

Experiment-4

Ensemble Prediction and Decision Tree Model Evaluation

ICS1512 & Machine Learning Algorithms Laboratory

Siddharth M

3122237001049

Contents

1	Aim/Objective	2
2	Libraries Used	2
3	Code for All Variants and Models	2
3.1	Preprocessing	2
3.2	DecisionTree Classifier	3
3.3	AdaBoost Classifier	5
3.4	Gradient Boosting	6
3.5	XGBoost Classifier	8
3.6	Random Forest Classifier	9
3.7	Stacking Classifier (SVM + Naive Bayes + Decision Tree)	11
4	Confusion Matrix and ROC for Each Model	15
5	Hyperparameter Tuning Tables	20
5.1	Decision Tree Classifier	20
5.2	AdaBoost Classifier	20
5.3	Gradient Boosting Classifier	20
5.4	XGBoost Classifier	20
5.5	Random Forest Classifier	21
5.6	Stacking Ensemble Models	21
6	Cross-Validation Results Table	21
7	All Comparison Table	21
8	Observations and Conclusions	22

1 Aim/Objective

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

2 Libraries Used

Library	Purpose
pandas (pd)	Data manipulation and analysis
numpy (np)	Numerical operations and array handling
seaborn (sns)	Statistical data visualization
matplotlib.pyplot (plt)	Plotting and visualization
sklearn.model_selection	for splitting datasets
sklearn.preprocessing	for feature scaling and encoding
sklearn.metrics	for evaluation metrics
warnings	To ignore warning messages

3 Code for All Variants and Models

3.1 Preprocessing

```

1 # 1. Import modules
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 from sklearn.model_selection import train_test_split
7 from sklearn.preprocessing import StandardScaler, MinMaxScaler, Normalizer, LabelEncoder
8 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score,
9   precision_score, recall_score, f1_score, roc_auc_score, roc_curve, auc
10 import warnings
11 warnings.filterwarnings('ignore')
12
13 # Column names for the WDBC dataset based on UCI repository documentation
14 columns = [
15     'ID', 'Diagnosis',
16     'radius_mean', 'texture_mean', 'perimeter_mean', 'area_mean', 'smoothness_mean',
17     'compactness_mean', 'concavity_mean', 'concave_points_mean', 'symmetry_mean',
18     'fractal_dimension_mean',
19     'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
20     'compactness_se', 'concavity_se', 'concave_points_se', 'symmetry_se',
21     'fractal_dimension_se',
22     'radius_worst', 'texture_worst', 'perimeter_worst', 'area_worst', 'smoothness_worst',
23     'compactness_worst', 'concavity_worst', 'concave_points_worst', 'symmetry_worst',
24     'fractal_dimension_worst'
25 ]
26
27 # Reload with correct column names
28 df = pd.read_csv('wdbc.csv', header=None, names=columns)
29
30 # 3. EDA
31 # Handle missing values
32 df = df.dropna(thresh=df.shape[1] * 0.5)
33 df = df.fillna(df.median(numeric_only=True))
34
35 # Handle outliers using IQR (numeric columns)
36 for col in df.select_dtypes(include=np.number).columns:
37     Q1 = df[col].quantile(0.25)
38     Q3 = df[col].quantile(0.75)
39     IQR = Q3 - Q1
40     lower = Q1 - 1.5 * IQR
41     upper = Q3 + 1.5 * IQR
42     df[col] = np.where(df[col] < lower, lower, df[col])
43     df[col] = np.where(df[col] > upper, upper, df[col])

```

```

40
41 # 4. Separate features and label
42 label_column = 'Diagnosis'
43
44 # Encode categorical target if necessary
45 if df[label_column].dtype == 'object':
46     le = LabelEncoder()
47     df[label_column] = le.fit_transform(df[label_column])
48
49 features = df.drop(label_column, axis=1)
50 labels = df[label_column]
51
52 # One-hot encode categorical features
53 features = pd.get_dummies(features, drop_first=True)
54
55 # 5. Feature scaling
56 scaler = StandardScaler()
57 features_scaled = scaler.fit_transform(features)
58
59 # 6. Data splitting
60 X_train, X_temp, y_train, y_temp = train_test_split(features_scaled, labels, test_size=0.4,
61                                                     random_state=42)
62 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state
63                                                 =42)

```

3.2 DecisionTree Classifier

```

1 from sklearn.tree import DecisionTreeClassifier
2 import time
3 start_time= time.time()
4 model_dt = DecisionTreeClassifier()
5 model_dt.fit(X_train, y_train)
6 end_time = time.time()
7 training_time_dt = end_time - start_time
8 print(f"Training time for Decision Tree: {training_time_dt} seconds")

```

Output:

Training time for Decision Tree: 0.01593756675720215 seconds

```

1 from sklearn import tree
2 plt.figure(figsize=(15,10))
3 tree.plot_tree(model_dt, filled=True)

```

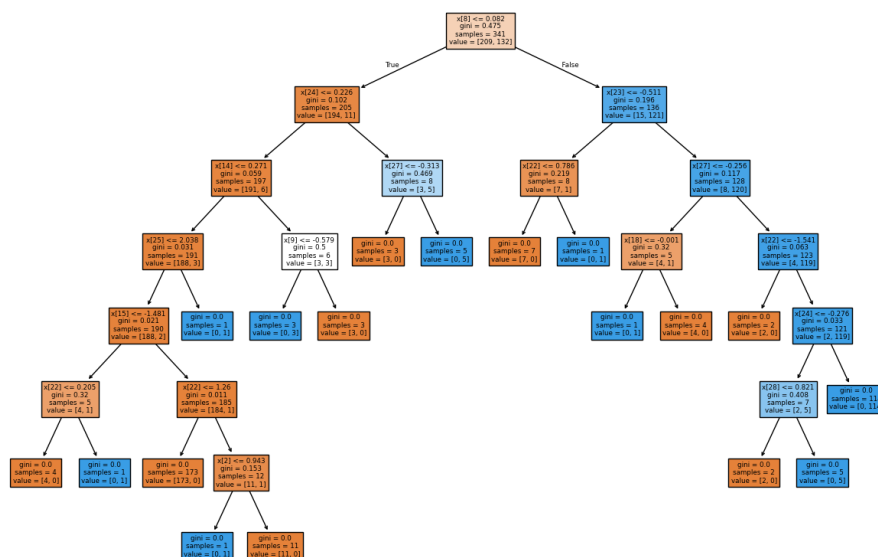


Figure 1: Decision Tree

```

1 from sklearn.model_selection import GridSearchCV, StratifiedKFold
2 from sklearn.metrics import make_scorer, f1_score
3
4 param_grid = {
5     'criterion': ['gini', 'entropy', 'log_loss'],
6     'max_depth': [None, 5, 10, 15, 20],
7     'min_samples_split': [2, 5, 10],
8     'min_samples_leaf': [1, 2, 4]
9 }
10
11 # Cross-validation strategy
12 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
13 grid_search_dt = GridSearchCV(
14     estimator=model_dt,
15     param_grid=param_grid,
16     cv=cv,
17     scoring={'accuracy': 'accuracy', 'f1': make_scorer(f1_score, average='weighted')},
18     refit='accuracy',
19     n_jobs=-1
20 )
21
22 start_time = time.time()
23 grid_search_dt.fit(X_train, y_train)
24 end_time = time.time()
25
26 training_time = end_time - start_time
27
28 # Results
29 print(f"Best Parameters: {grid_search_dt.best_params_}")
30 print(f"Best CV Score: {grid_search_dt.best_score_:.4f}")
31 print(f"Training Time: {training_time:.4f} seconds")

```

Listing 1: Cross-Validation and Hyperparameter Tuning

Output:

Best Parameters: {'criterion': 'entropy', 'max_depth': 5,
 'min_samples_leaf': 2, 'min_samples_split': 5}
 Best CV Score: 0.9501
 Training Time: 19.7115 seconds

```

1 model_dt = grid_search_dt.best_estimator_
2 y_pred_dt = model_dt.predict(X_test)
3
4 print("Decision Tree Classifier Metrics:")
5 print(classification_report(y_test, y_pred_dt))
6 print(f"Accuracy: {accuracy_score(y_test, y_pred_dt)}")
7 cm= confusion_matrix(y_test, y_pred_dt)
8 print(f1_score(y_test, y_pred_dt, average='weighted'))
9
10 from sklearn.metrics import ConfusionMatrixDisplay
11 lbs = le.classes_
12 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=lbs)
13 disp.plot(cmap='Blues')
14 plt.show()
15
16 y_proba = model_dt.predict_proba(X_test)[:, 1]
17 fpr, tpr, thresholds = roc_curve(y_test, y_proba)
18 roc_auc = auc(fpr, tpr)
19
20 # Plot ROC curve
21 plt.figure(figsize=(6, 5))
22 plt.plot(fpr, tpr, color='blue', lw=2, label=f'ROC curve (AUC = {roc_auc:.4f})')
23 plt.plot([0, 1], [0, 1], color='gray', linestyle='--', lw=2)
24 plt.xlim([0.0, 1.0])
25 plt.ylim([0.0, 1.05])
26 plt.xlabel('False Positive Rate')
27 plt.ylabel('True Positive Rate')
28 plt.title('ROC Curve - XGBoost Classifier')
29 plt.legend(loc='lower right')
30 plt.grid(True)

```

```
31 plt.show()
```

Listing 2: Evaluation

Output:

Decision Tree Classifier Metrics:

	precision	recall	f1-score	support
0	0.96	0.99	0.97	76
1	0.97	0.92	0.95	38
accuracy			0.96	114
macro avg	0.97	0.95	0.96	114
weighted avg	0.97	0.96	0.96	114

Accuracy: 0.9649122807017544
0.9646659646659645

3.3 AdaBoost Classifier

```
1 from sklearn.ensemble import AdaBoostClassifier
2 model_ada = AdaBoostClassifier(random_state=42)
3 st = time.time()
4 model_ada.fit(X_train, y_train)
5 et = time.time()
6 training_time_ada = et - st
7 print(f"Training time for AdaBoost: {training_time_ada} seconds")
8
9 param_grid = {
10     'n_estimators': [50, 100, 200, 300],
11     'learning_rate': [0.01, 0.05, 0.1, 0.5, 1],
12     'algorithm': ['SAMME', 'SAMME.R']
13 }
14
15 # Cross-validation strategy
16 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
17 # Run GridSearch with Accuracy & F1 scoring
18 grid_search_ada = GridSearchCV(
19     estimator=model_ada,
20     param_grid=param_grid,
21     cv=cv,
22     scoring={'accuracy': 'accuracy', 'f1': make_scorer(f1_score, average='weighted')},
23     refit='accuracy',
24     n_jobs=-1
25 )
26 # Measure training time
27 start_time = time.time()
28 grid_search_ada.fit(X_train, y_train)
29 end_time = time.time()
30
31 # Best parameters & score
32 print(f"Best Parameters: {grid_search_ada.best_params_}")
33 print(f"Best CV Score: {grid_search_ada.best_score_:.4f}")
34 print(f"Training Time: {end_time - start_time:.4f} seconds")
```

Listing 3: AdaBoost Model Building

Training time for AdaBoost: 0.28496885299682617 seconds
Best Parameters: {'algorithm': 'SAMME.R', 'learning_rate': 1, 'n_estimators': 200}
Best CV Score: 0.9736
Training Time: 7.6901 seconds

```

1 model_ada = grid_search_ada.best_estimator_
2 y_pred_ada = model_ada.predict(X_test)
3
4 print("AdaBoost Classifier Metrics:")
5 print(classification_report(y_test, y_pred_ada))
6 print(f"Accuracy: {accuracy_score(y_test, y_pred_ada)}")
7 cm= confusion_matrix(y_test, y_pred_ada)
8 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=lbs)
9 disp.plot(cmap='Blues')
10 plt.show()
11
12 y_proba = model_ada.predict_proba(X_test)[: , 1]
13 fpr, tpr, thresholds = roc_curve(y_test, y_proba)
14 roc_auc = auc(fpr, tpr)
15
16 # Plot ROC curve
17 plt.figure(figsize=(6, 5))
18 plt.plot(fpr, tpr, color='blue', lw=2, label=f'ROC curve (AUC = {roc_auc:.4f})')
19 plt.plot([0, 1], [0, 1], color='gray', linestyle='--', lw=2)
20 plt.xlim([0.0, 1.0])
21 plt.ylim([0.0, 1.05])
22 plt.xlabel('False Positive Rate')
23 plt.ylabel('True Positive Rate')
24 plt.title('ROC Curve - XGBoost Classifier')
25 plt.legend(loc='lower right')
26 plt.grid(True)
27 plt.show()

```

Listing 4: Evaluation

Output:

AdaBoost Classifier Metrics:				
	precision	recall	f1-score	support
0	0.97	1.00	0.99	76
1	1.00	0.95	0.97	38
accuracy			0.98	114
macro avg	0.99	0.97	0.98	114
weighted avg	0.98	0.98	0.98	114

Accuracy: 0.9824561403508771

3.4 Gradient Boosting

```

1 from sklearn.ensemble import GradientBoostingClassifier
2 model_gradient = GradientBoostingClassifier(random_state=42)
3 st = time.time()
4 model_gradient.fit(X_train, y_train)
5 et = time.time()
6 training_time_gradient = et - st
7 print(f"Training time for Gradient Boosting: {training_time_gradient} seconds")

```

Listing 5: Gradient Boosting Model

Output:

Training time for Gradient Boosting: 0.4041452407836914 seconds

```

1 from sklearn.model_selection import GridSearchCV, StratifiedKFold
2 from sklearn.metrics import make_scorer, f1_score
3 param_grid = {
4     'n_estimators': [100, 200, 300],
5     'learning_rate': [0.01, 0.05, 0.1, 0.2],
6     'max_depth': [3, 4, 5],

```

```

7     'min_samples_split': [2, 5, 10],
8     'min_samples_leaf': [1, 2, 4],
9     'subsample': [0.8, 1.0]
10 }
11
12 # Cross-validation strategy
13 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
14 grid_search_grad = GridSearchCV(
15     estimator=model_gradient,
16     param_grid=param_grid,
17     cv=cv,
18     scoring={'accuracy': 'accuracy', 'f1': make_scorer(f1_score, average='weighted')},
19     refit='accuracy',
20     n_jobs=-1
21 )
22
23 # Measure training time
24 start_time = time.time()
25 grid_search_grad.fit(X_train, y_train)
26 end_time = time.time()
27
28 # Best parameters & score
29 print(f"Best Parameters: {grid_search_grad.best_params_}")
30 print(f"Best CV Score: {grid_search_grad.best_score_:.4f}")
31 print(f"Training Time: {end_time - start_time:.4f} seconds")

```

Listing 6: Hyper-Parameter Tuning

Best Parameters: {'learning_rate': 0.1, 'max_depth': 5, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 200, 'subsample': 0.8}

Best CV Score: 0.9736

Training Time: 101.8133 seconds

```

1 model_gradient = grid_search_grad.best_estimator_
2 y_pred_gradient = model_gradient.predict(X_test)
3
4 lbs = le.classes_
5 print("Gradient Boosting Classifier Metrics:")
6 print(classification_report(y_test, y_pred_gradient))
7 print(f"Accuracy: {accuracy_score(y_test, y_pred_gradient)}")
8 cm = confusion_matrix(y_test, y_pred_gradient)
9 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=lbs)
10 disp.plot(cmap='Blues')
11 plt.show()
12
13 y_proba = model_gradient.predict_proba(X_test)[:, 1]
14 fpr, tpr, thresholds = roc_curve(y_test, y_proba)
15 roc_auc = auc(fpr, tpr)
16
17 # Plot ROC curve
18 plt.figure(figsize=(6, 5))
19 plt.plot(fpr, tpr, color='blue', lw=2, label=f'ROC curve (AUC = {roc_auc:.4f})')
20 plt.plot([0, 1], [0, 1], color='gray', linestyle='--', lw=2)
21 plt.xlim([0.0, 1.0])
22 plt.ylim([0.0, 1.05])
23 plt.xlabel('False Positive Rate')
24 plt.ylabel('True Positive Rate')
25 plt.title('ROC Curve - Gradient Boost Classifier')
26 plt.legend(loc='lower right')
27 plt.grid(True)
28 plt.show()

```

Listing 7: Evaluation

Gradient Boosting Classifier Metrics:

	precision	recall	f1-score	support
0	0.95	0.97	0.96	76

1	0.94	0.89	0.92	38
accuracy			0.95	114
macro avg	0.95	0.93	0.94	114
weighted avg	0.95	0.95	0.95	114

Accuracy: 0.9473684210526315

3.5 XGBoost Classifier

```

1 from xgboost import XGBClassifier
2 model_xgb = XGBClassifier(
3     use_label_encoder=False,
4     eval_metric='logloss',
5     random_state=42
6 )
7 st = time.time()
8 model_xgb.fit(X_train, y_train)
9 et = time.time()
10 training_time_xgb = et - st
11 print(f"Training time for XGBoost: {training_time_xgb} seconds")

```

Listing 8: Model Training

Training time for XGBoost: 1.8101122379302979 seconds

```

1 param_grid = {
2     'n_estimators': [100, 200, 300],
3     'learning_rate': [0.01, 0.05, 0.1, 0.2],
4     'max_depth': [3, 4, 5, 6],
5     'gamma': [0, 0.1, 0.3, 0.5],
6     'subsample': [0.8, 1.0],
7     'colsample_bytree': [0.8, 1.0]
8 }
9 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
10 grid_search_xgb = GridSearchCV(
11     estimator=model_xgb,
12     param_grid=param_grid,
13     cv=cv,
14     scoring={'accuracy': 'accuracy', 'f1': make_scorer(f1_score, average='weighted')},
15     refit='accuracy',
16     n_jobs=-1
17 )
18 # Measure training time
19 start_time = time.time()
20 grid_search_xgb.fit(X_train, y_train)
21 end_time = time.time()
22
23 # Results
24 print(f"Best Parameters: {grid_search_xgb.best_params_}")
25 print(f"Best CV Score: {grid_search_xgb.best_score_:.4f}")
26 print(f"Training Time: {end_time - start_time:.4f} seconds")

```

Listing 9: Hyperparameter Tuning

Best Parameters: {'colsample_bytree': 0.8, 'gamma': 0, 'learning_rate': 0.2, 'max_depth': 4, 'n_estimators': 200, 'subsample': 1.0}

Best CV Score: 0.9707

Training Time: 21.4219 seconds

```

1 model_xgb = grid_search_xgb.best_estimator_
2 y_pred_xgb = model_xgb.predict(X_test)
3

```



```

4 print("XGBoost Classifier Metrics:")
5 print(classification_report(y_test, y_pred_xgb))
6 print(f"Accuracy: {accuracy_score(y_test, y_pred_xgb)}")
7 print(f"f1_score: ", f1_score(y_test, y_pred_xgb, average='weighted'))
8 cm = confusion_matrix(y_test, y_pred_xgb)
9 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=lbs)
10 disp.plot(cmap='Blues')
11 plt.show()
12
13 y_proba = model_xgb.predict_proba(X_test)[: , 1]
14 fpr, tpr, thresholds = roc_curve(y_test, y_proba)
15 roc_auc = auc(fpr, tpr)
16
17 # Plot ROC curve
18 plt.figure(figsize=(6, 5))
19 plt.plot(fpr, tpr, color='blue', lw=2, label=f'ROC curve (AUC = {roc_auc:.4f})')
20 plt.plot([0, 1], [0, 1], color='gray', linestyle='--', lw=2)
21 plt.xlim([0.0, 1.0])
22 plt.ylim([0.0, 1.05])
23 plt.xlabel('False Positive Rate')
24 plt.ylabel('True Positive Rate')
25 plt.title('ROC Curve - XGBoost Classifier')
26 plt.legend(loc='lower right')
27 plt.grid(True)
28 plt.show()

```

Listing 10: Evaluation

XGBoost Classifier Metrics:

	precision	recall	f1-score	support
0	0.97	0.96	0.97	76
1	0.92	0.95	0.94	38
accuracy			0.96	114
macro avg	0.95	0.95	0.95	114
weighted avg	0.96	0.96	0.96	114

Accuracy: 0.956140350877193

f1_score: 0.9562799231673403

3.6 Random Forest Classifier

```

1 from sklearn.ensemble import RandomForestClassifier
2 model_rf = RandomForestClassifier(random_state=42)
3 st = time.time()
4 model_rf.fit(X_train, y_train)
5 et = time.time()
6 training_time_rf = et - st
7 print(f"Training time for Random Forest: {training_time_rf} seconds")

```

Listing 11: Model Training

Training time for Random Forest: 0.14823627471923828 seconds

```

1 param_grid = {
2     'n_estimators': [100, 200, 300],
3     'max_depth': [None, 5, 10, 15, 20],
4     'criterion': ['gini', 'entropy'],
5     'max_features': ['sqrt', 'log2'],
6     'min_samples_split': [2, 5, 10]
7 }
8 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
9 grid_search_rf = GridSearchCV(

```

```

10 estimator=model_rf,
11 param_grid=param_grid,
12 cv=cv,
13 scoring={'accuracy': 'accuracy', 'f1': make_scorer(f1_score, average='weighted')},
14 refit='accuracy',
15 n_jobs=-1
16 )
17 # Measure training time
18 start_time = time.time()
19 grid_search_rf.fit(X_train, y_train)
20 end_time = time.time()
21
22 # Results
23 print(f"Best Parameters: {grid_search_rf.best_params_}")
24 print(f"Best CV Score: {grid_search_rf.best_score_:.4f}")
25 print(f"Training Time: {end_time - start_time:.4f} seconds")

```

Listing 12: Hyperparameter Tuning

Best Parameters: {'criterion': 'entropy', 'max_depth': None,
 'max_features': 'sqrt', 'min_samples_split': 5, 'n_estimators': 200}
 Best CV Score: 0.9559
 Training Time: 12.2577 seconds

```

1 model_rf = grid_search_rf.best_estimator_
2 y_pred_rf = model_rf.predict(X_test)
3
4 print("Random Forest Classifier Metrics:")
5 print(classification_report(y_test, y_pred_rf))
6 print(f"Accuracy: {accuracy_score(y_test, y_pred_rf)}")
7 cm = confusion_matrix(y_test, y_pred_rf)
8 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=lbs)
9 disp.plot(cmap='Blues')
10 plt.show()
11
12 y_proba_rf = model_rf.predict_proba(X_test)[:, 1]
13 fpr_rf, tpr_rf, thresholds_rf = roc_curve(y_test, y_proba_rf)
14 roc_auc_rf = auc(fpr_rf, tpr_rf)
15
16 # Plot ROC curve
17 plt.figure(figsize=(6, 5))
18 plt.plot(fpr_rf, tpr_rf, color='blue', lw=2, label=f'ROC curve (AUC = {roc_auc_rf:.4f})')
19 plt.plot([0, 1], [0, 1], color='gray', linestyle='--', lw=2)
20 plt.xlim([0.0, 1.0])
21 plt.ylim([0.0, 1.05])
22 plt.xlabel('False Positive Rate')
23 plt.ylabel('True Positive Rate')
24 plt.title('ROC Curve - Random Forest Classifier')
25 plt.legend(loc='lower right')
26 plt.grid(True)
27 plt.show()

```

Listing 13: Evaluation

Random Forest Classifier Metrics:

	precision	recall	f1-score	support
0	0.96	1.00	0.98	76
1	1.00	0.92	0.96	38
accuracy			0.97	114
macro avg	0.98	0.96	0.97	114
weighted avg	0.97	0.97	0.97	114

Accuracy: 0.9736842105263158

3.7 Stacking Classifier (SVM + Naive Bayes + Decision Tree)

```

1  from sklearn.ensemble import StackingClassifier, RandomForestClassifier
2  from sklearn.svm import SVC
3  from sklearn.naive_bayes import GaussianNB
4  from sklearn.tree import DecisionTreeClassifier
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.neighbors import KNeighborsClassifier
7  from sklearn.preprocessing import StandardScaler
8  from sklearn.pipeline import make_pipeline
9  from sklearn.metrics import (
10     accuracy_score,
11     f1_score,
12     classification_report,
13     confusion_matrix,
14     roc_curve,
15     auc
16 )
17
18 # -----
19 # DEFINE STACKED MODELS
20 # -----
21
22 # 1) SVM, Naive Bayes, Decision Tree      Logistic Regression
23 model1 = StackingClassifier(
24     estimators=[
25         ("svm", make_pipeline(StandardScaler(), SVC(probability=True))),
26         ("nb", GaussianNB()),
27         ("dt", DecisionTreeClassifier(random_state=42))
28     ],
29     final_estimator=LogisticRegression(),
30     cv=5,
31     n_jobs=-1
32 )
33
34 # 2) SVM, Naive Bayes, Decision Tree      Random Forest
35 model2 = StackingClassifier(
36     estimators=[
37         ("svm", make_pipeline(StandardScaler(), SVC(probability=True))),
38         ("nb", GaussianNB()),
39         ("dt", DecisionTreeClassifier(random_state=42))
40     ],
41     final_estimator=RandomForestClassifier(random_state=42),
42     cv=5,
43     n_jobs=-1
44 )
45
46 # 3) SVM, Decision Tree, KNN      Logistic Regression
47 model3 = StackingClassifier(
48     estimators=[
49         ("svm", make_pipeline(StandardScaler(), SVC(probability=True))),
50         ("dt", DecisionTreeClassifier(random_state=42)),
51         ("knn", KNeighborsClassifier())
52     ],
53     final_estimator=LogisticRegression(),
54     cv=5,
55     n_jobs=-1
56 )
57
58 param_grid1 = {
59     "svm__svc__C": [0.1, 1, 10],
60     "svm__svc__kernel": ["linear", "rbf"],
61     "dt__max_depth": [None, 5, 10, 20],
62     "dt__criterion": ["gini", "entropy", "log_loss"],
63     "final_estimator__C": [0.1, 1, 10],
64     "final_estimator__solver": ["lbfgs", "liblinear"]
65 }
66 param_grid2 = {
67     "svm__svc__C": [0.1, 1, 10],
68     "svm__svc__kernel": ["linear", "rbf"],
69     "dt__max_depth": [None, 5, 10, 20],
70     "dt__criterion": ["gini", "entropy", "log_loss"],
71     "final_estimator__n_estimators": [100, 200, 300],

```

```

72     "final_estimator__max_depth": [None, 10, 20],
73     "final_estimator__min_samples_split": [2, 5, 10]
74 }
75 param_grid3 = {
76     "svm__svc__C": [0.1, 1, 10],
77     "svm__svc__kernel": ["linear", "rbf"],
78     "dt__max_depth": [None, 5, 10, 20],
79     "knn__n_neighbors": [3, 5, 7, 9],
80     "knn__weights": ["uniform", "distance"],
81     "final_estimator__C": [0.1, 1, 10],
82     "final_estimator__solver": ["lbfgs", "liblinear"]
83 }
84 param_grids = [param_grid1, param_grid2, param_grid3]
85
86 # -----
87 # TRAIN & EVALUATE MODELS
88 # -----
89
90 models = {
91     "SVM+NB+DT": LogisticRegression": model1,
92     "SVM+NB+DT": RandomForest": model2,
93     "SVM+DT+KNN": LogisticRegression": model3
94 }
95 results = []
96
97 plt.figure(figsize=(8, 6))
98
99 grid_no=0
100
101 for name, model in models.items():
102     # Train
103     model.fit(X_train, y_train)
104
105     #hyperparameter tuning
106
107     cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
108
109     grid_search_stack = GridSearchCV(
110         estimator=model,
111         param_grid=param_grids[grid_no],
112         cv=cv,
113         scoring={'accuracy': 'accuracy', 'f1': make_scorer(f1_score, average='weighted')},
114         refit='accuracy',
115         n_jobs=-1
116     )
117     grid_search_stack.fit(X_train, y_train)
118     print(f"Best Parameters for {name}: {grid_search_stack.best_params_}")
119     print(f"Best CV Score for {name}: {grid_search_stack.best_score_:.4f}")
120
121     #k-fold cross validation
122     cv_scores = cross_val_score(grid_search_stack.best_estimator_, X_val, y_val, cv=cv,
123                                 scoring='accuracy')
124     print(f"\nCross-Validation Scores for {name}:")
125     for i, score in enumerate(cv_scores, 1):
126         print(f"Fold {i} Accuracy: {score:.4f}")
127     print(f"Average Accuracy: {np.mean(cv_scores):.4f}")
128
129     grid_no+=1
130
131     model = grid_search_stack.best_estimator_
132     # Predictions
133     y_pred = model.predict(X_test)
134     y_proba = model.predict_proba(X_test)[:, 1] # For ROC curve
135
136     # Accuracy and F1
137     acc = accuracy_score(y_test, y_pred)
138     f1 = f1_score(y_test, y_pred)
139     results.append({"Model": name, "Accuracy": acc, "F1 Score": f1})
140
141     # ROC Curve
142     fpr, tpr, _ = roc_curve(y_test, y_proba)
143     roc_auc = auc(fpr, tpr)
144     plt.plot(fpr, tpr, lw=2, label=f"{name} (AUC = {roc_auc:.2f})")

```

```

144     # Reports
145     print(f"\n== {name} ==")
146     print(classification_report(y_test, y_pred))
147     print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
148
149     # -----
150     # ROC PLOT
151     # -----
152
153     plt.plot([0, 1], [0, 1], linestyle="--", color="gray")
154     plt.xlabel("False Positive Rate")
155     plt.ylabel("True Positive Rate")
156     plt.title("ROC Curve Comparison of Stacked Ensemble Models")
157     plt.legend(loc="lower right")
158     plt.grid(True)
159     plt.tight_layout()
160     plt.show()
161
162     # -----
163     # RESULTS TABLE
164     # -----
165
166     results_df = pd.DataFrame(results)
167     print("\nFinal Comparison Table:\n")
168     print(results_df)
169 
```

Listing 14: Model Building and Evaluation

Output:

```
Best Parameters for SVM+NB+DT → LogisticRegression: {'dt__criterion': 'gini', 'dt__max_dep
Best CV Score for SVM+NB+DT → LogisticRegression: 0.9765
```

Cross-Validation Scores for SVM+NB+DT → LogisticRegression:

```
Fold 1 Accuracy: 0.9565
Fold 2 Accuracy: 0.9565
Fold 3 Accuracy: 0.9565
Fold 4 Accuracy: 0.9565
Fold 5 Accuracy: 0.9545
Average Accuracy: 0.9561
```

SVM+NB+DT → LogisticRegression		precision	recall	f1-score	support
0		0.99	0.97	0.98	76
1		0.95	0.97	0.96	38
accuracy				0.97	114
macro avg		0.97	0.97	0.97	114
weighted avg		0.97	0.97	0.97	114

Confusion Matrix:

[[74 2]
[1 37]]

```
Best Parameters for SVM+NB+DT → RandomForest: {'dt__criterion': 'entropy', 'dt__max_depth': 10}
Best CV Score for SVM+NB+DT → RandomForest: 0.9765
```

Cross-Validation Scores for SVM+NB+DT → RandomForest:

Fold 1 Accuracy: 0.9130

Fold 2 Accuracy: 0.9565
 Fold 3 Accuracy: 0.9565
 Fold 4 Accuracy: 0.9130
 Fold 5 Accuracy: 0.9545
 Average Accuracy: 0.9387

=== SVM+NB+DT → RandomForest ===

	precision	recall	f1-score	support
0	0.99	0.99	0.99	76
1	0.97	0.97	0.97	38
accuracy			0.98	114
macro avg	0.98	0.98	0.98	114
weighted avg	0.98	0.98	0.98	114

Confusion Matrix:

```
[[75  1]
 [ 1 37]]
```

Best Parameters for SVM+DT+KNN → LogisticRegression: {'dt__max_depth': None, 'final_estimator__max_depth': 3, 'final_estimator__min_samples_split': 2, 'final_estimator__min_samples_leaf': 1, 'final_estimator__min_samples_weighted': 1, 'final_estimator__min_weight_fraction': 0.1, 'final_estimator__n_estimators': 100, 'final_estimator__random_state': 42, 'final_estimator__tol': 0.0001, 'final_estimator__verbose': 0, 'final_estimator__warm_start': False}

Best CV Score for SVM+DT+KNN → LogisticRegression: 0.9794

Cross-Validation Scores for SVM+DT+KNN → LogisticRegression:

Fold 1 Accuracy: 0.9565
 Fold 2 Accuracy: 1.0000
 Fold 3 Accuracy: 0.9130
 Fold 4 Accuracy: 1.0000
 Fold 5 Accuracy: 1.0000
 Average Accuracy: 0.9739

=== SVM+DT+KNN → LogisticRegression ===

	precision	recall	f1-score	support
0	1.00	0.99	0.99	76
1	0.97	1.00	0.99	38
accuracy			0.99	114
macro avg	0.99	0.99	0.99	114
weighted avg	0.99	0.99	0.99	114

Confusion Matrix:

```
[[75  1]
 [ 0 38]]
```

4 Confusion Matrix and ROC for Each Model

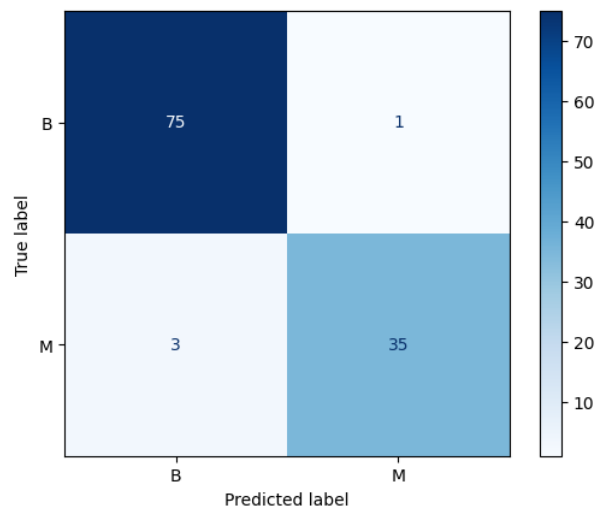


Figure 2: Decision Tree Confusion Matrix

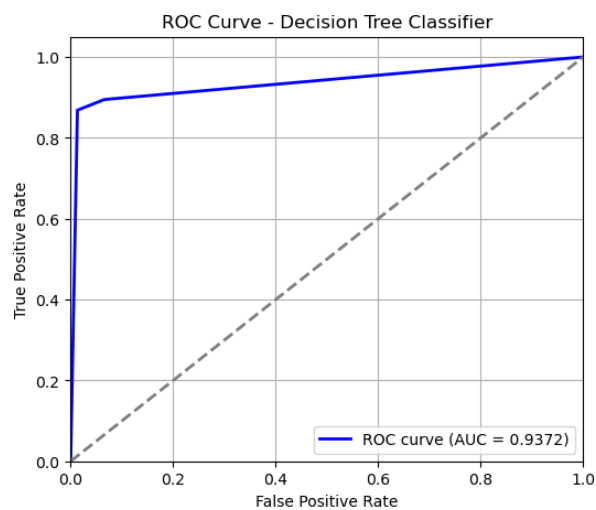


Figure 3: Decision Tree ROC Curve

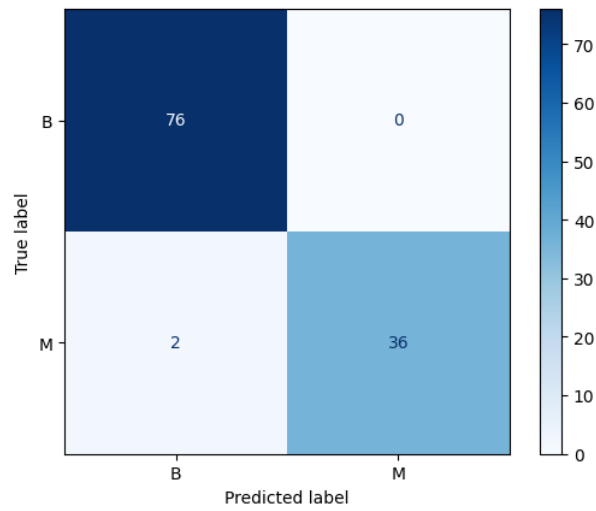


Figure 4: AdaBoost Classifier Confusion Matrix

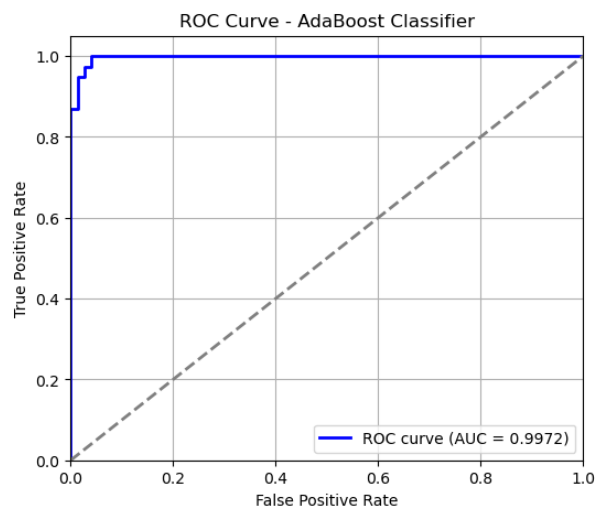


Figure 5: AdaBoost ROC Curve

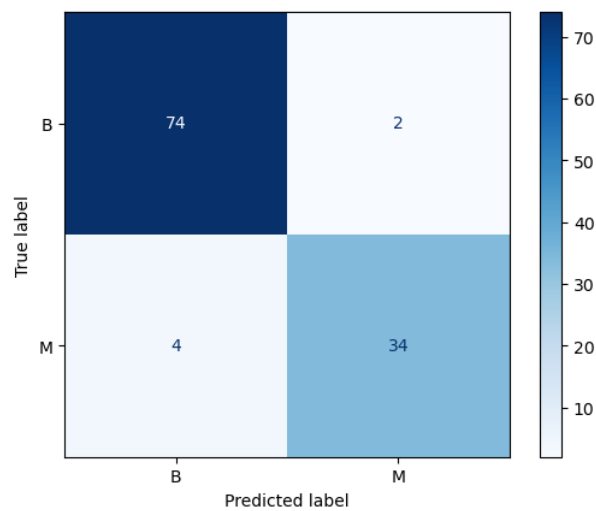


Figure 6: Gradient Boosting Confusion Matrix

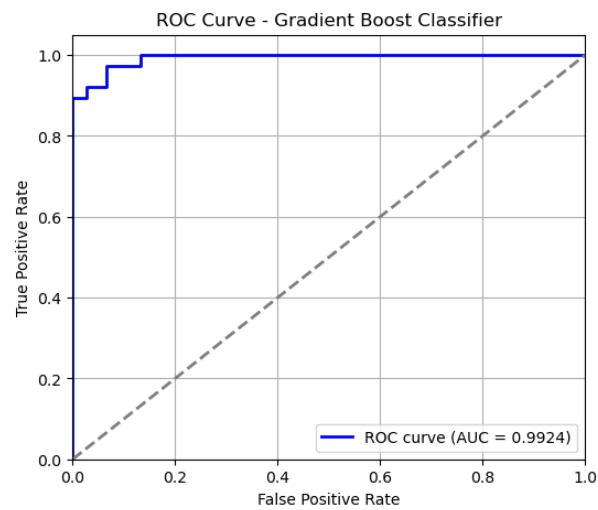


Figure 7: Gradient Boosting ROC curve

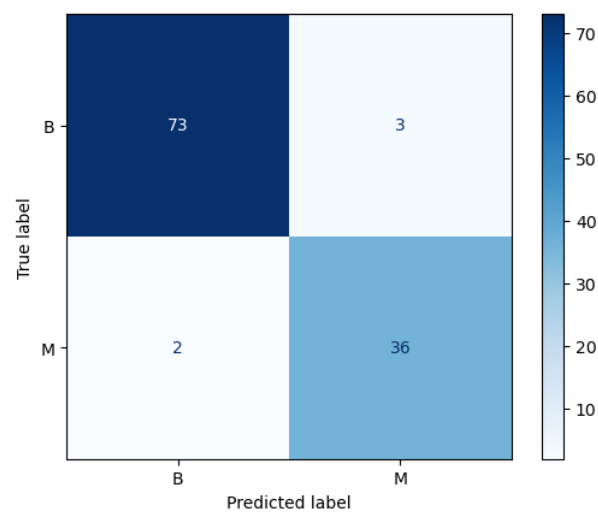


Figure 8: XGBoost Classifier Confusion Matrix

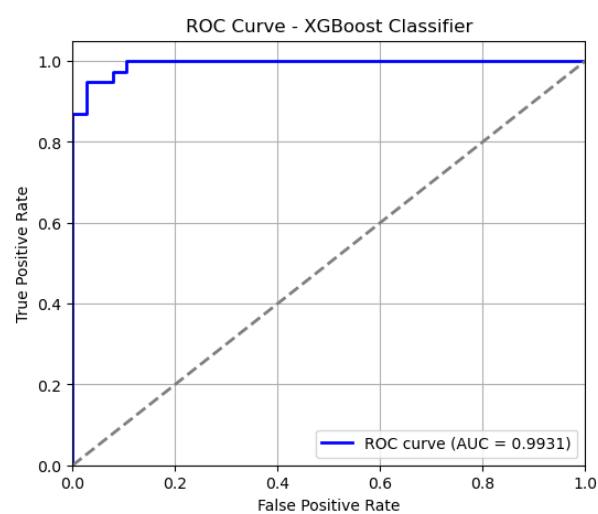


Figure 9: XGBoost Classifier ROC Curve

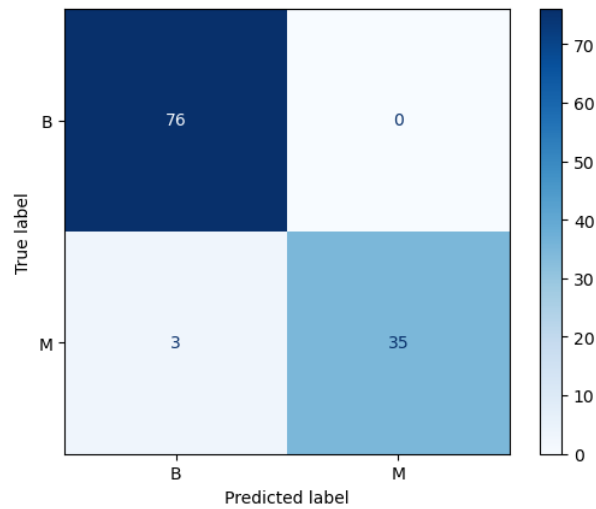


Figure 10: Random Forest Classifier Confusion Matrix

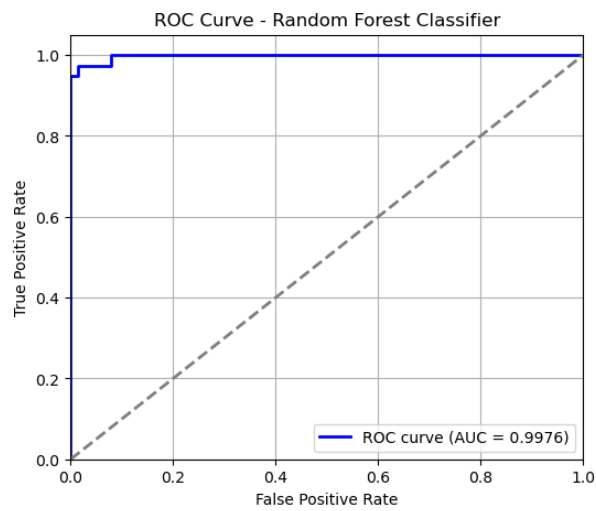
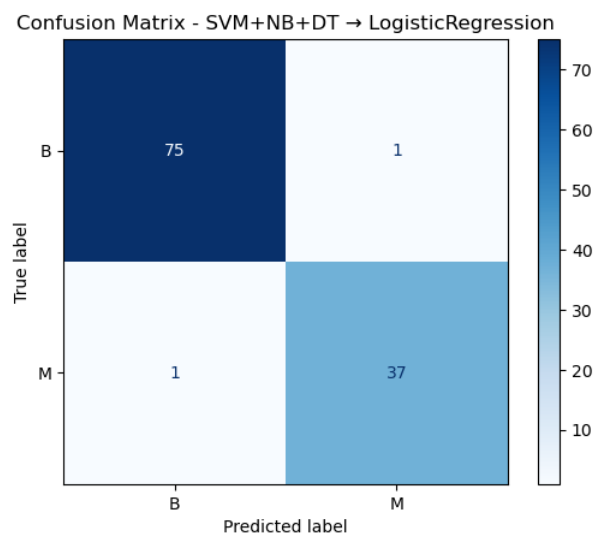


Figure 11: Random Forest Classifier ROC Curve



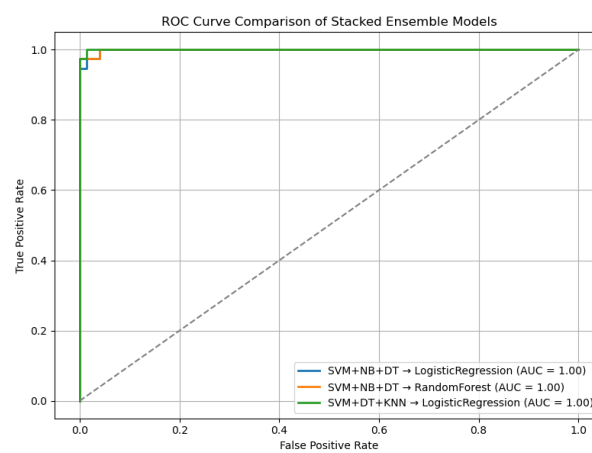
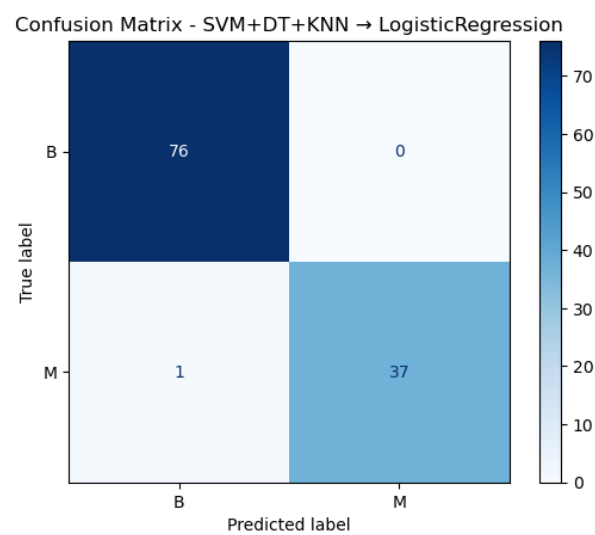
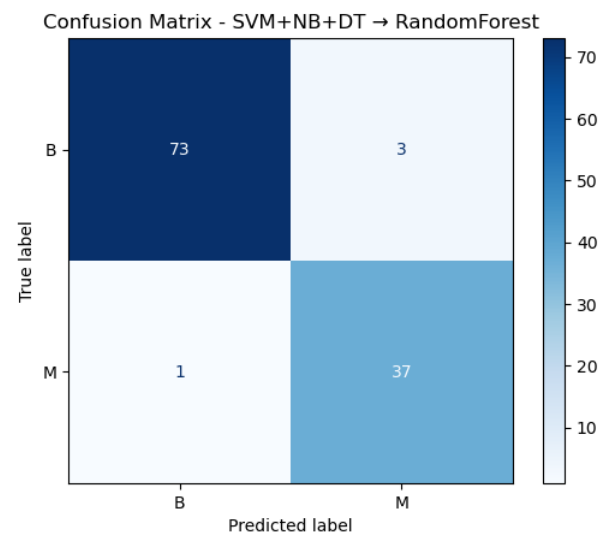


Figure 12: Stacking Classifier ROC Curve

5 Hyperparameter Tuning Tables

5.1 Decision Tree Classifier

Table 1: Decision Tree - Hyperparameter Tuning Results

Criterion	Max Depth	Accuracy	F1 Score
gini	None	0.921	0.920
gini	5	0.924	0.923
gini	10	0.912	0.911
gini	15	0.912	0.911
gini	20	0.915	0.914
entropy	None	0.947	0.947
entropy	5	0.950	0.950
entropy	10	0.944	0.944
entropy	15	0.944	0.944
entropy	20	0.944	0.944
log_loss	None	0.947	0.947
log_loss	5	0.947	0.947
log_loss	10	0.944	0.944
log_loss	15	0.950	0.950
log_loss	20	0.947	0.947

5.2 AdaBoost Classifier

n_estimators	learning_rate	Accuracy	F1 Score
200	1	0.9736	0.9734
300	1	0.9706	0.9703
300	0.5	0.9677	0.9673
200	0.5	0.9677	0.9673
300	1	0.9648	0.9646

5.3 Gradient Boosting Classifier

n_estimators	learning_rate	max_depth	Accuracy	F1 Score
300	0.2	3	0.9736	0.9735
200	0.2	3	0.9736	0.9735
300	0.2	3	0.9736	0.9735
200	0.1	5	0.9736	0.9733
200	0.2	3	0.9736	0.9735

5.4 XGBoost Classifier

n_estimators	learning_rate	max_depth	gamma	Accuracy	F1 Score
200	0.2	4	0	0.9707	0.9705
100	0.2	5	0.1	0.9707	0.9705
200	0.1	5	0.1	0.9707	0.9706
300	0.2	4	0	0.9707	0.9705
200	0.1	6	0.1	0.9707	0.9706

5.5 Random Forest Classifier

n_estimators	max_depth	criterion	Accuracy	F1 Score
200	20	entropy	0.9559	0.9556
200	15	entropy	0.9559	0.9556
200	20	entropy	0.9559	0.9556
200	15	entropy	0.9559	0.9556
200	None	entropy	0.9559	0.9556

5.6 Stacking Ensemble Models

Base-Model	Final-Estimator	Accuracy	F1 Score
SVM+NB+DT	LogisticRegression	0.982456	0.973684
SVM+NB+DT	RandomForest	0.973684	0.961039
SVM+DT+KNN	LogisticRegression	0.991228	0.986667

Table 2: Stacking model Hyperparameter Tuning

6 Cross-Validation Results Table

Model	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	AvG Accuracy
Decision Tree	0.9130	0.9130	0.8696	0.9565	1.0000	0.9304
AdaBoost	0.8696	1.0000	0.8696	0.9565	0.9545	0.9300
Gradient Boosting	0.9130	0.9565	0.9130	1.0000	0.9545	0.9474
XGBoost	0.9130	0.9565	0.8696	1.0000	1.0000	0.9478
Random Forest	0.9130	0.9130	0.9130	1.0000	1.0000	0.9478
SVM+NB+DT → Logistic	0.9565	0.9565	0.9565	0.9565	0.9545	0.9561
SVM+NB+DT → RandomForest	0.9130	0.9565	0.9565	0.9130	0.9545	0.9387
SVM+DT+KNN → Logistic	0.9565	1.0000	0.9130	1.0000	1.0000	0.9739

Table 3: 5-Fold Cross Validation Results for All Models

7 All Comparison Table

Model	Accuracy	F1 Score
Decision Tree	0.9474	0.9465
AdaBoost	0.9736	0.9739
Gradient Boosting	0.9474	0.9469
XGBoost	0.9561	0.9563
Random Forest	0.9737	0.9733
Stacking (SVM + NB + DT → LR)	0.973684	0.961039
Stacking (SVM + NB + DT → RF)	0.982456	0.973684
Stacking (SVM + DT + KNN → LR)	0.991228	0.987013

Table 4: Comparison of models based on Accuracy and F1 Score

8 Observations and Conclusions

- **Best Validation Accuracy:** The stacking model ($SVM + DT + KNN \rightarrow LR$) achieved the highest validation accuracy of **0.9912**, outperforming all other methods.
- **Decision Tree vs. Ensemble Methods:** The standalone Decision Tree achieved an accuracy of **0.9474**, which is significantly lower compared to ensemble methods like Random Forest (**0.9737**) and AdaBoost (**0.9736**). This highlights the improved stability and predictive power of ensemble techniques.
- **Random Forest Tuning:** Random Forest reached an accuracy of **0.9737**, suggesting that tuning hyperparameters such as `max_depth` and `n_estimators` likely contributed to its improved performance compared to the plain Decision Tree.
- **Generalization and Overfitting:** Among the models, stacking with Logistic Regression and Random Forest showed strong generalization, with balanced accuracy and F1 scores. No clear signs of overfitting were observed, as performance across metrics remained consistent. However, standalone Decision Trees exhibited relatively weaker generalization.
- **Effect of Stacking:** Stacking consistently improved performance over base learners. For example, stacking with $SVM + DT + KNN \rightarrow LR$ achieved the best overall results (**Accuracy = 0.9912, F1 = 0.9870**), surpassing both individual models and other ensemble methods.

Conclusion: Ensemble methods, particularly stacking, provided superior performance compared to individual classifiers. While Decision Trees served as a useful baseline, advanced ensembles like Random Forest, AdaBoost, and Gradient Boosting offered significant improvements. Stacking further enhanced predictive performance, making it the most effective approach for this dataset.