



A Comparative Analysis of Methods for Detecting and Diagnosing Breast Cancer Based on Data Mining

Ahmed T. Alhasani ¹, Hussein Alkattan ^{*2}, Alhumaima Ali Subhi ^{*3}, El-Sayed M. El-Kenawy ⁴,
Marwa M. Eid ⁴

¹ Al-Furat Al-Awsat Technical University Computer Center Administrator, Najaf, Iraq

² Department of System Programming, South Ural State University, Chelyabinsk 454080, Russia

³ Electronic Computer Center University of Diyala, Diyala, Iraq

⁴ Faculty of Artificial Intelligence, Delta University for Science and Technology, Mansoura, Egypt

Emails: ahmed.alhasani@atu.edu.iq; alkattan.hussein92@gmail.com;

alhumaimaali@uodiyala.edu.iq; skenawy@ieee.org; mmm@ieee.org

Abstract

Breast cancer is a significant public health concern worldwide, and early detection is crucial for its treatment. Although breast cancer has been extensively studied, there is still room for improvement in its classification accuracy. This study aims to improve the classification accuracy of breast cancer by applying information gain feature selection and machine learning techniques to the Wisconsin Diagnostic Breast Cancer (WDBC) dataset. The information gain method is utilized to reduce feature characteristics, and machine learning algorithms such as support vector machine (SVM), naive Bayes (NB), and C4.5 decision tree are employed for breast cancer classification. The study also conducts a comparison analysis based on accuracy value. The proposed model achieves maximum classification accuracy (100%) and a weighted average for precision (100%) and recall (100%) using a C4.5 decision tree, while SVM accuracy (98.42%) and weighted average for precision (98.17%) and recall (98.58%) are achieved using a C4.5 decision tree. The NB algorithm attains an accuracy of 96%, with a weighted average for precision (18.57%) and recall (50%). The proposed model's results are compared to similar studies and demonstrate significant progress, indicating new opportunities for breast cancer detection.

Keywords: Information Gain Feature Selection; Machine learning; classifier support vector machine; classifier naive Bayes; classifier C4.5 decision tree; Performance evaluation tests

1. Introduction

Breast cancer is a widespread cancer among women worldwide and is the most common form of cancer, surpassing lung cancer. It accounted for 11.7% of all cancer cases in 2020, with approximately 2.3 million new cases. Breast cancer is also the leading cause of cancer-related deaths among females, resulting in around 685,000 fatalities in 2020. Data mining is a powerful tool for discovering patterns and relationships within vast datasets. In the medical field, classification and data mining techniques are widely used for diagnosis and analysis to aid decision-making. This study utilizes publicly available breast cancer tumor data from the University of Wisconsin Hospital and employs feature selection and machine learning algorithms, namely support vector machine, C4.5 decision tree, and Naïve Bayes, for breast cancer classification. The study aims to identify the primary characteristics that influence breast cancer categorization and accurately diagnose and detect the disease in its early stages. The study also compares the accuracy of various machine learning algorithms for classification. The paper is organized into sections covering an overview of relevant literature, methods and materials, the proposed model, simulation method, evaluation metrics, and experimental outcomes and discussion. Overall, the proposed work has the potential to advance breast cancer detection and classification.

Guyon and Elisseeff (2003) proposed that information gain (IG), also known as mutual information, can identify the most beneficial attribute in a given set of training feature vectors for discriminating between classes to be learned by searching for a subset of the original variables. IG is one of the following three selection strategies: filter, wrapper, and embedded [1-3]. Utilizing selection techniques allows for the selection of relevant and informative features, or the selection of features that aid in the development of an accurate predictor. IG relies on entropy, a measure of the unpredictability of information, and the rank class of the characteristics that influence data classification. Moreover, p_i represents the probability of i in the given set of attributes [4-6].

The goal of this paper is Using data mining and machine learning techniques, this paper appears to propose an efficient model for detecting and diagnosing breast cancer in its early stages. Using a proposed Map-Reduce technique and information acquisition feature selection method, this paper seeks to better the accuracy and prognosis of patients with early breast cancer by identifying significant features. In addition to comparing the robustness of the proposed model with accuracy measures, the paper provides a comprehensive comparison and analysis of three machine learning algorithms in terms of accuracy, recall, and sensitivity measurements. Using data mining and machine learning techniques, the overall objective is to improve the detection and diagnosis of breast cancer [7-9].

2. Methods

A. Naive Bayes Classifier

A probabilistic classifier known as the Naive Bayes classifier is one that is founded on the Bayes theorem and takes into account a robust (naive) independence assumption. As a result, a Naive Bayes classifier takes into account the notion that each characteristic or feature independently contributes to the likelihood of a particular conclusion [10]. When the characteristics of the underlying probability model are taken into account, the Naive Bayes classifier is capable of being trained in a supervised learning environment in an extremely effective manner. As a result, it performs significantly better in many difficult real-world scenarios, particularly in computer-aided diagnosis, than one might anticipate. Because it is assumed that the variables are independent of one another, it is sufficient to compute simply the variances of the variables for each class rather than the whole covariance matrix [11-14].

$$p(F_1 \dots F_n) = \frac{p(C)p(F_1 \dots F_n|C)}{p(F_1 \dots F_n)} \quad (1)$$

where P is the probability, C is the class variable and $F_1 \dots F_n$ are Feature variables F_1 through F_n . The denominator is independent of C .

B. C4.5 Decision Tree Classifier

The information gain ratio, which is measured by entropy, serves as the foundation for C4.5. The information gain ratio metric is applied to each node in the tree in order to pick the test features for that node. [15-18] A measure of this kind is referred to as a feature (attribute) selection measure. As the test feature for the present node, the characteristic that has been determined to have the highest information gain ratio is selected. Let us assume that D is a set that contains the data instances D_1 through D_j . Imagine that the class label attribute contains m unique values, each of which defines m separate classes C_i (where i might range from 1 to m). Let's say that the number of D samples in class C_i is denoted by D_j . Provide the information that is anticipated to be required for the classification of a particular sample.

$$Splitinfo_A(D) = - \sum (|D_j| / |D|) * \log \log (|D_j| / |D|) \quad (2)$$

$$Gainratio(A) = Gain(A) / Splitinfo_A(D) \quad (3)$$

$$Gain = Info(D) - info_A(D) \quad (4)$$

$$Info(D) = - \sum p_i p_i \text{ and} \tag{5}$$

$$Splitinfo_A(D) = - \sum (|D_j| / |D|) * info(D_j) \tag{6}$$

3. Proposed Model

In this paper, the researchers have presented a new model to improve accuracy of the breast cancer classification. This model basically consists of the several stage are: Load WDBC dataset, preprocessing include proposed Map-Reduce technique, feature selection using information gain algorithm, and finally diagnosis and detected of breast cancer based on three machine learning algorithms includes Naive Bayes NB, Support Vector Machine SVM, C 4.5 decision tree. The framework of the proposed breast cancer classification comes in Figure 2.

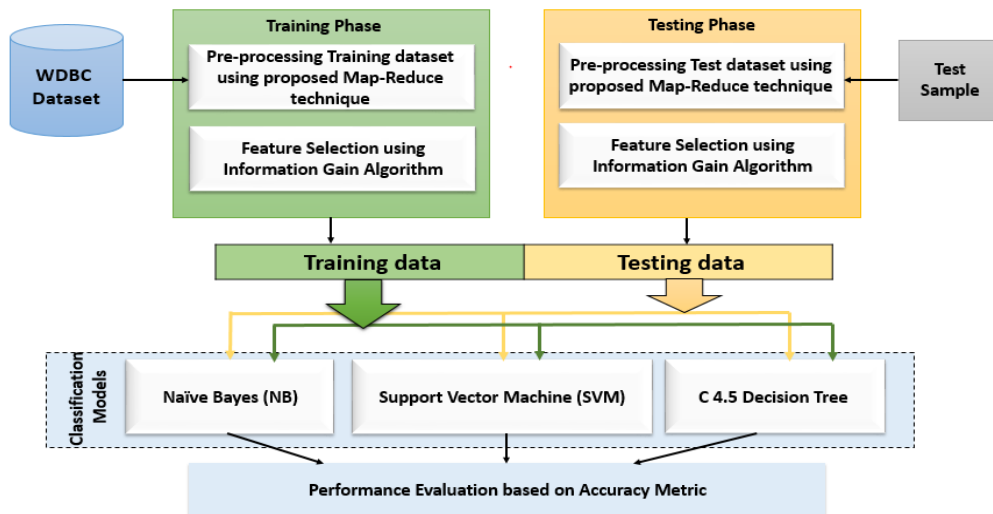


Figure 1: The Framework of the Proposed Breast Cancer Classification Model.

4. Pre-processing WDBC Dataset

After load data from WDBC dataset the proposed model prepare dataset by generation coding for each value in the input dataset using a propose Map-Reduce technique. The idea of the propose Map-Reduce technique is building a dictionary for the database includes (value (Feature), counter indicating the number of times a given value appears in Feature, index of dictionary) as illustrated in Figure 3. The proposed technique aims to eliminate duplicate values that may affect the accuracy of classification.

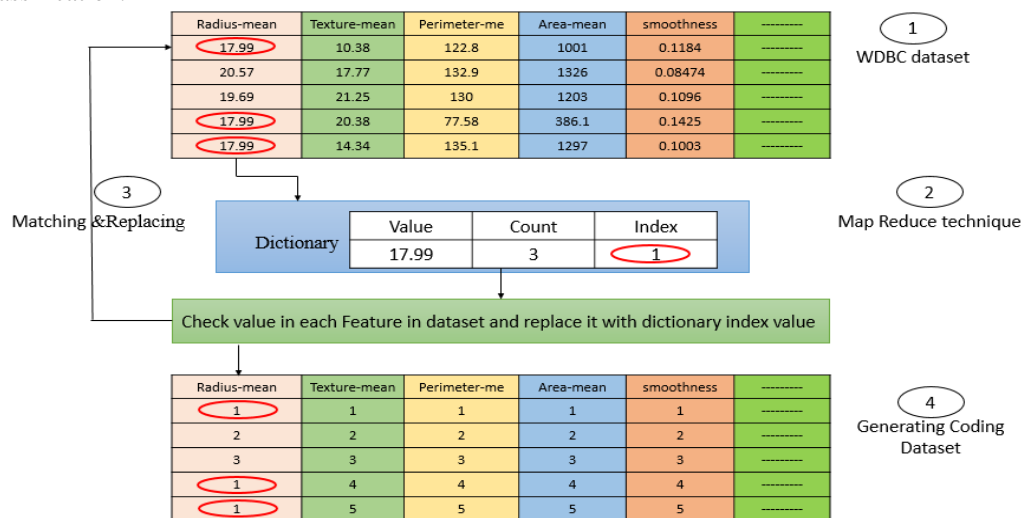


Figure 2: Main Steps of the Proposed Map-Reduce Technique to Generation Coding for Each Value in WDBC Dataset

For example as clarify in figure 3, the Map-Reduce technique take first feature: Radius-mean [17.99,2057,19.69, 17.99, 17.99,] from WDBC dataset .Now take each value in Radius –mean feature and find the number of times this value appears in this feature and give it a unique number keeping in mind the sequence of the index within the dictionary i.e. (value =17.99, count =3,index=1). Each value (17.99) in radius -mean is replaced by the index value (1) in the dictionary, So the index value represents the coding of the original value [19-21].

Figure 4 shows the impact of the Map-Reduce technique on 8 features from WDBC are: [radius mean, perimeter mean, area mean, compactness means, concavity mean, concave points mean, fractal dimension mean, perimeters]; Figure (4.a) presents original values without applied Map-Reduce technique and Figure (4.b) illustrates 1 values after applied Map-Reduce technique.

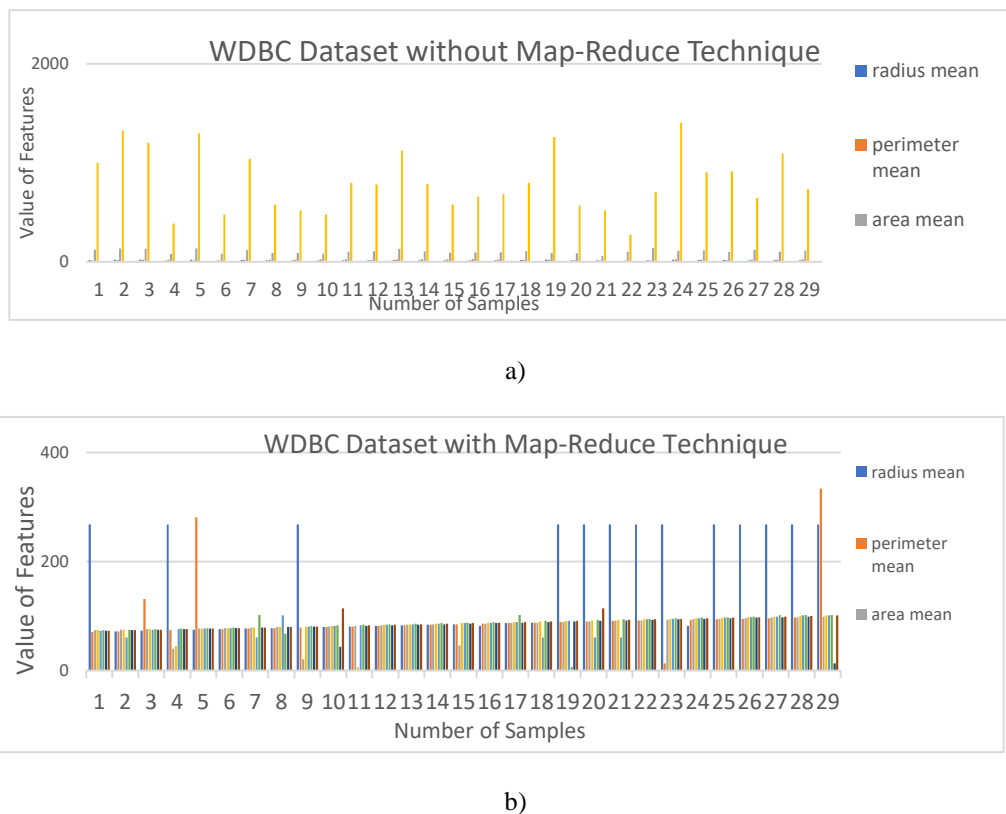


Figure 3: Application of Map-Reduce Technique to 8 Features of the WDBC Dataset; a) the Dataset without Coding Data, b) the Dataset with Coding Data.

5. Feature Selection using Information Gain Method

Table 1 provides a description of each of the 32 Features that are included in the WDBC dataset. If the proposed model takes on board all of the characteristics contained in the database and employs those characteristics in order to detect and diagnose breast cancer, then it is possible that the model will require a significant quantity of memory as well as a significant amount of time to perform the necessary computations. Additionally, it is likely that certain characteristics do not directly indicate whether or not the individual in question has breast cancer. As a result, the accuracy with which the disease is diagnosed may be negatively impacted [22-25]. Because of these factors, the model that was proposed resorted to making use of a method known as information gain in order to determine the characteristics that contribute the most to the considerable improvement in the performance of the classification algorithm. Calculating the entropy value of each feature in the data using Equation 1 is a prerequisite for employing the Information Gain methodology for feature selection. The value of entropy is applied when rating characteristics that have an impact on data classification. A feature that does not have a significant effect on the classification of the data has a very modest information gain, and it may be omitted without reducing the accuracy of the detection of a classed as indicated in Table 1.

Table 1: Information Gain Values for Each Feature in the WDBC Database.

#	Name Feature	information Gain	#	Name Feature	information Gain
1	concave points_mean	0.942090	16	compactness_mean	0.909129
2	smoothness_se	0.93506	17	radius_worst	0.900307
3	concavity_mean	0.93506	18	perimeter_worst	0.898584
4	area_worst	0.93506	19	fractal_dimension_worst	0.891554
5	radius_se	0.933733	20	texture_se	0.864296
6	perimeter_se	0.931545	21	texture_worst	0.863434
7	concavity_worst	0.930218	22	radius_mean	0.860781
8	concavity_se	0.928030	23	concave points_se	0.858593
9	area_mean	0.928030	24	symmetry_worst	0.84539
10	perimeter_mean	0.926703	25	fractal_dimension_mean	0.836177
11	area_se	0.925377	26	texture_mean	0.835711
12	compactness_se	0.923188	27	symmetry_se	0.818137
13	compactness_worst	0.918347	28	smoothness_mean	0.776178
14	fractal_dimension_se	0.917485	29	symmetry_mean	0.735036
15	concave points_worst	0.914832	30	smoothness_worst	0.719718

Table 1 displays the values arranged in descending order from the highest value to the lowest value. By applying the information gain method, keep the 25 features with the highest information gain values, that is, except for the 5 features that have the lowest values in the table. the results of the information gain are in range [0-1] and rearrange this value and remove the features that have lowest value which are [texture_mean =0.83571188, symmetry_se = 0.818137187, smoothness_mean =0.776178834, symmetry_mean = 0.735036638, smoothness_worst = 0.71971891].

6. Classification Model using Machine Learning

The Proposed model uses the most common and effective machine learning algorithms in the detection and diagnosis of breast cancer. Classification is two step processes: learning or training step where data is analyzed by a classification algorithm and testing step where data is used for classification and to estimate the accuracy of the Classification [26-28]. The input to classifier models C4.5 decision tree, Support Vector Machine (SVM), and Naïve Bayes (NB) is important features that selected in previous step and the output of these models is classified class breast cancer malignant or benign tumor [29].

7. Experimental Results and Discussion

This section will present the performance experiments of the proposed model and compare the classification algorithms used in this work based on the results of measures of accuracy, recall, and accuracy to discover which algorithm performed better in early detection and diagnosis of breast cancer. Take into account that the proposed model was applied to Wisconsin Diagnostic Breast Cancer (WDBC) dataset [30].

The experiment of the proposed model run under Windows 10 Professional operating system, an Intel(R) Core (TM) i3-7020U CPU @ 2.30GHz, 8 GB of random-access memory, and a 64-bit system type, with the proposed system running in a C# language environment.

In order to compare the performance of the C4.5 decision tree, Support Vector Machine (SVM), and Naïve Bayes (NB) classifier models, accuracy, precision, and recall measurements are used. Accuracy

metric is given by Eq. (11), where TP , TN , FP , and FN represent True Positive, True Negative, False Positive, and False Negative, respectively.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \times 100\% \quad (7)$$

Precision is the ratio of correct positive results divided by the number of all total predicted positive observations. Mathematically, it can be expressed as :

$$Precision = \frac{TP}{TP + FP} \times 100\% \quad (8)$$

Recall is also named sensitivity, is the ratio of correct positive results to all observations in actual class. Mathematically, it can be expressed as:

$$Precision = \frac{TP}{TP + FN} \times 100\% \quad (9)$$

After passing the dataset through Map-Reduce Technique to generate code for each data in the WDBC dataset to prepare this dataset to feed the Information Gain feature selection method to determine 25 important features in diagnosis cancer directly. The proposed model divided the dataset into 80% training dataset and 20% testing dataset. Breast cancer data contains tumors which represents the severity of the disease. To classify the tumors correctly from the training data set, the error rates and accuracy are calculated using classifiers. Table 3. shows test results based on accuracy metrics using Equation 11 for each machine learning algorithms. The best performance of the proposed model with C4.5 decision tree has accuracy ratio (100%) while accuracy ration of SVM (98.42%) and accuracy ratio value of the NB (96.6%).

Table 2: Test Accuracy Results for C4.5 decision tree, SVM, and NB classifier models.

Metric	Classification Methods with WDBC Dataset		
	C4.5 decision tree	SVM	NB
Accuracy	100%	98.42%	69.6%

Figure 4 shows compare between classifier models based on accuracy value that present in table 3. This compare proved the proposed model has ability to detection and diagnosis of malignant and benign tumors with perfect accuracy and there is no error rate, as the proposed model was able to classify cancer with an accuracy of 100% with C4.5 decision tree.

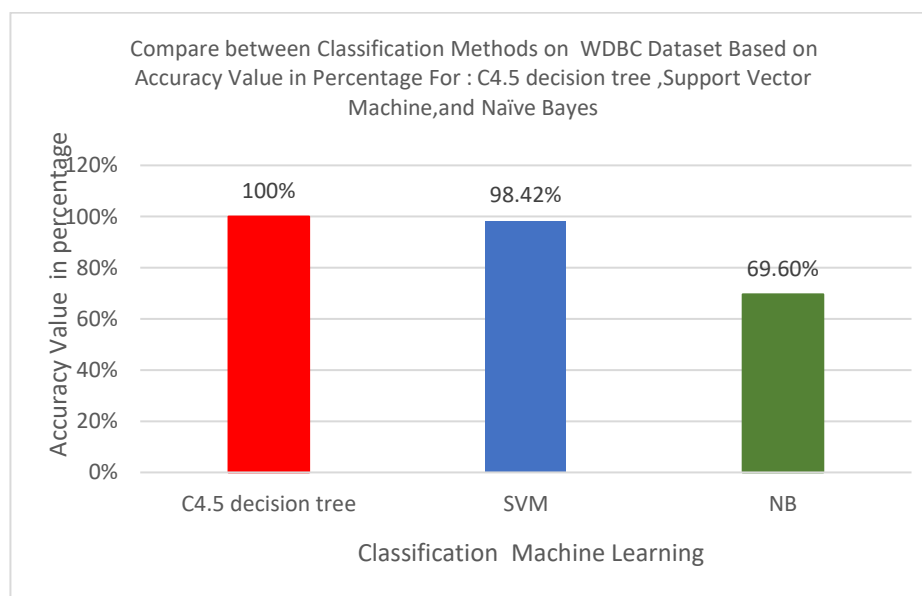


Figure 4: Performance Comparison of Classification Algorithm based on Accuracy Value

The performance of algorithms C4.5 decision tree, SVM, and NB are given the Tables 4. The classification precision using Equation (12) and Recall using Eqn (13) of three algorithms C4.5 decision tree, SVM, and NB are observed from Tables 4 via values of weighted average, which is available in the last row of each table.

Table 3: Results for C4.5 decision tree, SVM, and NB classifier models.

Class	Classification Models					
	C4.5 decision tree		SVM		NB	
	precision	Recall	precision	Recall	precision	Recall
Malignant	100%	100%	97.20%	98.58%	37.15%	100%
Benign	100%	100%	99.15%	98.32%	0%	0%
weighted average	100%	100%	98.17%	98.58%	18.57%	50%

The Figure 5 represents the comparison of the C4.5 decision tree, SVM, and NB classifier models based on the Table 4 values of weighted average Precision and Recall.

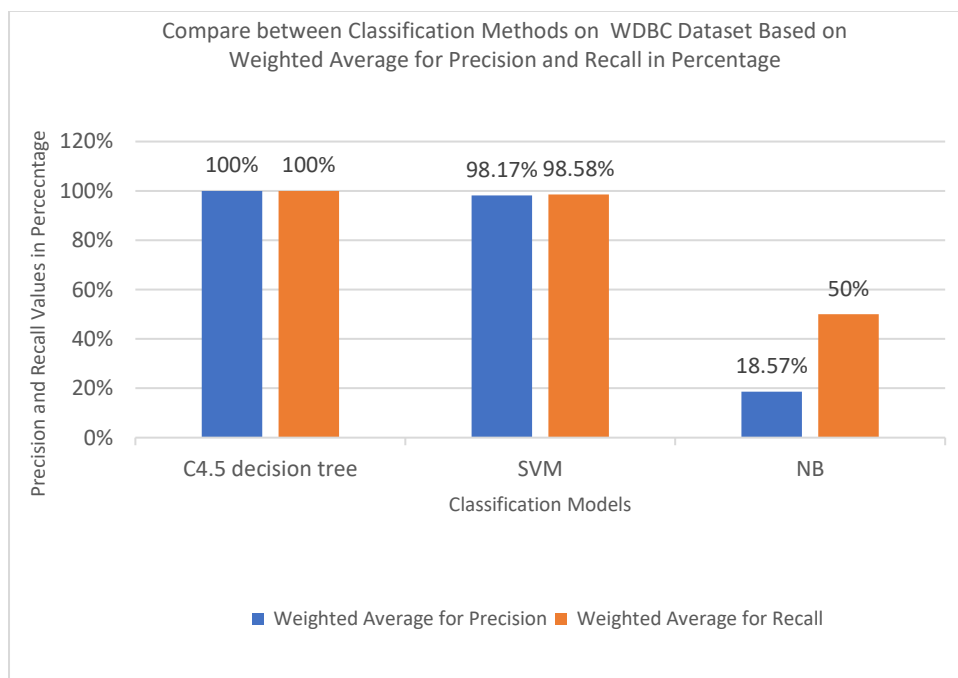


Figure 5: Performance Comparison of Classification Algorithm based on Weighted Average Precision and Recall Value

The best performance of the proposed model over all classification models based on weighted average precision and recall values as shown in table 4 and figure 6 with C4.5 decision tree that has best results of the weighted average precision (100%) and Recall (100%) while the results of the weighted average precision (98.17%) and Recall (98.58%). These results proved that the proposed model achieved excellent results for detected and diagnoses breast cancer.

Breast Cancer detection is a subject of considerable interest in the scientific literature because of the high accuracy of cancer detection by classification technique. The suggested approach delivered optimal results about accuracy. In table 5 shows the accuracy of the suggested model when compared to prior approaches as illustrated in section 2. The technique under consideration has a high success

rate, with the C45 decision tree classification algorithm achieving 100% accuracy. Thus, the proposed method could effectively diagnosis Breast Cancer.

Table 5: Analyses of the suggested method's performance on the WDBC dataset.

No	Authors & Reference	Classification Methods	Accuracy Value in Percent
1	Shokoufeh Aalaei et al.[7]	Artificial Neural Network (ANN)	97.3 %
2	Kui Liu et al.[8]	Fully-Connected Layer First Convolutional Neural Networks (FCLF-CNN)	99.28%
3	Laila Khairunnahar et al.[10]	Modified Logistic Regression	96.83%
4	Muhammet Fatih Ak .[4]	Logistic Regression	98.1%
5	Md Akizur Rahman & Ravie Chandren Muniyandi [11]	15-neuron neural network (NN)	99.4%
6	Ali Idri et al.[12]	Grid Search multilayer perceptron (GSMLP)	98.07 %
7	Our Proposed Method	C4.5 decision tree	100%

8. Conclusion

Breast cancer is a fatal disease if not detected and diagnosed early and can be fatal for the patient's life. The diagnosis and detection of breast cancer can be made accurate by using different data mining methods. The proposal of this paper was to use the information acquisition feature selection to select the most important features in referring to benign or malignant tumors directly. Also, in this work, three types of classification methods were based: C4.5 decision tree and support vector machine (SVM), and Naive Bayes (NB) in the diagnosis of breast cancer in the Wisconsin Diagnostic Breast Cancer (WDBC) dataset and compared its performance results. The precision, accuracy, and recall values obtained using the C 4.5 decision tree were the optimal results. This can be leveraged in the healthcare system to make diagnosis quick, accurate, and error-free. The proposed model can prove to be very useful in this process.

Funding: “This research received no external funding”

Conflicts of Interest: “The authors declare no conflict of interest.”

References

- [1] Hyuna Sung, Jacques Ferlay, MSc, ME2; Rebecca L. Siegel, Mathieu Laversanne, Isabelle Soerjomataram, Ahmedin Jemal, Freddie Bray, “Global Cancer Statistics 2020: GLOBOCAN Estimates of Incidence and Mortality Worldwide for 36 Cancers in 185 Countries,” CA CANCER J CLIN 2021;71:209–249.
- [2] Samuel O. Azubuiké, Colin Muirhead, Louise Hayes and Richard McNally, Rising global burden of breast cancer: the case of sub-Saharan Africa (with emphasis on Nigeria) and implications for regional development,” Azubuiké et al. World Journal of Surgical Oncology, 16-63, 2018.
- [3] Agarap, A. F. M., On breast cancer detection: an application of machine learning algorithms on the wisconsin diagnostic dataset. In Proceedings of the 2nd international conference on machine learning and soft computing, 5-9, 2018.
- [4] Ak, M. F., A comparative analysis of breast cancer detection and diagnosis using data visualization and machine learning applications. Healthcare, 8(2), 1-11, 2020.

- [5] Wang, H., & Yoon, S. W., Breast cancer prediction using data mining method. In IIE Annual Conference. Proceedings, 8-18, 2015.
- [6] Street, W. N., Wolberg, W. H., & Mangasarian, O. L., Nuclear feature extraction for breast tumor diagnosis. In Biomedical image processing and biomedical visualization, 1905, 861-870, 1993.
- [7] Aalaei, S., Shahraki, H., Rowhanimanesh, A., & Eslami, S., Feature selection using genetic algorithm for breast cancer diagnosis: experiment on three different datasets. Iranian journal of basic medical sciences, 19(5), 476, 2016.
- [8] Liu K., Kang G., Zhang N., Hou B., Breast cancer classification based on fully-connected layer first convolutional neural networks. IEEE Access, 6, 23722-23732, 2018.
- [9] Al Bataineh, A., A comparative analysis of nonlinear machine learning algorithms for breast cancer detection. International Journal of Machine Learning and Computing, 9(3), 248-254, 2019.
- [10] Khairunnahar L., Hasib M. A., Rezanur R. H. B., Islam M. R., Hosain M. K., Classification of malignant and benign tissue with logistic regression. Informatics in Medicine Unlocked, 16, 100189, 2019.
- [11] Rahman M. A., Muniyandi R. C., An Enhancement in cancer classification accuracy using a two-step feature selection method based on artificial neural networks with 15 neurons. Symmetry, 12(2), 271, 2020.
- [12] Idri A., Bouchra E. O., Hosni M., Abnane, I., Assessing the impact of parameters tuning in ensemble based breast cancer classification. Health and Technology, 10(5), 1239-1255, 2020.
- [13] Alharbi AH et al., Diagnosis of Monkeypox Disease Using Transfer Learning and Binary Advanced Dipper Throated Optimization Algorithm. Biomimetics, 8(3), 313, 2023.
- [14] Ibrahim Abdelhameed, El-Sayed M. El-kenawy, Applications and datasets for superpixel techniques: A survey. Journal of Computer Science and Information Systems, 15(3), 1-6, 2020.
- [15] M. Saber, Efficient phase recovery system, Indonesian Journal of Electrical Engineering and Computer Science (IJECS), 5(1), 123-129, 2017.
- [16] M Saber, Y Jitsumatsu, MTA Khan, A simple design to mitigate problems of conventional digital phase locked loop, Signal Processing: An international journal (SPIJ), 6(2), 65-77, 2012.
- [17] Mohamed Saber, A novel design and Implementation of FBMC transceiver for low power applications, Indonesian Journal of Electrical Engineering and Informatics (IJEI), 8(1), 83-93, 2020.
- [18] Amin Samy, Sayed A. Ward, Mahmud N Ali, Conventional Ratio and Artificial Intelligence (AI) Diagnostic methods for DGA in Electrical Transformers. International Electrical Engineering Journal, 6, 2096-2102, 2015.
- [19] Al-Salihy, N. K., & Ibrikci, T., Classifying breast cancer by using decision tree algorithms. In Proceedings of the 6th International Conference on Software and Computer Applications, 144-148, 2017.
- [20] Mohamed A. Abouelatta, et al. , Measurement and assessment of corona current density for HVDC bundle conductors by FDM integrated with full multigrid technique. Electric Power Systems Research, 199, 2021.
- [21] Venkatesan E., Velmurugan T., Performance analysis of decision tree algorithms for breast cancer classification. Indian Journal of Science and Technology, 8(29), 1-8, 2015.
- [22] T. Makarovskikh, A. Salah, A. Badr, A. Kadi, H. Alkattan and M. Abotaleb, Automatic classification Infectious disease X-ray images based on Deep learning Algorithms, 2022 VIII International Conference on Information Technology and Nanotechnology (ITNT), Samara, Russian Federation, pp. 1-6, 2022.
- [23] E. Akbari et al., Improved Salp Swarm Optimization Algorithm for Damping Controller Design for Multimachine Power System. IEEE Access, 10, 82910-82922, 2022.
- [24] M. Abotaleb, T. Makarovskikh, A. Ali Subhi, H. Alkattan and A. O. Adebayo, Forecasting and modeling on average rainwater and vapor pressure in Chelyabinsk Russia using deep learning models," 6th Smart Cities Symposium (SCS 2022), Hybrid Conference, Bahrain, 362-367, 2022.
- [25] H. Alkattan, M. Abotaleb, A. Ali Subhi, O. A. Adelaja, A. Kadi and H. K. Ibrahim Al-Mahdawi, The prediction of students' academic performances with a classification model built using data mining techniques. 6th Smart Cities Symposium (SCS 2022), Hybrid Conference, Bahrain, 353-356, 2022.
- [26] H. K. I. Al-Mahdawi, M. Abotaleb, H. Alkattan, A.-M. Z. Tareq, A. Badr, and A. Kadi, Multigrid Method for Solving Inverse Problems for Heat Equation," Mathematics, 10(15), 2022.
- [27] Doaa S. Khafaga, Hussein Alkattan, Alhumaima A. Subhi, Evaluating the Effect of Optimized Voting Using Hybrid Particle Swarm and Grey Wolf Algorithm on the Classification of the Zoo Dataset, Journal of Journal of Artificial Intelligence and Metaheuristics, 2(1), 2022.

- [28]Louloua M. AL-Saedi,Methaq Talib Gaata,Mostafa Abotaleb,Hussein Alkattan, New Approach of Estimating Sarcasm based on the percentage of happiness of facial Expression using Fuzzy Inference System, Journal of Journal of Artificial Intelligence and Metaheuristics, 1(1), 2022.
- [29]Rawat D., et al., Modeling of rainfalltime series using NAR and ARIMA model over western Himalaya, India. Arab. J. Geosci. 2022, 15, 1696.
- [30]Eid Marwa M, Fawaz Alassery, Abdelhameed Ibrahim, and Mohamed Saber, Metaheuristic optimization algorithm for signals classification of electroencephalography channels. Computers, Materials & Continua, 71(3), 4627-4641, 2022.