Medical science has greatly benefited from the application of machine learning and data mining in technology. Numerous studies have been conducted to apply machine learning in the medical field to classify medical datasets. This section provides a brief overview of several studies conducted in the field of breast cancer detection.

MM Islam et al. [10] proposed a novel approach that used two classification algorithms named Support Vector Machine (SVM) and K-Nearest Neighbors (K-NN) for the detection of breast cancer. The Wisconsin breast cancer diagnosis dataset that was used in their implementation was taken from the UCI machine learning repository. For accurate results, 10-fold cross-validation was performed. The approach's accuracy was 98.57% for SVM and 97.14% for K-NN, respectively, and their model's specificity for the testing phase was 95.65% for SVM and 92.31% for K-NN. However, their proposed model used a single dataset and only two machine learning techniques.

In another study, K Sivakami et al. [11] presents a disease status prediction strategy. This strategy is divided into two parts. 1. Information Treatment and Option Extraction, and 2. Decision Tree-Support Vector Machine (DT-SVM) Hybrid Model for predictions. They used Weka Software tools for data preparation, data analysis, and result comparison. In this study, three classification techniques are compared, and DT-SVM (91%) outperforms Instance-based Learning (IBL), Sequential Minimal Optimization (SMO), and Nave-based classifiers. They also used the Wisconsin breast cancer diagnosis dataset from the UCI machine learning repository. Furthermore, they only used single dataset, and their prediction results are poor. In addition, their dataset distribution (60% training and 40 testing) is not perfect because the sample dataset is limited.

Similarly, AA Bataineh et al. [14] based on performance metrics (such as accuracy, recall, and precision), they compared the effectiveness and efficiency of five nonlinear ML algorithms on the Wisconsin Breast Cancer Diagnostic (WBCD) dataset: Multilayer Perceptron (MLP), K-Nearest Neighbors (KNN), Classification and Regression Trees (CART), Gaussian Nave Bayes (NB), and Support Vector Machines (SVM). The Multilayer Perceptron (MLP) algorithm achieved the highest accuracy of the five nonlinear ML algorithms, with 99.12%. Despite having more proposed algorithms than the previous two and offering higher performance, the dataset problem persists.

Moreover, Mohammed et al. [15] proposed an approach that improves the accuracy and enhances the performance of three different classifiers: Decision Tree (J48), Naïve Bayes (NB), and Sequential Minimal Optimization (SMO). They conducted two different datasets as Wisconsin Breast Cancer (WBC) and the Breast Cancer dataset. The authors of this paper mainly focus to deal on imbalanced data. Data imbalance is a big problem in the classification field. The resampling techniques are used to deal with imbalanced data. In addition, they employed the 10-fold crossed validation method to solve the Data Unbalance issue. In particular, a resample filter was used to improve classifier performance. SMO outperformed the other two classifiers in terms of accuracy among the three. To assess the efficiency of the classifiers, they considered accuracy, standard deviation, ROC curve, and so on.

In addition, S Sharma et al. [17] compare the performance of three machine learning algorithms (ML) on the Wisconsin Diagnosis Breast Cancer dataset: Random Forest, K-NN(K-Nearest-Neighbor), and Naive Bayes. To compare the model performances, they have considered accuracy and precision. The K-Nearest-Neighbor(K-NN) (95.90%) outperforms Random Forest and Naïve Bayes. Furthermore, K-Nearest-Neighbor (K-NN) had the highest precision (98.27%) and f1-Score (94.20%). This study also made use of fewer ML algorithms and a single dataset.

Some researchers optimize the model to improve performance. To detect breast cancer, Assegie et al. [18] proposed an optimized K-Nearest Neighbor (K-NN) model. They used grid search techniques to find the best value of k for the K-NN model. This study also compares the effect of the hyper-parameter tuning model to the effect of the default hyper-parameter model. Hyper-parameter tuning has a significant impact on the performance of the KNN model. The optimized hyper-parameter tuning model then achieved 94.35%, while the default hyper-parameter achieved 90.10%. However, the model's performance is insufficient, and it only used a single dataset.

On the other hand, MF Aslan et al. [12] employed four machine learning algorithms for the early identification of breast cancer: the Artificial Neural Network (ANN), the conventional Extreme Learning Machine (ELM), the Support Vector Machine (SVM), and the K-Nearest Neighbor algorithm (k-NN). Their objective is to process the results based on the routine blood analysis and determine how well these techniques work to find breast cancer. They considered age, body mass

index (BMI), glucose, insulin, homeostasis model assessment (HOMA), leptin, adiponectin, resistin, and chemokine monocyte chemoattractant protein 1 (MCP1) attributes across the entire blood analysis dataset to complete this study. The blood analysis dataset was collected from papers M Patrício et al. [28]. They also tune the model's hyper-parameter for better results. Finally, standard Extreme Learning Machine (ELM) methods obtained the best results (80%) with the shortest training time (0.0075s) compared to the other three. Despite the fact that this study used a different dataset, model performance issues persist.

In order to identify the subtype of breast cancer, AA Bataineh et al. [7] present a performance comparison study. Five nonlinear machine learning methods are compared in this study: Gaussian Nave Bayes (NB), Classification and Regression Trees (CART), Multilayer Perceptron (MLP), and Support Vector Machines (SVM). Wisconsin Breast Cancer Diagnostic is their conducted dataset (WBCD). Finally, they evaluate the model performance with respect to the effectiveness and efficiency of each algorithm in terms of precision, recall, and accuracy. The Multilayer Perceptron (MLP) model outperformed the others in terms of accuracy (96.70%), precision (100%), and recall (97%). They used 10-fold cross-validation for more accurate results.

Table: Comparison of Previous Study

| Authors | Models | Dataset | Results (Best) | Weakness |
|---------|--------|---------|----------------|----------|
| MM Islam et al. [10] | 1. Support Vector Machine (SVM) 2. K-Nearest Neighbors (K-NN) | Wisconsin breast cancer diagnosis (WBCD) dataset | SVM (98.57%) | 1. Single dataset Only two 2. Machine learning techniques. |
| K Sivakami et al. [11] | 1. DT-SVM 2. Instance-based Learning (IBL) 3. Sequential Minimal Optimization (SMO) 4. Nave-based classifiers. | Wisconsin breast cancer diagnosis (WBCD) dataset | DT-SVM (91%) | 1. Single dataset 2. Results are poor 3. Dataset distribution (60% training and 40 testing) is not perfect 4. Fewer ML algorithms |

| | | | | |
|---|---|---|---|---|
| AA Bataineh et al. [14] | 1. Multilayer Perceptron (MLP)<br>2. K-Nearest Neighbors (KNN)<br>3. Classification and Regression Trees (CART)<br>4. Gaussian Nave Bayes (NB)<br>5. Support Vector Machines (SVM) | Wisconsin breast cancer diagnosis (WBCD) dataset | MLP (99.12%) | 1. Single dataset |
| Mohammed et al. [15] | 1. Decision Tree (J48)<br>2. Naïve Bayes (NB)<br>3. Sequential Minimal Optimization (SMO) | 1. Wisconsin Breast Cancer (WBC)<br>2. Breast Cancer dataset | 1. For Wisconsin Breast Cancer (WBC) SMO (99.56%)<br>2. For Breast Cancer dataset J48 (98.20%) | 1. Fewer ML algorithms |
| S Sharma et al. [17] | 1. Random Forest<br>2. K-NN(K-Nearest-Neighbor)<br>3. Naive Bayes | Wisconsin breast cancer diagnosis (WBCD) dataset | K-NN (95.90%) | 1. Fewer ML algorithms<br>2. Single dataset |
| Assegie et al. [18] | 1. K-Nearest Neighbor (K-NN) | Wisconsin breast cancer (WBC) dataset | K-NN (94.35%) | 1. Model performance is insufficient<br>2. Single dataset |
| MF Aslan et al. [12] | 1. Artificial Neural Network (ANN)<br>2. conventional Extreme<br>3. Learning Machine (ELM)<br>4. Support Vector Machine (SVM)<br>5. K-Nearest Neighbor algorithm (k-NN). | Blood analysis Dataset | ELM (80%) | 1. Model performance is insufficient |
| AA Bataineh et al. [7] | 1. Gaussian Nave Bayes (NB)<br>2. Classification and Regression Trees (CART)<br>3. Multilayer Perceptron (MLP)<br>4. Support Vector Machines (SVM) | Wisconsin breast cancer diagnosis (WBCD) dataset | MLP (96.70%) | 1. Low performance<br>2. Single dataset |

We can observe from the summary above that several studies employed few machine learning (ML) approaches, had poor performance, and only used one or two datasets. In this study, we attempt to address these problems and demonstrate how these ML approaches compare to one another. From the prior research, we chose two datasets and twelve machine learning approaches to conduct this study.