



# Analytical Overview of Machine Learning Algorithms in Breast Cancer Screening: Clinical Workflow, Applications and Research Gaps

Nishant Tripathi

Dept. of Cyber Physical System, School of Computer Science and Engineering, Faculty of Engineering and Technology, JAIN (Deemed-to-be University), Bangalore, Karnataka

**Abstract:** Breast cancer remains one of the leading causes of mortality among women worldwide. Early detection through screening is pivotal in improving survival rates and treatment outcomes. Over recent years, machine learning (ML) algorithms have emerged as powerful tools in enhancing breast cancer screening processes. This review paper provides a comprehensive analysis of the state-of-the-art ML algorithms employed in breast cancer screening. We explore various supervised, unsupervised, and reinforcement learning techniques, assessing their effectiveness in image analysis, risk prediction, and diagnostic accuracy. Key contributions include a detailed examination of convolutional neural networks (CNNs) in mammogram analysis, the role of support vector machines (SVMs) and random forests (RFs) in feature extraction and classification, and the application of ensemble methods in improving prediction robustness. Additionally, we discuss the integration of ML algorithms with clinical workflows, highlighting challenges such as data heterogeneity, interpretability, and ethical considerations. Through this analytical review, we aim to provide insights into the current landscape of ML applications in breast cancer screening, identify gaps in existing research, and suggest directions for future studies to enhance the efficacy and reliability of these technologies in clinical practice.

**Index Terms:** Deep Learning, Breast Cancer, Machine Learning (ML) Algorithms, Convolutional Neural Network, Support vector Machine

## I. INTRODUCTION

Deep learning is a subcategory of technology for artificial intelligence in which techniques process massive amounts of data to perform detection, gain knowledge from them, and execute activities spontaneously without being commanded. In current history, the increasing prevalence of convincing equipment and utility computing has resulted in a widespread implementation of Machine Learning in various spheres of human existence, tend to range from applying it for social networking guidelines to implementing it for workflow robotization in factories. And its prominence will broaden.[1]

Breast cancer is the frequently occurring illness in women. In the event of Indian women, there is a 50% chance of death as 1 of every 2 women detected with breast cancer dies.[2] Breast cancer ranks among the deadliest diseases and diverse diseases of our time, killing an a staggering amount of women all over the world After illness, it is the another reason to cause of mortality throughout women.[3] There are several kinds of ML. [4] and Carcinoma is anticipated using content extraction methodologies.

Another of the most crucial functions seems to be to find the most suitable and applicable automated system for cancer prediction. Carcinoma of breast generates from cancerous tumors, if there is overgrowth becomes uncontrollable. Breast forming a tumor when a high proportion of fatty and collagenous tissues in the breast begin to multiply atypically. Cancer cells scattered all across tumours, going to result in cancer stages. [3]

## II. MACHINE LEARNING ALGORITHMS FOR BREAST CANCER SCREENING

Automatic gaining of knowledge is an aspect of ML [6], the methodologies are intended to learn from prior data points; we insight a huge quantity of information, the ML framework explores that information, and based on that train model, we can make estimates of upcoming. [7], [8], [9] The following are the most critical supervised ml methodologies for foretelling breast cancer:



2.1 Artificial Neural Network (ANN)

A common data resource extraction algorithm is the Deep Convolutional Network. A neural network comprises three aspects: fully connected units, concealed units, and target output. This method is used to identify trends that are overly complex.

The automated system employs virtualization technology, memory distribution, network architecture, and a collective solution. [13] The ANN framework is the greatest extensively applied in CAD for tomosynthesis perception and biopsy making a choice. ANN is used in mammography interpretation in two ways: first, by implementing a classifier to the coloured region of interest (ROI) picture, and 2nd, by comprehending the circumstances utilizing data obtained from which was before picture signals[14].

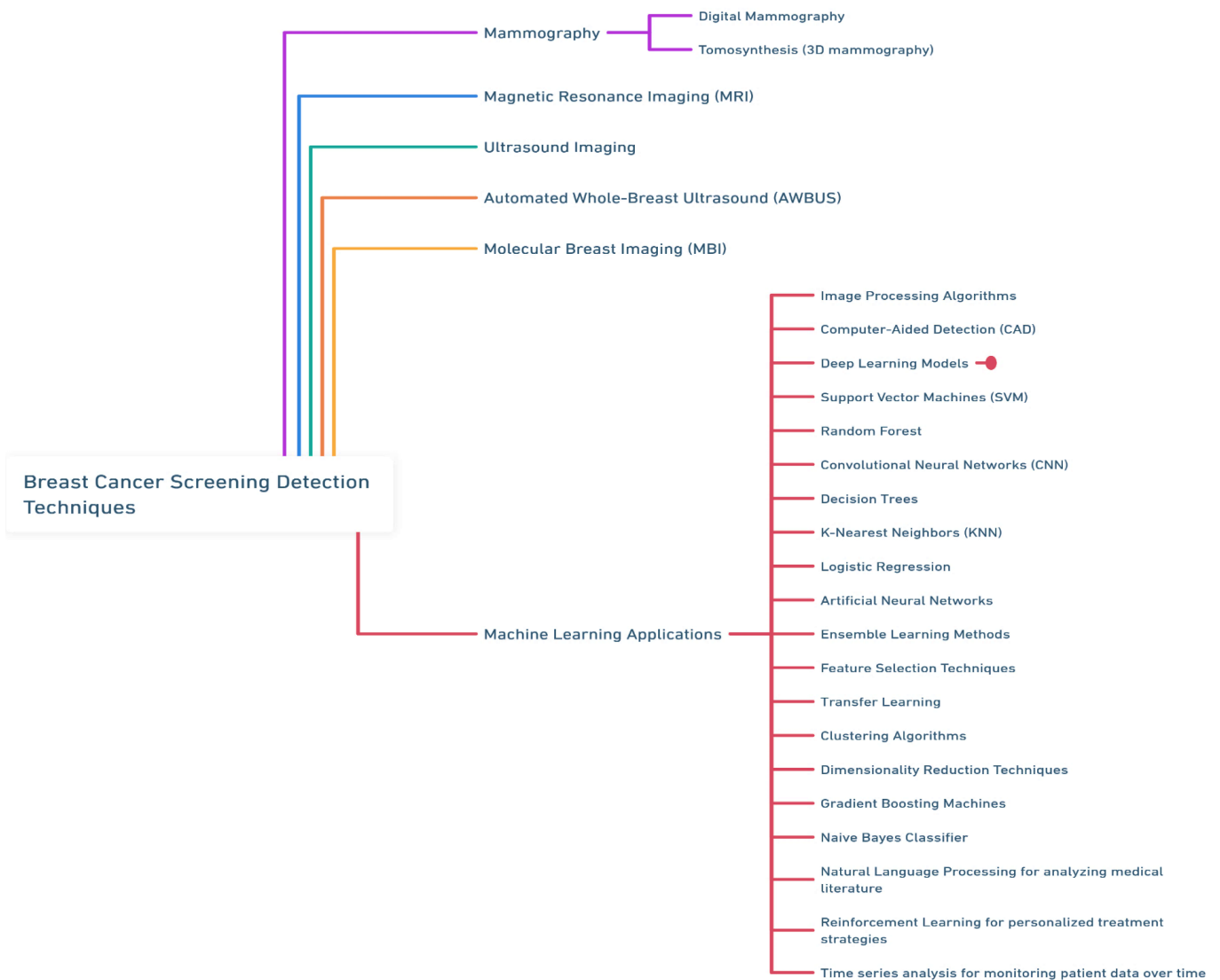
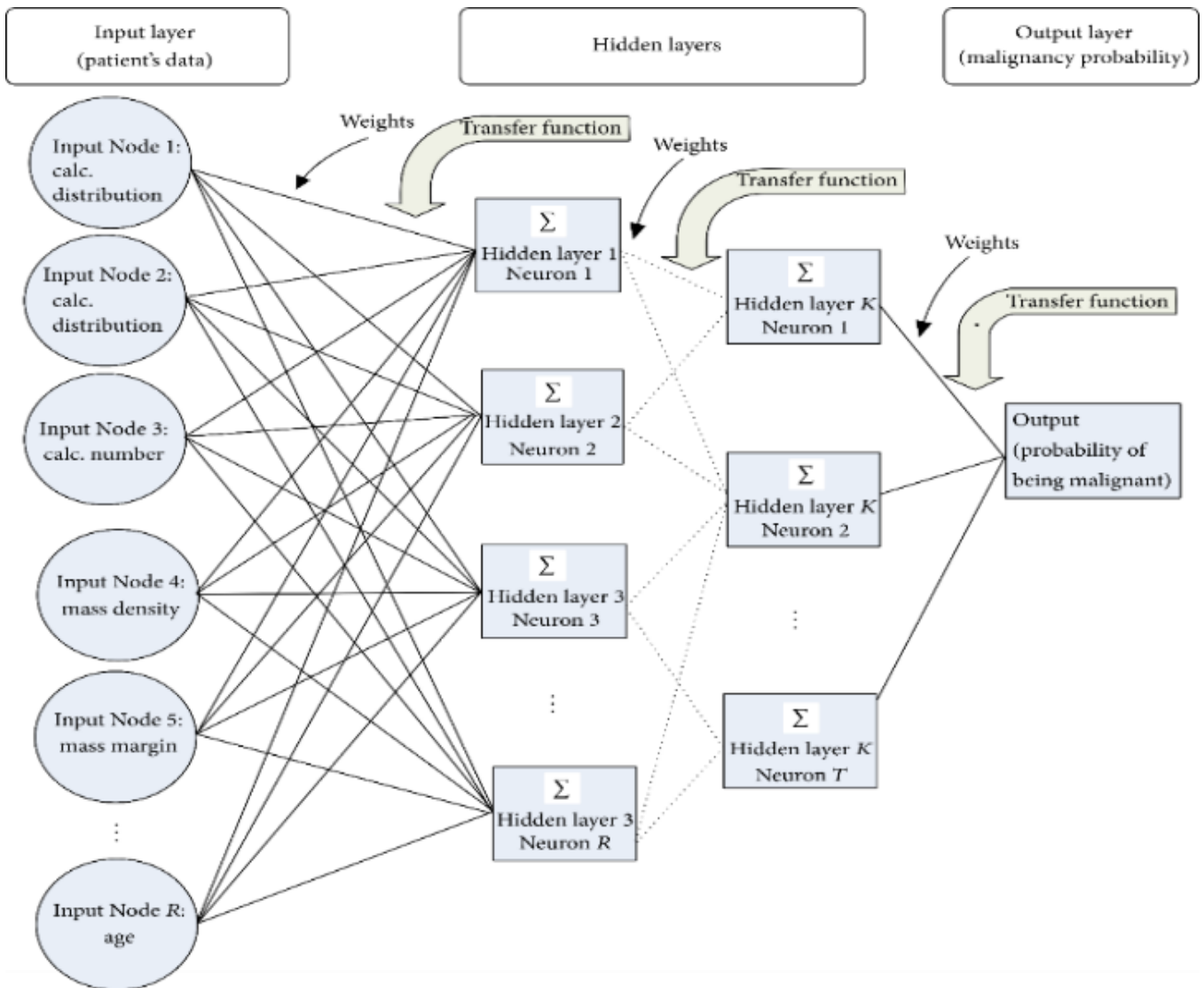


Figure 1: Display of massive kinds of breast cancer.[5]



**Figure 2:** A classic ANN formation for mammography malignant lesions categorization [14]

**2.2 Logistics Regression (LR)**

It is an algorithm for supervised learning with a higher proportion of dependent variables. This algorithm generates a binary value. Logistic support regression [15] can develop a continuous result from a given data set. A mathematical analysis and binary variables encompass this algorithm.

[16] Logistic is a type of regression model that anticipates the risk is the possibility that a specific piece of information or entrance falls into one of several sections. [17]. The assertion beside logistic regression is that the data follows a weight vector. In logistic regression, the sigmoid activation function is used to analyze the data.

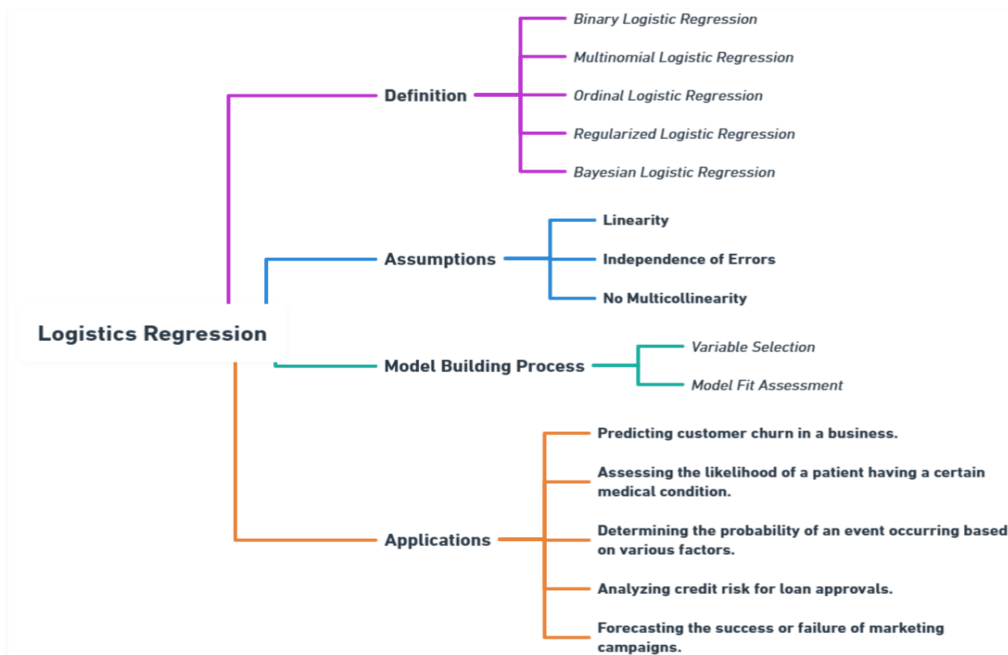


Figure 3 Implementation of Logistics Regression

2.3 K-Nearest Neighbor (KNN)

This technique is utilized to identify patterns. It is a valuable predictor of breast cancer. Identifying the pattern has been given equal weight in each class. K Closest Neighbor [18] Extract comparable depicted data from an extensive dataset. We employ feature similarity to classify a large dataset. [16] K can be regarded as an acknowledgment of the training information remarks comparable to the screening piece of data that information remarks to figure out the class. A k-nearest-neighbor method computes where a training dataset relates based on the information sets around it. The technique is a supervised training that employs regression and grouping. KNN gathers all nearby pieces of information prior to actually having to process a new one. The distance is established by features with a slightly elevated degree of variation. [2]

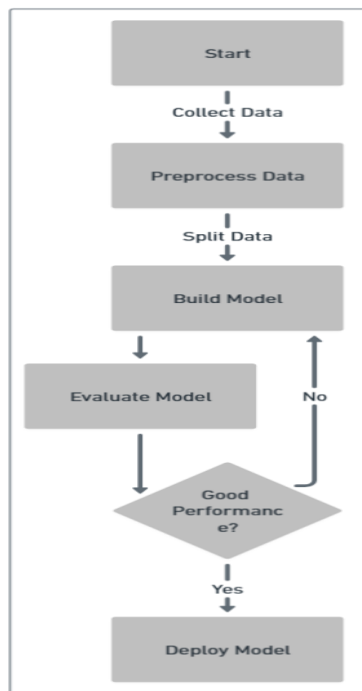


Figure 4 Flowchart of KNN Implementation



**2.4 Decision Tree (DT)**

A classification and a regression design are utilized in the DT [19]. The original data is splitted into smaller groupings. These smaller sets of information permit for the most reliable estimates. CART [20], C4.5 [21], C5.0 [22], and conditional tree [16], [23] are all decision tree methods. Decision trees are a powerful technique used in several disciplines, along with machine learning, image recognition, and analytical thinking [24]. DT are a collection of phases that proficiently and harmoniously unite a sequence of basic judgements in which each test matches a quantitative attribute to a predicted values [25].

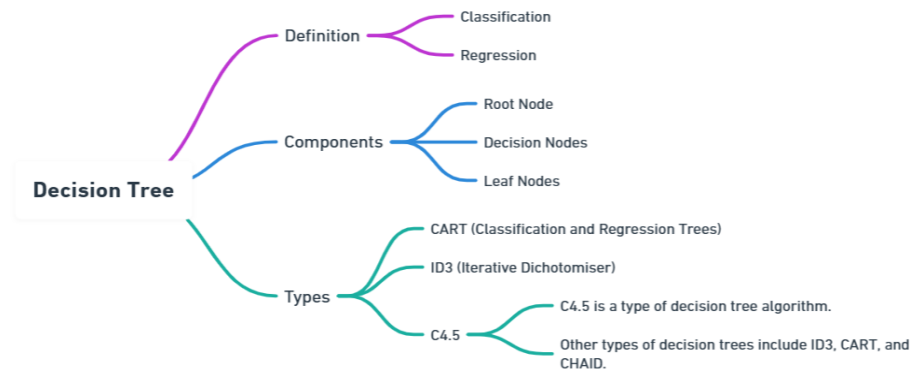


Figure 5 Implementation of Decision Tree Algorithm

**2.5 Naive Bayes Algorithm (NB)**

This pattern suggests a huge training source data. The methodology employs the Bayesian method to calculate probability.[26] It provides the best performance when estimating the probabilities of high dimensionality as an input [27]. It is a metaphor classifier that relates the training dataset to the training tuple [16]. Bayesian Networks are refers generally that utilize the Bayes theorem to think critically. It is delusional because it presumes that all attributes are distinct from each other, which is rarely the case in practical systems, but Nave Bayes is suitable for an extensive spectrum of machine learning problems. [2]

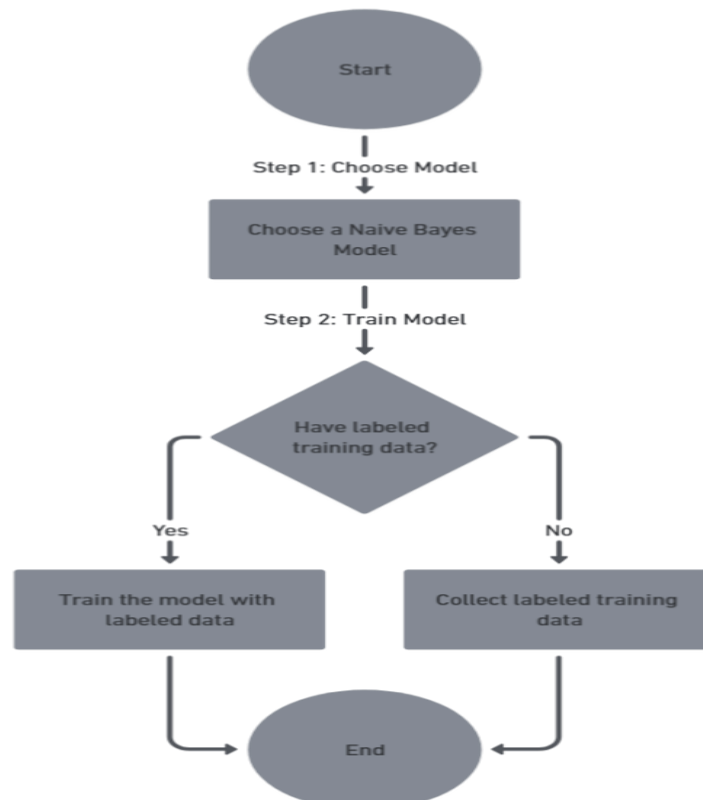


Figure 6 Flowchart of Naïve Bayes Algorithm



2.6 Support Vector Machine (SVM)

It is a supervised learning algorithm that can fix regression and classification issues [28]. It is composed of conceptual and statistical functions that are used to solve the linear model. It has the best accuracy when forecasting huge data. It is an impactful machine learning technique premised on three - dimensional and two-dimensional modelling [16], [29]. A separating hyperplane defines the SVM, a discriminative classifier. The concept of hyperplane contains a generalized statement of the maximal margin classifier. In a n-dimensional space, the centroid has (n-1) dimension and a plain feature space that does not require to cross over the origin. Even though the hyperplane is not viewed in extra dimensionality, the theory of a (n-1) dimensional plain dimension continues to remain. If no time series has a lower dimensional higher dimensional space, a linear descriptor can be formed. To generate a nonlinear classifier, the kernel trick must be implemented to maximum-margin hyperplanes. [30]

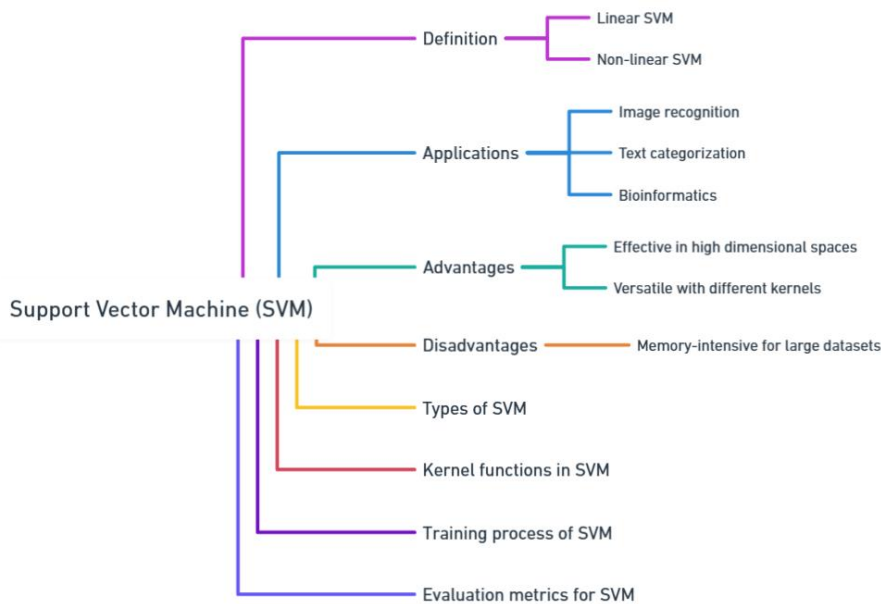


Figure 7 Implementation of SVM

Algorithm	Parameter	Description	Contribution
Support Vector Machines (SVMs)	Kernel Functions	Functions like linear, polynomial, and RBF kernels used to transform input data into higher-dimensional space.	Enhances the capability to handle non-linear relationships in the data, improving classification accuracy.
	Hyperparameters	Parameters such as C (regularization) and gamma (kernel coefficient).	Optimization of these parameters fine-tunes the model for better performance and generalization.
	Support Vectors	Data points that lie closest to the decision boundary.	These vectors define the decision boundary, crucial for accurate classification.
	Margin	Distance between the closest support vectors and the decision boundary.	Maximizing this margin improves the model's robustness and generalization ability.
Role in Feature Extraction	Feature Scaling	Techniques like normalization or standardization applied to input data.	Ensures that features contribute equally to the decision boundary, improving model performance.
	Feature Selection	Techniques such as Recursive Feature Elimination (RFE) used to select the most relevant features.	Reduces dimensionality, enhancing model efficiency and interpretability.



Algorithm	Parameter	Description	Contribution
Role in Classification	Binary and Multi-Class	SVMs can be extended to handle multi-class classification through strategies like One-vs-One or One-vs-Rest.	Versatile in addressing various classification problems within breast cancer screening.
	Decision Boundary	Hyperplane that separates different classes in the feature space.	Provides a clear and interpretable separation between benign and malignant cases.
Random Forests (RFs)	Number of Trees	Total number of decision trees in the forest, commonly ranging from 100 to 1000.	Increasing the number of trees typically enhances prediction accuracy and robustness.
	Max Depth	Maximum depth of each tree in the forest.	Controls model complexity and helps prevent overfitting.
	Feature Selection	Random selection of features at each split during tree construction.	Ensures diversity among trees, improving overall model robustness and reducing overfitting.
Role in Feature Extraction	Feature Importance	Measurement of the significance of each feature in predicting the target outcome.	Identifies the most influential features, aiding in interpretability and further feature selection processes.
	Handling Missing Values	Robust against missing data, as not all trees are affected by missing values in certain features.	Enhances model robustness and reliability when dealing with incomplete datasets.
Role in Classification	Ensemble Voting	Aggregation of predictions from all trees in the forest, typically through majority voting.	Improves overall prediction accuracy and stability.
	Class Probability Estimates	Ability to provide class probabilities rather than just binary outcomes.	Offers nuanced insights into the confidence of predictions, aiding in clinical decision-making.
Integration with Clinical Workflows	Model Interpretability	Techniques like feature importance and partial dependence plots.	Enhances transparency and trust in model predictions among clinicians.
	Scalability	Ability to handle large datasets and high-dimensional data efficiently.	Suitable for large-scale screening programs, ensuring practical utility in clinical settings.
Ethical Considerations	Bias Mitigation	Strategies to address biases in training data and model predictions.	Ensures fairness and equity in screening outcomes across diverse patient populations.
	Privacy and Security	Measures to protect patient data privacy and ensure secure handling of sensitive information.	Maintains compliance with regulations and fosters patient trust.

Table.1 outlines the critical parameters and roles of SVMs and RFs in feature extraction and classification, emphasizing their contributions to breast cancer screening.

## 2.7 Random Forest (RF)

Based on supervised learning [31], the Random Forest classifier can solve categorisation and recurrence problems. It is a ML basic foundation that estimates additional knowledge using preceding set of data. [16]The dataset was segmented into two sections, learning and evaluating, with 398 inferences in the learning phase but rather 171 inferences in the testing phase.

The number of estimators is set to 72 to ensure that each analysis is anticipated at least once. It is straightforward that the effects of treatment, diameter mean, surface mean, and exterior mean are compared, while the other factors have a considerable effect but ought to not be overshadowed in order to boost the model's validity. [2]



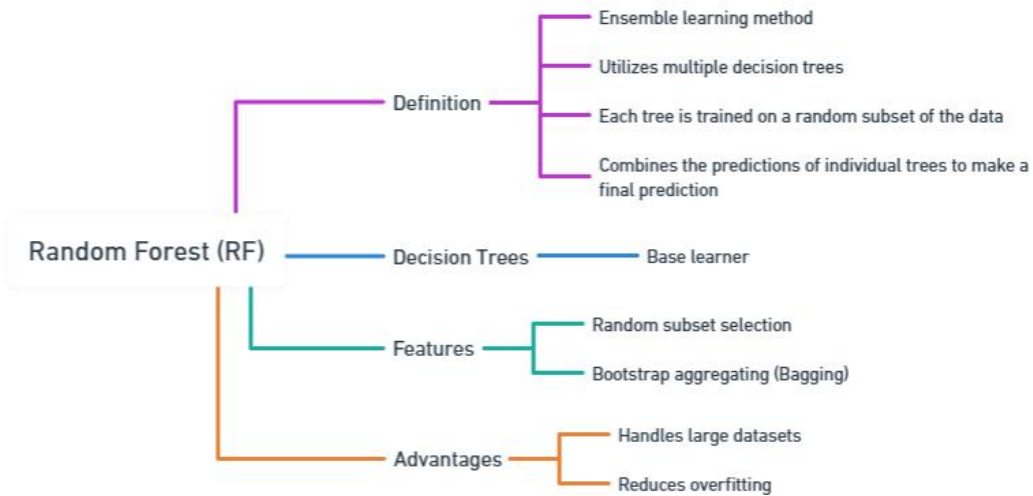


Figure 8 Characteristics of Random Forest Algorithm

**2.8 K Mean Algorithm**

The K mean is a clustering algorithm being used divide data into small factions. To determine the correlation of information extracted, an algorithm is used. Data points include at least yet another group that is especially suitable for processing complex datasets. [32]. Moreover, other studies, such as [14, 15], have proposed traditional methods for categorising breast cancer tumours using K-means and other techniques. [33]

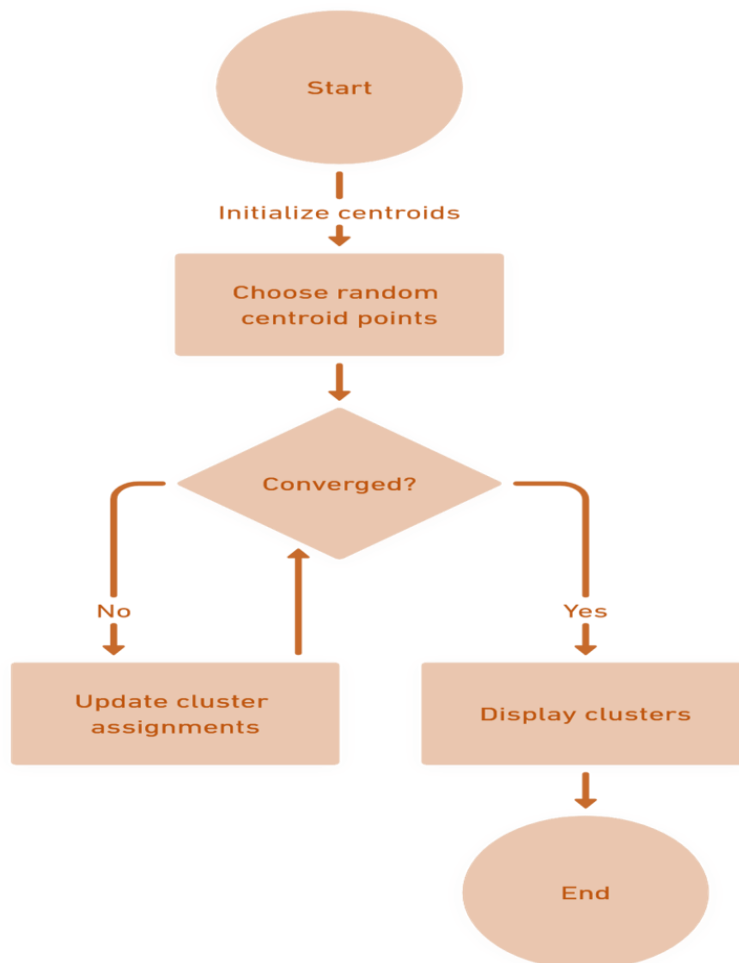


Figure 9 Flowchart of K Means Algorithm





**2.9 C Mean Algorithm**

Patterns are revealed via stability. A family is composed of a collection of equivalent observations. Conversely, each information statement in the C mean method is assigned to a single cluster. It is widely used in medical image segmentation and illness prediction. [34] Fuzzy C-means (FCM) is a data segmentation method that allows a mono piece of information to be linked with 2 or several groups. [35]. It is focused on the reductions of the factual feature to accomplish a decent categorization.

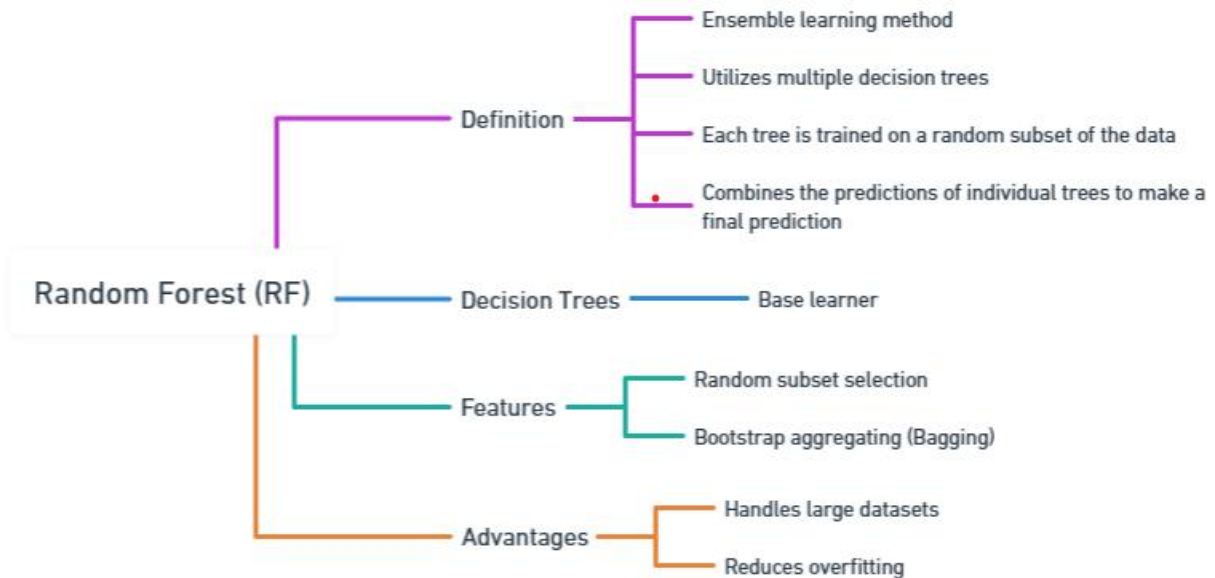


Figure 10 Characteristics of C Means Algorithms

**2.10 Hierarchical Algorithm**

The majority of the time, hierarchical algorithms are used to evaluate crude information in the form of matrices. A hierarchy isolates each group from the others. Every group is made up of data points that are equivalent. A probabilistic model is used to calculate the distance between each cluster.[36] Hierarchical RBF networks (HRBF) are made up of multiple RBF networks that have been assembled in different levels or cascade configurations to divide and solve an issue in even more least 1 phase. Mat Isa et al. used Hierarchical Radial Basis Function (HRBF) to enhance RBF performance in the diagnosis of cervical cancer.[37]

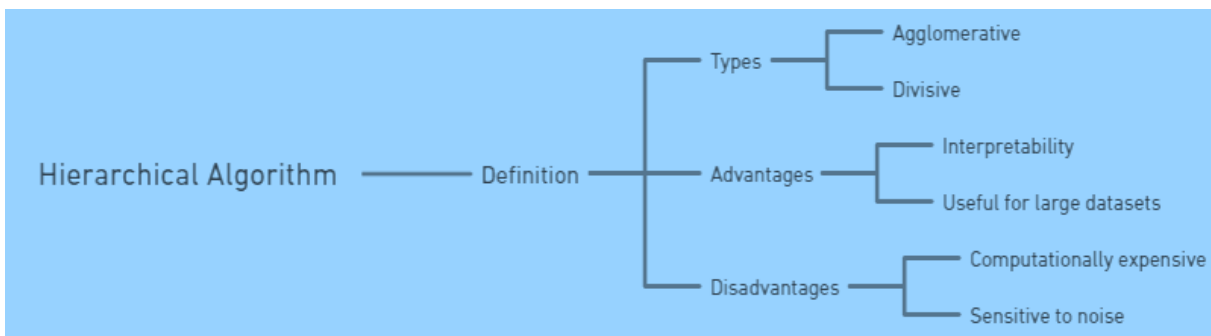


Figure 11 Approach of Hierarchical Algorithm

**2.11 Gaussian Mixture Algorithm**

It is probably the most familiar unaccompanied learning technique. Soft data aggregation is an approach for determining the possibility of varying sorts of data distribution. This algorithm's implementation is based on expectation maximization.[38] A GMM (Gaussian mixture model) is used for modelling data from one of the numerous groups; the groups may differ, but data points within the same group can be modelled by a Gaussian distribution. The image is a matrix with each element representing a pixel. The value is nothing more than a quantity that signifies the strength or appearance of the source images. [39]

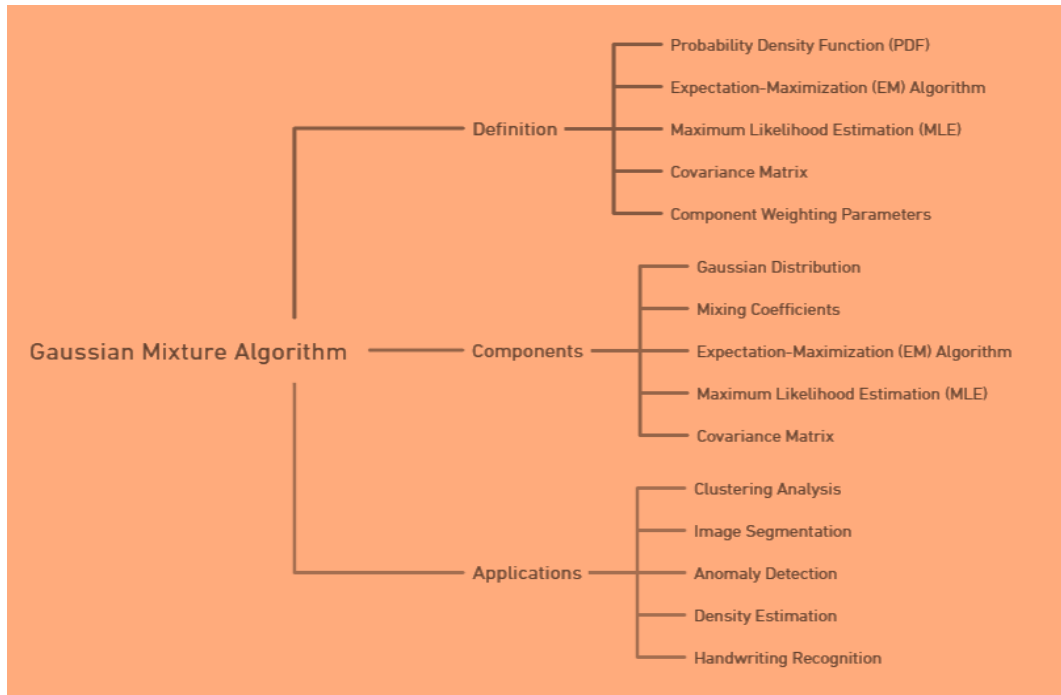


Figure 12 Characterization of Gaussian Mixture Algorithm

Algorithm	Year of Introduction	Complexity of Implementation	Year Introduced in Medical Industry
Artificial Neural Network (ANN)	1958	High: Requires significant computational resources, hyperparameter tuning, and large datasets for training.	1980s
Logistic Regression (LR)	1958	Low: Simple to implement with straightforward mathematical principles.	1980s
K-Nearest Neighbor (KNN)	1951	Medium: Requires tuning of the number of neighbors (k) and can be computationally intensive for large datasets.	1990s
Decision Tree (DT)	1986	Medium: Relatively easy to implement but requires pruning to prevent overfitting.	1990s
Naive Bayes (NB)	1960s	Low: Simple to implement and fast to train, especially for text classification tasks.	1990s
Support Vector Machine (SVM)	1992	High: Computationally intensive, especially with large datasets and requires careful selection of kernel functions.	2000s
Random Forest (RF)	1995	Medium: Requires tuning of the number of trees and other hyperparameters, but generally robust and scalable.	2000s
K-Means Algorithm	1957	Low: Simple and efficient to implement, but requires the number of clusters (k) to be specified beforehand.	1980s
C-Mean Algorithm	1973	Medium: More complex than K-Means due to the fuzziness parameter, but offers more flexible clustering.	1990s
Hierarchical Algorithm	1980s	High: Computationally intensive and less scalable for large datasets due to its hierarchical nature.	1990s
Gaussian Mixture Algorithm	1960s	High: Requires estimation of multiple parameters and is computationally expensive, especially for large datasets.	2000s

Table. 2. Comparison Table of Brest Cancer Screening based on Machine Learning algorithm Implementation Complexity



Algorithm	Type	Key Parameters	Advantages	Disadvantages	Applications in Breast Cancer
<b>Artificial Neural Network (ANN)</b>	Supervised Learning	Number of Layers, Number of Neurons, Learning Rate	Capable of identifying complex patterns; Effective in large datasets	Requires a large amount of data; Computationally intensive	CAD systems for mammogram analysis and biopsy decisions
<b>Logistic Regression (LR)</b>	Supervised Learning	Regularization Parameter, Learning Rate	Simple to implement; Provides probability estimates	Assumes linear relationship between features; Sensitive to outliers	Estimating probability of cancer based on mammographic features
<b>K-Nearest Neighbor (KNN)</b>	Supervised Learning	Number of Neighbors (k), Distance Metric	Simple and intuitive; Effective in low-dimensional data	Computationally expensive; Sensitive to noise in data	Classifying tumors based on similarity to known cases
<b>Decision Tree (DT)</b>	Supervised Learning	Depth of Tree, Split Criterion	Easy to interpret; Handles both numerical and categorical data	Prone to overfitting; Can create biased trees	Identifying key features in mammogram analysis
<b>Naive Bayes (NB)</b>	Supervised Learning	Smoothing Parameter (Laplace)	Works well with small datasets; Fast and scalable	Assumes feature independence; May underperform on complex data	High-dimensional data classification, such as gene expression profiles
<b>Support Vector Machine (SVM)</b>	Supervised Learning	Kernel Type, Regularization Parameter, Gamma	High accuracy; Effective in high-dimensional spaces	Not suitable for large datasets; Choice of kernel can be complex	Classifying mammographic findings into benign or malignant
<b>Random Forest (RF)</b>	Supervised Learning	Number of Trees, Max Depth	Robust to overfitting; Handles large datasets well	Can be slow to train; Less interpretable than single decision trees	Feature importance ranking, ensemble predictions for mammogram analysis
<b>K-Means Algorithm</b>	Unsupervised Learning	Number of Clusters (k), Initialization Method	Simple and scalable; Efficient on large datasets	Assumes clusters are spherical; Sensitive to initial conditions	Categorizing breast cancer tumors based on imaging characteristics
<b>C-Mean Algorithm</b>	Unsupervised Learning	Number of Clusters, Fuzziness Parameter	Allows data points to belong to multiple clusters; Flexible	Computationally expensive; Sensitive to initial conditions	Medical image segmentation, tumor categorization
<b>Hierarchical Algorithm</b>	Unsupervised Learning	Linkage Criteria, Distance Metric	Creates a comprehensive hierarchy of clusters; No need to specify number of clusters initially	Computationally intensive for large datasets; Less scalable	Analyzing hierarchical relationships in gene expression data
<b>Gaussian Mixture Algorithm</b>	Unsupervised Learning	Number of Components, Covariance Type	Models complex data distributions; Flexible clustering	Computationally expensive; Sensitive to initial parameter settings	Modeling distribution of tumor characteristics in medical images

Table.3. comparison of the key parameters, advantages, disadvantages, and applications of different machine learning algorithms used in breast cancer screening



Algorithm	Parameter	Description	Contribution
Convolutional Neural Networks (CNNs)	Layers	Multiple layers including convolutional, pooling, and fully connected layers.	Enhanced feature extraction from mammogram images, improving accuracy in detecting tumors.
	Filters	Filters applied in convolutional layers to detect edges, textures, and patterns in images.	Automatically learn and identify critical features for cancer detection.
	Activation Functions	Functions like ReLU, Sigmoid, and Softmax applied to introduce non-linearity.	Improve model's ability to capture complex patterns and make binary or multi-class classifications.
	Training Data	Large datasets of labeled mammogram images.	Better performance with larger, more diverse datasets, leading to improved generalization.
Support Vector Machines (SVMs)	Kernel Functions	Functions like linear, polynomial, and RBF kernels used to transform input data into higher-dimensional space.	Effective in handling non-linear relationships and improving classification accuracy.
Random Forests (RFs)	Hyperparameters	Parameters like C (regularization) and gamma (kernel coefficient).	Optimization of these parameters enhances model performance.
	Number of Trees	Total trees in the forest, commonly ranging from 100 to 1000.	Increasing number of trees typically improves prediction accuracy and robustness.
	Max Depth	Maximum depth of each tree in the forest.	Controls model complexity and overfitting.
Ensemble Methods	Voting Mechanism	Techniques like majority voting, weighted voting.	Combining multiple models to improve overall prediction accuracy and reliability.
Reinforcement Learning	Base Learners	Different types of algorithms used as base learners (e.g., decision trees, SVMs).	Diversity in base learners can lead to better generalization and robustness in predictions.
	Reward Function	Function that assigns rewards or penalties based on actions taken.	Guides the learning process towards optimal decision-making strategies.
Data Integration and Preprocessing	Exploration vs. Exploitation	Balance between exploring new actions and exploiting known profitable actions.	Crucial for improving learning efficiency and model performance over time.
	Data Augmentation	Techniques like rotation, flipping, and zooming applied to images.	Enhances model's ability to generalize from limited datasets.
	Normalization	Scaling data to a standard range (e.g., 0-1).	Improves convergence rate and model performance.
Clinical Workflow Integration	Interoperability	Ability to integrate ML models with existing clinical systems and workflows.	Ensures seamless adoption and practical utility of ML models in clinical settings.
	Interpretability	Techniques like saliency maps, LIME, and SHAP to explain model predictions.	Builds trust and transparency, aiding in clinical decision-making and adoption.
Ethical Considerations	Bias Mitigation	Strategies to address biases in training data and model predictions.	Ensures fairness and equity in screening outcomes across diverse patient populations.
	Privacy and Security	Measures to protect patient data privacy and ensure secure handling of sensitive information.	Maintains compliance with regulations and fosters patient trust.

Table. 4. Summarizing key parameters and contributions of various machine learning algorithms used in breast cancer screening, with a focus on convolutional neural networks (CNNs) for mammogram analysis



### III. INTEGRATION OF ML ALGORITHMS WITH CLINICAL WORKFLOWS

Integrating machine learning (ML) algorithms into clinical workflows for breast cancer screening holds promise for enhancing diagnostic accuracy and efficiency. However, this integration faces significant challenges. Data heterogeneity, stemming from diverse sources and formats, complicates model training and validation. Ensuring interpretability is critical for clinician trust and adoption, necessitating models that provide transparent and understandable predictions. Ethical considerations, including bias mitigation and patient privacy, must be addressed to ensure equitable and secure application. Successful integration demands robust data preprocessing, transparent algorithms, and adherence to ethical standards, fostering clinician collaboration and patient trust while leveraging ML's potential to improve screening outcomes.

#### a. Data Heterogeneity

Data heterogeneity poses a significant challenge in integrating ML algorithms into clinical workflows. Breast cancer screening data come from various sources and formats, leading to inconsistencies that complicate model training and validation. Effective integration requires robust data preprocessing and standardization techniques to ensure that the ML models can handle diverse and heterogeneous datasets.

#### b. Interpretability

Interpretability is crucial for the adoption of ML algorithms in clinical settings. Clinicians need to understand and trust the model's predictions to incorporate them into their decision-making processes. Transparent models, which provide clear and understandable predictions, along with tools like saliency maps and feature importance plots, are essential for building this trust.

#### c. Ethical Considerations

Ethical considerations are paramount when integrating ML algorithms into clinical workflows. Addressing bias in training data and ensuring fairness across diverse patient populations is critical. Additionally, protecting patient privacy and ensuring secure handling of sensitive information are necessary to maintain compliance with regulations and foster patient trust in these advanced technologies.

### IV. ML APPLICATIONS IN BREAST CANCER SCREENING

Machine learning (ML) algorithms have revolutionized breast cancer screening by enhancing the accuracy, efficiency, and reliability of diagnostic processes. Key applications include:

**1. Image Analysis:** Convolutional Neural Networks (CNNs) are extensively used for analyzing mammograms, ultrasounds, and MRI scans. They excel in detecting tumors, classifying them as benign or malignant, and identifying patterns that may be overlooked by human radiologists.

**2. Risk Prediction:** ML models, including logistic regression and decision trees, are utilized to predict the likelihood of developing breast cancer based on patient history, genetic factors, and lifestyle attributes.

**3. Automated Reporting:** Natural Language Processing (NLP) algorithms can generate automated radiology reports, summarizing findings and recommendations, thus saving time and reducing human error.

**4. Personalized Treatment Plans:** ML models analyze vast amounts of clinical data to recommend personalized treatment plans, predicting how patients will respond to different therapies.

### V. GAPS IN EXISTING RESEARCH

**1. Data Quality and Availability:** High-quality, annotated datasets are essential for training robust ML models. However, there is a lack of standardized, large-scale datasets in breast cancer screening, leading to potential biases and overfitting.

**2. Model Interpretability:** Many ML models, especially deep learning models, operate as "black boxes" with limited transparency in their decision-making processes. This hampers clinician trust and acceptance.





**3. Generalizability:** Models trained on specific datasets may not perform well across diverse populations or different clinical settings. This lack of generalizability limits the widespread adoption of these technologies.

**4. Integration with Clinical Workflows:** Seamless integration of ML models into existing clinical workflows remains a challenge. There is often a disconnect between model outputs and actionable insights that clinicians can use.

**5. Ethical and Legal Concerns:** Issues related to patient privacy, data security, and algorithmic bias need to be rigorously addressed to ensure ethical and legal compliance.

## VI. DIRECTIONS FOR FUTURE STUDIES

**1. Developing Standardized Datasets:** Establishing large, annotated, and standardized datasets that represent diverse populations and clinical settings will improve model robustness and generalizability.

**2. Improving Model Interpretability:** Research should focus on developing interpretable ML models and tools that provide insights into the decision-making processes, fostering clinician trust and transparency.

**3. Enhancing Generalizability:** Techniques like transfer learning, domain adaptation, and federated learning can be explored to improve the generalizability of models across different populations and clinical environments.

**4. Seamless Workflow Integration:** Designing ML models with a focus on integration into clinical workflows, including user-friendly interfaces and decision support tools, will facilitate their practical adoption.

**5. Addressing Ethical and Legal Concerns:** Establishing guidelines and frameworks for addressing ethical and legal issues, such as patient consent, data anonymization, and bias mitigation, is crucial for the responsible deployment of ML in healthcare.

## VII. CONCLUSION

At present, breast cancer is reportedly the most commonly diagnosed form of cancer. This article emphasizes various ways in which deep learning, machine learning and data mining are used in predicting breast cancer. The major benefit of this paper is how it shows that machine learning methods can be used to identify and facilitate early detection of breast cancer. Currently deep learning and machine learning techniques are widely employed when detecting malignancy. These algorithms have precision highly influenced by the dataset. As an organization, we seek innovative approaches and frameworks for instance deep learning or ML techniques that can apply to any dataset and give the best performance in terms of making projections on breast cancers.

Reported by Sara Al Ghunaim et al., SVM on Weka showed 98.03% reliability [41]. Even Flash tool when used had a higher truthfulness than other ML methods with a reliability of 99.68%. The dataset was obtained from Dataset, and its interpretation using methodologies such as CNN, SAE, SSAE resulted in a precision value of 98.9%. The tabular comparison highlights various algorithms used in AI for breast cancer applications, each with distinct strengths and weaknesses. Supervised learning algorithms such as Artificial Neural Networks (ANNs) and Support Vector Machines (SVMs) excel in handling large and high-dimensional datasets, making them suitable for tasks like mammogram analysis and classification of findings into benign or malignant categories. However, their computational intensity and need for extensive data are notable drawbacks.

Logistic Regression (LR) and Naive Bayes (NB) offer simplicity and speed, making them effective for probability estimation and classification tasks, though they may struggle with complex patterns and outliers. Decision Trees (DTs) and Random Forests (RFs) provide interpretable models and robustness to overfitting, respectively, but face challenges such as bias and training time.

Unsupervised learning algorithms, including K-Means, C-Means, and Hierarchical algorithms, are valuable for clustering and categorizing tumors based on imaging characteristics without prior labels. While these methods can be computationally intensive and sensitive to initial conditions, they offer flexibility and scalability for large datasets. Overall, the choice of algorithm depends on the specific requirements of the breast cancer application, considering factors like dataset size, complexity, and the need for interpretability. Understanding the trade-offs between different AI algorithms is crucial for optimizing their use in breast cancer detection, diagnosis, and treatment planning.



## REFERENCES

- [1]. Piluta, R., "Machine Learning in Healthcare: 12 Real-World Use Cases – NIX United," NIX United – Custom Software Development Company in US, Oct. 06, 2021. [<https://nix-united.com/blog/machine-learning-in-healthcare-12-real-world-use-cases-to-know/>](<https://nix-united.com/blog/machine-learning-in-healthcare-12-real-world-use-cases-to-know/>) (accessed Sep. 07, 2022).
- [2]. Sharma, S., Aggarwal, A., & Choudhury, T., "Breast Cancer Detection Using Machine Learning Algorithms," in 2018 International Conference on Computational Techniques, Electronics and Mechanical Systems (CTEMS), Dec. 2018, pp. 114–118. doi: 10.1109/CTEMS.2018.8769187.
- [3]. Sun, Y.-S., et al., "Risk Factors and Preventions of Breast Cancer," Int J Biol Sci, vol. 13, no. 11, pp. 1387–1397, 2017, doi: 10.7150/ijbs.21635
- [4]. Khourdifi, Y., & Bahaj, M., "Applying Best Machine Learning Algorithms for Breast Cancer Prediction and Classification," in 2018 International Conference on Electronics, Control, Optimization and Computer Science (ICECOCS), Dec. 2018, pp. 1–5. doi: 10.1109/ICECOCS.2018.8610632.
- [5]. Fatima, N., Liu, L., Hong, S., & Ahmed, H., "Prediction of Breast Cancer, Comparative Review of Machine Learning Techniques, and Their Analysis," IEEE Access, vol. 8, pp. 150360–150376, 2020, doi: 10.1109/ACCESS.2020.3016715
- [6]. Tuggener, L., et al., "Automated Machine Learning in Practice: State of the Art and Recent Results," in 2019 6th Swiss Conference on Data Science (SDS), Jun. 2019, pp. 31–36. doi: 10.1109/SDS.2019.00-11.
- [7]. Dhall, D., Kaur, R., & Juneja, M., "Machine Learning: A Review of the Algorithms and Its Applications," in Proceedings of ICRIC 2019, Cham, 2020, pp. 47–63. doi: 10.1007/978-3-030-29407-6\_5
- [8]. "Predicting factors for survival of breast cancer patients using machine learning techniques | Springer Link."(<https://link.springer.com/article/10.1186/s12911-019-0801-4>) (accessed Sep. 19, 2022).
- [9]. Tabra, Y., & Nidhal, F., "Reduced hardware requirements of deep neural network for breast cancer diagnosis," IAES International Journal of Artificial Intelligence (IJ-AI), vol. 11, pp. 1362–1372, Dec. 2022, doi: 10.11591/ijai.v11.i4.pp1362-1372.
- [10]. Dai, Q., Xu, S.-H., & Li, X., "Parallel Process Neural Networks and Its Application in the Predication of Sunspot Number Series," in 2009 Fifth International Conference on Natural Computation, Aug. 2009, vol. 1, pp. 237–241. doi: 10.1109/ICNC.2009.335.
- [11]. Tsai, W. K., Parlos, A., & Fernandez, B., "ASDM-a novel neural network model based on sparse distributed memory," in 1990 IJCNN International Joint Conference on Neural Networks, Jun. 1990, pp. 771–776 vol.1. doi: 10.1109/IJCNN.1990.137662.
- [12]. Hindawi, "Artificial Neural Networks in Mammography Interpretation and Diagnostic Decision Making."(<https://www.hindawi.com/journals/cmmm/2013/832509/>)(<https://www.hindawi.com/journals/cmmm/2013/832509/>) (accessed Sep. 20, 2022).
- [13]. Peng, C.-Y. J., Lee, K. L., & Ingersoll, G. M., "An Introduction to Logistic Regression Analysis and Reporting," The Journal of Educational Research, vol. 96, no. 1, pp. 3–14, Sep. 2002, doi: 10.1080/00220670209598786.
- [14]. Majumder, S., & Kehtarnavaz, N., "Vision and Inertial Sensing Fusion for Human Action Recognition: A Review," IEEE Sensors Journal, vol. 21, no. 3, pp. 2454–2467, Feb. 2021, doi: 10.1109/JSEN.2020.3022326.
- [15]. "Machine learning approaches for breast cancer diagnosis and prognosis | IEEE Conference Publication | IEEEXplore."(<https://ieeexplore.ieee.org/abstract/document/8280082>)(<https://ieeexplore.ieee.org/abstract/document/8280082>) (accessed Sep. 20, 2022).
- [16]. Baf, S., Im, E., & Bol, M., "Application of K-Nearest Neighbor (KNN) Approach for Predicting Economic Events: Theoretical Background."
- [17]. Sharma, H., & Kumar, S., "A Survey on Decision Tree Algorithms of Classification in Data Mining," International Journal of Science and Research (IJSR), vol. 5, Apr. 2016.
- [18]. Mahmood, A. M., Imran, M., Satuluri, N., Kuppa, M. R., & Rajesh, V., "An Improved CART Decision Tree for Datasets with Irrelevant Feature," in Swarm, Evolutionary, and Memetic Computing, Berlin, Heidelberg, 2011, pp. 539–549. doi: 10.1007/978-3-642-27172-4\_64.
- [19]. "Performance of Decision Tree C4.5 Algorithm in Student Academic Evaluation | SpringerLink." ([https://link.springer.com/chapter/10.1007/978-981-10-8276-4\\_36](https://link.springer.com/chapter/10.1007/978-981-10-8276-4_36)) (accessed Sep. 19, 2022).
- [20]. Pandya, R., & Pandya, J., "C5. 0 Algorithm to Improved Decision Tree with Feature Selection and Reduced Error Pruning," IJCA, vol. 117, no. 16, pp. 18–21, May 2015, doi: 10.5120/20639-3318.
- [21]. Song, Y., & Lu, Y., "Decision tree methods: applications for classification and prediction," Shanghai Arch Psychiatry, vol. 27, no. 2, pp. 130–135, Apr. 2015, doi: 10.11919/j.issn.1002-0829.215044.
- [22]. "Decision tree classifier for network intrusion detection with GA-based feature selection | Proceedings of the 43rd annual Southeast regional conference - Volume 2."





[https://dl.acm.org/doi/abs/10.1145/1167253.1167288](https://dl.acm.org/doi/abs/10.1145/1167253.1167288) (accessed Sep. 20, 2022).

- [23]. “Decision Tree Optimization in C4.5 Algorithm Using Genetic Algorithm - IOPscience.” [https://iopscience.iop.org/article/10.1088/17426596/1255/1/012012/meta](https://iopscience.iop.org/article/10.1088/1742-6596/1255/1/012012/meta) (accessed Sep. 20, 2022).
- [24]. Wu, W., Nagarajan, S., & Chen, Z., “Bayesian Machine Learning: EEG/MEG signal processing measurements,” *IEEE Signal Processing Magazine*, vol. 33, no. 1, pp. 14–36, Jan. 2016, doi: 10.1109/MSP.2015.2481559.
- [25]. Ibrahim, A. A., Hashad, A. I., & Shawky, N. E. M., “A Comparison of Open-Source Data Mining Tools for Breast Cancer Classification,” *Handbook of Research on Machine Learning Innovations and Trends*, 2017. [https://www.igi-global.com/chapter/a-comparison-of-open-source-data-mining-tools-for-breast-cancer-classification/www.igi-global.com/chapter/a-comparison-of-open-source-data-mining-tools-for-breast-cancer-classification/180964](https://www.igi-global.com/chapter/a-comparison-of-open-source-data-mining-tools-for-breast-cancer-classification/www.igi-global.com/chapter/a-comparison-of-open-source-data-mining-tools-for-breast-cancer-classification/180964) (accessed Sep. 19, 2022).
- [26]. “Support Vector Machines: Theory and Applications | SpringerLink.” (https://link.springer.com/chapter/10.1007/3-540-44673-7\_12) (accessed Sep. 19, 2022).
- [27]. Yang, Y., Li, J., & Yang, Y., “The research of the fast SVM classifier method,” in *2015 12th International Computer Conference on Wavelet Active Media Technology and Information Processing (ICCWAMTIP)*, Dec. 2015, pp. 121–124. doi: 10.1109/ICCWAMTIP.2015.7493959.
- [28]. Islam, Md. M., Iqbal, H., Haque, Md. R., & Hasan, Md. K., “Prediction of breast cancer using support vector machine and K-Nearest neighbors,” in *2017 IEEE Region 10 Humanitarian Technology Conference (R10-HTC)*, Dec. 2017, pp. 226–229. doi: 10.1109/R10-HTC.2017.8288944.
- [29]. Zhang, Y., *New Advances in Machine Learning. BoD – Books on Demand*, 2010.
- [30]. Li, Y., & Wu, H., “A Clustering Method Based on K-Means Algorithm,” *Physics Procedia*, vol. 25, pp. 1104–1109, Jan. 2012, doi: 10.1016/j.phpro.2012.03.206.
- [31]. Patel, B., & Sinha, P. G., “An Adaptive K-means Clustering Algorithm for Breast Image Segmentation,” *International Journal of Computer Applications*, vol. 10, Nov. 2010, doi: 10.5120/1467-1982.
- [32]. Li, Y., & Wu, H., “A Clustering Method Based on K-Means Algorithm,” *Physics Procedia*, vol. 25, pp. 1104–1109, Jan. 2012, doi: 10.1016/j.phpro.2012.03.206.
- [33]. Thangavel, K., & Mohideen, A. K., “Semi-supervised k-means clustering for outlier detection in mammogram classification,” in *Trendz in Information Sciences & Computing(TISC2010)*, Dec. 2010, pp. 68–72. doi: 10.1109/TISC.2010.5714611.
- [34]. Bijral, R. K., Manhas, J., & Sharma, V., “Hierarchical Clustering Based Characterization of Protein Database Using Molecular Dynamic Simulation,” in *Recent Innovations in Computing*, Singapore, 2022, pp. 427–437. doi: 10.1007/978-981-16-8248-3\_35.
- [35]. Chen, Y., Wang, Y., & Yang, B., “Evolving Hierarchical RBF Neural Networks for Breast Cancer Detection,” in *Neural Information Processing*, Berlin, Heidelberg, 2006, pp. 137–144. doi: 10.1007/11893295\_16.
- [36]. Zhang, J., Hong, X., Guan, S.-U., Zhao, X., Xin, H., & Xue, N., “Maximum Gaussian Mixture Model for Classification,” in *2016 8th International Conference on Information Technology in Medicine and Education (ITME)*, Dec. 2016, pp. 587–591. doi: 10.1109/ITME.2016.0139.
- [37]. Singh, A. S., “Segmentation of Breast Images Using Gaussian Mixture Models,” p. 5, 2017.
- [38]. Kajala, A., & Jain, V. K., “Diagnosis of Breast Cancer using Machine Learning Algorithms-A Review,” in *2020 International Conference on Emerging Trends in Communication, Control and Computing (ICONC3)*, Feb. 2020, pp. 1–5. doi: 10.1109/ICONC345789.2020.9117320.
- [39]. Alghunaim, S., & Al-Baity, H. H., “On the Scalability of Machine-Learning Algorithms for Breast Cancer Prediction in Big Data Context,” *IEEE Access*, vol. 7, pp. 91535–91546, 2019, doi: 10.1109/ACCESS.2019.2927080.
- [40]. “Deep Learning Techniques for Breast Cancer Detection Using Medical Image Analysis | SpringerLink.” [https://link.springer.com/chapter/10.1007/978-3-319-61316-1\_8](https://link.springer.com/chapter/10.1007/978-3-319-61316-1\_8) (accessed Sep. 20, 2022).
- [41]. Kargbo, R. B., “PROTAC-Mediated Degradation of Estrogen Receptor in the Treatment of Cancer,” *ACS Med. Chem. Lett.*, vol. 10, no. 10, pp. 1367–1369, Oct. 2019, doi: 10.1021/acsmchemlett.9b00397.