

# **A COMPUTER AIDED DIAGNOSIS SYSTEM FOR LUNG CANCER DETECTION USING SVM**

**A THESIS SUBMITTED TO  
THE GRADUATE SCHOOL OF APPLIED SCIENCES  
OF  
NEAR EAST UNIVERSITY**

**by**

**ERKAN EMIRZADE**

**In Partial Fulfillment of the Requirements for  
the Degree of Master of Science  
in  
Computer Engineering**

**NICOSIA, 2016**

**A COMPUTER AIDED DIAGNOSIS SYSTEM FOR  
LUNG CANCER DETECTION USING SVM**

**A THESIS SUBMITTED TO  
THE GRADUATE SCHOOL OF APPLIED SCIENCES  
OF  
NEAR EAST UNIVERSITY**

**by**

**ERKAN EMİRZADE**

**In Partial Fulfillment of the Requirements for  
the Degree of Master of Science  
in  
Computer Engineering**

**NICOSIA, 2016**



**Erkan Emirzade: A Computer Aided Diagnosis System for Lung Cancer Detection  
Using SVM**

**Approval of the Graduate School of Applied  
Sciences**

**Prof. Dr. İlkey SALİHOĞLU  
Director**

**We hereby certify that this thesis is satisfactory for the award of the  
Degree of Master of Science in Computer Engineering**

**Examining Committee in charge:**

Prof.Dr. Rahib Abiyev, Committee Chairman, Computer Engineering Department,  
NEU

Assist.Prof.Dr. Ümit İlhan, Committee Member, Computer Engineering  
Department, NEU

Assist.Prof.Dr. Elbrus Imanov, Committee Member, Computer Engineering  
Department, NEU

Assist.Prof.Dr. Kamil Dimililer, Committee Member, Automotive Engineering  
Department, NEU

Assist.Prof.Dr. Boran Şekeroğlu, Supervisor, Committee Member, Computer  
Information Systems Engineering Department, NEU

I hereby declare that all information in this document has been obtained and presented in accordance with academic rules and ethical conduct. I also declare that, as required by these rules and conduct, I have fully cited and referenced all material and results that are not original to this work.

Name, Last Name:

Signature:

Date:

## **ACKNOWLEDGEMENTS**

I would like to thank my supervisor Assist. Prof. Dr. Boran ŞEKEROĞLU who has shown plenty of encouragement, patience, and support as he guided me throughout this endeavor fostering my development as a graduate student and scientist. His important ideas and guidance assisted and help to improve my work.

Also I would like to thank my committee members, for their helpful reviews and comments.

Finally, I would like to thank my family and my best friends for continuous encouragement and positive feedbacks.

This research was generously supported by the Department of Computer Engineering of the Near East University. I am deeply grateful to all supporters.

**To my wife, daughter and my parents...**

## ABSTRACT

Computer aided diagnosis is starting to be implemented broadly in the diagnosis and detection of many varieties of abnormalities acquired during various imaging procedures. The main aim of the CAD systems is to increase the accuracy and decrease the time of diagnoses, while the general achievement for CAD systems are to find the place of nodules and to determine the characteristic features of the nodule. As lung cancer is one of the fatal and leading cancer types, there has been plenty of studies for the usage of the CAD systems to detect lung cancer. Yet, the CAD systems need to be developed a lot in order to identify the different shapes of nodules, lung segmentation and to have higher level of sensitivity, specificity and accuracy. This challenge is the motivation of this study in implementation of CAD system for lung cancer detection. In the study, LIDC database is used which comprises of an image set of lung cancer thoracic documented CT scans. The presented CAD system consists of CT image reading, image pre-processing, segmentation, feature extraction and classification steps. To avoid losing important features, the CT images were read as a raw form in DICOM file format. Then, filtration and enhancement techniques were used as an image processing. Otsu's algorithm, edge detection and morphological operations are applied for the segmentation, following the feature extractions step. Finally, support vector machine with Gaussian RBF is utilized for the classification step which is widely used as a supervised classifier.

**Keywords:** CAD systems; lung cancer; image pre-processing; segmentation; feature extraction; classification; global threshold; support vector machines; SVM; ANN



## ÖZET

Bilgisayar destekli tanı sistemleri, farklı muayene işlemlerinden elde edilen medikal görüntülerdeki çeşitli anomalilere tanı koyma ve ortaya çıkarmada yaygın olarak kullanılmaya başlanmıştır. BDT sistemleri için genel başarı nodüllerin yerlerini bulmak ve nodülün karakteristik özelliklerini belirleme ile ölçülürken, BDT sistemlerinin temel amacı doğruluk oranını artırmak ve tanı süresini azaltmaktır. Akciğer kanseri, en çok görülen ölümcül kanser türlerinden biri olduğu için akciğer kanserini tespit etmek için geliştirilen BDT sistemlerine yönelik birçok çalışma yapılmıştır. Buna rağmen BDT sistemleri, farklı şekillerdeki nodüllerin algılanmasında, akciğer görüntüsünün segmentasyonunda, yüksek düzeyde duyarlılık, özgüllük ve doğruluk değerlerinin elde edilmesinde yetersiz kalıyor.

Akciğer kanseri tespiti için geliştirilen BDT sisteminin yapılan çalışmadaki motivasyonunu bu yetersizlikler oluşturmaktadır. Çalışmada LIDC veri tabanındaki düşük dozda çekilmiş hastaların dokümanite edilmiş göğüs BT görüntüleri kullanılmıştır. Sunulan BDT sistemi; BT görüntüsünden okuma, görüntü önileme, segmentasyon, öznetelik çıkarma ve sınıflandırma adımlarından oluşmuştur. Önemli özneteliklerin kaybını önlemek için, BT görüntülerinin işlenmemiş halleri, yani DICOM, dosya biçiminde okunmuştur. Daha sonra, görüntü geliştirme teknikleri ve görüntü işleme tekniklerinden filtreleme kullanıldı. Segmentasyon için Otsu algoritması, kenar algılama ve morfolojik operasyonlar kullanıldı. Bunları öznetelik çıkartma adımı takip etti. Son olarak, sınıflandırma aşaması için yaygın bir denetimli sınıflandırıcı olarak kullanılan Destek Vektör Makinesinin Gauss RBF kerneli kullanılmıştır.

**Anahtar Kelimeler:** BDT sistemleri; akciğer kanseri; görüntü önileme; segmentasyon; öznetelik çıkarma; sınıflandırma; global threshold; destek vektör makinesi; SVM; ANN

## TABLE OF CONTENTS

<b>ACKNOWLEDGEMENTS</b> .....	i
<b>ABSTRACT</b> .....	iii
<b>ÖZET</b> .....	iv
<b>TABLE OF CONTENTS</b> .....	v
<b>LIST OF TABLES</b> .....	viii
<b>LIST OF FIGURES</b> .....	ix
<b>LIST OF ABBREVIATIONS</b> .....	x
 <b>CHAPTER 1: INTRODUCTION</b>	
1.1 Motivation .....	1
1.2 Objectives .....	2
1.3 Report Structure .....	2
 <b>CHAPTER 2: REVIEW OF CAD SYSTEMS FOR LUNG CANCER</b>	
2.1 Review of CAD Systems .....	4
2.1.1 The challenges for CAD .....	5
2.2 What Is Lung Cancer? .....	6
2.2.1 Non-small cell lung cancer (NSCLC) .....	6
2.2.2 Small cell lung cancer (SCLC) .....	7
2.3 Screening for Lung Cancer .....	7
2.3.1 Chest x-ray .....	8
2.3.2 CT scan .....	8
2.3.3 MRI (magnetic resonance imaging) .....	9
2.3.4 PET scan (positron emission tomography scan) .....	10
2.4 Pitfalls and Challenges of Lung Cancer Imaging .....	11
2.5 Medical Image File Formats .....	12
2.6 Summary .....	14

## CHAPTER 3: MAIN STEPS OF COMPUTER AIDED DIAGNOSIS SYSTEMS

3.1 Image Enhancement .....	16
3.1.1 Histogram .....	16
3.1.2 Median filtering .....	17
3.1.3 Histogram equalization .....	18
3.1.4 Histogram specification .....	19
3.2 Segmentation .....	20
3.2.1 Edge-based segmentation techniques .....	21
3.2.2 Thresholding .....	22
3.2.2.1 Global thresholding.....	23
3.2.3 Region growing .....	23
3.2.4 Morphological image processing.....	24
3.2.4.1 Erosion and dilation .....	24
3.2.4.2 Opening and closing .....	25
3.3 Image Feature Extraction .....	27
3.4 Classification .....	28
3.4.1 Neural networks.....	28
3.4.1.1 Architecture of neural networks.....	28
3.4.1.2 Learning .....	29
3.4.1.3 Training.....	31
3.4.1.4 Testing (classification).....	31
3.4.1.5 Advantages of ANN.....	32
3.4.1.6 Disadvantages of ANN .....	32
3.4.2 Support vector machine (SVM).....	32
3.4.2.1 Defining the SVM classifier formally.....	34
3.4.2.2 RBF kernel .....	36
3.4.2.3 Grid-search and cross-validation .....	37
3.5 Summary.....	39

## **CHAPTER 4: IMPLEMENTATION OF CAD SYSTEM FOR LUNG CANCER DETECTION USING SVM**

4.1 Pre-Processing .....	42
4.2 Segmentation .....	42
4.3 Feature Extraction.....	46
4.4 Classification .....	48
4.4.1 Training phase .....	49
4.4.2 Feature selection phase .....	49
4.4.3 Testing phase .....	51
4.5 Related Works .....	56
<b>CHAPTER 5: CONCLUSION .....</b>	<b>58</b>
<b>REFERENCES .....</b>	<b>59</b>
<b>APPENDIX: Algorithm and Source Codes.....</b>	<b>66</b>

## LIST OF TABLES

<b>Table 4.1:</b> Relevant articles published in IEEE published from 2010 to 2015.....	47
<b>Table 4.2:</b> SVM using RBF kernel with entire dataset.....	54
<b>Table 4.3:</b> SVM using RBF kernel with test dataset .....	54
<b>Table 4.4:</b> SVM using Quadratic kernel with test dataset .....	55
<b>Table 4.5:</b> SVM using Linear kernel with test dataset.....	55
<b>Table 4.6:</b> The effectiveness of the SVM kernels.....	55
<b>Table 4.7:</b> Summary of recent works related with lung cancer .....	57

## LIST OF FIGURES

<b>Figure 2.1:</b> Schema of typical CAD system for lung cancer .....	5
<b>Figure 2.2:</b> Typical x-ray machine .....	8
<b>Figure 2.3:</b> Typical CT machine.....	9
<b>Figure 2.4:</b> Typical MRI machine .....	10
<b>Figure 2.5:</b> Typical PET machine.....	11
<b>Figure 3.1:</b> The main steps of CAD system .....	15
<b>Figure 3.2:</b> Median filtering .....	17
<b>Figure 3.3:</b> Histogram equalization .....	19
<b>Figure 3.4:</b> Combine dilation and erosion to form closing or opening. ....	26
<b>Figure 3.5:</b> Components of a typical ANN.....	30
<b>Figure 3.6:</b> Hyperplanes' separation.....	33
<b>Figure 3.7:</b> Separating hyperplane (in two dimensions).....	35
<b>Figure 3.8:</b> Overfit classifier and better classifier .....	38
<b>Figure 4.1:</b> The flowchart of implementation of CAD system.....	41
<b>Figure 4.2:</b> Lung image in pre-processing stage .....	42
<b>Figure 4.3:</b> Lung image after Otsu's global threshold.....	43
<b>Figure 4.4:</b> Forming a mask of CT lung image .....	44
<b>Figure 4.5:</b> Binary lung image with removed unnecessary perimeter lines .....	45
<b>Figure 4.6:</b> Final stage of CT lung image segmentation .....	46
<b>Figure 4.7:</b> Classification of features (area, perimeter) in training phase .....	50
<b>Figure 4.8:</b> Classification of features (area, eccentricity) in training phase.....	50
<b>Figure 4.9:</b> Classification of features (perimeter, eccentricity) in training phase .....	51
<b>Figure 4.10:</b> Classification of features (area, perimeter) in testing phase .....	52
<b>Figure 4.11:</b> Classification of features (area, eccentricity) in testing phase.....	53
<b>Figure 4.12:</b> Classification of features (perimeter, eccentricity) in testing phase .....	53

## **LIST OF ABBREVIATIONS**

<b>CAD:</b>	Computer Aided Diagnosis
<b>SVM:</b>	Support Vector Machine
<b>ANN:</b>	Artificial Neural Network
<b>RBF:</b>	Radial Basis Function
<b>CT:</b>	Computed Tomography
<b>MRI:</b>	Magnetic resonance imaging
<b>PET:</b>	Positron Emission Tomography
<b>LDCT:</b>	The low-dose computed tomography
<b>NSCLC:</b>	Non-Small Cell Lung Cancer
<b>SCLC:</b>	Small Cell Lung Cancer
<b>DICOM:</b>	Digital Imaging and Communications in Medicine
<b>LIDC:</b>	Lung Image Database Consortium





# **CHAPTER 1**

## **INTRODUCTION**

### **1.1 Motivation**

Cancer is still the major cause of death in the world, further more lung cancer is the most frequently seen type of cancer among others (WHO, 2015). As there is no cancer registry system in TRNC, there is not any official data about cancer statistics. Yet, lung cancer is the leading cancer in males in Turkey as it is in the world male rates (T.C. Sağlık Bakanlığı, 2016). Early diagnosis and proper treatment may pull down the death rates, hence the CAD systems are increasingly becoming the preferred aid in diagnostic procedures by the doctors (Doi, 2007). CAD becomes a significant research topic in the diagnostic radiology and medical imaging. In fact, CAD systems help the doctors in interpreting the images of computed tomography (CT), magnetic resonance imaging (MRI), ultrasound, positron emission tomography (PET), conventional projection radiography as well as all other imaging methods.

Practically, diagnosis process incorporates the assistance of computers from medical imagery, lab work, and electronic medical records and more. When it comes to radiology, CAD is the essential system of procedures in medicine that help doctors in the medical image interpretation. The use of the digital processing and hybrid optical technologies afford the reduction in a processing time as well as enabling more enhancements in specificity and sensitivity. The computer aided diagnosis holds the great potential for the radiology and its utilization is based on its capability to speed up a diagnostic process as well as lessen probable errors.

The concept of the automated diagnosis exists from the year 1960, however the attempts in research and development were failed mostly (Doi, 2007).

Now, there are many institutions all over the world that involved in the development and research of CAD aspects. Day by day, CAD systems are giving more confidence in the medical area therefore CAD systems become a superior method for the cancer detection in interpreting X-ray, CT, MRI and other medical images. Using the outputs of CAD systems

as a reference, helps the doctors not only to accomplish their tasks more accurately and precisely but also in a shorter time.

The CAD systems ensure its reliability and efficiency to the integration of various scientific disciplines such as artificial intelligence, image processing, pattern recognition, etc. Although, CAD systems showed great improvement, it needs much to do in lung segmentation and in different shapes of nodule detection. CAD systems still have more false-positive results than experienced radiologist and have not achieved 100% accuracy, sensitivity and specificity which are very important measurements for the systems (Anshad and Kumar, 2014). This challenge is the motivation of this study in implementation of CAD system for lung cancer detection.

The main purpose of the CAD system is to enhance a diagnostic accuracy as well as radiologist's image interpretation consistency with the help of computer output. This output is highly useful, since the radiologist's diagnosis are based on the subjective judgment. Generally, there are two general approaches that can be applied in the computerized schemes for computer aided diagnosis. First, to identify the lesions location like lung nodule in the chest image by looking isolated abnormal pattern with the computer. Then, the next thing is to measure the features of image of abnormal or/and normal pattern like lung texture concerned to the interstitial infiltrate in vessel sizes and chest image related to the angiograms stenotic lesions.

## **1.2 Objectives**

The main objectives of this study are as follows;

- Exploring the increasing role of image processing and machine learning in computer aided diagnosis systems,
- Designing, implementing, and measuring the performance of computer aided diagnosis system for lung cancer detection using SVM,
- To contribute to the studies on computer aided diagnosis system for lung cancer detection.

## **1.3 Report Structure**

The rest of this report is organized as follows;

- Chapter 2 provides a brief review and the challenges of CAD systems. The chapter introduces the history of CAD systems and discusses the two notions of CAD; 1. Taking place of radiologists, 2. Taking computer outputs as a second opinion. Then follows with the description of lung cancer, the types of lung cancer and technical knowledge on imaging technologies that are used for diagnosis and treatment are reviewed.
- Chapter 3 introduces the main steps of CAD system. The chapter starts with the role of image enhancement and describes image enhancement techniques which are used in the implementation. Following this, the chapter moves to the details of commonly used segmentation techniques in image processing. Also, the importance of feature extraction and classification are described. In this chapter, an important portion is reserved to the explanation of artificial neural networks and support vector machines.
- Chapter 4 provides the details of design steps and illustrations of the implemented computer aided diagnosis system. Recent studies about lung cancer detection are summarized for comparison. Lastly, test results and statistical performance measurements of classification based on the features are presented.
- Finally Chapter 5 introduces the conclusions of the study and suggestions for the future work.

## **CHAPTER 2**

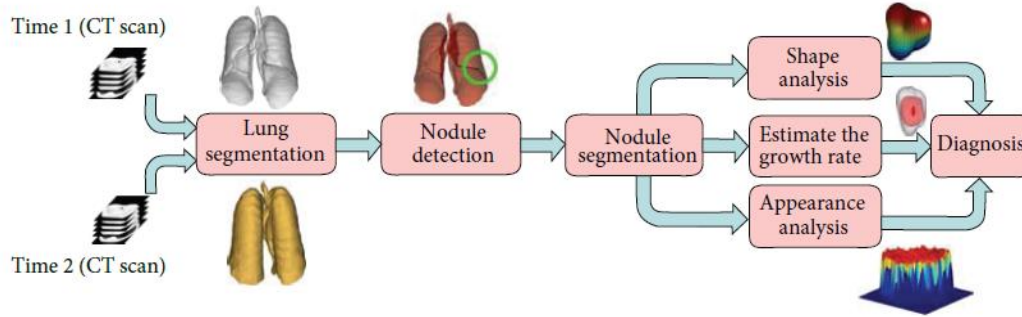
### **REVIEW OF CAD SYSTEMS FOR LUNG CANCER**

The chapter begins by the review and challenges of Computer Aided and Diagnosis system. Then, will introduce the brief description of lung cancer and the types of lung cancer. Finally, the chapter will present technical knowledge on imaging technologies which are using for diagnosis and treatment purposes.

#### **2.1 Review of CAD Systems**

The first steps of Computer Aided and Diagnosis (CAD) system go back to the years of 1950s. At the beginning, there was a thought that the CAD systems may take place of the radiologists in determining anomalous tissues, since computers perform given tasks much better than human. However, this thought could not find ground because of the ineligible computer sources and unsophisticated image-processing methods and techniques. In 1980s, systematic and large scale development and research about various CAD systems was started. The notion of taking computer outputs as a second opinion rather than taking place of radiologists gain more support, thus form the notion of computer aided diagnosis. The performance of automated diagnosis systems expected to be superior to physicians whereas the performance of CAD systems which is accepted as supplementary to the physicians do not have to be ahead of them (Doi, 2007). It took more than 30 years for the researches to become solid and in 1998 Food and Drug Administration approved ImageCheckerTM which was the first approved CAD system in mammography (Fujitaa et al., 2008). Studies keep going about both notions which are CAD and automated computer diagnosis. Although, the performance of CAD systems is not superior to the performance of physicians, many clinics in the Europe and USA benefit the aid of CAD systems in the early detection of breast cancer as a second opinion (Doi, 2007). Now, CAD has become part of the common clinical procedure for breast cancer detection in some clinics, thus CAD turns out to be the leading research area in machine learning of medical diagnosis and imaging (Lee et al., 2009).

A characteristic schema of CAD system for lung cancer is depicted in Figure 2.1 (El-Baz et al., 2013: p. 2)



**Figure 2.1:** Schema of typical CAD system for lung cancer

Several studies and findings on CAD systems point out that, by using CAD systems radiologists can get better results in nodule detection. And CAD does not only improve the performance and efficiency of the radiologists, but also improves the precision and accuracy of nodule detection rate as well (Fujitaa et al., 2008). Therefore, guidance and aid in nodule detection is the main objective of using CAD systems, which helps to increase the consistency and accuracy of radiologist's diagnosis. It is well known that anatomic structure of lungs makes nodule detection very hard mission thus radiologists may miss about 35% of the nodules on chest images (Anshad and Kumar, 2014). Whereas, CAD systems have still more false-positive thus sensitivity of nodule detection is low compared to experienced radiologists (Ginneken et al., 2011). It is understood that CAD systems advances the performance of radiologists in diagnosing process and nodule detection, however it does not become common in clinical usage. Regular usage in the clinics depends on meeting the following requirements: ensure in high sensitivity, decrease in the number of false positives, increase in the diagnosis speed, advance in automation, decrease in running cost, advance in the capability of detection of various shapes and types of nodules (Firmino et al., 2014).

### 2.1.1 The challenges for CAD

Some challenges for CAD system researchers and developers are as follows;

- Not easy to take needs and translate knowledge from radiologists (Ginneken et al., 2011).
- Radiologists are the critic stakeholder in CAD system implementation, they should point out the application fields and should provide sample databases with the ground truths for training and testing (Ginneken et al., 2011).
- The lack of homogeneity in the lung structure and affinity densities in veins, arteries, bronchi and bronchioles makes segmentation a very hard issue (El-Baz et al., 2013).
- The level of automation, the quickness of the system, the success rate of varied shape and size of nodule detection (El-Baz et al., 2013).
- For efficient validation of the CAD systems larger databases should be provided (El-Baz et al., 2013).
- There is a need for more advance techniques or improvement of the present techniques in segmentation of lungs in order to detect nodules less than or equal to 3 and ground-glass opacity (El-Baz et al., 2013).

## **2.2 What Is Lung Cancer?**

The unrestrained expansion of abnormal cells in lungs causes lung cancer. These abnormal cells disturb the smooth functioning and development of lung tissues. If this condition is left untreated, abnormal cells grow and form tumors and ultimately damage the lungs that are proving oxygen to the whole body via blood.

Two main type of lung cancers are non-small cell and small cell. Because of the large size of lungs, nodules can grow for a time untill detecting them (Lung Cancer, 2016).

### **2.2.1 Non-small cell lung cancer (NSCLC)**

The NSCLC is the commonly prevailing form of lung cancer, also according to American Cancer Society, NSCLC is responsible for 85 percent of the total lung cancers in America (American Cancer Society, 2016). The common tumors of lung cancers are following;

- Adenocarcinoma is the lung cancer in non-smokers, and equally found among men and women.

- Squamous cell carcinoma or epidermoid carcinoma is the lung cancer that is positively correlated with the tobacco smokers. This tumor is formed mainly in at the center of large bronchi. Males are more vulnerable to this type of tumor.
- Large cell carcinomas are the tumor cells that have comparatively larger size with excess cytoplasm. Unlike adenocarcinomas and epidermoid, these cells lack microscopic characteristics.

### **2.2.2 Small cell lung cancer (SCLC)**

The remaining 15 percent contribution in lung cancers is of SCLC. Tobacco smoking is the leading cause of SCLC and gets birth quickly as compared to NSCLC. In the body, this type of cancer is relatively spread quickly, higher growth rate and shorter multiplying time. Chemotherapy is a more effective treatment for the SCLC.

### **2.3 Screening for Lung Cancer**

The first step to diagnosing the lung cancer is the identification of symptoms. Symptoms are largely showing the damage to lungs and their functionality. Chest pain and cough are the most common symptoms of lung cancer. The cough gets worst on each passing day and also increases chest pain. Besides these, breath shortness, feeling weak, weight loss, blood in cough and fatigue are also commonly appeared symptoms among the lung cancer patients. Unfortunately, the scientific community has not developed any screening tool that could identify the lung cancer at early stage. Chest X-rays are commonly available tools for the screening, but they are not reliable enough yet. The development of screening tool is the necessity of time as many researchers have concluded that early-stage tumors are easy to cure. The low-dose computed tomography (LDCT) is recommended screening on the annual basis to smokers and those who quit smoking with last 15 years. According to American Society of Clinical Oncologists, people who are in the age group of 55-74 and smoked more than 30 years are at more risk of lung cancer. In addition to LDCT, following are the some imaging technologies that are used for diagnosis and treatment.

### 2.3.1 Chest x-ray

X-ray machine discharges radiation that goes into the body and imaging picture of the organs on the film. To diagnose lung cancer, x-ray imaging is used as step that helps in the identification of lung tumors. As mentioned, x-rays are not the final authority because they are unable to differentiate between the cancer and other lung diseases.



**Figure 2.2:** Typical x-ray machine, (Frederick Memorial Hospital, 2016)

### 2.3.2 CT scan

CT scan stands for computed tomography, and it is an extended version of X-ray in which computer is attached to the X-ray machine. Pictures that are taken from taken angles and distances are processed in the computer and presented in the 3-dimensionsal, cross-sectional (tomographic) and in slices form. In this way, bones, tissues, blood vessels, and organs are shown up clearly. The imaging of CT scan is useful for diagnosis, treatment and progress of medication. Recently, helical or multi-slice scanning is introduced that almost eliminated gaps in the collection of slides.





**Figure 2.3:** Typical CT machine, (Frederick Memorial Hospital, 2016)

### **2.3.3 MRI (magnetic resonance imaging)**

It is imaging technique in which radio waves and strong magnetic fields are used by the scanners to form the inside images of the body. The powerful magnet is aligned with the nuclei of atoms, and then magnetic field triggers atom to resonate. In this way, nuclei generate its own magnetic field and then the field is detected by the scanner for creating an image (Biederer et al., 2012). The advancement in technology has helped to take detailed pictures from different angles. MRI is particularly helpful when there is a need for the identification of soft tissues. This is the reason; this technique has more reliability.



**Figure 2.4:** Typical MRI machine, (Frederick Memorial Hospital, 2016)

#### **2.3.4 PET scan (positron emission tomography scan)**

If X-ray or CT scan diagnoses or doctor predicts any chances of lung cancer, PET scan is suggested for detailed results. In this imaging technique, tracer or radioactive glucose is injected and then scanners are rotating to take pictures which tissues or organs used tracer (Mac Manus et al., 2003). When malignant tumor cells use glucose, they are showing up brighter and more active in images. The integration of PET-CT scan is very useful for detecting the cancer. The CT scan gives a detailed view of tissues and organs, and PET gives pictures of abnormal activities and active cells. Researchers also concluded that PET-CT scan are producing more accurate results as compared to PET or CT scan alone.



**Figure 2.5:** Typical PET machine, (Frederick Memorial Hospital, 2016)

#### **2.4 Pitfalls and Challenges of Lung Cancer Imaging**

During the years 1985 to 1995 more than 100,000 medicolegal cases gathered which included all the medical properties. The collected cases showed that 90% of the professed errors happened on x-ray imaging, while the rest of the errors occurred respectively on CT scans and other studies (Mohammed et al., 2005). CT is more advance method in lung nodule detection with respect to analog radiography in which CT rate is 2.6 to 10 times higher than x-ray imaging (Firmino et al., 2014). Whereas, the increased number of CT image slices regards to x-ray per patient straightly increases the radiologist's amount of work in diagnosis (Shao et al., 2012). The increase in work per patient most likely effects the results negatively by increasing the errors in nodule detection and nodule characteristic determination. Thus, over diagnosis may cause loosing the focus whereas under diagnosis may cause missing lung cancer (Mohammed et al., 2005). The challenges of MRI of the lung is the limitations of low density of proton in the lung and signal corruption because of the sensitivity in air-tissue interfaces. Moreover, the quality of lung imaging in MRI depends on the patients' capability of attending the breath hold instructions compared with CT and x-ray imaging. MRI shows up worthy lung imaging technique, together with CT and x-ray. However, lung imaging with MRI still challenges being the most detailed and

devoid of any radiation load to the patient yet the least robust and most expensive of the three techniques (Biederer et al., 2012). Though it is not preferred as a main imaging technique for the staging and diagnosis of lung cancer, MRI has some superiority over other imaging techniques like distinguishing mass tissues from adjacent tissues of lungs. The lack of access and experience in MRI method are presumably the main hurdle of MRI to be a routine for lung cancer patients (Hochhegger et al., 2011). PET is the superior method in examining functional details of tissues whether it is benign or malignant whereas it is not precise on structural details (Sroubek et al., 2009). Those challenges demonstrate that each imaging method has their pros and cons, thus have complementary role in providing details of patient anatomy (Yelleswarapu et al., 2014). Thus, the complementary role and effects of utilization both CT and PET evolve PET/CT imaging and take place of PET-only scan imaging machines (Moon S-H et al., 2016).

## **2.5 Medical Image File Formats**

The development in medical instruments, provokes the hospitals to equip with different kind of digital imaging devices as well as image archiving and transmission servers called PACS. The standard in transmission and image data formats of different kind of imaging devices become more important since the PACS servers become a part of modern hospitals (Chen, 2012). The file formats of medical images which is created by medical devices divided into two broad categories that aim to standardize the format of images and to simplify and make stronger for analysis. The files of medical images consist of image data and metadata. Some file formats store image data and metadata in one file while some file formats store image data in one file and metadata on the other file. The first attempt of forming standardized medical imaging file was Interfile in 1980. Interfile format consists of two files, one with the image data and the other containing metadata. At the end of 1980s, Analyze 7.5 file format was created. It consists of a pair of files that one file contains image with extension “.img” and the other file with extension “.hdr” which contains metadata. By the beginning of 2000, the committee of National Institute of Health created Nifti file format to improve the frailty of Analyze format. Nifti file format authorizes the storage of image data and metadata in different files as well as a single file with extension “.nii”. National Electric Manufacturers Association and American College of

Radiology create the Dicom file format in 1993. It is accepted as the backbone of medical imaging divisions. Dicom file format is beyond being only the file format but existing as network communication protocol. The dicom file format merged image data and metadata into one file. Dicom file format supports non-dicom formatted files such as JPEG, RLE, JPEG-LS, JPEG-2000, MPEG2/MPEG4, etc. in order to enclose them in a Dicom file (Larobina and Murino, 2014). In spite of its broad acceptance in medical technique, large file sizes and a need of special software for viewing are the two disadvantages of DICOM format. Unlike to DICOM format, JPEG, TIFF and PNG image file formats allows users to view images without any special viewer software, thus they are more popular in daily life (Varma, 2012). JPEG compression algorithm was created in 1990s and become very popular because of its power in compressing larger files into smaller and allowing to transfer easily over internet. Using lossy compression technique may compress files more but causes loss of data while lossless compression keeps image quality high but costs a large file in size. And there is a possibility of distortion in radiologic images. In 1986, TIFF was specially developed to be compatible with varied devices. The full range in image colour depths, resolutions and size support is the the strength of TIFF whereas keeping high quality image causes files to be large in size which limits the portability. In 1995, file format PNG was developed by the objectives of flexibility in allowing lossless data compression, variable transparency and consistency in brightness. Quality of the image is not affected much because of the lossless data compression but this tends to have files in larger size with respect to JPEG (Wiggins, et al., 2001). It is seen that, each image file format has its own benefits and limitations thus using proper image file format gain more importance. The file formats in image domain frequently evolving. Having immense flexibility in complying development in technology, medical imaging and network infrastructure makes DICOM image file format as the preferable standard with regards to other image formats created by diagnostic techniques (Larobina and Murino, 2014). To be more convenient and practical in diagnostic procedures CAD systems should integrate with workflow of the clinics. That is, CAD systems will integrate the image archiving and transmission servers called PACS in daily workflow (Li and Nishikawa, 2015).

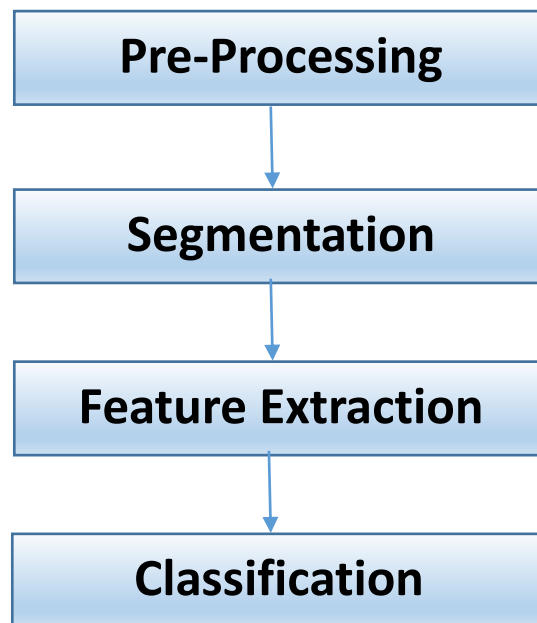
## **2.6 Summary**

The chapter give a start with discussion of two thoughts about the CAD systems. At the beginning, there was a thought that the CAD systems may take place of the radiologists in determining anomalous tissues, since computers perform given tasks much better than human. However, this thought could not find ground because of the ineligible computer sources and unsophisticated image-processing methods and techniques. The notion of taking computer outputs as a second opinion rather than taking place of radiologists gain more support, thus form the notion of computer aided diagnosis. Now, CAD has become part of the common clinical procedure for breast cancer detection in some clinics, thus CAD turns out to be the leading research area in machine learning of medical diagnosis and imaging. The pitfalls and challenges of lung cancer imaging is mentioned in imaging modality. Those challenges demonstrates that each imaging method has their pros and cons, thus have complementary role in providing details of patient anatomy.

## **CHAPTER 3**

### **MAIN STEPS OF COMPUTER AIDED DIAGNOSIS SYSTEMS**

This chapter will provide theoretical information about the methods applied at the implementation of CAD systems. The chapter starts with the role of image enhancement and describes applied image enhancement techniques for pre-processing, then moves to the details of commonly used segmentation techniques in image processing. Finally, importance of feature extraction and classification are described. The chapter will flow as Figure 3.1 below;



**Figure 3.1:** The main steps of CAD system

### **3.1 Image Enhancement**

Interference and other phenomena affect the quality of the medical images, which are caused by noise. This affects the process of measurement in imaging and the systems of data acquisition. Some improvements in appearance and visual quality of the images may assist in their interpretation by a medical specialist (Bankman, 2000: p. 1).

The purpose of image enhancement is to sharpen the features without changing the natural structure of the image, to be used in image analysis such as edge and boundary detection (Patel and Goswami, 2014).

Enhancement of the image is principally performed to ensure that the quality of the image is better than the original image in a subjective manner. Therefore, a method which is used for x-ray chest enhancement may not be very suitable for MRI brain image enhancement (Petrou and Petrou, 2010).

There are two broad categories of image enhancement, which can be named as the frequency and the spatial domain methods. Processing images based on the frequency approach is based on the transformation of the image based on the Fourier transformation. The image plain is referred to in the spatial approach, wherein the pixels of the image are directly manipulated. It is not unusual for practitioners to use a combination of these approaches in everyday usage. The evaluation of an image based on the visual interpretation of the quality of the image is a highly subjective process and therefore, a good image is a purely subjective term and these standards of evaluation are not always comparable to the performance of an algorithm (Gonzalez and Woods, 2002).

#### **3.1.1 Histogram**

Representation of an image in a histogram is a graphical presentation of the intensities and the intensity values of the probabilistic occurrences in the image (Patel and Goswami, 2014). Histograms show number of pixel values for each intensity value of an image and give detailed intensities of the image thus helps to realize the brightness and contrast of the image. The components of the graph are in the lower side when the picture is dark, pertaining to the scale of intensity. The components of the graph are on the higher side, similarly, when the image is bright or light in nature, pertaining to the scale of the intensity. Manipulation of the histogram can enhance the image quality. Apart from providing useful statistics of the image, this graphical representation can also be used for

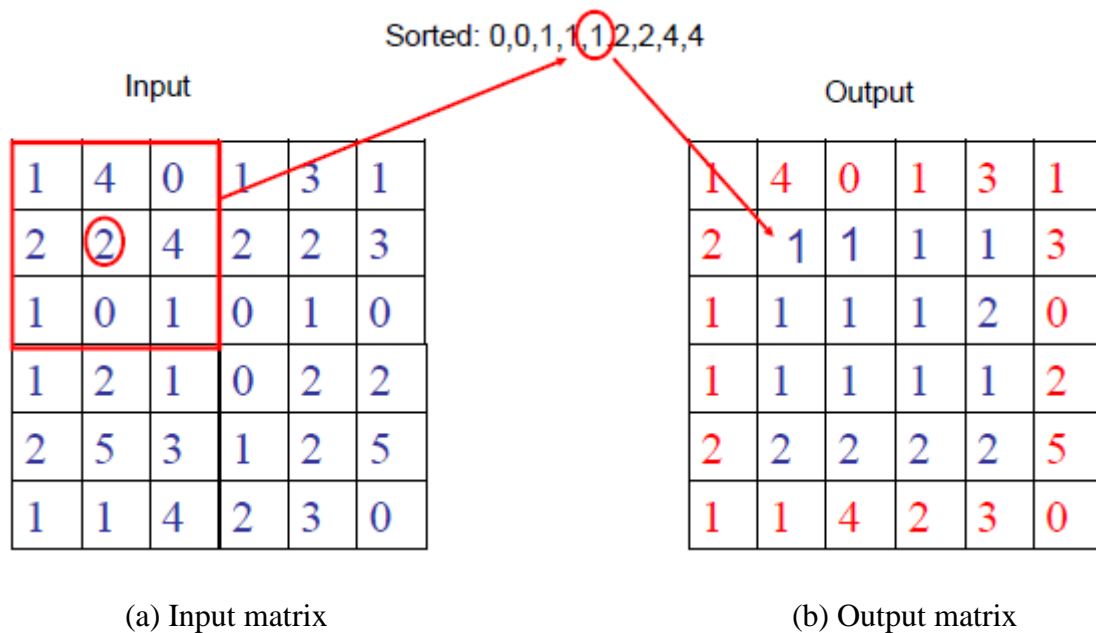


segmentation and image compression, which are other image processing applications in medical imaging.

### 3.1.2 Median filtering

In general, image filters categorized as linear or non-linear. Primary benefits of median filtering that is non-linear are the speed and the simplicity in usage. Having the property of adaptive usage and well edge conservation makes median filtering suitable for lots of application in image processing field that linear filters mostly fail (Pitas and Venetsanopoulos, 1990: p. 63). Median filtering is also powerful in eliminating particular sort of noise called 'salt and pepper'. If a pixel comprise an intense value, then the value alters with the median value within the neighbourhood (Russ, 2011: p. 217).

2D Median filtering example using 3x3 sampling windows which keeps edge values unchanged.



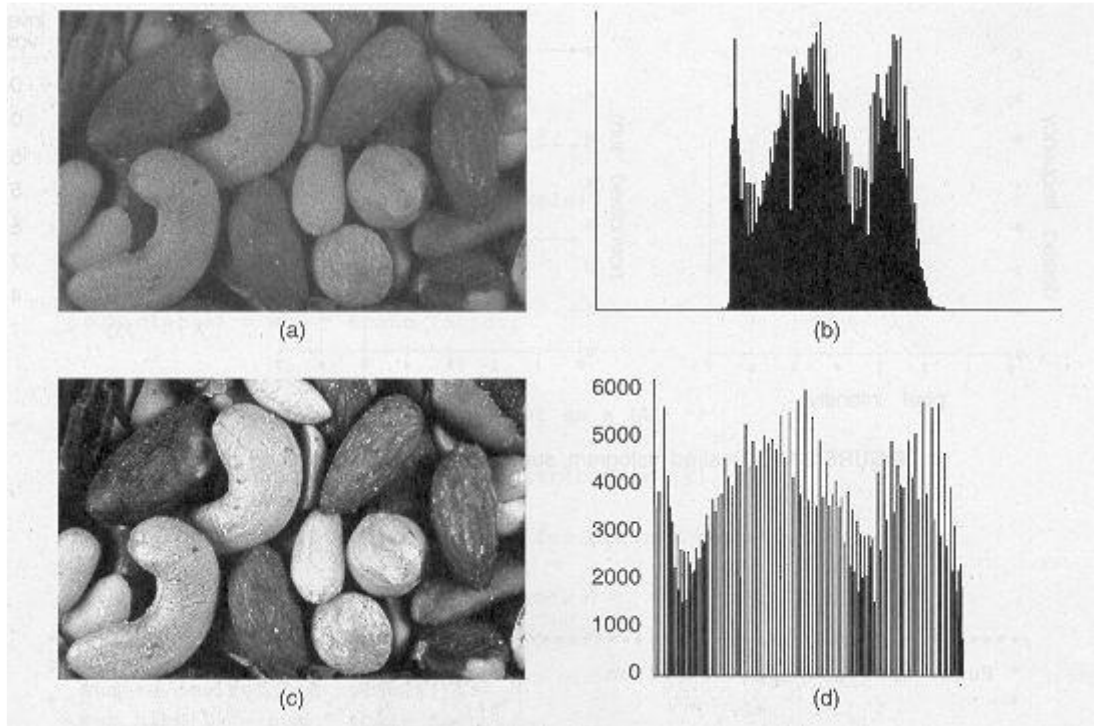
**Figure 3.2:** Median filtering

### 3.1.3 Histogram equalization

This is a technique by which the image intensities can be adjusted so that the contrast may be enhanced in the image (Senthilkumaran and Thimmiaraja, 2014). This can be applied to the entire image or on selected parts of the image to reduce or redistribute the intensity of distributions. In case the original histogram of the image has peaks and troughs, the same will also be there after the histogram equalization, but the same will be shifted.

Since histogram equalization is based by principle in being a point process, the image does not receive new intensities but new values will map to the existing values, but the actual number of points of intensity in the image will be equal or will be less in the processed image. Counting the pixel value of each point in the image is the starting point of the process and this can start with a zero valued array. 256 is the array size for 8 pixels (0-255). The image can be parsed and each array element can be incremented with respect to each processed pixel. In the second step, the histogram values will have to be stored in a separate array. Element 1, this array holds the values of the 1 and 0, the histogram elements. For example, the 255 element would contain the sum of the elements 255 to zero, counting in the reverse. By multiplying each element, the array can then be normalized, by the number of pixels and the maximum pixel value. For a 512 x 512, 8 bit image,  $255/262144$  would be that constant.

The completion of the Step two detailed above provides a Look Up table that may be used for the input image transformation. Steps 2 and 3 of this process are shown in Figure 3.3 from the Figure 3.3a, the normalized sum can be used to look up for the values by rounding off to the nearest whole number. Like number matching can be performed through the lookup table obtained through Step 2 (Luong, 2006: p. 22-23).



(a) Original image  
(c) Equalized image

(b) Histogram of original image  
(d) Histogram of equalized image

**Figure 3.3:** Histogram equalization

### 3.1.4 Histogram specification

This is a useful method to enhance the image. In contrast to the Histogram equalization method, this method transforms the image's histogram into histogram of a specific type to enhance the available gray scale ranges. While the histogram equalization technique creates an equated or uniform histogram, this may not be what is needed to enhance a certain image, which may be lacking in contrast or dynamic range. In this case, the histogram specification is useful wherein this method transforms the image's histogram into histogram of a specific type to enhance the available gray scale ranges (Chi-Chia et al., 2005). This technique can be performed by comparing a desirable histogram and the histogram of the input image and can be performed through the steps detailed below:

1. Equalize the histogram of the original image
2. Inversed histogram equalization is to be performed on the image that is equalized

The lookup table needs to be generated for the inverse histogram equalization which corresponds with the histogram that is the reference for this process and then the lookup table can be inverse transformed. The lookup table outputs are used to compute the lookup table inverse transform through analysis. For a particular input, the closest output becomes the inverse value (Luong, 2006: p. 23-24).

### **3.2 Segmentation**

The segmentation process aims to achieve the partitioning of the image into different regions which are called subsets or classes and these subsets are homogenous within themselves for one or more characteristics (Bankman, 2000).

In several image processing applications and computer vision, the segmentation is a very important process to be performed. This is because of the fact that it is the initiation of the image processing into more complicate approaches (Lucchese and Mitra, 2001).

Feature extraction, image display and image measurements in medical imaging predominantly use the segmentation operations. Classifying image pixels into regions of anatomy may be useful in some applications like blood vessels, muscles and bones. In some other applications, pathological regions may have to be separated, for example based on tissue deformation, cancer and lesions of multiple sclerosis. The division of the entire image into gray, white and spaces in the brain for cerebrospinal fluids may have to be performed in certain applications. In some cases, a specific feature has to be separated from an image, like for example, the cancerous growth from a CT image of the lungs (Memon et al., 2006).

Since the boundaries between the segments are not very clear, there is no universal rule for dividing these sections in the lungs and manual annotation of the segments of the pulmonary region may be difficult. Segmentation remains a problem in the absence of a universal rule, even though there are several methods proposed and used (Mesanovic et al., 2012).

There is a number of segmentation approaches proposed and yet there is no standard technique which can be used universally across all application. The type of the image data and the objective of the study of the image form the definition of the goal of segmentation.

Different algorithms result from the differences in the assumptions based on which the study is performed.

Segmentation can be classified through a variety of approaches (Bankman, 2009: p. 73) and these are listed as follows:

- Automatic, semi-automatic and manual.
- Local methods or pixel based methods and global methods or region based methods
- Thresholding, region growing, etc. (low level segmentation), manual delineation and multispectral or feature map techniques, contour following, dynamic programming, etc. (model based segmentation).
- Neural network techniques, statistical techniques, fuzzy techniques or classical techniques (edge-based, thresholding and region-based techniques).

The following are the most commonly used segmentation techniques (Bankman, 2009: p. 74):

1. Looking for the regions that satisfy a criterion of homogeneity or region segmentation
2. Looking for edges between regions with different characteristics or edge based segmentation

Thresholding is a low level segmentation method, wherein the selected threshold is used to divide the image into pixel groups having greater or lower values than the threshold values. Global methods which are based on histograms of the gray scale or based on local properties or selection of local thresholds and dynamic approaches can be used in the thresholding. Algorithms of clustering can achieve region segmentation through dividing the image into pixel clusters which have high resemblance in the feature space. Each pixel has to be studied and assigned to the best representative cluster that suits the pixel. Region segmentation can also be done using the technique of region growing, wherein algorithms assign adjacent pixels or spaces to a region based on the similarity of the previous or accompanying pixel closeness.

### **3.2.1 Edge-based segmentation techniques**

These techniques benefit from the changes in the gray tones in the images and transform the images into edge images. Edges represent a flaw or a lack of contiguity in the image.

Because the image is studied and edges drawn out, there is no alteration to the original image. Numerous parts of the image of different color levels form the objects of the image. There are varying gray levels in an image. There is a significant change in the intensity of the image in the edges and this can be associated with discontinuity in the feature of what is being studied through the image. Step edge can happen, where there is an abrupt change in the intensity of the image from one value to an entirely opposite value. Line edges can happen where the intensity value changes abruptly, but returns to the original value in a short space within the image. Due to the low frequency equipment or the smoothing that is caused through the sensors, there is a low probability of step and line edges in reality. Line edges can be transformed to roof edges and step edges can be transformed to ramp edges when the change in intensity is not rapid but occurs gradually (Senthilkumaran and Rajesh, 2009).

The principal disadvantage of edge based segmentation techniques is that the edges often do not enclose the object and therefore, a linking of the edges needs to be done to identify the specific regions that form contiguous segments in the image. Edge linking is the simplest technique, wherein small areas of the image are studied and linked based on the magnitude of the edge or the direction of the edge. This is an unreliable technique and it is often very costly (Bankman, 2009).

### **3.2.2 Thresholding**

The gray levels of the pixels of the background and the gray levels of the pixels of the object are totally different in many applications of pattern recognition and image processing. To separate the object from the background, thresholding is the easiest and simplest approach (Hossain et al., 2011).

Thresholding converts a multiple level image into a binary image. The value of 1 is assigned to the object and 0 assigned to the background pixels. The assignment is based on  $T$  – a threshold value, which may be the value of a color or an intensity that is desired to be separated. This technique is called global thresholding if the value  $T$  is held constant; else it is called local thresholding (Varshney et al., 2009). Global thresholding is an intuitive approach, where one threshold value is selected and applied to the entire image being processed and therefore, the thresholding is stated to be global in nature.

### 3.2.2.1 Global thresholding

This technique works on the basic assumption that there is a bimodal histogram in the image hence, the desired object being looked for can be separated by comparing a universal T value with each pixel of the image (Bankman, 2009: p. 74). Consider that the image  $f(x,y)$  is to be evaluated. The background pixels and the pixels of the object gave different gray levels which are grouped into two predominant gray scale levels. The simplest way to try to extract the object would be to identify and use the T value that separates the modes distinctly. Image  $g(x,y)$  representing the threshold is defined as

$$g(x,y) = \begin{cases} 1 & \text{if } f(x,y) \geq T \\ 0 & \text{if } f(x,y) < T \end{cases} \quad (3.1)$$

Thresholding converts a multiple level image into a binary image. The value of 1 is assigned to the object and 0 assigned to the background pixels. Converting a gray scale image into a black and white image through the binarization process has a number of applications, with the predominant purpose being the separation of the dark background and the light object and vice versa making to image to possess black pixel regions and white pixel regions and these regions are called connected components. When the white connected components are counted, target objects which are bright can be counted. By observing the connected component, it is also possible to understand the size and shape of the object being studied.

### 3.2.3 Region growing

While individual pixel intensities is the focus of the thresholding technique, this technique focuses on groups of pixels for similar intensities.

The first step is to find a starting point for each segment needed in the form of a seed pixel. Based on a predetermined formula, the pixels with the same property are merged together with the seed pixel domain. The new pixels so formed are used as new seed pixels as long as there are no more similar pixels to be merged and by then the region would have grown to represent a segment in the image (Tang, 2010).

In practice, there are a few questions that need to be answered when using this technique:

1. How to identify the seed pixels?
2. How to determine the formula for grouping the pixels?
3. How to determine when to stop the growth process?

The ease of computation and completion is the advantage of the region growing algorithm. As with the thresholding method, this technique is also very rarely used alone, but in conjunction with other advanced techniques. The disadvantages of the region growing techniques are (Tang, 2010).

1. Human intervention is needed to provide a seed point for every region which needs to be extracted
2. The region development is prone to being affected by noise where there are empty spaces of links to other local regions which could lead to error.

### **3.2.4 Morphological image processing**

Due to the nonlinear nature and the strong mathematical base, morphological image processing has been a highly preferred image processing technique. Image enhancement, shape analysis, segmentation, image analysis, computer vision problems, texture analysis, etc. are just some of the major applications of this technique (Chen and Haralick, 1995; Demin et al., 1995). The choice of the size and shape of the structuring elements determine the success of the application of this technique (Chen and Haralick, 1995).

#### **3.2.4.1 Erosion and dilation**

Binary images are processed through a number of techniques which are referred in total as morphological techniques, which includes the dilation and erosion of the binary images and this process is called the morphological procedure. The pixel values in the black and white binary images are just 0 and 1, making the process simple based on counting and sorting of pixels as per predetermined rules depending on the characteristics of the neighboring pixels. The original image is used to analyze the value of each pixel in the image and in practice, the original image is replaced after every few lines of binary codes are assigned back to the original image. The pixel values are used to evaluate the rest of the image and not the entire image is used up for the process (Russ, 2011: p. 468).

Erosion is the transformation of morphology which uses the vector subtraction of the set elements to form the dual to dilation (Haralick et al., 1987). This works by the process of



switching OFF pixels in a feature that were originally ON and this is aimed at removing the pixels that do not fall in the range of interest of the study. The simplest approach is to select a threshold value and switching off pixels which are different from this threshold value and thereby isolating the values or pixels which conform to the object being studied. The advantage is that when there is bright point in a dark region, simple enumeration can isolate this pixel, while this technique uses the average brightness of the region to determine whether a pixel is to remain ON or is to be switched OFF.

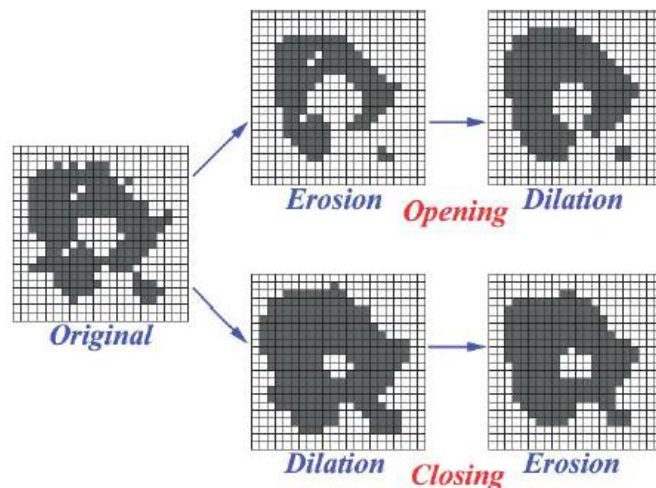
Since such pixels cannot be read through simple thresholding techniques as their gray scale property is similar to the regions of interest and therefore, they cannot be removed by Boolean logic. Even if one criterion is the gradient and the other one is the gray scale, analyzing using such a scale also provides extraneous pixels which will affect the quality of the study. A simple approach can be to remove any pixel that is in contrast and is touching the background. Performing this can remove a layer of pixels around the object being studied and this can have effects. It can cause dimensional shrinkage of the object or even cause a feature to be broken into pieces.

Instead of removing pixels, the operation to add pixels to the image is called dilation in the morphology universe of image processing (Haralick et al., 1987). The rule of dilation is analogous to the rule of erosion in that it adds background pixels which touch another pixel in the foreground or the object of the analysis. This operation adds a layer of pixels around the object of the study, which could cause the object to appear to be expanded and sometimes merge the features of the different segments. Since the processes of erosion and dilation respectively subtract or add to the foreground, they are also called plating and etching or growing and shrinking. Specific sets of rules of various kinds are available to determine which pixels to merge or add or remove in the formation of a combination of dilation and erosion.

#### **3.2.4.2 Opening and closing**

When an erosion process is followed by a dilation, it is called an opening since this combination of actions has the potential to create gaps between features which are touching (just touching), as indicated in the Figure 3.4 isolating pixel noise and removing fine lines in binary images is usually performed through the process of opening. When the same operations are performed in the opposite order it can produce a different result. This

is called closing, since it can close openings in the features of the images. Adjustment of the erosion and dilation can be performed through a number of parameters, especially pertaining to the adding or removing neighboring pixels and how many times the operations have to be performed. Mostly, the erosion and dilation sequences are kept in equal numbers, but this again is dependent on the nature and composition of the image. Opening of separate touching features can sometimes be done using multiple erosion operations till the point where the features have been separated and yet stopping short of completely erasing the pixels. When the separation is completed, the dilation process can be used to grow back the object to the original size and this has to be stopped before they eventually merge to negate the impact of the erosion process and this is accomplished by the rule which stops the closing process from merging a pixel that has been separated so that the separation is maintained. Additional logic is needed at each stage of the operation, since feature identification of the pixels must be performed at each stage of dilation. An additional rule also has to prevent the increase in the size of the object beyond the original size. If there are multiple objects in the image and with different sizes, this technique can make some features disappear for the sake of identifying some other features (Russ, 2011: p. 471).



**Figure 3.4:** Combine dilation and erosion to form closing or opening

In summation, it can be stated that dilation expands objects while erosion shrinks them. Opening operations can suppress islands, cuts out narrow isthmuses and capes smaller than

the structuring element used, while the closing operations fills the thin gulf and the small holes (Demin et al., 1995).

### **3.3 Image Feature Extraction**

An attribute of the image or a primitive characteristic of the image is called an image feature. These are important in separating different regions within an image or the identification of specific properties within the image (Pratt, 2007).

The individual features in images on which the analysis can be performed can be location (absolute position and relative position with respect to the other features), brightness (brightness parameters such as color density and color values), shape and size. Within each of these features, different measurements can be made of a pinpointed nature through a number of different means. A few measures each class is provided by default by every image analysis system and therefore, the most important aspect in image analysis is the understanding of which features are the most important to solve the problem at hand (Russ, 2011: p. 547).

A feature set that is good contains distinguishable elements that can help in clearly identifying and segregating the features. Objects in the same class must not be categorized as different and therefore, the system of understanding must be robust in nature. A small set whose values can accurately represent the different features of the image must be developed (Mesanovic et al., 2012). This is important in performing image classification of any kind, as it can affect the results of the classification significantly (Reza et al., 2011).

The most crucial step is the identification of the meaningful features in the image, because (Mesanovic et al., 2012):

1. From the initial set, all the possible subsets must be found and this is a laborious task.
2. Some of the discriminations apply to at least some of the subsets.
3. Variations between intra and inter class features within the image is narrow.
4. The inclusion of more and more features can reduce the utility of the model.

### **3.4 Classification**

The purpose of the classification is to label similar pixels in an image thus form several classes according to their similarities. Classification which is the significant component of image analysis, forms the classes by analyzing the image features. Two major type of classifications are; unsupervised classification and supervised classification (Thomas and Kumar, 2014). The following classifiers have plenty of implementation areas including medical image classification like lung cancer detection.

#### **3.4.1 Neural networks**

Learning systems that are biological in nature in the form of complex interconnected neurons have inspired the study of artificial neural networks (ANNs). In laymen's terms, artificial neural networks are composed of simple units that are densely interconnected in sets, producing a single real-valued output while taking many real-valued inputs (Mitchell, 1997).

Such neural networks can be used to segregate remote sensor data of various types and have also been produced better and more accurate results as compared to the statistical methods that are used commonly.

The two advantage of neural networks stated below contribute to their success (Jensen et al., 2001):

1. No need for normal distribution
2. Simulation of non-linear and complex patterns in an adaptive manner within the given structures of topology

##### **3.4.1.1 Architecture of neural networks**

Three types of neuron layers form the basic architecture of neural networks, namely hidden, input and output layers. Based on the architecture (connection pattern), two categories can be formed in ANNs (Jain et al., 1996). In networks that are feed-forward, flow of the signal is in the feed-forward direction or from the input to the output units. While the processing of the data can move over many layers of units, there are no feedback connections available. Networks that are incorporating feedback connection are recurrent networks. The properties of the network that are dynamic in nature are very important, which is in contrast to the feed-forward networks. In some instances there is a relaxation

process for the activation values of the units so that the evolution of the network to the stable states wherein there are no changes in the activations. In some other uses, the activate values change in the neurons of the output values are high so that the output of the network is constituted by the dynamical behavior. The multispectral reflectance values of the individual pixels combined with their surface roughness, texture, slope, terrain elevation, aspect, etc. The application of the hidden layer neurons helps the non-linear simulation of the neurons in the data that is input (Abraham, 2005).

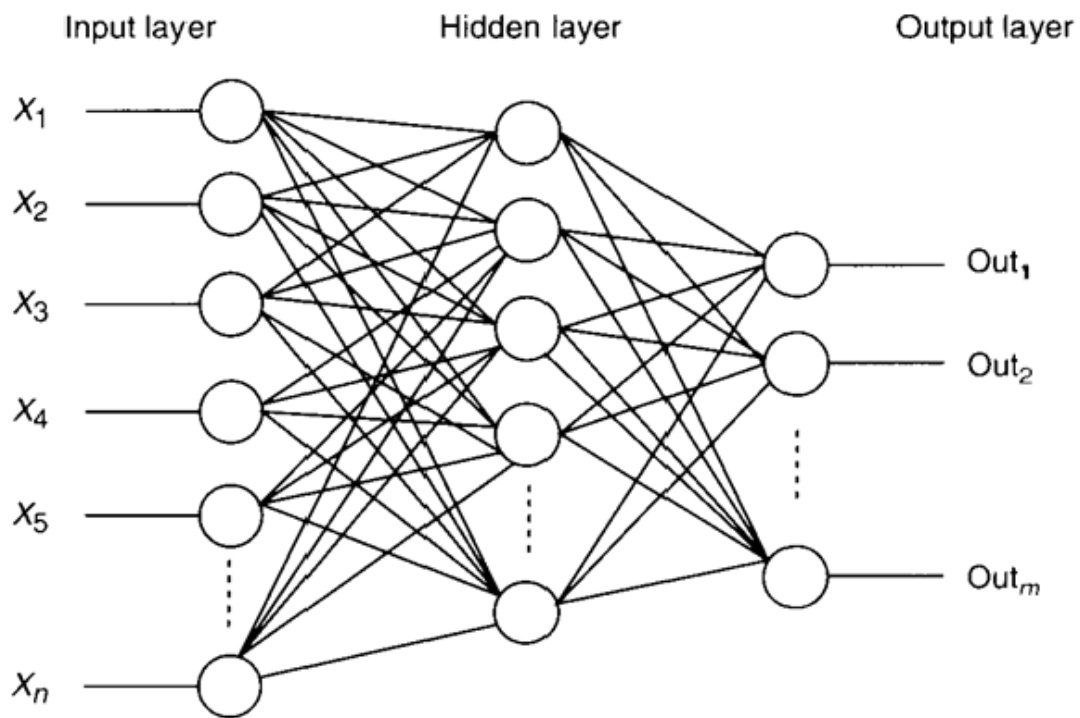
#### **3.4.1.2 Learning**

In the ANN context, the learning process can be stated to involve the intricacies of studying the updates of the architecture of the network and the weights of the connection so that tasks can be performed efficiently by the network. The connection weights must be learnt by the network from the patterns that are available. When there is iterative updating of the connection weights, the network must automatically learn to improve performance. The ability of ANNs to learn intuitively through experience enables them to be considered exciting. ANNs learn through their own experience by learning the underlying rules, like for example, the input-output relationships, instead of following a fixed set of rules that are coded by the humans. This is the major advantage of ANNs over traditional systems (Jain et al., 1996).

There are three different learning situations in ANNs, namely, reinforcement, supervised and unsupervised learning. In supervised learning, the inputs are provided in the form of an input vector along with the responses that are desired at the output layer, one for each node. The actual and the desired responses are compared and the discrepancies are found through a forward pass. This learning is used to then determine changes in the weights as per the subscribed rule. Since the desired signals in the output is provided by an external source (teacher) and therefore, this is called supervised learning (Abraham, 2005).

ANNs need testing and training like supervised classification to ensure that the information extracted from the ancillary and remotely sensed data to be useful information (Jensen, 2005).

Various classification problems are widely solved using a supervised ANN using the back-propagation technique. Figure 3.5 depicts the structure of the topology of the back-propagation ANN.



**Figure 3.5:** Components of a typical ANN

Back-propagation ANN comprises of the typical components of hidden, input and output layers. Individual training pixels comprising spectral reflection throughout the spectrum and other descriptors like slope, elevation, etc. are present in the input layer. Each layer comprises of interconnected nodes. This state of connectedness provides for the flow of information in different directions at the same time or in other words, back propagation is allowed. The weight of the connections is stored and then learnt by the ANN. During the classification procedure these weights are used. As the representativeness of the training data increases, the probability that the ANN will better mirror the reality and produce

classification that is more accurate. The individual thematic map classes such as forest or water may be represented by the output layer (Jensen et al., 2001).

#### **3.4.1.3 Training**

During training, x, y locations are specified by the analyst in the input image with attributes that are known used as the training sites (Jensen et al., 2001). The spectral information per pixel and the surrounding contextual information for the sites of training are accumulated and conveyed to the ANN input layer. The class value or the true target value is relayed to the output layer by assigning the value of 1 to this class at the same time and the rest of the neurons are denoted by the value 0.

Training of the neural network from an image as an example at a certain point of time and place may represent the state of things at the vicinity of the image and may be so for a particular reason and therefore, this cannot be extended through time and space.

When the weights are adjusted through the algorithm of back propagation, the learning is accomplished. Each time training happens, the true target value and the output of the network are compared. The difference between these values is considered as an error value and the feedback of the same is passed to the previous layers for the updating of the connection weights. The adjustment that is made to the network is proportional to the level of the error. Further improvements in the network will not be possible when the root mean square (RMS) error diminishes after a number of iterations of such feedback happen (Jensen et al., 1999) and at this stage it can be considered that the training process is accomplished and the network has achieved convergence. The inherent rules are stored in the network as example weights and they are used in the testing phase.

#### **3.4.1.4 Testing (classification)**

During this stage, the textual and/or the spectral characteristics of the scene in the form of individual pixels are passed on as input neurons irrespective of whether they originate from the rural or urban geography. The weights stored in the network are compared with the input neurons to produce an output value for the output layers. The fuzzy membership grade between 0 and 1 is assigned to every pixel in the output neuron representing the class of the neuron. The fuzzy classification map of the entire study is obtained through the value of every output neuron of every class. When these maps are defuzzified, a hard

classification map is obtained through using a local maximum function by unique classification of each pixel through the fuzzy membership highness (Jensen et al., 2001).

#### **3.4.1.5 Advantages of ANN**

- Adaptive learning from existing examples is possible, which is able to create objective classification.
- As there is more training data available, ANNs adjust the weights progressively, thus learning continuously.

#### **3.4.1.6 Disadvantages of ANN**

- Comprehensive explanation of how a classification or output has been obtained from an ANN can be difficult to trace or explain.
- ANN suffers from the inherent disadvantage of inability to depict the knowledge gained by the network in simple if-then relationships or rules.
- Image classification and interpretation rules of the ANN are hidden in the neuron weights and interpreting them can be complex (Qui and Jensen, 2004).

### **3.4.2 Support vector machine (SVM)**

SVMs are powerful machine learning techniques for classification and regression (Ganesan et al., 2014). Together with regression estimation and linear operation inversion, the SVMs are capable of providing a novel approach to pattern recognition problems and can establish connections with learning theories from statistics very clearly (Burges, 1998).

In a variety of areas, SVMs have provided a number of successful applications which include pattern recognition, supervised classification techniques, biometrics, image analysis and bioinformatics (Ma and Guo, 2014).

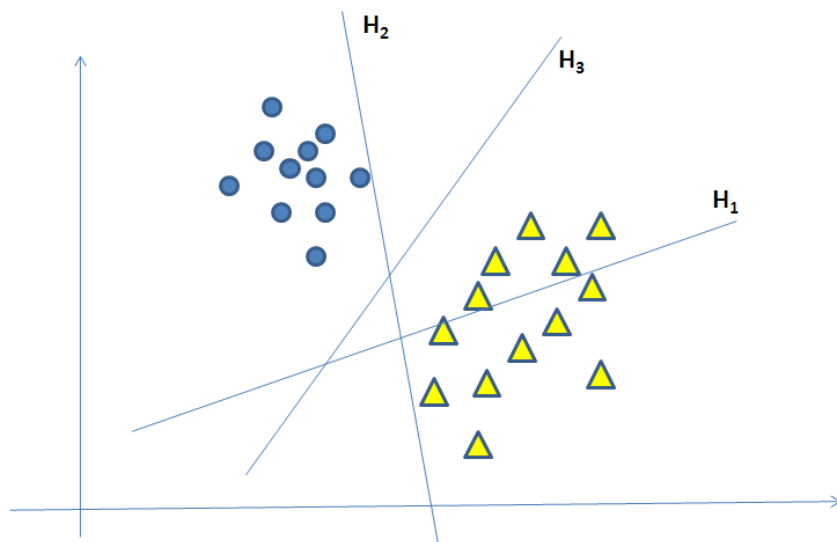
There is a newer approach to the other classification techniques that are supervised in nature. They are based on Vladimir Naumovich Vapnik's (a Russian scientist) statistical learning techniques, which were formulated in 1962. Those initial approaches have been improved and perfected by a number of researchers through new algorithms and techniques. Theoretical and algorithmic analysis of SVM has interested many researchers since then and this has merged different disciplines like functional analysis, statistics,



machine learning and optimization. Cartes and Vapnik later on introduced the soft margin classifier and the regression case came into the purview in 1995 (Cristianini and Shawe-Taylor, 2005).

Binary classification is the purpose of SVMs. This binary approach can also be extended to include scenarios which are multi class in nature encountered commonly in remote sensing. This is achieved through the dividing the problem (multiclass) into sequential analyses which are binary and the same can be solved using SVM (binary) using the one-against-all or the one-against-one approach (Mathur and Foody, 2008).

Constructing one hyperplane or many in a combination in an infinite and a high dimensional space is done by the support vector machines. Separations between separable classes can be done through any misclassification free hyperplane of the classes on all data points and this is called linear classification. On the other hand, many other hyperplanes as shown in Figure 3.6 below can classify the same data set. The objective of SVM is to find an approach to find the most optimal plane or the plane which provisions the largest possible distance margin of the two classes' nearest points. This is called the functional margin. In general, this approach ensures that when the margin is larger, the classifier returns a lower generalization error (Thome, 2012).



**Figure 3.6:** Hyperplanes' separation

H3 separates the planes with the highest margin, while H2 does it with a lower margin, H1 does not. When there is such a hyperplane in existence, the best border of separation is provided over these classes and this hyperplane is known as the hyperplane of maximum-margin and the corresponding maximum margin classifier (Thome, 2012).

### 3.4.2.1 Defining the SVM classifier formally

The hyperplane is the separation used by the SVM in the surface model. If we consider  $W$  to be the vector normal and  $b$  to be its relative in displacement for the origin, we derive the function of decision for the corresponding input as shown in equation 3.2 below. (Thome, 2012).

$$D(x) = W * x - b \quad (3.2)$$

where,

$$x \in \begin{cases} A & \text{if } D(x) > 0 \\ B & \text{if } D(x) < 0 \end{cases} \quad (3.3)$$

As depicted in Figure 3.7 below, the distance between the hyperplane and the  $x$  (with signal) is given by equation 3.4.

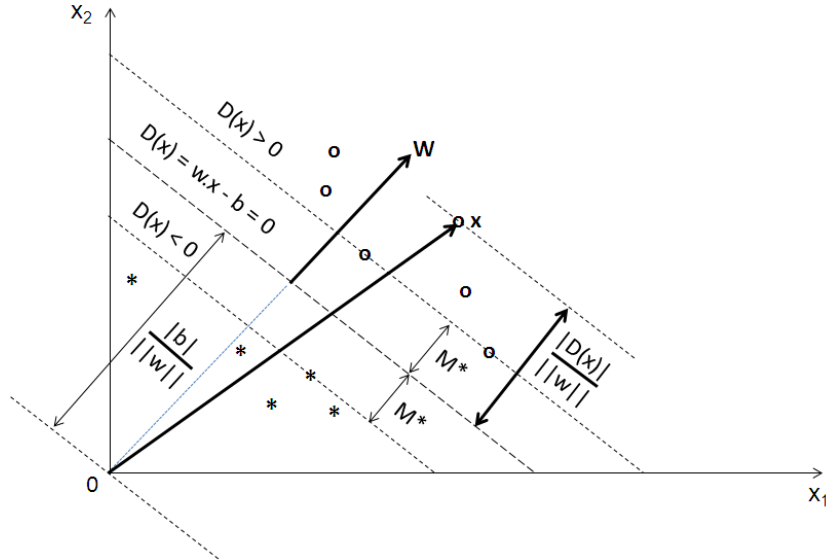
$$\frac{D(x)}{\|W\|} \quad (3.4)$$

Therefore,  $D(x_1)$  and  $D(x_2)$  which belong to different sets will have opposite signs and this is possible if and only the hyperplane's opposite sides are situated by if  $x_1$  and  $x_2$ .

$$\frac{|b|}{\|W\|} \quad (3.5)$$

As depicted in Figure 3.7, the Vector  $W$  and the hyperplane are perpendicular and the parameter in equation 3.5 explains the hyperplane that is offset in the form of a vector from

the origin. The margin  $M$  can be maximized through the choice of  $b$  and  $W$  and this will situate them as far as possible while the data sets are still separated. The following equations can describe these two hyperplanes (Cristianini and Shawe-Taylor, 2005).



**Figure 3.7:** Separating hyperplane (in two dimensions)

$$\begin{aligned} W * x - b &= +1 \\ \text{and} \\ W * x - b &= -1 \end{aligned} \tag{3.6}$$

Let  $x_1, \dots, x_p$  represent the set of sample points and  $y_1, \dots, y_p$  be their respective group classification where

$$y_i = \begin{cases} +1 & \text{if } x_i \in A \\ -1 & \text{if } x_i \in B \end{cases} \tag{3.7}$$

The training data comprising of the two samples being separable linearly, the two hyperplanes can be selected in such a fashion so that there are no points separating them and then we can attempt maximization of the distance between the planes (Bhavsar and Panchal, 2012).

$$\frac{2}{||W||} \tag{3.8}$$

Equation 3.8 represents the distance between the two hyperplanes and when this is maximized,  $W$  is minimized and so as to prevent the margin  $M$  attracting the data points, the following constraint is added to each equation (Thome, 2012).

$$\begin{aligned} W * x_i - b &\geq +1 \quad \forall i, y_i = +1 \\ \text{and} \\ W * x_i - b &\leq -1 \quad \forall i, y_i = -1 \end{aligned} \quad (3.9)$$

The maximum-margin hyperplane is found through the algorithm of this approach into the feature space that is transformed. There may be a non-linear transformation and/or high dimension may characterize the transformed space. The hyperplane is drawn by the classifier in the feature space which is a curve that is non-linear separation in the space that was original.

If a Gaussian radial basis function is based as the ideal kernel in this case, Hilbert space of indefinite dimension can be formed as a feature space. Well regularized maximum margin classifiers can be infinite dimension and therefore, the results are intact and effective. Some of the common kernels are (Thome, 2012).

Polynomial (homogeneous) kernel, Radial Basis Function kernel and Gaussian Radial Basis Function kernel.

### 3.4.2.2 RBF kernel

The RBF kernel can be a reasonable choice by default, often as it can map non-linearly the samples into a space of higher dimension and hence is different from the linear kernel in being able to handle relations between attributes and class labels that are non-linear state that the linear kernel is to be taken as a variant of the RBF kernel since penalty parameter  $C$  in the RBF and linear models have similar performance with the added parameters ( $C$ ,  $\gamma$ ). Additionally, the RBF and the sigmoid kernels behave similarly under certain conditions (Hsu, et al., 2010).

The hyperparameters' number that affects the model complexity of the selection is a secondary reason. The polynomial has greater number of hyperparameters than the RBF model and the latter has lesser difficulties that are numerical in nature. The important point is that  $0 < K_{ij} \leq 1$  which is different from the polynomial kernel where the values can

stretch to infinity ( $x_i^T x_j + r > 1$ ) or potentially to zero ( $x_i^T x_j + r < 1$ ) when there is a large degree. However, it is known that under some parameters the sigmoid kernel is invalid (Vapnik, 1995).

However, the RBF kernel cannot be applied universally, especially when there a large number of features where the linear kernel may be more useful (Hsu, et al., 2010).

### 3.4.2.3 Grid-search and cross-validation

$C$  and  $\gamma$  are the basic parameters for the RBF (kernel) and anonymous at the beginning the best  $C$  and  $\gamma$  for a problem and therefore, model selection of some kind must be performed in the form of a parameter search. Good  $C$  and  $\gamma$  must be identified so that the unknown data (testing data) can be predicted accurately by the classifier. It needs to be note that it may not always ideal to try to make the training very accurate or the accurate prediction of the training data by known classifiers. As stated before, the most useful strategy is the separation into two parts the available data from which one must be treated as unknown. The accuracy of prediction obtained from the set that is considered unknown can more predict the accuracy of the performance and when this approach is advanced, this is called as cross validation.

In v-fold cross validation, the training set is divided into v subsets of an equal size and one is tested through the classifiers which are trained from the subsets v-1. Thereby, whole training for each instance is provided with every other subset at least once and hence the accuracy of the cross validation is bound to be higher and this can prevent the problem of overfitting and this is represented in Figure 3.9 presenting a binary classification problem. Training data is represented by the filled circles and triangles and the testing data is represented by the hollow circles and triangles.

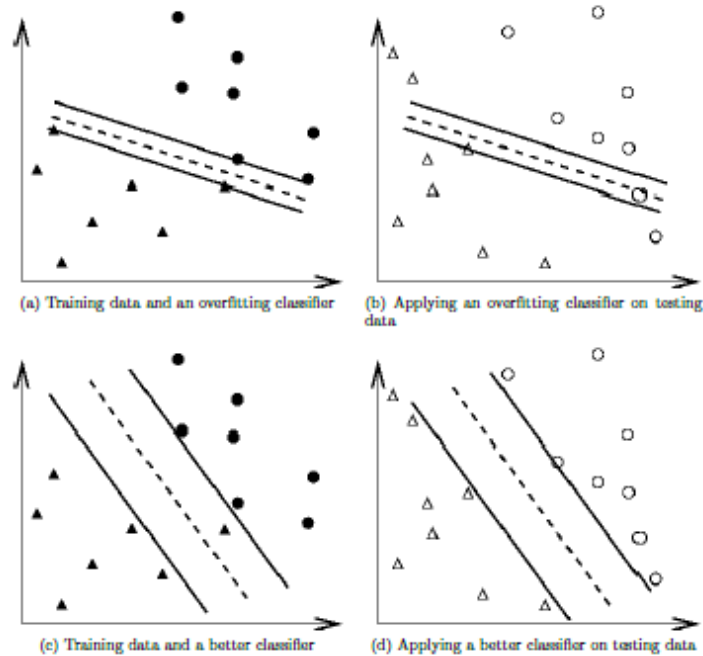
The classifiers' accuracy in testing shown in Figure 3.9a and 3.9b are not accurate as there is some overfitting of the training data. If it is considered that the data sets in 1a and 1b are for testing and training, as the validation and training ones during cross validation, high accuracy cannot be observed. Alternatively, as shown in Figures 3.9c and Figure 3.9d, there is no overfitting of the training data and this provides better testing accuracy and cross validation.

It is recommended that a “grid search” is performed through cross validation on  $C$  and  $\gamma$ . Different ( $C, \gamma$ ) pair values can be experimented with and the specific set with the cross

validation accuracy which is highest can be found. It was observed that the combination of  $C$  and  $\gamma$  in exponentially growing sequences produced the best results for good parameter identification.

The grid search methodology appears too straightforward and hence appears naïve. However, there are several advanced methodologies that can be used to save costs like cross validation rate approximation and it needs to be stated that there are two principal reasons for preferring the simplistic approach of grid search,

The first reason is that we may be assured of the accuracy of the results that are obtained without an exhaustive parameter search or using approximations, from a psychological perspective. The second reason is that grid search in this case has to contend with only two parameters and hence using this technique and the advanced techniques will consume more or less the same amount of time. The other advantage is that since each  $(C, \gamma)$  is independent, the grid search can be run parallelly, while many of the advanced approaches cannot be parallelized (Hsu, et al., 2010).



**Figure 3.8:** Overfit classifier and better classifier

### 3.5 Summary

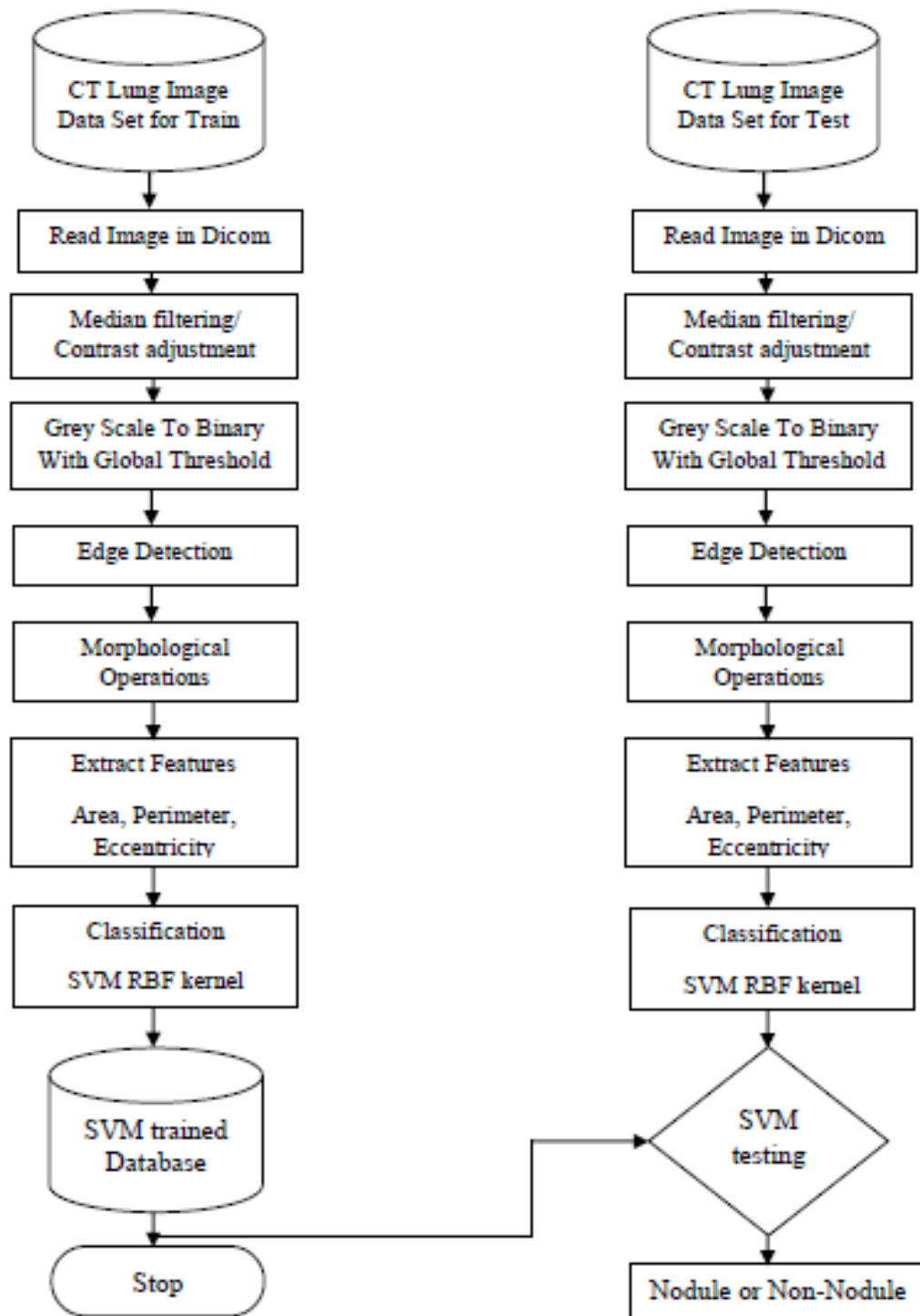
Enhancement of the image is principally performed to ensure that the quality of the image is better than the original image in a subjective manner. Therefore, a method which is used for x-ray chest enhancement may not be very suitable for MRI brain image enhancement. In several image processing applications and computer vision, the segmentation is a very important process to be performed. This is because of the fact that it is the initiation of the image processing into more complicate approaches. There is a number of segmentation approaches proposed and yet there is no standard technique which can be used universally across all application. An attribute of the image or a primitive characteristic of the image is called an image feature. A feature set that is good contains distinguishable elements that can help in clearly identifying and segregating the features. Objects in the same class must not be categorized as different and therefore, the system of understanding must be robust in nature. A small set who values can accurately represent the different features of the image must be developed. Classification which is the significant component of image analysis, forms the classes by analysing the image features. Two major type of classifications are; unsupervised classification and supervised classification. Learning systems that are biological in nature in the form of complex interconnected neurons have inspired the study of artificial neural networks (ANNs). Neural networks can be used to segregate remote sensor data of various types and have also been produced better and more accurate results as compared to the statistical methods that are used commonly. SVMs are powerful machine learning techniques for classification and regression. Together with regression estimation and linear operation inversion, the SVMs are capable of providing a novel approach to pattern recognition problems and can establish connections with learning theories from statistics very clearly. In a variety of areas, SVMs have provided a number of successful applications which include pattern recognition, supervised classification techniques, biometrics, image analysis and bioinformatics.

## **CHAPTER 4**

### **IMPLEMENTATION OF CAD SYSTEM FOR LUNG CANCER DETECTION USING SVM**

This chapter will provide the technical steps of implementation of CAD system for lung cancer detection using SVM. The main steps of proposed CAD system comprised of pre-processing, segmentation, feature extraction and classification. Matlab R2015a for windows operating system is utilized in the implementation of CAD system for lung cancer detection, which is a programming language developed by MathWorks. In the study, subset of public database called LIDC is used. The aim of Lung Image Database Consortium is to aid and support the institutions which formed consortium to build up guidelines for CT lung image and to establish a database of CT lung images. The image collection contains 1018 documented cases of CT scans, and it is accessible through the internet for the researchers, for teaching and training purposes as well as for CAD evaluation purposes (LIDC, 2015). The subset comprises of an image set of 271 documented whole-lung CT scans. The images were in DICOM which is a standard format in medical imaging. The flowchart of the implementation is shown in Figure 4.1.

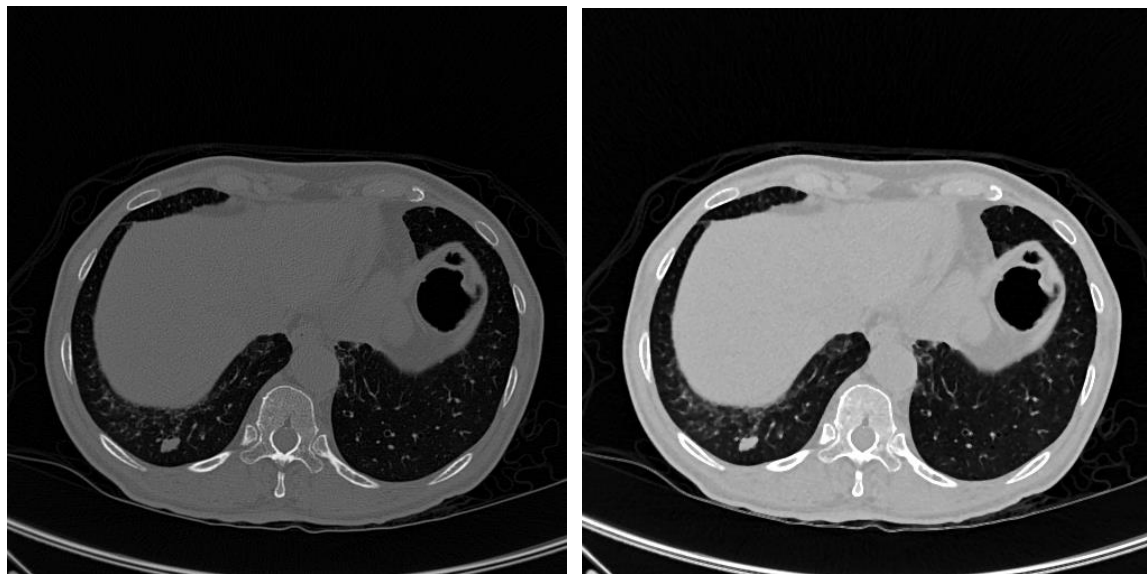




**Figure 4.1:** The flowchart of implementation of CAD system

#### 4.1 Pre-Processing

Interference and other phenomena affect the quality and contrast of the medical images, which are caused by noise and poor illumination. In the study, CT images were acquired as a raw DICOM format and interference is determined in the images. Therefore, median filtering image processing technique applied in order to reduce “salt and pepper” noise whilst preserving the edges. And contrast adjustment is carried out on the gray scale image. In Figure 4.2, the CT lung images are depicted. This process is one of the crucial steps in CAD system, since improving the enhancement of the medical image advance the further analysis including image analysis, feature detection and so on. Losing important information about the image would negatively impact the success of further analysis, hence it is taken in to consideration to keep the original structure of the image due to enhancement of the image.



(a) DICOM raw CT lung image

(b) Image after median filtering and contrast adjustment

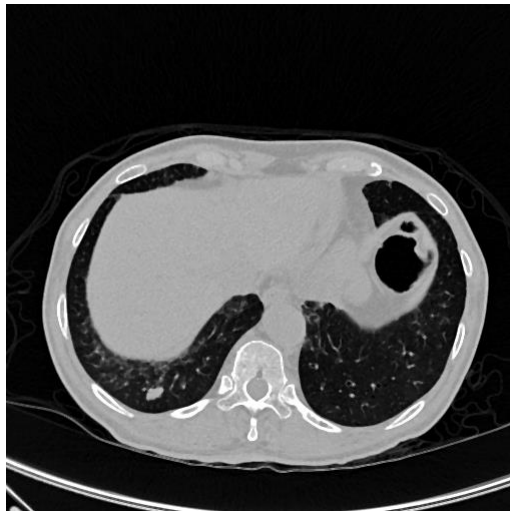
**Figure 4.2:** Lung image in pre-processing stage

#### 4.2 Segmentation

The segmentation process aims to achieve the partitioning of the image into different regions. Feature extraction, image display and image measurements in medical imaging predominantly use the segmentation operations. There is a number of segmentation

approaches proposed and yet there is no standard technique which can be used universally across all applications. Each segmentation technique has its own pros and cons. For some applications it is wisely to use various segmentation techniques together in order to have better results. For the particular requirements of this implementation, binarization process with global thresholding, edge detection and morphological operations that has major application in image enhancement and segmentation are commonly utilized.

The studied CT lung images have two types of pixels with distinct density. By using global image threshold (Otsu's method) which is a robust tool for image segmentation, the gray-scale images were converted to binary. The product image has the advantages of having smaller space in the storage and a significant increase in processing speed comparing with gray-scale image. In matlab gray thresh function uses Otsu's global image threshold method. The density of pixels in CT images which were greater than the threshold converted to white and the rest converted to black as in Figure 4.3b.



(a) Image after median filtering and contrast adjustment



(b) Image after Otsu's global threshold

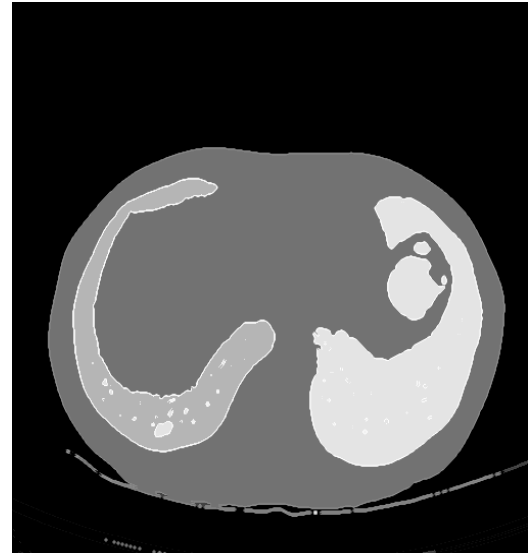
**Figure 4.3:** Lung image after Otsu's global threshold

After global threshold segmentation, edge detection utilized for identifying the boundaries of the CT lung image. As a precaution not to loose the nodules attached to the borders, gradient operator with prewitt method is applied to detect, highlight edges and the image values outside the borders accepted as uniform area. Figure 4.4a followed by clear border

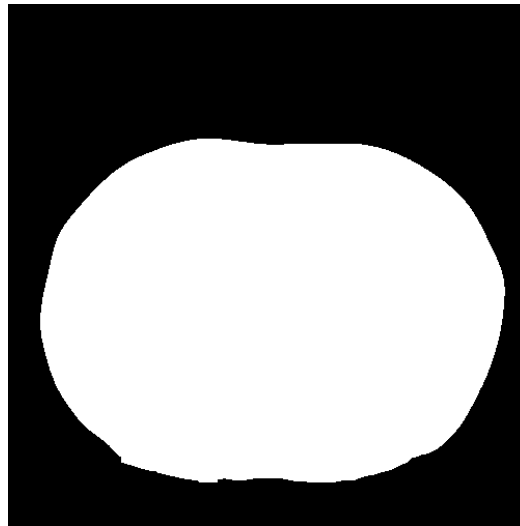
and filling detected regions as in Figure 4.4b, lastly to form a mask, gray-scale images were converted to binary while removing the unnecessary perimeter lines by applying erosion and dilation which are morphological operations as in Figure 4.4c.



(a) Edge detection gradient operator with prewitt method



(b) Clear borders and fill binary image



(c) Mask of the CT lung image

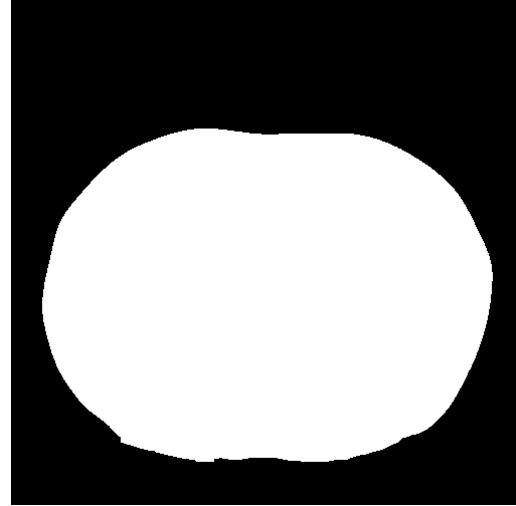
**Figure 4.4:** Forming a mask of CT lung image

Applying various image processing techniques to the CT lung image, a binary image with global threshold generated as in Figure 4.5a. Then the binary image was multiplied by the

mask in Figure 4.5b which produced in previous steps in order to obtain binary lung image with removed unnecessary perimeter lines as in Figure 4.5c.



(a) Image after Otsu's global threshold



(b) Mask of the CT lung image



(c) Image with removed perimeter lines

**Figure 4.5:** Binary lung image with removed unnecessary perimeter lines

Those steps followed with the morphological operations like `bwareaopen`, `bwmorph`, respectively until removing the undesired regions like the veins and soft tissues and nodule

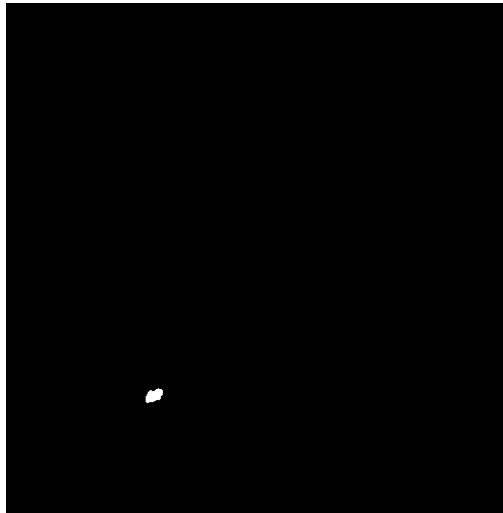
like objects. As for observation on the dataset binary objects having pixels except the range 0 and 2100 are removed to attain Figure 4.6c.



(a) Image with removed perimeter lines



(b) Segmented lung image



(c) Nodule or nodule like object

**Figure 4.6:** Final stage of CT lung image segmentation

### 4.3 Feature Extraction

Image features extraction step is very important in detecting and separating the desired region from CT lung images, which uses image processing techniques and algorithms.

In the study, the literature related to “CT images” and “feature extraction” techniques were obtained from IEEEXplore published from 2010 to 2015. For the feature extraction and selection the following methods were used in those articles as in Table 4.1.

**Table 4.1:** Relevant articles published in IEEE published from 2010 to 2015

Features Extracted	Authors	Year
Four features based on nodule shape were extracted compactness, circularity, two second, central moments.	(Xiaomin et al.)	2010
Three features were extracted; area, perimeter and eccentricity	(Chaudhary and Singh)	2012
Statistical texture features from the histogram such as average, gray level, standard deviation, smoothness, third moment, uniformity and entropy for each CT image were extracted.	(Orozco et al.)	2012
Gray features include gray average and gray variance and morphological feature includes area, perimeter, centroid, eccentricity, circularity and compactness.	(Shao et al.)	2012
Statistical texture features from the histogram and the Gray Level Co-occurrence Matrix (GLCM) in four different directions for each CT image were extracted.	(Orozco et al.)	2013
Two intensity & seven multiresolution based features were extracted such as mean and standard deviation, energy, entropy, mean, std, max probability, inverse difference moment, homogeneity	(Assefa et al.)	2013
area, energy, eccentricity, entropy, mean and standard deviation	(Tariq et al.)	2013
Three features were extracted; area, perimeter and eccentricity.	(Kulkarni and Panditrao)	2014
The GLCM is used as second order texture for feature extraction, mean standard deviation, skewness, kurtosis, entropy	(Thomas and Kumar)	2014
roundness, circularity, compactness, ellipticity and eccentricity	(Parinaz and Jamshid)	2015

Extraction step was done by regionprops and graycoprops functions which are Image Processing Toolbox in Matlab in order to predict the existence of lung cancer probability. By using those functions, Contrast, Correlation, Energy, Homogeneity, Area, MajorAxisLength, MinorAxisLength, Eccentricity, ConvexArea, EquivDiameter, Solidity, Extent, Perimeter, CentroidX, CentroidY are obtained.

Area, eccentricity and perimeter were selected as feature set. The following feature definitions are derived from matlab.

- **Area:** Gives a scalar value which defines the real number of pixels in the region of interest.
- **Eccentricity:** Gives a scalar value that defines the eccentricity of the ellipse. Region and eccentricity both have the same second-moments. It is the ratio of the distance between the major axis length and the center of the ellipse. The scalar value range is between 0 and 1. It can be conclude as it is a circle if the value is 0, while it is a line if the value is 1.
- **Perimeter :** Gives a scalar value which can be defines as the distance around the boundary of the region of interest. Perimeter can be found out by the distance between each adjoining pair of pixels around the border of the region of interest.

#### 4.4 Classification

The final step of CAD system is classification. The purpose of this step is to group the nodules and non-nodules based on the selected features. Support vector machine (SVM) is utilized which is powerful supervised machine learning techniques for classification and regression. In a variety of areas, SVMs have provided a number of successful applications which include pattern recognition, supervised classification techniques, biometrics, image analysis and bioinformatics. Recent studies point out that RBF gives the best performance, among other SVM kernels in lung nodule detection (Gomathi and Thangaraj, 2010; Aarthi and Ragupathy 2012; Gomathi and Thangaraj, 2012). Therefore, Gaussian radial basis function kernel (RBF) is selected as a SVM kernel in this step. The basic phases of supervised classification such as SVM contains training, feature selection, classification and test.



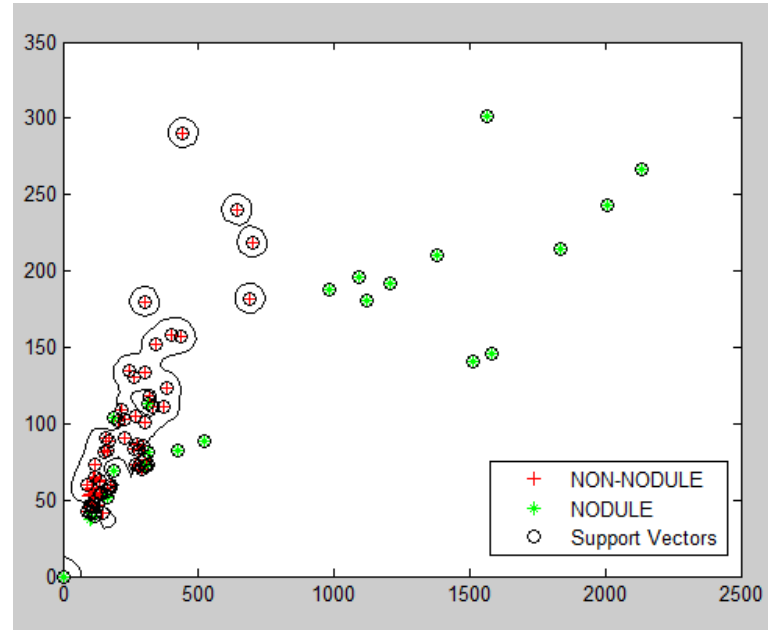
#### **4.4.1 Training phase**

LIDC database is used in the training phase. The database subset, which is formed randomly consists of 271 CT lung image scans. The grayscale images are in DICOM format that is the standard for medical images with having a size of 512x512 pixels. Nodule locations as a ground truth information are also provided which were detected by the radiologist. The total number of nodules are 49 in the database. In order to ensure an unbiased result, all images randomly selected from the database. The subset database which totally consists of 271 images are partitioned in to two groups. 50% of the CT lung images that have nodule and 50% of the CT lung images from non-nodule are used to derive training data set. Likewise, the other 50 % the CT lung images that have nodule and 50% of the CT lung images from non-nodule are used to derive the test data. The training is an iterative process in classification but for the training phase 135 lung CT images are used which consists of 24 images with nodule and 111 images with non-nodule.

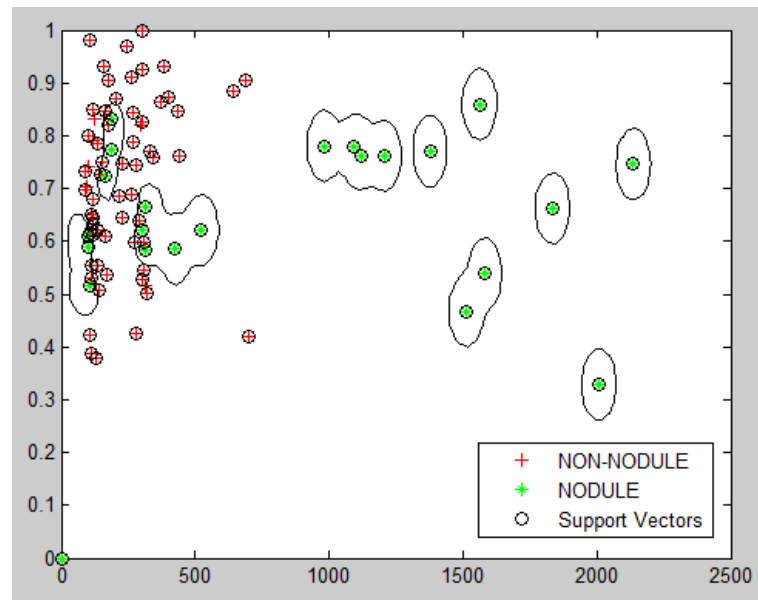
#### **4.4.2 Feature selection phase**

Area, eccentricity and perimeter were selected as main feature set according to the recent studies as indicated on Table 4.1. The goal of this step is to select the subset of features which gave better classification. Hence, by eliminating irrelevant features accuracy may increase. It also helps to reduce training time, prevent over fitting and accomplish better generalization. The features grouped as (area, perimeter), (area, eccentricity), (perimeter, eccentricity), (area, perimeter, eccentricity) and used in the training to measure the effectiveness of the classification with the selected features.

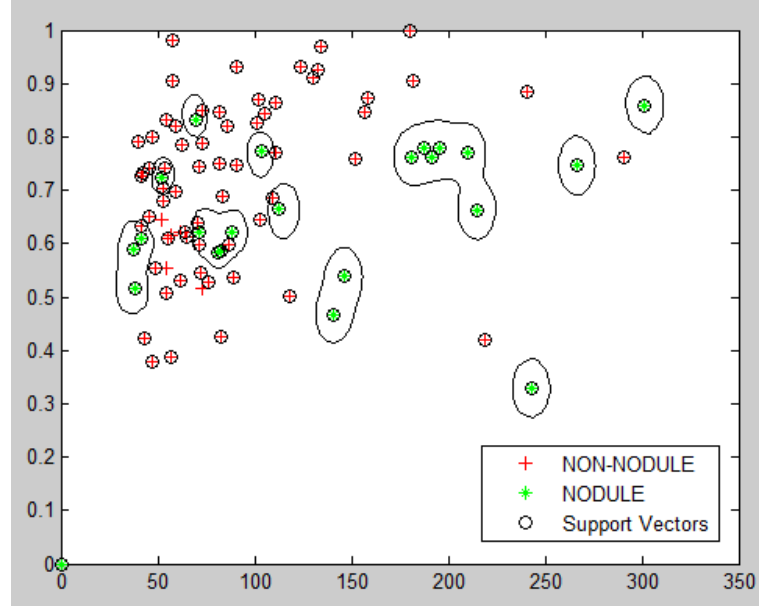
The following figures related with the selected feature classification is generated.



**Figure 4.7:** Classification of features (area, perimeter) in training phase



**Figure 4.8:** Classification of features (area, eccentricity) in training phase



**Figure 4.9:** Classification of features (perimeter, eccentricity) in training phase

#### 4.4.3 Testing phase

LIDC database is used in the testing phase. The subset database which totally contains 271 images are partitioned equally in to two groups. 136 lung CT images are used for the testing phase which consists of 25 images with nodule and 111 images with non-nodule.

The system that is trained with 135 lung CT images was tested with whole dataset as well as with 136 lung CT images different from those used for training. The performance of implementation and classification were tested by the statistical measures which are sensitivity, specificity and accuracy. The formulae of the statistical measures are as follows;

$$\text{Sensitivity} = \frac{\text{TPs}}{(\text{TPs} + \text{FNs})} \quad (4.1)$$

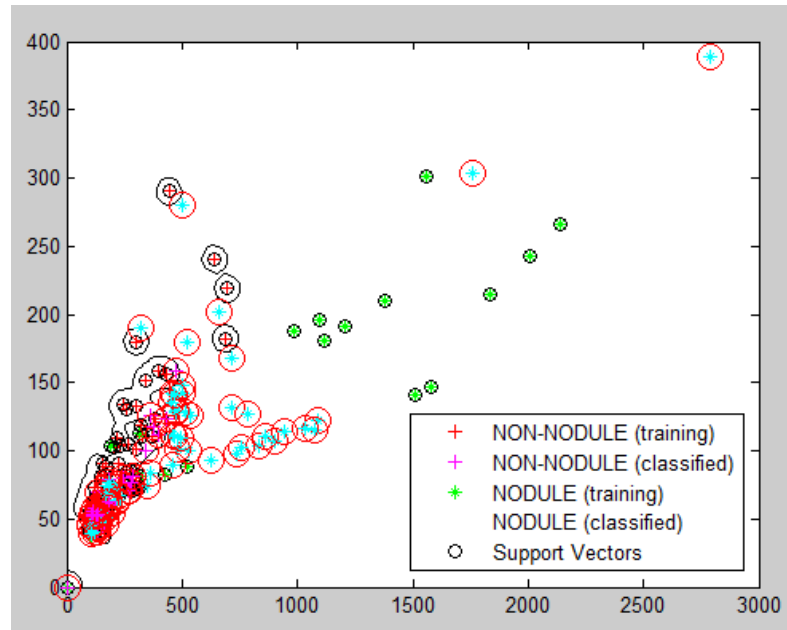
$$\text{Specificity} = \frac{\text{TNs}}{(\text{FPs} + \text{TNs})} \quad (4.2)$$

$$\text{Accuracy} = \frac{(\text{TP} + \text{TN})}{(\text{P} + \text{N})} \quad (4.3)$$

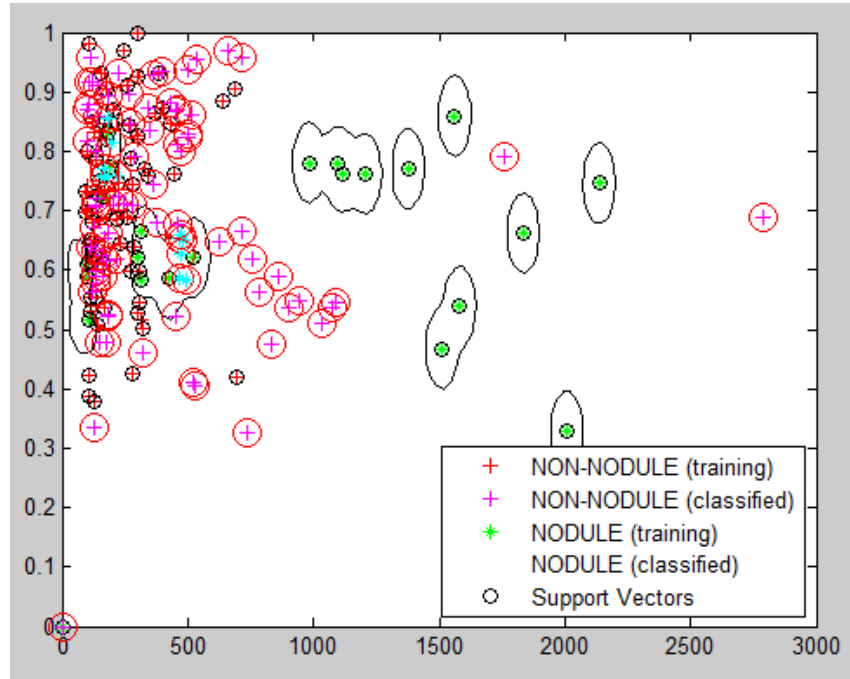
- **True positive (TP):** Nodule exists on the image and correctly diagnosed that Nodule is present.
- **False positive (FP):** Nodule does not exist on the image but diagnosed that Nodule is present.
- **True negative (TN):** Nodule does not exist on the image and correctly diagnosed that Nodule is not present.
- **False negative (FN):** Nodule exists on the image but diagnosed that Nodule is not present.

The features grouped as (area, perimeter), (area, eccentricity), (perimeter, eccentricity), (area, perimeter, eccentricity) and used in the testing phase to measure the effectiveness of the classification with the selected features.

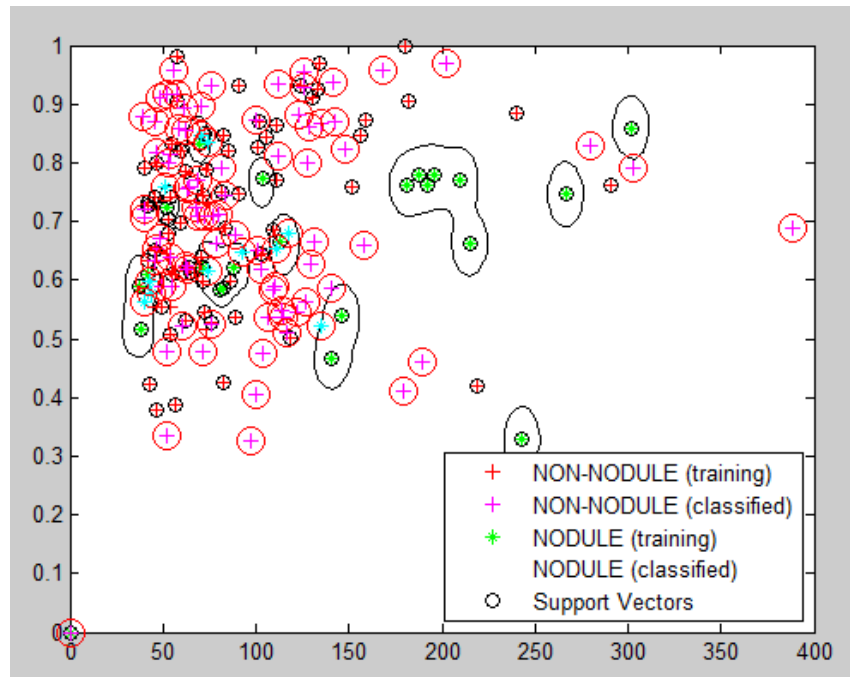
The following figures related with the selected feature classification is generated.



**Figure 4.10:** Classification of features (area, perimeter) in testing phase



**Figure 4.11:** Classification of features (area, eccentricity) in testing phase



**Figure 4.12:** Classification of features (perimeter, eccentricity) in testing phase

Table 4.2 indicates the confusion matrix for SVM using RBF kernel with whole dataset. Test has been carried out with entire dataset, which consists of 49 nodules and 222 non-nodules. The probability of system detecting positive, when there is nodule is 95.92% while the probability of detecting negative is 98.65%, when nodule is not present.

**Table 4.2:** SVM using RBF kernel with entire dataset

TP (a)	FN (b)	FP (c)	TN (d)	Sensitivity $a / (a+b)$	Specificity $d / (c+d)$	Accuracy $(a+d) / (P+N)$
47	2	3	219	95.92	98.65	98.15

*Note:* Total Positive (P): 49 Total Negative (N): 222

In Table 4.3 the confusion matrix for SVM using RBF kernel with test dataset takes part. Test has been carried out after training process only with test dataset, which consists of 25 nodules and 111 non-nodules. It is find out that the probability of detecting positive, when nodule presents is 92%, whereas the detecting rate of negative is 97.30% when there is no nodule.

**Table 4.3:** SVM using RBF kernel with test dataset

TP (a)	FN (b)	FP (c)	TN (d)	Sensitivity $a / (a+b)$	Specificity $d / (c+d)$	Accuracy $(a+d) / (P+N)$
23	2	3	108	92	97.30	97

*Note:* Total Positive (P): 25 Total Negative (N): 111

In Table 4.4 the confusion matrix for SVM using Quadratic kernel with test dataset takes part. Test has been carried out after training process only with test dataset, which consists of 25 nodules and 111 non-nodules. It is find out that the probability of detecting positive, when nodule presents is 60%, whereas the detecting rate of negative is 93.69% when there is no nodule.

**Table 4.4:** SVM using Quadratic kernel with test dataset

TP	FN	FP	TN	Sensitivity	Specificity	Accuracy
(a)	(b)	(c)	(d)	$a / (a+b)$	$d / (c+d)$	$(a+d) / (P+N)$
15	10	7	104	60	93.69	87.5

Note: Total Positive (P) :25 Total Negative (N): 111

Table 4.5 indicates the confusion matrix for SVM using Linear kernel with test dataset. Test has been carried out after training stage merely with test dataset which consists of 25 nodules and 111 non-nodules. The probability of system detecting positive, when there is nodule is 72% while the probability of detecting negative is 94.59%, when nodule is not present.

**Table 4.5:** SVM using Linear kernel with test dataset

TP	FN	FP	TN	Sensitivity	Specificity	Accuracy
(a)	(b)	(c)	(d)	$a / (a+b)$	$d / (c+d)$	$(a+d) / (P+N)$
18	7	5	106	72	94.59	91

Note: Total Positive (P): 25 Total Negative (N): 111

The following Table 4.6 shows the statistical measurements of selected features with respect to SVM kernels which is found out in the testing phase. The tests showed that, RBF gave best results for the classification. It is determined that lung vessels make up the majority of the false positive whereas small nodules make up the majority of the false negative.

**Table 4.6:** The effectiveness of the SVM kernels

SVM kernel	Sensitivity	Specificity	Accuracy
RBF	92	97.30	97
Quadratic	60	93.69	87.5
Linear	72	94.59	91

#### 4.5 Related works

As a result of the fatal and one of the leading cancer types, there has been plenty of studies to detect lung cancer. The recently, proposed paper named “Lung Nodule Detection Using Multi-Resolution Analysis” implemented a CAD system to detect lung nodules. The authors selected 165 CT images which consists of nodules from ELCAP database. As a first step, wiener filtering applied to the dataset to accomplish noises and a “novel circular object detection scheme” is used for the pre-processing step for the nodule detection. For the false positive reduction phase, Multi-resolution and Intensity based features were extracted. Nodules that are low in contrast and attached to the lung vessels were missed. The accuracy rate of the system is 81.21% (Assefa, et al., 2013). Another study called “Lung Nodule Classification in CT Thorax Images using Support Vector Machines” remove the segmentation step as a novel method. The authors selected 128 CT lung images from ELCAP and NBIA databases. After acquiring the CT images, statistical texture features like gray level average, standard deviation, smoothness, third moment, uniformity and entropy were computed, followed by 16 GLCM as a second order texture features. In the study, SVM is utilized as a classifier which gives reliability index of 84% (Orozco, et al., 2013). Parinaz, et al., introduced nodule detection system by using SVM. Initially, the dataset is pre-processed to get 49 nodule and 98 non-nodule images. At the segmentation step, a function is used which applies statistical region merging and segment feature extractor for converting numeric values into binomial. Gradient, geometric and intensity feature groups extracted. It is find out that SVM performed ideally in gradient feature group thus the system reached 89.9% of sensitivity. Ferzad, et al., firstly enhanced the dataset which consists of 60 CT images, followed by the region growing and thresholding. Following segmentation, post processing is applied in order to remove the unwanted regions and objects. After extracting five features, multiple classifiers named MLP, KNN and SVM is used. Majority voting system generated 84.16% in sensitivity. Yang, et al., evaluated the suggested system with 24 patient’s dataset consisting of 59 nodules. The raw images is segmented with combined thresholding and region growing methods and ROIs extracted by the selective enhancement filter. Followed by the feature extraction process which 4 different feature function is used. Hidden Conditional Random field is utilized as a classifier and gives 89.3% sensitivity with 1.2 FPs/patient.



**Table 4.7:** Summary of recent works related with lung cancer

<b>Dataset Used</b>	<b>Authors</b>	<b>Success Rate</b>
The authors selected 165 CT images, which consists of nodules from ELCAP database.	(Assefa et al., 2013)	Accuracy 81.21%
Totally 128 CT images used and 66 of them were consists of nodule images from ELCAP and NBIA database.	(Orozco et al., 2013)	Accuracy 84%
Dataset consists of 59 nodules from LIDC. 31 nodules used for training whereas 28 nodules used for testing. (LIDC)	(Yang et al., 2014)	Sensitivity 89.3%
Evaluated with 49 nodule and 98 non-nodule images. (LIDC)	(Parinaz and Jamshid, 2015)	Sensitivity 89.9%
The CAD system evaluated with 60 patients. Authors did not mention about the number of nodules (LIDC)	(Farzad et al., 2015)	Sensitivity 84.16%
The proposed CAD system evaluated with 49 nodule and 222 non-nodule images (LIDC)	(Proposed CAD System, 2016)	Sensitivity 92% Accuracy 97%

## **CHAPTER 5**

### **CONCLUSION**

A CAD system for lung cancer detection using SVM is performed. Proposed system helps the doctors to decide in evaluation of CT lung images whether there is a nodule or not. In the study, LIDC database is used which consists of documented lung CT images. The implemented CAD system consists of image pre-processing, segmentation, feature extraction and classification steps. Binarization process using global thresholding, edge detection and morphological operations that has major application in image enhancement and segmentation are commonly utilized at the segmentation step. SVM with RBF, which several studies also indicate, gives the best performance applied for the classification step. The subset database which totally consists of 271 images are partitioned into two groups. 50% of the CT lung images containing nodule and 50% of the CT lung images from non-nodule are used to derive training data set. Likewise, the other 50% the CT lung images containing nodule and 50% of the CT lung images from non-nodule are used to derive the test data. The performance of implementation and classification were tested by the statistical measures which are sensitivity, specificity and accuracy. The accuracy of classifier with regard to the fundamental truth is measured by confusion matrix. The CAD system performed at a rate of 97% accuracy, and achieved 92% sensitivity. It is determined that lung vessels make up the majority of the false positive whereas small nodules make up the majority of the false negative. It is possible to increase the accuracy of the implemented system by using different and larger image database for training. As a further suggestion, it can be also suggest the radiologists to cooperate with the engineers of computer science. For future studies, the proposed system will be tried by using different classifier like artificial neural networks.

## REFERENCES

- Aarhy, K.P., and Ragupathy, U.S. (2012). Detection of lung nodule using multiscale wavelets and support vector machine. *International Journal of Soft Computing and Engineering*, 2(3), 32-36.
- Abraham, A. (2005). Artificial neural networks. In P. Sydenham, and R.Thorn (Eds.), *Handbook of measuring system design* (pp. 901-908). doi: 10.1002/0471497398mm421
- American Cancer Society. (2016). About lung cancer. Retrived February 1 , 2016 from <http://www.cancer.org/cancer/lungcancer-non-smallcell/>
- Anshad, P.Y.M., and Kumar, S.S. (2014). Recent methods for the detection of tumor using computer aided diagnosis — A review. In *Proceedings of IEEE International Conference on Control, Instrumentation, Communication and Computational Technologies* (pp. 1014 – 1019). Kanyakumari: IEEE.
- Assefa, M., Faye, I., Malik, A.S., and Shoaib, M. (2013). Lung nodule detection using multi-resolution analysis. In *Proceedings of IEEE International Conference on Complex Medical Engineering* (pp. 457 - 461). Beijing: IEEE.
- Bankman, I.N. (2000). Handbook of medical image processing and analysis. Burlington, MA: Elsevier, Academic Press.
- Bankman, I.N. (2009). Handbook of medical image processing and analysis. (2<sup>nd</sup> Ed.). Burlington, MA: Elsevier Inc.
- Bhavsar, H., and Panchal, M.H. (2012). A review on support vector machine for data classification. In *Proceedings of the International Journal of Advanced Research in Computer Engineering & Technology*, 1(10), 185-189.
- Biederer, J., Beer, M., Hirsch, W., Wild, J., Fabel, M., Puderbach, M., and Van Beek, E. J. R. (2012). MRI of the lung (2/3). Why... when... how?. *NCBI Insights into imaging*, 3(4), 355-371.
- Biederer, J., Mirsadraee, S., Beer, M., Molinari, F., Hintze, C., Bauman, G., Both, M., Beek, E.J.R.V., Wild, J., and Puderbach, M. (2012). MRI of the lung (3/3)—current applications and future perspectives. *Insights into Imaging*, 3(4), 373-386.
- Burges, J.C.C. (1998). A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, 2, 121–167.
- Chaudhary, A., and Singh, S.S. (2012). Lung cancer detection on CT images by using image processing. In *Proceedings of IEEE International Conference on Computing Sciences* (pp. 142 - 146). Phagwara: IEEE.

- Chen, P. (2012). Study on medical image processing technologies based on DICOM. *Journal of Computers*, 7(10), 2354-2361.
- Chen, S., and Haralick, R.M. (1995). Recursive erosion, dilation, opening, and closing transforms. *In Proceedings of IEEE Transactions on Image Processing*, 4(3), 335-345.
- Chi-Chia, S., Shanq-Jang, R., Mon-Chau, S., and Tun-Wen, P. (2005). Dynamic contrast enhancement based on histogram specification. *In Proceedings of IEEE Transactions on Consumer Electronics* (pp. 1300-1305).
- Cristianini, N. and Shawe-Taylor, J. (2005). An introduction to support vector machines and other kernel-based learning methods. Cambridge: Cambridge University Press.
- Demin, W., Haese-Coat, V., Bruno, A., and Ronsin, J. (1995). Some statistical properties of mathematical morphology. *In Proceedings of IEEE Transactions on Signal Processing*, 43(8), 1955- 1965.
- Doi, K. (2007). Computer- aided diagnosis in medical imaging: historical review, current status and future potential. *Comput Med Imaging Graph*, 31(4-5), 198–211.
- El-Baz, A., Beache, G.M., Farb, G.G., Suzuki, K., Okada, K., Elnakib, A., Soliman, A., and Abdollahi, B. (2013). Review article computer-aided diagnosis systems for lung cancer: challenges and methodologies. *International Journal of Biomedical Imaging*, 2013, 1-46.
- Farzad, V.F, Abbas, A., and M.H. Fazel, Z. (2015). Lung nodule diagnosis from CT images based on ensemble learning. *In Proceedings of IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology* (pp. 1 - 7). Niagara Falls, ON: IEEE.
- Firmino, M., Morais, A.H., Mendoça, R.M., Dantas, M.R., Hekis, H.R., and Valentim, R. (2014). Computer-aided detection system for lung cancer in computed tomography scans: review and future prospects. *BioMedical Engineering*, doi:10.1186/1475-925X-13-41
- Frederick Memorial Hospital. (2016). X-ray, CT, MRI, PET. Retrived February 2, 2016 from <http://www.fmh.org/body.cfm?id=176>
- Fujitaa, H., Uchiyamaa, Y., Nakagawaa, T., Fukuokab, D., Hatanakac, Y., Haraa, T., Leea, G.N., Hayashia, Y., Ikedoa, Y., Gaoa, X., and Zhoau, X. (2008). Computer-aided diagnosis: the emerging of three CAD systems induced by Japanese health care needs. *Computer Methods and Programs in Biomedicine*, 92(3), 238-248.
- Ganesan, S., Subashini, T.S., and Jayalakshmi, K. (2014). Classification of x-rays using statistical moments and SVM. *In Proceedings of IEEE International Conference on Communications and Signal Processing* (pp. 1109- 1112). Melmaruvathur: IEEE.

- Ginneken, B., Prokop, C.M.S., and Prokop, M. (2011). Computer-aided diagnosis: how to move from the laboratory to the clinic. *RSNA Radiology*, 261(3), 719-732.
- Gomathi, M., and Thangaraj, P. (2010). A computer aided diagnosis system for lung cancer detection using support vector machine. *American Journal of Applied Sciences*, 7(12), 1532-1538.
- Gomathi, M., and Thangaraj, P. (2012). An effective classification of benign and malignant nodules using support vector machine. *Journal of Global Research in Computer Science*, 3(7), 6-9.
- Gonzalez, R.C. and Woods, R.E. (2002). Digital image processing. (2<sup>nd</sup> Ed.). Washington, DC: Prentice Hall.
- Haralick, R.M., Sternberg, S.R., and Zhuang, X. (1987). Image analysis using mathematical morphology. In *Proceedings of IEEE Transactions on Pattern Analysis and Machine Intelligence*, 9(4), 532-550.
- Hochegger, B., Marchiori, E., Sedlacek, O., Irion, K., Heussel, C.P., Ley, S., Ley, Z.J., Soares, S.A.Jr., and Kauczor, H.U. (2011). MRI in lung cancer: a pictorial essay. *The British Journal of Radiology*, 84(1003), 661–668.
- Hossain, S.S., Maiti, A., and Chaki, N. (2011). Image binarization using iterative partitioning: A global thresholding approach. In *Proceedings of IEEE International Conference on Recent Trends in Information Systems* (pp. 281-286). Kolkata: IEEE.
- Hsu, C.W., Chang, C.C., and Lin, C.J. (2010). A practical guide to support vector classification. Retrieved March 2, 2016 from <http://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>
- Jain, A.K., Mao, J., and Mohiuddin, K.M. (1996). Artificial neural networks: A tutorial, In *Proceedings of IEEE Computer Society*, 29(3), 31-44.
- Jensen, J.R. (2005). Introductory digital image processing: A remote sensing perspective. Washington, DC: Prentice Hall.
- Jensen, J.R., Qiu, F., and Ji, M. (1999). Predictive modelling of coniferous forest age using statistical and artificial neural network approaches applied to remote sensor data. *International Journal of Remote Sensing*, 20(14), 2805-2822.
- Jensen, J.R., Qiu, F. and Patterson, K. (2001). A neural network image interpretation system to extract rural and urban land use and land cover information from remote sensor data. *Geocarto International, A Multi-disciplinary Journal of Remote Sensing and GIS*, 16(1), 19-28.

- Kulkarni, A., and Panditrao, A. (2014). Classification of lung cancer stages on CT scan images using image processing. *In Proceedings of IEEE International Conference on Advanced Communication Control and Computing Technologies* (pp. 1384-1388). Ramanathapuram: IEEE.
- Larobina, M., and Murino, L. (2014). Medical image file formats. *Journal of Digital Imaging*, 27(2), 200–206.
- Lee, N., Laine, A.F., Marquez, G., Levsky, J.M. and Gohagan, J.K. (2009). Potential of computer-aided diagnosis to improve CT lung cancer screening. *In Proceedings of IEEE Reviews in Biomedical Engineering*, 2, 136-146.
- Li, Q., and Nishikawa, R.M. (2015). Computer-aided detection and diagnosis in medical imaging. Boca Raton, FL: CRC Press.
- LIDC. (2015). Image subset of 271 low-dose documented whole-lung CT scans. Retrived December 28, 2015 from <https://wiki.cancerimagingarchive.net/display/Public/LIDC-IDRI>.
- Lucchese, L. and Mitra, S.K. (2001). Color image segmentation: A state – of – the - art survey. Retrieved March 1, 2016 from [http://www.dli.gov.in/rawdataupload/upload/insa/INSA\\_2/20005aa4\\_207.pdf](http://www.dli.gov.in/rawdataupload/upload/insa/INSA_2/20005aa4_207.pdf)
- Lung Cancer. (2016). Lung Cancer 101. Retrived February 1, 2016 from [http://www.lungcancer.org/find\\_information/publications/163-lung\\_cancer\\_101/265](http://www.lungcancer.org/find_information/publications/163-lung_cancer_101/265)
- Luong, C.M. (2006). Introduction to image processing and computer vision. Retrieved March 1, 2016 from [https://campusvirtual.univalle.edu.co/moodle/pluginfile.php/78606/mod\\_resource/content/0/ChiL - Introduction to Image Processing and Computer Vision book .pdf](https://campusvirtual.univalle.edu.co/moodle/pluginfile.php/78606/mod_resource/content/0/ChiL_-_Introduction_to_Image_Processing_and_Computer_Vision_book_.pdf)
- Ma, Y., and Guo, G. (2014). Support vector machines applications. Lausanne: Springer Science and Business Media.
- Mac Manus, M. P., Hicks, R. J., Matthews, J. P., McKenzie, A., Rischin, D., Salminen, E. K., and Ball, D. L. (2003). Positron emission tomography is superior to computed tomography scanning for response-assessment after radical radiotherapy or chemoradiotherapy in patients with non–small-cell lung cancer. *Journal of Clinical Oncology*, 21(7), 1285-1292.
- Mathur, A., and Foody, G.M. (2008). Multiclass and binary SVM classification: implications for training and classification users. *IEEE Geoscience and Remote Sensing Letters*, 5(2), 241- 245.
- Memon , N.A., Mirza, A.M., and Gilani, S.A.M. (2006). Segmentation of lungs from CT scan images for early diagnosis of lung cancer. *World Academy of Science, Engineering and Technology*, 20, 113-118.

- Mesanovic, N., Mujagic, S., Huseinagic, H. and Kamenjakovic, S. (2012). Application of lung segmentation algorithm to disease quantification from CT images. *In Proceedings of the International Conference on System Engineering and Technology* (pp. 1-7). Bandung, Indonesia: IEEE.
- Mitchell, M.T. (1997). Machine learning. Boston, MA: McGraw-Hill.
- Mohammed, T.L.H., White, C.S., and Pugatch, R.D. (2005). The imaging manifestations of lung cancer. *Elsevier Inc*, 40(2), 98–108.
- Moon, S-H., Yang, B.Y., Kim, Y.J., Hong, M.K., Lee, Y.S., Lee, D.S., Chung, J.K., and Jeong, J.M. (2016). Development of a complementary PET/MR dual-modal imaging probe for targeting prostate-specific membrane antigen (PSMA). *Nanomedicine: Nanotechnology, Biology, and Medicine*, 12(4), 871–879.
- Orozco, H.M., Villegas, O.O.V., Dominguez, H.J.O, and Sanchez, V.G. (2013). Lung nodule classification in CT thorax images using support vector machines. *In Proceedings of IEEE 12<sup>th</sup> Mexican International Conference on Artificial Intelligence* (pp. 277 - 283). Mexico City: IEEE.
- Orozco, H.M., Villegas, O.O.V., Maynez, L.O., Sancez, V.G.C., and Dominguez, H.J.O. (2012). Lung nodule classification in frequency domain using support vector machines. *In Proceedings of IEEE 11<sup>th</sup> International Conference on Information Science, Signal Processing and their Applications* (pp. 870 - 875). Montreal, QC: IEEE.
- Parinaz, E., and Jamshid, B. (2015). Computer-aided detection of pulmonary nodules based on SVM in thoracic CT images. *In Proceedings of the 7<sup>th</sup> Conference on Information and Knowledge Technology* (pp. 1 - 6). Urmia: IEEE.
- Patel, S., and Goswami, M. (2014). Comparative analysis of histogram equalization techniques. *In proceedings of the International Conference on Contemporary Computing and Informatics* (pp. 167-168). Mysore: IEEE.
- Petrou, M. and Petrou, C. (2010). Image processing: The fundamentals. (2<sup>nd</sup> Ed.). West Sussex: John Wiley & Sons.
- Pitas, I., and Venetsanopoulos, A.N. (1990). Nonlinear digital filters principles and applications. New York, NY: Springer.
- Pratt, W.K. (2007). Digital image processing. (4<sup>th</sup> Ed.). New York, NY: John Wiley & Sons.
- Qui, F., and Jensen, J.R. (2004). Opening the black box of neural networks for remote sensing image Classification. *International Journal of Remote Sensing*, 25(9), 1749-1768.

- Reza, Z.M., Mueen, A., Seng, W.C., and Awedh, M.H. (2011). Combined feature extraction on medical x-ray images. *In Proceedings of the Third International Conference on Computational Intelligence, Communication Systems and Networks* (pp. 264-268). Bali: IEEE.
- Russ, J.C. (2011). The image processing handbook. (6<sup>th</sup> Ed.). Boca Raton, FL: CRC Press.
- Senthilkumaran, N. and Rajesh, R. (2009). Edge detection techniques for image segmentation – A survey of soft computing approaches. *International Journal of Recent Trends in Engineering*, 1(2), 250-254.
- Senthilkumaran, N., and Thimmiraja, J. (2014). Histogram equalization for image enhancement using MRI brain images. *In Proceedings of the World Congress on Computing and Communication Technologies* (pp. 80-83). Trichirappalli: IEEE.
- Shao, H., Cao, L., and Liu, Y. (2012). A detection approach for solitary pulmonary nodules based on CT images. *In Proceedings of the 2nd International Conference on Computer Science and Network Technology* (pp. 1253-1257). Changchun, China: IEEE.
- Sroubek, F., Sorel, M., Boldys, J., and Sroubek, J. (2009). PET image reconstruction using prior information from CT or MRI. *In Proceedings of the 16<sup>th</sup> International Conference on Image Processing* (pp. 2493 - 2496). Cairo: IEEE.
- Tang, J. (2010). A color image segmentation algorithm based on region growing. *In Proceedings of the 2nd International Conference on Computer Engineering and Technology* (pp. 634-637). Chengdu: IEEE.
- Tariq, A., Akram, M.U., and Javed, M.Y. (2013). Lung nodule detection in CT images using neuro fuzzy classifier. *In Proceedings of the Fourth International Workshop on Computational Intelligence in Medical Imaging* (pp. 49 - 53). Singapore: IEEE.
- T.C. Sağlık Bakanlığı. (2016). Türkiye kanser istatistikleri. Retrived February 1 , 2016 from <http://kanser.gov.tr/daire-faaliyetleri/kanser-istatistikleri/1793-2013y%C4%B1%C4%B1-t%C3%BCrkiye-kanser-istatistikleri.html>
- Thomas, R.A., and Kumar, S.S. (2014). Automatic detection of lung nodules using classifiers. *In Proceedings of IEEE International Control, Instrumentation, Communication and Computational Technologies* (pp. 705- 710). Kanyakumari:IEEE.
- Thome, A.C.G. (2012). SVM classifiers – concepts and applications to character recognition, advances in character recognition. doi: 10.5772/52009
- Vapnik, V. (1995). Machine learning. Retrived March 2, 2016 from <http://homepages.rpi.edu/~bennek/class/mmld/papers/svn.pdf>
- Varma, D.R. (2012). Managing DICOM images: Tips and tricks for the radiologist *Indian J Radiol Imaging*, 22(1), 4-13.



- Varshney, S.S., Rajpal, N., and Purwar, R. (2009). Comparative study of image segmentation techniques and object matching using segmentation. *In Proceedings of IEEE International Conference on Methods and Models in Computer Science* (pp. 1-6). Delhi: IEEE.
- WHO. (2015). Key Facts. Retrived February 1, 2016 from <http://www.who.int/mediacentre/factsheets/fs297/en>
- Wiggins, R.H., Davidson, H.C., Harnsberger, H.R., Lauman, J.R., and Goede, P.A. (2001). Image file formats: past, present, and future. *Radio Graphics infoRAD*, 21(3), 789-798.
- Xiaomin, P., Hongyu, G., and Jianping, D. (2010). Computerized detection of lung nodules in CT images by use of multiscale filters and geometrical constraint region growing. *In Proceedings of the 4<sup>th</sup> International Conference on Bioinformatics and Biomedical Engineering* (pp. 1 - 4). Chengdu: IEEE.
- Yang, L., Jinzhu, Y., Dazhe, Z., and Jiren, L. (2009). Computer aided detection of lung nodules based on voxel analysis utilizing support vector machines. *In Proceedings of IEEE International Conference on Future BioMedical Information Engineering* (pp. 90 - 93). Sanya: IEEE.
- Yang, L., Zhongqiu, W., Maozu, G., and Ping, L. (2014). Hidden conditional random field for lung nodule detection. *In Proceedings of IEEE International Conference on Image Processing* (pp. 3518 - 3521). Paris: IEEE.
- Yelleswarapu, V.R., Liu, F., Cong, W., and Wang, G. (2014). TOP-level designs of a hybrid low field MRI-CT system for pulmonary imaging. *In Proceedings of IEEE International Symposium on Biomedical Imaging* (pp. 975 - 978). Beijing: IEEE.

## APPENDIX

### Algorithm and Source Codes

Step 1: Start %%%PRE-PROCESSING +SEGMENTATION+ TRAINING %%%%

Step 2: Read the image data from DICOM file and set to Img1

Step 3: Apply median filtering to Img1 and set new image to Img2

Step 4: Adjust/ increase the contrast of Img2 and set to Img3

Step 5: Calculate the global threshold of Img3 and set to gt

%%If g(x, y) is a thresholded version of f(x, y) at global threshold gt %%%

Step 6: convert grey image Img3 to binary and set binary image to Img4

if f(x,y)> gt then g(x,y)=1

else

g(x,y)=0

Step 7: find the gradient magnitude, Gmag, and the gradient direction, Gdir, for binary image

Img4 using prewitt method

Step 8: set Img5 with gradient magnitude Gmag

Step 9: clear pixels on the border of the image Img5 using 8-connectivity and set new image

To Img6

Step 10: fill holes in the binary image Img6 and set to Img7

Step 11: Calculate the global threshold of Img7 and set to gt1

%%If g(x, y) is a thresholded version of f(x, y) at global threshold gt1 %%%%

Step 12: convert grey image Img7 to binary and set binary image to Img8

if f(x,y)> gt then g(x,y)=1

else

g(x,y)=0

Step 13: create disk-shaped structure of radius 3 and set to se1

Step 14: create disk-shaped structure of radius 3 and set to se2

Step 15: apply the binary erosion of Img8 by se1 and return the eroded binary image to Img9

%%The *binary erosion* of *Img8* by *se1*, denoted  $Img8 \ominus se1$ , is defined as the set operation  $Img8 \ominus se1 = \{z | (se1_z \subseteq Img8)\}$ .

Step 16: apply the binary dilation of Img9 by se2 and return the eroded binary image to Img10

%%The *binary dilation* of *A* by *B*, denoted  $A \oplus B$ , is defined as the set operation %%%5

%%%%  $A \oplus B = \{z | (\cap B_z) \cap A \neq \emptyset\}$  %%%

Step 17: multiply each element in array Img4 by the corresponding element in array Img10 and return the product output array to Img11.

Step 18: remove all objects of Img11 which is more than 2100 pixels using 8-connectivity  
And set to Img12

Step 19: remove all objects of Img12 which is less than 50 pixels using 8-connectivity and  
Set to Img13

Step 20: apply morphological opening operation to the binary image Img13 and set the product to Img14 %%% erosion followed by dilation %%%

Step 21: apply morphological closing operation to the binary image Img14 and set the product to Img15 %%%dilation followed by erosion%%%

Step 22: remove all objects of Img15 which is less than 50 pixels using 8-connectivity and  
Set to Img16

Step 23: Obtain features by regionprops & graycoprops functions & assign to  
FeatureListTestALL  
Contrast, Correlation, Energy, Homogeneity, Area, MajorAxisLength,  
MinorAxisLength, Eccentricity, ConvexArea, EquivDiameter, Solidity, Extent,  
Perimeter, CentroidX, CentroidY

Step 24: Read NoduleNonNoduleTable for training and set to NoduleTrnData as one column

Step 25: Transpose FeatureListTestALL to TrainDataSet to create one row for each object

Step 26: Train the system by using svmtrain function with RBF kernel and feed  
TrainDataSetFeature and NoduleTrnData and assign to svmStruct

Step 27 Stop

Step 1: Start %%%PRE-PROCESSING +SEGMENTATION+ TESTING %%%%

Step 2: Read the image data from DICOM file and set to Img1

Step 3: Apply median filtering to Img1 and set new image to Img2

Step 4: Adjust/ increase the contrast of Img2 and set to Img3

Step 5: Calculate the global threshold of Img3 and set to gt  
%%If  $g(x, y)$  is a thresholded version of  $f(x, y)$  at global threshold  $gt$  %%%

Step 6: convert grey image Img3 to binary and set binary image to Img4  
if  $f(x,y) > gt$  then  $g(x,y)=1$   
else  
 $g(x,y)=0$

Step 7: find the gradient magnitude, Gmag, and the gradient direction, Gdir, for binary image  
Img4 using prewitt method

Step 8: set Img5 with gradient magnitude Gmag

Step 9: clear pixels on the border of the image Img5 using 8-connectivity and set new image  
To Img6

Step 10: fill holes in the binary image Img6 and set to Img7

Step 11: Calculate the global threshold of Img7 and set to gt1  
%%If g(x, y) is a thresholded version of f(x, y) at global threshold gt1 %%%%

Step 12: convert grey image Img7 to binary and set binary image to Img8  
if f(x,y)> gt then g(x,y)=1  
else  
g(x,y)=0

Step 13: create disk-shaped structure of radius 3 and set to se1

Step 14: create disk-shaped structure of radius 3 and set to se2

Step 15: apply the binary erosion of Img8 by se1 and return the eroded binary image to Img9  
%%The *binary erosion* of *Img8* by *se1*, denoted  $Img8 \ominus se1$ , is defined as the set  
operation  $Img8 \ominus se1 = \{z | (se1)_z \subseteq Img8\}$ .

Step 16: apply the binary dilation of Img9 by se2 and return the eroded binary image to Img10  
%%The *binary dilation* of *A* by *B*, denoted  $A \oplus B$ , is defined as the set operation %%%5  
%%%%  $A \oplus B = \{z | (\hat{B}_z) \cap A \neq \emptyset\}$  %%%

Step 17: multiply each element in array Img4 by the corresponding element in  
array Img10 and return the product output array to Img11.

Step 18: remove all objects of Img11 which is more than 2100 pixels using 8-connectivity  
And set to Img12

Step 19: remove all objects of Img12 which is less than 50 pixels using 8-connectivity and  
Set to Img13

Step 20: apply morphological opening operation to the binary image Img13 and set the  
product to Img14 %%% erosion followed by dilation %%%

Step 21: apply morphological closing operation to the binary image Img14 and set the  
product to Img15 %%%dilation followed by erosion%%%

Step 22: remove all objects of Img15 which is less than 50 pixels using 8-connectivity and  
Set to Img16

Step 23: Obtain features by regionprops & graycoprops functions & assign to  
FeatureListTestALL  
Contrast, Correlation, Energy, Homogeneity, Area, MajorAxisLength,  
MinorAxisLength, Eccentricity, ConvexArea, EquivDiameter, Solidity, Extent,  
Perimeter, CentroidX, CentroidY

Step 24: Transpose FeatureListTestALL to TestDataSet to create one row for each object

Step 25: select features from TestDataSet and set to testDataMatrix

Step 26: Classify the test data by using svmclassify function and feed svmStruct

And testDataMatrix

Step 27: Stop

```
%%% TRAIN PROCESS %%%%%%%%%%

%%% Importing Original Image %%%%

cd('D:\Nodules\Train');
imagelist = dir('*.dcm');
num = numel(imagelist);
imdata = cell(1,numel(imagelist));

ImgSizeTmp = size(imdata);

ImgSize = ImgSizeTmp(2);
index = 0;

for m=1:ImgSize,
    DispMessage1 = ['Image Number: ',num2str(m)];
    disp(DispMessage1);
    imdata{m} = dicomread(imagelist(m).name);    %% Read DICOM imagelist
    ImportedImg = imdata{m};

    Img1 = ImportedImg;
    %figure(1), imshow(Img1,[]), title('Image 1');    %% Assign DICOM image to Img1

    Img2 = medfilt2(Img1); %% Median Filter Noise Removal and For Edge Detection
    %figure(2), imshow(Img2,[]), title('Image 2');

    Img3 = imadjust(Img2);    %% Contrast Strecting For Thresholding
    %figure(3), imshow(Img3,[]), title('Image 3');

    gt = graythresh(Img3);    %% Thresholding using graythresh
    Img4 = im2bw(Img3,gt);
    %figure(4);imshow(Img4,[]);title('Image 4');

    [Gmag, Gdir] = imgradient(Img4,'prewitt');    %% Prewitt Edge Detection
    Img5 = Gmag;
    %figure(5);imshow(Img5,[]);title('Image 5');

    Img6 = imclearborder(Img5);    %% Clear borders
    %figure(6);imshow(Img6,[]);title('Image 6');
```

```

Img7 = imfill(Img6,'holes');          %% Fills holes
figure(7);imshow(Img7,[]);title('Image 7');

gt = graythresh(Img7);                %% Thresholding using graythresh
Img8 = im2bw(Img7,gt);
figure(8);imshow(Img8,[]);title('Image 8');
se1 = strel('disk',3);
se2 = strel('disk',3);
Img9 = imerode(Img8,se1);
figure(9);imshow(Img9,[]);title('Image 9');

Img10 = imdilate(Img9,se2);
figure(10);imshow(Img10,[]);title('Image 10');

Img11 = immultiply(Img4,Img10);
figure(11);imshow(Img11,[]);title('Image 11');

LB = 0; UB = 2100;
Img12 = xor(bwareaopen(Img11,LB), bwareaopen(Img11,UB)); %% Removes pixels greater than
2100
figure(12);imshow(Img12,[]);title('Image 12');

Img13 = bwareaopen(Img12,15); %%% Adjust According to Database Nodule Areas  %%
Removes pixels less than 15
figure(13);imshow(Img13,[]);title('Image 13');

Img14 = bwmorph(Img13,'open');
%% Remove lung structure lines using opening and closing.
figure(14);imshow(Img14,[]);title('Image 14');

Img15 = bwmorph(Img14,'close');
figure(15);imshow(Img15,[]);title('Image 15'); %% Remove lung structure lines using opening
and closing.

Img16 = bwareaopen(Img15,15); %%% Adjust According to Database Nodule Areas  %%
Removes pixels less than 15
figure(16);imshow(Img16,[]);title('Image 13');

RegPropStat = regionprops(Img16,'All');
%% Regionprop Properties for ROI, objects after segmentation
figure(99);imshow(RegPropStat(1).Image,[]);title('Image 17');

SizeTemp = size(RegPropStat);
%% Number of objects after segmentation
RegPropSize = SizeTemp(1);
RegPropSizeArr(m) = RegPropSize;

GLCMPropStat = struct('Contrast',0,'Correlation',0,'Energy',0,'Homogeneity',0); %%
Initialize GLCMPropStat Structure with Fields

for i=1:RegPropSize

```

```

        GLCMPPropStat(i) = graycoprops((RegPropStat(i).Image),'All');
        %% GLCMCOProps Properties for ROI, objects after segmentation
    end

    if RegPropSize == 0

        FeatureListTestALL(j+index).LogicalRef = j+index;
        FeatureListTestALL(j+index).ImageNo = m;
        FeatureListTestALL(j+index).ImageFileName = imagelist(m).name;
        FeatureListTestALL(j+index).ObjectNumber = 0;
        FeatureListTestALL(j+index).Contrast = 0;
        FeatureListTestALL(j+index).Correlation = 0;
        FeatureListTestALL(j+index).Energy = 0;
        FeatureListTestALL(j+index).Homogeneity = 0;
        FeatureListTestALL(j+index).Area = 0;
        FeatureListTestALL(j+index).MajorAxisLength = 0;
        FeatureListTestALL(j+index).MinorAxisLength = 0;
        FeatureListTestALL(j+index).Eccentricity = 0;
        FeatureListTestALL(j+index).ConvexArea = 0;
        FeatureListTestALL(j+index).EquivDiameter = 0;
        FeatureListTestALL(j+index).Solidity = 0;
        FeatureListTestALL(j+index).Extent = 0;
        FeatureListTestALL(j+index).Perimeter = 0;
        FeatureListTestALL(j+index).CentroidX = 0;
        FeatureListTestALL(j+index).CentroidY = 0;

        index = index + 1;
    end

    if RegPropSize > 0

        for j=1:RegPropSize
            FeatureListALL(j+index).LogicalRef = j+index;
            FeatureListALL(j+index).ImageNo = m;
            FeatureListALL(j+index).ImageFileName = imagelist(m).name;
            FeatureListALL(j+index).ObjectNumber = j;
            FeatureListALL(j+index).Contrast = GLCMPPropStat(j).Contrast;
            FeatureListALL(j+index).Correlation = GLCMPPropStat(j).Correlation;
            FeatureListALL(j+index).Energy = GLCMPPropStat(j).Energy;
            FeatureListALL(j+index).Homogeneity = GLCMPPropStat(j).Homogeneity;
            FeatureListALL(j+index).Area = RegPropStat(j).Area;
            FeatureListALL(j+index).MajorAxisLength = RegPropStat(j).MajorAxisLength;
            FeatureListALL(j+index).MinorAxisLength = RegPropStat(j).MinorAxisLength;
            FeatureListALL(j+index).Eccentricity = RegPropStat(j).Eccentricity;
            FeatureListALL(j+index).ConvexArea = RegPropStat(j).ConvexArea;
            FeatureListALL(j+index).EquivDiameter = RegPropStat(j).EquivDiameter;
            FeatureListALL(j+index).Solidity = RegPropStat(j).Solidity;
            FeatureListALL(j+index).Extent = RegPropStat(j).Extent;
            FeatureListALL(j+index).Perimeter = RegPropStat(j).Perimeter;
            FeatureListALL(j+index).CentroidX = RegPropStat(j).Centroid(1);
            FeatureListALL(j+index).CentroidY = RegPropStat(j).Centroid(2);
        end
    end

```

```

end

index = j + index;

end
disp('End of Image Segmentation'); disp(' ');

[~, ~, NoduleNonNoduleTable] = xlsread('D:\Nodules\Nodule-NonNoduleTable.xlsx','Sheet1');
NoduleNonNoduleTable(cellfun(@(x) ~isempty(x) && isnumeric(x) &&
    isnan(x),NoduleNonNoduleTable)) = {''}; %% Reads the documented nodules table
DataSize = size(NoduleNonNoduleTable);
DataSize = DataSize(1);
Status = NoduleNonNoduleTable(:,2);
BinStat= NoduleNonNoduleTable(:,3);
for h=1:DataSize
    NoduleTrnData(h)=Status(h);
    NoduleTrnDataBin(h)=BinStat(h);
end

for z=1:index

    TrainDataSet(z,1) = FeatureListALL(z).Contrast;
    TrainDataSetNew(z,1) = FeatureListALL(z).Contrast;

    TrainDataSet(z,2) = FeatureListALL(z).Correlation;
    TrainDataSetNew(z,2) = FeatureListALL(z).Correlation;

    TrainDataSet(z,3) = FeatureListALL(z).Energy;
    TrainDataSetNew(z,3) = FeatureListALL(z).Energy;

    TrainDataSet(z,4) = FeatureListALL(z).Homogeneity;
    TrainDataSetNew(z,4) = FeatureListALL(z).Homogeneity;

    TrainDataSet(z,5) = FeatureListALL(z).Area;
    TrainDataSetNew(z,5) = FeatureListALL(z).Area;

    TrainDataSet(z,6) = FeatureListALL(z).MajorAxisLength;
    TrainDataSetNew(z,6) = FeatureListALL(z).MajorAxisLength;

    TrainDataSet(z,7) = FeatureListALL(z).MinorAxisLength;
    TrainDataSetNew(z,7) = FeatureListALL(z).MinorAxisLength;

    TrainDataSet(z,8) = FeatureListALL(z).Eccentricity;
    TrainDataSetNew(z,8) = FeatureListALL(z).Eccentricity;

    TrainDataSet(z,9) = FeatureListALL(z).ConvexArea;
    TrainDataSetNew(z,9) = FeatureListALL(z).ConvexArea;

    TrainDataSet(z,10) = FeatureListALL(z).EquivDiameter;
    TrainDataSetNew(z,10) = FeatureListALL(z).EquivDiameter;

```



```

TrainDataSet(z,11) = FeatureListALL(z).Solidity;
TrainDataSetNew(z,11) = FeatureListALL(z).Solidity;

TrainDataSet(z,12) = FeatureListALL(z).Extent;
TrainDataSetNew(z,12) = FeatureListALL(z).Extent;

TrainDataSet(z,13) = FeatureListALL(z).Perimeter;
TrainDataSetNew(z,13) = FeatureListALL(z).Perimeter;

TrainDataSet(z,14) = cell2mat(NoduleTrnDataBin(z));

end

TrainDataSetFeature = TrainDataSet(:,5:13);

figure(20);
svmStruct = svmtrain(TrainDataSetFeature,NoduleTrnData,'Kernel_Function','rbf','RBF_Sigma',
    0.1, 'BoxConstraint', 1,'method','SMO','showplot',true);

%% %% TEST PROCESS %% %% %% %% %% %% %% %% %% %% %%

cd('D:\Nodules\Test');
imagelist = dir('*.dcm');
num = numel(imagelist);
imdata = cell(1,numel(imagelist));

ImgSizeTmp = size(imdata);

ImgSize = ImgSizeTmp(2);
index = 0;

for m=1:ImgSize
DispMessage1 = ['Image Number: ',num2str(m)];
disp(DispMessage1);
imdata{m} = dicomread(imagelist(m).name);          %% Read DICOM imagelist
ImportedImg = imdata{m};

Img1 = ImportedImg;                                %% Read Original Image
%figure(1), imshow(Img1,[]), title('Image 1');

Img2 = medfilt2(Img1); %% Median Filter Noise Removal and For Edge Detection
%figure(2), imshow(Img2,[]), title('Image 2');

Img3 = imadjust(Img2);          %% Contrast Strechting For Thresholding
%figure(3), imshow(Img3,[]), title('Image 3');

gt = graythresh(Img3);          %% Thresholding using graythresh
Img4 = im2bw(Img3,gt);
%figure(4);imshow(Img4,[]);title('Image 4');

```

```

[Gmag, Gdir] = imgradient(Image4,'prewitt');      %% Prewitt Edge Detection
Image5 = Gmag;
figure(5);imshow(Image5,[]);title('Image 5');

Image6 = imclearborder(Image5);                  %% Clear borders
figure(6);imshow(Image6,[]);title('Image 6');

Image7 = imfill(Image6,'holes');                 %% Fills holes
figure(7);imshow(Image7,[]);title('Image 7');

gt = graythresh(Image7);                        %% Thresholding using graythresh
Image8 = im2bw(Image7,gt);
figure(8);imshow(Image8,[]);title('Image 8');

se1 = strel('disk',3);
se2 = strel('disk',3);
Image9 = imerode(Image8,se1);
figure(9);imshow(Image9,[]);title('Image 9');

Image10 = imdilate(Image9,se2);
figure(10);imshow(Image10,[]);title('Image 10');

Image11 = immultiply(Image4,Image10);           %% Gets lung structure from IMG4
figure(11);imshow(Image11,[]);title('Image 11');

LB = 0; UB = 2100;
Image12 = xor(bwareaopen(Image11,LB), bwareaopen(Image11,UB));    %% Removes pixels
                                                                    greater than 2100
figure(12);imshow(Image12,[]);title('Image 12');

Image13 = bwareaopen(Image12,15); %%% Adjust According to Database Nodule Areas  %%
                                                                    Removes area less than 15
figure(13);imshow(Image13,[]);title('Image 13');

Image14 = bwmorph(Image13,'open');
                                                                    %% Remove lung structure lines using opening and closing.
figure(14);imshow(Image14,[]);title('Image 14');

Image15 = bwmorph(Image14,'close');
figure(15);imshow(Image15,[]);title('Image 15'); %% Remove lung structure lines using opening
                                                                    and closing.

Image16 = bwareaopen(Image15,15); %%% Adjust According to Database Nodule Areas  %%
                                                                    Removes pixels less than 15
figure(16);imshow(Image16,[]);title('Image 13');

RegPropStat = regionprops(Image16,'All');
                                                                    %% Regionprop Properties for ROI, objects after segmentation
figure(99);imshow(RegPropStat(1).Image,[]);title('Image 17');

SizeTemp = size(RegPropStat);
                                                                    %% Number of objects after segmentation

```

```

RegPropSize = SizeTemp(1);
RegPropSizeArr(m) = RegPropSize;

GLCMPropStat = struct('Contrast',0,'Correlation',0,'Energy',0,'Homogeneity',0);    %%
    Initialize GLCMPropStat Structure with Fields

for i=1:RegPropSize
    GLCMPropStat(i) = graycoprops((RegPropStat(i).Image),'All');
    %% GLCMCOProps Properties for ROI, objects after segmentation
end

if RegPropSize == 0

    FeatureListTestALL(j+index).LogicalRef = j+index;
    FeatureListTestALL(j+index).ImageNo = m;
    FeatureListTestALL(j+index).ImageFileName = imagelist(m).name;
    FeatureListTestALL(j+index).ObjectNumber = 0;
    FeatureListTestALL(j+index).Contrast = 0;
    FeatureListTestALL(j+index).Correlation = 0;
    FeatureListTestALL(j+index).Energy = 0;
    FeatureListTestALL(j+index).Homogeneity = 0;
    FeatureListTestALL(j+index).Area = 0;
    FeatureListTestALL(j+index).MajorAxisLength = 0;
    FeatureListTestALL(j+index).MinorAxisLength = 0;
    FeatureListTestALL(j+index).Eccentricity = 0;
    FeatureListTestALL(j+index).ConvexArea = 0;
    FeatureListTestALL(j+index).EquivDiameter = 0;
    FeatureListTestALL(j+index).Solidity = 0;
    FeatureListTestALL(j+index).Extent = 0;
    FeatureListTestALL(j+index).Perimeter = 0;
    FeatureListTestALL(j+index).CentroidX = 0;
    FeatureListTestALL(j+index).CentroidY = 0;

    index = index + 1;

end

if RegPropSize > 0

    for j=1:RegPropSize
        FeatureListTestALL(j+index).LogicalRef = j+index;
        FeatureListTestALL(j+index).ImageNo = m;
        FeatureListTestALL(j+index).ImageFileName = imagelist(m).name;
        FeatureListTestALL(j+index).ObjectNumber = j;
        FeatureListTestALL(j+index).Contrast = GLCMPropStat(j).Contrast;
        FeatureListTestALL(j+index).Correlation = GLCMPropStat(j).Correlation;
        FeatureListTestALL(j+index).Energy = GLCMPropStat(j).Energy;
        FeatureListTestALL(j+index).Homogeneity = GLCMPropStat(j).Homogeneity;
        FeatureListTestALL(j+index).Area = RegPropStat(j).Area;
        FeatureListTestALL(j+index).MajorAxisLength =
RegPropStat(j).MajorAxisLength;
        FeatureListTestALL(j+index).MinorAxisLength =
RegPropStat(j).MinorAxisLength;
    end
end

```

```

        FeatureListTestALL(j+index).Eccentricity = RegPropStat(j).Eccentricity;
        FeatureListTestALL(j+index).ConvexArea = RegPropStat(j).ConvexArea;
        FeatureListTestALL(j+index).EquivDiameter = RegPropStat(j).EquivDiameter;
        FeatureListTestALL(j+index).Solidity = RegPropStat(j).Solidity;
        FeatureListTestALL(j+index).Extent = RegPropStat(j).Extent;
        FeatureListTestALL(j+index).Perimeter = RegPropStat(j).Perimeter;
        FeatureListTestALL(j+index).CentroidX = RegPropStat(j).Centroid(1);
        FeatureListTestALL(j+index).CentroidY = RegPropStat(j).Centroid(2);
    end
end

index = j + index;

end
disp('End of Image Segmentation'); disp(' ');

[~, ~, NoduleNonNoduleTableTest] = xlsread('D:\Nodules\NoduleNon-
    NoduleTable.xlsx','Sheet1');
NoduleNonNoduleTableTest(cellfun(@(x) ~isempty(x) && isnumeric(x) &&
    isnan(x),NoduleNonNoduleTableTest)) = {''}; %% Read documented nodule table for test
DataSize = size(NoduleNonNoduleTableTest);
DataSize = DataSize(1);
Status = NoduleNonNoduleTableTest(:,2);
BinStat= NoduleNonNoduleTableTest(:,3);
for h=1:DataSize
    NoduleTestData(h)=Status(h);
    NoduleTestDataBin(h)=BinStat(h);
end

for z=1:index

    TestDataSet(z,1) = FeatureListTestALL(z).Contrast;
    TestDataSetNew(z,1) = FeatureListTestALL(z).Contrast;

    TestDataSet(z,2) = FeatureListTestALL(z).Correlation;
    TestDataSetNew(z,2) = FeatureListTestALL(z).Correlation;

    TestDataSet(z,3) = FeatureListTestALL(z).Energy;
    TestDataSetNew(z,3) = FeatureListTestALL(z).Energy;

    TestDataSet(z,4) = FeatureListTestALL(z).Homogeneity;
    TestDataSetNew(z,4) = FeatureListTestALL(z).Homogeneity;

    TestDataSet(z,5) = FeatureListTestALL(z).Area;
    TestDataSetNew(z,5) = FeatureListTestALL(z).Area;

    TestDataSet(z,6) = FeatureListTestALL(z).MajorAxisLength;
    TestDataSetNew(z,6) = FeatureListTestALL(z).MajorAxisLength;

    TestDataSet(z,7) = FeatureListTestALL(z).MinorAxisLength;
    TestDataSetNew(z,7) = FeatureListTestALL(z).MinorAxisLength;

```

```

    TestDataSet(z,8) = FeatureListTestALL(z).Eccentricity;
    TestDataSetNew(z,8) = FeatureListTestALL(z).Eccentricity;

    TestDataSet(z,9) = FeatureListTestALL(z).ConvexArea;
    TestDataSetNew(z,9) = FeatureListTestALL(z).ConvexArea;

    TestDataSet(z,10) = FeatureListTestALL(z).EquivDiameter;
    TestDataSetNew(z,10) = FeatureListTestALL(z).EquivDiameter;

    TestDataSet(z,11) = FeatureListTestALL(z).Solidity;
    TestDataSetNew(z,11) = FeatureListTestALL(z).Solidity;

    TestDataSet(z,12) = FeatureListTestALL(z).Extent;
    TestDataSetNew(z,12) = FeatureListTestALL(z).Extent;

    TestDataSet(z,13) = FeatureListTestALL(z).Perimeter;
    TestDataSetNew(z,13) = FeatureListTestALL(z).Perimeter;

    TestDataSet(z,14) = cell2mat(NoduleTestDataBin(z));

end
RowSize=size(TestDataSetNew);
RowSize=RowSize(1);
for q=1:RowSize
    result(q) = trainedClassifier.predict(TestDataSetNew(q,:));
    FinalResults(1,q)=str2double(NoduleTestDataBin(q));
    FinalResults(2,q)=result(q);
end

FinalResultsTest = transpose(FinalResults);

x = index;

for j=1:x
    testDataMatrix=[TestDataSet(j,11),TestDataSet(j,12)];
    LungCancerTestList = svmclassify(svmStruct,testDataMatrix,'showplot',true);
    hold on;
    plot(testDataMatrix(1),testDataMatrix(2),'ro','MarkerSize',12);
    hold off;
    LungCancerTestList1(j)= LungCancerTestList(1,1);
end
figure(20);

TP=0;
TN=0;
FP=0;
FN=0;
SizeFinalResultsTest=size(FinalResultsTest);
row=SizeFinalResultsTest(1);
col=SizeFinalResultsTest(2);
for statistics = 1:row

```

```

        if FinalResultsTest(statistics,1)==1
        if FinalResultsTest(statistics,2)==1
            TP=TP+1;
        else
            FN=FN+1;
        end
        elseif FinalResultsTest(statistics,1)==0
        if FinalResultsTest(statistics,2)==0
            TN=TN+1;
        else
            FP=FP+1;
        end
    end
end

Specificity = (TN / (FP+TN))*100;
Sensitivity = (TP / (TP+FN))*100;
Accuracy = ((TP+TN)/row)*100;
end

```

## ABSTRACT

Computer aided diagnosis is starting to be implemented broadly in the diagnosis and detection of many varieties of abnormalities acquired during various imaging procedures. The main aim of the CAD systems is to increase the accuracy and decrease the time of diagnoses, while the general achievement for CAD systems are to find the place of nodules and to determine the characteristic features of the nodule. As lung cancer is one of the fatal and leading cancer types, there has been plenty of studies for the usage of the CAD systems to detect lung cancer. Yet, the CAD systems need to be developed a lot in order to identify the different shapes of nodules, lung segmentation and to have higher level of sensitivity, specificity and accuracy. This challenge is the motivation of this study in implementation of CAD system for lung cancer detection. In the study, LIDC database is used which comprises of an image set of lung cancer thoracic documented CT scans. The presented CAD system consists of CT image reading, image pre-processing, segmentation, feature extraction and classification steps. To avoid losing important features, the CT images were read as a raw form in DICOM file format. Then, filtration and enhancement techniques were used as an image processing. Otsu's algorithm, edge detection and morphological operations are applied for the segmentation, following the feature extractions step. Finally, support vector machine with Gaussian RBF is utilized for the classification step which is widely used as a supervised classifier.

**Keywords:** CAD systems; lung cancer; image pre-processing; segmentation; feature extraction; classification; global threshold; support vector machines; SVM; ANN

## ÖZET

Bilgisayar destekli tanı sistemleri, farklı muayene işlemlerinden elde edilen medikal görüntülerdeki çeşitli anomalilere tanı koyma ve ortaya çıkarmada yaygın olarak kullanılmaya başlanmıştır. BDT sistemleri için genel başarı nodüllerin yerlerini bulmak ve nodülün karakteristik özelliklerini belirleme ile ölçülürken, BDT sistemlerinin temel amacı doğruluk oranını artırmak ve tanı süresini azaltmaktır. Akciğer kanseri, en çok görülen ölümcül kanser türlerinden biri olduğu için akciğer kanserini tespit etmek için geliştirilen BDT sistemlerine yönelik birçok çalışma yapılmıştır. Buna rağmen BDT sistemleri, farklı şekillerdeki nodüllerin algılanmasında, akciğer görüntüsünün segmentasyonunda, yüksek düzeyde duyarlılık, özgüllük ve doğruluk değerlerinin elde edilmesinde yetersiz kalıyor.

Akciğer kanseri tespiti için geliştirilen BDT sisteminin yapılan çalışmadaki motivasyonunu bu yetersizlikler oluşturmaktadır. Çalışmada LIDC veri tabanındaki düşük dozda çekilmiş hastaların dokümanede edilmiş göğüs BT görüntüleri kullanılmıştır. Sunulan BDT sistemi; BT görüntüsünden okuma, görüntü önileme, segmentasyon, öznitelik çıkarma ve sınıflandırma adımlarından oluşmuştur. Önemli özniteliklerin kaybını önlemek için, BT görüntülerinin işlenmemiş halleri, yani DICOM, dosya biçiminde okunmuştur. Daha sonra, görüntü geliştirme teknikleri ve görüntü işleme tekniklerinden filtreleme kullanıldı. Segmentasyon için Otsu algoritması, kenar algılama ve morfolojik operasyonlar kullanıldı. Bunları öznitelik çıkartma adımı takip etti. Son olarak, sınıflandırma aşaması için yaygın bir denetimli sınıflandırıcı olarak kullanılan Destek Vektör Makinesinin Gauss RBF kerneli kullanılmıştır.

**Anahtar Kelimeler:** BDT sistemleri; akciğer kanseri; görüntü önileme; segmentasyon; öznitelik çıkarma; sınıflandırma; global threshold; destek vektör makinesi; SVM; ANN