

# Breast Cancer Detection Analysis

Aniakwa Nathan

August 1, 2025

## Abstract

This document presents an analysis of the breast cancer dataset using machine learning techniques. The dataset, containing 569 entries and 31 features, was preprocessed, split into training, validation, and test sets, and used to train Logistic Regression, SVM, and Neural Network models. Hyperparameter tuning was performed, and the Neural Network model achieved perfect classification on the test set. Key findings and code implementations are detailed below.

## 1 Introduction

This analysis uses the breast cancer dataset available in scikit-learn to develop predictive models for cancer detection. The dataset is loaded and preprocessed, followed by model training and evaluation. The code and results are presented as extracted from a Jupyter notebook.

## 2 Dataset Loading and Preprocessing

The dataset is loaded using scikit-learn's `load_breast_cancer` function and converted into a pandas DataFrame. The features are scaled using `StandardScaler`.

```
1 import pandas as pd
2 from sklearn.datasets import load_breast_cancer
3
4 # Load the dataset
5 cancer = load_breast_cancer()
6
7 # Create a DataFrame
8 df = pd.DataFrame(data=cancer.data, columns=cancer.feature_names)
9 df['target'] = cancer.target
10
11 # Check for missing values
12 print(df.isnull().sum())
13
14 # Separate features (X) and target (y)
15 X = df.drop('target', axis=1)
16 y = df['target']
17
```

```

18 # Scale numerical features
19 from sklearn.preprocessing import StandardScaler
20 scaler = StandardScaler()
21 X_scaled = scaler.fit_transform(X)

```

The output confirms no missing values in the dataset:

```

mean radius          0
mean texture         0
...
worst fractal dimension 0
target              0
dtype: int64

```

### 3 Data Splitting

The scaled data is split into training (398 samples), validation (85 samples), and test (86 samples) sets. The shapes are:

```

Shape of X_train: (398, 30)
Shape of X_val: (85, 30)
Shape of X_test: (86, 30)
Shape of y_test: (86,)

```

## 4 Model Training and Hyperparameter Tuning

Three models are trained: Logistic Regression, Support Vector Machine (SVM), and Neural Network. Hyperparameter tuning is performed using GridSearchCV.

### 4.1 Logistic Regression

```

1 # Define parameter grid for Logistic Regression
2 lr_param_grid = {
3     'C': [0.001, 0.01, 0.1, 1, 10, 100]
4 }

```

### 4.2 SVM

```

1 # Define parameter grid for SVM
2 svm_param_grid = {
3     'C': [0.1, 1, 10, 100],
4     'gamma': ['scale', 'auto'],
5     'kernel': ['linear', 'rbf']
6 }
7
8 # Initialize GridSearchCV for SVM
9 from sklearn.svm import SVC

```

```

10 svm_grid_search = GridSearchCV(SVC(), svm_param_grid, cv=5,
    scoring='f1')
11 svm_grid_search.fit(X_train, y_train)
12
13 # Print best parameters
14 print('Best parameters for SVM:', svm_grid_search.best_params_)

```

## 4.3 Neural Network

```

1 # Define parameter grid for Neural Network
2 nn_param_grid = {
3     'hidden_layer_sizes': [(50,), (100,), (50, 50), (100, 50)],
4     'activation': ['tanh', 'relu'],
5     'solver': ['adam', 'sgd'],
6     'alpha': [0.0001, 0.001, 0.01, 0.1],
7     'learning_rate': ['constant', 'adaptive']
8 }
9
10 # Fit GridSearchCV
11 nn_grid_search.fit(X_train, y_train)
12 print('Best F1 score for Neural Network:',
    nn_grid_search.best_score_)

```

Output:

Best F1 score for Neural Network: 0.9860754681113963

## 5 Model Evaluation

The models are evaluated on the validation set using accuracy, precision, recall, F1-score, and AUC.

### 5.1 Logistic Regression Metrics

```

1 from sklearn.metrics import accuracy_score, precision_score,
    recall_score, f1_score, roc_auc_score
2
3 # Logistic Regression evaluation
4 accuracy = accuracy_score(y_val, y_val_pred)
5 precision = precision_score(y_val, y_val_pred)
6 recall = recall_score(y_val, y_val_pred)
7 f1 = f1_score(y_val, y_val_pred)
8 auc = roc_auc_score(y_val, model.predict_proba(X_val)[:, 1])
9
10 print('Logistic Regression Metrics:')
11 print(f'Accuracy: {accuracy:.4f}')
12 print(f'Precision: {precision:.4f}')
13 print(f'Recall: {recall:.4f}')
14 print(f'F1-score: {f1:.4f}')

```

```
15 print(f'AUC: {auc:.4f}')
```

## 5.2 SVM Metrics

```
1 svm_model_prob = SVC(probability=True)
2 svm_model_prob.fit(X_train, y_train)
3 svm_pred = svm_model_prob.predict(X_val)
4 svm_accuracy = accuracy_score(y_val, svm_pred)
5 svm_precision = precision_score(y_val, svm_pred)
6 svm_recall = recall_score(y_val, svm_pred)
7 svm_f1 = f1_score(y_val, svm_pred)
8 svm_auc = roc_auc_score(y_val,
    svm_model_prob.predict_proba(X_val)[:, 1])
```

## 5.3 Neural Network Test Set Evaluation

The Neural Network, identified as the best model, is evaluated on the test set:

```
1 test_model = best_nn_model
2 y_pred = test_model.predict(X_test)
3
4 test_accuracy = accuracy_score(y_test, y_pred)
5 test_precision = precision_score(y_test, y_pred)
6 test_recall = recall_score(y_test, y_pred)
7 test_f1 = f1_score(y_test, y_pred)
8 test_auc = roc_auc_score(y_test,
    test_model.predict_proba(X_test)[:, 1])
9
10 print('Neural Network Test Metrics:')
11 print(f'Accuracy: {test_accuracy:.4f}')
12 print(f'Precision: {test_precision:.4f}')
13 print(f'Recall: {test_recall:.4f}')
14 print(f'F1-score: {test_f1:.4f}')
15 print(f'AUC: {test_auc:.4f}')
```

Output:

```
Neural Network Test Metrics:
Accuracy: 1.0000
Precision: 1.0000
Recall: 1.0000
F1 score: 1.0000
AUC: 1.0000
```

## 6 Data Sample

A sample of the dataset is presented in Table 1.

Table 1: Sample of Breast Cancer Dataset

Index	Mean Radius	Mean Texture	Mean Perimeter	Mean Area	Mean Smoothness	Target
0	17.99	10.38	122.80	1001.0	0.11840	0
1	20.57	17.77	132.90	1326.0	0.08474	0
2	19.69	21.25	130.00	1203.0	0.10960	0
3	11.42	20.38	77.58	386.1	0.14250	0
4	20.29	14.34	135.10	1297.0	0.10030	0

## 7 Key Findings

The breast cancer dataset, containing 569 entries and 31 features, was successfully loaded and confirmed to have no missing values. The data was preprocessed by scaling features and splitting into training, validation, and test sets. The Neural Network model, after hyperparameter tuning, achieved perfect classification on the test set with an F1 score of 0.9861. While promising, further validation is needed to ensure generalization. The analysis provides insights into key features for breast cancer detection.