

Universal ML Dashboard: A Comprehensive Machine Learning and Data Analysis Tool

Abstract

This project presents the **Universal ML Dashboard**, an interactive machine learning and data analysis application developed as the AI Lab Final Project. The dashboard is designed to work with **any CSV dataset** and focuses primarily on **regression-based machine learning tasks**. It provides an end-to-end pipeline including data upload, preprocessing, exploratory data analysis, model training, evaluation, comparison, prediction, and result export. The system is implemented using Python with Streamlit for the user interface and Scikit-learn for machine learning. The project emphasizes modular design, best practices in machine learning, and user-friendly interaction for both technical and non-technical users.

1. Introduction

Machine learning workflows often require multiple tools for data cleaning, visualization, model training, and evaluation. This fragmentation makes it difficult for beginners and students to understand the complete pipeline. The **Universal ML Dashboard** addresses this problem by providing a **single integrated platform** that allows users to perform the

entire regression-based machine learning workflow through an interactive web interface.

The project was developed as part of the **AI Lab Final Project** with the objective of demonstrating practical implementation of machine learning concepts, object-oriented programming, and data visualization techniques.

2. Problem Statement

Many existing machine learning tools are either too complex for beginners or limited to specific datasets and use cases. There is a need for a **generic, reusable, and interactive dashboard** that can:

- Work with any CSV dataset
 - Automatically handle numeric and categorical features
 - Train and compare multiple machine learning models
 - Provide clear visualizations and evaluation metrics
-

3. Objectives

The main objectives of this project are:

- To design a **universal machine learning dashboard** for regression problems
- To automate data preprocessing including cleaning and encoding

- To implement and compare multiple regression models
 - To visualize data patterns, model performance, and feature importance
 - To follow **machine learning best practices** such as proper train-test splitting and scaling
 - To allow users to export predictions and evaluation results
- Cleaning missing and invalid values
 - Automatically detecting numeric and categorical columns
 - Encoding categorical variables
 - Applying feature scaling using StandardScaler

4. System Overview and Architecture

The system follows a **modular and object-oriented architecture**, consisting of the following main components:

4.1 Streamlit Interface

Streamlit is used to build the interactive web-based user interface. The application is divided into multiple tabs:

- Data Overview
- Exploratory Data Analysis
- Machine Learning
- Model Comparison
- Export Results
- About

4.2 DataProcessor Module

Responsible for:

- Loading CSV datasets

4.3 MLModelManager Module

Handles all machine learning operations, including:

- Training multiple regression models
- Hyperparameter updates
- Cross-validation
- Performance evaluation
- Overfitting detection
- Model saving and loading

4.4 VisualizationManager Module

Manages all plots and visual outputs such as:

- Correlation matrix
- Distribution plots
- Model comparison charts
- Prediction vs actual plots
- Residual plots
- Feature importance graphs

5. Dataset Handling and Preprocessing

The dashboard accepts **any CSV file** uploaded by the user. Once uploaded:

- The dataset is cleaned automatically
- Numeric and categorical columns are detected
- Missing values are identified and reported
- Categorical features are encoded for machine learning
- Feature scaling is optionally applied

A key best practice followed is that **scaling is fitted only on training data** and then applied to test data, avoiding data leakage.

6. Exploratory Data Analysis (EDA)

The EDA module helps users understand the dataset before modeling. It includes:

- Dataset preview and column-level information
- Statistical summary of numeric features
- Summary of categorical features
- Correlation matrix for numeric variables
- Distribution plots for numeric features
- Frequency plots for categorical features

- Missing value analysis with counts and percentages

These visualizations help in identifying trends, relationships, and potential issues in the data.

7. Machine Learning Models

The dashboard supports training of multiple regression models:

- **Linear Regression**
- **Decision Tree Regressor**
- **Random Forest Regressor**
- **Gradient Boosting Regressor**

7.1 Training Process

- User selects target and feature columns
- Dataset is split into training and testing sets
- Optional feature scaling is applied
- Models are trained using the training set
- Performance is evaluated on both training and test sets

7.2 Evaluation Metrics

The following metrics are used:

- R^2 Score
- Mean Absolute Error (MAE)
- Root Mean Squared Error (RMSE)

Cross-validation (K-Fold) is optionally performed to ensure robust evaluation.

7.3 Overfitting Detection

The system automatically checks for overfitting by comparing training and testing R² scores and displays warnings when significant differences are detected.

8. Model Comparison and Analysis

All trained models are compared using:

- Performance comparison charts
- Detailed metric tables
- Cross-validation statistics

The **best model** is selected based on the highest test R² score, with RMSE used as a tie-breaker. Additional analysis includes:

- Prediction vs Actual plots
- Residual plots
- Feature importance visualization for tree-based models

9. Prediction and User Interaction

After training, users can:

- Select any trained model
- Enter feature values through an interactive form
- Generate real-time predictions

The system ensures that preprocessing steps such as encoding and scaling are consistently applied during prediction.

10. Export and Model Persistence

The dashboard allows exporting:

- Prediction results (CSV)
- Model evaluation metrics (CSV)
- Feature importance values (CSV)

Additionally, trained models can be saved to disk and loaded later using joblib, enabling model reuse without retraining.

11. Technologies Used

- **Python** (Object-Oriented Programming)
- **Streamlit** (Interactive Dashboard)
- **Scikit-learn** (Machine Learning)
- **Pandas & NumPy** (Data Processing)
- **Matplotlib** (Visualization)

12. Challenges and Limitations

- The dashboard is limited to **regression tasks only**
- Large datasets may affect performance due to in-memory processing

- Deep learning models are not included

Despite these limitations, the system is highly flexible for academic and learning purposes.

13. Conclusion and Future

Enhancements

The Universal ML Dashboard successfully demonstrates an end-to-end machine learning pipeline in an interactive and user-friendly manner. It fulfills all objectives of the AI Lab Final Project and showcases practical implementation of machine learning concepts.

Future Enhancements:

- Support for classification problems
 - Integration of deep learning models
 - Automated feature selection
 - Cloud deployment
 - Report generation directly from the dashboard
-