

Early Breast Cancer Prediction using Machine Learning and Deep Learning Techniques

¹Ms. Swati B. Patil, Dr. Ritesh V. Patil², Dr. Parikshit N. Mahalle³

¹Department of Computer Engineering
Vishwakarma Institute of Information Technology
Pune, India
swati.patil@viit.ac.in

²Department of Computer Engineering
PEDA's College Of Engineering, Manjari(BK)
Research Guide
Department of Computer Engineering
Vishwakarma Institute of Information Technology
Pune, India
rvpatil3475@yahoo.com

³Department of Computer Engineering
Vishwakarma Institute of Information Technology
Pune, India
parikshit.mahalle@viit.ac.in

Abstract: Breast Cancer (BC) is considered as one of the utmost lethal diseases across the globe that has a very high morbidity and mortality rate. Accurate and early prediction along with diagnosis is one of the most crucial characteristics for the treatment of Breast Cancer. Doctors can have an edge over Breast cancer if they are able to predict it in its early stages using deep learning and machine learning techniques. This paper proposed consists of comparison between the accuracy of various machine learning models like Support vector machine (SVM), K-Nearest Neighbours (KNN), Naïve Bayes (NB), Logistic Regression (LR), Random Forest (RF), Decision Tree (DT), XGB Classifier and deep learning model of Artificial neural networks (ANN) for the precise detection of breast cancer.

The most crucial properties from the database have been chosen using one feature-selection technique. Correlation is also used to choose the most correlated features from the data. Implementing the ANN model consists of one input layer, two hidden layers, and one output layer. All Machine Learning models and ANN model are then applied to selected features. The results demonstrated that the SVM classifier achieved the highest performance with an accuracy of ~98.24%.

Keywords: Breast Cancer, Machine Learning, XGB Classifier, Decision Tree, Naïve Bayes, Logistic Regression, Support Vector Machine, Random Forest, Deep Learning, Artificial Neural Network.

I. INTRODUCTION

The most common malignant tumour, Breast Cancer is responsible for about 10.4% of all cancer-related deaths in females ageing between 20 and 50 [1][3]. As observed in the figure 1, statistics in India show that it accounts for the majority of newly diagnosed cases of cancer and cancer-related fatalities. In the modern world, it must be taken into consideration as a crucial health issue. Specialist physicians have discovered various elements, such as lifestyle, hormonal, and environmental factors, that may raise a person's risk of getting BC. A DNA mutation that has impacted numerous generations of the family affects more than 5%–6% of BC patients. Additional reasons of BC include old age, obesity, and abnormal postmenopausal hormone levels.

These problems have already been addressed by a number of ways, but the accuracy of those strategies has been constrained by noisy data. The dataset used in the proposed method is Wisconsin Diagnostic Breast Cancer (WDBC) that accurately classify BC. Because of this, the main objectives of this research are to use cancer disease features and a good pre-processing model to improve model accuracy for BC diagnosis and predict cancer affection at a preliminary phase. This study's main objective is to offer a simple way for detecting BC. This study meticulously examines current cancer detection techniques, producing amazingly precise and effective results. The structure for the following sections of the article is as follows.

According to data published in December 2020 by International Agency for Research on Cancer (IARC), BC took over Lung Cancer as the leading and most common cancer among women

worldwide. Overall number of cancer diagnoses has than doubled over the last two decades, rising from a projected 10 million in 2000 to 19.3 million in 2020 [2]. According to projections, there will be a further rise in the number of people that will be given cancer diagnoses in the upcoming years, reaching a level that is about 50% greater than in 2040 than it was in 2020. As of 2020, there will be 10 million cancer-related deaths worldwide, up from 6.2 million in 2000. Cancer is the cause of more than one in six fatalities. This highlights the need of supporting both the cancer research and cancer prevention efforts.

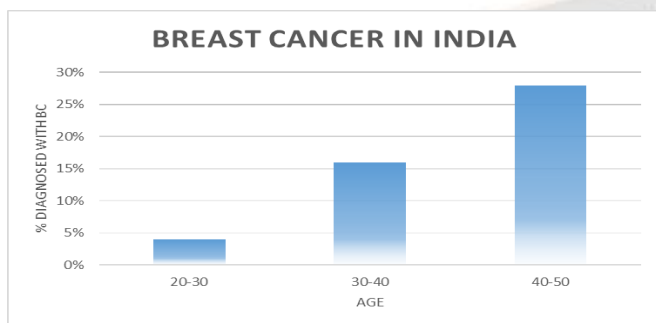


Fig. 1. Breast Cancer in India [3]

Section III of this paper describes the models and methodologies used for forecasting BC. Section IV displays experimental results obtained for very model used. The conclusion and references can be found in Section V and VI of the paper respectively.

II. RELATED WORKS

The research's authors [4] are especially interested in developing supervised prediction models using supervised machine learning techniques for accurately anticipating the consequences of actual illnesses. The results suggest that combining multivariate data with different selecting features, categorization, and dimensionality reduction procedures may result in tools that are beneficial for inference in this field. The 32 distinguishing characteristics of the dataset help condense the multidimensionally enormous dataset to just the few necessary dimensions. SVM, KNN, and LR are three of the pertinent methods, and the SVM has attained a maximum accuracy of 92.7% when compared to other computations. In this strategy, SVM resulted to be the best method for determining the likelihood of developing breast cancer in the presence of complicated data.

This study's goal was to assess a patient's risk of breast cancer using machine learning while taking into account a variety of characteristics, such as age, gender, family history, lumps, resistin, msp.1, leptin, adiponectin and genetics [5]. The patients are divided into groups based on whether or not they

are at risk of acquiring breast cancer. The Random Forest classification technique was used for prediction [6].

In an effort to reduce the fatality rate caused on by the disease, the author of the article [7] examines the previously reported techniques for the investigation for BC classification and diagnosis. The differences between these datasets are examined from a number of perspectives, including input image, size (views included), and classes contained. The CBIS-DDSM, an upgraded version of the DDSM, is said to be the best well-known and complete database.

The innovative models RetinaNet and YOLO that are currently in use and are thought to be more transparent an typical CNN networks, are shown to produce superior outcomes and more accurate performance for bulk identification and malignancy classification. Additionally, compared to models that employ unbalanced-sized data, models that use balanced numbers to determine the classes are more egalitarian and achieve higher levels of overall accuracy.

The development of a model that takes into account all known risk variables presents a significant challenge in the prediction of breast cancer [8]. Without considering other crucial elements, existing prediction models might merely analyse mammographic pictures or demographic risk factors. These models, which can be used to identify high-risk women, may also lead to multiple screenings and invasive ultrasound and MRI imaging sampling. Patients could feel the psychological and financial strain [9].

To improve the grouping reliability of the WDBC dataset, Lavanya and Usha Rani [10] used a half-and-half and dynamic process with 10-crease cross approval. This study demonstrates how machine learning algorithms have a major impact on the diagnosis and outlook for breast cancer. Finding the primary location of the cancer is the main objective of the present study. Therefore, using specialised methods is necessary for early breast cancer identification.

An accuracy rate of 90% or better was seen in the most of the investigations, which is considered to be extraordinary.

The approach in this paper is unique, which takes into account a variety of calculations and techniques to reach an accuracy of ~98% that is higher than that of preceding distributions. ANN, SVM, DT and LR, to name a few, each achieved ~97%, 98%, 95%, 96% percent F1-scores respectively. It is clear from the accuracy% of the models employed in this research that they are more reliable than the models utilised in prior studies. Numerous confirmed model correlations can be found, and the review research can be used to derive the method.

According to studies, the situation may become better if women can detect BC and start treatment at an early stage. To do this,

they must precisely predict how the ailment will progress from a benign state to BC. Machine learning technology can aid in making accurate predictions at an early stage. Although there are several machine learning systems, they all provide incorrect and imprecise predictions. Additionally, they worry about both overfitting and underfitting. As a result, the developed model is used as an aid to medical technicians in applying machine learning to detect cancer illnesses as soon as possible. If someone does have breast cancer, it will be confirmed and shown.

III. MATERIALS AND METHODS

ML and DL approaches are being utilized in proposed approach to predict BC. The BCWD data set is trained and evaluated using ML models, particularly the classifiers from the DT, SVM, RF, NB, KNN, LR, and XGB families. By using the batch normalisation technique, the DL method provides the ANN model by normalising the inputs to a layer for each mini-batch. Phases of the proposed system are as follows: feature selection, database spitting, layer construction, compilation, fitting, model training, and model evaluation.

3.1 Methodology

In the proposed methodology to predict BC, standard ML and DL approaches are used. The BCWD data set is trained and evaluated using the conventional ML approach using seven ML algorithms, including NB, RF, XGBoost, DT, LR, SVM, and KNN. Correlation feature selection technique is used to improve model performance. The DL method proposes a Keras-based ANN model. The proposed system's steps, which comprise the feature-selection method, the spitting database, modelling training, and model evaluation, are depicted in Figure 2.

of data. Breast cancer can be classified as benign (B) or malignant (M) on a scale of 0 to 1.

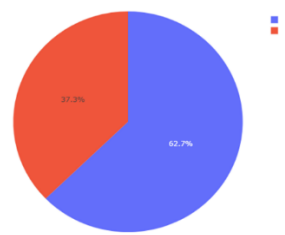


Fig. 3. Pie chart by Plotly

3.1.3 Selection of Features

Although dataset frequently comprises a large number of features but it is not necessary that all of them may be useful in building a machine-learning model that can make the necessary prediction. Forecasts may be erroneous as a result of using some of the features. Therefore, while developing a machine-learning model, properly selecting characteristics is essential.

The correlation matrix was utilized to cut down on the number of features in this study. Correlation is a statistical concept that refers to the degree to which two variables have a linear connection with one another. High correlation characteristics are more linearly dependent and influence the dependent variable equally. When there is a significant association between the two characteristics, omitting one of features will have no effect. As shown in Figure 3, after utilizing the correlation matrix on a dataset, the seven characteristics that correlate greater than 0.95 were deleted. After eliminating the highly correlated parameters, 26 attributes were chosen from the dataset for further examination.

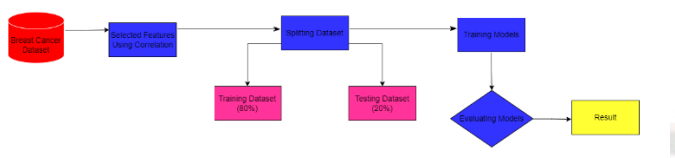


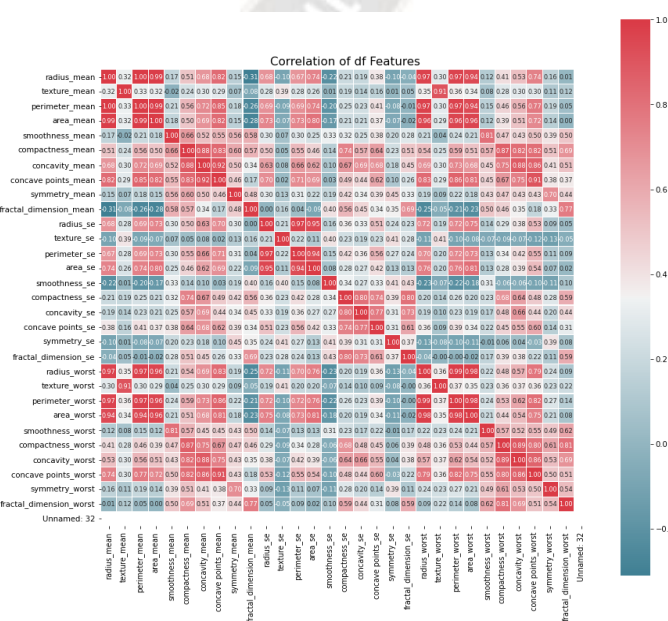
Fig. 2. Methodology

3.1.1 Loading BC Dataset

In this step, the data from the BWCD Data Set are loaded into a Jupyter notebook for further analysis.

3.1.2 Dataset

For the creation of a new and effective method for the identification of BC, a well labelled dataset is a crucial necessity. The Kaggle dataset of WBCD is used in this research work. Features in the dataset were calculated from a digital image of a breast mass obtained via a fine needle aspiration (FNA). They describe the features of the visible cell nuclei in the image. 30 features and 569 samples make up the entire set



3.1.4 Splitting dataset

Training and testing datasets have been created from the dataset. The sklearn model selection has imported the train test split. We split the data in half, with 80% being used for training and 20% being used for testing. The models are trained using training datasets, and the results of each model's CV are then given.

3.1.5 Selection of Model

Model selection involves picking the best machine learning models to train the dataset. In our paper, a variety of classification models, including KNN, NB, DT and SVM, which are employed with the model selection technique. In our dataset, a dependant variable (Y) only has two sets of values: Benign(B) and Malignant(M). Our model utilizes these data to project the outcome based on the different classification methods.

3.2 Proposed Algorithms

3.2.1 Logistic Regression

Logistic regression is one of the ML classification methods for analysing datasets with one or more relationship between the independent variable (IVs) that forecast a result and a dependent variable (DV). The result of logistic regression is translated using the logistic sigmoid function to produce a probability that may be shifted to two or more distinct groups. Logistic regression limits the cost function's range to be between 0 and 1.

3.2.2 Random Forest

An integrated classifier known as Random Forest combines bags with a random selection of features.

Data can be handled by the random forest without pre-processing. It generates the most precise classification of various data sets. A representation of the RF algorithm is shown in Figure 4.

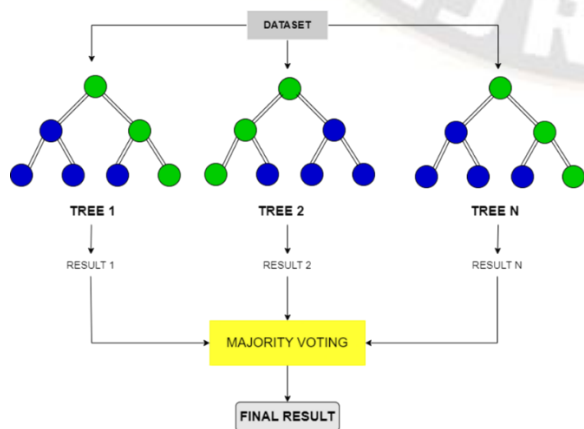


Fig. 5. An illustration of a Random Forest that consists of n different decision trees

3.2.3 Decision Tree

The decision tree is one of the first and most significant machine learning algorithms (DT). A decision tree is a process framework of choosing, testing and aligning the results of grouping data into a tree-like framework. Figure 5 depicts the DT diagram along with its attributes and guidelines.

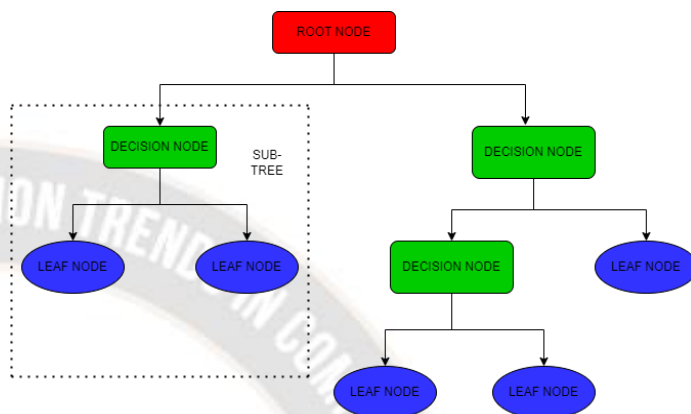


Fig. 6. An illustration of a Decision tree

3.2.4 Support Vector Machine

SVM, is a popular supervised learning technique for problems involving classification and regression. However, Machine Learning Classification difficulties are where it is most commonly applied. In order to facilitate rapid future categorization of new data points, the SVM technique looks to locate the optimal line or distance measure that may split n-dimensional space into classes. This ideal boundary for making decisions is referred to as a hyperplane.

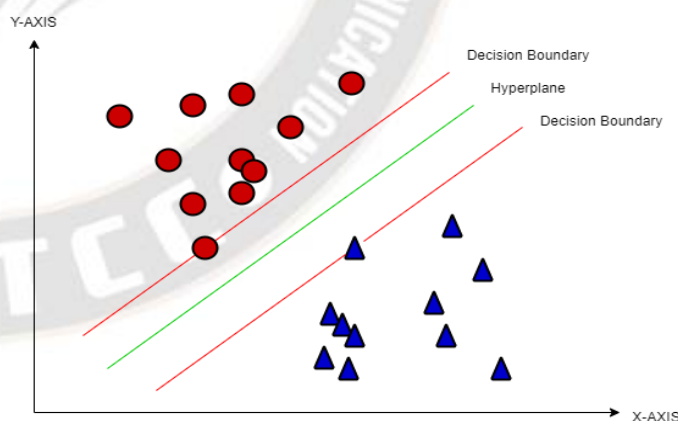


Fig. 7. An illustration of SVM Classifier

3.2.5 Naïve Bayes

Based on the Bayes theorem, the NB machine learning algorithm is supervised learning method used in classification issues. With a large training set, text categorization is where it is most frequently utilised. One of the most straightforward and efficient classification methods available today is the NB. It

facilitates the quick creation of machine learning models that produce predictions. It provides rough approximation based on the odds that an event will occur since it is a classification method. One of the simplest yet most reliable classification systems is the NB classifier. However, it has a number of drawbacks, including the precise categories assigned to the training data[11][13].

3.2.6 K-Nearest Neighbours

A fundamental machine learning method built on the supervised learning approach is K-Nearest Neighbour. The K-NN method, which assumes that the new iteration and the preceding examples are identical, puts the test instances in the group which is the most comparable to the existing categories. Once all previous data has been captured, an unique data point is classified using the K-NN technique based on similarity. This indicates that the K-NN approach can swiftly classify new data.

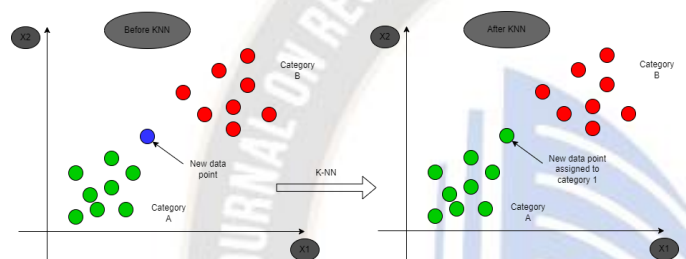


Fig. 8. An illustration of KNN Algorithm

3.2.7 XGB Classifier

The collection of gradient boosting methods known as eXtreme Gradient Boosting, or XGBoost, is intended for usage with current data science tools and issues. The major characteristics of XGBoost include being incredibly parallelizable and scalable, quick to execute, and frequently outperforming competing algorithms. Additionally, in order to decrease overfitting and improve performance, it uses a more regularised model formalisation.

3.2.8 Artificial Neural Network

3.2.8.1 Encoding

The encoding approach offers the fully-connected neural network (FCNN) training speedup. LabelEncoder is used to handle categorical variables in ANN. For this label encoder is imported from the sci-kit learn library.

3.2.8.2 Scaling the Data

Since the node weights and input data are multiplied by the machine learning algorithm over a long period of time, the data must be scaled down when building the artificial neural network. The data is scaled to reduce that period of time. To scale the data, StandardScaler is used.

3.2.8.3 Creating Layers for ANN model

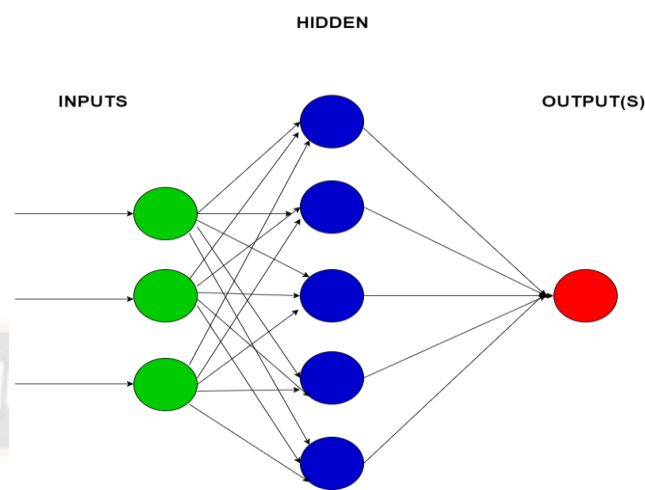


Fig. 9. ANN Architecture

Computer models of biological neural networks' structure and operation are known as "artificial neuron networks" (ANN). Information travelling through the network has an impact on how the ANN is built since a neural network adapts or learns based on input and output. ANNs are thought of as instruments for modelling or identifying patterns in the complex interactions between inputs and outputs in nonlinear statistical data analysis. An ANN is also known as a neural network. The Keras and its packages have therefore been imported. After importing the Keras libraries, we create three different layers types: one output and input, and two hidden layers.

The ANN network consists of one output layer and output layer each along with two hidden layers. Activation function was employed as the relu for both the hidden and input layers. In the output layer, sigmoid activation function is utilised. The sigmoid activation function is employed when dealing with classification problems involving two different types of results. The sigmoid function is used when there are three or more categorization results.

3.2.8.4 Compiling and Fitting

The optimizer is chosen as adam for gradient descent. Binary_crossentropy is the loss function used.

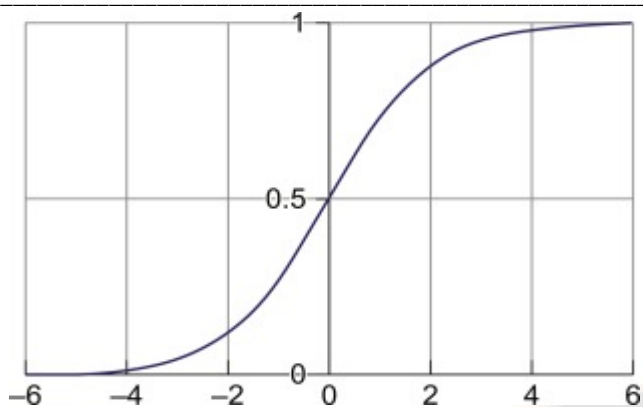


Fig. 10. Compiling and Fitting of ANN

IV. RESULTS

In result evaluation, the accuracies of various classification methods are assessed which are used for prediction of BC in this research paper. JUPYTER notebook was used for this study with an Intel Core i5 dual-core CPU operating with 8GB of RAM and 1.6 GHz. Machine Learning algorithms: SVM, K-NN, NB, LR, RF, DT, XGB Classifier with an ANN Deep Learning algorithm were used to make the prediction. The data was split into 80-20 models for testing and training and used the train-test splitting technique. Keras library was used to implement the ANN model. The models for ML were implemented using sci-kit-learn package. A test dataset (20%) was used to assess the model, while a training data set (80%) was utilised to optimise model performance and record the findings of cross-validation (CV). First, the correlation between the features was examined and elimination of any parameters that had a high correlation (above 95%) with other features was carried out. The chosen features were then subjected to correlation using both the standard ML models and the ANN models. Many ANN parameters were changed for each testing of 32 batch size and epochs. In-depth documentation will be kept of each experiment's testing procedure and CV outcomes.

Table 1 displays the model's accuracy as well as its Recall, F1 Score and Precision.

Model	ACU	PRE	REC	F1
SVM	98.24	97.40	100.00	0.98
ANN	96.92	96.00	99.00	0.97
KNN	95.61	93.75	100.00	0.96
NB	95.61	96.05	97.33	0.96
LR	95.61	97.33	96.00	0.96
RF	94.73	97.26	94.66	0.95

XGB Classifier	94.73	97.26	94.66	0.95
DT	93.85	97.22	93.33	0.95

Table 1: Results of an individual model with variations in the dataset

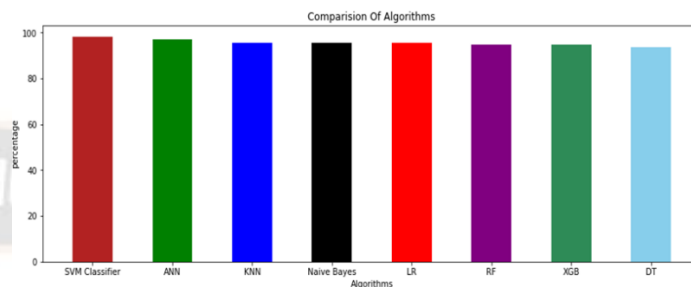


Fig. 11. Bar chart of algorithms using pyplot.

CV Results- The ANN model used for deep learning reported a score of 96.92. SVM Machine Learning approach registered the best performance, scoring 98.24%, while the Decision Tree scored 93.85 percent, which accounted for the lowest performance.

V. CONCLUSION

In this study, a number of ML models as well as a DL model are used to identify breast cancer. All the Machine Learning models that were exhibited worked well in identifying benign and malignant breast cancer. Logistic Regression, SVM, KNN, XGB Classifier, RF, NB are the machine learning algorithms that were utilised. An ANN model containing an one output layer along with two layers of hidden neurons, and an input layer was developed for deep learning. To consistently normalise inputs to these layers, batch normalisation method is utilized. In order to choose the important database properties, correlation matrix was used as a feature selection technique. The most important attributes are sorted out and are utilized while applying ML models and ANN. According to the findings, the SVM classifier had the optimum performance.

REFERENCES

- [1] R. L. Siegel, K. D. Miller, and A. Jemal, "Cancer statistics, 2019," CA: A Cancer Journal for Clinicians, vol. 69, no. 1, pp. 7–34, 2019.
- [2] 'WHO | Breast cancer', WHO. <http://www.who.int/cancer/prevention/diagnosis-screening/breast-cancer/en/> (accessed Feb. 18, 2020).
- [3] <https://cancerconsultindia.com/blog/breast-cancer-statistics-rise-of-breast-cancer-in-india>
- [4] Vyas, Sarthak & Chauhan, Abhinav & Rana, Deepak & Ansari, Mohd. (2022). Breast Cancer Detection Using Machine Learning Techniques. International Journal for Research in

- Applied Science and Engineering Technology. 10.10.22214/ijraset.2022.43055.
- [5] Cruz, J. A., & Wishart, D. S. (2006). Applications of machine learning in cancer prediction and prognosis. *Cancer informatics*, 2, 117693510600200030.
- [6] Rawal, R. (2020). Breast cancer prediction using machine learning. *Journal of Emerging Technologies and Innovative Research (JETIR)*, 13(24), 7.
- [7] Karthik Moudgalya Umesh, Abdul Rahman Bin S. Senathirajah, R. A. Sheedul Haque, Gan Connie. (2023). Examining Factors Influencing Blockchain Technology Adoption in Air Pollution Monitoring. *International Journal of Intelligent Systems and Applications in Engineering*, 11(4s), 334–344. Retrieved from <https://ijisae.org/index.php/IJISAE/article/view/2673>
- [8] Charbuty, Bahzad and Adnan Mohsin Abdulazeez. "Classification Based on Decision Tree Algorithm for Machine Learning." (2021).
- [9] Brédart A, Kop JL, Antoniou AC, Cunningham AP, De Pauw A, et al. Clinicians' use of breast cancer risk assessment tools according to their perceived importance of breast cancer risk factors: an international survey. *J Community Genet* . 2019;10(1):61–71. doi: 10.1007/s12687-018-0362-8.
- [10] Yala, A., Lehman, C., Schuster, T., Portnoi, T., & Barzilay, R. (2019). A deep learning mammography-based model for improved breast cancer risk prediction. *Radiology*, 292(1), 60–66.
- [11] Lavanya, D., & Rani, K. U. (2012). Ensemble decision tree classifier for breast cancer data. *International Journal of Information Technology Convergence and Services*, 2(1), 17–24.
- [12] Kathija and S. Shajun Nisha. "Breast Cancer Data Classification Using SVM and Naïve Bayes Techniques." (2017).
- [13] Bharati, S., Podder, P., & Mondal, M. (2020). Artificial neural network-based breast cancer screening: a comprehensive review. *arXiv preprint arXiv:2006.01767*.
- [14] Karabatak, M. (2015). A new classifier for breast cancer detection based on Naïve Bayesian. *Measurement*, 72, 32–36.
- [15] Sadoughi, F., Kazemy, Z., Hamedan, F., Owji, L., Rahmanikatifari, M., & Azadboni, T. T. (2018). Artificial intelligence methods for the diagnosis of breast cancer by image processing: a review. *Breast Cancer: Targets and Therapy*, 10, 219.
- [16] Roslidar, R., Rahman, A., Muharrar, R., Syahputra, M. R., Arnia, F., Syukri, M., ... & Munadi, K. (2020). A review on recent progress in thermal imaging and deep learning approaches for breast cancer detection. *IEEE Access*, 8, 116176–116194.
- [17] Ferroni, P., Zanzotto, F. M., Rioldino, S., Scarpato, N., Guadagni, F., & Roselli, M. (2019). Breast cancer prognosis using a machine learning approach. *Cancers*, 11(3), 328.
- [18] Chand, S. (2020). A comparative study of breast cancer tumor classification by classical machine learning methods and deep learning method. *Machine Vision and Applications*, 31(6), 1–10.
- [19] Chugh, G., Kumar, S., & Singh, N. (2021). Survey on machine learning and deep learning applications in breast cancer diagnosis. *Cognitive Computation*, 13(6), 1451–1470.
- [20] Battineni, G., Chintalapudi, N., & Amenta, F. (2020). Performance analysis of different machine learning algorithms in breast cancer predictions. *EAI Endorsed Transactions on Pervasive Health and Technology*, 6(23), e4–e4.
- [21] Ara, S., Das, A., & Dey, A. (2021, April). Malignant and benign breast cancer classification using machine learning algorithms. In *2021 International Conference on Artificial Intelligence (ICAI)* (pp. 97–101). IEEE.
- [22] Paul Garcia, Ian Martin, Laura López, Sigurðsson Ólafur, Matti Virtanen. *Predictive Analytics in Education: Leveraging Machine Learning for Student Success*. *Kuwait Journal of Machine Learning*, 2(1). Retrieved from <http://kuwaitjournals.com/index.php/kjml/article/view/164>
- [23] Osareh, A., & Shadgar, B. (2010, April). Machine learning techniques to diagnose breast cancer. In *2010 5th international symposium on health informatics and bioinformatics* (pp. 114–120). IEEE.
- [24] Benbrahim, H., Hachimi, H., & Amine, A. (2020). Comparative study of machine learning algorithms using the breast cancer dataset. In *International conference on advanced intelligent systems for sustainable development* (pp. 83–91). Springer, Cham.
- [25] Islam, M., Haque, M., Iqbal, H., Hasan, M., Hasan, M., & Kabir, M. N. (2020). Breast cancer prediction: a comparative study using machine learning techniques. *SN Computer Science*, 1(5), 1–14.
- [26] Fröhlich, H., Patjoshi, S., Yeghiazaryan, K., Kehrer, C., Kuhn, W., & Golubnitschaja, O. (2018). Premenopausal breast cancer: potential clinical utility of a multi-omics-based machine learning approach for patient stratification. *EPMA Journal*, 9(2), 175–186.
- [27] Kumar, A., & Poonkodi, M. (2019). Comparative study of different machine learning models for breast cancer diagnosis. In *Innovations in soft computing and information technology* (pp. 17–25). Springer, Singapore.