

# Survey and Comparison of Classification Algorithm for Medical Image

*Sweety Maniar, Jagdish S. Shah*  
Department of Computer/IT Engineering  
Gujarat Technological University

**Abstract**—Medical image classification and retrieval systems have been finding extensive use in the areas of image classification according to imaging modalities, body part and diseases. One of the major challenges in the medical classification is the large size images leading to a large number of extracted features which is a burden for the classification algorithm and the resources. In this paper, a novel approach for automatic classification of fundus images is proposed. The method uses image and data pre-processing techniques to improve the performance of machine learning classifiers. Some predominant image mining algorithms such as Classification, Regression Tree (CART), Neural Network, Naive Bayes (NB), Decision Tree (DT) K-Nearest Neighbor

## II RELETED WORK

**Index Terms**—classification, medical imaging, neural networks, texture and shape features.

## I. INTRODUCTION

Image mining is the process used to extract meaningful information from images. It deals with the embedded knowledge, extracting inherent data, image data relationship and other patterns that are not clearly found in the images [1]. There is a system which is Content Based Image Retrieval (CBIR) which aims at searching of images available in databases for any particular images so as to get a related image. The extracting images based on some features such as shape, texture, region and so on. On the other end, Retrieval of image is the fast developing and challenging research part in both unmoving and moving images. Especially, the medical image classification plays an important role in human diagnosis and treatment. It is also used for healthcare students in the educational domain and studies by explaining with these images. Medical images are mainly used to detect specific diseases occur in the human body.

Image matching is more important in the field of mining images. Frequently used technique is nearest neighborhood in which objects are represented as n dimensional vectors. In [6] the visual queries are represented in the retrieval process. So that the images mainly based on the user request and the mechanism is considered as query-by-example used to compare the target images to find the image indices present in the image database. For ease of access digital medical images stored in huge databases as well as Content based image retrieval (CBIR) which is mainly used in diagnostic cases like query medical image. The CBIR images is based on some features such as edge, shape and texture which are extracted automatically [7]. If there is empty in the image set or less than the total images then the system randomly chosen the image for creating the association rules. This paper gives a survey on several techniques in image mining which was already proposed method they are Neural Network, CART, Naive Bayes, KNN and Decision Tree. This paper provides best method in medical image classification based on the classification accuracy, processing time and error rates.

Deshpande et al [8] provides data mining approach which is used to identify the image content present in the association rules. The association rule algorithm helps to detect the regular item set with the help of some iterative methods. This algorithm helps to minimize the number of scans in Apriori algorithm. It is very essential to advance the image quality and make the extraction phase as simple and reliable.

Li-Hong Juang et al [9] focused on tracking tumor objects of (MRI) brain images by using K-means algorithm. The process which is also useful for detecting exact lesion objects in images. The main purpose of this algorithm is to resolve the MRI image by changing the gray-level image into colour image.

S.L.A. Lee et al [10] concentrated on lung nodule detection which is used to spot the lung abnormalities in CT lung images with the help of Random forest algorithm. This algorithm provides hybrid random forest based nodule classification. It is also used to detect 32 patients with 5721 images. The accuracy in proposed system is noted as 97.11 whereas in the developed system the high receiver operator characteristic is given 97.86% accuracy.

Mahnaz Etehad Tavakol et al [11] provide the high infrared cameras to diagnose the vascular changes of breasts by using the adaboost algorithm. The algorithm is used to classify the invisible images into benign, malignant and normal. In this system the accuracy of 83% is given which gives better performance than the proposed system of 66%.

Ming-Yih Lee et al [12] proposed an entropy based feature extraction and some other protocols for the breast cancer diagnosis using decision tree algorithm. The Morphological operations used in this system to detect the unified abnormal regions. This method gives 86% accuracy which is better than the proposed system of 59%.

Ye Chen et al [13] focused on the detection of brain structural changes from the Magnetic resonance images which helps to aid the treatment of neurological diseases with the help of

Support Vector Machine algorithm. In addition the algorithm which helps to analyze the MR images from the various datasets. The accuracy range between 70% and 87% are noted.

Wen-Jie Wu et al [14] suggested both the classification accuracy and the optimal classification model which helps to detect the ultrasound breast tumor images by using genetic algorithm. The algorithm is to calculate the near optimal parameters to differentiate the tumor as benign or malignant. The accuracy of proposed system is 95% which is improved better in the developing system by reducing the biopsies of benign lesions.

Daniel J. Evers et al [15] has given the study to evaluate whether the optimal spectroscopy improve the accuracy of transthoracic lung biopsies using Classification and regression tree (CART) algorithm. Based on the derived parameter the algorithm classifies the type of tissue present in the system. The overall accuracy is 91% sensitivity.

Daniel J. Evers et al [16] has given the study to evaluate whether the optimal spectroscopy improve the accuracy of transthoracic lung biopsies using Classification and regression tree (CART) algorithm. Based on the derived parameter the algorithm classifies the type of tissue present in the system. The overall accuracy is 91% sensitivity.

Min-Chun Yang et al [17] enhance the naïve bayes classification algorithm by separating the ultra sound images pixel- by-pixel then the image measured by gray scale is converted to binary image which is then evaluated by two-phase criteria. So, the detection sensitivity can be further developed.

Shengjun Zhou et al [18] suggested that in the medical applications the images are segmented. To manage the segmentation, fuzzy c-means clustering do the classification of pixels into some divisions. Then the algorithm assigns the membership values for those pixels to form the centroid.

Ravi Babu et al. [19] focused to determine the image classification rate for the purpose of digital image classification. The K-Nearest neighbor algorithm uses the learning technique to find out the classification time of those images. The lazy based and instance based are the two learning techniques. To compare the curves the algorithm is used which based on some comparison. Finally the nearest neighbor classifiers used to measure the distance of the two curves [20].

### III. WORKING OF CLASSIFICATION SYSTEM

To automatically categorize medical images, we have experimented on real mammograms with two data mining techniques, association rule mining and neural networks. In both cases, the problem consists of building a mammography classification model using attributes extracted from and attached to mammograms, then evaluating the effectiveness of the model using new images. The process of building the classification model (classifier) includes preprocessing and extraction of visual features from already labelled images (i.e. training set).

Figure 1 shows an overview of the categorization process adopted for both systems. The first step is represented by the image acquisition and image enhancement, followed by feature extraction. The last one is the classification part where the technique for supervised learning is different. All these parts of the classification systems are discussed in more detail later.

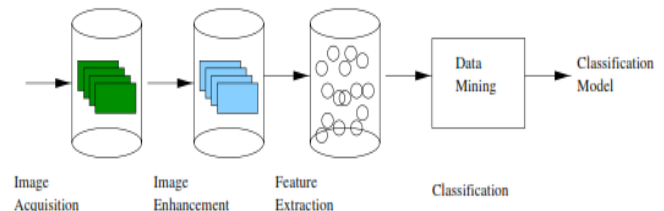


Figure 1. Classification System

#### 1. Image Acquisition

To have access to real medical images for experimentation is a very difficult undertaking due to privacy issues and heavy bureaucratic hurdles. The data collection that was used in our experiments was taken from the Mammographic Image Analysis Society (MIAS) [18]. This same collection has been used in other studies of automatic mammography classification.

#### 2. Image Enhancement

Mammograms are images difficult to interpret, and a preprocessing phase of the images is necessary to improve the quality of the images and make the feature extraction phase more reliable. Pre-processing is always a necessity whenever the data to be mined is noisy, inconsistent or incomplete and pre-processing significantly improves the effectiveness of the data mining techniques

Image enhancement helps in qualitative improvement of the image with respect to a specific application [10]. In order to diminish the effect of over brightness or over darkness in the images and accentuate the image features, we applied a widely used technique in image processing to improve visual appearance of images known as Histogram Equalization. Histogram equalization increases the contrast range in an image by increasing the dynamic range of grey levels (or colours) [10]. This improves the distinction of features in the image. The method proceeds by widening the peaks in the image histogram and compressing the valleys. This process equalizes the illumination of the image and accentuates the features to be extracted. That is how the different illumination conditions at the scanning phase are reduced.

#### 3. Feature Extraction

Feature extraction is a best form of dimensionality reduces. When the input to the various methods are too big to be give and it is believed to be disreputably unneeded (more data, but not more information) then the input data will be changed into a compact version with different number of features (also named as features vector). Storing the input data into the other

format of features is called features extraction. The numbers of techniques for feature extraction are given below.

#### A. Texture

Texture demonstration can be of different types: structural and statistical. Statistical features can be calculated with co occurrence matrices, principal component analysis. [13]The features like mean variance standard deviation, energy, entropy, correlation, inertia are calculated using co-occurrence matrix. Contrast is the compute of difference in the gray level for co occurrence matrix [9].

$$\text{Mean} = \sum_{i=1}^n \sum_{j=1}^m x_y / mn \quad (3)$$

$$\text{Variance} = \frac{1}{mn} \sum_{i=1}^n \sum_{j=1}^m x_y - \text{Mean} \quad (4)$$

$$\sigma = \sqrt{\text{Variance}} \quad (5)$$

$$\text{Correlation} = \sum_{i,j} \frac{(i-\mu_i)(j-\mu_j)P_{ij}}{\sigma_i \sigma_j} \quad (6)$$

$$\text{Entropy} = - \sum_i \sum_j P_d(i,j) \log P_d(i,j) \quad (7)$$

$$\text{Contrast} = \sum_i \sum_j (i-j)^2 P_d(i,j) \quad (8)$$

In Equation 1 . m and n are size of image , $X_{ij}$ . In Equation 5 and Equation 6 is the pixel at i and j position,  $P_d(i,j)$  is the probability distribution function.

#### B. Shape

Shape features have an significant role in primary group of medical images based on their content [2]. Features such as Area, Edge, Fourier Descriptor, Circularity, are used to retrieve medical images [14,8].

**Area:** Area of selection in square pixels or in calibrated square units.

**Edge:** Using canny edge detector, gradient, and other operators.

**Fourier Descriptor:** Fourier Descriptors (FDs) is a powerful feature for boundaries and objects representation.

$$a(k) = \sum_{n=0}^{N-1} z(n) \exp\left[\frac{-j2\pi kn}{N}\right], 0 \leq k \leq N-1 \quad (9)$$

Discrete Fourier Transform of  $z(n)$ (boundary point) gives value of Fourier Descriptor.

$$\text{Circularity} = 4\pi \left( \frac{\text{Area}}{\text{Perimeter}^2} \right) \quad (10)$$

Equivalence diameter (circle with same area as the region)

$$\text{Equivalence Diameter} = \sqrt{\frac{4 * \text{Area}}{\pi}} \quad (11)$$

### 4. Classification Algorithm

#### A. Neural Network

Artificial neural network models have been studied for many years in the hope of achieving human-like performance in several fields such as speech and image understanding. The networks are composed of many nonlinear computational elements operating in parallel and arranged in patterns reminiscent of biological neural networks.

Computational elements or nodes are connected in several layers (input, hidden and output) via weights that are typically adapted during the training phase to achieve high performance. Instead of performing a set of instructions sequentially as in a Von Neumann computer, neural network models explore simultaneously many hypotheses using parallel networks composed of many computational elements connected by links with variable weights

The architecture of the neural network consists of three layers: an input layer, a hidden one and an output layer. The number of nodes in the input layer is equal to the number of elements existing in one transaction in the database. While the output layer was consisting of one node.

The node of the output layer is the one that gives the classification for the image. It classifies it as normal or abnormal. In the training phase, the internal weights of the neural network are adjusted according to the transactions used in the learning process. For each training transaction the neural network receives in addition the expected output. This allows the modification of the weights. In the next step, the trained neural network is used to classify new images.

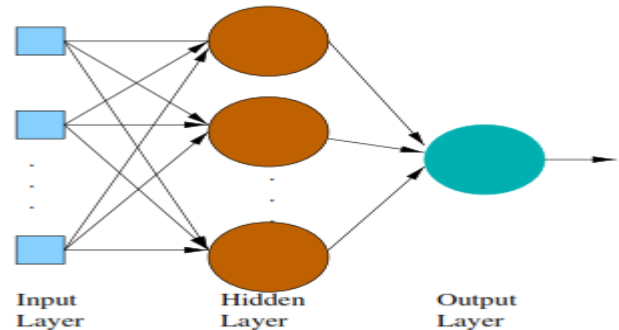


Figure 2. Neural Network

#### B. CART

The classification and regression tree (CART) algorithm is mainly used for the classification of different tissues in image mining, which is on the basis of several derived parameters. The recursive partitioning method used in the CART algorithm to introduce the tree based modelling which is later converted to the statistical mainstream. To select the optimal tree value the algorithm involves the cross validation scheme from some rigorous approaches. Based on the technique called surrogate splits the algorithm automatically handles the missing values. For example the variable ( $x=t1$ ) is selected then the greatest separation is produced so ( $x=t1$ ) is said to be split. If this variable X it sends to which is less than t1 then the data is send to left or else it sends to right. The process is repeated for all the nodes. So that it is easy to conclude that CART algorithm uses only the binary splits.

### C. K-Means

K-Means algorithm is said to be an unsupervised clustering algorithm. It works well for numerical data alone. The pixel-by-pixel image classification is possible by defining single and multiple thresholds. So that histogram statistics is used in this algorithm for the pixel based classification. The main work of this process is to check whether the histogram is bimodal or not. If it is then the gray value will be appeared otherwise the images get partitioned into several regions. The threshold of gray value can be determined using the peak values. However it converges only the local minimum values. So the algorithm involves number of clusters for the optimization.

### D. Naive Bayes

The Naive bayes algorithm is the most powerful technique. It does the testing process easily and the classification problems can be solved. It can be able to build a model fastly and giving better predictions. To find the missing data the naïve bayes algorithm plays a major role. The unseen data can be easily predicted by characterizing the problem in naïve bayes method. During the construction time and prediction time this algorithm separates the attributes value. The probability of each attributes in isolation process needs only the enough data. So, there is no need of more data collection in this algorithm. Finally, if the data has high correlated features the performance will be degraded.

### E. Decision Tree

Decision tree algorithm is one of the classifier technique which is in the form of tree structure. For classification and prediction, the powerful tools are available in this algorithm. It has four divisions such as Decision node, leaf node, edge and path. A single attribute is represented in the decision node. Leaf node defines the target attribute. Splitting of one attribute is edge and the path is a final decision. For continuous attribute this algorithm is not applicable.

## IV.COMPARITION OF CLASSIFICATION ALGORITHM

In this part, the comparative results and the datasets are listed for the data mining algorithms. The accuracy of various algorithms is clearly noted in this study.

### A. Dataset Description

Various image datasets helps to find the classification performance of data mining algorithms. The used data sets are shown in table 1.

S. NO	ALGORITHM	DATASET
1	<b>KNN</b>	Brain Images
2	<b>CART</b>	Lung Images
3	<b>K-Means</b>	Brain Lesion image Dataset
4	<b>Naive Bayes</b>	Breast Lesion Images
5	<b>Decision Tree</b>	Breast Images

Table 1. Dataset Comparison

### B. Comparison of Data Mining Algorithms

This part lists out the positive and negative aspects used in various algorithms present in this paper for the data mining algorithm.

S. NO	ALGORITHM	PURPOSE	LIMITATIONS
1	<b>KNN</b>	It is used to analyze the MR images from the heterogeneous dataset	Some of the features are not properly used in local image features.
2	<b>CART</b>	This algorithm is used to enable the accurate fraction estimation of the substances	Complex classification steps are followed.
3	<b>K-Means</b>	It is used to find exact lesion objects	Parameters are not sufficient for the detection process
4	<b>Naive Bayes</b>	It aims to enhance computer aided system to offer real time detection It improves the detection sensitivity.	The speckle noises present in images affects the pixel classification Low scan speed.
5	<b>Decision Tree</b>	Thermograph images was projected for the feature extraction	Credibility and sensitivity are not accurate.

Table 2. Algorithm Comparisons

## V.CLASSIFICATION PARAMETER

The confusion matrix can be used to determine the performance of the system. This matrix describes all possible outcomes of a prediction results in table structure. The possible outcomes of a two class prediction be represented as True positive (TP), True negative (TN), False Positive (FP) and False Negative (FN). The normal and abnormal images are correctly classified as True Positive and True Negative respectively. A False Positive is when the outcome is incorrectly classified as positive when it is a negative. False Positive is the False alarm in the classification process. A false negative is when the outcome is incorrectly predicted as negative when it should have been in fact positive.

In our system consider,

TP= Number of Abnormal images correctly classified

TN= Number of Normal images correctly classified

FP= Number of Normal images classified as Abnormal

FN= Number of Abnormal images classified as Normal.

**Precision:** The fraction of abnormal images with correct results.

$$TP/(FP+TP)$$

**Sensitivity (Recall):** The probability of the test finding the abnormal case among all abnormal cases.

$$TP/(FN+TP)$$

**Specificity:** The probability of the test finding the normal case among all normal cases.

$$TN/(FP+TN)$$

**Accuracy:** The fraction of test results those are correct.

$$(TN+TP)/FP+TN+FN+TP$$

## VI. CONCLUSION

In this study, the overall performance of various algorithms present in this paper was analyzed based on the classification accuracy. The accuracy in image classification is the main idea of evaluating the performance in data mining algorithms. The



overall result shown in this paper is step into further development in future technology. To evaluate the best indications clinical studies are more essential. This paper projected several parameters which is used by data mining methods. By comparing these algorithms, neural network results shows better performance among the other methods presented in this work and it is also gives the best classification accuracy as well as saves the computing time. In future, we have planned to enhance neural network method to improve the classification accuracy, sensitivity, specificity and as well as processing time.

## VII. REFERENCES

- [1] C. Lakshmi Devasena, T.Sumathi, Dr. M. Hemalatha “An Experiential Survey on Image Mining Tools, Techniques and Applications” International Journal on Computer Science and Engineering (IJCSSE), Vol. 3, No. 3, pp. 5061-5067, 2011.
- [2] E. Arnoldi, M. Gebregziabher, U.J. Schoepf, R. Goldenberg, L. Ramos-Duran, P.L. Zwerner, K. Nikolaou, M.F. Reiser, P. Costello, C. Thilo “Automated computeraided stenosis detection at coronary CT angiography: initial experience”, Elsevier, Vol. 20, No.5, pp. 1160-1167, 2010.
- [3] Lefkovits Sz, Leftkovits L. “Enhanced Gabor Filter Based Facial Feature Detector, The proceeding of the European Integration- Between Tradion and modernity”, Elsevier, computer science section, vol. 1, pp. 78738-8863, 2013.
- [4] R.Venkata Ramana Chary, Dr.D.Rajya Lakshmi and Dr. K.V.N Sunitha “Feature extraction methods for color image similarity”, Advanced Computing: An International Journal (ACIJ), Vol.3, No.2, pp. 147-157, 2012.
- [5] Ankit Agrawal, Sanchit Misra, Ramanathan Narayanan, Lalith Polepeddi, Alok Choudhary, “A Lung Cancer Outcome Calculator Using Ensemble Data Mining on SEER Data”, BIODDD 2011, No. 5, 2011.
- [6] S.Balan and T.Devi, “Design and Development of an Algorithm for Image Clustering In Textile Image Retrieval Using Color Descriptors”, International Journal of Computer Science, Vol.2, No.3, pp. 199-211, 2012.
- [7] Vamsidhar Enireddy, Kiran Kumar Reddi, “A Data Mining Approach for Compressed Medical Image Retrieval”, International Journal of Computer Applications (0975 – 887) Vol. 52, No.5, 2012.
- [8] D. Deshpande, “Association Rule Mining Based on Image Content”, International Journal of Information Technology and Knowledge Management, Vol. 4, No. 1, pp. 143-146, 2011.
- [9] Li-Hong Juang , Ming-Ni Wu “MRI brain lesion image detection based on color-converted K-means clustering segmentation”, elsevier, computer science section, vol.43 ,No.7 , pp. 941-949, 2010.
- [10] S.L.A. Leea, A.Z. Kouzania, E.J. Hub “Random forest based lung nodule classification aided by clustering”, elsevier, computer science section, vol.34, No.7, pp. 535-542, 2010.
- [11] Mahnaz EtehadTavakol , Vinod Chandran , E.Y.K. Ng , Raheleh Kafieh, “Breast cancer detection from thermal images using bispectral invariant features”, elsevier, International journal of thermal sciences, vol.69, pp. 21-36, 2013.
- [12] Ming-Yih Leea, Chi-Shih Yang, “Entropy-based feature extraction and decision tree induction for breast cancer diagnosis with standardized thermograph images”, elsevier, computer science meth ods and programs in biomedicine, vol.100, No.3, pp. 269-282, 2010.
- [13] Ye Chena, Judd Storrs, Lirong Tana, Lawrence J. Mazlackc, Jing-Huei Leeb, Long J. Lua, “Detecting brain structural changes as biomarker from magnetic resonance images using a local feature based SVM approach”, elsevier, Journal of Neuroscience methods, vol. 221 , No. 15 , pp. 22-31, 2014.
- [14] Wen-Jie Wua, Shih-Wei Lina, Woo Kyung Moonb, “Combining support vector machine with genetic algorithm to classify ultrasound breast tumor images”, elsevier, computer science section, vol. 36 , No. 8 , pp. 627-633, 2012.
- [15] Jarich W. Spliethoff , Daniel J. Ever, Houke M. Klomp, Johanna W. van Sandick, Michel W. Wouters, Rami Nachabe, Gerald W. Lucassen, Benno H.W. Hendriks, Jelle Wesseling, Theo J.M. Ruers,” Improved identification of peripheral lung tumors by using diffuse reflectance and fluorescence spectroscopy”, elsevier, computer science section, vol.80 , No. 2 , pp. 165-171, 2013.
- [16] Min-Chun Yang, Chiun-Sheng Huang, “Whole breast lesion detection using naive bayes classifier for portable ultrasound”, elsevier, computer science section, vol. 38, No.11, pp. 1870-1880, 2012.
- [17] Shengjun Zhou, Yuanzhi Cheng, Shinichi Tamura, “ Automated lung segmentation and smoothing techniques for inclusion of juxtapleural nodules and pulmonary vessels on chest CT images”, elsevier, computer science section, vol. 13, pp. 62-70, 2014.
- [18] U Ravi Babu, Y. Venkaswarlu, Aneel Kumar Chintha. “Handwritten Digit Recognition Using K-Nearest Neighbour Classifier” IEEE World Congress on Computing and Communication Technologies ISBN: 978-1-4799-2876-7 pp 60-65, 2014.
- [19] U Ravi Babu, Y. Venkaswarlu, Aneel Kumar Chintha. “Handwritten Digit Recognition Using K-

Nearest Neighbour Classifier” IEEE World Congress on Computing and Communication Technologies ISBN: 978-1-4799-2876-7 pp 60-65, 2014.

[20] Umut Konur, Fikret S. Gurgun, “ Computer aided detection of spina bifida using nearest neighbor

classification with curvature scale space features of fetal skulls extracted from ultrasound images, vol. 85, pp. 80-95, 2015.