# Deep Convolution Neural Network Based Breast Cancer Bigdata Analysis for Crowd Cloud Sourcing

**Rethinakumar, GopinathGanapathy, Jeong-Jin Kang**

*Abstract: Breast cancer is one of the dangerous diseases leads fast death among women. Several kinds of cancers are affecting people, but breast cancer affects highly women. In medical industry removal of women breasts or major surgery is taken forward as the solution, where it reoccurs after surgery also. Only solution to save women from breast cancer is to identify and detect the earlier stage of cancer and provide necessary treatment. Hence various research works have been focused on finding good solution for diagnosing and classifying the cancer stages as benign, malignant or severe malignant. Still the accuracy of classification needs to be improved on complex breast cancer datasets. Few of the earlier research works have proposed machine learning algorithms, which are semi-automatic and accuracy is also not high. Thus, to provide a better solution this paper aimed to use one of the deep learning algorithms such as Convolution Neural Networks for diagnosing various kinds of breast cancer dataset. From the experimental results, it is obtained that the proposed deep learning algorithms outperforms than the other algorithms.*

*Keywords: Breast Cancer, Convolution Neural Network, Deep Learning, Diagnosis, Prediction, Benign, Malignant.*

## I. INTRODUCTION

One of the most usual cancers is Breast Cancer, affects women highly worldwide, exemplifying the majority of other cancer types and related deaths according to global statistics, considering it as a significant public health issue in present's society. The amount of death can be controlled by earlier diagnosis and it save people life significantly. Accurate diagnosing finds the stages of the cancer as benign, malignant, and severe malignant; it leads to apply timely treatment to the cancer, and it avoid patients undergoing improper treatments. Therefore, accurate diagnosis of cancer classes is the matter of several recent researches. Numerous approaches have been proposed for diagnosing the breast cancer data whereas the accuracy of classification is less. Since the breast cancer dataset has various and unique advantages of clinical features, machine learning algorithms have been used. But the machine learning algorithms are not fully automatic. Hence to provide an automatic learning, detection, and classification of breast cancer dataset is applied using deep learning algorithms. Deep learning algorithms are recognized as better method and used for breast cancer classification and prediction. Generally, data analytics and data mining approaches are widely used classifying complex datasets. Particularly in medical industry those approaches are widely applied for diagnosing and decision making. The earlier

**Rethinakumar**, Assistant Professor, Dept of Information and Communication, Dong Seoul University, Korea.Research Scholar, Bharathidasan University, India.

**GopinathGanapathy**, Registrar, Bharathidasan University, India.

**Jeong-Jin Kang**, Professor, Dept of Information and Communication, Dong Seoul University, Korea.

methods extract the standard features for classifying the data. Though the classification accuracy is less.

The main objective of this paper is to analyse and observe the features used to predict the cancer class as benign or malignant. To extract a greater number of features and hidden features deep learning algorithm is used. Convolutional Neural Network is used as the deep learning algorithm for analysing and diagnosing the complex data. Deep learning algorithm learn the data using striding method. Convolution value is calculated from the stride function, where stride function can fetch the features from various sized images.

## II. RELATED WORKS

To understand the issues and challenges faced by the earlier research works, this section presents a detailed survey on various methods focused on breast cancer analysis. Author in [1] explained breast cancer is one of the most important kind of cancers among various cancers. Breast cancer is considered as serious cancer type, it is a hot research topic with great value [2]. Using data science and machine learning algorithms healthcare industries developing a great assistance for medical practitioners in decision making. Nowadays, it is a big challenge of diagnosing the pattern of breast cancer, since patterns are different is shape, texture and other clinical features. So, the healthcare industry is paying more attention in developing an efficient application using machine learning algorithms [2]. In the earlier works, some of the researchers have focused on detecting breast cancer using image analysis for analysing the cancers have spread beyond the breast, other organs and nearby lymph nodes [3-5], and cell biology [6-8] using selective but small datasets from algorithm evaluation challenges [9-12].

Some of the earlier research works have focused on machine learning algorithms for diagnosing cancer dataset [13]. The dataset used in the experimenting the machine learning algorithms is Wisconsin Diagnostic Breast Cancer dataset [14], and obtained significant output. Authors in [15] used Gated Recurrent Unit [16] combined with Support Vector Machine [17] algorithms for diagnosing breast cancer on WDBC dataset. Some of the research works have used multiple kernel learning [18] method for various healthcare problems [19-20] and provides the way to examine the learned model. One of the authors combined SVM [21] and Random optimization [22] methods for demographic, clinical and biochemical data prediction. One of the authors in [18] utilised MKL method for cancer-based

378

thrombosis risk analysis to predict the progression of cancer disease in an oncology setting of breast cancer patients. Also, the author has provided a proof of concept for assessing the MKL based decision support system is really useful for cancer diagnosis.

## III. LIMITATIONS AND MOTIVATION

Some of the limitations identified from the above discussion are, the number of cancer classes are more, whereas SVM can classify only two (positive and negative), so it needs a classifier can find more classes. Thinking about Gated Recurrent Unit, it is highly suitable for huge volume of numerical or time series, and dependent data analysis. Also, combining GRU with SVM creates more time and computational complexity. Cell biology is not a common method used in all the medical centers. In the initial stage the breast cancer analysis has done over breast images and clinical data of the breast cancer, not on tissue analysis. In addition to the above said limitations, the detection and classification accuracy is less and all the methods are tested on different datasets.

From the above discussion it is identified that still the medical / healthcare industry requires a fully automatic and efficient algorithm for diagnosing and classifying breast cancer stage. Hence, to improve the accuracy of classification this paper used a fully automatic and effective method as convolution neural network and prove the performance. To do that, the entire contribution of the work is carried out into two different stages such as data acquisition from IoT industry and data analysation using convolution neural network.
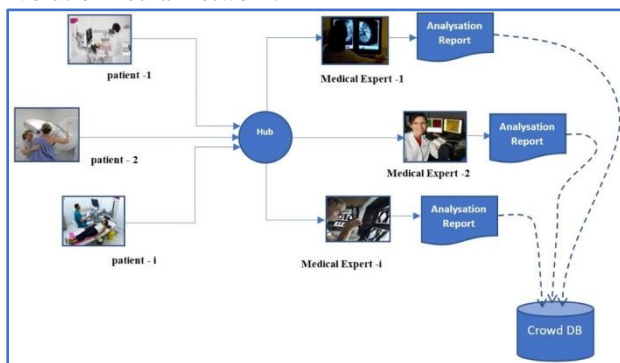


**Figure-1.IoT Based Breast Cancer Data Acquisition System**

The initial stage of the work has concentrated on data acquisition from medical network. The sample network is illustrated in Figure-1. The number of patients $P = \{P_1, P_2, ..., P_i, ... P_n\}$, are monitored using different IoT devices such as, mammogram, CT scan and other scanning devices interconnected in the healthcare. The set of all data monitored from the devices are directly feed into the internet hub (it is a software model) where it collects and transmit to the corresponding server. From the server various medical experts $ME = \{ME_1, ME_2, ..., ME_i, ..., ME_m\}$ access the data with authorisation and process the data using various data/image analytical methods to predict the health condition of the patients and the medical report/summary is uploaded in the crowd cloud DB. The final analyzation reports are sent to cloud and the volume of data is increased

into a high amount, which is unimaginable. Since, the data is coming from various medical industries, IoT devices and medical experts, it become crowd data. The crowd data has two major kinds such as breast cancer data and the appropriate medical report used for enhancing the detection and classification accuracy.

## IV. DEEP LEARNING ALGORITHMS

One of major part of Artificial Intelligent (AI) algorithms to imitate human brain regarding data processing, pattern creation, and decision making is deep learning. It is a subset of machine learning in AI, that has high ability of leaning unstructured data. Deep learning is also called as deep neural network or deep neural learning. Comparing with AI and machine learning algorithms, deep learning algorithms are highly capable to learn large volume of unstructured data automatically. Hence, mainly deep learning methods have been used for bigdata analytics. Every region of the world using deep learning method for evolving and exploring digital era. The bigdata is obtained from various resources like internet, social websites, online movies, e-commerce and search engines, and etc. These data are particularly available and can be accessed by any internet applications belongs to cloud computing. Since the data is unstructured, it is highly difficult to extract the relevant information from the data.

Most of the IT industries trust that AI methods can incredibly extract the information from the data and it will be an automated support. Most of the earlier data processing industries are using machine learning method for bigdata analytics, and they are not fully automatic and self-adaptive. Deep learning is fully automatic, learn the data deeply, uses a hierarchical level of ANN to complete a task based on the machine learning process. ANN is created in accordance to the human brain, with neuron nodes interconnected together like www. Conventional methods analyse the data linearly. But the hierarchical method of deep learning process enables machine to process data using a non-linear approach. In terms of complex dataset like fraud detection, spatial, temporal and other geo-space data, deep learning methods can fetch the important features like time, geographical locations, IP addresses, and the type of retailer. More number of layers in the deep learning method performs their functions efficiently and classify the data. Due to more advantages, this paper utilizes Convolution Neural Network to analyse the breast data.

## V. PROPOSED CONVOLUTION NEURAL NETWORK

The entire research work is carried out into three different stages such as, data generation, data analysis and data prediction, which is illustrated in Figure-2. In the first stage the data is generated from patient body suing various medical IoT devices like MRI, X-RAY, Mammogram and so on. The data collected from various resources may be in the form of images and/or alpha numerical, stored in the local admin system and transferred to cloud DB. The cancer data stored in the cloud DB can be accessed by only the

authorised medical experts from anywhere at any time. The accessed data is analysed by the medical experts and predict the cancer classes using various algorithms. Most of the cases the treatment and medicine information are also uploaded in the cloud DB. This data with the analysation report is updated regularly with the treatment and medicine details. This updated data is included in the cloud crowd data can be used for training the model. This paper proposed a CNN model for breast cancer disease diagnosing and predicting the classes. Before applying CNN algorithm, it is essential to understand the various functions involved in the overall process [20]. For example, initially the input image is stored in the form of matrix, hold pixel values. A 3 x 3 image matrix is also considered as a 9 x 1 vector. It is feed into multi-level perceptron to classify.
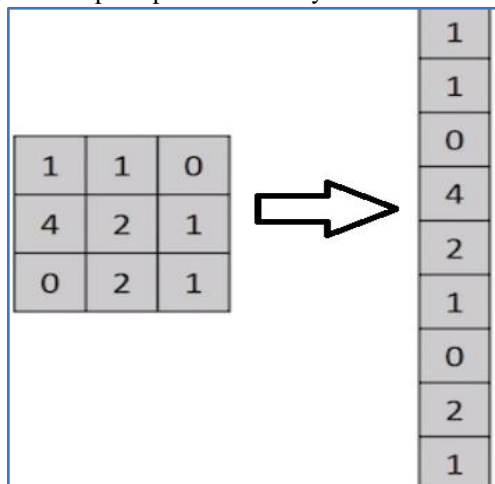


**Figure-2. Data Matrix Representation**

Figure-2 illustrates the matrix and vector form of an image. The CNN can obtain the spatio temporal information based on the appropriate filters. A training process is applied on the images in order to understand the entire image information in better manner. The 4 x 4 x 3 image is shown in Figure-3, where it shows the three layers such as R, G, and B, where the learning process uses 3 x 3 of the data from top-left corner to bottom-right corner using stride function.
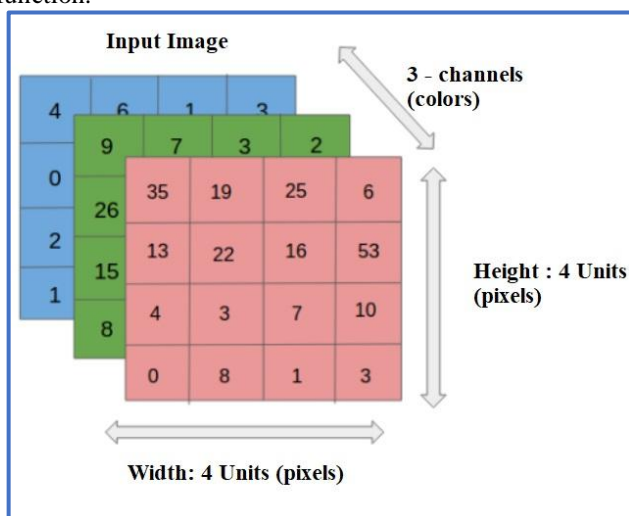


**Figure-3. Image Representation**

From this, it is understanding that, learning an image with high dimensions like 7680 x 4320. The CNN method reduces the dimensionality without any data loss for improving the predicting accuracy. Hence it is essential to

design and implement a novel architecture for learning the entire features. But it should be scalable for any size of dataset. To do this CNN used convolution layer as the kernel.
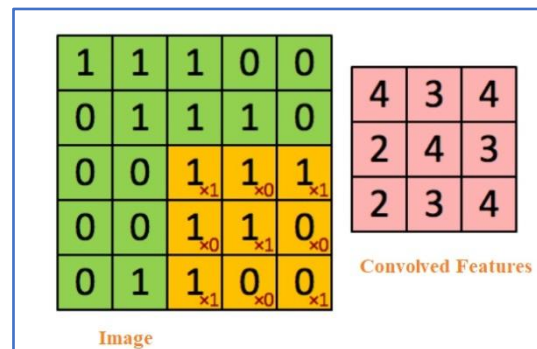


**Figure-4. Convoluting a 5x5x1 image using kernel 3x3x1 provides 3x3x1 convolved feature**

For example, a 5x5x1 image is carried out into a convolution operation (is the first operation carried out by the first set of convolution layer). Convolution layers is also called as kernel or filter (highlighted in yellow colour in Figure-4). The kernel size is 3 x 3, indicated by k=3. To learn the entire image, the kernel performs 9 times, since the stride length is 1. The kernel movement is illustrated in Figure-5.
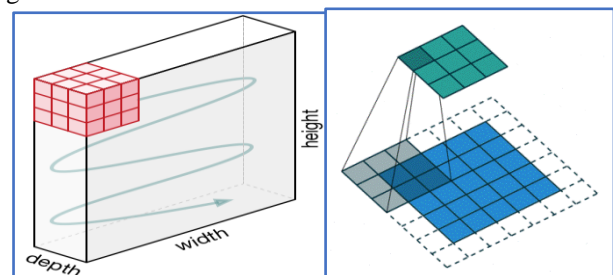


**Figure-5. Movement of the Kernel**

Initially the input image is feed into CNN-1 layer. The convoluted output from CNN-1 is the activation map. Various filters are applied in the CNN-1 layer for extracting the important features from the input image is used as the input to the next layer. Filters in the CNN-1 can extract only different features lead to obtain the correct the classes. In order to maintain the size of the image, image-padding (valid padding) process is applied, where it reduces the number of features. Following CNN, pooling layers are added to reduces the number of parameters. In this paper two stages of CNN (CNN-1, CNN-2) is added with pooling layers for predicting the accurate pattern. CNN layers help to learn deeply and extract the features. Since CNN learn the data deeply, high number of hidden features are also extracted and compared with other shallow network where the features are highly general. Output obtained from CNN-2 is feed into Fully Connected layer. The input of other layers is flatted and sent for transform the output as the number of classes as desired by the network.

The final output is generated from the output layer and compare the percentage of error. To do that a loss function is applied in the FC layer to calculate the mean square loss and gradient or error is calculated and backpropagated to

update the filters used in the CNN. One cycle of training process is completed using a single forward and backward pass. To improve the efficiency of cancer detection and classification a new CNN architecture is proposed (see Figure-3). CNN has four different layers as CNN-1 (Convolution-1) to CNN-5 (Convolution-5) with 96 kernels, and two FC-1/2 (Fully Connected) layers. The input image (size of 224 x 224 x 3) is feed into CNN-1 and consider 11 x 11 patch for pooling with 4 strides.
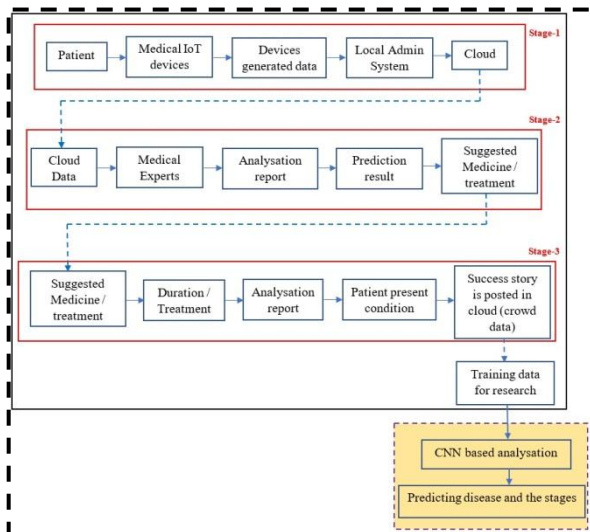


**Figure-2. Overall Architecture of the Research Work**

The maximum pooling layer is activated by ReLU method, to learn the data using pooling size of 2 x 2. The other CNN layers used 5 x 5, 3 x 3, and 3 x 3 patches for 96 kernels, 384 kernels, 384 kernels and 356 kernels respectively. Finally, the FC layers used 4096 neurons for classifying the cancer classes as normal, benign, malignant and severe malignant. All the images input to the training process are learned thoroughly and entire image features are extracted. Finally, the features are classified based on the ODD values.
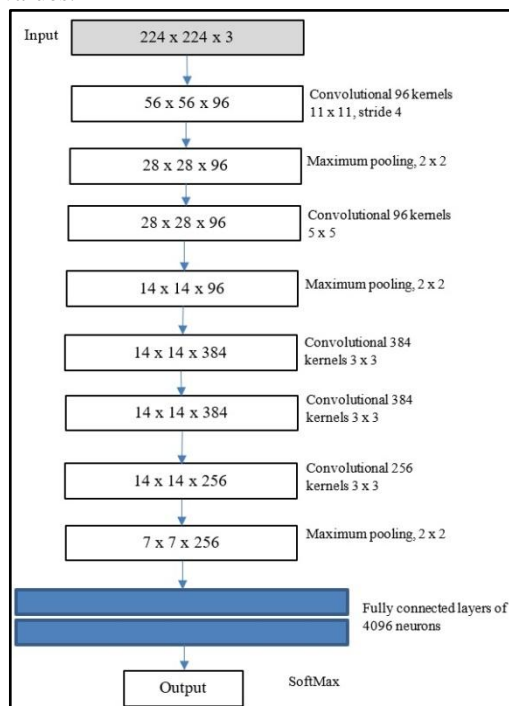


**Figure-3. CNN Architecture**

*Dataset*

To examine the performance of the proposed CNN different kinds of breast cancer dataset is used in the experiment. Some of the datasets used in the experiment are BreaKHis, MIAS, DDSM, and CBIS-DDSM. The total number of images used in the experiment is 12000, whereas it has three different classes as normal, benign and malignant. Normal defined as good breast images, benign defined as breast images with mass, not cancer, but it may become cancer, and malignant defined as the cancer. From the entire images 60% of the images is used for training the CNN model to test the remaining 40% of the images. To make the model as generalised classification model new patients' image is used for testing process and the classification results are verified. In this paper, the CNN has been introduced to the mammograms, which allow each picture to be either benign or malignant in one of the two groups. Three sets of Mammogram dataset BreaKHis, MIAS, DDSM, and CBIS-DDSM has been selected for the most commonly used verification by MATLAB. The complete dataset information is given in Table-1 for understanding the performance of the proposed CNN. Four different datasets with a number of normal and abnormal images are given. All the datasets are already used for classification and it is considered as benchmark dataset, and it is used here for experimenting and evaluating the performance.

**Table.1. Dataset Information**

| Dataset | Total Images | Normal | Abnormal |
|---|---|---|---|
| **BreaKHis** | 2000 | 1100 | 900 |
| **DDSM** | 4000 | 1400 | 2600 |
| **CBIS-DDSM** | 4000 | 1000 | 3000 |
| **MIAS** | 2000 | 600 | 1400 |
| **Total** | 12000 | 4100 | 7900 |

The proposed method has examined in the experiment using the dataset given in Table-1 and the performance measures are calculated. For example, Initially the classification accuracy is calculated using CNN and the results are given in Table-2. From the table, the basic performance values TP, TN, FP, and FN are calculated.

**Table.2. Classification Accuracy**

| Datasets | DB-Normal | DB-Abnormal | CNN-Normal | CNN-Abnormal |
|---|---|---|---|---|
| **BreaKHis** | 1100 | 900 | 1100 | 899 |
| **DDSM** | 1400 | 2600 | 1400 | 2598 |
| **CBIS-DDSM** | 1000 | 3000 | 1000 | 2999 |
| **MIAS** | 600 | 1400 | 600 | 1399 |

From the initial classification, it obtained that the proposed CNN architecture obtained an average of 100% accuracy in classifying

normal images and 99.99% of accuracy in classifying abnormal images. In terms of dataset, on BreaKHis, CNN obtained 99.98% of TP, on DDSM CNN obtained 99.92% of TP, on CBIS-DDSM CNN obtained 99.96% of TP, and in MIAS CNN obtained 99.92% of TP. The variation of TP shows the variation in the quality of the input images. The performance of the proposed CNN is calculated using the performance factors such as TP, TN, FP and FN, given in Table-3.

**Table-3. TP, TN, FP, and FN Obtained using ProposedAlexNet**

| Dataset | TN | FN | TP | FP |
|---------|----|----|----|----|
| **BreaKHis** | 1 | 0 | 0.998889 | 0.111111 |
| **DDSM** | 1 | 0 | 0.999231 | 0.076923 |
| **CBIS-DDSM** | 1 | 0 | 0.999667 | 0.033333 |
| **MIAS** | 1 | 0 | 0.999286 | 0.071429 |

From the values calculated values given in Table.3 (TP, TN, FP, FN), the set of all performance measures are calculated for DDSM dataset and given below for verification. The following calculation is obtained for DDSM dataset only. For other datasets, the values are given in Table.3.

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN} = \frac{99.9 + 1}{1 + 0 + 99.9 + .1} = 99.90$$

$$Sensitivity = \frac{TP}{TP + FN} = \frac{99.9}{99.9 + 0} = 1$$

$$Specificity = \frac{TN}{TN + FP} = \frac{1}{1 + 0.1} = 1$$

$$Precision = \frac{TP}{TP + FP} = \frac{99.9}{99.9 + 0.1} = 0.999$$

$$Recall = \frac{TP}{TP + FN} = \frac{99.9}{99.9 + 0} = 1$$

$$F1\ Score = \frac{2 * Recall * Precision}{Recall + Precision} = 0.999$$

Similar to the above calculation, the performance measures are calculated for all the four different datasets DDSM, CBIS-DDSM, MIAS, and BreaKHis image sets, and given in Table.4. From the table-4, it is noticed that the proposed approach outperforms in identifying, detecting and classifying breast mammograms collected from various data sources. To evaluate the performance once again, the obtained results in terms of AUC and accuracy is calculated and compared with various existing approaches, given in Table-4. From the table-4, it is concluded that proposed CNN outperforms than the other approaches and it is decided that it is highly suitable for medical image processing and analyzation.

**Table.4. Performance Measures Obtained from Proposed Deep Learning Algorithm**

| Dataset | Accuracy | Sensitivity | Specificity | Precision | Recall | F1-Score |
|---------|----------|-------------|-------------|-----------|--------|----------|
| **BreaKHis** | 99.88889 | 1 | 1 | 1 | 1 | 1 |
| **DDSM** | 99.92308 | 1 | 1 | 1 | 1 | 1 |
| **CBIS-DDSM** | 99.96667 | 1 | 0.333333333 | 0.98 | 1 | 0.989 |
| **MIAS** | 99.92857 | 1 | 0.37037037 | 0.983 | 1 | 0.991 |

**Table.5. Performance Evaluation**

| Research Work | Proposed Model | Dataset | Accuracy |
|---------------|----------------|---------|----------|
| Jain & Levy (2016) | DCNN | DDSM | 60% |
| Jiang (2017) | DCNN - GoogleNet | BCDR-F03 | 88% |
| Duraisamy&Emperumal (2017) | DCNN-VGG | MIAS & BCDR | 85% |
| Ragab et al. (2019) | DCNN-SVM | CBIS-DDSM | 87.2% |
| Proposed Approach | CNN | **BreaKHis** | 99.88889 |
| Proposed Approach | CNN | **DDSM** | 99.92308 |
| Proposed Approach | CNN | **CBIS-DDSM** | 99.96667 |
| Proposed Approach | CNN | **MIAS** | 99.92857 |

To evaluate the performance of the proposed CNN, the obtained accuracy is compared with the existing research works, have carried out the similar research works given in author in [17] obtained 60% on DDSM dataset, author in [18] obtained 88% on BCDR-F03 dataset, Duraisamy&Emperumal (2017) obtained 85% on MIAS and BCDR dataset and author in [19] obtained 60% on CBIS-DDSM dataset. From the overall existing system, author in [19] obtained the highest accuracy.

Comparing with the various existing method, the proposed CNN has experimented over four different datasets and the accuracy is calculated. For all the datasets, the proposed algorithm obtained the highest accuracy of 99.99%, which is higher than the other existing approaches.

From the above discussion, experimental and comparison of results it is clear that the proposed CNN architecture proves that it is highly suitable for analysing the breast cancer data.

## VI. CONCLUSION

The main objective of this paper is to design and implement a novel automatic algorithm for detecting and classifying the breast cancer data obtained from the cloud crowd data. To do that, a convolutional neural network model is implemented in MATLAB software and the results are verified. To verify the performance various kind of dataset is collected from different benchmark datasets and the experiment is carried out. From the results, it is identified that the proposed CNN outperforms in detecting and classifying breast cancer. The performance is evaluated by comparing the results with the other existing approaches and proved the proposed CNN is better than the others. Hence it is concluded that the proposed CNN is suitable for diagnosing and classifying breast cancer data anywhere at any time.

### FUTURE WORK

In future the proposed CNN model can be combined with mobile cloud and crowd computing architecture and it can also be compared with the other deep learning algorithms such as RNN, AlexNet, GRU and the performance will be evaluated.

## REFERENCES

1. 2017. Cancer Statistics.(Mar 2017), https://www.cancer.gov/about-cancer/ understanding / statistics.
2. [n. d.]. ([n. d.]). https://seer.cancer.gov/statfacts/html/breast.html.
3. Wang, D., Khosla, A., Gargeya, R., Irshad, H. & Beck, A. H. Deep learning for identifying metastatic breast cancer. arXiv preprint arXiv:1606.05718 (2016).
4. Nazeri, K., Aminpour, A. &Ebrahimi, M. Two-stage convolutional neural network for breast cancer histology image classification. In International Conference Image Analysis and Recognition, 717–726 (Springer, 2018).
5. Golatkar, A., Anand, D. &Sethi, A. Classification of breast cancer histology using deep learning. In International Conference Image Analysis and Recognition, 837–844 (Springer, 2018).
6. Albarqouni, S. et al. Aggnet: deep learning from crowds for mitosis detection in breast cancer histology images. IEEE transactions on medical imaging 35, 1313–1321 (2016).
7. Veta, M. et al. Assessment of algorithms for mitosis detection in breast cancer histopathology images. Med. image analysis 20, 237–248 (2015).
8. Rao, S. Mitos-rcnn: A novel approach to mitotic figure detection in breast cancer histopathology images using region based convolutional neural networks. arXiv preprint arXiv:1807.01788 (2018).
9. Bejnordi, B. E. et al. Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. Jama 318, 2199–2210 (2017).
10. Bándi, P. et al. From detection of individual metastases to classification of lymph node status at the patient level: the camelyon17 challenge. IEEE Transactions on Med. Imaging (2018).
11. Litjens, G. et al. 1399 h&e-stained sentinel lymph node sections of breast cancer patients: the camelyon dataset. GigaScience 7, giy065 (2018).
12. Aresta, G. et al. Bach: Grand challenge on breast cancer histology images. arXiv preprint arXiv:1808.04277 (2018).
13. Gouda I Salama, M Abdelhalim, and MagdyAbd-elghanyZeid. 2012. Breast cancer diagnosis on three different datasets using multi-classifiers. Breast Cancer (WDBC) 32, 569 (2012), 2.
14. William H Wolberg, W Nick Street, and Olvi L Mangasarian. 1992. Breast cancer Wisconsin (diagnostic) data set. UCI Machine Learning Repository [http://archive. ics. uci. edu/ml/] (1992).
15. Abien Fred Agarap. 2017. A Neural Network Architecture Combining Gated Recurrent Unit (GRU) and Support Vector Machine (SVM) for Intrusion Detection in Network Traffic Data. arXiv preprint arXiv:1709.03082 (2017).
16. Kyunghyun Cho, Bart Van Merriënboer, CaglarGulcehre, DzmitryBahdanau, FethiBougares, HolgerSchwenk, and YoshuaBengio, (2014), "Learning phrase representations using RNN encoder-decoder for statistical machine translation", arXiv preprint arXiv:1406.1078 (2014).
17. C. Cortes and V. Vapnik, (1995), "Support-vector Networks.Machine Learning 20.3", (1995), 273–297. https://doi.org/10.1007/BF00994018.
18. Gönen, M.; Alpaydın, E. Multiple kernel learning algorithms. J. Mach. Learn. Res. 2011, 12, 2211–2268.
19. Ferroni, P.; Zanzotto, F.M.; Scarpato, N.; Riondino, S.; Nanni, U.; Roselli, M.; Guadagni, F. Risk assessment for venous thromboembolism in chemotherapy treated ambulatory cancer patients: A precision medicine approach. Med. Dec. Mak. 2017, 37, 234–242.
20. Ferroni, P.; Roselli, M.; Zanzotto, F.M.; Guadagni, F. Artificial Intelligence for cancer-associated thrombosis risk assessment. Lancet Haematol.2018, 5, e391.
21. Cristianini, N.; Shawe-Taylor, J. An Introduction to Support Vector Machines and other kernel-based learning methods.Ai Magazine 2000, 22, 190.
22. Matyas, J. Random optimization. Automat. Rem. Control 1965, 26, 246–253.
23. 17. Jain A, Levy D. 2016. Breast mass classification using deep convolutional neural networks. In: 30th conference on neural information processing systems (NIPS 2016). Barcelona, Spain. 1_6.
24. 18. Jiang F. 2017. Breast mass lesion classification in mammograms by transfer learning. In: ICBCB '17. Hong Kong, 59_62 DOI 10.1145/3035012.3035022.
25. 19. Ragab DA, Sharkas M, Marshall S, Ren J. 2019. Breast cancer detection using deep convolutional neural networks and support vector machines.PeerJ 7:e6201 http://doi.org/10.7717/peerj.6201.
26. 20. https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b1164a53.

## AUTHORS PROFILE

**Rethinakumar,** Assistant Professor, Dept of Information and Communication, Dong Seoul University, Korea.Research Scholar, Bharathidasan University, India.

**GopinathGanapathy**, Registrar, Bharathidasan University, India.
**Jeong-Jin Kang,** Professor, Dept of Information and Communication, Dong Seoul University, Korea.

*Retrieval Number: C10810193S20/2020©BEIESP*
*DOI: 10.35940/ijitee.C1081.0193S20*

383

*Published By:*
*Blue Eyes Intelligence Engineering*
*& Sciences Publication*