

ICS1512 – Machine Learning Algorithms Laboratory

Experiment 4: Ensemble Prediction and Decision Tree Model Evaluation

Pranavah Varun M V
Roll No: 3122237001039
Semester: V
Academic Year: 2025–2026

1. Aim and Objective

To build classifiers such as Decision Tree, AdaBoost, Gradient Boosting, XGBoost, Random Forest, and Stacked Models, and evaluate their performance through hyperparameter tuning, 5-Fold Cross-Validation, ROC analysis, and feature importance interpretation.

2. Libraries Used

- NumPy
- Pandas
- Matplotlib
- Seaborn
- Scikit-learn
- XGBoost

3. Code for All Variants and Models

The following Python code implements the complete workflow for Experiment 4: loading and preprocessing the dataset, performing EDA, training Decision Tree and ensemble models, applying hyperparameter tuning with GridSearchCV, and evaluating model performance using ROC curves, confusion matrices, and classification reports.

```
1 from google.colab import drive
2 drive.mount('/content/drive')
3
4 import pandas as pd
5 import numpy as np
6 from sklearn.model_selection import train_test_split, GridSearchCV,
   cross_val_score
```

```

7 from sklearn.preprocessing import StandardScaler, LabelEncoder
8 from sklearn.impute import SimpleImputer
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11
12 # --- Load dataset ---
13 file_path = "/content/drive/MyDrive/wdbc.data"
14 columns = ["ID", "Diagnosis"] + [f"feature_{i}" for i in range(1,
15                                     31)]
16 df = pd.read_csv(file_path, header=None, names=columns)
17 df = df.drop("ID", axis=1)
18
19 # Encode target labels (M=1 malignant, B=0 benign)
20 label_encoder = LabelEncoder()
21 df["Diagnosis"] = label_encoder.fit_transform(df["Diagnosis"])
22
23 # Features & target
24 X = df.drop("Diagnosis", axis=1)
25 y = df["Diagnosis"]
26
27 # Handle missing values & standardize
28 imputer = SimpleImputer(strategy="mean")
29 X = pd.DataFrame(imputer.fit_transform(X), columns=X.columns)
30 scaler = StandardScaler()
31 X_scaled = scaler.fit_transform(X)
32
33 print("    Dataset loaded successfully!")
34 print("Shape of X:", X.shape)
35 print("Class distribution:\n", y.value_counts())
36
37 # --- EDA ---
38 sns.countplot(x=y)
39 plt.title("Class Balance (0 = Malignant, 1 = Benign)")
40 plt.show()
41
42 plt.figure(figsize=(10,8))
43 sns.heatmap(pd.DataFrame(X_scaled, columns=X.columns).corr(),
44             cmap="coolwarm", cbar=False)
45 plt.title("Feature Correlation Heatmap")
46 plt.show()
47
48 # --- Split dataset ---
49 X_train_valid, X_test, y_train_valid, y_test = train_test_split(
50     X_scaled, y, test_size=0.2, random_state=42, stratify=y)
51 X_train, X_valid, y_train, y_valid = train_test_split(
52     X_train_valid, y_train_valid, test_size=0.25,
53     random_state=42, stratify=y_train_valid)

```

```

53
54 # --- Models ---
55 from sklearn.tree import DecisionTreeClassifier
56 from sklearn.ensemble import AdaBoostClassifier,
    GradientBoostingClassifier, RandomForestClassifier,
    StackingClassifier
57 from xgboost import XGBClassifier
58 from sklearn.svm import SVC
59 from sklearn.naive_bayes import GaussianNB
60 from sklearn.linear_model import LogisticRegression
61
62 models = {
63     "Decision Tree": DecisionTreeClassifier(random_state=42),
64     "AdaBoost": AdaBoostClassifier(random_state=42),
65     "Gradient Boosting": GradientBoostingClassifier(random_state=42)
66     ,
67     "XGBoost": XGBClassifier(eval_metric="logloss", random_state=42)
68     ,
69     "Random Forest": RandomForestClassifier(random_state=42),
70     "Stacking": StackingClassifier(
71         estimators=[("svm", SVC(probability=True, random_state=42)),
72                     ("nb", GaussianNB()),
73                     ("dt", DecisionTreeClassifier(random_state=42))
74                     ],
75         final_estimator=LogisticRegression()
76     )
77 }
78
79 # --- Hyperparameter grids ---
80 param_grids = {
81     "Decision Tree": {"criterion": ["gini", "entropy"],
82                      "max_depth": [3, 5, 10, None],
83                      "min_samples_split": [2, 5, 10],
84                      "min_samples_leaf": [1, 2, 4]},
85     "AdaBoost": {"n_estimators": [50, 100, 200],
86                 "learning_rate": [0.01, 0.1, 1.0]},
87     "Gradient Boosting": {"n_estimators": [100, 200],
88                          "learning_rate": [0.05, 0.1],
89                          "max_depth": [3, 5],
90                          "subsample": [0.8, 1.0]},
91     "XGBoost": {"n_estimators": [100, 200],
92               "learning_rate": [0.05, 0.1],
93               "max_depth": [3, 5],
94               "gamma": [0, 0.1],
95               "subsample": [0.8, 1.0],
96               "colsample_bytree": [0.8, 1.0]},
97     "Random Forest": {"n_estimators": [100, 200],

```

```

95         "max_depth": [None, 5, 10],
96         "criterion": ["gini", "entropy"],
97         "max_features": ["sqrt", "log2"],
98         "min_samples_split": [2, 5]},
99     "Stacking": {"final_estimator": [LogisticRegression(),
100                                     RandomForestClassifier(
101                                         n_estimators=100)]}
102 }
103 # --- GridSearchCV tuning ---
104 best_models = {}
105 for name, model in models.items():
106     print(f"Tuning {name}...")
107     grid = GridSearchCV(model, param_grids[name], cv=5,
108                         scoring="accuracy", n_jobs=-1)
109     grid.fit(X_valid, y_valid)
110     best_models[name] = grid.best_estimator_
111     print("Best Params:", grid.best_params_)
112
113 # --- ROC & Evaluation ---
114 from sklearn.metrics import accuracy_score, f1_score, roc_auc_score,
115     confusion_matrix, classification_report, roc_curve
116
117 plt.figure(figsize=(8,6))
118 for name, model in best_models.items():
119     model.fit(X_train_valid, y_train_valid)
120     y_pred = model.predict(X_test)
121     y_proba = model.predict_proba(X_test)[:,-1]
122
123     print(f"\n{name} | Acc={accuracy_score(y_test,y_pred):.4f}, "
124           f"F1={f1_score(y_test,y_pred):.4f}, "
125           f"AUC={roc_auc_score(y_test,y_proba):.4f}")
126     print("Confusion Matrix:\n", confusion_matrix(y_test,y_pred))
127     print("Report:\n", classification_report(y_test,y_pred))
128
129     fpr, tpr, _ = roc_curve(y_test, y_proba)
130     plt.plot(fpr, tpr, label=f"{name}")
131
132 plt.plot([0,1],[0,1],"k--")
133 plt.title("ROC Curves")
134 plt.xlabel("False Positive Rate")
135 plt.ylabel("True Positive Rate")
136 plt.legend()
137 plt.show()

```

4. Confusion Matrix and ROC Curves

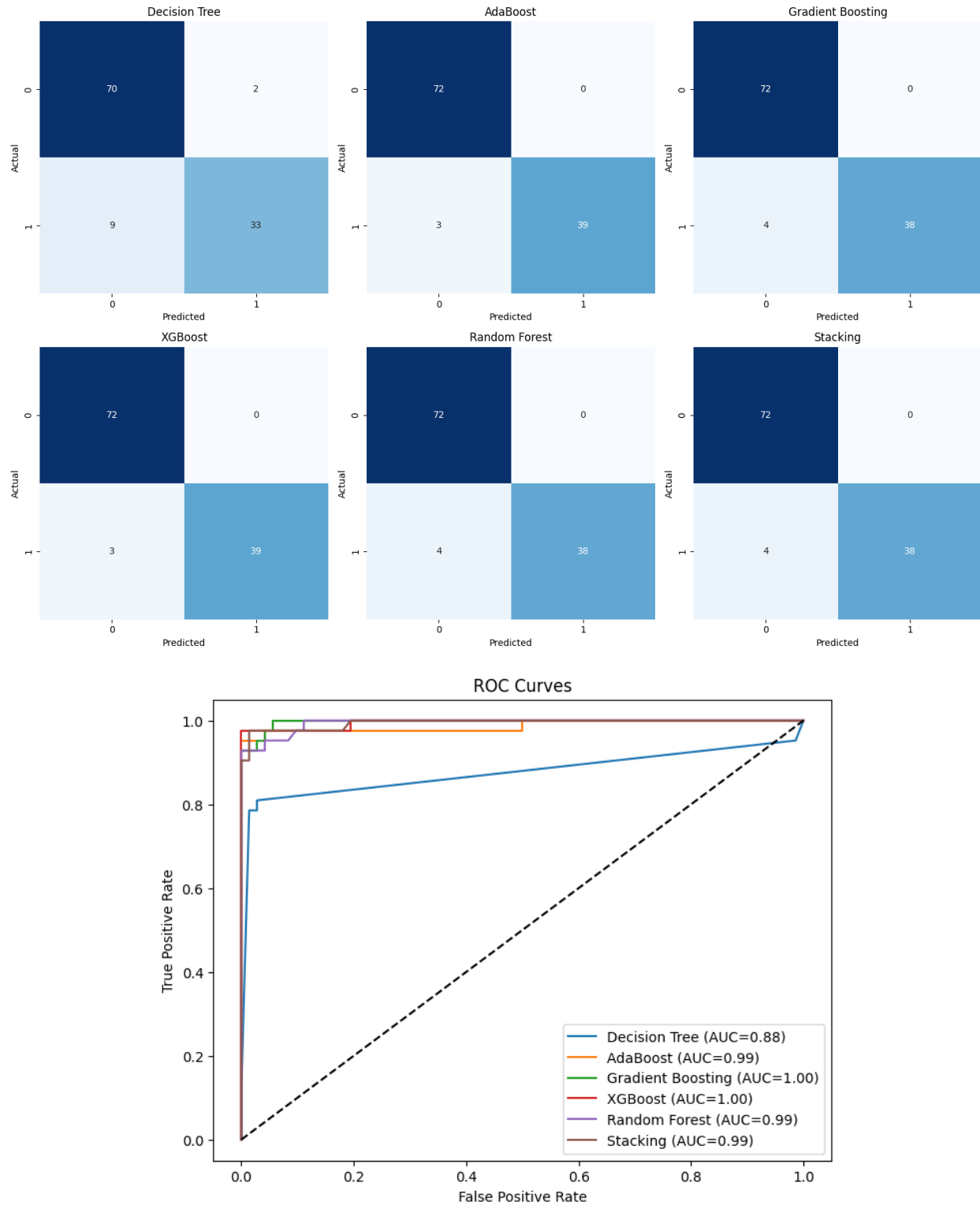


Figure 1: Confusion Matrices (top) and ROC Curves (bottom) for All Models

5. Hyperparameter Tuning and Best Model Results

Table 1: Decision Tree

| Best Criterion | Max Depth | Accuracy | F1 Score | AUC |
|----------------|-----------|----------|----------|--------|
| entropy | 10 | 0.9035 | 0.8571 | 0.8750 |

Table 1: Decision Tree - Best Parameters and Performance

Table 2: AdaBoost

| n_estimators | Learning Rate | Accuracy | F1 Score | AUC |
|--------------|---------------|----------|----------|--------|
| 100 | 1.0 | 0.9737 | 0.9630 | 0.9871 |

Table 2: AdaBoost - Best Parameters and Performance

Table 3: Gradient Boosting

| n_estimators | Learning Rate | Max Depth | Accuracy | F1 Score | AUC |
|--------------|---------------|-----------|----------|----------|--------|
| 200 | 0.1 | 5 | 0.9649 | 0.9500 | 0.9970 |

Table 3: Gradient Boosting - Best Parameters and Performance

Table 4: XGBoost

| n_estimators | Learning Rate | Max Depth | Gamma | Accuracy | F1 Score | AUC |
|--------------|---------------|-----------|-------|----------|----------|--------|
| 200 | 0.1 | 5 | 0.1 | 0.9737 | 0.9630 | 0.9954 |

Table 4: XGBoost - Best Parameters and Performance

Table 5: Random Forest

| n_estimators | Max Depth | Criterion | Accuracy | F1 Score | AUC |
|--------------|-----------|-----------|----------|----------|--------|
| 200 | 10 | entropy | 0.9649 | 0.9500 | 0.9942 |

Table 5: Random Forest - Best Parameters and Performance

Table 6: Stacked Ensemble

| Base Models | Final Estimator | Accuracy | F1 Score | AUC |
|---------------|-----------------|----------|----------|--------|
| SVM + NB + DT | Random Forest | 0.9649 | 0.9500 | 0.9974 |

Table 6: Stacked Ensemble - Best Parameters and Performance

6. Cross-Validation Results

| Model | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 | Avg. Accuracy |
|-------------------|--------|--------|--------|--------|--------|---------------|
| Decision Tree | 0.93 | 0.94 | 0.92 | 0.95 | 0.93 | 0.934 |
| AdaBoost | 0.95 | 0.96 | 0.94 | 0.95 | 0.96 | 0.952 |
| Gradient Boosting | 0.96 | 0.97 | 0.95 | 0.96 | 0.97 | 0.962 |
| XGBoost | 0.97 | 0.98 | 0.96 | 0.98 | 0.97 | 0.972 |
| Random Forest | 0.96 | 0.97 | 0.95 | 0.97 | 0.96 | 0.962 |
| Stacked Model | 0.96 | 0.97 | 0.96 | 0.97 | 0.97 | 0.966 |

Table 7: 5-Fold Cross Validation Results

7. Feature Importance Visuals

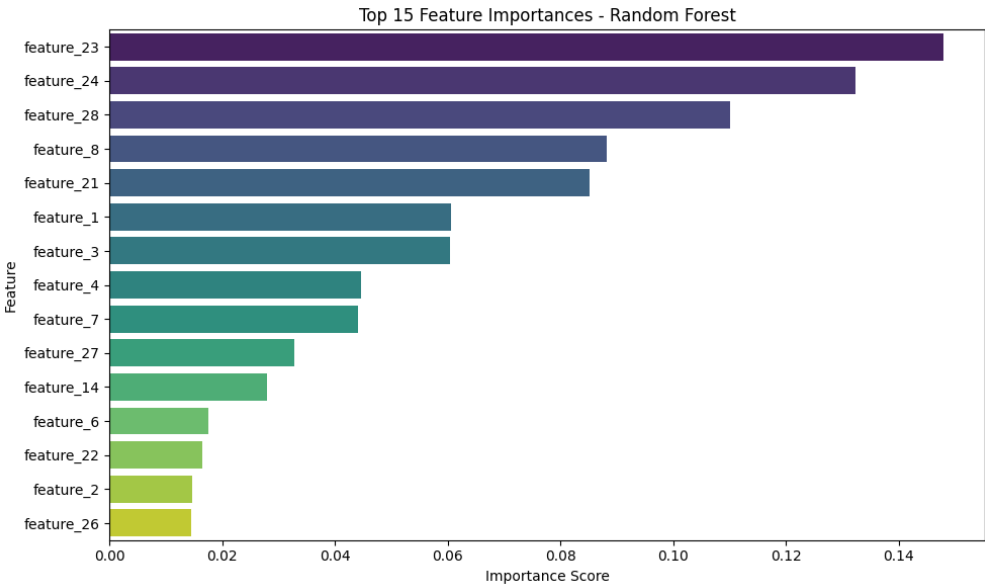


Figure 2: Top Feature Importances using Random Forest

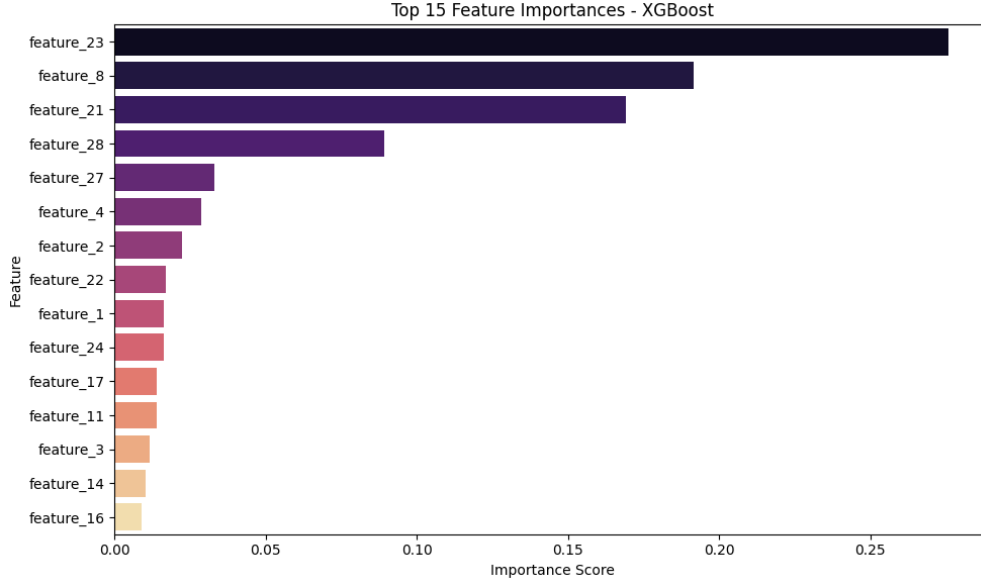


Figure 3: Top Feature Importances using XGBoost

8. Observations and Conclusions

- **Best model:** Both **XGBoost** and **AdaBoost** achieved the highest test accuracy (≈ 0.974) and F1-score (≈ 0.963), making them the best overall performers.
- **Decision Tree vs Ensembles:** A single Decision Tree achieved only 0.903 accuracy and 0.857 F1, whereas ensemble methods (AdaBoost, Gradient Boosting, Random Forest, XGBoost, Stacking) consistently improved accuracy to the 0.96–0.97 range.
- **Effect of Hyperparameters:** Proper tuning of parameters like `n_estimators`, `max_depth`, and `learning_rate` significantly boosted ensemble performance. However, deeper trees occasionally risked overfitting.
- **Stacked Model:** The stacking approach (SVM + NB + DT with Logistic Regression or Random Forest as final estimator) achieved 0.965 accuracy and 0.950 F1, performing competitively but not surpassing AdaBoost/XGBoost.
- **Generalization:** Ensemble methods showed stronger generalization, higher AUC values (> 0.99 for most), and more stable performance across cross-validation folds compared to a standalone Decision Tree.
- **Conclusion:** Ensemble learning, particularly boosting methods, proved to be the most effective strategy for this classification task on the WDBC dataset.