



Medical Image Segmentation and Classification Using MKFCM and Hybrid Classifiers

Kottayilkunnel Gopalapillai Satheesh^{1*} Alex Noel Joseph Raj¹

¹Vellore Institute of Technology University, India

* Corresponding author's Email: kgsatheeshphd2016@gmail.com

Abstract: Tuberculosis (TB) is a common infectious disease caused by bacteria named mycobacterium tuberculosis, which is preventable and curable if detected early. In feature extraction of medical images, any unwanted features extracted may lead to efficiency loss. To overcome this, the features are optimized using Orthogonal Learning Particle Swarm Optimization (OLPSO) technique, which is used to identify the specific set of features from the image and ranks the features based on decision task equation. Based on which the images are classified. In addition, this paper proposes a hybrid classification to differentiate the images as Cavitory TB and Miliary TB by nomination method of classification. The hybrid classifier is an integration of Support Vector Machine (SVM) and Artificial Neural Network (ANN) which are applied to CT scan lung images to provide results with high accuracy. This experiment results show that, it is possible to identify and classify TB images by using MATLAB classifiers.

Keywords: Artificial neural network, Orthogonal learning particle swarm optimization, Support vector machine, Tuberculosis.

1. Introduction

Tuberculosis (TB) is a contagious disease caused by the bacteria, which generally affects the lungs. TB is airborne, which means it can spread from the infected person to vulnerable people through the air, which makes it unpredictable. It can be treated if detected and treatment is undertaken in the early stage. It still remains the ruling cause of death in the developed and developing countries. As per WHO report, 1.5 million deaths occurred due to TB in the year 2013. There are many researches being conducted to develop a framework to medically diagnose and detect TB disease using CT scan images [1]. Most image processing techniques fail to extract the information from an image such as position, shape, etc., which becomes even more challenging when the technique should can handle features like quality, quantity. After medical data is processed, Classification of medical data is another important task in the diagnosis of this disease. To overcome such problems and with an aim to achieve promising results in the diagnosis of TB, in the proposed

technique we employ a hybrid classifier. The steps that are involved in processing of medical images are pre-processing, segmentation, feature extraction, classification [2]. The methods used for performing these processes in our system is explained below.

Initially, the pre-processing is performed to remove noise in input medical data, for this wiener filter technique is employed in this paper. It involve in converting the raw data into a suitable format for further processing. The segmentation process is carried out by Multiple Kernel Fuzzy C Means (MKFCM) clustering technique, which is one of the successful segmentation techniques [3-4]. It improves and manages the problems caused by over or under segmentation. After the segmentation process is completed, the meaningful information or features are to be extracted from the lung images, through which various ranks of features are collected. In the collected features, only the top ranked features are then forwarded, and the remaining are rejected. The feature extraction is achieved using Orthogonal Learning Particle Swarm Optimization (OLPSO) technique in this proposed system. This OLPSO leads

to present an optimizer which cooperatively optimizes the requirement problem and helps to place the correlated dimensions in the same swarm. Once the feature extraction phase is completed, the classification will be carried-out using a hybrid classifier [5-6]. Classification includes, two important processes that are, training phase and testing phase. In the training phase, the hybrid classifier is trained with extracted features by using training data and in testing, the features of the image which is segmented are fed into the trained classifier to detect whether the image has Cavitory TB or Miliary TB. This process of classification is carried out using the SVM classifier. The improvement in classifier performance can produce even better results in the classification accuracy [5]. Once the classification into Cavitory TB or Miliary TB has been carried out, it has also been proposed to classify the stages of TB accordingly as, Stage 1 or Stage 2 or Stage 3 based on the severity of infection by using ANN classifier. This hybrid classifier which is a combination of SVM classifier and ANN classifier is used to detect the TB in a short time with more accuracy, which perform well than the other existing methods [7-10]. The reduction of time required and complexities have been greatly reduced.

2. Literature review

Z. H. Zhan et al. [11] have proposed an orthogonal learning (OL) strategy for PSO. They entitled this PSO as orthogonal learning particle swarm optimization (OLPSO). Here, they applied to both global & local versions of PSO that produced the OLPSO-G & OLPSOL algorithms, correspondingly. This new algorithm and learning strategy were tested on a set of 16 benchmark functions & were compared with other existing PSO algorithms & some state of the art evolutionary algorithms. The test outcomes showed the efficiency and productivity of the proposed learning strategy & algorithms. The comparisons illustrated that OLPSO expressively enhanced the performance of PSO, provided a faster global convergence, greater solution quality & stronger robustness.

J. Kuruvilla and K. Gunavathi [12] have presented a computer based classification methodology in computed tomography (CT) images of lungs made by utilizing artificial neural network. The complete lung was segmented from the CT images. At the same time, the parameters were computed from the segmented image. The statistical parameters, for instance mean, standard deviation,

skewness, kurtosis, fifth central moment & sixth central moment were utilized to do the classification. This procedure was carried out by feed forward as well as feed forward back propagation neural networks. The feed forward back propagation network delivered improved classification in comparing to feed forward networks. The parameter skewness provided the highest classification precision. Amongst the existing thirteen training functions of back propagation neural network, the Traingdx function delivered the highest classification precision of 91.1%. They proposed two new training functions. The outcomes of the tests revealed that the proposed training function 1 gave precision of 93.3%, specificity of 100% & sensitivity of 91.4% & a mean square error of 0.998. The proposed training function 2 provided a classification precision of 93.3% & minimum mean square error of 0.0942.

H.C. Huang et al. [14] have proposed a multiple kernel fuzzy c-means (MKFC) algorithm that extends the fuzzy c-means algorithm with a multiple kernel learning setting. This proposed method was more resistant to ineffective kernels & irrelevant features through combining multiple kernels & automatically adjusting the kernel weights. This thing prepared the selection of kernels less critical. Tests on synthetic & real-world data showed the efficiency of the proposed MKFC algorithm.

Most of the existing methods has been encountered with various challenges in their perception. Removal of noises in the images especially in case of medical image is the common risky problem in all methods. This may leads to misclassification, various system complexities and low performance nature as a result, which should be noted. In addition, time requirement, sensitivity and over/ under segmentation are some of the issues seen in all the system. This system is designed in such a way by considering the noise problem as a major one. The use of OLPSO with the combination of hybrid classifier (SVM+ANN) involve in noise reduction, further improves the performance comparatively than the other existing methods described above.

3. Proposed methodology for using MKFCM and hybrid classifiers

Segmentation and deduction of Lung Tuberculosis using ranked based OLPSO and MKFCM method has been proposed in this paper. The paper consists of pre-processing, segmentation, feature extraction, classification which are discussed in the following sections.

3.1 Pre-processing using Wiener filter and region growing

Medical images for instance magnetic resonance imaging, x-ray, CT etc. all these images are collected from different sources, thus contaminating the image with different kinds of noises. One such noise is additive noise. The objective of wiener filter is to get rid of the additive noise from images using statistical method. Typically all the filters are intended to do the function for a preferred frequency response & known source noise (equipment noise). This filtering technique minimize the mean square error. The Wiener filter implemented by,

$$X(i, j) = \frac{H^*(i, j)p_s(i, j)}{|H(i, j)|^2 p_s(i, j) + p_n(i, j)} \quad (1)$$

Where, $X(i, j)$ = wiener filter for noise removal in an image. $H(i, j)$ = Degradation function on image. $H^*(i, j)$ = Complex conjugate of degradation function on image. $p_n(i, j)$ = Power Spectral Density of Noise on image. $p_s(i, j)$ = Power Spectral Density of un-degraded image. The image after removal of additive noise is forwarded to the next phase.

Once the image is ready for further process, the desired lung region is still with unwanted body portion, this has to be eliminated for faultless nodule extraction. Here we make use of the primary region removal algorithm, to separate the required portion. To perform this task we employ region growing for body portion stripping. The objective of the region growing is to map the input image data into sets of connected pixels, known as regions, with regards to a predefined criterion that usually inspects the characteristic of local groups of pixels. A set of connected region pixels has two portions, lung and extracted body portion. From this only lung portion is captured for nodule extraction, using MKFCM technique.

3.1.1 Nodule detection multiple kernel fuzzy c-means

The fuzzy c-means algorithm (FCM), have been broadly used in the medical image segmentation and such a wide range of use of it is because of the application of fuzziness to categorize every image pixel. FCM have noteworthy utilization in this area of study for the reason that, they could keep more information from the original image with compare to hard c-means segmentation techniques. The key benefit of FCM is it lets pixels to belong to several clusters with a reasonable level of membership. The motive of the FCM algorithm is to separate the vector

space into subspaces according to the distance measure. The FCM algorithm does not take into consideration about the local spatial property that makes it sensitive to noise.

3.1.2 FCM algorithm

Given a data set $X = \{x_1, x_2, \dots, x_n\}$, where the data point $x_j, R^p (j = 1, 2, \dots, n)$, n is the number of data points, and p is the input dimension of a data point, traditional FCM groups X into c clusters by minimizing the weighted sum of distances between the data and the cluster centres or prototypes defined in Eq. (2)

$$Q = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m \|x_j - o_i\|^2 \quad (2)$$

Here, $\| \cdot \|$ is the Euclidean distance. u_{ij} is the membership of data x_j belonging to cluster i , which is represented by the prototype o_i . The constraint on u_{ij} is $\sum_{j=1}^n u_{ij}^m = 1$ and m is the fuzzification coefficient.

3.1.3 Kernel fuzzy c-means

The fundamental concept of KFCM is to first map the input data into a feature space with greater dimension utilizing a non-linear transform & then attain FCM in that feature space. As a result the original complex & nonlinearly separable data structure in input space might turn into simple & linearly separable in the feature space after the nonlinear transforms. Therefore we will be able to get superior performance. One more benefit of KFCM is, unlike FCM which require the preferred number of clusters in advance, it can adaptively determine the number of clusters in the data under some criteria. In KFCM, the input data is the integration of the pixel intensity & the local spatial information of a pixel which is designated for clustering. KFCM restricts that the prototypes in the kernel space are in fact mapped from the original data space or the feature space. The objective function is presented in Eq. (3)

$$Q = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m \|\Phi(x_j) - \Phi(O_i)\|^2 \quad (3)$$

The function is reformulated in Eq. (4)

$$Q = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m (1 - k(x_j, O_i)) \quad (4)$$

Where, $(1 - k(x_j, O_i))$ is robust distance measurement derived in the kernel space.

3.1.4 Multiple kernel fuzzy c-means

In KFCM, the input data is the $Q = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m \|\Phi_{com}(x_j) - \Phi_{com}(O_i)\|^2$ group of the pixel intensity & the local spatial information of a pixel that is carefully chosen for clustering. Now, multiple kernels or composite kernels are taken into consideration in place of a single fixed kernel. With multiple kernels, the kernel techniques achieve more flexibility on kernel choices and additionally reflect & fuse data from numerous heterogeneous or homogeneous sources. Specially, in image segmentation problems, the input data includes characteristic of image pixels resultant from diverse sources. So, we can define various kernel functions deliberately for the intensity information & the texture information separately, & after that we merge these kernel functions & apply the composite kernel in MKFCM to get better image-segmentation outcomes. The combination of the collaborative kernel can be automatically adjusted in the learning of multiple-kernel FCM (MKFCM) or it can be settled by trial & error or cross-validation. The general context of MKFCM (Eq. (5)) purpose is to reduce the objective function.

$$Q = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m \|\Phi_{com}(x_j) - \Phi_{com}(O_i)\|^2 \quad (5)$$

Where $\Phi_{com}(x_j)$ represents the composite kernel data and $\Phi_{com}(O_i)$ represents the composite prototype in the kernel space, u_{ij} is the membership function. To enhance the Gaussian-kernel-based KFCM by adding a local information term in the objective function. The objective function is represented in Eq. (8)

$$Q1 = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m (1 - k(x_j, O_i)) \quad (6)$$

$$Q2 = \alpha \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m (1 - k(\bar{x}_j, O_i)) \quad (7)$$

$$Q = Q1 + Q2 \quad (8)$$

Where, x_j is the intensity of pixel j . In the new objective function, the additional term is the weighted sum of differences between the filtered intensity x_j (the local spatial information) and the clustering prototypes. It is worth pointing out that k_1 or k_2 in the first variant of MKFCM-K-based image segmentation can be changed to any other Mercer kernel function for the information related to image pixels a composite kernel that joins different shaped kernels can be defined as in Eq. (9).

$$k_{com} = k_1 + \alpha k_2 \quad (9)$$

Where α is the constant of Mercer kernel function, k_{com} represents the composite kernel and k_1 (Eq. (10)) is still the Gaussian kernel for pixel intensities

$$k_1(x_i, x_j) = \exp(-|x_i - x_j|^2 / r^2) \quad (10)$$

k_2 (Eq. (11)) is a polynomial kernel for the spatial information

$$k^2(x_i, x_j) = (x_i x_j + d)^2 \quad (11)$$

If $k_{com} = k_1 + \alpha k_2$ is the composite kernel, the minimized objective function of the MKFCM is derived and defined in Eq. (12)

$$Q = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m \|\Phi_{com}(x_j) - O_i\|^2 \quad (12)$$

Where O_i is the prototype of kernel metric function and x_j is the data belong to the composite kernel object that is represented as $\Phi_{com}(x_j)$. For example, the input image data x_j is set to be $x_j = [x_j, x_j, s_j]$ the same as the third variant of MKFCM, then the composite kernel (k_L) is designed as in Eq. (13)

$$k_L = w_1^b k_1 + w_2^b k_2 + w_3^b k_3 \quad (13)$$

To reduce the impact of outliers the MKFCM algorithm assess the centroids. The importance of MKFCM-based image-segmentation algorithms is the flexibility in selections & combinations of the kernel functions in diverse shapes & for diverse pieces of information. After integrating the different kernels in the kernel space, there is no requirement to alter the computational processes of MFKCM. This is one more benefits to reflect & fuse the image information from several heterogeneous or homogeneous sources. In the MKFCM method, we can simply fuse the texture information into segmentation algorithms by simply incorporating a kernel designed for the texture information in the composite kernel. The portion that is extracted is known as lung nodules. It comprises of both types of nodules, tuberculosis nodule & a regular nodule. There is a necessity for partitioning the class on the basis of the features which are extracted from both types of designated nodules.

3.2 Feature extraction using first order statistics

Statistical approaches used to measure the spatial distribution of pixel values and the following features

are first order statistics as discussed in calculated. Features are defined in Eqs. (14)-(20).

a. Variance:

Variance in the gray level in a neighbourhood region of a pixel is a measure of the texture and the histogram of intensity levels is used as a concise and summary of the statistical information of an image.

$$variance = \sum_{i=0}^{G-1} (i - \mu)^2 p(i) \quad (14)$$

Where $p(i)$ represents the pixels having the gray-scale intensity in the i^{th} level.

b. Dispersion:

The purpose of measures the dispersion is to find out how spread out data values are on the number line. Another term for these statistics is measures of spread. The measure of dispersion compared with Interquartile range, mean absolute deviation, Central moment of all orders, range, Standard deviation, Variance.

c. Skewness:

The skewness gives an asymmetric value of the probability distribution about its mean, which an image can be positive and negative.

$$skewness = \frac{\sum_{i=1}^M \sum_{j=1}^N [p(i,j) - \mu]^3}{(MN)\sigma^3} \quad (15)$$

d. Kurtosis:

It measures the flatness of the distribution with respect to the normal distribution.

$$kurtosis = \frac{\sum_{i=1}^M \sum_{j=1}^N [p(i,j) - \mu]^4}{(MN)\sigma^4} \quad (16)$$

f. Average energy:

It's calculated from the sum of the square of all the elements in GLCM and which gives the uniformity of the pixels. The value of the energy may be 0 and 1. The image constant value is one.

$$energy = \sum_{i,j} p(i,j)^2 \quad (17)$$

g. Inverse difference moment (IDM):

It is a measure of local homogeneity,

$$IDM = \frac{1}{R} \sum_{i=1}^M \sum_{j=1}^N \frac{1}{1+(i-j)^2} p(i,j) \quad (18)$$

h. Sum average:

The sum average calculated by using the following equation,

$$SM = \sum_{i=1}^{2M} i p_{x+y}(i) \quad (19)$$

Where $p_{x+y}(k) = \sum_{i=1}^M \sum_{j=1}^N p(i,j) \quad k = 2,3, \dots, 2M$

i. Sum entropy:

The sum entropy calculated by using the following equation,

$$SE = - \sum_{i=1}^{2M} p_{x+y}(i) \log p_{x+y}(i) \quad (20)$$

After the successful feature extraction, the set of features of both types of nodules are procured. These feature sets are further fed to OLPSO for feature ranking. The acquired ranked features are considered as a primary feature classification.

3.2.1 Feature selection using OLPSO

We initiated our research work focused on six first-order statistical features but, as our analysis has exposed, a number of them are irrelevant because of their mutual correlations. T test was utilized to evaluate the importance of features. To increase the performance of the classifier, it is essential to select the best feature or the feature combination. The productivity can be maximized when the best obtained features are classified utilizing hybrid classifiers.

The t value of the T test for independent samples is defined by relation (21), where $m1$ and $m2$ represent the mean values of the two independent samples and EE_{m1-m2} (22) is the standard error of the mean difference.

$$t = \frac{m1-m2}{EE_{m1-m2}} \quad (21)$$

$$EE_{m1-m2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad (22)$$

Where $s1$ and $s2$ are the standard deviations of the two samples and $n1$ and $n2$ are the number of samples.

The number of degrees of freedom is:

$$df = n_1 - n_2 - 1 \quad (23)$$

This decision task Eq. (23) is considered as the fitness equation for OLPSO technique. Fitness value is Maximum when we consider the feature combination as the best feature. Initialization of the particle population & updating of velocity in OLPSO as follows,

Particle Swarm Optimization (PSO) is population-based stochastic-optimization iterative learning algorithm which emulates the swarm behaviours for instance birds flocking & fish schooling. PSO finds for an optimal selected features over each particle flying in the search space and by adjusting its flying direction according to its experience of personal best and neighbourhood best. The information of neighbourhood's best experience & personal best experience is utilized in a simple way and therefore the flying is adjusted by learning of summation of two experiences. So, the challenge is to use both personal & neighbourhood best to guide the particle towards the global best. Orthogonal experimental design (OED) provides a capability to find the best combination levels for various features with a comparatively small number of experimental samples. We utilize the OED method to develop a promising learning exemplar. The OED is utilized to find the best combination of particle's best historical position & its neighbourhood's best historical positioned is additionally utilized to form an orthogonal learning (OL) strategy for PSO to find and keep useful information of a particle's best position & its neighbourhood's best position. The OL strategy deliver better productivity to PSO & hence yield better global optimization.

In a D-dimensional hyperspace, particle i , with velocity $V_i = [v_{i1}, v_{i2}, \dots, v_{iD}]$, position vector $X_i = [x_{i1}, x_{i2}, \dots, x_{iD}]$ to indicate its current state, where i is a positive value. Its personal historical best position $P_i = [p_{i1}, p_{i2}, \dots, p_{iD}]$. The best position of all particles in i^{th} particle neighbourhood is $P_i = [p_{n1}, p_{n2}, \dots, p_{nD}]$. V_i and X_i are initialized randomly, updated generation to generation.

Updating,

$$u_{id} = c_1 r_{id} (p_{id} - x_{id}) + c_2 r_{2d} (p_{nd} - x_{id}) \quad (24)$$

$$v_{id} = v_{id} + u_{id} \quad (25)$$

$$X_{id} = x_{id} + v_{id} \quad (26)$$

In Eq. (24) c_1 and c_2 are acceleration parameters, r_{id} and r_{2d} are random values. V_{MAXd} is used to clamp the updated velocity. If $|v_{id}|$ exceeds V_{MAXd} , then it is set to $\sin v_{id} V_{MAXd}$. In this way the particles are driven to p_i add p_n . Equation (26) states how velocity is updated in position vector.

To control or adjust the flying velocity, however, an inertia weight or a constriction factor ω in Eq. (27)

$$v_{id} = \omega v_{id} + u_{id} \quad (27)$$

In the traditional PSO, each particle updates its flying velocity and position according to its personal best position and its neighbourhood's best position. The concept is simple and appealing, but this learning strategy can cause the "oscillation" phenomenon and the "two steps forward, one step back" phenomenon. This velocity represents in Eq. (28)

$$v_{id} = (p_{id} - x_{id}) + (p_{nd} - x_{id}) \quad (28)$$

Using the OED method, the original PSO can be modified as an OLPSO with an OL strategy that combines information of p_i and p_n to form a better guidance vector p_o . The particle's flying velocity is thus changed as

$$v_{id} = \omega v_{id} + c r_d (p_{od} - x_{id}) \quad (29)$$

In Eq. (29) ω is the same as in Eq. (24) and c is fixed to be 2.0, the same as c_1 and c_2 , and r_d is a random value uniformly generated within the interval $[0, 1]$. The guidance vector p_o is constructed for each particle i , respectively, from p_i and p_n as

$$p_o = p_i \oplus p_n \quad (30)$$

In Eq. (30) p_o defines best particle for decided task selection which are further fed into a hybrid classification.

3.2.2 Hybrid classification using ANN and SVM

A feed forward neural network is a biologically motivated classification algorithm. This algorithm comprises of a number of simple neuron-like processing units, organized in layers. Data enters at the inputs and passes through the network, layer by layer, until it arrives at the outputs. At the time of normal operation, that is when it acts as a classifier, there is no feedback between layers. This is why they are also known as feed forward neural networks. The back propagation neural network is a multi-layered, feed forward neural network, which is best known example for training algorithm. It has following steps,

Step 1: select a network architecture

Step 2: randomly initialize weights

Step 3: while error is too large

Select training pattern and feed forward to find the actual network output

Calculate errors and back propagate error signals

Adjust weights

Step 4: Evaluate performance using the test set

The data sets are divided into three subsets.

Training set: the training of this model is using this set, by pairing this set with expected output set.

Validation set: in order to estimate how well the model is trained this set of data is utilized. Test set: this set of results are compared with the validation set to test for the validity of the obtained results. Classified Nodules are again processed using SVM classifier wherein, based on the selected features. But this stage of nodule cannot be predicted, so we employ SVM and some new features for to detect the stage of nodule. The new features (Eqs. (31)-(33)) are Area, Eccentricity and Equivalent Diameter (ED).

$$Area = \sum_{x=1}^m \sum_{y=1}^n fb(x, y) \quad (31)$$

$$Eccentricity = \frac{\{\sum_{x=1}^m fb(x, :)\}}{\{\sum_{y=1}^n fb(:, y)\}} \quad (32)$$

$$ED = \frac{2}{\sqrt{\pi}} \sqrt{Area} \quad (33)$$

These new features are fed as input for Multiple SVM clustering to predict the stage of the nodule. SVM classification is essentially a binary (two-class) classification technique, which has to be modified to handle the multiclass tasks in real world situations. For a k-class problem, these methods design a single objective function for training all k-binary SVMs simultaneously and maximize the margins from each class to the remaining ones. Here, we consider the method by Weston and Watkins wherein k represents stages of tuberculosis. Given a labelled training set represented by $\{(x_1, y_1), \dots, (x_l, y_l)\}$ of cardinality l , where $x_i \in R^d$ and $y_i \in \{1, \dots, k\}$, the formulation proposed is given as follows:

$$\begin{aligned} & \min \\ & w_m \in h, b \in R^k, \\ & \xi \in R^{l \times k} \frac{1}{2} \sum_{m=1}^k w_m^T w_m + C \sum_{i=1}^l \sum_{t \neq y_i} \xi_{i,t} \quad (34) \end{aligned}$$

Subjected to

$$\begin{aligned} (w_{y_i}^T \phi(x_i) + b_{y_i}) & \geq (w_t^T \phi(x_i) + b_t) + 2 - \xi_{i,t} \\ \xi_{i,t} & \geq 0, i = 1, \dots, l, t \in \{1, \dots, k\} \setminus y_i \end{aligned}$$

The resulting decision function is

$$\arg \max x_m f_m(x) = \arg(x) \quad (35)$$

$$\arg(x) = \arg \max_m (w_m^T \phi(x) + b_m) \quad (36)$$

Equation (36) states that the presence of tuberculosis in the image, and then the resulting vector identifies the class of input tuberculosis.

4. Experiments and results

In this work, Lung Image Database Consortium (LIDC) dataset is employed for the classification of nodule, non-nodule and overlapping of both nodule and non-nodule. This large database consists of 4,532 nodules extracted from CT images. Here, the nodules $\leq 5\mu m$ are utilized, because it shows an effective outcome in determination of both nodule and non-nodule features.

The multiple KFCM and their combination with the hybrid classifiers are evaluated by performing experiments on the medical data sets and the results are tabulated for verification that is shown in the following part.

The original image (Fig. 3) taken as an input from the dataset and N number of samples image used for processing. In first stage of process, we need to remove the outer region of lung by using region growing technique.

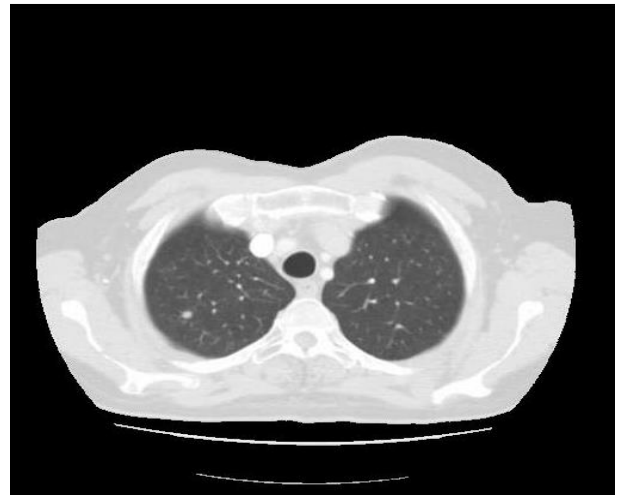


Figure.1 Input Image

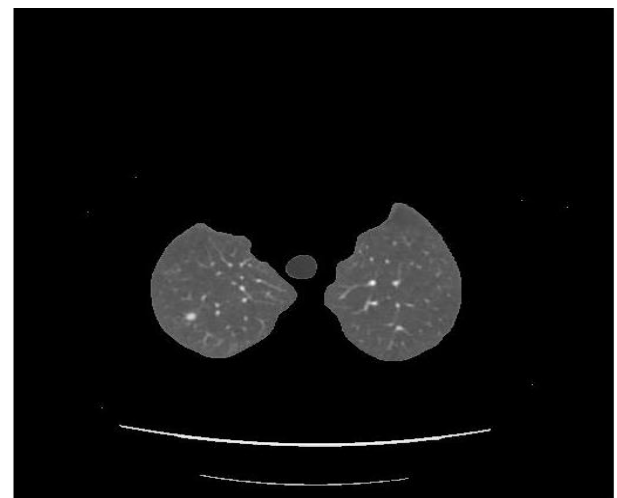


Figure.2 Identified lung portion

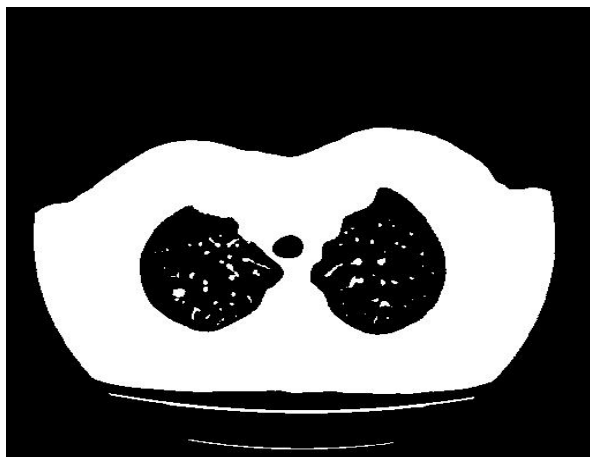


Figure.3 KFCM based segmentation

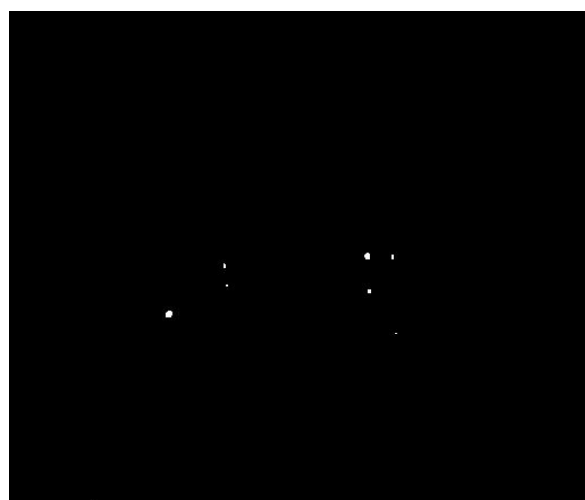


Figure.4 Detected major tuberculosis regions

Figure 4 shows the divided two sessions namely left and right lungs that are processed to apply clustering algorithm to group the similar region based on using threshold values as 3*3 then region are similar pixel values also segmented by using MKFCM algorithm.

After segmentation, that small size of nodules removed from the image, which is helps to next process of classification to the nodules as tuberculosis or not.

Figure 5 is used to classify the given image as a tuberculosis or not. Each nodule as to discriminate using a hybrid classifier. The Linear SVM used to identify the nodules are normal or abnormal. The features of the nodules from the above Fig.6 are subjected to OL-PSO algorithm for finding the primary rank based features which will be used by the hybrid classifiers.

The lungs image divided into two sessions, namely the lung right and lung left. After clustering, the first order statistics features processed for both lungs, calculated the values separately and both feature values are merged to find the average value.

The table 1 represents the values of nodule and nodule values of selected features from which the features which has notable changes are selected as primary features used for hybrid classifiers. Here, the selected primary features are listed as variance, dispersion, skewness and kurtosis.

Table 1. Nodule and non-nodule values of features

First Order statistics	Non-Nodule	Non-Nodule	Nodule	Nodule	Rank Selected Non-Nodule	Rank Selected Nodule
Variance	4.7010	5.5618	4.2807	4.1291	5.5618	4.2807
Dispersion	61.6901	106.8	93.4032	87.4010	106.8	93.4032
Skewness	0.0923	0.0010	0.0740	0.1002	0.0923	0.1002
Kurtosis	3.7053	2.4056	2.2450	3.1210	3.7053	3.1210
Avg. energy	0.5790	0.6146	0.6402	0.6510	0	0
IDM	0	0	0	0	0	0
Sum average	11.720	13.6	13.75	13.333	0	0
Sum entropy	0.0100	0.0008	0.007	0.008	0	0

Table 2. First order statistics for classified nodule and non- nodule

First order statistics	Krishnamurthy et al [17]	Non-Nodule	Nodule1	Nodule 2	Nodule 3
Area	49	40	37	32	26
ESI	0.098	0.552	0.589	0.63	0.67

1. The values variation occurs due to the probability of information gained from the investigation method done in particular features. Ranks are provided for the selected nodule and non- nodule from the segmented lung image.
2. The value tends to zero, if the probability of feature details become invalid.
3. IDM is a feature that is not easily accessible and it shows invalid output most of the times, thus the nodule and non- nodule values become zero.
4. In addition, the rank values takes the average value received from the combination of statistics or features respectively, in which the statistic values combined with Avg. energy, IDM, Sum average, Sum entropy is always results in unclear and provides zero value.
5. The ranks of selected nodule and non- nodule for those four features remains zero as their combination will not produce a valid outcome.

After ranking of features, the classification is performed using SVM classifier. The classification takes place on three factors Area, Eccentricity, and Edge Sharp Index (ESI).

From the above table 2, the area and ESI variation in the nodule and non- nodule image investigation. In this SVM algorithm, shape-based analysis is carried out to identify the Tuberculosis cell. The deviation in the shape is given by analysing the above table. The nodules taken for classification are divided into three nodules as nodule 1, 2, and 3 based on the tumour size.

1. In first phase, the experiment is conducted on a tiny small sized cell i.e. Nodule 1, in which

it does not provide a clear identification of the feature and shape variation is more due the increased area deviation, eccentricity and decreased ESI value.

2. The nodule 2 inspection is on the medium sized infected cell that resembles as that of a tumour cell. The variation in eccentricity is still considerably high as it struggles to identify whether it is a tumour or other diseased cell. It also indicates that the nodule 2 has irregularity in shape between consecutive slices compared with all other nodules.
3. The nodule 3 reveals the examination of tumours from a large sized tumours cell, which shows low area and eccentricity and high edge sharp index. The classification result is more successful in large sized tumour i.e. nodule 3.
4. In order to differentiate the characterization between nodule 2 and nodule 3 and their deviation, it requires a highly advanced algorithm to perform, which can be taken as a future work.

The statistic data taken for calculation leads to the performance result of the proposed algorithm. The formula for finding the statistic deals with the score values of True positive (TP), True negative (TN), False positive (FP) and False negative (FN) identification. The performance measure of this work is with the good sensitivity of 89.87%, specificity of 82.88%, positive likelihood ratio of 5.24, negative likelihood ratio of 0.122, and disease prevalence of 7.69% respectively. The occurrence of False Positive (FP) rate is also reduced to a great extent simultaneously.

Table 3. Statistical Results

Statistic	Formula	Gori et al. [15]	Messay et al. [16]	Krishnamurthy et al [17]	Proposed value
Sensitivity	TP/TP + FN	74.3%	82.66%	88%	89.87%
Specificity	TN/TN + FP	82.3%	83.27%	84.05%	82.88%
Positive likelihood ratio	Sensitivity/100-Specificity	4.197	4.940	5.51	5.24
Negative likelihood ratio	100-Sensitivity/Specificity	0.3122	0.208	0.142	0.122
Disease prevalence	TP+FN/TP + FN + TN + FP	9.46%	9.02%	8.87%	7.69%

5. Conclusion

The proposed hybrid classifier and OLPSO are evaluated by conducting the experiments in lung image. In this study it is identified that the lesion size of Tuberculosis cell in lung is 20mm for normal lung and greater than 20mm as abnormal lung Tuberculosis cell. Thus the resultant output of this technique is obtained and evaluated shows comparatively high performance than the other methods. The first order statistics and their variations from table1 and table2 reveal better evaluation results and the statistical resultant value based on the sensitivity, specificity, PLR, NLR and disease prevalence are also calculated. Based on this study it is more evident that hybrid KFCM segmentation technique is good for evaluation of lung tuberculosis cell region. But the detection and diagnosis of tuberculosis still requires more and more improvement in classifying the nodules can be taken as the future work. By selecting the suitable combination of classification or segmentation algorithm, it can be made possible in future.

References

- [1] A. Kulkarni and A. Panditrao, "Classification of Lung Cancer Stages on CT Scan Images Using Image Processing", In: *Proc. of IEEE International Conference on Advanced communication Control and Computing Technologies (ICACCCT)*, pp.384-1388, 2014.
- [2] U. Bağcıa, M. Brayb, J. Cabanc, J. Yaod, and D.J. Molluraa, "Computer-assisted detection of infectious lung diseases: A review", *Computerized Medical Imaging and Graphics*, Vol.36, No.1, pp.72-84, 2012.
- [3] H.B. Nandpuru, S.S. Salankar, and V.R. Bora, "MRI Brain Cancer Classification Using Support Vector Machine", In: *Proc. of IEEE Students' Conference on Electrical, Electronics and Computer Science*, pp.1-6, 2014.
- [4] P. Yugander, B.J. Sheshagiri, K. Sunanda, and E. Susmitha, "Multiple Kernel Fuzzy C-Means Algorithm with ALS method for Satellite and Medical Image Segmentation", In: *Proc. of International Conference on Devices, Circuits and Systems (ICDCS)*, pp.244-248, 2012.
- [5] S. A. Patil, V. R. Udipi, C. D. Kane, A. I. Wasif, and J. V. Desai, "Geometrical and Texture Features Estimation of Lung Cancer and TB Images Using Chest X-ray Database", In: *Proc. of International Conference on Biomedical and Pharmaceutical Engineering*, pp.58-75, 2009.
- [6] M. Seeraa and C.P. Limb, "A hybrid intelligent system for medical data classification", *Expert Systems with Applications*, Vol.41, No.5, pp.2239-2249, 2014.
- [7] S.M. Jadhav, "Generalized Feed forward Neural Network based cardiac arrhythmia classification from ECG signal data", In: *Proc. of 6th international conference on Advanced Information Management and Service (IMS)*, pp. 351-356, 2010.
- [8] B. Shanna and K. Venugopalan, "Classification of Hematomas in Brain CT Images using Neural Network", In: *Proc. Of Issues and Challenges in Intelligent Computing Techniques (ICICT)*, pp.41-46, 2014.
- [9] T. Sun, J. Wang, X. Li, P. Lv, F. Liu, Y. Luo, Q. Gao, H. Zhu, and X. Guo, "Comparative evaluation of support vector machines for computer aided diagnosis of lung cancer in CT based on a multi-dimensional data set", *Computer Methods and Programs in Biomedicine*, Vol.111, No.2, pp.519-524, 2013.
- [10] W. Shen, M. Zhou, F. Yang, C. Yang, and J. Tian, "Multi-scale Convolutional Neural Networks for Lung Nodule Classification", In: *Proc. of International Conference on Information Processing in Medical Imaging*, pp.588-599, 2015.
- [11] Z.H. Zhan, J. Zhang, Y. Li, and Y.H. Shi, "Orthogonal Learning Particle Swarm Optimization", *IEEE transactions on evolutionary computation*, Vol.15, No.6, pp.832-847, 2011.
- [12] J. Kuruvilla and K. Gunavathi, "Lung cancer classification using neural networks for CT images", *Computer Methods and Programs in Biomedicine*, Vol.113, No.1, pp.202-209, 2014.
- [13] M. Keshani, Z. Azimifar, F. Tajeripour, and R. Boostani, "Lung nodule segmentation and recognition using SVM classifier and active contour modeling: A complete intelligent system", *Computers in Biology and Medicine*, Vol.43, pp.287-300, 2013.
- [14] H.C. Huang, Y.Y. Chuang, and C.S. Chen, "Multiple kernel fuzzy clustering", *IEEE Transactions on Fuzzy Systems*, Vol.20, No.1, pp.120-134, 2012.
- [15] S. Krishnamurthy, G. Narasimhan, and U. Rengasamy, "Three dimensional lung nodule segmentation and shape variance analysis to detect lung cancer with reduced false positives", In: *Proc. of the Institution of Mechanical Engineers, Part H: Journal of Engineering in Medicine*, Vol.230, No.1, pp.58-70, 2016.

- [16] S. Thawkar and R. Ingolikar, "Automatic Detection and Classification of Masses in Digital Mammograms", *International Journal of Intelligent Engineering and Systems*, Vol.10, No.1, pp.65-74, 2017.
- [17] D. Mahammad Rafi and C.R. Bharathi, "Optimal Fuzzy Min-Max Neural Network (FMMNN) for Medical Data Classification Using Modified Group Search Optimizer Algorithm", *International Journal of Intelligent Engineering and Systems*, Vol.9, No.3, pp.1-10, 2016