# Data Analyzer & Cleaner
**A PROJECT REPORT**
for
**Mini Project-I (K24MCA18)**
**Session (2024-25)**
**Submitted by**
**Chetanya Bedi**
**202410116100053**
**Deepak Sharma**
202410116100055
**Devansh Kumar**
202410116100059
**Dhruv Bathla**
202410116100062

**Submitted in partial fulfilment of the**
**Requirements for the Degree of**

# MASTER OF COMPUTER APPLICATION

**Under the Supervision of**
**Dr. Vipin Kumar**



**Submitted to**

**DEPARTMENT OF COMPUTER APPLICATIONS**
**KIET Group of Institutions, Ghaziabad**
**Uttar Pradesh-201206**

# CERTIFICATE

Certified that **Chetanya Bedi (202410116100053)**, **Deepak Sharma (202410116100055)**, **Devansh Kumar(202410116100059), Dhruv Bathla (202410116100062)** have carried out the project work having "**Data Analyzer and Cleaner**" (**Mini Project-I, K24MCA18**) for **Master of Computer Application** from Dr. A.P.J. Abdul Kalam Technical University (AKTU**)** (formerly UPTU), Lucknow under my supervision. The project report embodies original work, and studies are carried out by the student himself/herself and the contents of the project report do not form the basis for the award of any other degree to the candidate orto anybody else from this or any other University/Institution.

**Dr. Vipin Kumar**                    **Dr. Arun Kr. Tripathi**
**Assistant Professor**                **Dean**
**Department of Computer Applications**    **Department of Computer Applications**
**KIET Group of Institutions, Ghaziabad**    **KIET Group of Institutions, Ghaziabad**

# ABSTRACT

In the age of big data, organizations rely heavily on accurate and well-structured data for decision-making, analytics, and machine learning applications. However, raw data often contains inconsistencies, missing values, duplicate records, and outliers, making it unsuitable for direct analysis. The **"Data Analyzer and Cleaner"** project aims to develop an automated system that efficiently processes raw datasets, improving their quality and usability.

This system will incorporate **data cleaning, transformation, and statistical analysis** to enhance data reliability. Key functionalities include **handling missing values** through imputation techniques, **detecting and removing duplicate entries**, **identifying outliers**, and **standardizing formats**. Additionally, the system will provide **descriptive statistics and data visualization** to help users gain insights into the dataset's structure and quality. Machine learning models and rule-based techniques will be employed to automate error detection and correction.

The project will be implemented using technologies such as **Python, Pandas, NumPy, and Scikit-learn**, ensuring efficient data processing. A **user-friendly interface** will allow users to upload datasets, apply cleaning operations, and export refined data for further analysis. This tool will be beneficial for **data scientists, researchers, and business analysts** who require high-quality datasets for their work.

By automating data preprocessing, this system reduces manual effort, minimizes errors, and accelerates analytical workflows. The **Data Analyzer and Cleaner** project will play a crucial role in streamlining data preparation, ultimately leading to **more accurate predictions, better insights, and improved decision-making** in various industries such as finance, healthcare, and e-commerce.

# ACKNOWLEDGEMENTS

# TABLE OF CONTENTS

# Chapter 1
# INTRODUCTION

## 1.1 Project Description

In today's data-driven world, organizations and individuals rely on accurate and well-structured data for analysis, forecasting, and decision-making. However, raw data often contains inconsistencies, missing values, duplicate records, and outliers, making it unsuitable for direct use. The **"Data Analyzer and Cleaner"** project aims to develop an automated system that streamlines the process of cleaning, analyzing, and preparing data for further processing. This system will focus on essential data preprocessing tasks such as handling missing values using imputation techniques, removing duplicate entries, detecting and managing outliers, and standardizing data formats to ensure consistency. Additionally, it will provide statistical analysis and data visualization features to help users understand the dataset's structure and quality.

The project will be developed using **Python**, leveraging powerful data-processing libraries such as **Pandas, NumPy, and Scikit-learn**, along with visualization tools like **Matplotlib and Seaborn**. A **user-friendly interface** will enable users to upload raw datasets, apply cleaning operations, and export the refined data for further analysis. By automating these critical tasks, the system reduces the need for manual intervention, minimizes errors, and enhances the efficiency of data preprocessing workflows.

The **"Data Analyzer and Cleaner"** will be an invaluable tool for **data scientists, researchers, business analysts, and professionals** who work with large datasets. High-quality, well-processed data leads to more accurate models, better insights, and improved decision-making across various industries, including finance, healthcare, and e-commerce. This project will play a crucial role in optimizing data workflows, ensuring that users have access to clean, structured, and reliable data, ultimately making the process of data analysis more efficient, accurate, and accessible.

## 1.2  Project Scope

The **"Data Analyzer and Cleaner"** project focuses on developing an automated system for preprocessing raw datasets, ensuring data accuracy, consistency, and usability. It will include functionalities such as **missing value imputation, duplicate removal, outlier detection, data standardization, and statistical analysis**. The system will support multiple data formats and provide **visual insights** to help users understand dataset quality. Implemented using **Python, Pandas, NumPy, and Scikit-learn**, it will feature a **user-friendly interface** for seamless interaction. The project aims to assist **data scientists, analysts, and businesses** by reducing manual data-cleaning efforts and improving efficiency. It will enhance **data-driven decision-making** in industries like finance, healthcare, and e-commerce. The system's scalability will allow it to handle large datasets while maintaining performance. Ultimately, the project will provide a **reliable and automated solution** for refining raw data, making it ready for further analysis and predictive modeling.

## 1.3  Project Overview

- **Data Cleaning** : Removes duplicates, corrects inconsistencies, and standardizes formats.
- **Missing Value Handling**: Uses imputation techniques to fill missing data and analyze all the values From the data.
- **Outlier Detection**: Identifies and removes anomalies from datasets.
- **Statistical Analysis**: Computes key metrics (mean, median, mode, variance, etc.).
- **User-friendly Interface**: Allows users to upload, process, and export cleaned data.

# Chapter 2

# Feasibility Study

The **feasibility study** of the **"Data Analyzer and Cleaner"** project evaluates its technical, economic, and operational viability. Technically, it leverages **Python, Pandas, NumPy, and machine learning** to automate data cleaning and analysis. Economically, it is cost-effective, reducing manual effort and improving efficiency. Operationally, its **user-friendly interface** ensures accessibility for researchers, analysts, and businesses. With increasing data-driven applications, this tool enhances data quality, supports better decision-making, and is highly feasible for deployment across industries like finance, healthcare, and e-commerce.

## Market Feasibility

The market feasibility of the **"Data Analyzer and Cleaner"** project is highly promising, given the increasing reliance on data-driven decision-making across industries. Businesses, researchers, and analysts require clean, accurate data for insights and predictions. With growing big data adoption, this tool addresses a critical need, ensuring efficiency, accuracy, and automation in data preprocessing.

## Technical Feasibility

The technical feasibility of the "Data Analyzer and Cleaner" project is high, as it utilizes well-established technologies like Python, Pandas, NumPy, and Scikit-learn for data processing. The system will efficiently handle large datasets, automate cleaning tasks, and provide insights through statistical analysis and visualization, ensuring seamless usability and integration for data professionals.

## Economical Feasibility

The economic feasibility of the "Data Analyzer and Cleaner" project is highly favorable, as it reduces the time and cost associated with manual data cleaning. By automating preprocessing tasks, businesses and researchers can save resources while improving data accuracy. The project's implementation using open-source tools further minimizes development and operational costs.

## Operational Feasibility

The operational feasibility of the "Data Analyzer and Cleaner" project is high, as it automates data preprocessing, reducing manual effort and errors. With a user-friendly interface and efficient algorithms, it ensures seamless integration into existing workflows. Its adaptability benefits researchers, analysts, and businesses, enhancing data quality for improved decision-making and analytics.

## Schedule Feasibility

The **schedule feasibility** of the **"Data Analyzer and Cleaner"** project ensures timely completion by following a structured development plan. The project will be divided into phases: **requirement analysis, design, implementation, testing, and deployment**. With efficient task management, milestone tracking, and iterative testing, the project is expected to be completed within the planned timeframe.

# Chapter - 3
# Project Objective

The objective of the **"Data Analyzer and Cleaner"** project is to develop an efficient and automated system for preprocessing raw datasets by identifying and rectifying inconsistencies, missing values, duplicate records, and outliers. The system aims to enhance data quality through **cleaning, transformation, and statistical analysis**, ensuring accuracy and reliability for further processing.

- **Automated Data Cleaning :** Detect and handle missing values, duplicates, and inconsistencies.

- **User-Friendly platform:** Enable users to upload, analyze, and export cleaned datasets easily.

- **Enhanced Data Insights :** Provide statistical summaries and visualizations for better understanding.

- **Data Transformation:** Standardize formats, normalize values, and improve data structure.

- **Performance Optimization & Automation :** Implement batch processing and scheduled automation for continuous data cleaning.

- **Descriptive Analysis and Insights :** Generate summary statistics like mean, median, standard deviation, etc. ,

- **Automated Data Cleaning :** Identify and fill missing values using various imputation techniques (mean, median, mode, or ML-based).
-

# Chapter 4

# Hardware and Software Requirement

## Hardware Requirements

- **Processor:** Intel Core i3 (8th Gen) or AMD Ryzen 3

- **RAM:** 4 GB

- **Storage:** 128 GB SSD or 500 GB HDD

- **Graphics Card:** Integrated GPU

- **Display:** 1366x768 resolution monitor

- **Network:** Basic internet connectivity for cloud-based data operations (if required)

## Software Requirements

1. **Operating System:** Windows 10/11, Linux (Ubuntu), or macOS
2. **Programming Language:** Python (Version 3.8 or later)
3. **Libraries:** Pandas, NumPy, Matplotlib, Seaborn, Scikit-learn (for data processing & visualization)
4. **Development Tools:** Jupyter Notebook, VS Code, or PyCharm
5. **Database (if needed):** MySQL, SQLite, or MongoDB
6. **Web Framework (if needed):** Flask or Django for web-based interface
7. **Other Tools:** Git (for version control), Anaconda (for managing libraries), and cloud services (if required)

# Chapter - 5
# Project Flow of "Data Analyzer and Cleaner"

The project follows a structured flow to ensure efficient data processing and cleaning. The key steps are:

1. **Data Input:**
   - **User uploads raw data files (CSV, Excel, JSON, etc.).**
   - **Data can also be fetched from a database or API.**

2. **Data Preprocessing:**
   - **Detect and handle missing values (imputation or removal).**
   - **Remove duplicate records.**
   - **Standardize formats (date, time, text, and categorical values).**

3. **Outlier Detection and Removal:**
   - **Identify anomalies using statistical methods (Z-score, IQR).**
   - **Provide options to remove or correct outliers.**

4. **Data Transformation:**
   - **Normalize and scale numerical data.**
   - **Convert categorical data into a uniform format.**
   - **Encode text-based data for analysis.**

5. **Data Analysis & Visualization:**
   - **Generate summary statistics (mean, median, standard deviation, etc.).**
   - **Provide graphical insights (charts, histograms, correlation plots).**

6. **Data Export & Storage:**
   - **Save cleaned data in multiple formats (CSV, Excel, JSON).**
   - **Option to store processed data in a database for future use.**

7. **User Interaction & Feedback:**
   - **User reviews processed data before final export.**
   - **System logs errors and improvements for future updates.**

# Flow Chart:

Flowchart is a diagrammatic representation of sequence of logical steps of a program. Flowcharts use simple geometric shapes to depict processes and arrows to show relationships and process/data flow.
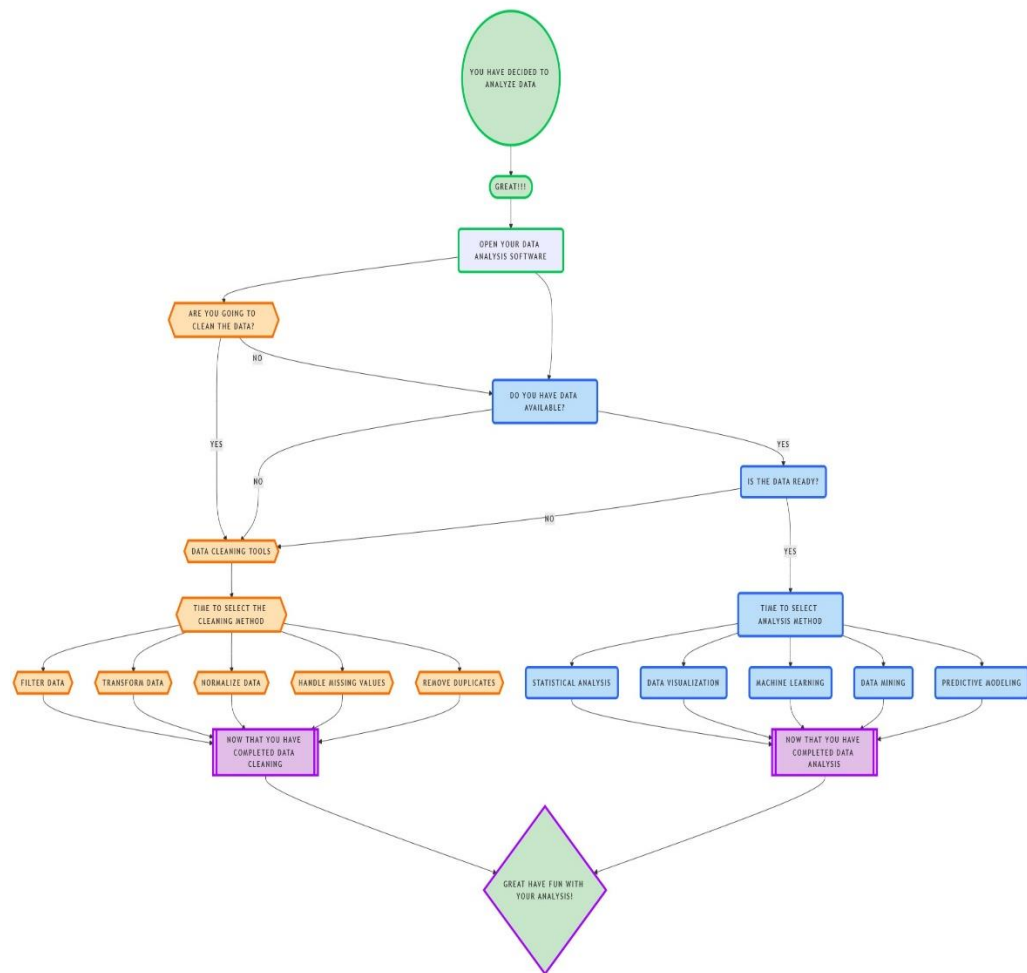


Fig 5.1  Flowchart representation  of         Data Analyzer and Cleaner

## Entity Relationship Diagram:

• ER model stands for an Entity-Relationship model. It is a high-level data model. This model is used to define the data elements and relationship for a specified system.

• It develops a conceptual design for the database. It also develops a very simple and easy to designview of data.

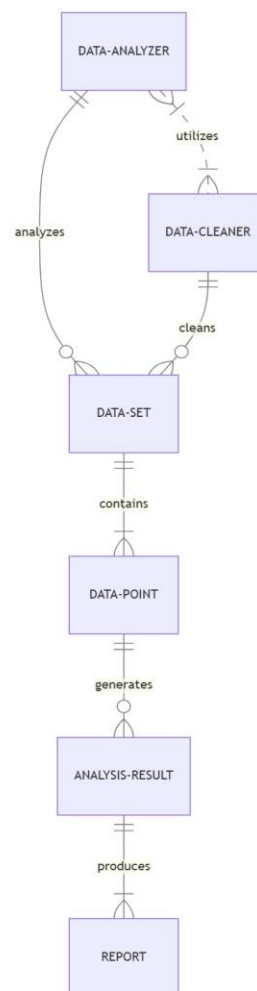• In ER modelling, the database structure is portrayed as a diagram called an entity-relationshipdiagram.



**Fig 5.2  ER Diagram representation  of Data Analyzer and cleaner**

# Chapter - 6
# Project  Outcome

The "Data Analyzer and Cleaner" project will result in a fully functional system that automates data preprocessing, improving accuracy, consistency, and reliability. The key outcomes include:

1. **Clean and Well-Structured Data:**

   - **Removal of duplicates, missing values, and inconsistencies.**
   - **Standardized and formatted datasets ready for analysis.**

2. **Improved Data Quality for Analysis & Decision-Making:**

   - **Enhanced data accuracy leads to better insights and predictions.**
   - **Reduction in errors for machine learning and business analytics.**

3. **Automated Data Processing:**

   - **Saves time and effort by reducing manual data cleaning.**
   - **Batch processing and scheduled automation for efficiency.**

4. **User-Friendly Interface:**

   - **Simple UI for uploading, analyzing, and exporting cleaned data.**
   - **Easy access to statistical summaries and visual insights.**

5. **Scalability and Integration:**

   - **Support for various data formats (CSV, Excel, JSON).**
   - **Compatibility with databases and cloud storage for future scalability.**

# REFERENCES

1. https://www.w3schools.com/html/html_css.asp
2. https://developer.mozilla.org/en-US/docs/Web/JavaScript
3. https://nodejs.org/en
4. https://www.geeksforgeeks.org/javascript/
5. https://www.mongodb.com/
6. https://github.com/mongodb/mongo
7. https://www.simplilearn.com/tutorials/nodejs-tutorial/what-is-nodejs