# Introduction

Real-world datasets are often incomplete, inconsistent, and noisy. Issues such as missing values, duplicate records, and outliers can significantly affect the accuracy and reliability of data analysis and machine learning models. Data preprocessing is therefore a critical step before performing any meaningful analysis or model training.

This project focuses on building a **Data Cleaning & Preprocessing Tool** using Python that automates essential preprocessing tasks. The tool prepares raw data into a clean, structured, and reliable format suitable for further analysis or machine learning applications.

# Problem Statement

Raw datasets collected from real-world sources frequently contain missing values, duplicate entries, and abnormal data points (outliers). These issues can lead to biased results, incorrect predictions, and poor model performance if not handled properly.

Manually cleaning such data is time-consuming and error-prone. Hence, there is a need for an automated preprocessing tool that efficiently cleans datasets and ensures data quality before analysis or modeling.

# Objectives
- Identify and handle missing values
- Detect and remove duplicate records
- Identify and treat outliers
- Convert raw data into a clean and consistent format
- Save the cleaned dataset for further use

# Tools & Libraries Used
- Python
- Pandas
- NumPy
- Matplotlib
- Seaborn
- Jupyter Notebook

# Methodology
1. Data loading from CSV
2. Missing value imputation using mean
3. Duplicate removal
4. Outlier detection using IQR and domain rules
5. Data type correction (Age flooring)
6. Saving cleaned data

# Results

The dataset was successfully cleaned by handling missing values, removing duplicates, correcting invalid values, and eliminating outliers. The final dataset is consistent and ready for analysis or machine learning.

## Applications
- Machine learning preprocessing
- Data analysis and visualization
- Academic projects
- Business analytics

# Conclusion

This project demonstrates the importance of data preprocessing in real-world datasets. Automating cleaning tasks improves data quality, reduces errors, and enhances the reliability of downstream analysis and models.

# Future Scope
- GUI using Streamlit or Flask
- Advanced imputation techniques
- Feature scaling and encoding
- Integration with ML pipelines