# Advanced Machine Learning Techniques for Early Detection and Classification of Breast Cancer

Sumit Kushwaha

Department of Computer Applications, University Institute of Computing, Chandigarh University, Mohali, India

## Abstract

Breast cancer is one of the most dangerous diseases, and it is the second leading cause of morbidity in women. Breast cancer arises when hazardous, malignant tumours grow in the mammary gland. Self-tests and routine clinical checks help in early detection and hence substantially improve survivorship. Breast cancer classification is a medical method that academics and specialists find difficult to implement. Several microarray studies have utilised gene signatures to create classifications that predict medical outcomes for various cancer patients. Signatures from diverse studies usually suffer from low consistency when used in the classification of databases, regardless of the study from which they were produced. By integrating the auto-encoder and Principal Component Analysis, the researchers provide an unsupervised feature training strategy for characterizing different qualities from variations in gene expression. An ensemble classifier based on the AdaBoost algorithm was created as the framework for the gathered attributes to anticipate medical outcomes in breast cancer. During the experiments, the researchers created an additional classifier using the same classifier learning strategy to act as a median for the suggested technique. Experiments reveal that the proposed system, which makes use of deep learning techniques, outperforms others.

## Keywords

Breast Cancer Detection, Machine Learning in Healthcare, Early Diagnosis and Treatment, SDG 3, SDG 9, AI, Deep Learning

## 1. Introduction

Breast cancer is the recurrent malignancy in women and the second foremost reason of mortality. Since the origins of breast cancer remain unknown, early detection is essential for effective treatments and a lower mortality rate. There are countless cells in the human body [1]. Cancer begins when biological alterations drive these cells to develop abnormally, resulting in a mass known as a malignant tumor. As per Cancer Statistics, roughly 2.25 million people worldwide are affected by cancer in 2018, as in figure 1 [2]. Every year, over 1,157,294 new people with cancer are recorded, and about 784, 821 cancer-related fatalities are documented. Males have a 7.34 percent likelihood of dying from illness, while females have a 6.28 percent chance. Oral mucosa and lung diseases contribute to 25% of cancer fatalities in men, while oral cavity and breast cancers are responsible for 25% of cancer deaths in women. Figure 1 depicts a comprehensive overview of cancer data for the year 2018. As it could be seen, breast cancer is the maximum commonly noted carcinoma in women, accounting for 14% of all malignancies. According to the Globocan 2018 survey, there were 162,468 new cases of breast cancer and 87,090 fatalities [3], [4].
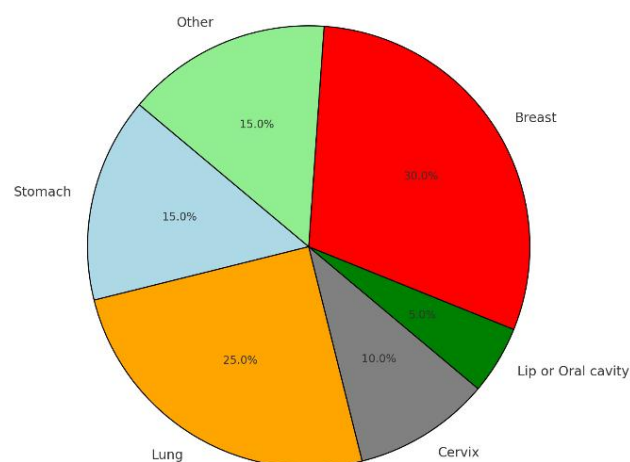


**Figure 1.** Cancer Statistics

Ultrasonography systems have long been utilized for breast illness diagnosis and screening because of their cost-effectiveness, accessibility, noninvasiveness, and real-time imaging. Ultrasonography has been demonstrated to accurately determine benign and malignant tumors in investigations. Elastography is a method that evaluates the stiffness of tissue using ultrasonography as the primary diagnostic tool. Elastography is beneficial in discriminating between benign and malignant breast cancers [5], [6] because there is a considerable difference in flexibility between benign and malignant tumors. It has introduced a new characteristic to assess the stiffness of tissue and has demonstrated that the elasticity of tissue is of significant utility for diagnostic purposes. Breast cancers in their initial phases produce minimal indications since the tumor remains tiny and so curable. As the result, early diagnosis and treatment are challenging but critical. Early diagnosis and treatment with mammography considerably enhance the likelihood of living, according to both randomized trials and inhabitants' analyses. Ultrasound imaging is by far the most efficient cancer screening technology commercially available since it could detect cancer numerous generations earlier than visible signs appear [7]. Nevertheless, approximately 5%–10% of mammogram readings are viewed as abnormal or inconclusive unless additional examinations such as ultrasound scans or breast biopsies reveal a final diagnosis of normal or benign breast tissue [8].

Breast cancer is a disease in which mammary cells are grown unnaturally, and certain malignant tumors have grown from these cells. Cancer stem cells can differentiate and infect healthy breast tissue or even other tissues by lymph nodes, leading to mature phase or metastatic breast cancers. Lymphoma nodes are tiny nodes that sieve lymphatic fluid and are among the primary sites where breast cancer had spread [9]. Breast cancer could spread to nearby liver, lymph, bone, lungs, and brain if it is not diagnosed and treated promptly. Whenever it spreads to the bones and brain, the patient is suffering from excruciating agony daily [10]. Artificially intelligent (AI) and machine learning (ML) systems have recently advanced and developed approaches to assist pathology and physicians in more precisely diagnosing, assessing, and treating patients. Several investigations have used and documented the different applications of derived modeling in medical research for earlier modeling of carcinogenic effects or speedy and precise assessment of clinical outcomes. Forecasting related to modern image processing methods and AI technologies were built for the earlier detection of malignancy results in cancer screening [11], [12].

As a result, a biopsy process is frequently used, in which a piece of cells from the suspected area is taken and sent to a laboratory to determine. A physician examines the specimen below a magnification, and a final report verifies or denies the existence of malignant cells in the blood. Although this strategy is perhaps the most definite way to make a diagnosis, it has a partiality problem where many diagnoses could be given on the same material, especially among non-specialized pathologists [13]. There was also the problem of histopathological inadequate staffing, which might keep tissue samples on wait for up to two months. This happens frequently in Norway. Because of the explosive rise of artificial intelligence, the computer-aided diagnosis (CAD) method of diagnostic imaging is quickly becoming one of the most reliable approaches to detecting probable diseases at an early stage and identifying them [14]. The CAD methods used by BC make use of highly advanced computer technologies to automatically determine whether breast findings are benign or malignant. The potential of computer-aided diagnosis (CAD) to enhance classification of malignant and benign lesions has been proven in preclinical investigations. A growing number of computer-aided diagnosis (CAD) systems are now being utilized for specific breast imaging modalities. Entire photos could now be digitized to build and generate massive media content thanks to the introduction of full slide digital scanners. This allowed academics and clinicians to perform the mathematical methodology of histopathology pictures by building accurate CAD methods to assist in the identification and treatment of various diseases.

Computer-aided detection and diagnosis (CADD) technologies are being utilized for more than two decades and are frequently mentioned to as "second thoughts." CAD devices work by training highly technical skill solutions utilizing machine learning and pattern detection procedure using radiographic pictures with conventional diagnostic features. These procedures could then detect the imaging features they were skilled with on test photos (i.e., images that have not been seen or utilized in training), enabling computers to assist in disease identification and treatment. This technology has the potential to be tremendously operative in oncology, secondary in the better detection and treatment of a wide range of tumor systems.The utilization of CAD models for breast cancer identification and treatment was being studied frequently in the literature. That research employed a variety of imaging techniques and machine learning techniques, with some undergoing a diagnosis process for acceptability testing. Given the significant phenotypic variability in tumors, a huge rate of false positives, and low diagnostic rates, the efficacy of these investigations have indeed been restricted. As a result, many studies have been devoted to enhancing these mechanisms. Recent developments in machine learning, particularly "deep learning," have shifted study in this discipline in a more positive direction. The primary goal of several studies was to employ machine learning techniques to investigate the prediction performance of effectiveness of treatment after 6 months of pharmaceutical utilize in patients with pediatric who poses asthma and to uncover the crucial factors for insights into the process behind such reactions. Using machine learning algorithms like those used in this work to detect non-responders vs. responders to medication is critical in asthma care, as forecasting response to therapy is difficult (if not unattainable) in an existing process. The research used an observation pediatrics asthma cohort that was culturally diverse, age diversified, and matched real-life clinical scenarios, with the majority of the kids having a moderately severe illness.The primary outcome measures were four independent indices of therapy reaction (after six months), as measured by improvements in pulmonary function and sense of control. Together with the other measures and diagnostics, all of these objectives were examined at the beginning and after 6 months of

therapy administration. Machine learning techniques, such as Random Forest and AdaBoost classification algorithms, were used to test the accurate predictions, with patient outcomes (as categorical data) specified as results to also be forecasted based on other clinically relevant information and evaluations. Recognizing Machine learning is critical in healthcare because such techniques are regarded as black-box approaches.

In radionics investigation, ensemble teaching techniques that aggregate the forecasts of numerous classifications to achieve greater efficiency than a unique estimation have piqued attention. In contrast to the traditional Machine learning techniques, such methodologies have turned out to be exceptionally suitable for modeling heterogeneous collections of any length and diversity, while also excelling at trading off approximations and prediction mistakes. In terms of detecting chemical subtypes, sentinel lymph node metastatic prognosis and early prognostication of therapeutic responsiveness during breast mpMRI mammography, Boosting Ensemble Classification approaches have been found to outperform standard classification methods and permit the passage. They employed four minimum probabilities, a maximum probability, the product of probabilities, and the average of probabilities as unweighted voting mechanisms for ensemble classification, in addition to majority-based voting. Furthermore, its prediction efficacy for distinguishing normal from abnormal breast tumors has shown potential in a recent survey using DCE MRI radiomics independently, although that's yet to be tested with the mpMRI dataset. A person's propensity can be disrupted by a variety of causes; among them, breast cancer is the foremost cause of death in women. As a result, the main consideration is to reduce the fatality rate through early disease detection. Breast cancer is sometimes defined as the uncontrolled proliferation of abnormal cells in the breast's milk-secreting ducts. Early diagnosis and treatment is the most efficient way to reduce breast cancer mortality. Additionally, identifying breast cancer requires a precise and practical detection equation that enables clinicians to differentiate between benign and malignant breast tumors without needing to perform an immediate surgical biopsy. The objective of these detection systems is to categories people into "benign" or "malignant" groups. A benign tumor rarely results in human death and is not harmful to the body. This particular sort of tumor only develops in one area of the body. A cancerous tumor is more harmful and can be fatal to people. The aberrant cell development that causes this form of tumor causes it to spread quickly. A fine spine subjective dissection is a simple and quick approach to detecting development in cancer patients. This procedure entails extracting a sample of cells or fluid from just a nodule, followed by a microscopic examination of the tissues. If the knob cannot be detected, visualization approaches are required to determine the precise location. The doctor uses an ultrasound to evaluate the spine and guides it around the spot in the ultrasonography procedure.

Deep learning algorithms for clinical diagnosis relying on pathological pictures were developed by many researchers. For instance, utilizing a database of histology images, the scientists constructed a system capable of categorizing malignant or benign tumors. During the training step of the program, the greatest efficiency was 0.99 percent. The study gives an approach the with intention of reducing the mortality rate from breast cancer and preserving women's lives in a previous experiment. Eight breast cancer subtypes are categorized into two parts: benign and malignant, for each subtype containing four sub-classes to help lower the incidence. The writers of this paper used a magnification of 40x to cover the area of concern, then magnified the images to concentrate just on that area. A convolutional neural network contributed to texture feature extraction data and implementationof image processing filtration, a max-pooling did computations to avoid over-fitting, and ultimately, fully-connected levels encapsulated the result. The information was split into training and testing sets sections, with 90 percent and 10% proportions, accordingly.

## 2. Related Works

Breast cancer kills more women in the United States than any other malignancy. Machine learning-based forecasting analytics offer faster breast clinical diagnostic procedures. Nevertheless, evaluating systems that could accurately identify cancer remains difficult. To enhance breast cancer detection performance, researchers introduced data exploratory techniques (DET) and constructed four alternative forecasting analytics. Before developing a design, the researchers conducted in-depth investigations into four-layered critical DET. These investigations included feature allocation, correlations, removal, and hyperparameters adjustment. The goal of these investigations was to find the most reliable characteristic categorization into malignant and benign classifications. We were able to improve the performance of the prediction model using these mining techniques, which resulted in a maximum F1 score and an accuracy score that was higher than it had been before. The Breast Cancer Coimbra Dataset (BCCD) and Wisconsin Diagnostic Breast Cancer (WDBC) datasets have been used to assess the hypothesized methodologies and classifications. To evaluate each classifier's effectiveness and preparation time, conventional performance indicators such as segmented images and K-fold cross-validation approaches have been used. Using approach DET, the systems' analytical abilities were enhanced, with polynomial Support Vector Machine achieving 99.3 percent, LR 98.06 percent, KNN 97.35 percent, and EC 97.61 percent accuracy with the WDBC dataset. In terms of effectiveness, researchers further contrasted the present research's important findings with previous research. The outcomes of the deployment approach to help clinicians establish an efficient mechanism for identifying and diagnosis of breast cancerous tissue. Except for SVM, the BCCD database did not offer successful outcomes with experimental projections; hence, these findings were omitted in this investigation. In the "Data Explanation"paragraph, researchers supplied the databases' findings. Because these databases are for American patients, the outcomes may not have been comparable or beneficial for Asian patients. It is one of the report's drawbacks, which might be addressed in the future by using a new dataset and implementing neural networks. Thus, the research based on the enhanced machine learning-based forecasting system for the cancer diagnosis failed in providing the high efficiency or performance which is proposed..

Artificial intelligence has permitted the automatic diagnosis of the disease on image data, and mammograms take a significant part in detecting breast cancer in women. The focus of this research was to create a deep learning approach for detecting breast cancer in mammographic images of varied densities, as well as to contrast the effectiveness of the algorithm with earlier research [15]. Mediolateral and Craniocaudal view mammography were incorporated and synthesized for each breast from 1502 participants who had digital mammography between the year 2007 to 2015, resulting in 3002 combined pictures [16]. On the integrated pictures, two convolutional neural networks were programmed to find any malignant tumor. The results are then associated to a meta-analysis that included 12 prior deep learning experiments using 301 combined photos from 284 people. DenseNet-169 had an AUC of $0.952 \pm 0.006$ for breast cancer detection in each integrated mammography and EfficientNet-B5 had an AUC of $0.953 \pm 0.020$ for breast cancer detection in each integrated mammography. As breast density rose, the efficiency for malignancy identification declined (concentration A, average Area Under Curve = 0.985 vs. saturation D, medium Area Under Curve = 0.902 using DenseNet-169). So, when the age of the individuals was employed as a covariate for malignancy identification, the results were similar (average AUC, $0.953 \pm 0.005$). The DenseNet-169's mean specificity and sensitivity (87 and 88 percent, individually) outperformed the meta-mean analysis's values (81 and 82 percent, combined). Deep learning will be effective in forecasting breast cancer in digital mammography of varying concentrations, with the best results coming from breasts with lower parenchyma thickness. The current research has some limitations. For machine learning research, the percentage of patients selected was limited. Our mammography, on the other hand, was rigidly categorized based on pathological verification and a 5-year follow-up duration to reduce the chance of intermediate malignancy. Secondly, the source photos consist of a single tertiary institution of higher learning, to whom secondary institutions sent more serious illnesses. Finally, mammography was created with a solitary piece of equipment. Once the methodology could be widely used, more research is required to confirm this across universities and vendors. Furthermore, digital mammography identification of breast cancer utilizing deep learning has the standard limitation in the research [17]. Breast cancer is commonly in women's breasts, according to a study conducted in 2016, which found that 2.8 million women globally were identified with breast cancer that same year. The therapeutic treatment of a patient with breast cancer is expensive, and given the expense and charge of preserving a citizen's health, breast cancer prevention becomes a significant public health concern. Several technologies have been developed for this aim over the last 20 years, including mammography, which is commonly utilized for breast cancer detection. Nevertheless, mammogram false - positive could arise when the individual is diagnosed positively by some other method. Furthermore, the possible disadvantage of the mammogram could lead patients and doctors to seek out alternative testing procedures. The research study began with infrared digital photography, which implies that a simple thermal assessment among a normal breast and a cancer breast usually indicates an increase in thermal activities in pre-cancerous regions and the regions around growing breast cancer. Moreover, researchers discovered via this research that a Computer-Aided Diagnostic (CAD) using IR image analysis could not be accomplished without the need for a model like the well-known hemisphere paradigm. The creation of a relative analysis of many breast cancer recognition strategies employing sophisticated computer vision technology and deep learning techniques is the paper's critical contributor [18].

Mammography is a first-line imagination evaluation method for detecting initial breast cancer. Deep-learning-based computational methods, including the convolutional neural network (CNN), are frequently utilized as classifications in mammography examinations for speedy automated breast tumor detection. A multilayered Convolutional neural network has numerous convolutional-pooling layers and fully linked systems for categorizing numerous feature maps on two-dimensional (2D) image data, which could also improve transmission correctness and minimize error margin. This multilayer structure, though, has numerous drawbacks, including computationally expensive, the need for a large-scale training database, and limited fit for real-time clinical uses. As a result, this research proposes a multilayer design made up of a flattening layer, three convolution layers, two pooling layers, and a classification algorithm for CNN-based classifiers for autonomous breast tumor detection [19]. The suggested approach uses a fractional-order convolution layer procedure in the first convolution layers to improve the contrast and eliminate undesirable interference to achieve the optimum object's limits; in the second and the two kernel convolutional, third convolutional-pooling layers, and pooling processes are being used to ensure the continued improvement and deepening of the feature patterns for further removal of the desirable characteristics at various scales and levels. Furthermore, the feature patterns' dimensions have been reduced. During the classification step, a multilayer network equipped with an adaptive moment estimate technique is applied to optimize the network configuration of classifiers for mammography categorization. This is accomplished by distinguishing tumor-free feature patterns from tumour feature sequences. An optimization approach that is able to handle sparse gradients and noisy situations is provided by the Adam algorithm, which is a combination of the best qualities of the AdaGrad and RMSProp algorithms. By taking into account an 'exponentially weighted average,' this approach is utilized to speed up the gradient descent algorithm. This is accomplished by taking into account the gradients. K-fold cross-validations are done on images picked from a controlled breast imaging subgroup of a digital folder for mammography transmission. In terms of recall (percent), F1 score, accuracy (percent), and Youden's index, precision (percent), the experimental findings show promising performance for automatic breast tumor screening [20].

The efficiency of a neural net in forecasting 5-15 year breast-cancer-specific mortality was investigated in the present study. A total of 951 breast cancer survivors were separated into two groups: 651 for training and 300 for verification. Tumor size, mitotic count, histological type, axillary nodal status, tubule formation, nuclear pleomorphism, tumor

necrosis, and age have all been supplied as inputs to the network. The area under the ROC curve (AUC) has been utilized to measure the prediction systems' efficiency in producing longevity predictions for individuals in the independent testing set. The AUC values for 5-15 year breast-cancer-specific lifespan in the neural network models were 0.90, 0.88, and 0.88, correspondingly. The area under curve values for regression models were 0.85, 0.86, and 0.897, respectively. With a precision of 72% and a sensitivity of 78%, the axillary lymph node condition predicted 5-year longevity. At this specificity level, the neural network model's sensitivity was 92%. At 5 years, the frequency of erroneous forecasts for nodal status was 82/300 and the neural network was 40/300. The frequency of erroneous guesses slightly increased to 49/300 once the nodal condition was removed from the neural network model (AUC 0.877). In the prognosis of 5 to 15-year breast cancer-specific survival, an artificial neural network is extremely accurate. The management of patients who pass of an intercurrent reason or have inadequate follow-up continues to be a concern in logistic regression using neural network models. Several approaches have been presented that could utilize the data from these individuals, but none has yet been acknowledged as the preferred way. The management of victims who pass of an intercurrent reason or have inadequate follow-up continues to be a concern in logistic regression using neural network models. Several approaches have been presented that could utilize the data from these individuals, but none has yet been acknowledged as the preferred way [32]. Thus this work collects breast cancer data points, and the suggested technique's MCC (Matthew's correlation coefficient) index, ACC (accuracy), AUC (area under the receiver operating characteristic curve), as well as other assessment variables have been evaluated and associated with prominent gene signature-based methodologies, such as the baseline method.

## 3. Materials and Method

To accomplish dimension reduction for high-dimensional data sets, a variety of selection of features and feature extraction approaches existed and are commonly employed. Inherently, all of these algorithms explore removing irrelevant and redundant data to improve the categorization of unique test instances. The Principal Component Analysis (PCA) method and the autoencoder neural network utilize only for dimension removal (also known as feature acquisition) are extensively explained in this research.

## 3.1 Principal Component Analysis

One of the most general technologies for linear dimensional decrease in multivariate technology is Principal Component Analysis. PCA determines the principal components of information that are statistically independent eigen values, each reflecting some amount of variance in the dependent variable, using the linear combination and its eigen values and eigenvectors.
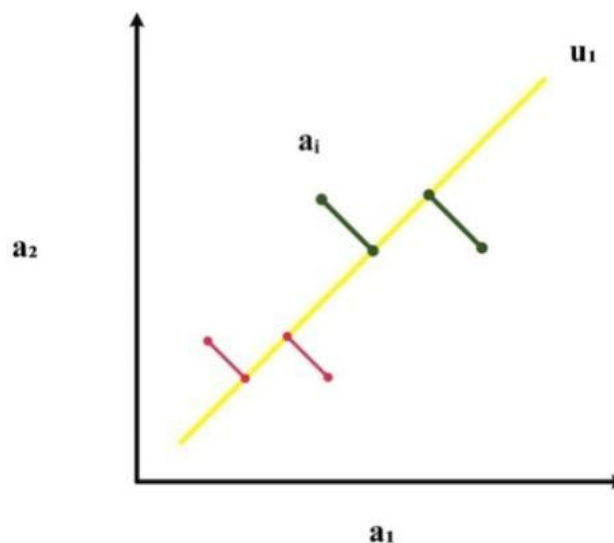


**Figure 2.** Orthogonal Projection

Let $A=\{a_i\}_{i=1}^m$ represent the set of the training dataset, $a_i$ denotes the set of variables with dimension D which holds the gene expression profile in this research.PCA has two goals: (a) extracting one of the most essential data from $a_m$; and (b) compressing the dimensionality of A by maintaining only the sensitive details. This is defined as an orthonormal basis of the initial D-dimensional source onto a new k-dimensional domain (k <D), with the variance of the planned outcomes as the goal to be reduced, as shown in Figure 2.

This part provides a quick overview of the autoencoder neural network, a nonlinear dimensionality reduction technique that is well-known for the extraction of features. An auto-encoder, as shown in Figure 3, is a feed-forward neural net that is frequently qualified to acquire interpretations or efficient encoding of the original input $A=\{a_i\}_{i=1}^m$. In this approach, it develops a functional $b(f(a_i)) \approx c_i$ that approximates the model parameters, which are made up of a small number of image activation functions and symbolized by the network's hidden neurons. Over 13 transcriptomic databases, researchers examined stacking encoder, Principal component analysis, and Gauss SVM. Their findings

revealed that the autoencoder outperforms the bulk of the data sources, which further inspires us in current research. The architecture consists of an autoencoder is classified in the following sections, according to the research: (1) the input units$a_i$ (2) a "code" or hidden recognition$x_i=f(a_i)$; 34) a decoder function g;; (4) an encoder component f; (5) the output systems, also known as "rebuilding" $b(f(a_i))$; and (6) a loss function $L(a_i, c_i)$ computing a scalar$\|a_i-c_i\|_2$which quantifies how useful the rebuilding $c_i$. The autoencoder's optimization goal is to reduce the predicted standards of L over the training instances A.
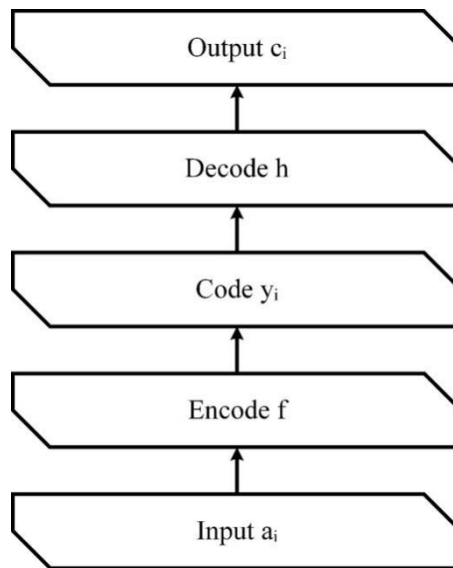


**Figure 3.** Autoencoder General Framework

The differential among outputs and inputs is assessed using mean squared error throughout each iteration process, and weight adjustments to multiple levels are made using backpropagation across the neural net. The encoder method achieves non-linear dimensionality minimization if the no. of hidden neurons is more than one. When the dimensions of the hidden units will be fewer than the dimension of input, the auto encoder will be able to discover the optimum feature reduction of the input on the hidden units since it will meet the qualifications necessary to do so. In the event that this is not the case, the auto encoder is trained to map the characteristic to a region with a greater dimension. It provides a helpful method for significantly reducing the noise of input data, which in turn makes the process of developing deep learning models significantly more effective. Anomalies can be found with their help, as can difficulties with unsupervised learning, and complexity can be removed with their assistance. The proposed technique in this research is presented in the following subsections. First, research goes over the five gene expression databases researchers utilized and how they processed them. Then researchers propose their solution to the issue.

The gene expression data were obtained from the NCBI GEO catalogue, which is accessible to the public. Each product includes 129,159 gene expression patterns from the Affymetrix microarray technology, including 22,268 probing for 978 benchmark genes and 21,290 gene expressions in each analysis. To assess the feasibility and robustness of the proposed technique, researchers used five distinct breast cancer databases from LINCS Cloud1. Except for GSE4922, which has been standardized using the method RMA, all five databases were regularized by their original authors utilizing the procedure MAS5.0. The statistics are detailed in Table 1 Database of Breast Cancer Patients

**Table 1.** Patients Dataset of Breast Cancer

| Info | Deprived Consequence | Virtuous Consequence | Overall | Samples Eliminated |
|------|----------------------|----------------------|---------|--------------------|
| GSE2034 | 94 | 180 | 274 | 12 |
| GSE4923 | 29 | 106 | 135 | 154 |
| GSE6533 | 20 | 78 | 98 | 225 |
| GSE6289 | 40 | 155 | 195 | 10 |
| GSE11122 | 27 | 155 | 182 | 16 |

Patients typically perform immunological and pathologic indicators in addition to affecting the resulting diagnosis in the multiple datasets. GSE6534, for instance, has exclusively ER-positive cases, whereas other databases have both ER-negative and ER-positive cases. The GSE4922 and GSE6534 datasets comprise two lymph node-negative and lymph node-positive individuals, whereas the additional data source solely cover lymphatic node-negative cases. With multiple datasets, researchers did two-step pretreatment to prepare them for the classification problem: To start, the researchers

utilized a dataset segmentation method in which all cancer patients were split into terrible prediction and great prediction when many distant metastases had occurred within the preceding five years. This was done so that the researchers could compare the two groups. People who had their treatment prematurely discontinued within the past five years or who had received adjuvant therapy were not included in the analysis. Second, because the gene expression levels of the microarray platform were determined by a number of distinct methods, the researcher's quantile normalized the various sources of data by employing the MAS 5.0 methodology. After that, every probe was linked to an Entrez Gene ID, and the results were averaged. With the help of our implementation of Affymetrix's MAS 5.0 expression measure, an instance of AffyBatch can be converted into an instance of Expression Set using the MAS 5.0 function. Entrez Gene offers one-of-a-kind integer identities for genes and other loci for a selection of the organisms that serve as models. Researchers aim to integrate variable assortment and feature removal techniques and deep learning approaches in this research in terms of learning additional significant properties from gene expression profiles and building a more effective classification for malignancy prognostic forecasting. The workflow of the suggested approach is represented in Figure 4.
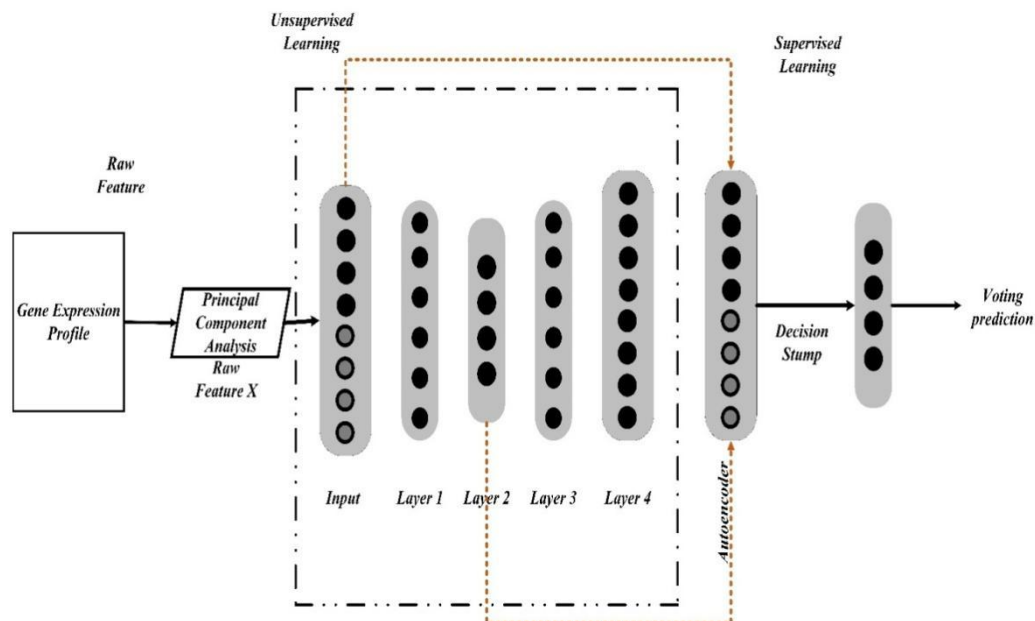


**Figure 4.** Flowchart based on Proposed Method

### 3.2 Unsupervised Learning

The feature learning technique is divided into two phases, as represented by the feature strategy instruction:

### 3.2.1 Principal Component Analysis

Researchers use Principal component analysis (as mentioned in the Section above) as the selection of feature approach to minimize the dimensions of the gene expression profiles because their dimension is quite high and they include redundant and incorrect information. In the meantime, PCA attempts a mathematical approach to the source information while retaining specific information.

### 3.2.2 Auto Encoder

Following Principal component analysis, the generated characteristics are essentially a linear function of the unique information. As a consequence of this, an enhanced version of the principal component analysis characteristics, in addition to raw attributes, are fed into an extraction of feature abstraction in order to acquire knowledge of high-level and complex features. These features will then be utilized in the subsequent categorization algorithm in order to acquire non-linear affiliations between the gestures of genetic variations. After this, the newly streamlined set of features ought to be in a position to condense the majority of the information that was included in the first set of features. For extracting features, researchers use an autoencoder neural net, and the setup specifics are below. As previously stated, both of the feature learning stages are independent of any particular labels, demonstrating unsupervised training.

### 4. Supervised Learning

The features created via the suggested two-phase unsupervised feature learning method are then attached with a set of attributes for classifier learning to accomplish the objective of forecasting treatment outcomes for people with cancer. In this research, they use a variation of the AdaBoost algorithm as the classifier's learning strategy because it performs well in classification techniques. The PCA-AE-Ada is the label given to this classifier. The process of classifier study was based on sample labels, indicating that it is supervised.

The following are the main characteristics of our suggested technique for forecasting treatment outcomes in breast cancer victims:

• At first, the set of gene expressions is presented that $\{a_i\}_{i=1}^{m}$

• Secondly, to acquire a compressed set of features, principal component analysis is used.

• And then the raw gene expression data and flattened feature are merged. The input of the autoencoder neural network is taken to acquire more intricate representations utilizing the deep learning approach.

• Lastly, the compressed feature and deep representation are concatenated in an inclusive manner that is utilized for training an ensemble classifier.

• In particular, a baseline classification termed PCA-Ada is built as a comparative in the assessment studies conducted in this research, utilizing PCA trampled finding as input features to compare with the classification built utilizing features created from the two-stage feature training architecture.

## 5. Experimental Analysis

### 5.1 Data Configuration

On the one side, whenever the principal component analysis approach is utilized to minimize the dimension in gene expression profiles, researchers could get compressed feature vectors with varied dimensions because each dataset is diverse in size, which is inconsistent with the extracting features neural network's input size. Researchers padded all the feature maps with zero values to constrain the selected features in almost the same dimensions while causing efficiency degradation, to conform the system researchers built to data with varying dimensions.

The principal component analysis AdaBoost algorithm model, on either side, has a feature dimensionality of (m-1)+32, whereas the baseline system PCA-Ada has a feature dimensionality of m-1, implying that all these two approaches are dealing with dissimilar hyper-parameter combinations. As more than just a result, in addition to making a much more valid comparison, researchers introduced 33 raw gene expression relevant information to the PCA-Ada inputs randomly, forcing these two approaches to extra comparable formations without adding data redundancy.

### 5.2 Feature Extraction Objective Function

The layered autoencoder was used in the second stage of the suggested feature learning strategy to form a deep neural technique by overlaying numerous auto-encoders framework.

The encoder and decoder components of the auto-encoder are made up of many non-linear modification layers that take the merged approximations of the original information $\widetilde{A}$ as input as expressed in equations (1) and (2):

$$y^{(1)} = \alpha\left(\varphi^{(1)}\widetilde{A} + u^{(1)}\right) \tag{1}$$

$$y^{(j)} = \alpha\left(\varphi^{(j)}y^{(j-1)} + u^{(j)}\right) j=2,\ldots.n \tag{2}$$

In the above equation n indicates the layer of number and $\alpha$ specifies the activation function. The hidden vector, bias vector, and weight matrix are indicated as $y^{(j)}$, $u^{(j)}$ and $\varphi^{(j)}$ among the j-th layer.

### 5.3 Activation Function

As an activation function, researchers used ELU (Exponential Linear Unit), which quickens deep neural network training and improves classification performance mentioned in Equations (3) and (4):

$$\omega\left(y^{(j)}\right) = \begin{cases} \delta\left(\exp\left(y^{(j)}\right) - 1\right), & if \, y^{(j)} \leq 0 \\ y^{(j)} & if \, y^{(j)} > 0 \end{cases} \tag{3}$$

$$\omega'(y^{(j)}) = \begin{cases} \omega\left(y^{(j)}\right) + \delta & if \, y^{(j)} \leq 0 \\ 1 & if \, y^{(j)} > 0 \end{cases} \tag{4}$$

The ELU hyperparameters $\omega$ (set = 1:0) determines the point where a net negative input ELU saturates.

### 5.4 Optimization

Researchers retrain the auto encoder computational model by minimizing the squared restoration error with such a sparsity penalty using a stochastic gradient-based optimization technique called Adam (adaptive moment estimation). Adam is an acronym for adaptive moment estimation. According to the findings of the research done, this is the most effective method for accelerating the gradient-based improvement process. Adam requires a low memory need for the first-order gradient during each iteration of the procedure. Adam is the optimizer to use in this situation if one wants to train the neural network more effectively and in less time. Because Adam executes gradient clipping implicitly coordinate-wise, it is able, in contrast to SGD, to deal with heavy-tailed noise. It combines the advantages of AdaGrad (this performs well enough with sparse slopes) and RMSProp (this performs well in online and non-stationary situations) to calculate adaptive learning rates individually for various variables from predictions of first and second values of the

gradients in this instance. The sample size in this research is 64, the repetition times are 10k, and the training error is 0.001. The default values for the other variables are used.

## 5.5 Adaboost Algorithm

The AdaBoost algorithm is utilized to determine the classifiers in the classification process of learning their suggested protocol, which is a supervised learning meant to generate a classifier model that best differentiates the positive and negative occurrences. As shown the sequence of the training phase, $\{(a_i', b_i)\}_{i=1}^m$ where $a_i'$, signifies the training dataset and $b_i$ is a Boolean value allocated during the database pre-processing step based on the clinical characteristics of cancer patients. AdaBoost is a useful technique for improving the classification accuracy of a simple training procedure by integrating a group of weak classifiers $\{h_j(a')\}$ into a stronger classifier $h(a')$. Researcher utilized decision stumps as the weak classifier learning algorithm in this paper. If $a'$ is categorized as a positive instance, $h(a')$ returns 1; alternatively, it returns 0.

They [22] presented a variation of the AdaBoost method, that the research uses below. This form limits the base classifiers to only using one information $v_j$. As a consequence, each weak classifier has a feature space $v_j$, a threshold j, and a equivalence parity $pa_j$ that is also 1 or -1, signifying the disparity's orientation mentioned in equation (5).

$$h_j(a')=\begin{cases} 1 & \text{if } pa_j v_j(a') < pa_j \theta_j \\ 0 & \text{otherwise} \end{cases} \tag{5}$$

$$(pa_j, \theta_j)=arg \min_{(pa_k, \theta_k)} \sum_{i=1}^m \left| hk(a_i') - b_i \right| \tag{6}$$

For every weak classification $h_j(a')$ , the boosting method creates the ideal values for $pa_j$, and $\theta_j$ so that the amount of misclassified training instances is reduced. To do so, it examines all potential groupings of the both $pa_j$, and $\theta_j$ the amount of which would be restricted because only an infinite number of iterations examples are provided in equation (6).

## 6. Result and Discussion

Four common approaches for distant metastatic forecasting were examined to determine the classification method: the 70-gene classification algorithm, the 76-gene classification algorithm, and two variants of gene set data classification model: set centroid and set-median. The 70-gene and 76-gene classifiers are very well gene signature-depending on the approaches for predicting cancer outcomes, however, gene set statistical classifiers were shown to function similar to the earlier classifications as well as being more reliable. A maximum of 70 genes were chosen as gene characteristics in the 70-gene classifier. The pattern of the different outcome categories (positive result and negative result) was determined to use the mean vector of the 70 gene expression profiles, and the data was allotted to more correlated categories using Person's correlations. A overall of 76 genes were chosen as 76 gene characteristics in the 76-gene classifier. The weighted linear combination of the 76 genes' predicting the value was used to construct a relapse score with each model, and each sample was then allocated to one of two performance categories based on whether the relapsed value was greater than a cutoff. The gene set statistics classifier first retrieved pre-specified genetic variants from the MSigDB dataset, then calculated statistical values to determine the best feature set obtained from the genetic variants, which then was utilized to build the centroid classification. PCA, set-median, Set-centroid, and t-test are some of the statistical approaches utilized to examine the gene sets. The set-centroid and set-median statistics approaches were chosen since they have been shown to perform much better than the others.

Moreover, researchers create a baseline classification using the same procedure as the suggested technique to compare the deep learning classifiers with the non-deep learning classification. In cancer databases, there is a significant disparity between the number of patients who have favorable results and those who have terrible outcomes. In database GSE11121, for instance, there are only 28 bad result cases contrasted to 154 good result cases. The MCC (Matthew's correlation coefficient, presented in Equation (7) and the AUC, which are said to be the greatest dependable quantify criteria when the set of data dispersion is highly unbalanced, have been used as two main evaluation measures in this situation.

$$Matthews\ Correlation\ Coefficient= \frac{T^+ T^- - F^+ F^-}{\sqrt{(T^++F^+)(T^++F^-)(T^-+F^+)(T^-+F^-)}} \tag{7}$$

$TE^+$ indicates true positive, $TE^-$ indicates true negative, $FE^+$ indicates false positive, and $FE^-$ indicates false negative. This part also includes accuracy (ACC), specificity (SP), and Recall (R), which were explained in the equations (8) (9) and (10).

$$Recall = \frac{TE^+}{TE^+ + FE^-} \tag{8}$$

$$Accuracy= \frac{TE^+ + TE^-}{TE^+ + TE^- + FE^+ + FE^-} \tag{9}$$

$$Specificity = \frac{T^{E^-}}{TE^- + FE^+} \tag{10}$$

To extract features from gene expression profiles and develop a classifier, researchers used GSE2034 as the combined training assessment datasets in this research. GSE2034, in particular, was subjected to ten times five-fold validation set. The entire dataset was randomized and classified into 5 groups each session, four of which were used to train a model and the remaining to evaluate the provided strategy. The final score was calculated by averaging all of the data and provided as the statistical assessment. During the test phase, we also ran multiple studies on the other four GEO databases. Table 2 and Table 3 contain all of the data for the two approaches that were suggested and the graphical representation based on Tables are presented in Figure 5 and Figure 6.

**Table 2.** PCA-AE Adaboost Classifier Performance

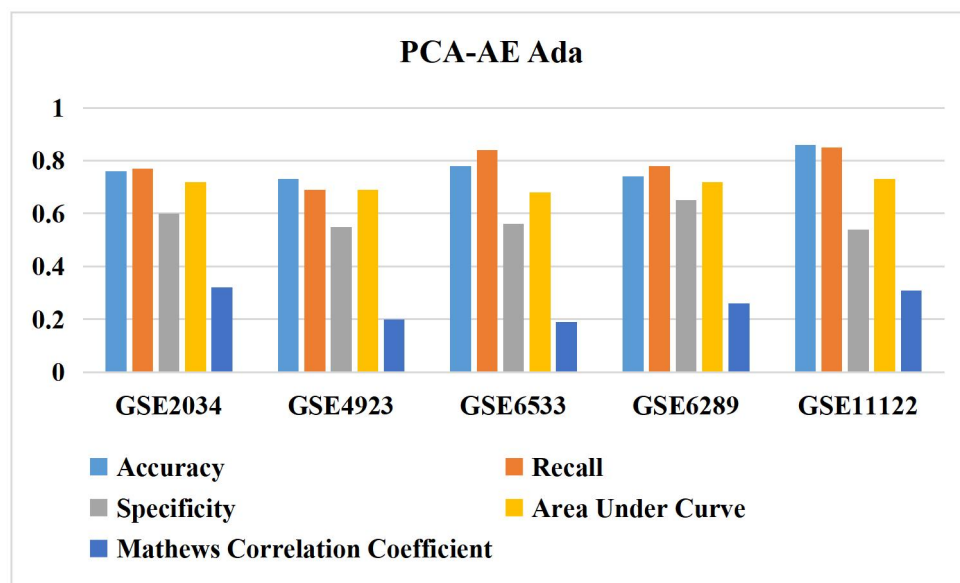| Info | Accuracy | Recall | Specificity | Area Under Curve | Mathews Correlation Coefficient |
|------|----------|--------|-------------|------------------|--------------------------------|
| GSE2034 | 0.76 | 0.77 | 0.60 | 0.72 | 0.32 |
| GSE4923 | 0.73 | 0.69 | 0.55 | 0.69 | 0.20 |
| GSE6533 | 0.78 | 0.84 | 0.56 | 0.68 | 0.19 |
| GSE6289 | 0.74 | 0.78 | 0.65 | 0.72 | 0.26 |
| GSE11122 | 0.86 | 0.85 | 0.54 | 0.73 | 0.31 |



**Figure 5.** PCA-AE Adaboost Classifier Performance

**Table 3.** PCA-Ada Classifier Performance

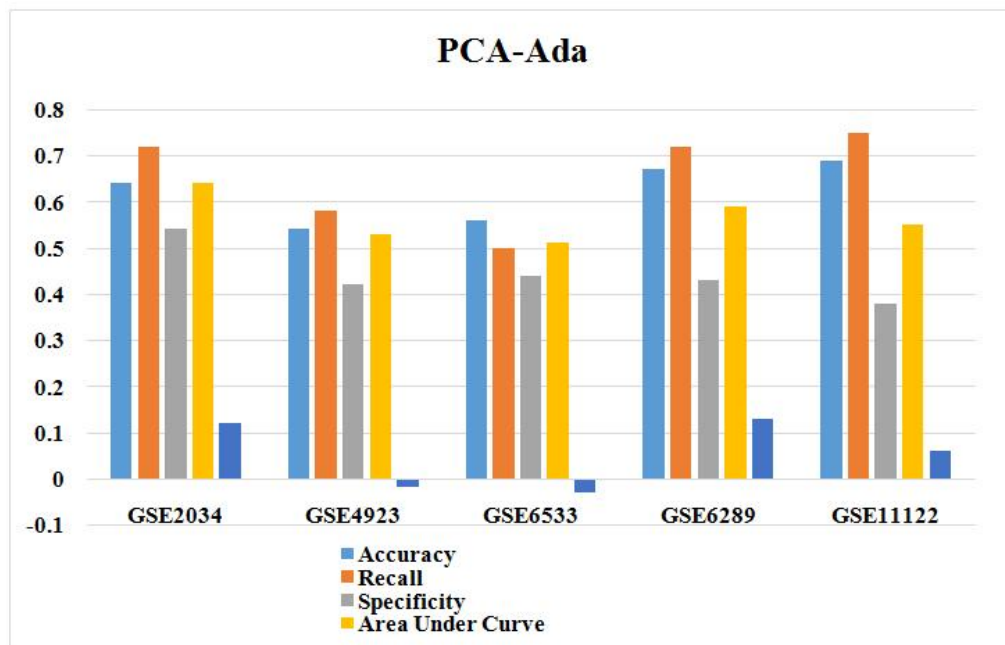| Info | Accuracy | Recall | Specificity | Area Under Curve | Mathews Correlation Coefficient |
|------|----------|--------|-------------|------------------|--------------------------------|
| GSE2034 | 0.64 | 0.72 | 0.54 | 0.64 | 0.12 |
| GSE4923 | 0.54 | 0.58 | 0.42 | 0.53 | -0.02 |
| GSE6533 | 0.56 | 0.5 | 0.44 | 0.51 | -0.03 |
| GSE6289 | 0.67 | 0.72 | 0.43 | 0.59 | 0.13 |
| GSE11122 | 0.69 | 0.75 | 0.38 | 0.55 | 0.06 |

**Figure 6.** PCA-Ada Classifier Performance

Table 2 demonstrates that the proposed PCAAE-Ada classifier efficiency well both on training validation and independent testing dataset, and it operates consistently. It has strong AUC ratings (almost 70%) and ACC rates (nearly 80%) and present superior efficiency in terms of accuracy in positively and negatively situations (with fairly good recall and specificity).When measuring the results of the PCA-AE-Ada classifier to that of the PCA-Ada classifier, researchers could see that PCA-AE-Ada with deep learning approaches outperforms PCA-Ada with compressed PCA characteristics for all assessment criteria [21]. Notably, the Recall and Specificity rates show that PCA-AE-Ada performs best for recognizing both positively and negatively occurrences, whereas PCA-Ada only performs well enough for positive examples. This shows that deep learning may efficiently address the issue of imbalanced training dataset distributions while also improving classifier generalization. On the multiple datasets, the Area under the curve and MCC values using this technique as well as the other four classifiers are displayed in Figure 7 and Figure 8, respectively. The specifics of the other four ways are not displayed.
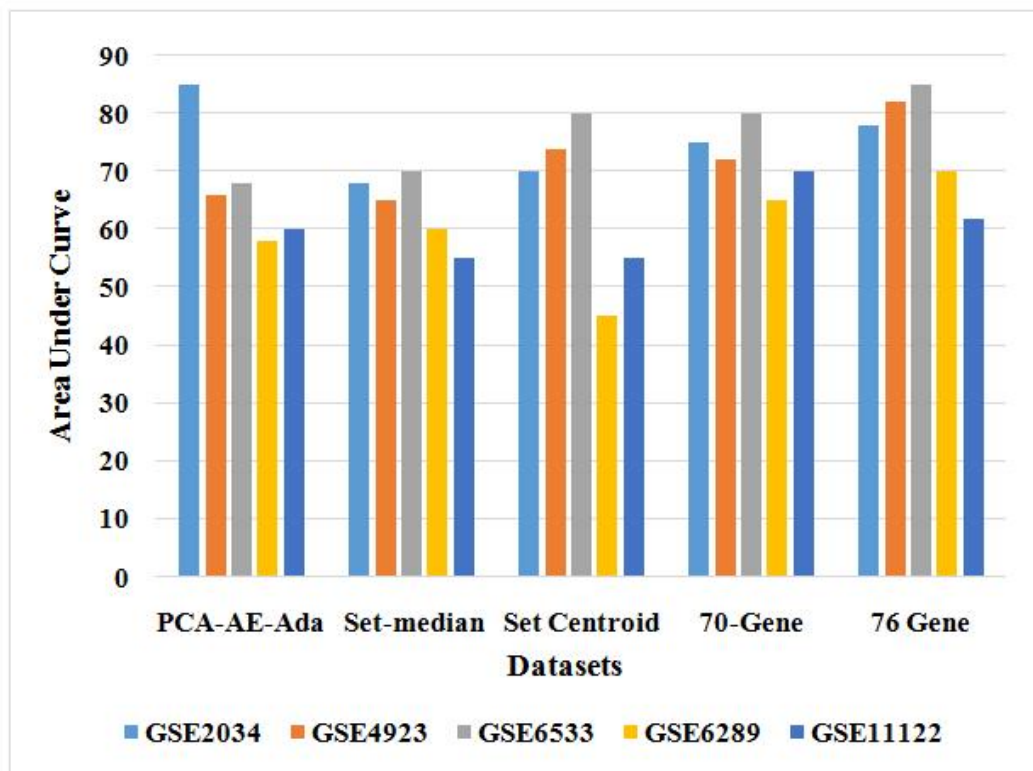


**Figure 7.** AUC Score based on Five Classifier on Five database

Over most databases based on Figure 7, the ensemble classifier has the best Area under the roc curve, albeit it doesn't function as well on GSE11121 as the two genetic set-based techniques. Then there are the gene set-based algorithms, which generate superior AUC results than the two-gene signature classifications. MCC scores demonstrate a similar tendency.
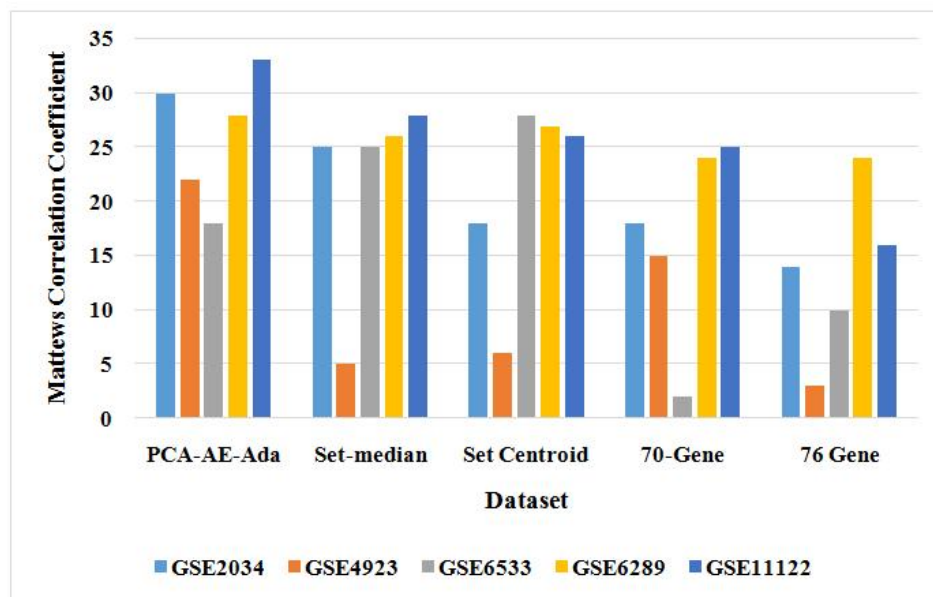


**Figure 8.** MCC Score based on Five Classifier on Five database

Except for research, four typical classifications (particularly those based on gene signatures)low performance on GSE4922 and GSE6532 than on the other three different data source (see Figure 8). This is attributed to the fact that these data sources include lymph node-positive and lymph node-negative patients, whereas the other databases exclusively include lymph node-negative victims. Interestingly, researcher ensemble classifier has a strong performance.

The suggested technique, in comparison to the other classification techniques, is much less sensitive to unbalanced statistics and is, therefore, more reliable. Unfortunately, the classification stops working on dataset GSE6532; this could be due to an inherent bias toward the right path, as they are the more difficult situations in the clinical world that necessitate adequate treatment preparation. Furthermore, to show the efficiency of the proposed method in a much more comprehensive way, the area under curve and Matthews correlation coefficient of the approach were aggregated over the five public NCBI databases, and the final findings were presented in Figure 9.
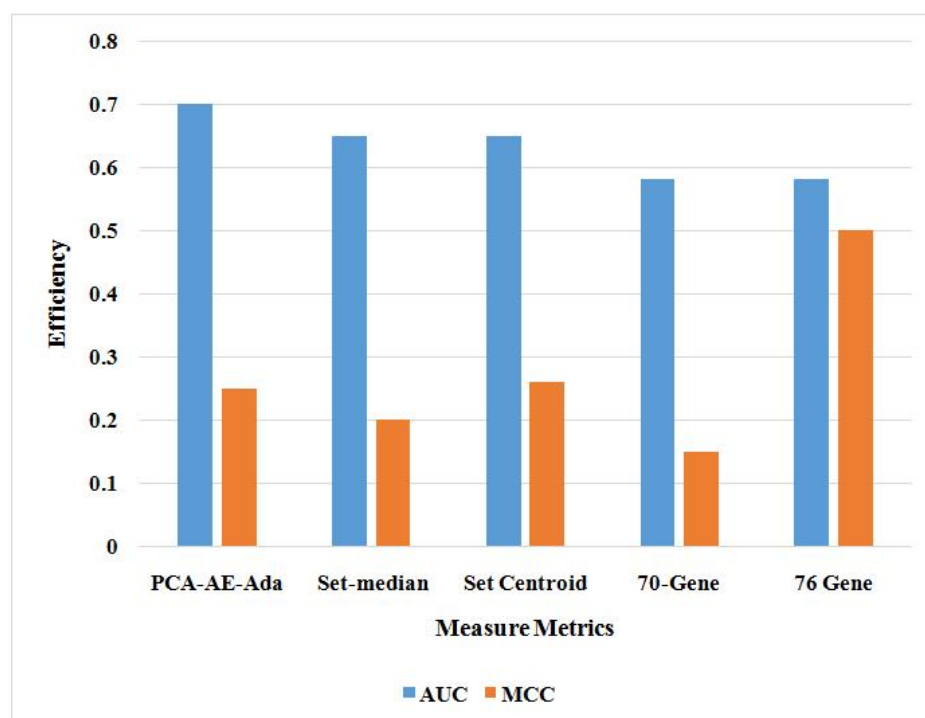


**Figure 9.** Overall efficiency

The predictor has an Area under curve of over 0.715, whereas the two gene set classifiers have a substantially lower outcome of around 0.55, according to a thorough examination. The two gene signature classifier, however, could also achieve an AUC of less than 0.6. The indices of the MCC show a similar pattern. In summary, the ensemble classifier based on PCA and autoencoder characteristics outperform other existing approaches in classification performance as well as generalization ability across diverse databases.

## 7. Conclusion

Researchers describe a new strategy for predicting the medical outcomes of people with cancer utilizing deep learning in this work. The auto-encoder and neural system and Principal Component Analysis are merged with the exportation of deep learning approaches to learning more suitable parameters from gene expression profiles in the feature learning stage. They use the AdaBoost technique to build an ensemble classifier for the completed prediction problem during the classifier process of learning. As evidenced by the assessment testing results, the suggested technique has a stronger forecasting ability, and the classifier built using deep learning techniques outperforms the others. The characteristics automatically generated by the neural system provided a good ability for fast generalization and explicitly increased the efficiency of diagnosis and forecast, according to the research will be discussed. Our classifier, on the other hand, has some flaws. To begin with, the model created is difficult to analyze—a common difficulty with neural networks. Furthermore, determining which traits are most essential for the forecasting task is challenging. Finally, due to the deep learning model's complicated structure, the quantity of storage is less capable of approximating. While the approach has shown promising results, more publicly available variables are needed to improve generalization capacity.

## Reference

[1] G. Chugh, S. Kumar, and N. Singh, "Survey on machine learning and deep learning applications in breast cancer diagnosis," Cognitive Computation, vol. 13, no. 6, pp. 1451–1470, 2021.

[2] Ali, L., Rahman, A., Khan, A., Zhou, M., Javeed, A., & Khan, J. A. (2019). An automated diagnostic system for heart disease prediction based on $\chi^2$ statistical model and optimally configured deep neural network. IEEE Access, 7, 34938–34945.

[3] World Health Organization. Cardiovascular diseases. https://www.who.int/health-topics/cardiovascular-diseases (Accessed April 19, 2025)

[4] Li, H., & Boulanger, P. (2020). A survey of heart anomaly detection using ambulatory electrocardiogram (ECG). Sensors, 20(5), 1461.

[5] Duncker, D.(2021). Smart wearables for cardiac monitoring—real-world use beyond atrial fibrillation. Sensors, 21(7), 2539.

[6] Isakadze, N., & Martin, S. S. (2020). How useful is the smartwatch ECG? Trends in Cardiovascular Medicine, 30(7), 442–448.

[7] Kushwaha, S., Chithras, T., Girija, S. P., Prasanth, K. G., Minisha, R. A., Dhanalakshmi, M., Jayanthi, A., Robin, C. R. R., & Rajaram, A. (2024). Efficient liver disease diagnosis using infrared image processing for enhanced detection and monitoring. Journal of Environmental Protection and Ecology, 25(4), 1266–1278.

[8] Kumar, V., & Kushwaha, S. (2023). Hybrid metaheuristic model based performance-aware optimization for map reduce scheduling. International Journal of Computers and Applications, 45(12), 776-788. Taylor and Francis Online.

[9] Sandhu, M., Kushwaha, S., & Arora, T. (2023). A comprehensive review of GAN-based denoising models for low-dose computed tomography images. International Journal of Image and Graphics, 2023, 1-38. World Scientific Publishing Co.

[10] Kushwaha, S. (2023). An effective adaptive fuzzy filter for speckle noise reduction. Multimedia Tools and Applications, 2023, 1-16. Springer.

[11] Dhanalakshmi, P., Venkatesh, V., Ranjit, P. S., Hemalatha, N., Divyapriya, S., Sandhiya, R., Kushwaha, S., Marathe, A., & Huluka, M. A. (2022). Application of machine learning in multi-directional model to follow solar energy using photo sensor matrix. International Journal of Photoenergy, 2022, 1-9. Hindawi.

[12] Hemanand, D., Mishra, N., Premalatha, G., Mavaluru, D., Vajpayee, A., Kushwaha, S., & Sahile, K. (2022). Applications of intelligent model to analyze the green finance for environmental development in the context of artificial intelligence. Computational Intelligence and Neuroscience, 2022, 1-8. Hindawi.

[13] Singh, A., Khan, R. A., Kushwaha, S., & Alshenqeeti, T. (2022). Roll of Newtonian and Non-Newtonian motion in analysis of two-phase hepatic blood flow in artery during jaundice. International Journal of Mathematics and Mathematical Sciences, 2022, 1-10. Hindawi.

[14] Singh, A., Kushwaha, S., Alarfaj, M., & Singh, M. (2022). Comprehensive overview of backpropagation algorithm for digital image denoising. Electronics, 11(10), 1590. MDPI.

[15] Kushwaha, S., & Singh, R. K. (2015). Study and analysis of various image enhancement method using MATLAB. International Journal of Computer Sciences and Engineering, 03(01), 15-20.

[16] Kushwaha, S., Kondaveeti, S., Vasanthi, S. M., W, T. M., Rani, D. L., & Megala, J. (2024). Graph-informed neural networks with green anaconda optimization algorithm based on automated classification of condition of mental health using alpha band EEG signal. 2024 4th International Conference on Sustainable Expert Systems (ICSES), 44–50.

[17] A, S. B., S, S., S, R. S., Nair, A. R., & Raju, M. (2022). Scalogram based heart disease classification using hybrid CNN-naive Bayes classifier. In 2022 International Conference on Wireless Communications Signal Processing and Networking (WiSPNET) (pp. 345–348). IEEE.

[18] Bulbul, H. I., Usta, N., & Yildiz, M. (2017). Classification of ECG arrhythmia with machine learning techniques. In 2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA) (pp. 546–549). IEEE.

[19] Sharma, P., & Gupta, D. V. (2018). Disease classification from ECG signal using R-peak analysis with artificial intelligence. International Journal of Signal Processing, Image Processing and Pattern Recognition, 11(3), 29–40.

[20] Ullah, A., Anwar, S. M., Bilal, M., & Mehmood, R. M. (2020). Classification of arrhythmia by using deep learning with 2-D ECG spectral image representation. Remote Sensing, 12(10), 1685.

[21] Khan, M. U., Aziz, S., Naqvi, S. Z. H., & Rehman, A. (2020). Classification of coronary artery diseases using electrocardiogram signals. In 2020 International Conference on Emerging Trends in Smart Technologies (ICETST). IEEE.