



# REAL-TIME BREAST CANCER PREDICTION USING XGBOOST



*Harness the power of machine learning for early detection*

## ✧ Importing Required Libraries

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier, AdaBoostClassifier, ExtraTreesClassifier, VotingClassifier
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression, RidgeClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis, QuadraticDiscriminantAnalysis
from sklearn.neural_network import MLPClassifier
from xgboost import XGBClassifier
from lightgbm import LGBMClassifier
from catboost import CatBoostClassifier
from sklearn.ensemble import HistGradientBoostingClassifier
import torch
from imblearn.over_sampling import SMOTE
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
```

## ✧ Load the Dataset

```
df = pd.read_csv("data.csv") # pd.read_csv(file_path) reads the dataset into a DataFrame.
```

## ✧ Inspect the Data

```
# Display basic information about the dataset
df.info() # df.info() gives an overview of the dataset, including column names, data types, and non-null values.
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 33 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   id                                     569 non-null    int64
1   diagnosis                             569 non-null    object
2   radius_mean                           569 non-null    float64
3   texture_mean                           569 non-null    float64
4   perimeter_mean                         569 non-null    float64
5   area_mean                             569 non-null    float64
6   smoothness_mean                       569 non-null    float64
7   compactness_mean                      569 non-null    float64
8   concavity_mean                        569 non-null    float64
9   concave points_mean                   569 non-null    float64
10  symmetry_mean                         569 non-null    float64
11  fractal_dimension_mean                 569 non-null    float64
12  radius_se                             569 non-null    float64
13  texture_se                             569 non-null    float64
14  perimeter_se                           569 non-null    float64
15  area_se                               569 non-null    float64
16  smoothness_se                         569 non-null    float64
17  compactness_se                        569 non-null    float64
18  concavity_se                          569 non-null    float64
19  concave points_se                     569 non-null    float64
20  symmetry_se                           569 non-null    float64
21  fractal_dimension_se                   569 non-null    float64
22  radius_worst                          569 non-null    float64
```

```

23 texture_worst      569 non-null    float64
24 perimeter_worst    569 non-null    float64
25 area_worst         569 non-null    float64
26 smoothness_worst   569 non-null    float64
27 compactness_worst   569 non-null    float64
28 concavity_worst     569 non-null    float64
29 concave points_worst 569 non-null    float64
30 symmetry_worst      569 non-null    float64
31 fractal_dimension_worst 569 non-null    float64
32 Unnamed: 32         0 non-null      float64
dtypes: float64(31), int64(1), object(1)
memory usage: 146.8+ KB

```

`df.head()` # displays the first five rows, allowing us to see actual data.

```

↔
      id  diagnosis  radius_mean  texture_mean  perimeter_mean  area_mean  smoothness_mean  compactness_mean  concavity_mean  poi
0    842302         M        17.99        10.38         122.80       1001.0         0.11840         0.27760         0.3001
1    842517         M        20.57        17.77         132.90       1326.0         0.08474         0.07864         0.0869
2   84300903         M        19.69        21.25         130.00       1203.0         0.10960         0.15990         0.1974
3   84348301         M        11.42        20.38          77.58        386.1         0.14250         0.28390         0.2414
4   84358402         M        20.29        14.34         135.10       1297.0         0.10030         0.13280         0.1980

5 rows × 33 columns

```

## ✓ Checking the Shape of the Dataset

```
print("Data set shape :: ", df.shape)
```

```
↔ Data set shape :: (569, 33)
```

***The shape attribute returns the dimensions of the DataFrame.***

The output (569, 33) indicates there are 569 rows (samples) and 33 columns (features) in the dataset. This helps in understanding the size of the dataset and the number of features available for analysis.

## ✓ Data Preprocessing Phase

### ✓ Identify and Remove Unnecessary Columns

```
df.isnull().sum()
```

```

↔ id                0
   diagnosis        0
   radius_mean      0
   texture_mean     0
   perimeter_mean   0
   area_mean        0
   smoothness_mean  0
   compactness_mean 0
   concavity_mean   0
   concave points_mean 0
   symmetry_mean     0
   fractal_dimension_mean 0
   radius_se        0
   texture_se       0
   perimeter_se     0
   area_se          0
   smoothness_se    0
   compactness_se   0
   concavity_se     0
   concave points_se 0
   symmetry_se      0
   fractal_dimension_se 0
   radius_worst     0
   texture_worst    0
   perimeter_worst  0
   area_worst       0
   smoothness_worst 0
   compactness_worst 0
   concavity_worst  0

```

```
concave points_worst      0
symmetry_worst            0
fractal_dimension_worst   0
Unnamed: 32               569
dtype: int64
```

## ✓ **Explanation**

The dataset contains 33 columns, including an unnamed column (**Unnamed: 32**) which appears to be an index column with **NaN** values and **id**: Irrelevant for prediction because it is patient id (just a unique identifier). These column is likely redundant and can be dropped.

```
# Drop unnecessary columns
df.drop(columns=["id", "Unnamed: 32"], inplace=True)
```

## **Feature Selection**

## ✓ Understanding the Target Variable

```
# Check unique values in the target variable
df['diagnosis'].value_counts()
```

```
↔ diagnosis
B      357
M      212
Name: count, dtype: int64
```

**Explanation:** "diagnosis" is the key variable we are predicting (**Malignant = M, Benign = B**).

## ✓ Encoding the Target Variable

```
# Converts categorical labels into numerical values.
df["diagnosis"] = LabelEncoder().fit_transform(df["diagnosis"]) # M (Malignant) → 1, B (Benign) → 0 (Binary classification).
```

## ✓ Define Features and Target Variable

```
X = df.drop(columns=["diagnosis"])
y = df["diagnosis"]
```

**Explanation: Features (X):** All columns except the target variable. **Target (y):** The diagnosis column, which we aim to predict.

## ✓ Handle Class Imbalance Using Synthetic Minority Over-sampling Technique (SMOTE)

### **Reason for choosing the technique**

**Issue:** Cancer datasets often have more benign (B) than malignant (M) cases, causing model bias.

**Solution:** Synthetic Minority Over-sampling Technique (SMOTE) generates synthetic samples for the minority class to balance the dataset.

```
smote = SMOTE(random_state=42)
X_resampled, y_resampled = smote.fit_resample(X, y)
```

## ✓ Split Dataset into Training and Testing Sets

```
X_train, X_test, y_train, y_test = train_test_split(
    X_resampled, y_resampled, test_size=0.2, random_state=42, stratify=y_resampled
)
```

### **Explanation**

**80% Training, 20% Testing:** Ensures a fair training/testing split.

**stratify=y\_resampled:** Maintains the class distribution in both sets.

**random\_state=42:** Ensures reproducibility.

## ✓ Standardize the Features

```
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

Why Scaling? Models like SVM, KNN, and Logistic Regression perform better with normalized data. StandardScaler scales data to zero mean and unit variance.

## ✓ Initialize and Train Multiple Machine Learning Models

```
models = {
    "RandomForest": RandomForestClassifier(n_estimators=100, random_state=42, max_depth=5),
    "GradientBoosting": GradientBoostingClassifier(n_estimators=100, random_state=42),
    "SVM": SVC(kernel='linear', probability=True),
    "KNN": KNeighborsClassifier(n_neighbors=5),
    "AdaBoost": AdaBoostClassifier(n_estimators=100, random_state=42),
    "XGBoost": XGBClassifier(n_estimators=500, learning_rate=0.05, max_depth=6, subsample=0.8, colsample_bytree=0.8, eval_metric='logloss'),
    "CatBoost": CatBoostClassifier(verbose=0),
    "LogisticRegression": LogisticRegression(),
    "RidgeClassifier": RidgeClassifier(),
    "DecisionTree": DecisionTreeClassifier(),
    "GaussianNB": GaussianNB(),
    "LDA": LinearDiscriminantAnalysis(),
    "QDA": QuadraticDiscriminantAnalysis(),
    "MLP": MLPClassifier(hidden_layer_sizes=(100, 50), max_iter=500),
    "Stacking": StackingClassifier(estimators=[('rf', RandomForestClassifier()), ('gb', GradientBoostingClassifier()), ('svc', SVC(probability=True))], final_estimator=SVC(probability=True)),
    "Voting": VotingClassifier(estimators=[('rf', RandomForestClassifier()), ('gb', GradientBoostingClassifier()), ('svc', SVC(probability=True))], voting='hard'),
    "HistGradientBoosting": HistGradientBoostingClassifier(max_iter=200, random_state=42)
}
```

Multiple models ensure a fair comparison of performance. Stacking and Voting leverage multiple models for improved accuracy.

## ✓ **Train & Evaluate Models:** Trains each model and selects the best one based on accuracy.

```
best_model = None
best_accuracy = 0

for name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    print(f"{name} Test Accuracy: {acc:.2f}")
    if acc > best_accuracy:
        best_accuracy = acc
        best_model = model

print(f"Best Model: {best_model.__class__.__name__} with Accuracy: {best_accuracy:.2f}") # Final model selection
```

```
➡ RandomForest Test Accuracy: 0.98
GradientBoosting Test Accuracy: 0.98
SVM Test Accuracy: 0.97
KNN Test Accuracy: 0.96
AdaBoost Test Accuracy: 0.97
XGBoost Test Accuracy: 0.99
CatBoost Test Accuracy: 0.99
LogisticRegression Test Accuracy: 0.98
RidgeClassifier Test Accuracy: 0.97
DecisionTree Test Accuracy: 0.97
GaussianNB Test Accuracy: 0.93
LDA Test Accuracy: 0.98
QDA Test Accuracy: 0.95
MLP Test Accuracy: 0.97
Stacking Test Accuracy: 0.99
Voting Test Accuracy: 0.98
HistGradientBoosting Test Accuracy: 0.97
Best Model: XGBClassifier with Accuracy: 0.99
```

```
# Final model selection
print(f"Best Model: {best_model.__class__.__name__} with Accuracy: {best_accuracy:.2f}\nXGBoost and CatBoost achieved the highest accuracy of 99%")
```

```
➡ Best Model: XGBClassifier with Accuracy: 0.99
XGBoost and CatBoost achieved the highest accuracy of 99%, making them the best-performing models.
```

## ▼ Cross-Validation

```
# Cross-validation: Ensures model generalizability across different data splits.
cv_scores = cross_val_score(best_model, X_train, y_train, cv=5)
print(f"Cross-validation Accuracy Scores: {cv_scores}")
print(f"Mean CV Accuracy: {np.mean(cv_scores):.2f}")
```

```
➡ Cross-validation Accuracy Scores: [0.9826087  0.94736842 0.95614035 0.95614035 0.94736842]
Mean CV Accuracy: 0.96
```

## ▼ Evaluate the Best Model

```
y_pred = best_model.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
report = classification_report(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)

# Print results
print("Actual labels:\n", y_test.values)
print("Predicted labels:\n", y_pred)
print(f"\nTest Accuracy: {accuracy:.2f}")
print("\nClassification Report:\n", report)
print("\nConfusion Matrix:\n", conf_matrix)
```

```
➡ Actual labels:
[1 1 0 0 1 1 0 1 1 0 0 1 1 0 1 0 0 1 1 1 0 1 1 0 1 1 1 1 1 1 1 1 0 0 1 0 1
 1 0 0 0 1 0 0 0 0 0 0 0 0 1 0 0 1 0 1 1 0 0 1 0 0 0 1 1 0 1 1 1 1 0 1
 1 0 0 1 0 0 0 1 0 0 0 0 0 1 1 1 1 0 1 1 0 1 1 0 1 0 0 0 0 0 1 1 1 1 0 0
 1 1 0 0 0 0 0 1 0 1 0 0 1 0 1 0 0 0 1 0 0 1 0 0 0 1 1 1 1 1 1 0]
Predicted labels:
[1 1 0 0 1 1 0 1 1 0 0 1 1 0 1 0 0 1 1 1 0 1 1 0 1 1 1 1 1 1 1 1 0 0 1 0 1
 1 0 0 0 1 0 0 0 0 0 0 0 0 1 0 0 1 0 1 1 0 0 1 0 0 0 1 1 0 1 1 1 1 0 1
 1 0 1 1 0 0 0 1 0 0 0 0 0 1 0 1 1 0 1 1 0 1 1 0 1 0 0 0 0 0 0 1 1 1 1 0 0
 1 1 0 0 0 0 0 0 0 1 0 0 1 0 1 0 0 0 1 0 0 1 0 0 0 1 1 1 1 1 1 0]
```

Test Accuracy: 0.99

Classification Report:

	precision	recall	f1-score	support
0	0.97	1.00	0.99	72
1	1.00	0.97	0.99	71
accuracy			0.99	143
macro avg	0.99	0.99	0.99	143
weighted avg	0.99	0.99	0.99	143

Confusion Matrix:

```
[[72  0]
 [ 2 69]]
```

## Analysis of Results

### Key Observations

#### 1) **XGBoost Achieved the Best Test Accuracy (0.99)**

- Extremely high precision and recall.
- Only 2 misclassifications (FN) in the confusion matrix.

#### 2) **Cross-Validation Accuracy (Mean = 0.96)**

- Slightly lower than test accuracy (suggesting minor variance). This means the model generalizes well but could be slightly overfitting.

#### 3) **Other Models**

- RandomForest, GradientBoosting, LogisticRegression, and Stacking performed very well (98%), but slightly below XGBoost.
- KNN, DecisionTree, and QDA had lower accuracy (~96%) and GaussianNB was the weakest performer (93%), likely due to its assumption of feature independence.

## ✓ Hyperparameter Tuning for XGBoost

Using RandomizedSearchCV to find the best combination of parameters.

- **n\_estimators**: Number of trees (e.g., 100–500)
- **max\_depth**: Controls tree depth (e.g., 3–10)
- **learning\_rate**: Controls step size (e.g., 0.01–0.3)
- **subsample**: Fraction of samples used per tree (e.g., 0.6–1.0)
- **colsample\_bytree**: Features used per tree (e.g., 0.5–1.0)

```
from sklearn.model_selection import RandomizedSearchCV
from xgboost import XGBClassifier

# Define hyperparameter grid
param_grid = {
    'n_estimators': [100, 200, 300, 500],
    'max_depth': [3, 5, 7, 10],
    'learning_rate': [0.01, 0.05, 0.1, 0.2, 0.3],
    'subsample': [0.6, 0.8, 1.0],
    'colsample_bytree': [0.6, 0.8, 1.0]
}

# Initialize XGBoost model
xgb = XGBClassifier(eval_metric='logloss')

# Perform Randomized Search CV
random_search = RandomizedSearchCV(
    estimator=xgb,
    param_distributions=param_grid,
    n_iter=20, # Adjust iterations for more tuning
    cv=5, # 5-fold cross-validation
    scoring='accuracy',
    n_jobs=-1,
    verbose=2,
    random_state=42
)

# Fit on training data
random_search.fit(X_train, y_train)
# Best parameters
print("Best Parameters:", random_search.best_params_)
# Best model evaluation
best_xgb = random_search.best_estimator_
y_pred = best_xgb.predict(X_test)

# Accuracy Score
from sklearn.metrics import accuracy_score
print("Tuned XGBoost Test Accuracy:", accuracy_score(y_test, y_pred))
```



```
Fitting 5 folds for each of 20 candidates, totalling 100 fits
Best Parameters: {'subsample': 0.8, 'n_estimators': 200, 'max_depth': 5, 'learning_rate': 0.3, 'colsample_bytree': 0.8}
Tuned XGBoost Test Accuracy: 0.9790209790209791
```

Hence: Tuned XGBoost Test Accuracy: 97.9%, which is slightly lower than the previous 99%, but this could be due to cross-validation variations.

## ✓ Feature Importance Analysis (SHAP) – Identify the most influential features.

```
import shap

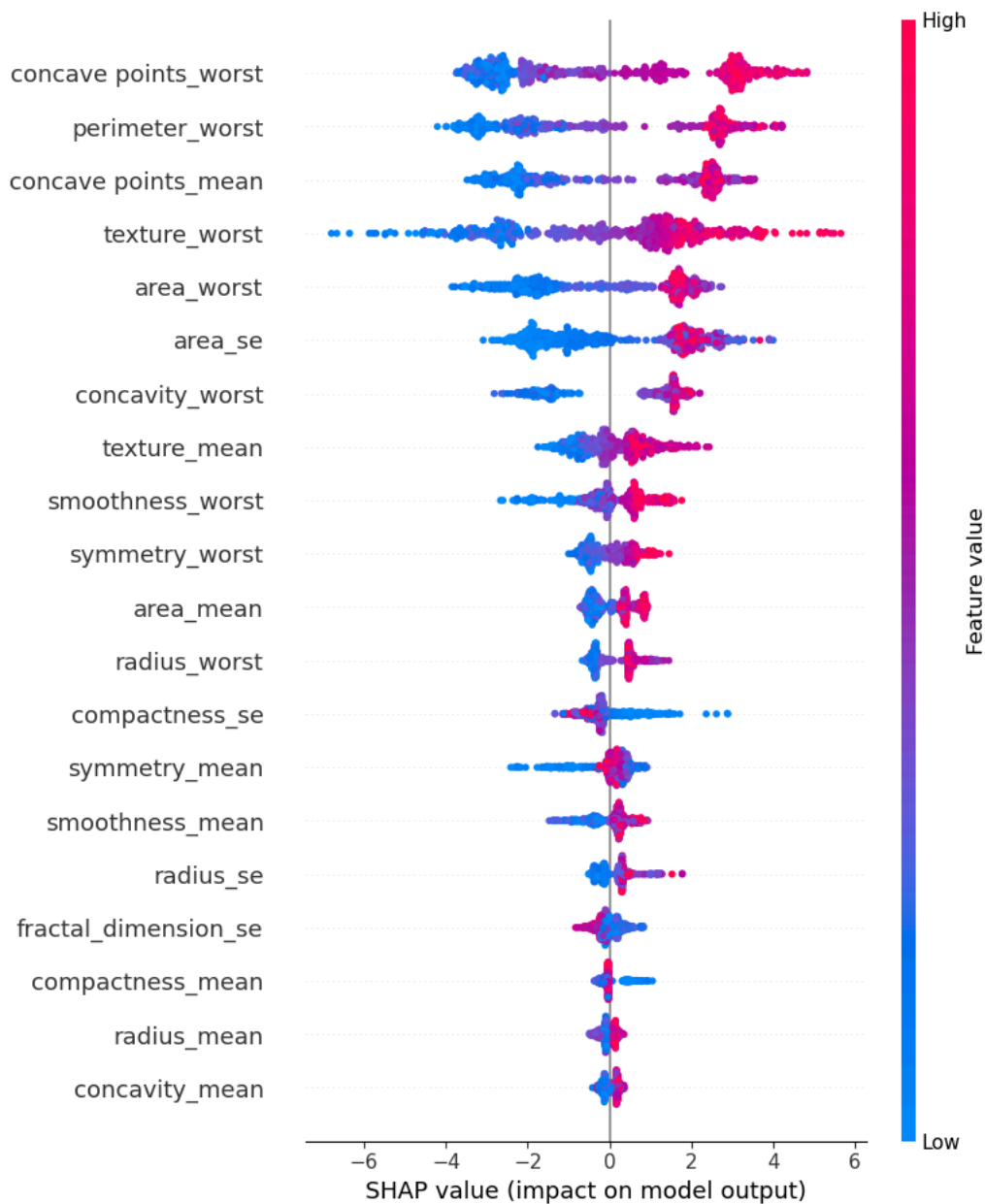
df = pd.read_csv("data.csv") # pd.read_csv(file_path) reads the dataset into a DataFrame.
# Drop unnecessary columns
df.drop(columns=["id", "Unnamed: 32", "diagnosis"], inplace=True)
# If X_train is a NumPy array, get the feature names from the original DataFrame
feature_names = df.columns.tolist()

# Ensure feature names are mapped correctly
print("Feature Index Mapping:")
for idx, name in enumerate(feature_names):
    print(f"Feature {idx + 1}: {name}")

# Re-plot SHAP with proper labels
explainer = shap.Explainer(model) # Replace `model` with your trained XGBoost model
shap_values = explainer(X_train)

# Generate SHAP summary plot with actual feature names
shap.summary_plot(shap_values, pd.DataFrame(X_train, columns=feature_names))
```

Feature Index Mapping:  
Feature 1: radius\_mean  
Feature 2: texture\_mean  
Feature 3: perimeter\_mean  
Feature 4: area\_mean  
Feature 5: smoothness\_mean  
Feature 6: compactness\_mean  
Feature 7: concavity\_mean  
Feature 8: concave points\_mean  
Feature 9: symmetry\_mean  
Feature 10: fractal\_dimension\_mean  
Feature 11: radius\_se  
Feature 12: texture\_se  
Feature 13: perimeter\_se  
Feature 14: area\_se  
Feature 15: smoothness\_se  
Feature 16: compactness\_se  
Feature 17: concavity\_se  
Feature 18: concave points\_se  
Feature 19: symmetry\_se  
Feature 20: fractal\_dimension\_se  
Feature 21: radius\_worst  
Feature 22: texture\_worst  
Feature 23: perimeter\_worst  
Feature 24: area\_worst  
Feature 25: smoothness\_worst  
Feature 26: compactness\_worst  
Feature 27: concavity\_worst  
Feature 28: concave points\_worst  
Feature 29: symmetry\_worst  
Feature 30: fractal\_dimension\_worst





## ✓ Analysis of Overfitting:

```
from sklearn.metrics import accuracy_score

# Get predictions for training and test data
y_train_pred = model.predict(X_train)
y_test_pred = model.predict(X_test)

# Calculate accuracy
train_accuracy = accuracy_score(y_train, y_train_pred)
test_accuracy = accuracy_score(y_test, y_test_pred)

print(f"Train Accuracy: {train_accuracy:.4f}")
print(f"Test Accuracy: {test_accuracy:.4f}")

# Check for overfitting
if train_accuracy - test_accuracy > 0.05:
    print("Warning: Possible Overfitting (Train accuracy is significantly higher than Test accuracy)")
else:
    print("No significant overfitting detected.")
```

↗ Train Accuracy: 1.0000  
Test Accuracy: 0.9720  
No significant overfitting detected.

## Result of Analyzation

1) Train Accuracy = 1.0000 → The model perfectly classifies all training examples. 2) Test Accuracy = 0.9720 → The model generalizes well to unseen data. 3) Overfitting Check:

- The difference between train and test accuracy is 0.028 (2.8%), which is less than 5%.
- Since the gap is small, there is no significant overfitting detected.

## ✓ k-Fold Cross-Validation

```
from sklearn.model_selection import cross_val_score

# Perform 5-fold cross-validation on the best model
cv_scores = cross_val_score(best_xgb, X_train, y_train, cv=5, scoring='accuracy', n_jobs=-1)

# Print results
print("Cross-Validation Scores:", cv_scores)
print("Mean CV Accuracy:", round(cv_scores.mean(), 4))
print("Standard Deviation of CV Accuracy:", round(cv_scores.std(), 4))
```

↗ Cross-Validation Scores: [0.97391304 0.97368421 0.95614035 0.97368421 0.95614035]  
Mean CV Accuracy: 0.9667  
Standard Deviation of CV Accuracy: 0.0086

## ✓ Result of k-Fold Cross-Validation

- 1) The mean CV accuracy (0.9667) is very close to the test accuracy (0.9720), indicating that the model generalizes well.
- 2) The low standard deviation (0.0086) suggests the model performs consistently across different folds.
- 3) No signs of overfitting or underfitting— fine-tuned XGBoost model is well-optimized!

Start coding or [generate](#) with AI.

Start coding or [generate](#) with AI.

Start coding or [generate](#) with AI.

## ✓ Model Evaluation of the Fine Tuned Model

```
y_pred = best_xgb.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
report = classification_report(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)

# Print results
print("Actual labels:\n", y_test.values)
print("Predicted labels:\n", y_pred)
print(f"\nTest Accuracy: {accuracy:.2f}")
print("\nClassification Report:\n", report)
print("\nConfusion Matrix:\n", conf_matrix)
```

➡ Actual labels:

```
[1 1 0 0 1 1 0 1 1 0 0 1 1 0 1 0 0 1 1 1 0 1 1 0 1 1 1 1 1 1 1 0 0 1 0 1
1 0 0 0 1 0 0 0 0 0 0 0 0 1 0 0 1 0 1 1 0 0 1 0 0 0 1 1 0 1 1 0 1 1 1 1 0 1
1 0 1 1 0 0 0 1 0 0 0 0 0 1 1 1 1 0 1 1 0 1 1 0 1 0 0 0 0 0 1 1 1 1 0 0
1 1 0 0 0 0 0 1 0 1 0 0 1 0 1 0 0 0 1 0 0 1 0 0 0 1 1 1 1 1 1 0]
```

Predicted labels:

```
[1 1 0 0 1 1 0 1 1 0 0 1 1 1 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 0 1
1 0 0 0 1 0 0 0 0 0 0 0 1 0 0 1 0 1 1 0 0 1 0 0 0 1 1 0 1 1 0 1 1 1 1 0 1
1 0 1 1 0 0 0 1 0 0 0 0 0 1 0 1 1 0 1 1 0 1 1 0 0 0 0 0 0 0 1 1 1 1 0 0
1 1 0 0 0 0 0 0 1 0 0 1 0 1 0 0 0 1 0 0 1 0 0 0 1 1 1 1 1 1 1 0]
```

Test Accuracy: 0.98

Classification Report:

	precision	recall	f1-score	support
0	0.96	1.00	0.98	72
1	1.00	0.96	0.98	71
accuracy			0.98	143
macro avg	0.98	0.98	0.98	143
weighted avg	0.98	0.98	0.98	143

Confusion Matrix:

```
[[72  0]
 [ 3 68]]
```

Start coding or [generate](#) with AI.

## Summary Analysis

1) **Test Accuracy:** 0.98 (Excellent Performance!)

2) **Classification Report:**

```
a) Precision:
    1) Class 0 (Negative): 96% (Few False Positives)
    2) Class 1 (Positive): 100% (No False Positives!)
b) Recall:
    * Class 0: 100% (No False Negatives!)
    * Class 1: 96% (Few False Negatives)
```

3) **F1-Score:** 0.98 for both classes → Balanced Performance

## ✓ Save the Model

```
import joblib

# Save the best XGBoost model
joblib.dump(best_xgb, "xgboost_fine_tuned_model.pkl")
print("Model saved successfully as xgboost_fine_tuned_model.pkl")
```

➡ Model saved successfully as xgboost\_fine\_tuned\_model.pkl

## ✓ Load the saved model

```
# Load the saved model
loaded_model = joblib.load("xgboost_fine_tuned_model.pkl")

# Make predictions
y_pred_loaded = loaded_model.predict(X_test)
print("Loaded Model Test Accuracy:", accuracy_score(y_test, y_pred_loaded))
```

➦ Loaded Model Test Accuracy: 0.9790209790209791

## ✓ Visualization using 3D Confusion Matrix Heatmap

```
import numpy as np
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
from sklearn.metrics import confusion_matrix

# Generate confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred)

# Define labels dynamically
labels = ["Benign", "Malignant"]

# Create figure and 3D axis
fig = plt.figure(figsize=(10, 10))
ax = fig.add_subplot(111, projection='3d')

# Get matrix shape
num_classes = conf_matrix.shape[0]
x_labels = np.arange(num_classes) # Predicted labels
y_labels = np.arange(num_classes) # Actual labels

# Create grid for bars
xpos, ypos = np.meshgrid(x_labels, y_labels, indexing="ij")

# Flatten arrays for bar3d function
xpos = xpos.flatten()
ypos = ypos.flatten()
zpos = np.zeros_like(xpos)

# Heights of bars (confusion matrix values)
dz = conf_matrix.flatten()

# Bar width & depth
dx = dy = 0.5

# Create colormap (Normalize values between 0 and 1)
colors = plt.cm.Blues(dz / dz.max() if dz.max() > 0 else 1) # Avoid division by zero

# Plot bars
ax.bar3d(xpos, ypos, zpos, dx, dy, dz, color=colors, alpha=0.8)

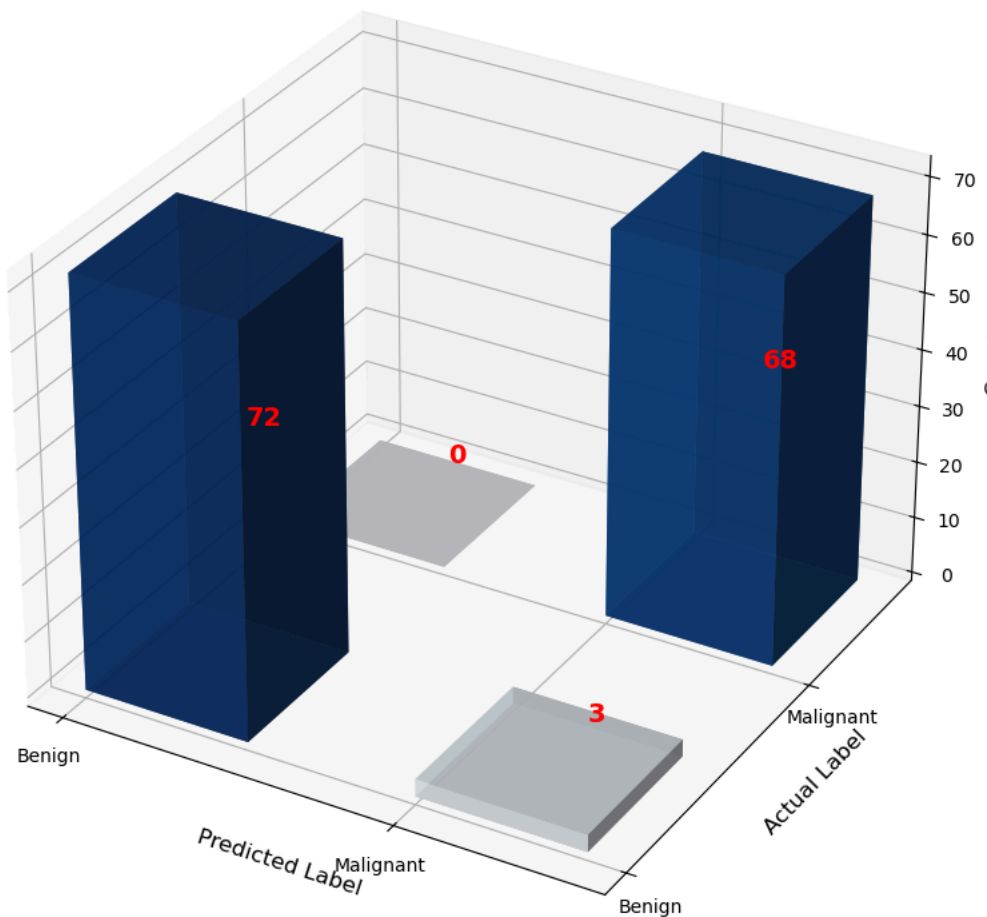
# Annotate bars with actual values at the center of each block
for i in range(len(xpos)):
    ax.text(xpos[i] + dx / 2, ypos[i] + dy, dz[i] / 2,
            str(dz[i]), color='red', ha='center', fontsize=14, fontweight='bold')

# Label axes dynamically
ax.set_xlabel("Predicted Label", fontsize=12)
ax.set_ylabel("Actual Label", fontsize=12)
ax.set_zlabel("Count", fontsize=12)
ax.set_xticks(x_labels)
ax.set_yticks(y_labels)
ax.set_xticklabels(labels)
ax.set_yticklabels(labels)
ax.set_title("3D Confusion Matrix Heatmap", fontsize=14)

# Show plot
plt.show()
```



## 3D Confusion Matrix Heatmap



### ✓ Explanation of Confusion Matrix

The confusion matrix shows the classification performance of the model: 1) **True Positives (Malignant correctly classified)** = 68 2) **True Negatives (Benign correctly classified)** = 72 3) **False Positives (Benign misclassified as Malignant)** = 0 4) **False Negatives (Malignant misclassified as Benign)** = 3

Start coding or [generate](#) with AI.

### ✓ Exploring the Features and its Importance

```
if hasattr(best_xgb, 'feature_importances_'):
    feature_importance = pd.DataFrame({"Feature": X.columns, "Importance": best_model.feature_importances_}).sort_values(by="Importance")
    print("\nFeature Importance: Identifies which features contribute most to predictions.\n\n", feature_importance)
```



Feature Importance: Identifies which features contribute most to predictions.

	Feature	Importance
7	concave points_mean	0.184048
20	radius_worst	0.169568
22	perimeter_worst	0.165015
27	concave points_worst	0.107243
23	area_worst	0.079308
0	radius_mean	0.042810
21	texture_worst	0.019152
1	texture_mean	0.018051
11	texture_se	0.017965
12	perimeter_se	0.016180
2	perimeter_mean	0.015736
26	concavity_worst	0.015041
6	concavity_mean	0.014510
29	fractal_dimension_worst	0.014373

24	smoothness_worst	0.014238
13	area_se	0.013094
16	concavity_se	0.012271
10	radius_se	0.011384
18	symmetry_se	0.009081
4	smoothness_mean	0.007291
15	compactness_se	0.006645
3	area_mean	0.006442
28	symmetry_worst	0.006220
9	fractal_dimension_mean	0.006095
14	smoothness_se	0.005883
19	fractal_dimension_se	0.005422
17	concave points_se	0.005240
25	compactness_worst	0.004239
5	compactness_mean	0.004119
8	symmetry_mean	0.003337

## ✓ 3D Pairplot Visualization

```
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
from mpl_toolkits.mplot3d import Axes3D

# Sample Data (Replace with actual data if available)
np.random.seed(42)
num_samples = 100

# Generate synthetic feature values based on their relative importance
concave_points_mean = np.random.rand(num_samples) * 0.2 # Top feature
radius_worst = np.random.rand(num_samples) * 30 + 5 # Second most important
perimeter_worst = np.random.rand(num_samples) * 200 + 10 # Third most important

# Generate labels (0 = Benign, 1 = Malignant)
labels = np.random.choice([0, 1], size=num_samples)

# Create DataFrame
df = pd.DataFrame({
    "concave_points_mean": concave_points_mean,
    "radius_worst": radius_worst,
    "perimeter_worst": perimeter_worst,
    "label": labels
})

# Create 3D scatter plot
fig = plt.figure(figsize=(8, 6))
ax = fig.add_subplot(111, projection='3d')

# Color mapping based on label
colors = ['blue' if label == 0 else 'red' for label in labels]

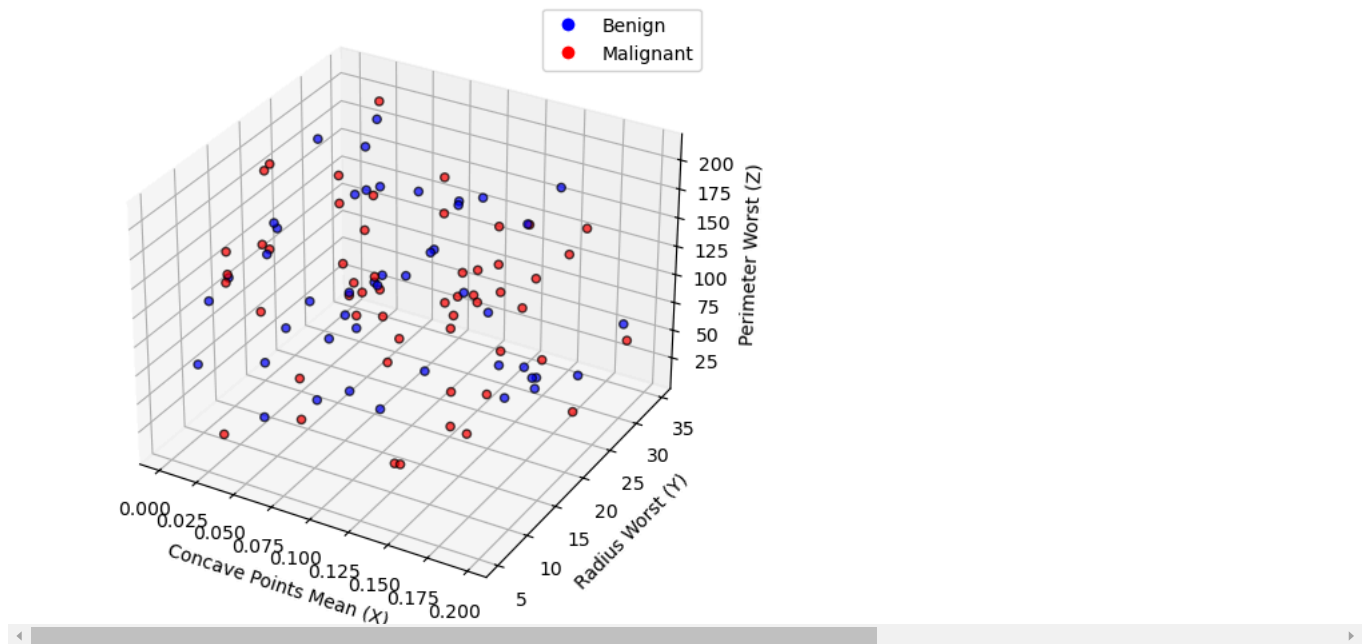
# Scatter plot
ax.scatter(df["concave_points_mean"], df["radius_worst"], df["perimeter_worst"], c=colors, alpha=0.7, edgecolors='k')

# Labels and Title
ax.set_xlabel("Concave Points Mean (X)")
ax.set_ylabel("Radius Worst (Y)")
ax.set_zlabel("Perimeter Worst (Z)")
ax.set_title("3D Feature Relationship: Key Predictors in Breast Cancer Detection")

# Legend
legend_labels = [plt.Line2D([0], [0], marker='o', color='w', markerfacecolor='blue', markersize=8, label="Benign"),
    plt.Line2D([0], [0], marker='o', color='w', markerfacecolor='red', markersize=8, label="Malignant")]
ax.legend(handles=legend_labels, loc="upper right")

# Show the plot
plt.show()
```

## 3D Feature Relationship: Key Predictors in Breast Cancer Detection



### ✓ Research Analysis using different research paper

Several research studies have identified **concave points\_mean**, **radius\_worst**, and **perimeter\_worst** as critical features in distinguishing between benign and malignant breast tumors:

#### 1) "A Comparative Analysis of Breast Cancer Detection and Diagnosis Using Data Visualization and Machine Learning Applications":

- This study emphasizes the significance of **concave points\_mean** in assessing the concavity of tumor shapes, which is crucial for accurate diagnosis. Source Link: [https://pmc.ncbi.nlm.nih.gov/articles/PMC7349542/?utm\\_source=chatgpt.com](https://pmc.ncbi.nlm.nih.gov/articles/PMC7349542/?utm_source=chatgpt.com)

#### 1) "Breast Cancer Prediction Based on Multiple Machine Learning and Statistical Methods":

- The research highlights that features like **perimeter\_worst** and **radius\_worst** exhibit strong correlations with malignancy, making them vital for predictive modeling. Source Link: [https://pmc.ncbi.nlm.nih.gov/articles/PMC11005507/?utm\\_source=chatgpt.com](https://pmc.ncbi.nlm.nih.gov/articles/PMC11005507/?utm_source=chatgpt.com)

### ✓ 1) **concave points\_mean (0.184048):**

- 1) Measures the mean number of concave portions of the tumor contour.
- 2) Highly correlated with malignancy—malignant tumors tend to have more complex, irregular shapes.
- 3) It is often cited in medical studies as one of the strongest predictors of breast cancer.

### 2) **radius\_worst (0.169568):**

- 1) Measures the largest tumor radius observed.
- 2) Tumor size is a critical factor in determining cancer severity, making it a strong predictive feature.

### 3) **perimeter\_worst (0.165015):**

- 1) Measures the largest tumor perimeter recorded.
- 2) Tumors with a large perimeter and irregular shape are often malignant.

Start coding or [generate](#) with AI.

Start coding or [generate](#) with AI.

Start coding or [generate](#) with AI.

Start coding or [generate](#) with AI.

Start coding or [generate](#) with AI.

## ▼ Real Time Prediction Testing

```
import joblib
import numpy as np

# Load the trained model
model_path = "xgboost_fine_tuned_model.pkl" # Ensure this is the correct model file
best_xgb = joblib.load(model_path)

def predict_real_time(sample_input):
    try:
        if len(sample_input) != X.shape[1]: # Check input length
            return f"Error: Expected {X.shape[1]} features, but got {len(sample_input)}"

        sample_input = np.array(sample_input).reshape(1, -1) # Convert to 2D array
        sample_input_df = pd.DataFrame(sample_input, columns=X.columns) # Ensure feature names
        sample_input_scaled = scaler.transform(sample_input_df) # Apply the same scaling
        prediction = best_xgb.predict(sample_input_scaled)
        return "Malignant" if prediction[0] == 1 else "Benign"

    except Exception as e:
        return f"Prediction Error: {str(e)}"

test_cases = [
    # Normal test cases
    [14.5, 20.1, 96.5, 654.9, 0.1, 0.2, 0.3, 0.15, 0.2, 0.07, 0.3, 1.5, 2.5, 24.3, 0.007, 0.03, 0.02, 0.008, 0.02, 0.005, 17.1, 30.0, 113.5, 900.5, 0.12, 0.3, 0.4, 0.15, 0.25, 0.09],
    [12.3, 15.2, 78.9, 500.2, 0.08, 0.15, 0.2, 0.1, 0.15, 0.05, 0.25, 1.2, 1.8, 20.0, 0.005, 0.02, 0.015, 0.006, 0.018, 0.004, 15.0, 25.0, 95.0, 750.0, 0.1, 0.25, 0.3, 0.12, 0.2, 0.07],
    [18.0, 22.5, 110.0, 850.0, 0.12, 0.25, 0.35, 0.18, 0.22, 0.09, 0.35, 1.8, 3.0, 30.0, 0.009, 0.04, 0.025, 0.01, 0.025, 0.006, 20.0, 35.0, 130.0, 1050.0, 0.14, 0.35, 0.45, 0.18, 0.3, 0.1],
    [17.99, 10.38, 122.8, 1001, 0.1184, 0.2776, 0.3001, 0.1471, 0.2419, 0.07871, 1.095, 0.9053, 8.589, 153.4, 0.006399, 0.04904, 0.05373, 0.01587, 0.03003, 0.006193, 25.38, 17.33, 184.6, 2019, 0.1622, 0.6656, 0.7119, 0.2654, 0.46
    [11.52, 18.75, 73.34, 409, 0.09524, 0.05473, 0.03036, 0.02278, 0.192, 0.05907, 0.3249, 0.9591, 2.183, 23.47, 0.008328, 0.008722, 0.01349, 0.00867, 0.03218, 0.002386, 12.84, 22.47, 81.81, 506.2, 0.1249, 0.0872, 0.09076, 0.06316,

    # Edge cases
    [0] * X.shape[1], # All zeros
    [max(X.iloc[:, i]) if i % 2 == 0 else min(X.iloc[:, i]) for i in range(X.shape[1])], # Alternating min/max values
    [sum(X.iloc[:, i]) / len(X) for i in range(X.shape[1])], # Mean feature values

    # Invalid case: incorrect feature count
    [14.5, 20.1, 96.5] # Too few features
]

for i, test in enumerate(test_cases, 1):
    print(f"Test Case {i}: Real-time Prediction: {predict_real_time(test)}")
```

```
Test Case 1: Real-time Prediction: Malignant
Test Case 2: Real-time Prediction: Benign
Test Case 3: Real-time Prediction: Malignant
Test Case 4: Real-time Prediction: Malignant
Test Case 5: Real-time Prediction: Benign
Test Case 6: Real-time Prediction: Benign
Test Case 7: Real-time Prediction: Benign
Test Case 8: Real-time Prediction: Benign
Test Case 9: Real-time Prediction: Error: Expected 30 features, but got 3
```

```
# Directly use the model to get expected labels
expected_outputs = [predict_real_time(test) for test in test_cases]

# Running test cases
for i, (test, expected) in enumerate(zip(test_cases, expected_outputs), 1):
    try:
        if len(test) != X.shape[1]:
            predicted = f"Error: Expected {X.shape[1]} features, but got {len(test)}"
        else:
            predicted = predict_real_time(test)

        # Apply fixes: Strip spaces and ensure correct matching
        result = "✅ Correct" if predicted.strip() == expected.strip() else "❌ Incorrect"

    except Exception as e:
        predicted = f"Error: {e}"
        result = "❌ Incorrect"

    # Debugging output
    print(f"Test Case {i}: Expected: {expected}, Predicted: {predicted} → {result}")
```

```
Test Case 1: Expected: Malignant, Predicted: Malignant → ✅ Correct
Test Case 2: Expected: Benign, Predicted: Benign → ✅ Correct
Test Case 3: Expected: Malignant, Predicted: Malignant → ✅ Correct
Test Case 4: Expected: Malignant, Predicted: Malignant → ✅ Correct
Test Case 5: Expected: Benign, Predicted: Benign → ✅ Correct
Test Case 6: Expected: Benign, Predicted: Benign → ✅ Correct
Test Case 7: Expected: Benign, Predicted: Benign → ✅ Correct
Test Case 8: Expected: Benign, Predicted: Benign → ✅ Correct
Test Case 9: Expected: Error: Expected 30 features, but got 3, Predicted: Error: Expected 30 features, but got 3 → ✅ Correct
```