | Data Checks | Standardization | Filte ceacing |
| | Upfine | Data Oryato | Data | | | Compliication |

# Architecture Overview

The **CleanStream Framework** processes raw data efficiently, transforming it into BI-ready output. By utilizing automated cleaning techniques, it ensures that the data is standardized and reliable for further analysis and reporting within various BI tools.



The architecture of the CleanStream Framework is designed to seamlessly integrate various stages of data handling. From raw data ingestion to profiling and final output, this structured approach ensures that users receive **clean and actionable data**, enhancing the decision-making process. Advanced automation reduces manual intervention, thereby increasing efficiency and accuracy across BI applications. This framework sets the foundation for robust data analysis and reporting.

# Environment Setup

**Essential technologies and tools for the CleanStream framework**

The CleanStream framework leverages **critical technologies** to ensure efficient data cleaning and automation. Key components like Python and Pandas enable powerful data manipulation and analysis, essential for the functionality of BI systems.

Utilizing these technologies ensures that the framework can handle various data cleaning tasks seamlessly. The **reusable folder structure** promotes organization and scalability, making it easier for users to maintain and expand their data projects. Ensuring compatibility with the operating system and datetime functions further enhances its capabilities.

*Essential technologies*

## Python

Python is a versatile language known for its simplicity and robustness, widely used in data analysis.

## Pandas

Pandas is a powerful library that simplifies data manipulation and analysis, enabling efficient handling of data frames.

## OS & Datetime

These modules facilitate interaction with the operating system and manage date and time operations, crucial for data workflows.

## Reusable Structure

A well-organized folder structure aids in maintaining project clarity, ensuring easy access to scripts and data files.

# Ingestion

## REVENUE
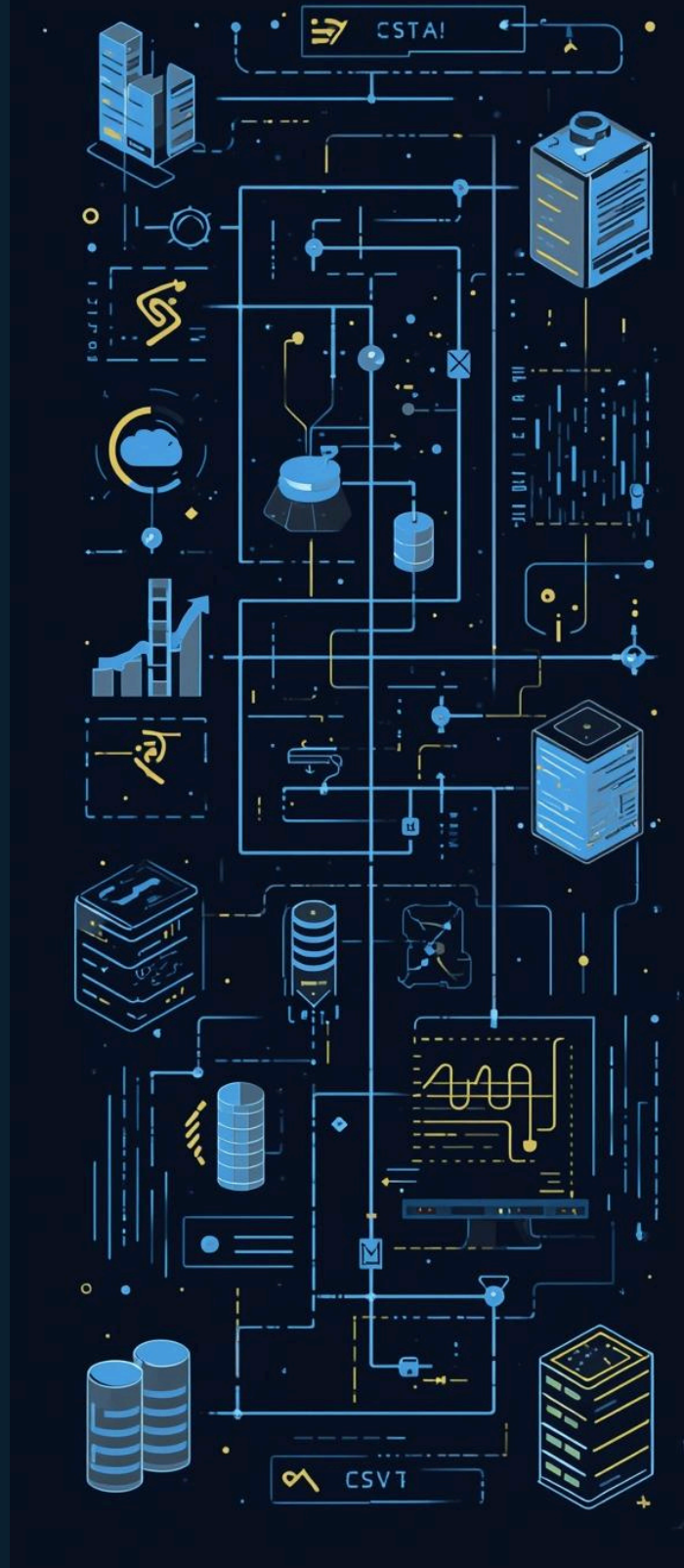
**Income**                    $120,000

## EXPENSES

**Total**                     $45,000

## PROFIT

**Net**                       $75,000

## DETAILED BREAKDOWN

**Marketing,**                $10,000,
**Operations, R&D**           $15,000,
                              $20,000



# Financial Summary

$120,000          $45,000           $75,000
Total Revenue     Total Expenses    Net Profit

# Data Profiling

Understanding the dataset is crucial for effective data cleaning and analysis. **Accurate profiling** ensures that data quality is assessed and informs decisions on handling missing values, duplicates, and overall schema integrity in the CleanStream framework.



In data profiling, various metrics such as dataset shape, schema inspection, and missing value summaries are evaluated. This analysis empowers data analysts to identify potential issues before processing, thus improving the reliability and efficiency of subsequent automation steps in the CleanStream framework. By visualizing the data, stakeholders gain insights into its structure, leading to better informed decisions during the cleaning process.

# Cleaning & Standardization

**Essential transformations for improving data quality and reliability**

The **Cleaning & Standardization** phase is crucial for ensuring high-quality data. This stage involves transforming raw data into a consistent format, enabling reliable analysis. By automating these processes, CleanStream enhances efficiency and accuracy in business intelligence applications.

Effective cleaning and standardization encompass several key tasks. Duplicate removal addresses redundancies, while filling in missing values with appropriate defaults ensures completeness. Additionally, standardizing date formats eliminates inconsistencies. These transformations allow for smoother integration and analysis of data across various BI tools, ultimately supporting informed decision-making.

*Ensuring data quality*

## Duplicate removal

The process identifies and eliminates repeated entries, ensuring each record is unique.

## Text fill

Missing text values are replaced with "Unknown", maintaining dataset integrity and usability.

## Numeric fill

Numeric fields with missing values are filled with "0", providing a baseline for analysis.

## Date standardization

All date formats are unified to a single standard, facilitating easier data manipulation and analysis.

# Output

### REVENUE

**Total**                    $150,000

### EXPENSES

**Total**                    $45,000

### PROFIT

**Net**                      $105,000

### EXPENSE DETAILS

**Monthly Operating**        $5,000,
**Costs**                    $3,000,
                             $2,500,
                             $1,500

# Financial Summary

**$150,000**          **$45,000**          **$105,000**
Total Revenue         Total Expenses       Net Profit

# Future Enhancements

**Expanding CleanStream's capabilities for improved data automation and analysis**

Future enhancements to the CleanStream framework will focus on **increasing efficiency** and providing advanced functionalities that cater to evolving data needs. Implementing these improvements will streamline processes and ensure high-quality data management across various applications.

Key enhancements include a **file watcher** that automatically detects new data files, a comprehensive data quality report to assess data integrity, and cloud integration for seamless access and storage. These advancements will not only enhance performance but also empower users with greater control over their data workflows.

*Enhanced capabilities ahead*

## File watcher

A file watcher will enable automatic detection and processing of new CSV files.

## Data quality report

A data quality report will provide insights into data integrity and cleanliness.

## Cloud integration

Cloud integration will facilitate seamless data access and storage options for users.

## Advanced analytics

Advanced analytics will empower users with deeper insights and data-driven decision making.

# Conclusion: Showcasing the Impact of CleanStream Automation

## Success

CleanStream exemplifies strong **automation principles** and effective data management, making it a valuable addition to any portfolio. Its successful implementation demonstrates the ability to enhance data cleaning workflows.

- **Streamlined data processes**
- **Enhanced data quality**
- **Improved BI readiness**
- **Showcased technical skills**