

# Data Leakage Detection And Intelligent Data Preprocessing System

Pakirathan K<sup>1</sup>, Murshith Ahamed Eithirish<sup>2</sup>, Gokul<sup>3</sup>

<sup>1, 2, 3</sup>Dept of Artificial Intelligence And Data Science

<sup>1, 2, 3</sup>Sir Issac newton college of engineering and technology, nagapattinam, tamilnadu, india

**Abstract-** Data leakage is one of the most critical challenges in machine learning systems, leading to unrealistic model performance and poor generalization in real-world applications. Leakage occurs when information from outside the training dataset is inadvertently used during the model training process, causing biased predictions and overly optimistic evaluation metrics. Detecting such leakage before model development is essential for building reliable and robust machine learning systems. This paper proposes a Data Leakage Detection and Intelligent Data Preprocessing System that automatically identifies potential leakage sources in datasets prior to model training. The system integrates dataset profiling, leakage detection, preprocessing techniques, and visualization tools within a Flask-based web application. Users can upload datasets, analyze data quality, detect different types of leakage such as target leakage and temporal leakage, and apply safe preprocessing operations. The system also provides interactive data visualizations and exports a cleaned dataset ready for machine learning tasks. By combining leakage detection with automated preprocessing, the proposed solution improves model reliability, reduces human error, and enhances the overall machine learning workflow.

**Keywords:** Data Leakage Detection, Intelligent Data Preprocessing, Machine Learning, Dataset Profiling, Target Leakage, Temporal Leakage, Data Visualization, Feature Engineering, Flask Web Application, Data Quality, Model Reliability, Automated Preprocessing.

## I. INTRODUCTION

Machine learning models rely heavily on the quality, integrity, and reliability of the data used during training. Poor data quality or incorrect preprocessing practices can significantly affect the performance of predictive models. One of the most common yet often overlooked problems in machine learning pipelines is **data leakage**. Data leakage occurs when information from the testing dataset or future observations is unintentionally included in the training data. This can happen during preprocessing steps such as feature scaling, feature engineering, encoding, or incorrect data

splitting. When leakage occurs, models may show extremely high accuracy during training but fail to perform effectively when deployed in real-world scenarios. Traditional machine learning workflows rely heavily on manual data inspection and developer expertise to identify leakage risks. However, this approach is time-consuming and prone to errors. Therefore, there is a growing need for automated systems that can detect leakage patterns and ensure safe data preprocessing. To address these challenges, this research proposes a **Data Leakage Detection and Intelligent Data Preprocessing System** that automates the process of dataset analysis, leakage detection, and data preprocessing. The system helps data scientists and researchers prepare reliable datasets for machine learning applications.

## II. RELATED WORK

Several studies have explored the impact of data quality and leakage prevention in machine learning systems. Researchers have identified data leakage as a significant factor that affects model reliability and evaluation accuracy. Improper preprocessing steps such as applying scaling before data splitting or using future information during training are common causes of leakage. Machine learning frameworks such as **Scikit-learn** provide various preprocessing tools including scaling, encoding, and feature transformation. However, these libraries rely on users to apply preprocessing steps correctly. If used incorrectly, these operations can introduce leakage into the dataset. Visualization libraries such as **Matplotlib**, **Seaborn**, and **Plotly** help in understanding data distributions and relationships between features. However, they do not provide direct mechanisms for identifying leakage risks. Recent research has proposed statistical techniques such as **correlation analysis**, **feature importance analysis**, and **dataset validation methods** to detect suspicious relationships between features and target variables. Despite these advancements, there is still a lack of integrated systems that combine leakage detection, preprocessing, and visualization in a single platform. The proposed system addresses this limitation by providing a unified framework that performs dataset profiling, leakage detection, preprocessing, and visualization automatically.

### III. PROPOSED SYSTEM

The proposed system is designed to automatically detect data leakage and perform intelligent data preprocessing before machine learning model development. The system integrates multiple modules that analyze dataset structure, detect leakage patterns, perform preprocessing operations, and generate visualization reports. The architecture ensures that the dataset used for machine learning models is clean, validated, and free from potential leakage risks.

#### A. Dataset Upload Module

The dataset upload module allows users to upload CSV datasets through the web interface. The system validates the uploaded file to ensure it is in the correct format and free from corruption. Once uploaded, the dataset is stored securely on the server for further analysis. This module acts as the entry point of the system and ensures that only valid datasets are processed.

#### B. Dataset Profiling Module

The dataset profiling module analyzes the structure of the uploaded dataset. It identifies the number of rows and columns, detects numerical and categorical features, and calculates statistics such as mean, median, and standard deviation. The module also detects missing values and duplicate records, providing an overview of the dataset quality before further processing.

#### C. Data Leakage Detection Module

The data leakage detection module is the core component of the proposed system. It analyzes the relationship between features and the target variable using statistical techniques such as correlation analysis. Features with extremely high correlation with the target variable are flagged as potential leakage sources. The module also detects temporal leakage by analyzing date-based features that may reveal future information.

#### D. Data Preprocessing Module

The preprocessing module performs several data cleaning operations to improve dataset quality. These operations include handling missing values, removing duplicate records, detecting and removing outliers, and applying feature scaling techniques. Encoding methods such as label encoding and one-hot encoding are also applied to

convert categorical variables into numerical form suitable for machine learning models.

#### E. Data Visualization Module

The visualization module provides graphical representations of the dataset to help users understand patterns and relationships within the data. The system generates various visualizations such as histograms, bar charts, scatter plots, box plots, and correlation heatmaps. These visualizations allow users to identify trends, outliers, and feature relationships within the dataset.

#### F. Clean Dataset Export Module

After preprocessing and leakage detection, the system generates a final cleaned dataset. The dataset is validated to ensure that no major leakage patterns remain. The cleaned dataset is then exported as a CSV file, which can be directly used for machine learning model training and analysis.

### IV. IMPLEMENTATION

The implementation of the proposed system is carried out using a combination of Python-based technologies and web development frameworks. The system is developed as a web application that allows users to upload datasets, analyze them for leakage risks, perform preprocessing operations, and visualize dataset characteristics. The implementation focuses on building a modular architecture where each component performs a specific function in the data analysis pipeline.

#### A. Development Environment

The system is implemented using **Python** as the primary programming language. The backend framework used for developing the web application is **Flask**, which provides a lightweight and flexible environment for building data-driven applications. The application is executed on a local server environment and can be accessed through a web browser. The development environment includes tools such as **Visual Studio Code** for coding and debugging, and Python libraries for data analysis and machine learning preprocessing.

#### B. Data Processing Libraries

Several Python libraries are used to handle dataset processing and analysis: These libraries provide efficient methods for analyzing dataset structures and performing statistical calculations required for leakage detection.

#### C. Data Visualization Tools

The data visualization module is implemented using Python libraries such as Matplotlib and Seaborn to create meaningful visual representations of dataset characteristics. It includes histograms for analyzing data distribution, scatter plots for understanding relationships between features, box plots for detecting outliers, and correlation heatmaps for examining dependencies among features. These visualizations enable users to gain deeper insights into data patterns and effectively identify anomalies.

#### D. Web Interface Design

The frontend of the system is designed using **HTML, CSS, and Bootstrap**. Bootstrap provides responsive components that allow the application interface to adapt to different screen sizes. The web interface provides options for dataset upload, dataset analysis, preprocessing selection, and visualization display. The user interface communicates with the backend Flask server to process user requests and display the results dynamically.

#### E. System Workflow

The system workflow begins when the user uploads a dataset through the web interface. The dataset is then processed by the profiling module to analyze its structure. The leakage detection module identifies potential leakage features using statistical analysis. After identifying leakage risks, the preprocessing module performs cleaning operations on the dataset. Finally, visualization charts are generated, and a cleaned dataset is exported for machine learning model development.

### V. EXPERIMENTAL RESULTS

The experimental evaluation of the proposed **Data Leakage Detection and Intelligent Data Preprocessing System** was conducted using publicly available datasets obtained from Kaggle. These datasets were selected because they contain common data quality issues such as missing values, duplicate records, inconsistent feature distributions, and potential data leakage patterns. The purpose of the experiments was to evaluate the effectiveness of the system in identifying leakage risks, performing preprocessing operations, and preparing datasets suitable for machine learning model development. During the experimental phase, datasets were uploaded into the developed Flask-based web application. The system first performed dataset profiling to analyze the structure of the uploaded data. It successfully identified the number of rows, number of columns, feature types, and missing value distributions. The system also generated statistical summaries such as mean, median, and

standard deviation for numerical features, providing a clear overview of dataset characteristics.

### VI. CONCLUSION

Data leakage is a significant challenge in machine learning workflows that can lead to misleading model performance and unreliable predictions in real-world applications. When information from outside the training dataset is unintentionally used during model development, it creates overly optimistic evaluation results that do not reflect actual performance after deployment. Therefore, detecting and preventing data leakage is an essential step in building reliable machine learning systems. This paper presented a **Data Leakage Detection and Intelligent Data Preprocessing System** designed to automatically identify leakage patterns and improve dataset quality before model training. The proposed system integrates dataset profiling, statistical leakage detection, preprocessing techniques, and visualization tools within a single web-based platform. By analyzing feature relationships with the target variable and detecting abnormal correlations, the system can identify potential leakage sources effectively. The system also provides intelligent preprocessing operations such as handling missing values, removing duplicates, detecting outliers, and performing feature scaling and encoding. These preprocessing techniques help improve dataset quality and ensure compatibility with machine learning algorithms. Additionally, visualization tools provide graphical insights into dataset distributions and feature relationships, allowing users to better understand the structure of the data.

### REFERENCES

- [1] J. Gama, I. Žliobaitė, A. Bifet, M. Pechenizkiy, and A. Bouchachia, "A Survey on Concept Drift Adaptation," *ACM Computing Surveys*, vol. 46, no. 4, pp. 1–37, 2014.
- [2] S. Rabanser, S. Günemann, and Z. Lipton, "Failing Loudly: An Empirical Study of Methods for Detecting Dataset Shift," *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 1396–1408, 2019.
- [3] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [4] W. McKinney, "Data Structures for Statistical Computing in Python," *Proceedings of the Python in Science Conference*, pp. 51–56, 2010.
- [5] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, 2016.
- [6] J. Brownlee, *Machine Learning Mastery with Python*, Machine Learning Mastery Publishing, 2016.

- [7] F. Chollet, Deep Learning with Python, Manning Publications, 2018.
- [8] Kaggle Datasets Repository. Available: <https://www.kaggle.com>
- [9] Python Software Foundation, “Python Language Reference,” Available: <https://www.python.org>
- [10] Matplotlib Development Team, “Matplotlib: Visualization with Python,” Available: <https://matplotlib.org>
- [11] M. Waskom, “Seaborn: Statistical Data Visualization Library,” Available: <https://seaborn.pydata.org>