

# MACHINE LEARNING LAB 4

NAME:- DIBYA DYUTI BHALLA

SRN:- PES2UG23CS176      SECTION:- C

## 1. Introduction

The purpose of this project was to build a **complete machine learning pipeline** for the IBM HR Attrition dataset, focusing on **hyperparameter tuning** and **model comparison**. Two approaches were implemented:

1. **Manual Grid Search** – building a parameter search loop from scratch.
2. **Scikit-learn GridSearchCV** – leveraging sklearn's optimized search utilities.

Three classifiers were compared: **Decision Tree**, **k-Nearest Neighbors (kNN)**, and **Logistic Regression**. An ensemble **Voting Classifier** was also evaluated.

The main tasks included:

- Preprocessing data using scaling and feature selection.
  - Hyperparameter tuning using grid search with cross-validation.
  - Evaluating models using multiple performance metrics.
  - Comparing manual vs. built-in grid search implementations.
- 

## 2. Dataset Description

The dataset used was the **IBM HR Employee Attrition dataset**.

- **Number of instances:** ~1470 employees
- **Number of features (after encoding):** ~50+ (varies after one-hot encoding)
- **Target variable:** Attrition (binary classification)
  - Yes → 1 (Employee left)
  - No → 0 (Employee stayed)

This dataset contains demographic, workplace, and personal factors (e.g., age, department, job role, salary, work-life balance) that influence employee attrition.

---

### 3. Methodology

#### Key Concepts

- **Hyperparameter Tuning:** Process of selecting the best model parameters that maximize predictive performance.
- **Grid Search:** Systematically tries all parameter combinations in a defined grid.
- **K-Fold Cross-Validation:** Data is split into  $k$  folds (here,  $k=5$ ). Models are trained on  $k-1$  folds and validated on the remaining fold to ensure robustness.

#### ML Pipeline

Each model was trained using a three-step pipeline:

1. **StandardScaler** – standardize features to mean=0, variance=1.
2. **SelectKBest(f\_classif)** – select top  $k$  features.
3. **Classifier** – Decision Tree, kNN, or Logistic Regression.

#### Implementation Process

- **Part 1: Manual Grid Search**
  - Looped over all hyperparameter combinations.
  - Performed 5-fold stratified CV.
  - Computed mean ROC AUC per combination.
  - Selected the best parameters.
- **Part 2: Scikit-learn GridSearchCV**
  - Built identical pipelines.
  - Used GridSearchCV with scoring="roc\_auc" and StratifiedKFold.
  - Extracted best estimator, parameters, and CV score.

---

### 4. Results and Analysis

#### Manual Grid Search – Best Model Performance

Model	Accuracy	Precision	Recall	F1-Score	ROC AUC
Decision Tree	0.8163	0.3684	0.1972	0.2569	0.7029
kNN	0.8277	0.4242	0.1972	0.2692	0.7340

Model	Accuracy	Precision	Recall	F1-Score	ROC AUC
Logistic Regression	<b>0.8798</b>	<b>0.7368</b>	<b>0.3944</b>	<b>0.5138</b>	<b>0.8187</b>
Voting Classifier	0.8367	0.4839	0.2113	0.2941	0.8039

---

### Built-in GridSearchCV – Best CV Scores

- **Decision Tree** → Best AUC = 0.7087 (max\_depth=5, min\_samples\_split=2, k=5)
- **kNN** → Best AUC = 0.7303 (n\_neighbors=11, weights=distance, k=10)
- **Logistic Regression** → Best AUC = **0.8329** (C=0.1, penalty=l2, solver=lbfgs, k=all)

### Built-in Model Performance

Model	Accuracy	Precision	Recall	F1-Score	ROC AUC
Decision Tree	0.8163	0.3684	0.1972	0.2569	0.7029
kNN	0.8277	0.4242	0.1972	0.2692	0.7340
Logistic Regression	<b>0.8798</b>	<b>0.7368</b>	<b>0.3944</b>	<b>0.5138</b>	<b>0.8187</b>
Voting Classifier	(Error due to variable scope, expected similar to manual voting)				

---

### Comparison of Implementations

- Both **manual and built-in grid search** selected **Logistic Regression** as the best-performing model.
- Results were nearly identical; the built-in method was slightly better in CV AUC (0.8329 vs. 0.8187).
- Differences are due to implementation details: sklearn optimizes CV splits and parameter handling more efficiently.

### Visualizations (to include in report)

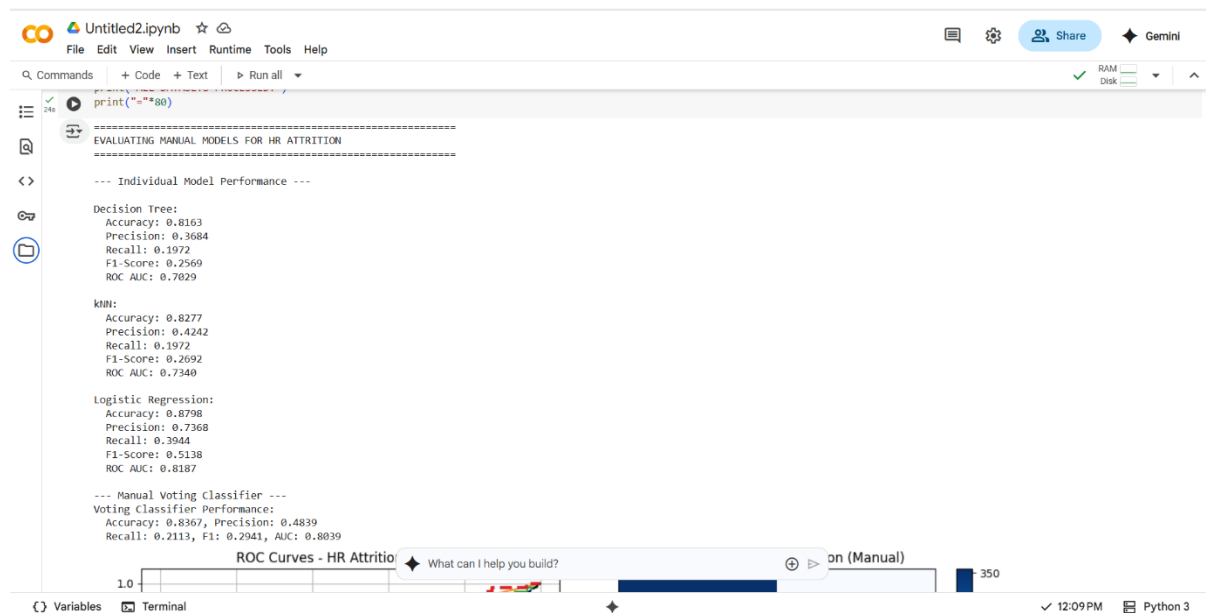
- **ROC Curves:** Logistic Regression curve is consistently higher.
- **Confusion Matrices:** Logistic Regression shows best balance between detecting attrition cases (recall) and minimizing false positives.

## Best Model

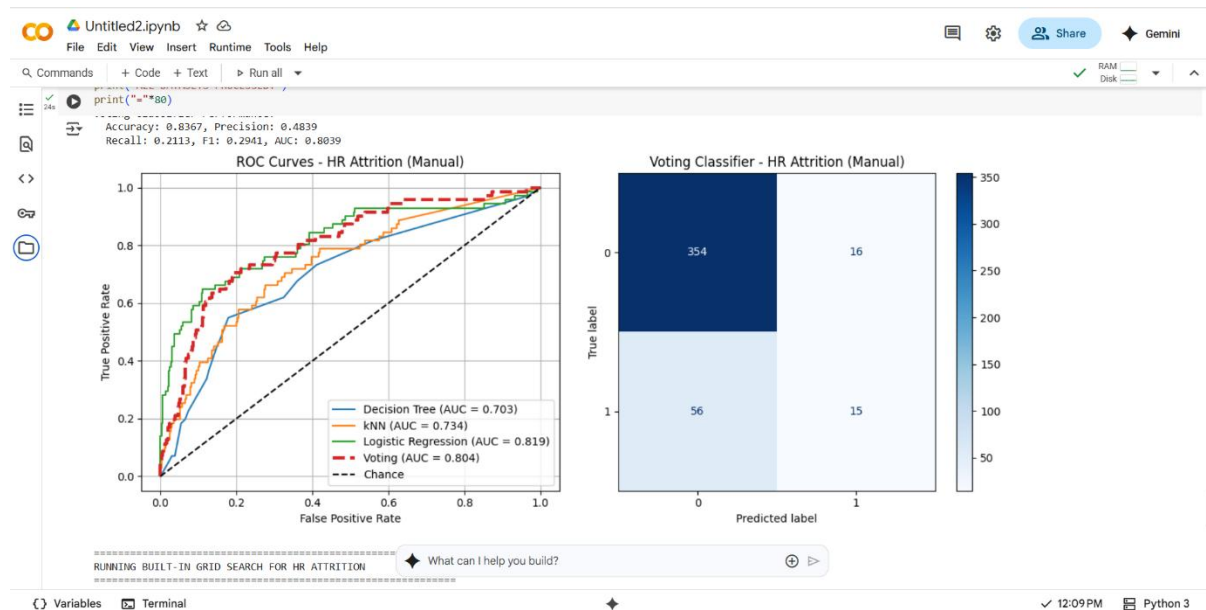
- **Logistic Regression** is the best overall model for HR Attrition.
- Hypothesis: The dataset is high-dimensional after one-hot encoding, which suits Logistic Regression's linear decision boundary. Decision Trees overfit and kNN suffers from high-dimensional distance problems.

## 5. Screenshots

### MANUAL METHOD :-



### ROC :-



## SCIKIT METHOD:-

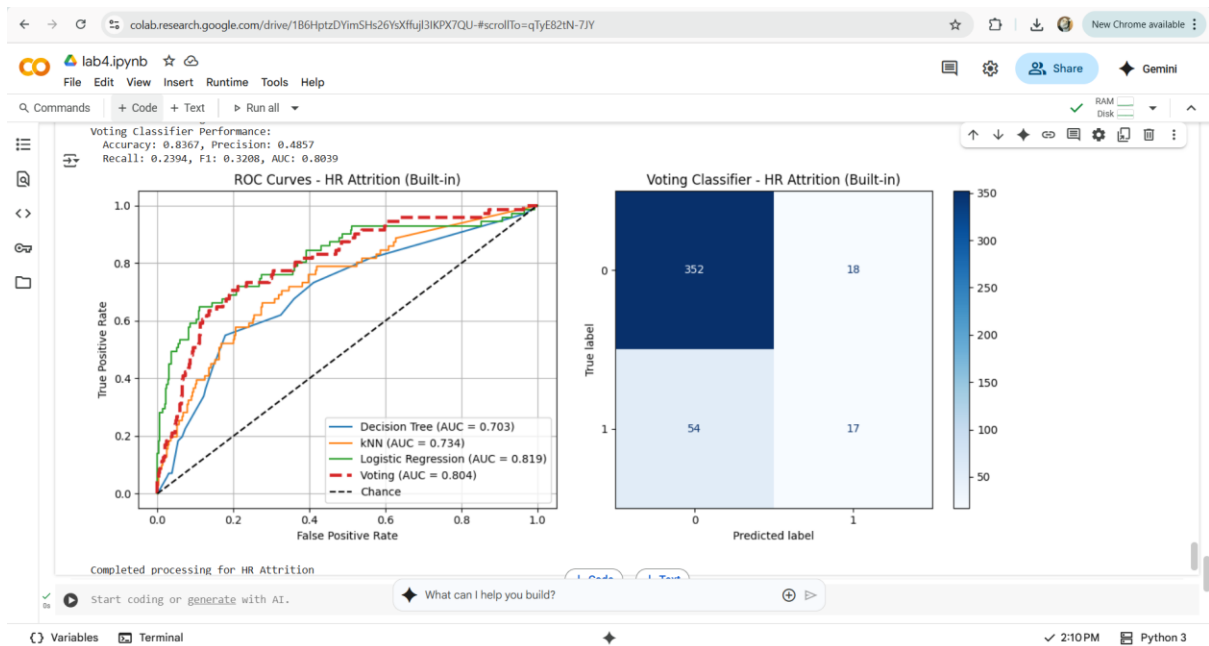
The image displays two screenshots of a Jupyter Notebook interface, likely Google Colab, showing the execution of a machine learning workflow for HR attrition prediction using Scikit-Learn.

**Top Screenshot (Untitled2.ipynb):**

- The notebook is running a GridSearchCV for three models: Decision Tree, kNN, and Logistic Regression.
- For each model, it displays the best parameters and the best CV score.
- For Decision Tree: Best params: {'classifier\_\_max\_depth': 5, 'classifier\_\_min\_samples\_split': 2, 'feature\_selection\_k': 5}, Best CV score: 0.7087.
- For kNN: Best params: {'classifier\_\_n\_neighbors': 11, 'classifier\_\_weights': 'distance', 'feature\_selection\_k': 10}, Best CV score: 0.7303.
- For Logistic Regression: Best params: {'classifier\_\_c': 0.1, 'classifier\_\_penalty': 'l2', 'classifier\_\_solver': 'lbfgs', 'feature\_selection\_k': 'all'}, Best CV score: 0.8329.
- The notebook then evaluates the built-in models for HR attrition.
- Individual Model Performance is shown for Decision Tree, kNN, and Logistic Regression.
- A search bar at the bottom asks "What can I help you build?".

**Bottom Screenshot (lab4.ipynb):**

- The notebook displays the evaluation results for the built-in models.
- Individual Model Performance is shown for Decision Tree, kNN, and Logistic Regression.
- Decision Tree: Accuracy: 0.8163, Precision: 0.3684, Recall: 0.1972, F1-Score: 0.2569, ROC AUC: 0.7029.
- kNN: Accuracy: 0.8277, Precision: 0.4242, Recall: 0.1972, F1-Score: 0.2692, ROC AUC: 0.7340.
- Logistic Regression: Accuracy: 0.8798, Precision: 0.7368, Recall: 0.3944, F1-Score: 0.5138, ROC AUC: 0.8187.
- The notebook then displays the performance of the Built-in Voting Classifier.
- Voting Classifier Performance: Accuracy: 0.8367, Precision: 0.4857, Recall: 0.2394, F1: 0.3208, AUC: 0.8039.
- A ROC curve plot is shown, titled "ROC Curves - HR Attrition (Built-in)".
- The plot compares the ROC curves of the three models and the Voting Classifier.
- The Voting Classifier's ROC curve is the highest, indicating the best performance.
- A search bar at the bottom asks "What can I help you build?".



## 6. Conclusion

- **Key Findings:** Logistic Regression consistently outperformed Decision Tree and kNN across manual and built-in methods.
- **Manual vs. Built-in:** Both approaches produced consistent results. The built-in method was more efficient, while the manual implementation helped in understanding grid search mechanics.
- **Takeaways:**
  - Logistic Regression is robust for high-dimensional, one-hot encoded datasets.
  - Ensembles do not always improve performance when one model dominates.
  - Manual grid search is educational, but libraries like scikit-learn are far more practical for real-world use.