# Machine Learning Assignment 2

## Classification Models

| | |
|---|---|
| **Student Name:** | Aravindan B |
| **Student ID:** | 2025aa05026 |
| **Program:** | M.Tech (AIML), BITS Pilani WILP |
| **Course:** | Machine Learning |
| **Submission Date:** | February 2026 |
| **Email:** | 2025aa05026@wilp.bits-pilani.ac.in |

# 1. GitHub Repository Link

## Repository URL:

```
https://github.com/algoyog/mlassignment2
```

## Repository Contents:

- **app.py**   -- Streamlit web application
- **requirements.txt**   -- Python dependencies
- **README.md**   -- Complete documentation
- **.gitignore**   -- Git ignore rules
- **model/model_training.ipynb**   -- Jupyter notebook for model training
- **input/adult_income.csv**   -- UCI Adult Income dataset
- **output/**   -- Generated results and charts
- **utils/ml_utils.py**   -- Shared ML utilities

# 2. Live Streamlit App Link

## Application URL:

```
https://mlassignment2-kmepmbozytsyeyf3yerkyb.streamlit.app/
```

## Implemented Features:

- CSV file upload with dataset preview
- Target column selection with class distribution
- Model selection dropdown (6 classification models)
- All 6 evaluation metrics displayed (Accuracy, AUC, Precision, Recall, F1, MCC)
- Confusion matrix visualization
- Classification report
- Model comparison table with best model highlighted
- Results download as CSV
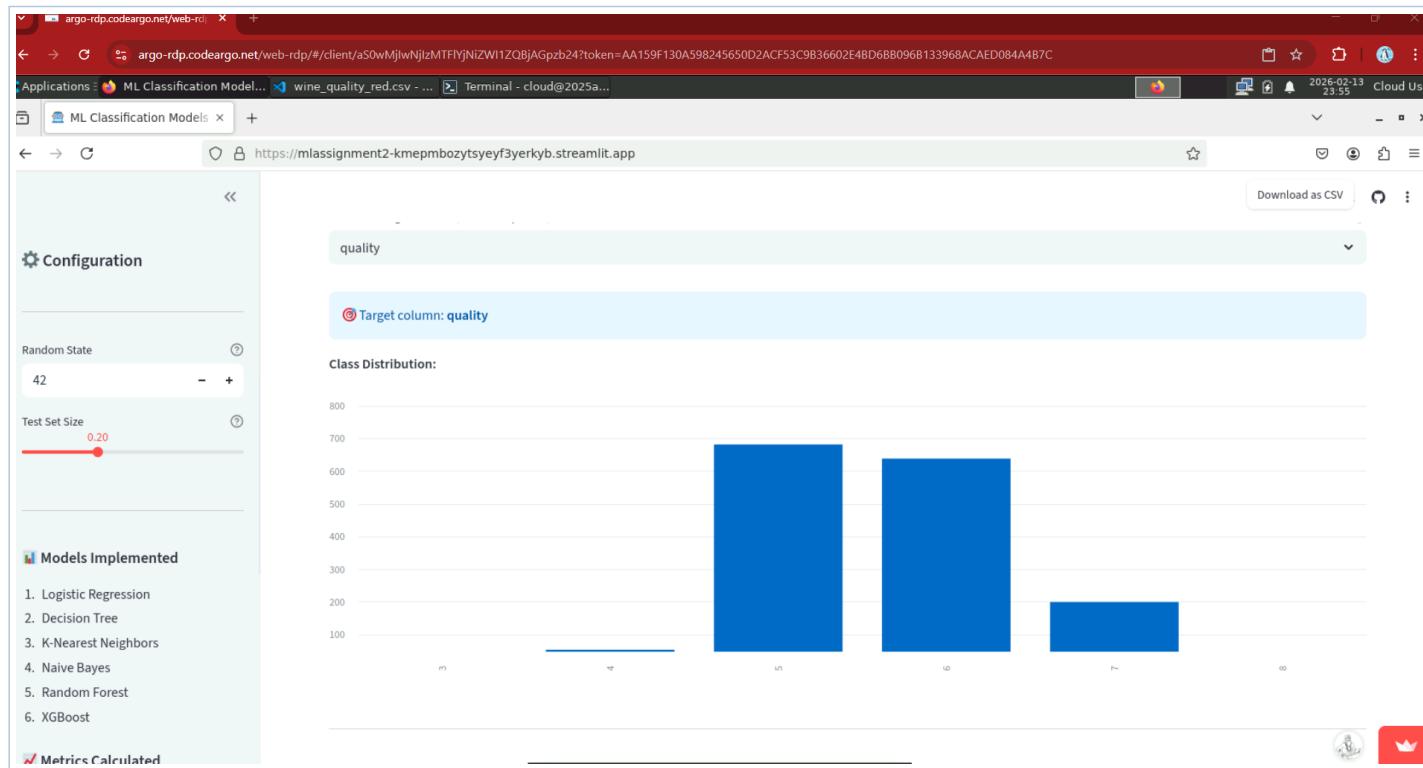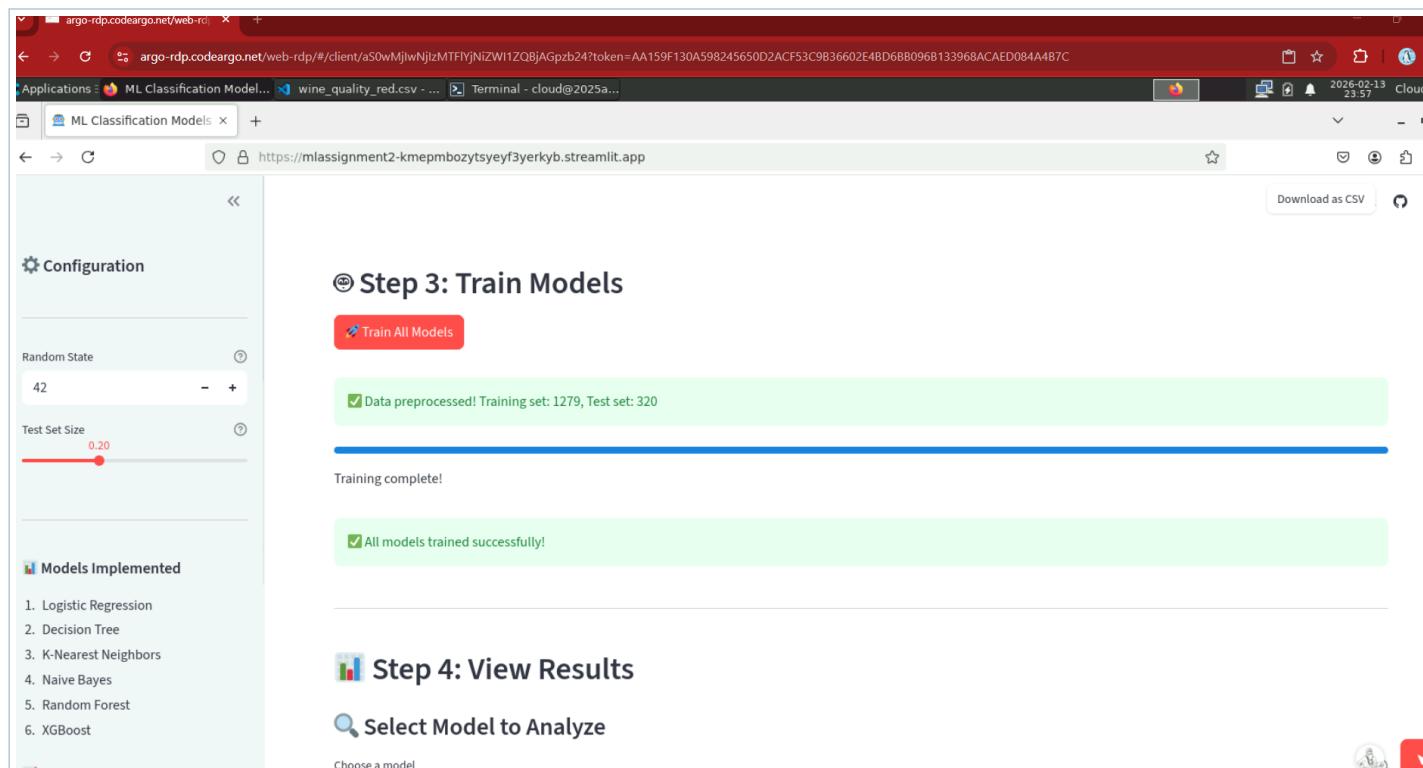
# 3. BITS Virtual Lab Screenshots

The following screenshots demonstrate the Streamlit application running on the BITS Virtual Lab (argo-rdp.codeargo.net) using the Wine Quality Red dataset (1599 instances, 12 features) as a test dataset.

**Screenshot 1: Step 1: CSV Upload - wine_quality_red.csv loaded successfully (1599 rows, 12 features)**



**Screenshot 2: Step 2: Dataset Preview - Feature table and dataset information displayed**

**Screenshot 3: Step 2: Target Column Selection - Class distribution visualization for "quality"**



**Screenshot 4: Step 3: Model Training - All 6 models trained successfully on BITS Virtual Lab**

**Screenshot 5: Step 4: Model Comparison Table - All 6 models with all 6 evaluation metrics**

## Model Comparison Table

| | Accuracy | AUC | Precision | Recall | F1 | MCC |
|---|---|---|---|---|---|---|
| Logistic Regression | 0.590600 | 0.755500 | 0.569500 | 0.590600 | 0.567300 | 0.325000 |
| Decision Tree | 0.593800 | 0.708000 | 0.590800 | 0.593800 | 0.592100 | 0.363900 |
| K-Nearest Neighbors | 0.609400 | 0.747600 | 0.584100 | 0.609400 | 0.595900 | 0.373300 |
| Naive Bayes | 0.562500 | 0.737700 | 0.574500 | 0.562500 | 0.568100 | 0.329900 |
| Random Forest | 0.662500 | 0.833800 | 0.637700 | 0.662500 | 0.646200 | 0.454700 |
| XGBoost | 0.678100 | 0.817100 | 0.665700 | 0.678100 | 0.668700 | 0.486700 |

🏆 **Best Model:** XGBoost (Accuracy: 0.6781)

## 💾 Download Results

⬇ Download Results as CSV

### Configuration

**Random State** ⑦

42 − +

**Test Set Size** ⑦
0.20

### 📊 Models Implemented

1. Logistic Regression
2. Decision Tree
3. K-Nearest Neighbors
4. Naive Bayes
5. Random Forest
6. XGBoost

# 4. README Content

## Problem Statement

This project implements a comprehensive machine learning classification pipeline featuring six different classification algorithms applied to the UCI Adult Income dataset. The goal is to predict whether an individual's annual income exceeds $50K based on census attributes, and to compare the performance of traditional ML models (Logistic Regression, Decision Tree, K-Nearest Neighbors, Naive Bayes) against ensemble methods (Random Forest and XGBoost) using six evaluation metrics.

## Dataset Description

**Dataset Name:**  UCI Adult Income Dataset (Census Income)

**Source:**  UCI Machine Learning Repository

**Type:**  Binary Classification

**Instances:**  30,162 (Requirement >= 500: MET)

**Features:**  14 (Requirement >= 12: MET)

**Target:**  income (<=50K or >50K)

Class Distribution: Class 0 (<=50K): 22,654 samples (75.1%) | Class 1 (>50K): 7,508 samples (24.9%)

Numerical Features (6): age, fnlwgt, education_num, capital_gain, capital_loss, hours_per_week

Categorical Features (8): workclass, education, marital_status, occupation, relationship, race, sex, native_country

## Data Preprocessing Steps

- Label Encoding on all 8 categorical features
- Target encoded: <=50K to 0, >50K to 1
- StandardScaler normalization on all features
- 80-20 stratified train-test split

## Model Performance Comparison Table

| ML Model Name | Accuracy | AUC | Precision | Recall | F1 | MCC |
|---|---|---|---|---|---|---|
| Logistic Regression | 0.8175 | 0.8501 | 0.8060 | 0.8175 | 0.8018 | 0.4613 |
| Decision Tree | 0.8508 | 0.8855 | 0.8446 | 0.8508 | 0.8451 | 0.5789 |
| K-Nearest Neighbors | 0.8190 | 0.8498 | 0.8133 | 0.8190 | 0.8154 | 0.4993 |
| Naive Bayes | 0.7978 | 0.8498 | 0.7830 | 0.7978 | 0.7697 | 0.3798 |
| Random Forest | 0.8589 | 0.9136 | 0.8534 | 0.8589 | 0.8526 | 0.6003 |
| **XGBoost** | **0.8671** | **0.9243** | **0.8624** | **0.8671** | **0.8624** | **0.6269** |

**Best Model: XGBoost (Accuracy = 0.8671, AUC = 0.9243)**

## Model Performance Observations

| ML Model Name | Observation about Model Performance |
|---|---|
| **Logistic Regression** | Achieves 81.75% accuracy with AUC 0.8501, serving as a solid linear baseline. It performs well because features like education_num, age, and hours_per_week have approximately linear relationships with income. The model converges reliably with lbfgs solver and is the most interpretable among the six, though it falls short on non-linear patterns compared to tree-based models. |
| **Decision Tree** | Achieves 85.08% accuracy with AUC 0.8855 and MCC 0.5789. With max_depth=10, it successfully captures non-linear interactions between occupation, marital_status, and education without extreme overfitting. It outperforms linear models significantly while remaining interpretable through decision rules. However, it is surpassed by ensemble methods which reduce its variance. |
| **K-Nearest Neighbors** | KNN with k=5 achieves 81.90% accuracy with AUC 0.8498, slightly edging out Logistic Regression. Feature scaling via StandardScaler is essential for KNN distance calculations and was correctly applied. On 30,162 instances, it is computationally heavier at prediction time but benefits from the demographic clustering present in the Adult Income dataset. |
| **Naive Bayes** | Gaussian Naive Bayes achieves the lowest accuracy at 79.78% with MCC 0.3798. The feature independence assumption does not hold well here as features like education, occupation, and marital_status are correlated, limiting its performance. Despite this, it achieves a competitive AUC of 0.8498, indicating reasonable probability calibration, and is the fastest model to train. |
| **Random Forest** | Achieves 85.89% accuracy with AUC 0.9136 and MCC 0.6003, ranking second overall. The ensemble of 100 trees reduces overfitting through bagging and random feature subsets. Its AUC of 0.9136 reflects excellent discriminative ability. It consistently outperforms all traditional models across every metric and provides feature importance insights for interpretability. |
| **XGBoost** | Best performing model with 86.71% accuracy, highest AUC (0.9243), F1 (0.8624), and MCC (0.6269). Its gradient boosting framework iteratively corrects errors from previous trees, capturing complex non-linear patterns. With learning_rate=0.1, max_depth=6, and 100 estimators, it demonstrates the clear advantage of advanced ensemble techniques over traditional classifiers on this tabular dataset. |

## Overall Insights

- Best performing model: XGBoost with 86.71% accuracy, AUC 0.9243, F1 0.8624, MCC 0.6269
- Ensemble methods (Random Forest and XGBoost) outperform all traditional algorithms across every metric
- Naive Bayes is weakest (79.78%) due to its feature independence assumption not holding for this correlated dataset
- All models achieved AUC > 0.84, indicating good discriminative ability despite class imbalance (75%/25%)

## END OF SUBMISSION