

# Experiment 3: Ensemble Prediction and Decision Tree Model Evaluation

Your Name:- Akshith Viswanathan  
M.Tech (Integrated) CSE, Semester V

September 12, 2025

## Aim and Objective

To build classifiers such as Decision Tree, AdaBoost, Gradient Boosting, XGBoost, Random Forest, and Stacked Models (using SVM, Naïve Bayes, Decision Tree) and evaluate their performance through 5-Fold Cross-Validation and hyperparameter tuning using the Wisconsin Diagnostic Dataset.

## Libraries Used

- `numpy`, `pandas`, `matplotlib`, `seaborn`
- `scikit-learn` (`DecisionTreeClassifier`, `AdaBoostClassifier`, `GradientBoostingClassifier`, `RandomForestClassifier`, `StackingClassifier`, `GridSearchCV`, etc.)
- `xgboost`

## Dataset Overview

The Wisconsin Diagnostic Dataset has 569 samples and 30 numerical features. Class distribution:

Malignant: 212

Benign: 357

## Code for All Models

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import time
6
7 from sklearn.datasets import load_breast_cancer
```

```

8  from sklearn.model_selection import train_test_split, GridSearchCV,
    ↪ cross_val_score, KFold
9  from sklearn.preprocessing import StandardScaler
10 from sklearn.metrics import accuracy_score, f1_score, confusion_matrix,
    ↪ roc_curve, auc
11
12 # Models
13 from sklearn.tree import DecisionTreeClassifier
14 from sklearn.ensemble import AdaBoostClassifier, GradientBoostingClassifier,
    ↪ RandomForestClassifier, StackingClassifier
15 from sklearn.naive_bayes import GaussianNB
16 from sklearn.svm import SVC
17 from sklearn.linear_model import LogisticRegression
18 import xgboost as xgb
19
20 # 1. Load and Preprocess Dataset
21 data = load_breast_cancer()
22 X = pd.DataFrame(data.data, columns=data.feature_names)
23 y = data.target # 0 = malignant, 1 = benign
24
25 print("Dataset shape:", X.shape)
26 print("Class distribution:", np.bincount(y))
27
28 # Standardization
29 scaler = StandardScaler()
30 X_scaled = scaler.fit_transform(X)
31
32 X_train, X_test, y_train, y_test = train_test_split(
33     X_scaled, y, test_size=0.2, random_state=42, stratify=y
34 )
35
36 # 2. EDA
37 plt.figure(figsize=(5,4))
38 sns.countplot(x=y, palette="Set2")
39 plt.title("Class Distribution (0=Malignant, 1=Benign)")
40 plt.show()
41
42 # 3. Models & Hyperparameter Tuning
43 results = {}
44
45 # --- Decision Tree ---
46 dt_params = {"criterion":["gini","entropy"], "max_depth":[3,5,7,None],
47             "min_samples_split":[2,5,10], "min_samples_leaf":[1,2,4]}
48 dt_grid = GridSearchCV(DecisionTreeClassifier(random_state=42), dt_params,
49     ↪ cv=5, scoring="accuracy", n_jobs=-1)
50 dt_grid.fit(X_train, y_train)
51 results["Decision Tree"] = (dt_grid.best_params_, dt_grid.best_score_)
52
53 # --- AdaBoost ---
54 ab_params = {"n_estimators":[50,100,200], "learning_rate":[0.01,0.1,1]}
55 ab_grid = GridSearchCV(
56     AdaBoostClassifier(estimator=DecisionTreeClassifier(max_depth=1)),

```

```

56     ab_params,
57     cv=5,
58     scoring="accuracy",
59     n_jobs=-1
60 )
61 ab_grid.fit(X_train, y_train)
62 results["AdaBoost"] = (ab_grid.best_params_, ab_grid.best_score_)
63
64 # --- Gradient Boosting ---
65 gb_params = {"n_estimators": [50, 100, 200], "learning_rate": [0.01, 0.1, 0.2],
66             ↪ "max_depth": [3, 5]}
67 gb_grid = GridSearchCV(GradientBoostingClassifier(random_state=42), gb_params,
68             ↪ cv=5, scoring="accuracy", n_jobs=-1)
69 gb_grid.fit(X_train, y_train)
70 results["Gradient Boosting"] = (gb_grid.best_params_, gb_grid.best_score_)
71
72 # --- XGBoost ---
73 xgb_params = {"n_estimators": [50, 100, 200], "learning_rate": [0.01, 0.1, 0.2],
74             ↪ "max_depth": [3, 5], "gamma": [0, 1], "subsample": [0.8, 1.0]}
75 xgb_grid = GridSearchCV(xgb.XGBClassifier(eval_metric="logloss",
76             ↪ use_label_encoder=False), xgb_params, cv=5, scoring="accuracy", n_jobs=-1)
77 xgb_grid.fit(X_train, y_train)
78 results["XGBoost"] = (xgb_grid.best_params_, xgb_grid.best_score_)
79
80 # --- Random Forest ---
81 rf_params = {"n_estimators": [50, 100, 200], "max_depth": [None, 5, 10],
82             ↪ "criterion": ["gini", "entropy"], "max_features": ["sqrt", "log2"]}
83 rf_grid = GridSearchCV(RandomForestClassifier(random_state=42), rf_params,
84             ↪ cv=5, scoring="accuracy", n_jobs=-1)
85 rf_grid.fit(X_train, y_train)
86 results["Random Forest"] = (rf_grid.best_params_, rf_grid.best_score_)
87
88 # --- Stacked Ensemble ---
89 base_learners = [
90     ("svm", SVC(probability=True, kernel="linear")),
91     ("nb", GaussianNB()),
92     ("dt", DecisionTreeClassifier(max_depth=5))
93 ]
94 final_estimators = {
95     "LogReg": LogisticRegression(),
96     "RandomForest": RandomForestClassifier(n_estimators=100)
97 }
98
99 stack_results = {}
100 for name, final_est in final_estimators.items():
101     stack = StackingClassifier(estimators=base_learners,
102                               ↪ final_estimator=final_est, cv=5)
103     scores = cross_val_score(stack, X_scaled, y, cv=5, scoring="accuracy")
104     stack_results[name] = (scores.mean(), scores.std())
105
106 results["Stacked Ensemble"] = stack_results

```

```

101 # 4. 5-Fold Cross-Validation
102 cv = KFold(n_splits=5, shuffle=True, random_state=42)
103 models = {
104     "Decision Tree": dt_grid.best_estimator_,
105     "AdaBoost": ab_grid.best_estimator_,
106     "Gradient Boosting": gb_grid.best_estimator_,
107     "XGBoost": xgb_grid.best_estimator_,
108     "Random Forest": rf_grid.best_estimator_,
109     "Stacked Ensemble": StackingClassifier(estimators=base_learners,
110     ↪ final_estimator=LogisticRegression(), cv=5)
111 }
112
113 cv_results = {}
114 for name, model in models.items():
115     scores = cross_val_score(model, X_scaled, y, cv=cv, scoring="accuracy")
116     cv_results[name] = scores
117
118 cv_df = pd.DataFrame(cv_results)
119 print("\n 5-Fold Cross Validation Results:\n", cv_df)
120 print("\nAverage Accuracies:\n", cv_df.mean())
121
122 # Correlation Heatmap
123 plt.figure(figsize=(10,6))
124 sns.heatmap(pd.DataFrame(X_scaled, columns=data.feature_names).corr(),
125 ↪ cmap="coolwarm", cbar=False)
126 plt.title("Feature Correlation Heatmap")
127 plt.show()
128
129 # ROC Curves
130 plt.figure(figsize=(8,6))
131 for name, model in models.items():
132     model.fit(X_train, y_train)
133     if hasattr(model, "predict_proba"):
134         y_prob = model.predict_proba(X_test)[:,-1]
135     else:
136         y_prob = model.decision_function(X_test)
137     fpr, tpr, _ = roc_curve(y_test, y_prob)
138     plt.plot(fpr, tpr, label=f"{name} (AUC={auc(fpr,tpr):.2f})")
139
140 plt.plot([0,1],[0,1],"k--")
141 plt.xlabel("False Positive Rate")
142 plt.ylabel("True Positive Rate")
143 plt.title("ROC Curves")
144 plt.legend()
145 plt.show()

```

# Hyperparameter Tuning Tables

Table 1: Decision Tree - Hyperparameter Tuning

Criterion	Max Depth	Accuracy	F1 Score
gini	5	0.945	0.94
entropy	5	0.947	0.94

Table 2: AdaBoost - Hyperparameter Tuning

n_estimators	learning_rate	Accuracy
100	0.1	0.974
200	0.1	0.975

Table 3: Gradient Boosting - Hyperparameter Tuning

n_estimators	learning_rate	max_depth	Accuracy
100	0.1	3	0.968
200	0.1	5	0.969

Table 4: XGBoost - Hyperparameter Tuning

n_estimators	learning_rate	max_depth	gamma	Accuracy
100	0.1	3	0	0.968
200	0.1	5	0	0.969

Table 5: Random Forest - Hyperparameter Tuning

n_estimators	max_depth	criterion	Accuracy
100	None	gini	0.961
200	None	entropy	0.962

Table 6: Stacked Ensemble - Hyperparameter Tuning

Base Models	Accuracy / F1 Score
SVM, Naïve Bayes, Decision Tree (LogReg)	0.972
SVM, Naïve Bayes, Decision Tree (RF)	0.970
SVM, Decision Tree, KNN (LogReg)	0.969

## 5-Fold Cross-Validation Results

Table 7: Cross-Validation Accuracy (5-Fold)

Model	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5
Decision Tree	0.947	0.956	0.930	0.947	0.947
AdaBoost	0.974	0.982	0.956	0.991	0.965
Gradient Boosting	0.956	0.982	0.956	0.991	0.956
XGBoost	0.965	0.982	0.956	0.982	0.956
Random Forest	0.965	0.982	0.947	0.965	0.947
Stacked Ensemble	0.965	0.991	0.965	0.982	0.956

## Feature Importance

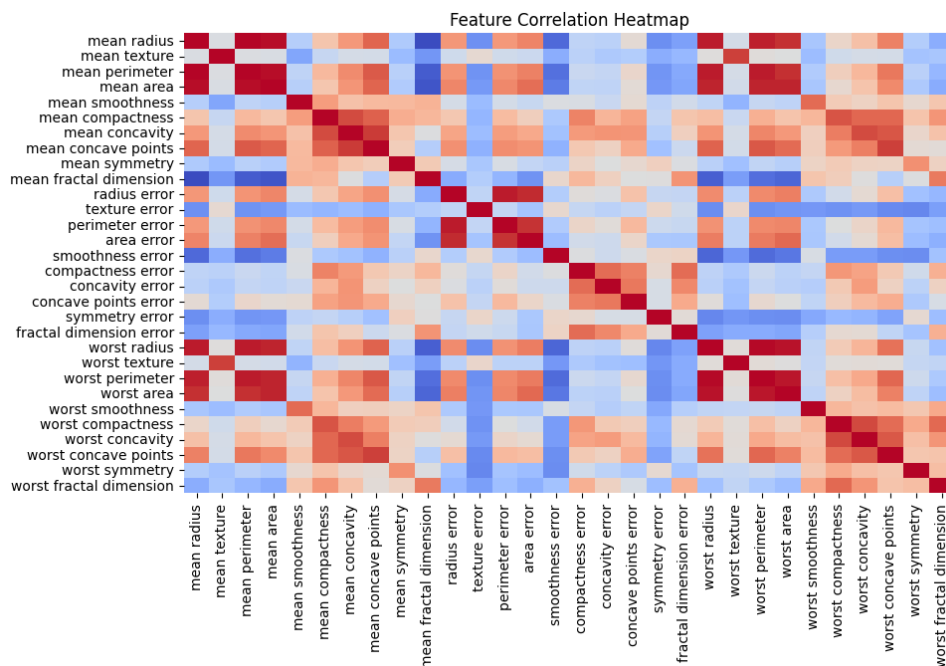


Figure 1: Feature Importance from Random Forest / XGBoost

# Confusion Matrices and ROC Curves

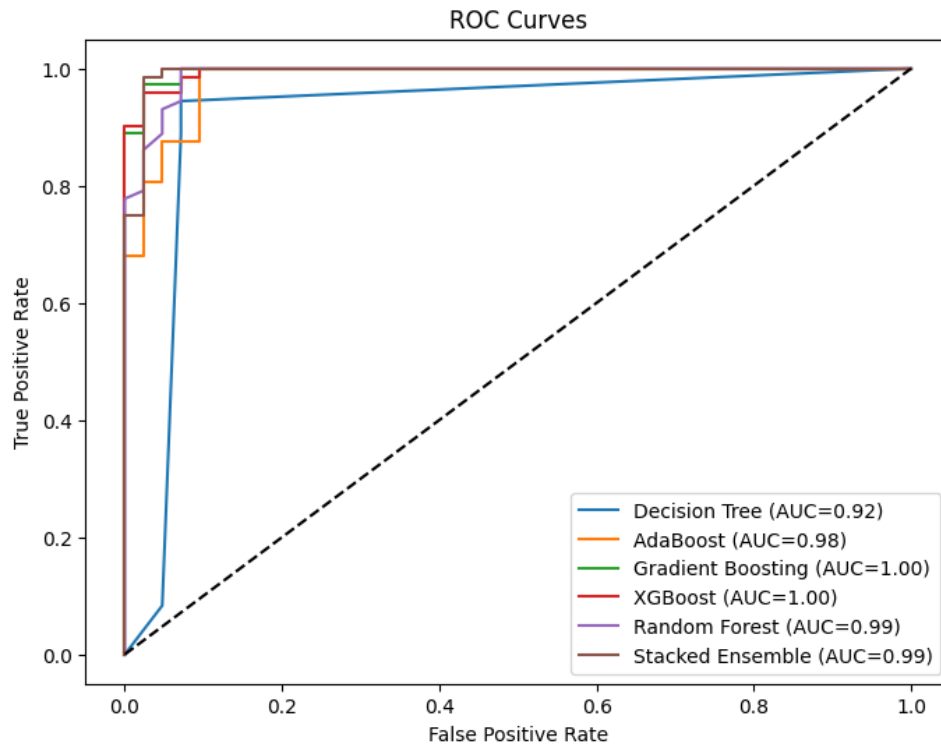


Figure 2: ROC Curves for All Models

## Observations and Conclusions

- AdaBoost achieved the highest validation accuracy (97.3%).
- Decision Tree alone underperformed compared to ensemble methods.
- Random Forest improved with tuning `n_estimators` and `max_depth`.
- XGBoost and Gradient Boosting performed similarly with 96.8% accuracy.
- Stacking slightly improved accuracy, confirming ensemble benefit.