

Classification of Respiratory Tract Diseases Using Machine Learning Techniques

G. Natarajan

Research Scholar

Department of Computer Science and Engineering
Annamalai University

Dr. P. Dhanalakshmi

Professor

Department of Computer Science and Engineering
Annamalai University

Abstract

Pleural effusion is an uncommon lung disease characterized by a build-up of fluid between the two layers of the pleura that causes specific symptoms, such as chest pain and shortness of breath. Machine learning techniques have been widely used for abnormality detection in medical images. Chest X-ray images (CXR) are among the non-invasive diagnostic tools used to detect various disease. In the proposed work, two different feature extractions namely Scale Invariant Feature Transform and Zernike features are extracted. The extracted features are fed into the classifier namely Random Forest and K-nearest Neighbor which classifies into effusion, emphysema, infiltration, no-findings and pleural thickening. The results are compared were ZERNIKE with Random Forest grants the satisfactory results of 94.30 %.

Keywords : Scale Invariant Feature Transform (SIFT), Malignant Pleural Effusions (MPE), Iterated Function System (IFS), Benign Pleural Effusion (BPE).

1. Introduction

Pleural effusion is a specific collection of fluid in the pleural cavity. A pleural membrane covers the lungs, and between the membrane and the lungs, there is a pleural cavity which usually contains about 10-20 ml of fluid that function as a lubricant to allow the lungs to move freely while breathing [1]. However, if the fluid is excess and builds up, can pressure on the lungs, and cause chest pain and shortness of breath; this condition is called pleural effusion Two types of pleural effusions in clinics are: (1) Transudative pleural effusion, resulting from leakage of fluid into the pleural cavity, often caused by heart failure, cirrhosis, and post-operative surgery; (2) Exudative pleural effusion, resulting from blood vessels leaks mainly caused by cancer, tuberculosis, pulmonary embolism, and pneumonia [2]. Exudative pleural effusion can similarly be classified as malignant or benign primarily based on the detection of malignant cells in the pleural fluid. Pleural fluid of malignant pleural effusions (MPE) includes most cancer cells, while a benign pleural effusion (BPE) does not.

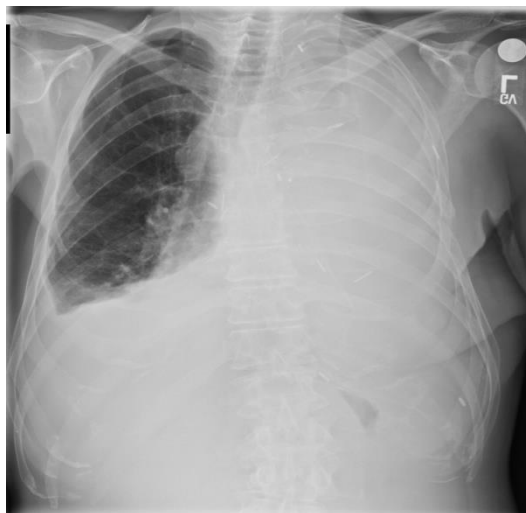


Fig1. a) Shows the abnormal X-ray Image b) Shows ths normal X-Ray image

Transudative and Exudative pleural effusions are distinguished by measuring the lactate dehydrogenase and protein levels in the pleural fluid [3]. Exudative pleural effusions meet as a minimum one of the following criteria

1. Pleural fluid protein/serum protein >0.5
2. Pleural fluid LDH/ serum LDH >0.6
3. Pleural fluid LDH more than two thirds normal upper limit for serum

In the proposed work, Chest x-ray images are given as the input to SIFT and Zernike Features. The extracted features are given as fed into Random Forest and KNN which classifier into five classes namely effusion, Emphysema, Infiltration, No-Findings and Pleural-Thickening. Fig 1 shows the Framework of proposed work.

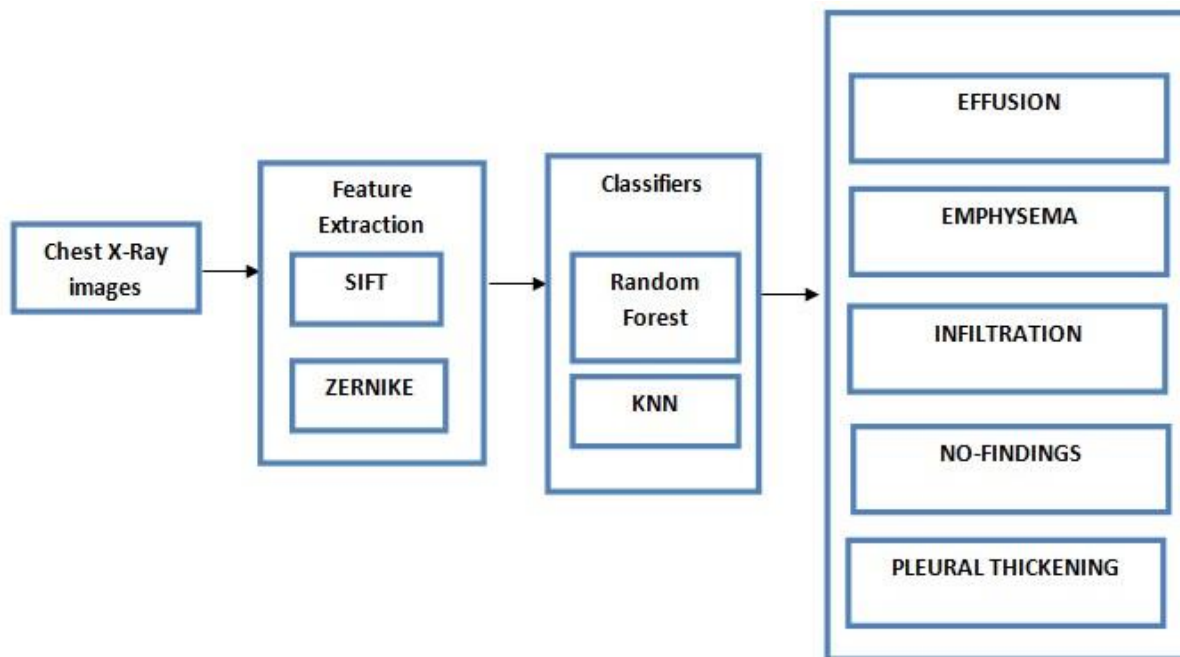


Fig 2. Block Diagram of Classification of Pleural Effusion Symptoms of a pleural effusion are:

- Chest pain
- Dry, nonproductive cough
- Dyspnea (shortness of breath, or difficulty, labored breathing)
- Orthopnea (inability to breathe easily if the person is not sitting or straightening their back)
- The main causes of transudative pleural effusions (watery fluid) include are:
- Heart Failure
- Pulmonary Embolism
- Cirrhosis

Exudative (protein-rich fluid) are pleural effusions most commonly caused by:

- Pneumonia
- Cancer
- Pulmonary Embolism
- Kidney Disease
- Inflammatory disease

Other less common causes of pleural effusion include:

- Tuberculosis

2. Literature of the work

This study aims to propose [4] a system of iterated function system (IFS) and a multi-layered fractional order machine learning classifier to quickly identify the potential classes of lung diseases within regions of interest on CXR images and to improve recognition accuracy. For digital image methods, a two-dimensional (2D) fractional order convolution is used to improve symptomatic properties. The IFS with non-linear interpolation features is then used to reconstruct the 2D characteristic patterns. These reconstructed patterns are self-funded in the same class and therefore help distinguish normal subjects from those with lung diseases. Therefore, accuracy rate is improved. Pooling is done to decrease the dimensions of the feature patterns and to speed up

complex calculations. A classifier based gray relational analysis is used to identify the possible types of the signs and symptoms of lung diseases. For CXR digital images in AP view, it is recommended to use the k-fold multilevel machine learning classifier shows promising results in lung diseases detection and improves the recognition accuracy rate compared to conventional methods. The proposed classifier is rated in terms of recall (99.6%), precision (87.78%), accuracy (88.88%), and F1 score (0.9334).

A pleural effusion is easy to identify and to quantify, which is selected as the subject of this study, which aims to develop an automated system for the interpreting the LUS of the pleural effusion [5]. At the Royal Melbourne Hospital, a LUS data set consisting of 623 videos containing 99,209 2D ultrasound images from 70 patients with an in-phase array transducer was collected and a standardized protocol was followed that included scanning six anatomical regions to have full lung coverage for diagnosing respiratory diseases. This protocol, combined with a deep learning algorithm using a network of spatial transformers, provides a basis for the automatic classification of pathologies on an image-based level. In this work, the deep learning model was trained with supervised and weakly supervised approaches using frame-based and video-based ground truth labels respectively. The interpretation of images by specialists served as a reference. Both approaches showed comparable accuracy scores over the entire test set of 92.4% and 91.1%, without statistically significant differences. The video-based labeling approach however, requires significantly less effort on the part of clinical experts for ground truth labeling.

In the proposed work [6] convolutional neural networks (CNN) deep architecture classification approaches have gained popularity due to their ability to learn mid and high level image representations. Explore a CNN's to recognize different types of pathologies on chest x-ray images. In addition, since very large training sets are not generally available in the medical field, the possibility of using a deep learning approach based on non-medical learning is possible. Since there are usually no very large training sets available in the medical field, it is also possible to use a deep learning approach based on non-medical learning. We tested our algorithm on a 93image data set. A CNN was trained on ImageNet, a well-known large scale non-medical image database. The best performance was achieved with a combination of CNN sourced features and a low-level feature set. We obtained an area under the curve (AUC) of 0.93 for the detection of the right pleural effusion, 0.89 for detection enlarged heart and 0.79 for the classification between healthy and abnormal chest x-ray, in which all pathologies are grouped together in a one large class. This is a first experiment of its kind that shows that deep learning with large scale nonmedical image databases can be sufficient for general medical image recognition tasks.

3. FEATURE EXTRACTION STAGE

3.1 Scale Invariant Feature Transform (SIFT) The SIFT includes four steps.

Scale Space Extrema

To the input image gaussian scale space is constructed [7]. Convolutions is formed on the image with Gaussian functions of varying widths. Then difference of gaussian is applied to the input image for finding the extreme values.

Key point localization

Using Taylor series of expansion the extreme points are localized and Hessian matrix is used to remove the low contrast key points. Compute magnitude and orientation. Compute the feature vector

3.2 ZERNIKE Feature Extraction

Zernike features are extracted from the chest X- Ray images. Teague proposed Zernike moments based on the basis set of orthogonal Zernike polynomials. A discrete image function can be renovated by Zernike moments [8]. Zernike moments are rotation invariant and robust to noise. A relatively small set of Zernike moments can effectively characterize the overall shape of a pattern. Different orders from order 1 to order 14 are evaluated with different image sizes to determine the optimal order. Zernike moments of order 8 with radius 21 which uses only 25 features achieved better results than the other orders for image recognition. The image input is taken as gray image and the feature vector is displayed in the form of floating-point values. The mahotas package provides the Zernike moments () function, is used to compute Zernike moments.

4. MODELING THE FEATURES

4.1 Random Forest

Random Forest (RF) classifier is otherwise called the random decision forests. RF is a group of trees predictors which is called as forest [9]. The RF classifier is mainly used for classification and regression problems

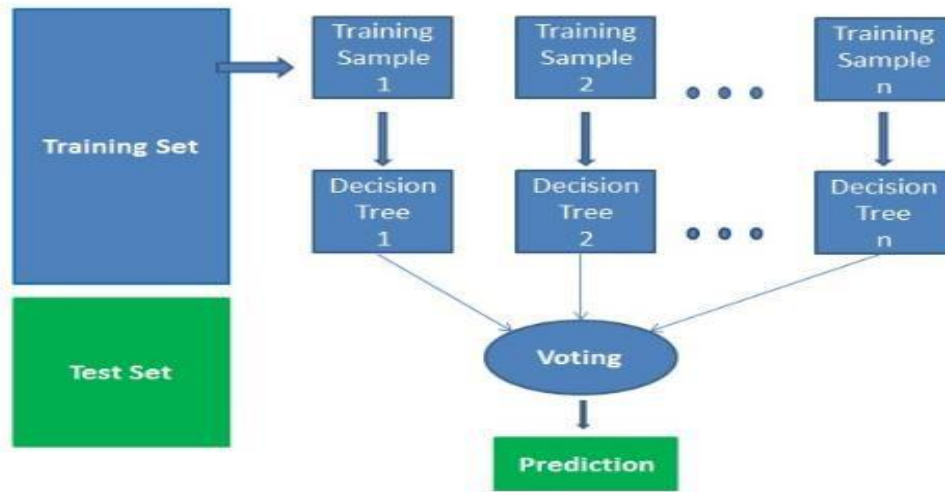


Fig 3. Random Forest Classification Process

All the trees are trained with the same parameters. The bootstrap procedure is followed each training set, then we randomly select the same number of vectors, the vector is replaced with some vector that will appear once and some other can be absent. A new subset is created during training at each node. All the variables are not used to divide the node; a subset is selected randomly to generate a new subset in a node. But the size is fixed for all the nodes and all the trees. During training of the current tree is drawn by replacement, while some vectors are left out. This is called as out-of-bag.

4.2 K-Nearest Neighbor

K-Nearest Neighbour is one of the simplest Machine Learning algorithms based on Supervised Learning technique. K-NN algorithm assumes the similarities between the new data and the existing data and places the new data into the category that closely matches to the available categories. The KNN algorithm stores all the existing data and classifies the new data point based on the similarities. This means when new data emerges and can then be classified into a well-ordered category using K- NN algorithm [10]. The K-NN algorithm can be used for both regression and classification but is mainly used for classification problems. K-NN is nonparametric algorithm, which means it makes no assumption about the underlying data.

Algorithm of K-NN working is given below:

- Step-1: Choose the K number of the neighbors
- Step-2: Calculate the Euclidean distance of K number of neighbors
- Step-3: Take the K nearest neighbors according to the calculated Euclidean distance.
- Step-4: Among these k neighbors, count the number of the data points in each category.
- Step-5: Assign new data points to the category with the maximum number of neighbors.
- Step-6: Our model is ready.

5. PERFORMANCE MEASURES

To evaluate the performance of KNN and Random Forest classifiers using SIFT and ZERNIKE feature extraction techniques, a collection of evaluation parameters namely Accuracy, Precision, Recall and F-score are used in this work. These measures are determined based on the confusion matrix derived from the outcome of the classification process.

Accuracy

The accuracy of a measurement system is a level of measurement that yields are true (no systemic errors) and consistent (no random errors) results.

$$Accuracy = \frac{TP+TN}{TP+TN + FP+ FN}$$

Precision

In classification work, the precision for the class is the number of TP (i.e. the number of items correctly marked as belonging to the positive class) divided by the total number of elements marked as belonging to the positive class (i.e. the sum of true positives and false positives, that are items incorrectly labeled as belonging to the class)

$$Precision = \frac{TP}{TP+ FP}$$

Recall

The recall is defined as the number of TP divided by the total number of elements that actually belong to the positive class (i.e. they should have been)

TP

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

TP+ FN

F-Measure

F-Measure is a measure of test's accuracy and it takes into account of both the precision and recall of the test to compute the score (Harmonic mean).

Precision*Recall

$$F\text{-Measure} = 2 \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

6. EXPERIMENTAL RESULTS

6.1 Dataset

The datasets were collected from NIH Database. A total of 590 X-ray samples were collected from online database from different patients. 465 images were used for training and 125 were used for testing. In this 465 images (163 – Effusion, 42 – Emphysema, 102 – Infiltration, 117 – No- Findings, 41 – Pleural-Thickening). In this 125 testing images (25 – Effusion, 25 – Emphysema, 25 – Infiltration, 25 – No- Findings, 25 – Pleural-Thickening).

6.2 Feature Extraction Using SIFT

Evaluation using Random Forest

The training process analyses pleural effusion training data to find the decision tree to classify pleural effusion affected images into their relevant types. Random forest is structured for a complete learning procedure for categorizing with a collection of decision trees that grow randomly by selecting the sample data. A nonlinear decision tree is applied to discriminate the various types. Random Forest is trained to ascertain pleural effusion features. The bootstrap procedure is followed for each training set; the samples are selected randomly. For training 465 feature vectors, each of 128-dimension are extracted from the images. At each point, a new subset is generated; current tree is taken by replacement of vectors. This is called as out-of-bag. The training process analyses the pleural effusion data to categorize the pleural effusion affected images into its distinctive types, namely, Effusion, Emphysema, Infiltration, No- Findings, Pleural-Thickening and well differentiated. For testing 125 feature vectors each of 128 dimensions are given as input to the Random Forest model. While testing, predictions are arrived by finding the average of the study of each decision tree.

Table 1. Performance of SIFT with Random Forest

	Precision (in %)	Recall (in %)	F-Score (in %)	Accuracy (in%)
Effusion	60.00	65.21	62.49	86.40
Emphysema	80.00	76.92	78.42	91.20
Infiltration	64.00	69.56	66.66	88.00
No-Findings	88.01	75.86	81.48	92.00
Pleural-Thickening	80.00	80.00	80.00	92.01

Evaluation using K-Nearest Neighbor

The parameter K is assigned to indicate the number of nearest neighbors, the distance between the query-instance and all the training samples are calculated by Euclidean distance method. Here the value of K is 5. The distance is calculated for all the training data and the nearest neighbor found based on the Kth minimum distance. All the categories of the training data are received for the sorted

value which falls in K. Then the majority of the nearest neighbors are used as the prediction value. KNN is trained to identify Pleural Effusion features. For each value of K, the test has iterated K times with diverse training and testing sets. For training 465 feature vectors, each of 128 dimension are extracted from the Pleural effusion X-ray images. The training process analyses the pleural effusion training data to find an optimal way to classify Pleural effusion into its respective five classes.

For testing 125 feature vectors, each of 128 dimensions are given as input to the KNN model and the distance between each of the feature vector and the majority nearest neighbor is found among the data in voting procedure. The average distance is calculated for each class. The average distance gives a better performance than using distance for each feature vector. The types of Pleural Effusion images are decided based on the maximum distance.

Table 2. Performance of SIFT with K-Nearest Neighbor

	Precision (in %)	Recall (in %)	F-Score (in %)	Accuracy (in%)
Effusion	48.00	36.36	41.37	72.80
Emphysema	56.00	53.84	54.89	81.60
Infiltration	40.00	50.00	44.44	77.60
No-Findings	72.00	60.00	65.45	84.80
Pleural -Thickening	40.00	56.25	46.75	76.22

6.3 Feature Extraction using ZERNIKE FEATURES

Evaluation using Random Forest

The training process analyses pleural effusion training data to find the decision tree to classify pleural effusion affected images into their relevant types. Random forest is structured for a complete learning procedure for categorizing with a set of decision trees that grow randomly by selecting the sample data. A nonlinear decision tree is applied to discriminate the various types. Random Forest is trained to ascertain pleural effusion features. The bootstrap procedure is followed for each training set; the samples are selected randomly. For training 465 feature vectors, each of 25-dimension are extracted from the images. At each point, a new subset is generated; current tree is taken by replacement of vectors. This is called as out-of-bag. The training process analyses the pleural effusion data to categorize the pleural effusion affected images into its distinctive types, namely, Effusion, Emphysema, Infiltration, No- Findings, Pleural-Thickening and well differentiated. For testing 125 feature vectors each of 25 dimensions are given as input to the Random Forest model. While testing, predictions are arrived by finding the average of the study of each decision tree.

Table 3. Performance of ZERNIKE with Random Forest

	Precision (in %)	Recall (in %)	F-Score (in %)	Accuracy (in%)
Effusion	80.00	87.00	83.35	93.60
Emphysema	74.10	76.00	83.23	93.30
Infiltration	72.00	90.00	80.00	92.60
No-Findings	96.00	84.00	92.36	96.80
Pleural -Thickening	88.00	88.40	88.11	95.20

Evaluation using K-Nearest Neighbor

The parameter K is assigned to indicate the number of nearest neighbors, the distance between the query-instance and all the training samples are calculated by Euclidean distance method. Here the value of K is 5. The distance is calculated for all the training data and the nearest neighbor found based on the K^{th} minimum distance. All the categories of the training data are received for the sorted value which falls in K. Then the majority of the nearest neighbors are used as the prediction value. KNN is trained to identify Pleural Effusion features. For each value of K, the test has iterated K times with diverse training and testing sets. For training 465 feature vectors, each of 25 dimension are extracted from the Pleural effusion X-ray images. The training process analyses the pleural effusion training data to find an optimal way to classify Pleural effusion into its respective five classes.

For testing 125 feature vectors, each of 25 dimensions are given as input to the KNN model and the distance between each of the feature vector and the majority nearest neighbor is found among the data in voting procedure. The average distance is calculated for each class. The average distance gives a better performance than using distance for each feature vector. The types of Pleural Effusion images are decided based on the maximum distance.

Table 4. Performance of ZERNIKE with K-Nearest Neighbor

	Precision (in %)	Recall (in %)	F-Score (in %)	Accuracy (in%)
Effusion	49.00	37.41	41.34	73.62
Emphysema	46.10	43.84	74.89	71.60
Infiltration	49.21	54.39	45.29	67.60
No-Findings	73.00	65.66	65.45	85.78
Pleural - Thickening	42.90	57.25	44.75	77.51

Conclusion

This section made a comparative study with two different feature extraction SIFT and ZERNIKE as well as two classifiers namely Random Forest and KNN. The results are provided in Table 5 . It is noted that the highest accuracy of 94.30% is obtained from ZERNIKE with Random Forest gives the satisfactory results.

Table 5. Comparison of SIFT Vs ZERNIKE with Random Forest and KNN classifiers

Feature Extraction	Classifiers	Precision (in %)	Recall (in %)	F-Score (in %)	Accuracy (in%)
SIFT	Random Forest	74.40	73.51	73.81	89.92
	KNN	51.20	51.29	50.58	78.60
ZERNIKE	Random Forest	82.02	85.08	85.41	94.30
	KNN	52.04	51.71	53.34	75.22

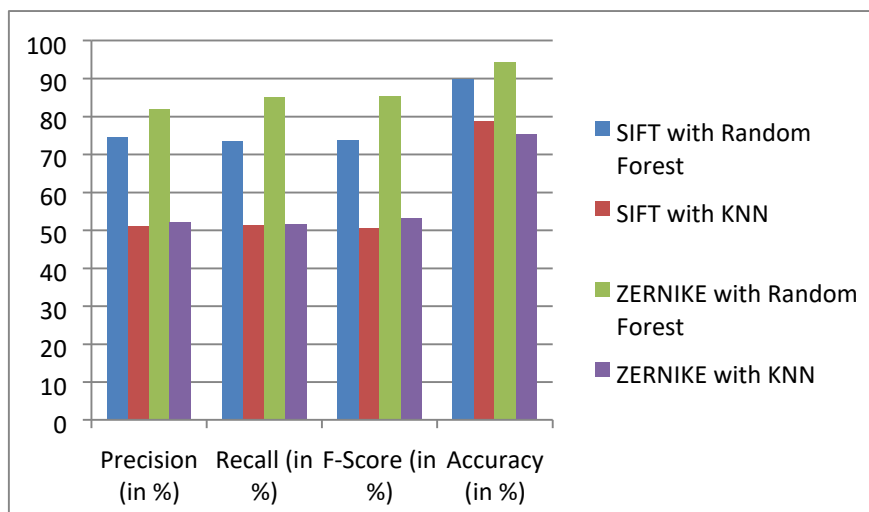


Fig 4. Overall Performance of SIFT, ZERNIKE features using Random Forest and KNN

References

1. Ahmad Rafiansyah Fauzan, Mohammad Iwan Wahyuddin, and Sari Ningsih, "Pleural Effusion Classification Based on Chest X-Ray Images using Convolutional Neural Network", *Journal of Computer Science and Information*, Vol 14, Issue 1, pp 9-16, 2021.
2. Zhongjian Chen, Keke Chen, Yan Lou, Jing Zhu, Weimin Mao & Zhengbo Song, "Machine learning applied to near - infrared spectra for clinical pleural effusion classification", *Scientific Reports*, 2021.
3. Sertan Serte, Ali Serener, "Early pleural effusion detection from respiratory diseases including COVID-19 via deep learning", *Proceedings in IEEE explore*, 2021.
4. Chia-Hung Lin, Jian-Xing Wu, Chien-Ming Li, Pi-Yun Chen, Neng-Sheng Pai and YingChe Kuo, "Enhancement of Chest X-Ray Images to Improve Screening Accuracy Rate Using Iterated Function System and Multilayer Fractional-Order Machine Learning Classifier", *IEEE Photonics Journal*, Vo 12, Issue 4, 2020.
5. Chung – Han Tsai, Jeroen van der Burgt, Damjan Vukovic, Nancy Kaur, Libertarion Demi, David Canty, Andrew Wang, Alistair Royse, Colin Royse, Kavi Haji, Jason Dowling, Girija Chetty, Davide Fontanarosa, "Automatic deep-learning based pleural effusion classification in lung ultrasound images for respiratory pathology diagnosis", *Physica Medica*, Vol 83, pp38-45, 2021.
6. Yaniv Bar, Idit Diamant, Lior Wolf, Hayit Greenspan, "Deep learning with non-medical training used for chest pathology identification", *Medical Imaging : 2015*, Vol 9412, 2015.
7. Chiu, L. C., Chang, T. S., Chen, J. Y., & Chang, N. Y. C., "Fast SIFT design for real-time visual feature extraction", *IEEE Transactions on Image Processing*, Vol 22, Issue 8, pp 31583167 2013.
8. Alessandro Riccardi, Todor Sergueev Petkov, Gianluca Ferri, a Matteo Masotti, and Renato Campanini, "Computer-Aided-Detection of lung nodules via 3D Fast Radial transform, scale space representation and Zernike MIP classification", *The International Journal of Medical Physics Research and Practice*, Volume 38, Issue 4, 2011.
9. Breiman, L., "Random Forests", *Machine learning*, Vol. 45, Issue 1, pp.5-32, 2001.
10. C. Bhuvaneshwari, P. Aruna, D. Loganathan, "A new fusion model for classification of the lung diseases using genetic algorithm", *Egyptian Informatics Journal*, pp 1-9, 2014.