# Prediction of Breast Cancer Using Mammogram Images by Data mining on the Extracted Features of the processed Image

Tarun Gangil[1] Sneha Prajwal[2]

Research Scholar, Department of Medical Electronics, Dr. Ambedkar Institute of Technology, Bangalore, Karnataka,

India

Research Scholar, Department of Computer Science and Engineering, Bapuji Institute of Engineering and Technology ,Davangere, Karnataka, India

**ABSTRACT:** The paper presents the algorithm to detect the tumour from the mammogram images. Two major problems are handled in this paper. One is how to detect tumours as suspicious regions with a very weak contrast to their background and another is how to extract features which categorize tumours. The tumour detection method follows the scheme of (a) mammogram enhancement. (b) The segmentation of the tumour area. (c) The extraction of features from the segmented tumour area and followed by implementing clustering techniques on the extracted features such as PAM  CLARA and BIRCH. First Stage we will perform the image enhancement which is required to improve the quality of the image such that the noise can be reduced from the image and the area of interest where the maximum probability of finding the tumour can be enhanced. For performing such enhancement we will use the techniques such as filtering, top hat operation, DWT. Contrast stretching is used to enhance the properties of the image. The most common segmentation method used is thresholding, K means and Otsu. The features are extracted from the segmented breast area. Next stage includes the implementation of data mining algorithms from the extracted features and comparing its results. The method was tested on 100 mammographic images, from the mini-MIAS database. The algorithm was developed using a licenced version of MATLAB version R2015a. The Aim is to detect the tumour in the mammogram images, extract its spatial features and cluster them using data mining techniques.

**KEYWORDS**: PAM, CLARA, BIRCH, kernel function, separating hyper plane, mammography, contrast stretching, segmentation, image enhancement, discrete wavelet transform.

## I. INTRODUCTION

Breast Cancer is the most widely recognized growth in women, however it can likewise show up in men also [1]. India is likely to have over 17.3 lakh new cases of cancer and over 8.8 lakh deaths due to the disease by 2020 with cancers of breast, lung and cervix topping the list. In its projection, the Indian Council of Medical Research (ICMR) said in 2016 the total number of new cancer cases is expected to be around 14.5 lakh and the figure is likely to reach nearly 17.3 lakh new cases in 2020. Over 7.36 lakh people are expected to succumb to the disease in 2016 while the figure is estimated to shoot up to 8.8 lakh by 2020. Data also revealed that only 12.5 per cent of patients come for treatment in early stages of the disease. Among females, breast cancer topped the list and among males mouth cancer. The northeast reported the highest number of cancer cases in both males and females [2]. Breast Cancer is brought about by the improvement of malignant cells in the breasts. The malignant cells originate in the lining of the milk glands or ducts of the breast (ductal epithelium).

Tumour is uncontrolled growth of cells which can be either benign or malignant. Breast Tumour is a tumour present in Breast. There are 3 types in breast cancer such as malignant tumour, benign tumours and non-malignant tumour. Malignant tumours are cancerous and can invade and destroy nearby tissue, it spreads to other parts of the body. Benign tumours are not cancerous and may grow larger but do not spread to other parts of the body. Tumour can be

easily identified in mammogram because tumour part is highly bright (having high intensity) compared to other part (background) of the mammogram image.

Breast cancer is the most common non skin malignancy in women and the second leading cause of female cancer mortality. Breast tumours and masses usually appear in the form of dense regions in mammograms. A typical benign mass has a round, smooth and well circumscribed boundary; on the other hand, a malignant tumour usually has a speculated, rough, and blurry boundary.

**MAMMOGRAPHY**

Mammography is a uniquely important type of medical imaging used to screen for breast cancer. All women at risk go through mammography screening procedures for early detection and diagnosis of tumour. A typical mammogram is an intensity X-ray image with gray levels showing levels of contrast inside the breast that which characterize normal tissue and different calcifications and masses. Mammography is believed to reduce mortality from breast cancer. No other imaging technique has been shown to reduce risk, but remaining aware of breast changes and physician examination are considered essential parts of regular breast care.

Detection and diagnosis of breast cancer in its early stage increases the chances for successful treatment and complete recovery of the patient. Micro calcifications and masses are two most common types of suspected abnormalities in mammogram images. Screening mammography is currently the best available radiological technique for early detection of breast cancer. A mammogram is an x-ray of the breast tissue which is designed to identify abnormalities which could be benign or malignant tumours. Based on the level of suspicion of the abnormality, radiologists usually recommend a routine follow up to confirm the result. The earlier breast cancer cannot be detected by the physical examination, because the tumour is too small to be picked up from these reasons the study of the automatic image processing of X-ray mammography is strongly desired.

Important visual clues of breast cancer include preliminary signs of masses and calcification clusters. A mass is a localized collection of tissue seen in two different projections, and calcifications are small calcium deposits. If a potential mass is seen in only a single projection it should be called 'Asymmetry' or 'Asymmetric Density' until its three-dimensionality is confirmed. Masses have different density (fat containing masses, low density, and high density), different margins (circumscribed, obscured, indistinct, and speculated) and different shape (round, oval, lobular, irregular). Round and oval shaped masses with smooth and circumscribed margins usually indicate benign changes. On the other hand, a malignant mass usually has a speculated, rough and blurry boundary. Unusually smaller and clustered calcifications are associated with malignancy while there are other calcifications (diffuse, regional, segmental and linear) that are typically benign. Such calcifications are termed as **Microcalcifications (MC)**. In the early stages of breast cancer, these signs are subtle and hence make diagnosis by visual inspection difficult. Figure 1 shows a Mammogram image containing mass.



**Figure 1: Mammogram showing a mass.**

Typically a Radiologist evaluates the X-ray of the breast to determine the presence or absence of masses or tumours. These growths are generated by radical changes on the cell structure, due to the presence of cancer. After the evaluation, there are three possible outcomes for the mammogram:

- A normal one indicates no presence of high density bodies within the breast,
- A "benign" mammogram suggests the presence of a non-cancerous tumour
- A malign one adverts the presence of a cancerous tumour on the patient's breast.

In this paper, we propose a new classification approach for breast cancer using statistical features and an event processing engine. By means of a set of rules we classify and then isolate significant occurrences, correlate and enhance them, to continuously improve the classification rate. We define a group of queries to map these rules into the model and classify a collection of 200 images from a commonly used database of mammograms. Our results display considerable classification rate accuracy and an inexpensive and easy implementation, which can be integrated with other systems.

## WORKING OF MAMMOGRAPHY

The basic principle behind mammography is the same that all X-ray imaging is based upon. X-ray radiations, or roentgen radiations, are produced when an electron beam is focused on a special type of material usually called the target material. Tungsten and molybdenum are some of the most common examples of target materials. However, in mammography, because higher resolution requirements, molybdenum is typically used. A spectrum of roentgen radiations emerges from the target and is directed through breast. On the opposite side of the breast, a special _lm is placed that is sensitive to roentgen radiations. Once inside the breast, each ray of X-ray photons is attenuated by varying degrees dependent on the variation in density of the tissue within it.

The amount of attenuated radiations due to the higher-density tissue will be considerably higher than the radiations attenuated by the surrounding fatty tissue. Fibrous tissue and masses appear lighter in an X-ray image. Figure 2, shows a typical mammographic setup and shows a typical mammographic image. Figure 3 shows mammogram with large portions of dense fibrous tissue. This dense tissue is the region which is having highest probability of occurrence of cancerous cells.
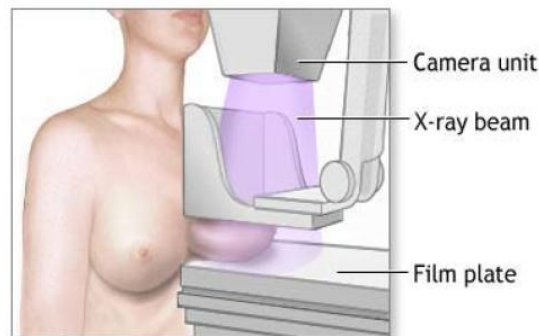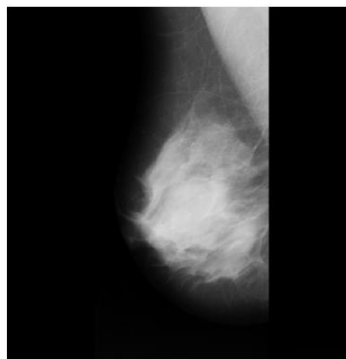


Figure. 2 Mammography Setup



Figure 3: Mammogram with large portions of dense fibrous tissue

To ensure the detection of cancerous masses, radiologists require more than one mammographic image of a patient's breast. For the same reason, it is also important that different angles are used to achieve different views of a patient's breast. Medio-lateral oblique (MLO) view and cranio-caudal (CC) view, used in combination, have become the international standard views in mammographic screening.

## II. RELATED WORK

In the year 2013 P.Spandana worked on Novel Image Processing Techniques for Early Detection of Breast Cancer, Matlab and Lab View implementation. Here the Image enhancement is done using wavelets and adaptive histogram equalization technique. Segmentation of masses is done using region growing technique and Extraction of border of the mass using canny edge detection and morphological operations. [3]. In the year 2014 Mario Mustra worked on Detection of Areas Containing Micro calcifications in Digital Mammograms. Here results of detection of suspicious areas which could contain micro calcifications are shown to be good in comparison to other methods which do not use manual selection of suspicious areas. For background suppression a combination of wavelet filtering and grayscale morphology is used. [4]. In the year 2013 Y.Ireaneus Anna Rejani, worked on Early detection of Breast Cancer using SVM Techniques. Here the Image is segmented and the features are obtained from the region of interest. The extracted features were classified using Support vector machine (SVM) for the detection of image to be cancerous or non-cancerous.[5]. In the year 2015, Yousef Nasiri worked on Breast cancer detection in mammograms using wavelets and counterlet transforms. Here SVM Technique is used for classifying the extracted features from the mammogram images and the feature vectors are created using wavelet and counterlet transform on the mammogram images. [6]. In the year 2014, Nadia El Atlas worked on Computer Aided Breast cancer detection using Mammograms. Here the survey of deaths due to breast cancer is represented and an overview of digital image processing and pattern analysis techniques to address several areas in CAD of breast cancer, including the four stages of CAD system: image preprocessing, image segmentation, features extraction and selection and image classification.[7]. In the year 2013, Sharanya Padmanabhan worked on Enhanced accuracy of breast cancer detection in digital mammograms using wavelet analysis. the process of object detection, recognition and classification of mammograms with the aim of differentiating between normal and abnormal (benign or cancerous) cells. Here the MATLAB numerical analysis software is used for processing the image to extract out features for the cancerous images.[8].

## III. PROPOSED ALGORITHM

Detection of tumours in mammogram is divided into three main stages. The first step involves an enhancement procedure, image enhancement techniques are used to improve an image, where to increase the signal to noise ratio and to make certain features easier to see by modifying the colors or intensities. Then the intensity adjustment is an image's intensity values to a new range. After the mammogram enhancement segment the tumour area. Then the features are extracted from the segmented mammogram. Figure 4 and 5 represents the proposed algorithm.

A) *Image enhancement:*

Image enhancement can be defined as conversion of the image quality to a better and more understandable level. The enhancement procedure is (a) the mammogram images are filtered by Gaussian smoothing filter based on standard deviation. (b) Perform morphological top hat filtering on the gray scale input image using the structuring element. (c) The top hat output is decomposed into two scales using discrete wavelet transform and then the image is reconstructed as shown in figure 4.
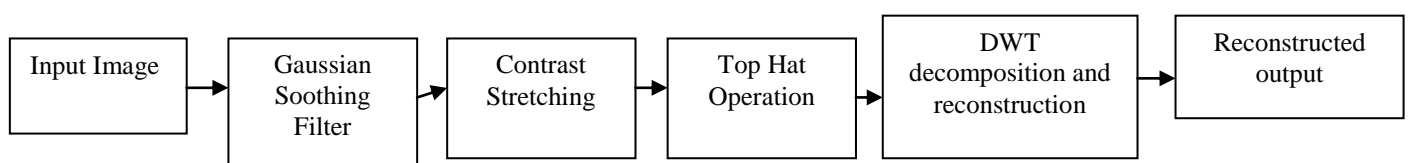
Input Image → Gaussian Soothing Filter → Contrast Stretching → Top Hat Operation → DWT decomposition and reconstruction → Reconstructed output

Figure 4: Block diagram of image enhancement

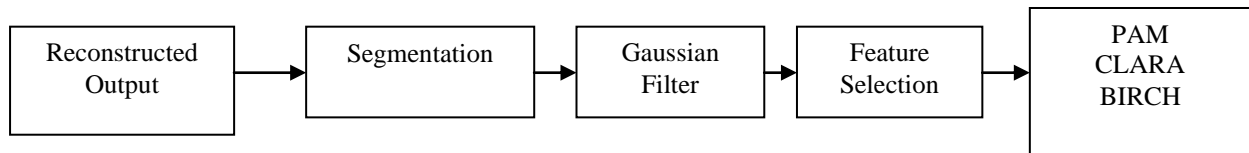| Reconstructed Output | → | Segmentation | → | Gaussian Filter | → | Feature Selection | → | PAM CLARA BIRCH |
|---|---|---|---|---|---|---|---|---|

Figure 5: Block diagram of the proposed method.

*B)   Noisy Image & Denoising*

Image noise is random (not present in the object imaged) variation of brightness or color information in images, and is usually an aspect of electronic noise. It can be produced by the sensor and circuitry of a scanner or digital camera. Image noise can also originate in film grain and in the unavoidable shot noise of an ideal photon detector. Image noise is an undesirable by-product of image capture that adds spurious and extraneous information. The original meaning of "noise" was and remains "unwanted sound"; unwanted electrical fluctuations in signals received by AM radios caused audible acoustic noise ("static"). By analogy unwanted electrical fluctuations themselves came to be known as "noise". Image noise is, of course, inaudible.

Images taken with both digital cameras and conventional film cameras will pick up noise from a variety of sources. Many further uses of these images require that the noise will be (partially) removed - for aesthetic purposes as in artistic work or marketing, or for practical purposes such as computer vision. In salt and pepper noise (sparse light and dark disturbances), pixels in the image are very different in color or intensity from their surrounding pixels; the defining characteristic is that the value of a noisy pixel bears no relation to the color of surrounding pixels. Generally this type of noise will only affect a small number of image pixels. When viewed, the image contains dark and white dots, hence the term salt and pepper noise. Typical sources include flecks of dust inside the camera and overheated or faulty CCD elements.

In Gaussian noise, each pixel in the image will be changed from its original value by a (usually) small amount. A histogram, a plot of the amount of distortion of a pixel value against the frequency with which it occurs, shows a normal distribution of noise. While other distributions are possible, the Gaussian (normal) distribution is usually a good model, due to the central limit theorem that says that the sum of different noises tends to approach a Gaussian distribution. In either case, the noise at different pixels can be either correlated or uncorrelated; in many cases, noise values at different pixels are modelled as being independent and identically distributed, and hence uncorrelated.
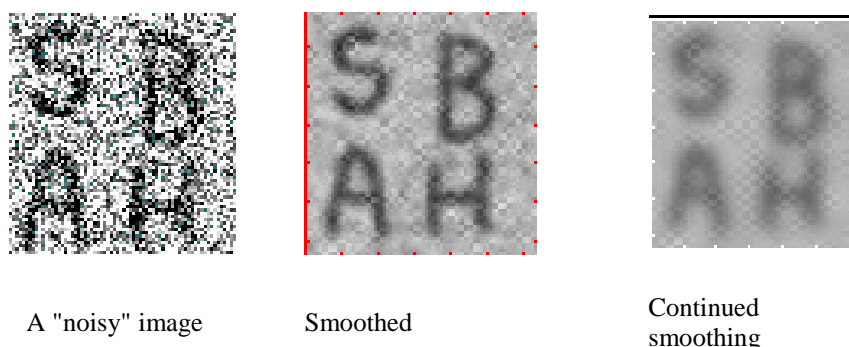


A "noisy" image          Smoothed          Continued smoothing

Figure 6: Image Smoothening

One goal in image restoration is to remove the noise from the image in such a way that the "original" image is discernible. Of course, "noise" is in the eye of the beholder; removing the "noise" from a Jackson Pollack painting would considerably reduce its value. Nonetheless, one approach is to decide that features that exist on a very small scale in the image are noise, and that removing these while maintaining larger features might help "clean things up".

One well-traveled approach is to smooth the image. The simplest such version is replace each pixel it by the average of the neighboring pixel values. If we do this a few times we get the image in the middle above; if we do it many times, we get the image on the right. On the plus side, much of the spotty noise has been muted out. On the downside, the sharp boundaries that make up the letters have been smeared due to the averaging. While many more sophisticated approaches exist, the goal is the same: to remove the noise, and keep the real image sharp. The trick is to not do too much, and to "know when to stop".

## C) Image Restoration

Image restoration is the operation of taking a corrupted/noisy image and estimating the clean original image. Corruption may come in many forms such as motion blur, noise, and camera misfocus. Image restoration is different from image enhancement in that the latter is designed to emphasize features of the image that make the image more pleasing to the observer, but not necessarily to produce realistic data from a scientific point of view. Image enhancement techniques (like contrast stretching or de-blurring by a nearest neighbor procedure) provided by "Imaging packages" use no a priori model of the process that created the image.

With image enhancement noise can be effectively be removed by sacrificing some resolution, but this is not acceptable in many applications. In a Fluorescence Microscope resolution in the z-direction is bad as it is. More advanced image processing techniques must be applied to recover the object. Deconvolution is an example of image restoration method. It is capable of: Increasing resolution, especially in the axial direction Removing noise Increasing contrast.

The purpose of image restoration is to "compensate for" or "undo" defects which degrade an image. Degradation comes in many forms such as motion blur, noise, and camera misfocus. In cases like motion blur, it is possible to come up with an very good estimate of the actual blurring function and "undo" the blur to restore the original image. In cases where the image is corrupted by noise, the best we may hope to do is to compensate for the degradation it caused. In this project, we will introduce and implement several of the methods used in the image processing world to restore images.

## D) Top-Hat Transform

In mathematical morphology and digital image processing, **top-hat transform** is an operation that extracts small elements and details from given images. There exist two types of top-hat transform: The *white top-hat transform* is defined as the difference between the input image and its opening by some structuring element; The *black top-hat transform* is defined dually as the difference between the closing and the input image. Top-hat transforms are used for various image processing tasks, such as feature extraction, background equalization, image enhancement, and others. The **Top-Hat Transform** or peak detector is another composite operation: the image opened by the structuring element is subtracted from the original image. The brightest spots on the original image are highlighted using this transformation.

The white top-hat transform returns an image, containing those "objects" or "elements" of an input image that:

- Are "smaller" than the structuring element (i.e., places where the structuring element does not fit in), and
- Are brighter than their surroundings.

The black top-hat returns an image, containing the "objects" or "elements" that:

- Are "smaller" than the structuring element, and
- Are darker than their surroundings.

The size, or width, of the elements that are extracted by the top-hat transforms can be controlled by the choice of the structuring element. The bigger the latter, the larger the elements extracted. Both top-hat transforms are images that contain only non-negative values at all pixels. On the left side there is the original image, on the right side is the image after top-hat transform as shown in figure 7.
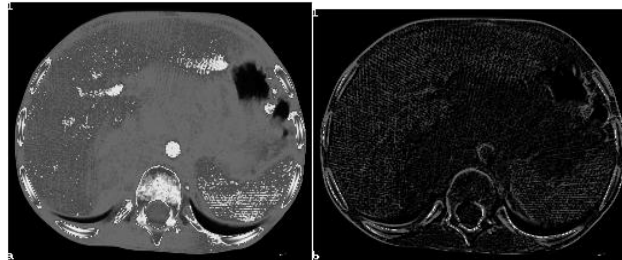
Figure 7: Example of Top Hat Transform using MRI Images

*E) ROI Extraction Using Segmentation*

The main logical step in the system is to perform image segmentation, this involves separating the suspected areas they may contain abnormalities from the image. The suspicious area is an area that is brighter than its surroundings, has almost uniform density, has a regular shape with varying size, and has fuzzy boundaries. A number of different segmentation techniques were investigated and are outlined in this section. A *Region Of Interest* (ROI) is a portion of an image that we want to filter or perform some other operation on. ROI defined by creating a *binary mask*, which is a binary image that is the same size as the image we want to process with pixels that define the ROI set to 1 and all other pixels set to 0. We can define more than one ROI in an image. The regions can be geographic in nature, such as polygons that encompass contiguous pixels, or they can be defined by a range of intensities. In the latter case, the pixels are not necessarily contiguous. In this stage the potential tumour locations are extracted from the segmented image obtained in the previous stage.

In our project ROI means mass present in the Mammogram image. ROI extraction is the extraction of mass from the input Mammogram Image. By calculating the threshold intensity value of the image, ROI is extracted, because tumour part usually has higher brightness than other part of the Mammogram. In the figure 8, a Region of Interest is defined at near top left of the image. This is useful, for example when we want to crop an object from an image, or when we want to perform template matching within sub image.
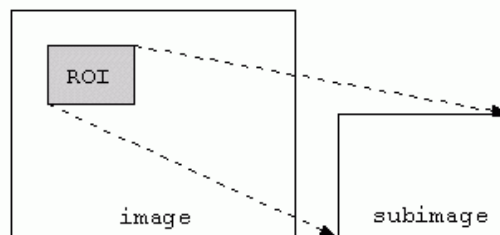


**Figure 8. An image with Region of Interest defined.**

For ROI extraction segmentation technique is used. Segmentation refers to the process of partitioning a digital image into multiple segments (sets of pixels, also known as super pixels). The goal of segmentation is to simplify and/or change the representation of an image into something that is more meaningful and easier to analyze. Image segmentation is typically used to locate objects and boundaries (lines, curves, etc.) in images. More precisely, image segmentation is the process of assigning a label to every pixel in an image such that pixels with the same label share certain visual characteristics.

The result of image segmentation is a set of segments that collectively cover the entire image, or a set of contours extracted from the image. Each of the pixels in a region are similar with respect to some characteristic or computed property, such as color, intensity, or texture. Adjacent regions are significantly different with respect to the same characteristic(s).

In analyzing mammogram image, it is important to distinguish the suspicious region from its surroundings. The methods used to separate the region of interest from the background are usually referred as the segmentation process. Different segmentation methods are available. The block diagram for segmentation is shown in figure 9.
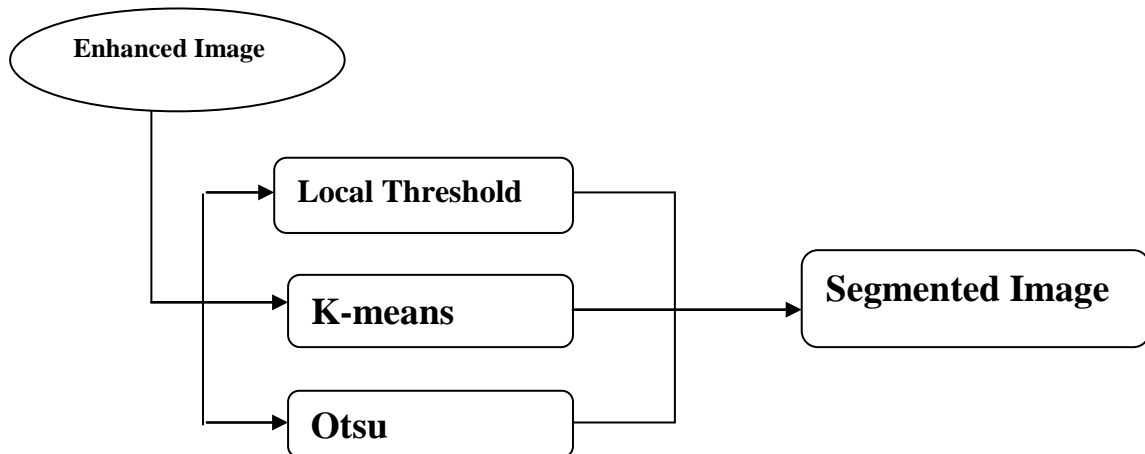
Figure 9: Segmentation Block Diagram

### Local Thresholding

This technique has been proven to provide an easy and convenient way to perform the segmentation on digital mammogram. The segmentation is determined by a single value known as the intensity threshold value. Then, each pixel in the image is compared with the threshold value. Pixel intensity values higher than the threshold will result in a white spot in the output image.

### K-means Segmentation

K-means segmentation is a method of cluster analysis which aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean. This is a region clustering method that does not need prior information or starts point and is based on an iterative process. This method only requires a stop function, which is the number of clusters, k, in the segmented image. The higher the k value, the clearer the segmentation but the processing time will increase. The algorithm iterates over two steps:

1. Compute the mean of each cluster.

2. Compute the distance of each point from each cluster by computing its distance from the corresponding cluster mean. Assign each point to the cluster it is nearest to.

### Otsu Segmentation

Otsu's segmentation method is used to automatically perform histogram shape-based image thresholding, or, the reduction of a gray level image to a binary image. The algorithm assumes that the image to be threshold contains two classes of pixels (e.g. foreground and background) then calculates the optimum threshold separating those two classes so that their combined spread (intra-class variance) is minimal.

Otsu segmentation method has shown a more satisfactory performance in the medical image segmentation. It has been found to perform well compared to other thresholding methods in segmenting the masses in digital mammogram. The output of a segmentation process is a binary image. In order to retrieve the texture information, the segmented image is masked with a 16-bit quantization image. Instead of using the original image, a quantized image is used. In the quantized image, the amount of represented intensities is visible to humans. By reducing quantization level to 16 bits, the area with masses can be identified on the mammogram. The masked image is then used as input for the features extraction process. In our project, we have used this segmentation technique.

### F) Wavelets

Wavelets are a relatively new mathematical tool which has contributed significantly to image and signal analysis over the past twenty years. A wavelet can be defined as a mathematical function used to divide a given function into

different scale components and each scale component can then be studied with a resolution that matches its scale. A wavelet transform is then a representation of a function by wavelets. There are two types of wavelet transform: the continuous wavelet transform (CWT) and the discrete wavelet transform (DWT). The continuous wavelet transform of a function *f* using a wavelet function basis is defined by equation 1:

$$f(a, b) = \int f(x) \psi x_{a,b}(x) dx$$

............................... 1)

Where $\psi(x)$ is the mother wavelet function. The basis of the wavelet function is obtained by scaling and shifting a signal mother wavelet function given by equation 2.

$$\psi_{a,b}(x) = \frac{1}{\sqrt{a}} \psi\left(\frac{x - b}{a}\right) \; ; \; a > 0$$

........................................ 2)

Where *a* is a scale factor and *b* is the shift value. The DWT is obtained by taking *a = 2* and *b =Z*. In both cases the transform, transforms the function into a function of scale and translation. While the Fourier transform uses sinusoids of infinite duration to decompose the wavelet transform uses wavelets of finite duration to perform the same operation. When it comes to dealing with images the discrete wavelet transform is the transform chosen, the reasons are outlined in the coming section. The type of wavelet mother function chosen is that of the Haar wavelet, its mother function $\psi(t)$ can be described as follows, equation 3 and scaling function given by equation 4:

$$\psi(t) = \begin{cases} 1 & 0 \le t < \frac{1}{2} \\ -1 & \frac{1}{2} \le t < 1 \\ 0 & otherwise. \end{cases}$$

.................................................. 3)

Its scaling function $\varphi(t)$ can be described as

$$\varphi = \begin{cases} 1 & 0 \le t < 1 \\ 0 & otherwise. \end{cases}$$

.................................................. 4)

**Discrete Wavelet Transform (DWT)**

The discrete wavelet transform is used for image processing as it gives a detailed insight to an image's spatial and frequency characteristics, unlike the Fourier transform or CWT which deal only with an image's frequency characteristics. The first stage of image de-noising is image decomposition. The DWT of an image is most efficiently calculated by passing it through a series of filters, called a filter bank, as shown in Figure 3.5. In this figure, the image is represented by x[n], the low pass filter is represented by G[n], the high pass filter is represented by H[n] and the down sampling operator represented by ↓. At each level, the high pass filter produces detail coefficients for that level while the low pass filter produces approximation coefficients which are fed into the next level, and the process continues for as many levels as are required. The approximation coefficients are the high scale, low frequency components of the image whereas the detail coefficients (vertical, horizontal and diagonal) are the low scale, high frequency components. Image decomposition of the noisy mammogram can be seen in Figure 10.

*G) Thresholding*

The next stage in de-nosing the image is to threshold the coefficients. This involves selecting and applying a threshold to the detail coefficients for each level from 1 to N. This can be achieved using hard thresholding, setting to zero all the elements whose absolute values are lower than the threshold or by soft thresholding which first sets to zero all the elements whose absolute values are less than the threshold, and then shrinking the non-zero coefficients towards zero. Soft thresholding is used in this project, although no visual difference was observed when hard thresholding was used.
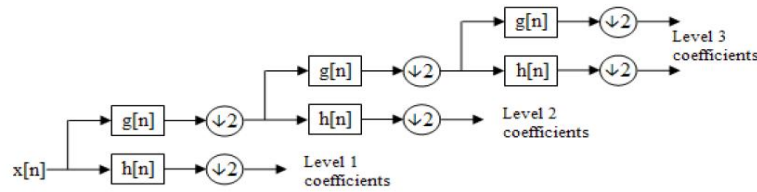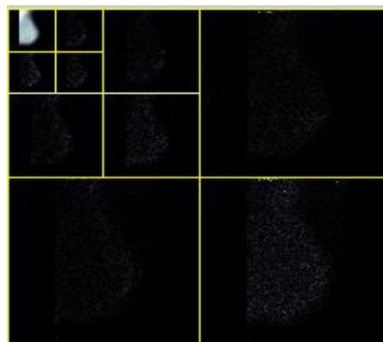
Figure - Three level filter bank



Figure - Three level decomposition of mammogram

Figure 10: Discrete Wavelet Transform

### H) *Wavelet Reconstruction*

The final step is to perform wavelet reconstruction using the original approximation coefficients of level N and the modified detail coefficients of levels 1-N. This is achieved by using the inverse discrete wavelet transform (IDWT). Where wavelet decomposition uses down sampling and filtering the reconstruction process consists of up-sampling and filtering. The noisy image is shown alongside the final de-noised image in Figure.

### I) *Feature Extraction*

In pattern recognition, image classification, and Bio-medical applications and in image processing, feature extraction is a special form of dimensionality reduction. When the input data to an algorithm is too large to be processed and it is suspected to be notoriously redundant (e.g. the same measurement in both feet and meters) then the input data will be transformed into a reduced representation set of features (also named features vector). Transforming the input data into the set of features is called feature extraction. If the features extracted are carefully chosen it is expected that the features set will extract the relevant information from the input data in order to perform the desired task using this reduced representation instead of the full size input.

Feature extraction involves simplifying the amount of resources required to describe a large set of data accurately. When performing analysis of complex data one of the major problems stems from the number of variables involved. Analysis with a large number of variables generally requires a large amount of memory and computation power or a classification algorithm which over fits the training sample and generalizes poorly to new samples. Feature extraction is a general term for methods of constructing combinations of the variables to get around these problems while still describing the data with sufficient accuracy. Many data analysis software packages provide for feature extraction and dimension reduction. Common numerical programming environments such as MATLAB, SciLab, NumPy and the R language provide some of the simpler feature extraction techniques (e.g. principal component analysis) via built-in commands. More specific algorithms are often available as publicly-available scripts or third-party add-ons.

Feature is used to denote a piece of information which is relevant for solving the computational task related to a certain application. More specifically, features can refer to

- The result of a general neighbourhood operation (feature extractor or feature detector) applied to the image,
- Specific structures in the image itself, ranging from simple structures such as points or edges to more complex structures such as objects.

Many features has been extracted for the abnormalities of mammograms. The extraction methods of texture feature play very important role in detecting abnormalities of mammograms because of the nature of mammograms. Texture features have been proven to be useful in differentiating masses and normal breast tissues. Texture features are able to

isolate normal and abnormal lesion with masses and micro calcification. Feature extraction block diagram shown in figure 11.
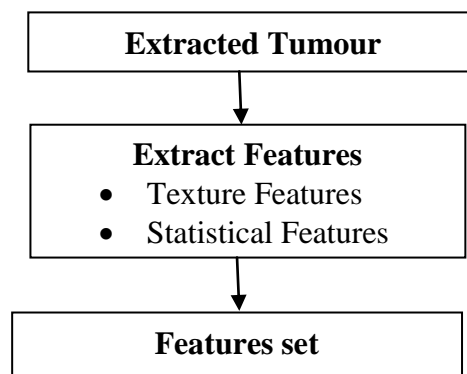


**Figure 11: Extraction of Features**

**Texture and Statistical Features**

The features are extracted for each ROI. The extracted features are then feed to SVM classifier to train it for the recognition of a particular ROI of similar nature. These features are 1. Area, 2. Centriod, 3.major axis length, 4.minor axis length, 5.eccentricity, 6.orientation, 7.filled area, 8.extrema, 9.solidity, 10.equivdiameter. The area is the scalar value; it computes the actual number of pixels in the region. Then the centroid is the vector and it computes the centre of the tumour region.

*J) Data mining Algorithms*

#### *PAM METHOD*

The PAM(Partition around medoids) algorithm was developed by Leonard Kaufman and Peter J. Rousseeuw, which is similar to k-means, where the dataset is broken into number of groups called clusters. A cluster is group of similar items. PAM –a Partitional clustering technique it create a one-level partitioning of the data points .This algorithm takes 2 inputs, dataset and K as number of clusters to be formed. It partitions the dataset into k clusters consists of n objects. It finds a sequence of objects called medoids that are centrally located in clusters. An actual object mediod is center which is accessible to all data objects within a group (cluster). The goal is to find a set of k clusters that minimizes the sum of the dissimilarities of all the objects to their nearest medoids.[9].
Steps:
1. Initially it takes the dataset and k as input and it randomly chooses k data items as the initial medoids.
2. Assigns each nearest mediod to the cluster.
3. Randomly select a non-medoid data item X. For the medoid M data item, it computes the total swapping cost S.
4. If the total cost swapping Cost S is less than zero(S<0), then perform the swap operation to generate the new set of k-medoids. i.e it swaps the non mediod with old mediod.
5. Repeat steps 2, 3 and 4 till the medoids stabilize their locations, actual mediod is selected which is accessible to all data points in the clusters.

#### *CLARA METHOD*

The CLARA (Clustering for Large Applications) is a clustering algorithm based on Randomized Search method was developed by Leonard Kaufman and Peter J. Rousseeuw. It takes k as number of clusters to be formed and dataset as an input.

It handles the large dataset where as PAM handles small dataset.

Steps:

1. It randomly selects data point from the dataset. i.e It draws multiple samples of the data set, applies PAM on each sample, and gives the best clustering as the output.
2. PAM generates a set of Mediods. The quality of mediods is measured by taking the average dissimilarity between every object in the entire data set D and the medoid of its cluster given by equation 5.

$$Cost(M, D) = \frac{\sum_{i=1}^{n} dissimilarity(O_i, rep(M, O_i))}{n}$$ ……………  5)

Where M is a set of selected medoids, dissimilarity(Oi, Oj) is the dissimilarity between objects Oi and Oj, and rep(M, Oi) returns a medoid in M which is closest to Oi.

3. Calculate the mean (or the sum) of the dissimilarities of the observations to their closest medoid. This is used as a measure of the goodness of the clustering.Retain the sub-dataset for which the mean (or sum) is minimal.[10].

### ♣ *BIRCH METHOD*

Balanced Iterative Reducing and Clustering Using Hierarchies, it works on large dataset, it's a unsupervised data mining algorithm. BIRCH can typically find a good clustering with a single scan of the dataset. Dataset consists of multidimensional data points with noise. The main advantage of birch is, it removes the outliers efficiently. It builds a tree called the Characteristic Feature Tree (CFT) for the given data and generates CF nodes. The CF nodes consists of CF subclusters.CF sub clusters consists of necessary information for clustering. For a Cluster X with N-data points($X_i$), Clustering Feature(CF) is a triple summarizing information such as (N,LS,SS). Where N is Number of data points, LS is Linear Sum of Data points and SS is Square sum of Data points in that cluster. This information is sufficient to compute the distance between two clusters.CF tree takes two parameters the threshold distance TD and the branching factor BF. Branching Factor includes B and L where B is Maximum number of CF in non-leaf Node and L is Maximum of CF in Leaf Node. TD is maximum radius of CF. CF tree built dynamically as data is scanned and inserted. The tree size depends on the parameter TD. There are 3 phases in BIRCH

**Phase 1:** Load the data into the memory by building CF Tree
Input: Dataset D with data points and Threshold Distance TD into tree.
Output: Initial CF Tree is constructed.
Dataset fits into memory, tree is constructed by removing outliers. If the dataset is more and doesn't fit into memory increase the value of TD so that it fits the memory.

**Phase 2:** Global Clustering
Input: Initial CF tree/Smaller CF tree
Output: good cluster
In this step initial CF tree is loaded and algorithm scans the leaf nodes of initial CF tree to generate smaller CF tree. An existing algorithm i.e standard algorithm (global algorithm) Example agglomerative cluster/Kmeans method are applied on the smaller CF tree to cluster all leaf entries. In this step we can specify the number of cluster required or the diameter of the cluster. global clustering step labels these subclusters into global clusters (labels) The root of Sub cluster that has smallest radius are merged with root node, tis done by obtaining the values of threshold distance and Branching factor. Then this is done repeatedly till it reaches a leaf. In this step cluster quality is improved and good clusters are formed.

**Phase 3:** Cluster refining
Cluster are refined by scanning the entire dataset once again to label the data points. i.e. assign the data points to the cluster. Centroid for the clusters are assigned and redistribution of data points closest to centroid are obtained. This is repeated to obtain new set of clusters. Quality cluster are obtained by removing outliers.[11].

## IV. SIMULATION RESULTS

Figure 12 shows the output GUI for the processing of Mammogram images using MATLAB platform. The input image is set to Gaussian smoothing, contrast stretching, top hat operation, DWT and the image will be reconstructed. The reconstructed image is segmented to obtain the desired region of interest and various geometrical features are extracted out of it. These features are being classified using three data mining algorithms : PAM, CLARA and BIRCH. The results of these data mining algorithms are compared and graphical results are being shown. Finally if the image is classified to be as the cancerous image then it is further classified to predict, whether the image is of type Malignant or Benign. The samples of extracted features are used as an input for data mining algorithms. Here at this stage, the features are used to predict for the occurrence of image which contains cancerous lesions or not. If the cancerous image is predicted, then algorithm is trained to predict for the stage of cancer, i.e. Malignant or Benign. The comparison of the output of the three data mining algorithms is produced graphically, with respect to the comparison in its percentage of accuracy, sensitivity and specificity.
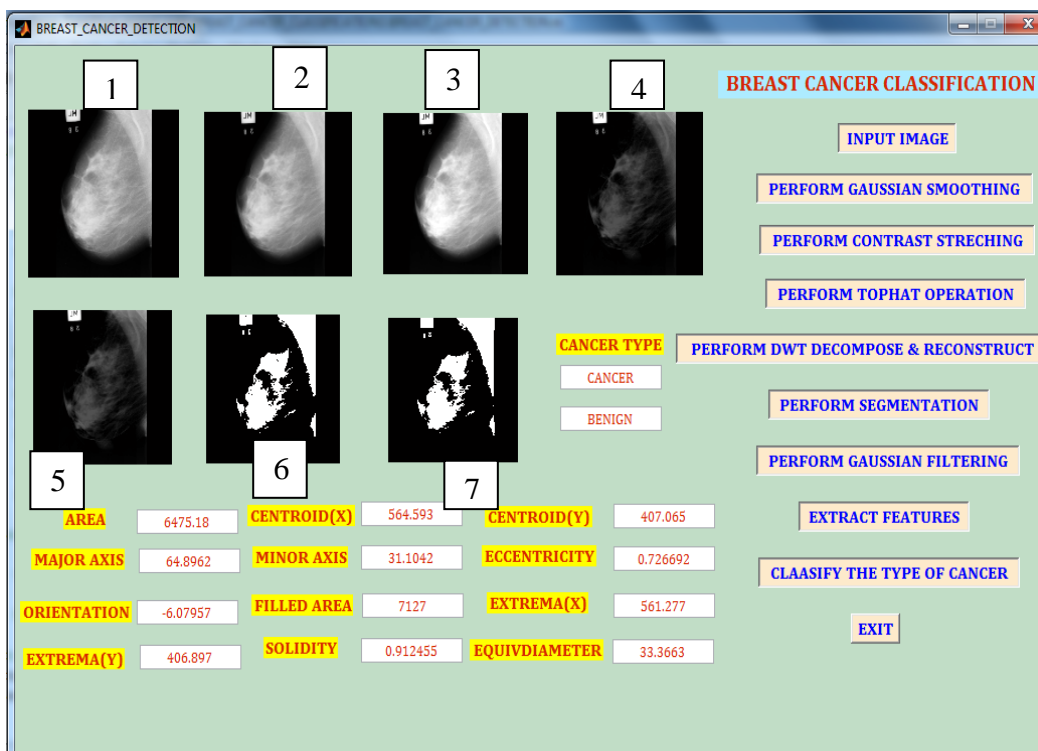


Figure 12: Output screenshot of image processing and Feature extraction of Mammogram images using MATLAB
1) Input Image 2) Gaussian Smoothing 3) Contrast Stretching 4) Top Hat Operation 5) DWT Decomposition and reconstruction 6) Segmentation 7) Gaussian Filtering 8) Feature Extraction 9) Classify as Cancerous (Benign or Malignant) or non cancerous

*Data Mining Results*

A total sample size of 200 is taken for the analysis. 100 samples are taken for training the algorithm and 100 samples were taken for testing the algorithm. Where each sample represents the extracted features out of the Region of interest of the Mammogram Image, where highest possibility of cancerous tissue is suspected to be present.

**1) Results of PAM:**

Here the confusion matrix is shown by testing with the features of sample size 100 as sown in Table 2. The percentage of accuracy is given by the equation 6 and accuracy is shown in Table 1.

$$\% \ accuracy = \frac{Number \ of \ samples \ correctly \ classified}{Total \ number \ of \ samples} X \ 100$$

..........................................................................6)

Table 1: Results of PAM algorithm

| Number of Cancerous Samples | Number of Non Cancerous Samples | Correctly Classified | Incorrectly Classified |
|---|---|---|---|
| 64 | 36 | 63 | 37 |

**Accuracy %= 63%**

Table 1 represents the number of cancerous and non classified cancerous samples taken for the testing of PAM algorithm. By using the accuracy percentage formula in equation number 6 gives the percentage of correctly samples using this algorithm. Hence the calculated percentage of accuracy is 63%, with respect to the number of samples correctly classified.

Table 2: Confusion Matrix from PAM Algorithm

| | Positive | Negative |
|---|---|---|
| Cancerous (True) | 53 | 11 |
| Non Cancerous (False) | 11 | 25 |

Table 2 represents the confusion matrix, according to the classification performed. True positive is the case where the given sample is from the cancerous image, correctly classified. True Negative is the case where the given sample is from the non cancerous image, correctly classified. False positive is the case where the classifier fails to classify it as a cancerous image, whereas it is a cancerous image only. False negative is the case where the classifier fails to classify it as a non cancerous image, whereas it is a non cancerous image only.

The specificity percentage and sensitivity percentage are calculated from the above confusion matrix using equation 7 and 8 respectively.

$$\% \ Specificity = \frac{True \ Negative}{False \ Positive + True \ Negative} X \ 100$$

.............................................................................7)

$$\% \ Sensitivity = \frac{True \ Positive}{True \ Positive + False \ Negative} X 100$$

.............................................................. 8)

Calculated values of Sensitivity and Specificity for PAM algorithm are as follows using equation 8 and 7 respectively:

Sensitivity= 67.94%

Specificity= 17.18 %

### 2) Results of CLARA:

Here the confusion matrix is shown by testing with the features of sample size 100 as sown in Table 4. The accuracy is given by the Table 3.

Table 3: Results of CLARA algorithm

| Number of Cancerous Samples | Number of Non Cancerous Samples | Correctly Classified | Incorrectly Classified |
|---|---|---|---|
| 72 | 28 | 85 | 15 |

**Accuracy %= 85%**

Table 3 represents the number of cancerous and non classified cancerous samples taken for the testing of CLARA algorithm. By using the accuracy percentage formula in equation number 6 gives the percentage of correctly samples using this algorithm. Hence the calculated percentage of accuracy is 85%, with respect to the number of samples correctly classified.

Table 4: Confusion Matrix from PAM Algorithm

|  | Positive | Negative |
|---|---|---|
| Cancerous (True) | 53 | 8 |
| Non Cancerous (False) | 19 | 20 |

Table 4 represents the confusion matrix, according to the classification performed. True positive is the case where the given sample is from the cancerous image, correctly classified. True Negative is the case where the given sample is from the non cancerous image, correctly classified. False positive is the case where the classifier fails to classify it as a cancerous image, whereas it is a cancerous image only. False negative is the case where the classifier fails to classify it as a non cancerous image, whereas it is a non cancerous image only.

Calculated values of Sensitivity and Specificity for CLARA algorithm are as follows using equation 8 and 7 respectively:
Sensitivity= 72.60%
Specificity= 29.62 %

### 3)  Results of BIRCH:

Here the confusion matrix is shown by testing with the features of sample size 100 as sown in Table 6. The accuracy is given by the Table 5.

Table 5: Results of BIRCH algorithm

| Number of Cancerous Samples | Number of Non Cancerous Samples | Correctly Classified | Incorrectly Classified |
|---|---|---|---|
| 68 | 32 | 72 | 28 |

**Accuracy %= 72%**

Table 5 represents the number of cancerous and non classified cancerous samples taken for the testing of BIRCH algorithm. By using the accuracy percentage formula in equation number 6 gives the percentage of correctly samples using this algorithm. Hence the calculated percentage of accuracy is 72%, with respect to the number of samples correctly classified.

Table 6: Confusion Matrix from BIRCH Algorithm

|  | Positive | Negative |
|---|---|---|
| Cancerous (True) | 53 | 11 |
| Non Cancerous (False) | 15 | 21 |

Table 6 represents the confusion matrix, according to the classification performed. True positive is the case where the given sample is from the cancerous image, correctly classified. True Negative is the case where the given sample is from the non cancerous image, correctly classified. False positive is the case where the classifier fails to classify it as a cancerous image, whereas it is a cancerous image only. False negative is the case where the classifier fails to classify it as a non cancerous image, whereas it is a non cancerous image only.

Calculated values of Sensitivity and Specificity for BIRCH algorithm are as follows using equation 8 and 7 respectively:

Sensitivity= 71.62 %
Specificity= 42.30 %

**Graphical Representation**: Figure 13 represents the comparison of the results of accuracy, sensitivity and specificity of PAM, CLARA and BIRCH algorithms.
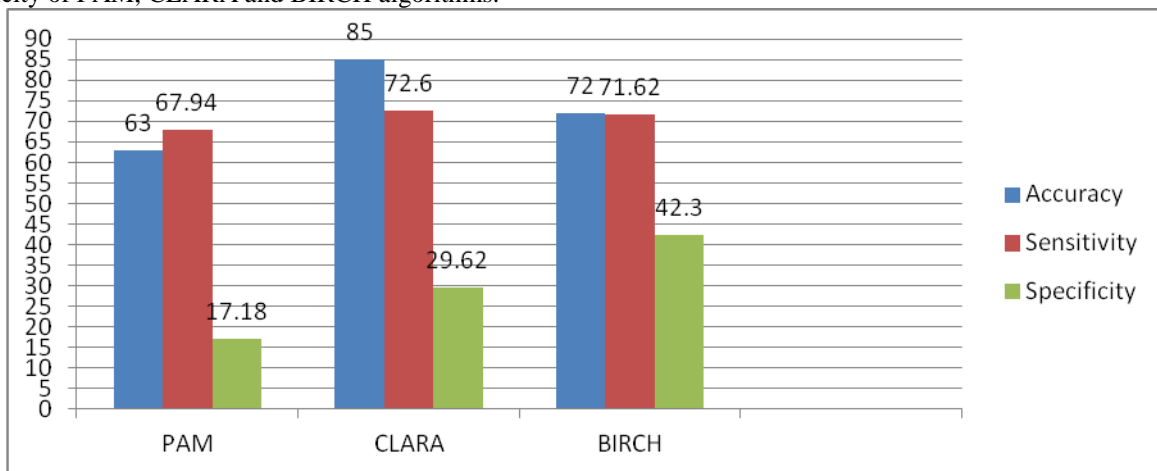


Figure 13: Graphical Comparison of PAM, CLARA and BIRCH

The graphical representation in Figure 13 represents the comparison of the results of the 3 algorithms, PAM, CLARA and BIRCH, with respect to accuracy, sensitivity and specificity. Here the accuracy percentage of CLARA algorithm is found to be highest as 85 %. The accuracy of BIRCH is found to be 72%. Whereas the accuracy of PAM algorithm is found to be least, 65%. The sensitivity is calculated on the basis of true positive and false negative, which is also known as true positive rate. It is 67.94 % in PAM, which is least. 72.6% in CLARA, whereas in BIRCH algorithm it is found out to be 71.62 which is highest among the three algorithms. The specificity is calculated on the basis of true negative and false positive, which is also known as true negative rate. It is 17.18% in PAM, 29.62% in CLARA and 42.3 % in BIRCH. Specificity is found to be highest in BIRCH algorithm and least in PAM algorithm.

## V.  CONCLUSION AND FUTURE WORK

The paper describes about the modality of mammography for imaging the breasts for the presence of cancerous tissue. Here the poor contrast image is being enhanced by using morphological operators, Top hat operation. The contrast of the image is improved by using contrast stretching and the features are extracted from the image using discrete wavelet transform. The image is segmented and the desired region of interest is being taken out. From this region of interest, various parameters are being extracted whose values are being fed to different data mining algorithms such as PAM, CLARA and BIRCH. The accuracy, sensitivity and specificity are calculated for the three algorithms. The accuracy of PAM algorithm is 63%. The accuracy of CLARA algorithm is 85% and the accuracy of BIRCH is 72%. Hence it can be concluded that the accuracy of CLARA algorithm is found to be highest for the dataset of 200 samples where 100 is used for training and 100 is used for testing.

In future work, the mammogram images can be compared with the images for breast cancer tissue screened with other modalities such as Ultrasound, MRI etc and the Image fusion algorithms can be applied to the images obtained from different modalities and hence obtaining the further details of the anatomical structure.

## REFERENCES

1.      Amberly K. Windisch A. Patrick Schneider, Christine M. Zainer, Christopher Kevin Kubat, Nancy K. Mullen, "The breast cancer epidemic: 10 facts." [Online]. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4135458/. , , Date Accessed on [28-May-2017]

2.   ICMR News, "Over 17 lakh new cancer cases in India by 2020: ICMR." [Online]. Available: http://icmr.nic.in/icmrsql/archive/2016/7.pdf , Date Accessed on [28-May-2017]
3.   P. Spandana, K. M. M. Rao, and P. B. V. V. S. N. Prabhakar, "Novel Image Processing Techniques for Early Detection of Breast Cancer , Mat lab and Lab view implementation," *Point-of-Care Healthc. Technol.*, vol. 1, no. 1, pp. 16–18, 2013.
4.   A. Strauss, A. Sebbar, S. Desarnaud, and P. Mouillard, "Automatic detection and segmentation of microcalcifications on digitized mammograms Anne Strauss, Abdeljalil Sebbar, Serge Desarnaud, Pierre Mouillard Michele," *Cent. Morphol. Math. Ec. des Mines Paris, Fontainebleau*, vol. 1, no. 5, pp. 5–6, 2014.
5.   Y.Ireaneus Anna Rejani, "Early Detection of Breast Cancer using SVM Classifier Technique," IJCSE, vol. 1, no. 3, pp. 127–130, 2013.
6.   X.-W. Zhu and L. Xiao, "Research of Blind Watermark Detection Algorithm Based on Wavelet and Contourlet Transform Domain," *2009 Int. Conf. E-bus. Inf. Syst. Secur.*, vol. 1, no. 1, pp. 1–5, 2009.
7.   K. Ganesan, U. R. Acharya, C. K. Chua, L. C. Min, and K. T. Abraham, "Computer-Aided Breast Cancer Detection Using Mammograms : A Review," *IEEE Rev. Biomed. Eng.*, vol. 6, no. 1, pp. 77–98, 2013.
8.   S. Padmanabhan and R. Sundararajan, "Enhanced accuracy of breast cancer detection in digital mammograms using wavelet analysis," *2012 Int. Conf. Mach. Vis. Image Process.*, vol. 1, no. 1, pp. 153–156, 2012.
9.   L. Kaufman and P.J. Rousseeuw, "Finding Groups in Data: an Introduction to Cluster Analysis", John Wiley & Sons, 1990. ISBN-0-471-73578-7
10.  Raymond T. Ng and Jiawei Han CLARANS: A Method for Clustering Objects for Spatial Data Mining, *IEEE Conf.Inf.Syst,*vol.14,no. 5, 2012.
11.  Tian Zbau.g, Ragbu Rau)akt-ishuan, aud Mirou Liv,,y, BIRCH: An Effcient Data Clustering Mtthod for VPTV Largr Databases, *Technical Report, Computer Scieuces Dept.,Univ. of Wisconsiu-Madison*, 1995.

## BIOGRAPHY

**Mr. Tarun Gangil** is a Research Scholar at Department of Medical Electronics, Dr. Ambedkar institute of Technology, Bangalore, Karnataka, India. He received Masters in Technical Education(M.Tech) degree in 2015 from Manipal University, Manipal, Karnataka India, His area of Interest is Medical Image Processing, Data Mining, Machine Learning.

**Ms. Sneha Prajwal**, Research Scholar, Dept of CSE at Bapuji institute of engineering and technology Karnataka, India. She received her Master of Technology degree in 2014 from BIET, Karnataka, India. Her research area are Data mining, image processing, Big Data analytics, cloud computing.