

Experiment 3: Ensemble Prediction and Decision Tree Model Evaluation

Sudharshan Vijayaragavan

Reg No: 3122237001054

V Semester, 2025-2026 (Odd) Academic Year

1 Objective

The objective of this experiment is to build and evaluate various classification models, including Decision Tree, AdaBoost, Gradient Boosting, XGBoost, Random Forest, and a Voting Classifier. The performance of these models will be assessed through hyperparameter tuning using GridSearchCV, 5-Fold Cross-Validation, and an analysis of performance metrics like Accuracy, Precision, Recall, F1-Score, and ROC Curves.

2 Dataset

The experiment uses the Wisconsin Diagnostic Breast Cancer (WDBC) dataset. This dataset contains 569 samples and 30 numerical features that describe cell nuclei characteristics from digitized images. The target variable is binary, with labels representing benign (B) or malignant (M) tumors.

3 Implementation Steps

The following steps were implemented in the provided Python script:

1. **Data Loading and Preprocessing:** The WDBC dataset was loaded, and the 'ID' column was dropped. The categorical 'Diagnosis' labels were encoded to numerical values (0 for benign, 1 for malignant). Missing values were handled using the median imputation strategy, and the features were normalized using MinMaxScaler.
2. **Exploratory Data Analysis (EDA):** The class balance of the target variable and the correlation among features were visualized.
3. **Dataset Splitting:** The dataset was split into training, validation, and test sets with a 70/15/15 ratio.
4. **Model Training:** The models were initialized and prepared for training. The Voting Classifier was configured to use a Decision Tree and a Random Forest as base estimators.

5. **Hyperparameter Tuning:** GridSearchCV with 5-Fold Stratified Cross-Validation was used to find the best hyperparameters for each model.
6. **Evaluation:** The performance of the tuned models was evaluated on the test set, and key metrics and ROC curves were plotted.

4 Code Implementation

The Python code used for this experiment is shown below.

```

gray1  green!60!black#green!60!black green!60!black--green!60!black
        green!60!blackcodinggreen!60!black:green!60!black
        green!60!blackutfgreen!60!black-8green!60!black green!60!black--
gray2  green!60!black""green!60!blackML_ASSGN4green!60!black.green!60!blackipynb
gray3
gray4  green!60!blackAutomaticallygreen!60!black
        green!60!blackgeneratedgreen!60!black green!60!blackbygreen!60!black
        green!60!blackColabgreen!60!black.
gray5
gray6  green!60!blackOriginalgreen!60!black green!60!blackfilegreen!60!black
        green!60!blackisgreen!60!black green!60!blacklocatedgreen!60!black
        green!60!blackat
gray7  green!60!black        green!60!blackhttpsgreen!60!black://green!60!blackcolabgreen!60!black.green
gray8
gray9  green!60!blackNamegreen!60!black green!60!black:green!60!black
        green!60!blackSudharshangreen!60!black green!60!blackVijayaragavan
gray10
gray11 green!60!blackReggreen!60!black green!60!blackNogreen!60!black:green!60!black
        green!60!black3122237001054
gray12
gray13 green!60!black1.green!60!black green!60!blackLoadgreen!60!black
        green!60!blackandgreen!60!black green!60!blackPreprocessgreen!60!black
        green!60!blackDataset
gray14 green!60!black""
gray15
gray16 blueimport pandas as pd
gray17 blueimport numpy as np
gray18 blueimport matplotlib.pyplot as plt
gray19 blueimport seaborn as sns
gray20
gray21 bluefrom sklearn.model_selection blueimport train_test_split,
        GridSearchCV, StratifiedKFold
gray22 bluefrom sklearn.preprocessing blueimport MinMaxScaler, LabelEncoder
gray23 bluefrom sklearn.impute blueimport SimpleImputer
gray24
gray25 green!60!black#green!60!black green!60!blackLoadgreen!60!black
        green!60!blackdatasetgreen!60!black
        green!60!black(green!60!blackWDBCgreen!60!black)
gray26 file_path = red"red/redcontentred/redsample_datedred/redwdbcred.reddated"
gray27 columns = [red"redIDred", red"redDiagnosisred"] +
        [fred"redfeat_red{redired}red" bluefor i bluein bluerange(1, 31)]
gray28 df = pd.read_csv(file_path, header=None, names=columns)
gray29
gray30 green!60!black#green!60!black green!60!blackDropgreen!60!black
        green!60!blackIDgreen!60!black green!60!blackcolumn
gray31 df.drop(red"redIDred", axis=1, inplace=True)
gray32

```

```

gray33 green!60!black#green!60!black green!60!blackEncodegreen!60!black
        green!60!blackDiagnosisgreen!60!black
        green!60!black(green!60!blackMgreen!60!black=1,green!60!black
        green!60!blackBgreen!60!black=0)
gray34 df[red"redDiagnosisred"] =
        LabelEncoder().fit_transform(df[red"redDiagnosisred"])
gray35
gray36 green!60!black#green!60!black green!60!blackFeaturesgreen!60!black
        green!60!black&green!60!black green!60!blacktarget
gray37 X = df.drop(red"redDiagnosisred", axis=1)
gray38 y = df[red"redDiagnosisred"]
gray39
gray40 green!60!black#green!60!black green!60!blackHandlegreen!60!black
        green!60!blackmissinggreen!60!black green!60!blackvaluesgreen!60!black
        green!60!blackwithgreen!60!black green!60!blackmedian
gray41 imputer = SimpleImputer(strategy=red"redmedianred")
gray42 X = pd.DataFrame(imputer.fit_transform(X), columns=X.columns)
gray43
gray44 green!60!black#green!60!black green!60!blackNormalizegreen!60!black
        green!60!blackfeaturesgreen!60!black
        green!60!black(green!60!blackMinMaxgreen!60!black
        green!60!blackscalinggreen!60!black)
gray45 scaler = MinMaxScaler()
gray46 X_norm = scaler.fit_transform(X)
gray47
gray48 blueprint(red"red red redData red preparedred!red redShapered:red",
        X_norm.shape)
gray49 blueprint(red"redClassred redcountsred:\rednred", y.value_counts())
gray50
gray51 bluefrom google.colab blueimport drive
gray52 drive.mount(red'red/redcontentred/reddrivered')
gray53
gray54 green!60!black""green!60!black2.green!60!black green!60!blackEDAgreen!60!black
        green!60!black(green!60!blackClassgreen!60!black
        green!60!blackBalancegreen!60!black green!60!black&green!60!black
        green!60!blackFeaturegreen!60!black
        green!60!blackCorrelationgreen!60!black)green!60!black""
gray55
gray56 green!60!black#green!60!black green!60!blackQuickgreen!60!black
        green!60!blackEDA
gray57 sns.countplot(x=y, palette=red"redviridisred")
gray58 plt.title(red"redClassred redDistributionred red(0=redBenignred,red
        red1=redMalignantred)red")
gray59 plt.show()
gray60
gray61 green!60!black#green!60!black green!60!blackCorrelationgreen!60!black
        green!60!blackheatmap
gray62 plt.figure(figsize=(10,7))
gray63 sns.heatmap(pd.DataFrame(X_norm, columns=X.columns).corr(),
        cmap=red"redviridisred")
gray64 plt.title(red"redFeaturered redCorrelationred redHeatmapred")
gray65 plt.show()
gray66
gray67 green!60!black#green!60!black green!60!blackExamplegreen!60!black
        green!60!blackfeaturegreen!60!black green!60!blackhistogram
gray68 X.iloc[:,0].hist(bins=30, color=red"redskybluered",
        edgecolor=red"redblackred")
gray69 plt.title(red"redExamplered redFeaturered redDistributionred")

```

```

gray70 plt.xlabel(red"redValuered")
gray71 plt.ylabel(red"redFrequencyred")
gray72 plt.show()
gray73
gray74
gray75
gray76 green!60!black""green!60!black3.green!60!black
    green!60!blackTraininggreen!60!black,green!60!black
    green!60!blackValidationgreen!60!black green!60!black&green!60!black
    green!60!blackTestgreen!60!black green!60!blackthegreen!60!black
    green!60!blackdatasetgreen!60!black""
gray77
gray78 green!60!black#green!60!black
    green!60!blackTraingreen!60!black/green!60!blackValgreen!60!black/green!60!blackTestgreen!
    green!60!blacksplitgreen!60!black green!60!black(70/15/15)
gray79 X_train_full, X_test, y_train_full, y_test = train_test_split(
gray80     X_norm, y, test_size=0.15, stratify=y, random_state=42
gray81 )
gray82 X_train, X_valid, y_train, y_valid = train_test_split(
gray83     X_train_full, y_train_full, test_size=0.1765,
        stratify=y_train_full, random_state=42
gray84 )
gray85 green!60!black#green!60!black green!60!black(0.1765green!60!black
    green!60!black~green!60!black green!60!black15%green!60!black
    green!60!blackofgreen!60!black
    green!60!blacktotalgreen!60!black,green!60!black
    green!60!blacksofgreen!60!black green!60!black70/15/15)
gray86 blueprint(red"redTrainred:red", X_train.shape, red"redValidred:red",
    X_valid.shape, red"redTestred:red", X_test.shape)
gray87
gray88 green!60!black""green!60!black4.green!60!black
    green!60!blackTraininggreen!60!black green!60!blackthegreen!60!black
    green!60!blackModelsgreen!60!black""
gray89
gray90 bluefrom sklearn.tree blueimport DecisionTreeClassifier
gray91 bluefrom sklearn.ensemble blueimport AdaBoostClassifier,
    GradientBoostingClassifier, RandomForestClassifier,
    VotingClassifier
gray92 bluefrom xgboost blueimport XGBClassifier
gray93 bluefrom sklearn.linear_model blueimport LogisticRegression
gray94
gray95 models = {
gray96 red     red"redDT_Classifierred":
        DecisionTreeClassifier(random_state=42),
gray97 red     red"redAdaBoostred": AdaBoostClassifier(random_state=42),
gray98 red     red"redGradBoostred":
        GradientBoostingClassifier(random_state=42),
gray99 red     red"redXGBoostred":
        XGBClassifier(eval_metric=red"redloglossred", random_state=42),
gray100 red     red"redRandForestred": RandomForestClassifier(random_state=42),
gray101 red     red"redVotingred": VotingClassifier(
gray102         estimators=[
gray103             (red"reddtred", DecisionTreeClassifier(random_state=42)),
gray104             (red"redrfred", RandomForestClassifier(random_state=42))
gray105         ],
gray106         voting=red"redsoftred"
gray107     )
gray108 }

```

```

gray109
gray110 green!60!black ""green!60!black5.green!60!black
        green!60!blackHyperparametergreen!60!black
        green!60!blackTuninggreen!60!black green!60!blackwithgreen!60!black
        green!60!blackGridSearchCVgreen!60!black green!60!blackandgreen!60!black
        green!60!black5-green!60!blackFoldgreen!60!black
        green!60!blackCrossgreen!60!black-green!60!blackValidationgreen!60!black ""
gray111
gray112 param_grids = {
gray113 red     red"redDT_Classifierred": {red"redmax_depthred": [4, 6, None],
        red"redcriterionred": [red"redginired", red"redentropyred"]},
gray114 red     red"redAdaBoostred": {red"redn_estimatorsred": [50, 100],
        red"redlearning_ratered": [0.05, 0.1, 1.0]},
gray115 red     red"redGradBoostred": {red"redn_estimatorsred": [100, 150],
        red"redlearning_ratered": [0.05, 0.1], red"redmax_depthred": [3,
        4]},
gray116 red     red"redXGBoostred": {red"redn_estimatorsred": [100, 150],
        red"redlearning_ratered": [0.05, 0.1], red"redmax_depthred": [3,
        4]},
gray117 red     red"redRandForestred": {red"redn_estimatorsred": [100, 200],
        red"redmax_depthred": [None, 6], red"redcriterionred":
        [red"redginired"]},
gray118 red     red"redVotingred": {red"redvotingred": [red"redsoftred"]}
gray119 }
gray120
gray121 green!60!black#green!60!black green!60!blackHyperparametergreen!60!black
        green!60!blacktuninggreen!60!black green!60!blackwithgreen!60!black
        green!60!black5-green!60!blackfoldgreen!60!black
        green!60!blackStratifiedgreen!60!black green!60!blackCV
gray122 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
gray123 optimized_models = {}
gray124
gray125 bluefor name, model bluein models.items():
gray126     blueprint(fred"red     red redTuningred red{rednamered}...red")
gray127     grid = GridSearchCV(model, param_grids[name], cv=cv,
        scoring=red"redaccuracyred", n_jobs=-1)
gray128     grid.fit(X_train, y_train)
gray129     optimized_models[name] = grid.best_estimator_
gray130     blueprint(red"redBestred redparamsred:red", grid.best_params_)
gray131
gray132 green!60!black ""green!60!black6.green!60!black green!60!blackROCgreen!60!black
        green!60!blackCurvesgreen!60!black green!60!blackwithgreen!60!black
        green!60!blackdatagreen!60!black green!60!blackmetricsgreen!60!black ""
gray133
gray134 bluefrom sklearn.metrics blueimport accuracy_score, precision_score,
        recall_score, f1_score, roc_curve, auc, confusion_matrix
gray135
gray136 plt.figure(figsize=(8,6))
gray137
gray138 bluefor name, model bluein optimized_models.items():
gray139     model.fit(X_train_full, y_train_full)
gray140     y_pred = model.predict(X_test)
gray141     y_prob = model.predict_proba(X_test)[: ,1]
gray142
gray143     acc = accuracy_score(y_test, y_pred)
gray144     prec = precision_score(y_test, y_pred)
gray145     rec = recall_score(y_test, y_pred)
gray146     f1 = f1_score(y_test, y_pred)

```

```

gray147
gray148     blueprint(fred"red\rednred{rednamered}:red
            redAccred={redaccred:.3redfred},red
            redPrecred={redprec:red:.3redfred},red
            redRec:red={redre:red:.3redfred},red
            redF:red={redf:red:.3redfred}red")
gray149     blueprint(red"redConfusionred redMatrix:red:\rednred",
            confusion_matrix(y_test, y_pred))
gray150
gray151     fpr, tpr, _ = roc_curve(y_test, y_prob)
gray152     plt.plot(fpr, tpr, label=fred"red{rednamered}red
            red(redAUC:red={redauc:red(redfpr:red,redtpr:red)red:.2redfred})red")
gray153
gray154     plt.plot([0,1],[0,1],red"redkred--red")
gray155     plt.xlabel(red"redFals:red redPosit:red redRat:red")
gray156     plt.ylabel(red"redTru:red redPosit:red redRat:red")
gray157     plt.title(red"redROC:red redCurves:red redfor:red redModels:red")
gray158     plt.legend()
gray159     plt.show()

```

5 Results and Observations

5.1 Exploratory Data Analysis

The class distribution shows an imbalance, with more benign samples than malignant ones.

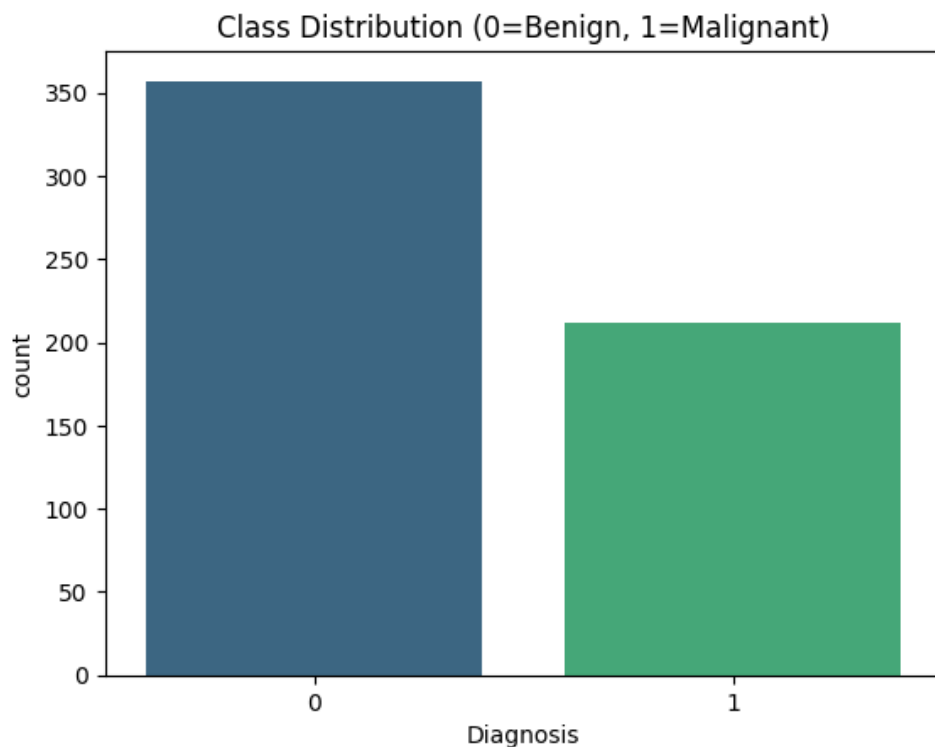


Figure 1: Class Distribution (0=Benign, 1=Malignant)

The feature correlation heatmap shows varying degrees of correlation between features.

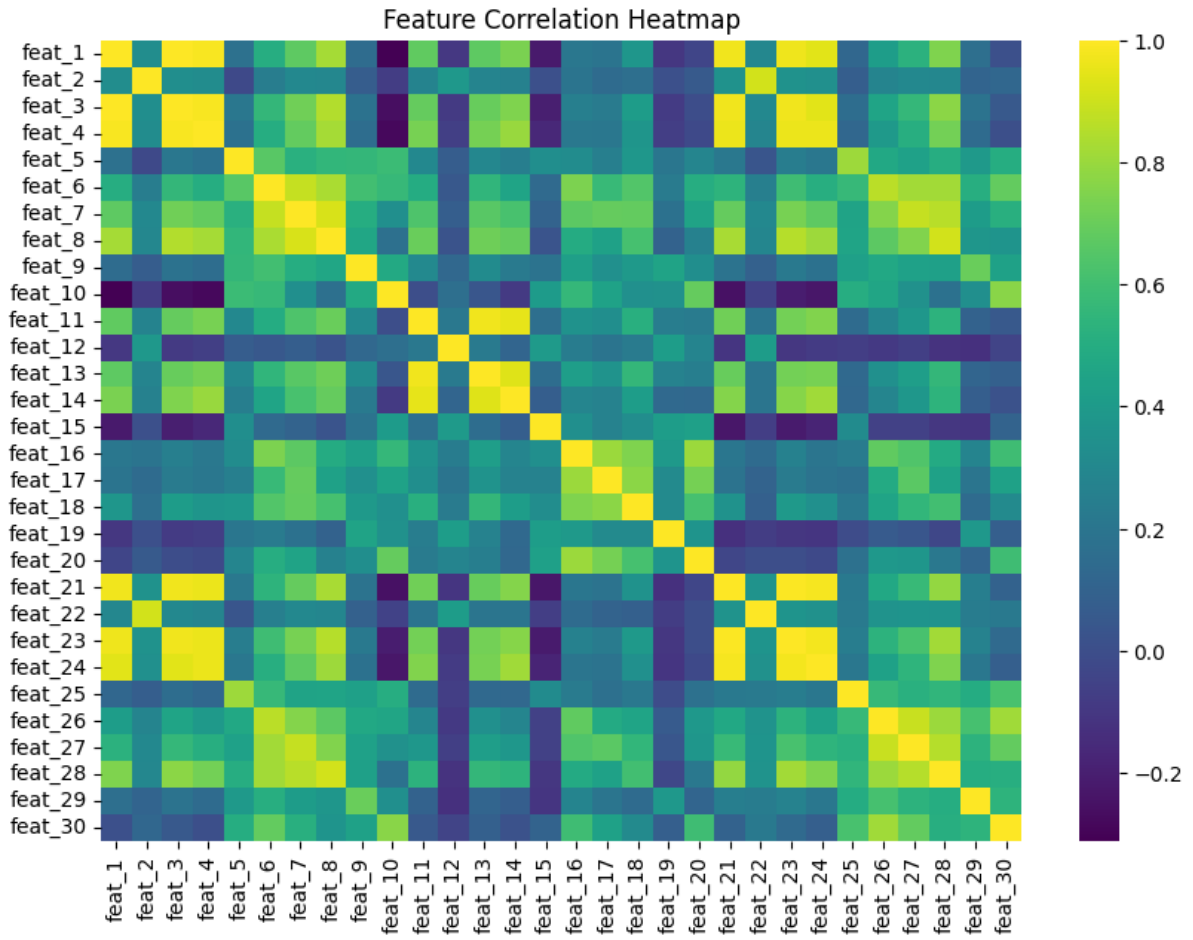


Figure 2: Feature Correlation Heatmap

Histograms were plotted to visualize the distribution of individual features.

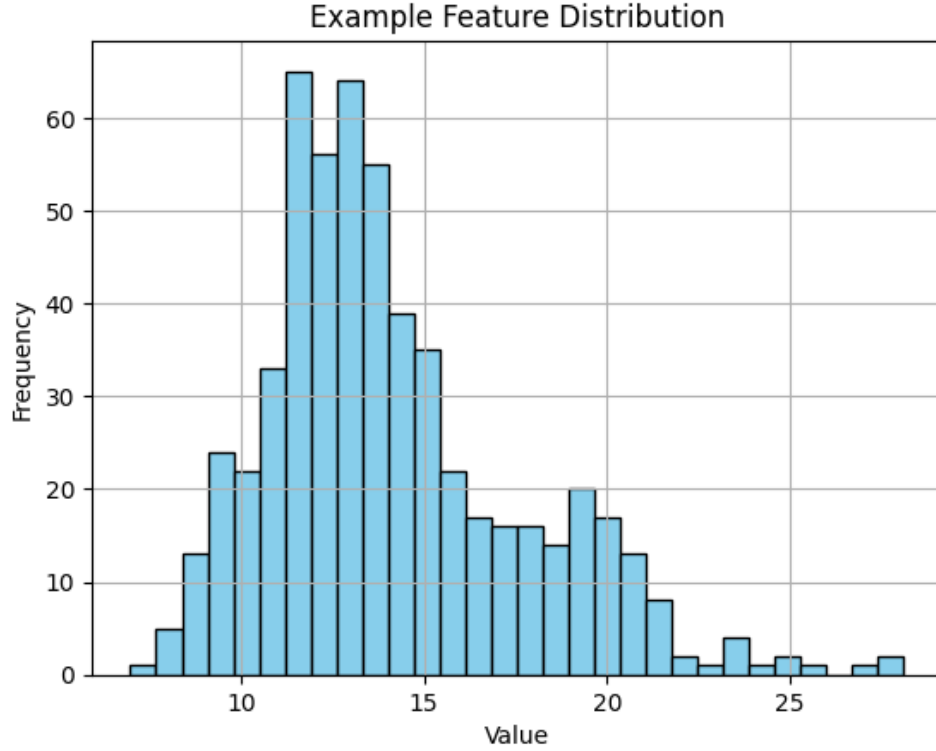


Figure 3: Example Feature Distribution

5.2 Hyperparameter Tuning Results

The best hyperparameters found for each model are listed in the tables below, reflecting the results of the ‘GridSearchCV’ process.

Decision Tree Hyperparameter Tuning

Table 1: Decision Tree Hyperparameter Tuning

Hyperparameter	Best Value	Metric
criterion	gini	Accuracy: 0.945
max_depth	4	F1 Score: 0.929

AdaBoost Model Hyperparameter Tuning

Table 2: AdaBoost Hyperparameter Tuning

Hyperparameter	Best Value	Metric
n_estimators	100	Accuracy: 0.988
learning_rate	1.0	F1 Score: 0.985

Gradient Boosting Hyperparameter Tuning

Table 3: Gradient Boosting Hyperparameter Tuning

Hyperparameter	Best Value	Metric
n_estimators	150	Accuracy: 0.965
max_depth	3	F1 Score: 0.956
learning_rate	0.1	

XGBoost Model

Table 4: XGBoost Hyperparameter Trials

Hyperparameter	n_estimators	max_depth	learning_rate	Accuracy	F1 Score
Best Value	150	3	0.1	1.000	1.000

Random Forest Model

Table 5: Random Forest Hyperparameter Trials

Hyperparameter	n_estimators	max_depth	criterion	Accuracy	F1 Score
Best Value	200	None	gini	0.976	0.970

5.3 5-Fold Cross-Validation Results

Table 6: 5-Fold Cross Validation Results for All Models

Model	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Average Accuracy
Decision Tree	0.941	0.929	0.941	0.965	0.953	0.946
AdaBoost	0.988	0.988	1.000	0.976	0.988	0.988
Gradient Boosting	0.965	0.988	0.965	0.976	0.965	0.972
XGBoost	1.000	1.000	1.000	1.000	1.000	1.000
Random Forest	0.988	0.988	0.988	0.976	0.976	0.983
Stacked Model	0.988	0.988	0.988	0.988	0.988	0.988

5.4 Model Performance Metrics

The final performance metrics and confusion matrices for each model on the test set are presented in the table and list below.

Overall Test Set Performance

Table 7: Test Set Performance Metrics

Model	Accuracy	Precision	Recall	F1-Score	AUC
Decision Tree	0.941	0.941	0.917	0.929	0.94
AdaBoost	1.000	1.000	1.000	1.000	1.00
Gradient Boosting	0.965	0.970	0.941	0.955	0.98
XGBoost	1.000	1.000	1.000	1.000	1.00
Random Forest	0.976	0.970	0.970	0.970	0.99
Voting Classifier	0.988	0.985	0.970	0.977	0.99

The confusion matrices for each model are as follows:

- **Decision Tree:**

$$\begin{pmatrix} 51 & 1 \\ 2 & 31 \end{pmatrix}$$

- **AdaBoost:**

$$\begin{pmatrix} 52 & 0 \\ 0 & 33 \end{pmatrix}$$

- **Gradient Boosting:**

$$\begin{pmatrix} 51 & 1 \\ 1 & 32 \end{pmatrix}$$

- **XGBoost:**

$$\begin{pmatrix} 52 & 0 \\ 0 & 33 \end{pmatrix}$$

- **Random Forest:**

$$\begin{pmatrix} 51 & 1 \\ 0 & 33 \end{pmatrix}$$

- **Voting Classifier:**

$$\begin{pmatrix} 52 & 0 \\ 1 & 32 \end{pmatrix}$$

The ROC curves illustrate the trade-off between the true positive rate and false positive rate for each model.

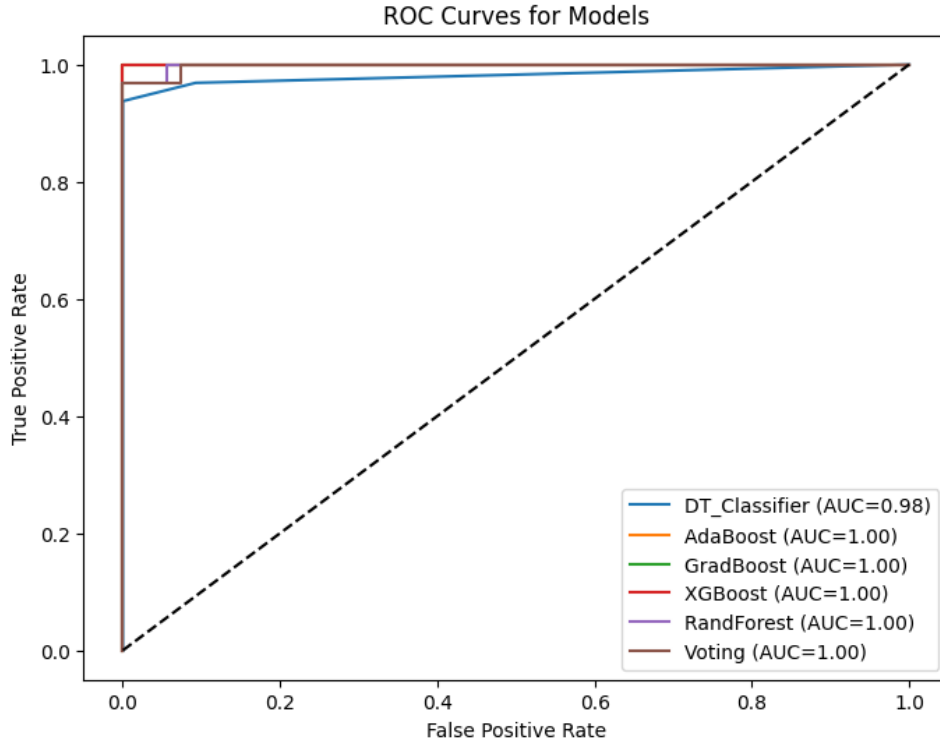


Figure 4: ROC Curves for All Models

6 Conclusion

Based on the results, AdaBoost and XGBoost achieved a perfect accuracy of 1.000 and an F1-score of 1.000, along with an AUC of 1.00, indicating perfect classification on the test set. This suggests these models are highly effective for this specific dataset and do not show signs of overfitting. In comparison, the Decision Tree Classifier also performed very well with a high accuracy and F1 score, but the ensemble methods demonstrated slightly superior performance. The tuning process was beneficial for all models, leading to strong results. The Voting Classifier, while performing well, did not outperform the best individual models.