

A MASS DETECTION ALGORITHM FOR MAMMOGRAM IMAGES

A THESIS SUBMITTED TO  
THE GRADUATE SCHOOL OF NATURAL AND APPLIED SCIENCES  
OF  
MIDDLE EAST TECHNICAL UNIVERSITY

BY

MUHAMMED YEŞİLKAYA

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS  
FOR  
THE DEGREE OF MASTER OF SCIENCE  
IN  
ELECTRICAL AND ELECTRONICS ENGINEERING

SEPTEMBER 2014



Approval of the thesis:

**A MASS DETECTION ALGORITHM FOR MAMMOGRAM IMAGES**

submitted by **MUHAMMED YEŞİLKAYA** in partial fulfillment of the requirements for the degree of **Master of Science in Electrical and Electronics Engineering Department, Middle East Technical University** by,

Prof. Dr. Canan Özgen \_\_\_\_\_  
Dean, Graduate School of **Natural and Applied Sciences**

Prof. Dr. Gönül Turhan Sayan \_\_\_\_\_  
Head of Department, **Electrical and Electronics Engineering**

Prof. Dr. Gözde Bozdağı Akar \_\_\_\_\_  
Supervisor, **Electrical and Electronics Eng., Dept., METU**

Prof. Dr. Nevzat Güneri Genç \_\_\_\_\_  
Co-supervisor, **Electrical and Elect., Eng., Dept., METU**

**Examining Committee Members:**

Assoc. Prof. Dr. Yeşim Serianağaoğlu \_\_\_\_\_  
Electrical and Electronics Engineering Dept., METU

Prof. Dr. Gözde Bozdağı Akar \_\_\_\_\_  
Electrical and Electronics Engineering Dept., METU

Prof. Dr. Nevzat Güneri Genç \_\_\_\_\_  
Electrical and Electronics Engineering Dept., METU

Assoc. Prof. Dr. Ali Rıza Sever \_\_\_\_\_  
Radiology Dept., Hacettepe University

Assist. Prof. Dr. Fatih Kamışlı \_\_\_\_\_  
Electrical and Electronics Engineering Dept., METU

**Date:** \_\_\_\_\_

**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: MUHAMMED YEŞİLKAYA

Signature :

# ABSTRACT

A MASS DETECTION ALGORITHM FOR MAMMOGRAM IMAGES

YEŞİLKAYA, MUHAMMED

M.S., Department of Electrical and Electronics Engineering

Supervisor : Prof. Dr. Gözde Bozdağı Akar

Co-Supervisor : Prof. Dr. Nevzat Güneri Gençer

September 2014, 77 pages

Breast cancer is the most common cancer type encountered among woman in the world and causes many deaths. In order to prevent mastectomies, decrease the probability of return and reduce mortality, early detection of cancer lesion is crucial. Mammography is a frequently used screening technique to detect and diagnose lesions. However, sometimes it is difficult for radiologists to see and diagnose lesions due to low contrast of mammograms. Computer Aided Detection / Diagnosis (CAD / CADx) systems have been developed to help radiologists.

In this thesis, we propose a method for classification of mass regions in MLO (Mediolateral oblique) view mammograms. The suspicious regions are first determined by Iris filtering with variable window sizes applied on the breast region without pectoral muscle. Then classification is applied to textural features obtained using Gabor filter applied on these suspicious regions. We reduced false detection ratio nearly 50 percent with a cost of missing 9 percent of true mass regions with classification. For pectoral muscle region determination a novel algorithm is also proposed. This algorithm is based on average derivative calculation and line fitting with least square solution. Our algorithm outperforms other algorithms given in the literature in terms of FP (False positive) pixel percentage and FN (False negative) pixel percentage metrics.

Keywords: Mass Detection in Mammogram Images, Iris Filter, Gabor Filter Bank

# ÖZ

## MAMMOGRAM GÖRÜNTÜLERİ İÇİN BİR KİTLE TESPİT ALGORİTMASI

YEŞİLKAYA, MUHAMMED

Yüksek Lisans, Elektrik ve Elektronik Mühendisliği Bölümü

Tez Yöneticisi : Prof. Dr. Gözde Bozdağı Akar

Ortak Tez Yöneticisi : Prof. Dr. Nevzat Güneri Gençer

Eylül 2014 , 77 sayfa

Meme kanseri, dünyada kadınlar arasında en yaygın kanser tipidir ve ölümlere neden olmaktadır. Meme ameliyatlarını önlemek, nüksetme olasılığını düşürmek ve ölüm oranını azaltmak için kanser lezyonunun erken tespiti çok önemlidir. Mamografi, lezyon tespit ve teşhisinde sıkça kullanılan bir görüntüleme tekniğidir; fakat meme filmlerinin düşük kontrastından, lezyonların tespit ve teşhisi bazen radyologlar için zordur. Bu sebepten radyolaglara yardım etmek için yakın geçmişte, Bilgisayar Destekli Tespit / Teşhis sistemleri (CAD / CADx) geliştirilmiştir.

Bu tezde, MLO (yandan) görüntülü mammogramlarda kitle bölgelerinin sınıflandırılması için bir yöntem sunmaktayız. Öncelikle, göğüs kası olmayan meme bölgesine değişik pencere boyutlarıyla uygulanan Iris süzgeciyle şüpheli bölgeler belirlenmektedir. Sonra Gabor filtrenin bu şüpheli bölgelere uygulanması ile elde edilen doku-sal (textural) özelliklere sınıflandırma uygulanmaktadır. Sınıflandırma ile doğru kitle bölgelerinin yüzde dokuzu kaçırılmasına karşın, yanlış tespit edilen bölgeleri yaklaşık yüzde elli oranında azaltmaktayız. Ayrıca göğüs kası bölgesinin belirlenmesi için yeni bir algoritma sunulmaktadır. Bu algoritma ortalama türev hesabı ve en küçük kare çözümü ile doğru uydurmaya dayanmaktadır. Bizim algoritmamız, literatürdeki diğer algoritmalarından FP (Yanlış pozitif) piksel yüzdesi ve FN (Yanlış pozitif) piksel yüzdesi açısından daha iyi sonuç vermektedir.

Anahtar Kelimeler: Mammogram Görüntülerinde Kitle Tespiti, Iris Süzgeci, Gabor Süzgeç Kümesi

*To My Mother*

## ACKNOWLEDGMENTS

Firstly, I would like to express my sincere thanks to my supervisor Prof. Dr. Gözde Bozdağı Akar and co-supervisor Prof. Dr.Nevzat Güneri Gençer for their supervision and guidance throughout this study.

I would like to thank Assoc. Dr. Ali Rıza Sever for his contributions to this thesis.

I would like to thank my director Seçkin Öz Saraç in my work for his understanding and his encouragement throughout my thesis study.

I am indebted to my family for their continuous love, encouragement, and support during my academic studies.



# TABLE OF CONTENTS

ABSTRACT . . . . .	v
ÖZ . . . . .	vi
ACKNOWLEDGMENTS . . . . .	viii
TABLE OF CONTENTS . . . . .	ix
LIST OF TABLES . . . . .	xii
LIST OF FIGURES . . . . .	xiv
LIST OF ABBREVIATIONS . . . . .	xxi
CHAPTERS	
1 INTRODUCTION . . . . .	1
1.1 Mammography . . . . .	1
1.2 Breast Lesion Types . . . . .	3
1.2.1 Microcalcifications . . . . .	3
1.2.2 Masses . . . . .	4
1.2.3 Architectural distortions . . . . .	5
1.2.4 Bilateral asymmetry . . . . .	5
1.3 CAD System for Mammograms . . . . .	6
1.3.1 CAD system evaluation . . . . .	6
1.3.2 Available databases . . . . .	7

1.4	Literature Review . . . . .	8
1.4.1	Pectoral muscle segmentation . . . . .	8
1.4.2	Mass detection . . . . .	9
1.4.3	Classification of textural features . . . . .	10
1.5	Contributions . . . . .	10
1.5.1	Pectoral muscle segmentation . . . . .	11
1.5.2	Mass detection . . . . .	11
1.6	Outline of the Thesis . . . . .	12
2	BACKGROUND SUBTRACTION AND PECTORAL MUSCLE REGION SEGMENTATION . . . . .	13
2.1	Image Processing Techniques Used in Background Subtraction	13
2.2	Background Subtraction . . . . .	14
2.3	Pectoral Muscle Segmentation . . . . .	15
2.3.1	Pectoral muscle segmentation algorithm . . . . .	15
2.4	Evaluation of pectoral muscle segmentation algorithm . . . . .	23
2.5	Contrast Enhancement and Mass detection in the Pectoral Regions . . . . .	26
3	DETECTION OF POSSIBLE MASS REGIONS WITH IRIS FILTER	29
3.1	Iris Filter . . . . .	29
3.2	Suspicious ROI Detection Algorithm . . . . .	31
3.2.1	Iris filter implementation . . . . .	35
3.2.2	Outputs of Iris filters . . . . .	37
3.2.3	Designation of suspicious regions . . . . .	40
3.3	Results of Suspicious ROI Detection Algorithm . . . . .	44

4	SVM CLASSIFICATION USING FEATURES, EXTRACTED WITH GABOR FILTERS . . . . .	47
4.1	Feature Extraction with Gabor Filter Bank . . . . .	47
4.1.1	Gabor filter bank . . . . .	48
4.1.2	Partition of ROIs . . . . .	49
4.1.3	Feature extraction . . . . .	50
4.2	Classification . . . . .	51
4.2.1	Determination of SVM parameters . . . . .	52
4.2.2	Classification performance comparison of regions with different dimensions including same masses . . . . .	60
4.2.3	Classification of suspicious regions . . . . .	62
5	RESULTS . . . . .	65
6	CONCLUSION . . . . .	69
	REFERENCES . . . . .	73

## LIST OF TABLES

### TABLES

Table 1.1 Grading of imaging reports of microcalcifications according to risk of malignancy [1]. . . . .	4
Table 1.2 Performance comparison of different textural approaches ([2]). . . . .	10
Table 2.1 Mean and standard deviation of FP, FN pixel percentages for mammograms from mini-MIAS database. . . . .	28
Table 3.1 $R_{min}$ , $R_{max}$ values of Iris filters. . . . .	35
Table 3.2 Mean FPPi for each Iris filter. . . . .	46
Table 4.1 Patch numbers creating sub-windows. . . . .	50
Table 4.2 MIAS statistical data in terms of mass existence. . . . .	53
Table 4.3 ROI selection statistic. . . . .	55
Table 4.4 Statistic of training and test sets. . . . .	55
Table 4.5 Maximum truth rate and Az values for regions with dimensions. . . . .	60
Table 4.6 Number of regions in training set. . . . .	63
Table 5.1 Chosen SVM parameters. . . . .	65
Table 5.2 Mean FPPi after classification with respect to Iris filters. . . . .	66
Table 5.3 Mean FPPi after classification with respect to region sizes. . . . .	66

Table 5.4 True mass regions detected after classification. . . . . 66

## LIST OF FIGURES

### FIGURES

Figure 1.1 Left:MLO view. Right: CC view [3]. . . . .	2
Figure 1.2 CC mammogram of a breast couple from DDSM [4]. . . . .	2
Figure 1.3 MLO mammogram of a breast couple from DDSM [4]. . . . .	3
Figure 1.4 Microcalcifications in mammograms from mini-MIAS [5]. . . . .	4
Figure 1.5 Mass Shapes in Mammograms [6]. . . . .	4
Figure 1.6 Two malignant mass ROIs cropped from full-sized mammograms with their pathological characteristics from DDSM [4]. . . . .	5
Figure 1.7 ROC curve. . . . .	7
Figure 1.8 Benign circumscribed mass from mini-MIAS. . . . .	7
Figure 1.9 Conventional CAD mass detection algorithm [7]. . . . .	8
Figure 1.10 Proposed mass detection algorithm. . . . .	11
Figure 2.1 Left:Erosion operation and illustration, Right:Dilation operation and illustration [8]. . . . .	14
Figure 2.2 Left:Determination of connected components, Right:Labeling of con- nected components [9]. . . . .	14

Figure 2.3 Background subtraction, Step1:Binary image obtainment after threshold, Step2:Erosion operation to eliminate mammogram label, Step3:Dilation operation to get back regions lost on mammograms after erosion step, Step4:Connected component labeling to discriminate breast region from other unnecessary parts. . . . .	15
Figure 2.4 Pectoral muscle segmentation algorithm. . . . .	16
Figure 2.5 ROI determination. . . . .	17
Figure 2.6 Average derivative calculation for each raw and a typical average derivative plot . . . . .	17
Figure 2.7 Minimum derivative points for each raw in ROI. . . . .	18
Figure 2.8 Minimum derivative points for each raw in ROI after first step moving average operation. . . . .	18
Figure 2.9 Minimum derivative points for each raw in ROI after second step moving average operation. . . . .	19
Figure 2.10 Angle of curvature values for the points, on pectoral muscle boundary.	20
Figure 2.11 Pectoral muscle boundary points after the removal of saliencies. . .	20
Figure 2.12 Imaginary line for horizontal derivative search muscle. . . . .	21
Figure 2.13 Minimum derivative calculation on the horizontal lines passing through imaginary line piece added. . . . .	21
Figure 2.14 Least square solution applied to pixels. . . . .	23
Figure 2.15 Minimum derivative calculation on the horizontal lines passing through imaginary line piece added (Second step). . . . .	23
Figure 2.16 Least square solution applied to pixels (Second step). . . . .	24
Figure 2.17 Minimum derivative calculation on the horizontal lines passing through imaginary line piece added (Third step). . . . .	24

Figure 2.18 Least square solution applied to pixels (Third step). . . . .	25
Figure 2.19 Calculated pectoral muscle boundaries with different line piece sizes for the same mammogram; First row:8 , 12, 16, 24, Second row:32, 36, 52, 60	26
Figure 2.20 Calculated pectoral muscle boundaries for some mammograms, from mini-MIAS database. . . . .	27
Figure 2.21 Mammogram image, taken with a camera. . . . .	27
Figure 3.1 Iris filter definition. . . . .	30
Figure 3.2 Left:Gradient map of the mass region with one local maximum, Right:Mass region. . . . .	31
Figure 3.3 Left:Gradient map of the mass region with two local maximums, Right:Mass region. . . . .	32
Figure 3.4 Iris filter's support region and mass detection. Left:Mass edge falls between $R_{min}$ circle and $R_{max}$ circle. Middle:Mass is smaller than $R_{min}$ circle. Right:Mass region covers $R_{out}$ circle. . . . .	33
Figure 3.5 Left:Pixel of interest on the mammogram, Right:Convergence index map for each pixel in the region of support for chosen pixel of interest ( $R_{min} = 2, R_{max} = 62$ ). . . . .	33
Figure 3.6 Left:Convergence index map for region of support, with $R_{min} = 2$ and $R_{max} = 22$ , Right:Convergence index map for region of support with $R_{min} = 32$ and $R_{max} = 52$ . . . . .	34
Figure 3.7 Left:Convergence index map for region of support, with $R_{min} = 12$ and $R_{max} = 32$ , Right:Pixel of interest and largest convergence index region. . . . .	34
Figure 3.8 Left:Convergence index map for each pixel in the region of support, with $R_{min} = 22$ and $R_{max} = 42$ , for chosen pixel of interest, Right:Pixel of interest and largest convergence index region. . . . .	35
Figure 3.9 Iris filter implementation. . . . .	36



Figure 3.10 Original breast region and its histogram. . . . .	37
Figure 3.11 Histogram equalized breast region and its histogram. . . . .	38
Figure 3.12 Left:Output of Iris filter 1, Right: Threshold level = 0,90 applied to filter output added to original mammogram. . . . .	39
Figure 3.13 Left:Output of Iris filter 2, Right: Threshold level = 0,75 applied to filter output added to original mammogram. . . . .	39
Figure 3.14 Left:Output of Iris filter 3, Right: Threshold level = 0,75 applied to filter output added to original mammogram. . . . .	40
Figure 3.15 Left:Output of Iris filter 4, Right: Threshold level = 0,75 applied to filter output added to original mammogram. . . . .	41
Figure 3.16 Left:Output of Iris filter 5, Right: Threshold level = 0,75 applied to filter output added to original mammogram. . . . .	41
Figure 3.17 Left:Suspicious pixels obtained after filtering and threshold operation for Iris filter 2, Right:Connected component labeling and center determina- tion result. . . . .	42
Figure 3.18 Left:Suspicious pixels obtained after filtering and threshold operation for Iris filter 3, Right:Connected component labeling and center determina- tion result. . . . .	43
Figure 3.19 Region determination for centers obtained by Iris filters 1, 2 cases. .	43
Figure 3.20 Region determination for centers obtained by Iris filters 3, 4, 5 cases.	44
Figure 3.21 Region determination when centers are close to each other. . . . .	45
Figure 3.22 Left:Gradient map of spiculated malignant mass region, Right:Spiculated malignant mass region. . . . .	45
Figure 4.1 Feature extraction stages with Gabor filter bank. . . . .	47

Figure 4.2 Gabor filter bank: filters in the same column have the same orientation, filters in the same row has the same frequency. . . . .	49
Figure 4.3 Left:Combined frequency response of Gabor filters, Right: Combined frequency response of Gabor filters without frequency shift. . . . .	50
Figure 4.4 Left:Patches and sub-windows for ROIs with size 128 x 128, Right:Patches and sub-windows for ROIs with size 256 x 256. . . . .	51
Figure 4.5 Left:SVM illustration, Right:ROC curve obtained by adding offsets to hyper-plane. . . . .	53
Figure 4.6 Selected regions of dimension 128 x 128 with mass. . . . .	54
Figure 4.7 Selected regions of dimension 128 x 128 without mass. . . . .	54
Figure 4.8 Selected regions of dimension 256 x 256 with mass. . . . .	54
Figure 4.9 Selected regions of dimension 256 x 256 without mass. . . . .	54
Figure 4.10 Left:Az values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Right: Az values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 1 run. . . . .	55
Figure 4.11 Left:Truth rate values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Right:Truth rate values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 1 run. . . . .	56
Figure 4.12 Left:Sensitivity values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Right:Sensitivity values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 1 run. . . . .	56

Figure 4.13 (1 - specificity) values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Right:(1 - specificity) values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 1 run. . . . . 57

Figure 4.14 Top Left:Sensitivity values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Top Right:(1 - specificity) values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Bottom:Performance points for different SVM parameters. . . . . 58

Figure 4.15 Left:Mean Az values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 100 runs, Right:Mean Az values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 100 runs. . . . . 58

Figure 4.16 Left:Mean truth rate values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 100 runs, Right:Mean truth rate values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 100 runs. . . . . 59

Figure 4.17 Left:Mean sensitivity values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 100 runs, Right:Mean sensitivity values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 100 runs. . . . . 59

Figure 4.18 Mean (1 - specificity) values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after one run, Right:Mean (1 - specificity) values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after one run. . . . . 60

Figure 4.19 Top Left:Mean Sensitivity values obtained with respect to SVM parameters for selected regions of dimension 128 x 128 after 100 runs, Top Right:Mean (1 - specificity) values obtained with respect to SVM parameters for selected regions of dimension 128 x 128 after 100 runs, Bottom:Performance points for different SVM parameters after 100 runs. . . .	61
Figure 4.20 Left:ROC curve obtained with SVM parameters satisfying maximum truth rate condition for selected regions of dimension 128 x 128 after 100 runs, Right:ROC curve obtained with SVM parameters satisfying maximum truth rate condition for selected regions of dimension 256 x 256 after 100 runs.	61
Figure 4.21 Left:ROC curve comparison of different selected region sizes with SVM parameters satisfying maximum Az condition, Right:ROC curve comparison of different selected region sizes with SVM parameters satisfying maximum truth rate condition. . . . .	62
Figure 4.22 SVM implementation for classification of suspicious ROIs, training set is obtained from selected ROIs. . . . .	63
Figure 5.1 Left:Sensitivity values obtained with respect to SVM parameters after classification of suspicious regions with dimension 128 x 128, Right:Sensitivity values obtained with respect to SVM parameters after classification of suspicious regions with dimension 256 x 256. . . . .	67
Figure 5.2 Left:(1 - specificity) values obtained with respect to SVM parameters after classification of suspicious regions with dimension 128 x 128, Right:(1 - specificity) values obtained with respect to SVM parameters after classification of suspicious regions with dimension 256 x 256. . . . .	67
Figure 5.3 Top Left:Sensitivity values obtained with respect to SVM parameters after classification of suspicious regions with dimension 128 x 128, Top Right:(1 - specificity) values obtained with respect to SVM parameters for suspicious regions of dimension 128 x 128, Bottom:Performance points for different SVM parameters. . . . .	68

## LIST OF ABBREVIATIONS

Az	Area Under Receiving Operating Characteristic Curve
B	Benign
BI-RADS	Breast Imaging Reporting and Data System
CAD	Computer Aided Diagnosis
CC	Cranio Caudel
CP	Completeness
CR	Correctness
FFDM	Full Field Digital Mammography
FN	False Negative
FP	False Positive
FPr	False Positive Rate
FPpI	False Positive per Image
FPF	False Positive Fraction
FROC	Free Response ROC
IARC	International Agency for Reserach on Cancer
M	Malignant
MIAS	Mammography Image Analysis Society
MLO	Medio-Lateral Oblique
O	Orientation
ROC	Receiving Operating Characteristic
ROI	Region of Interest
S	Scaling
SEER	Surveillance, Epidemiology, and End Results
SVM	Support Vector Machine
TN	True Negative
TP	True Positive
TPF	True Positive Fraction
LBP	Local Binary Patterns
LDA	Linear Discriminant Analysis
NNet	Neural Network
PCA	Principal Component Analysis
2DPCA	2 Dimensional Principal Component Analysis
NN	Nearest Neighbor Classifier



# CHAPTER 1

## INTRODUCTION

Breast cancer is the most common form of cancer that affects women all over the world and is considered a major health problem. According to the statistics of National Cancer Institute, Surveillance, Epidemiology, and End Results (SEER) program, lifetime risk of developing breast cancer among American women is 12.2 % [10]. In the European Community, breast cancer represents 19 % of cancer deaths and 24 % of all cancer cases [11] [12]. Women diagnosed between ages 40-49 years are the major victims having about 25 % of all breast cancer deaths. The World Health Organization's International Agency for Research on Cancer (IARC) has estimated more than one million cases of breast cancer to occur annually and reported that more than 400, 000 women die each year from this disease [13]. Therefore, imaging techniques such as mammography are used for early detection of breast cancer.

### 1.1 Mammography

Mammography is a particular form of radiography. It uses radiation levels between specific intervals with a purpose to acquire breast images to diagnose an eventual presence of structures that indicates a disease, especially cancer. Early detection of mammary pathologies is extremely important. The technological advances in imaging have contributed for the increase in the successful detection of breast cancer cases. In this area, mammography has an important role to detect lesions in early stages and make a favorable prognosis [6].

Mammography procedure is similar to the other X-Ray procedures. However, low doses, which presents high quality images with low noise, are used. [14]. It is desirable to use lowest radiation dose compatible with excellent image quality [6].

In terms of sensitivity and specificity mammography has better performance for fatty breasts. Dense breast tissue in young women is particularly difficult to assess. Mammography is also used in assisting needle core biopsies and for localization of non-palpable lesions [15]. In screening mammography breasts are compressed uniformly; because it is important to ensure image contrast. Thus, these tools have to be highly

sensitive to identify, as correctly as possible, those tumors that could be malignant.

Breast tissue image acquisition is done using two views in order to assess differences in density between the breast tissue: a cranio caudal (CC) and mediolateral oblique (MLO), Figure 1.1.

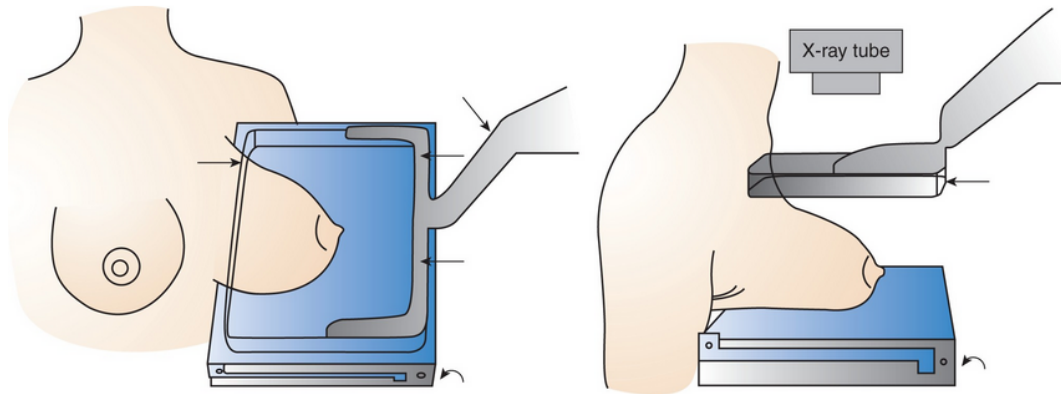


Figure 1.1: Left:MLO view. Right: CC view [3].

Figure 1.2 shows CC view of a left and right breast couple from Digital Database for Screening Mammography (DDSM) [4].

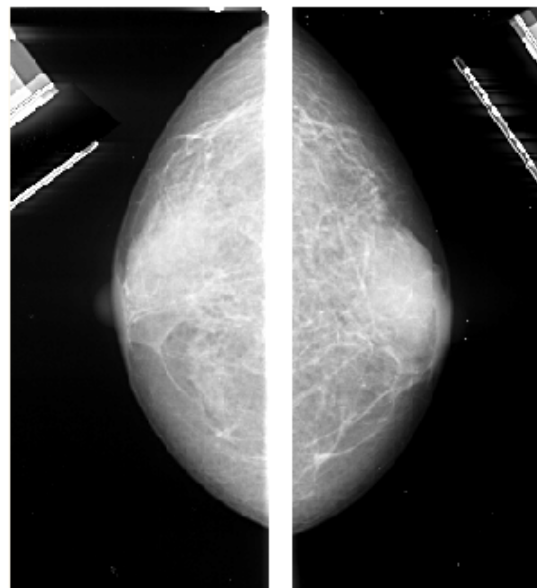


Figure 1.2: CC mammogram of a breast couple from DDSM [4].

Figure 1.3 is the MLO view of a left and right breast couple. Generally, on the MLO view, more breast tissue can be projected than on the CC view because of the slope and curve of the chest wall [16]. The image should include the free margin of the pectoral major muscle to ensure that the tail of the breast is imaged [17].



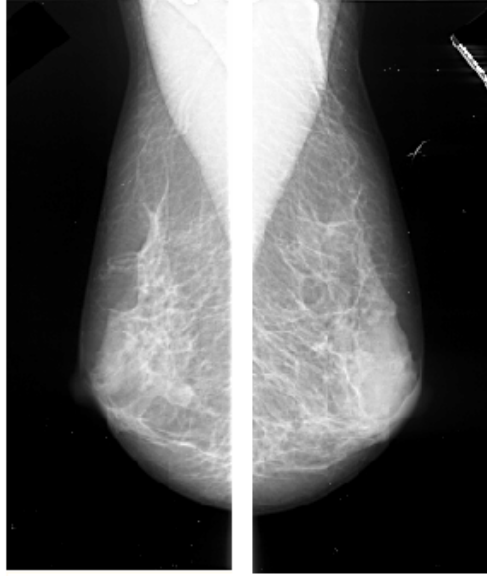


Figure 1.3: MLO mammogram of a breast couple from DDSM [4].

## 1.2 Breast Lesion Types

Breast cancer characteristic lesions are: microcalcifications, masses, architectural distortions and bilateral asymmetry [6]. In the following subsections, each of these lesions are introduced.

### 1.2.1 Microcalcifications

Microcalcifications are small deposits of calcium. They are brighter than surrounding tissues. Size of microcalcifications change from 0,33 to 0,7 mm [6]. Although they have high inherent attenuation properties, it is difficult to detect them in mammography due to their low contrast. Associated with extra cell activity in breast tissue microcalcifications may show up in clusters or in patterns [18]. A typical mammogram from mini-MIAS ([5]) database with microcalcifications is shown in Figure 1.4 [19].

A microcalcification cluster normally is more detectable than an isolated microcalcification, and contributes for the diagnosis of early stages of breast cancer. These clusters may have three or more microcalcifications present in a mammogram region with an area around  $1 \text{ cm}^2$  [18]. it is important to be able to distinguish benign and malignant microcalcifications, Table 1.1 presents the grade, degree of suspicion and mammographic appearance [1].

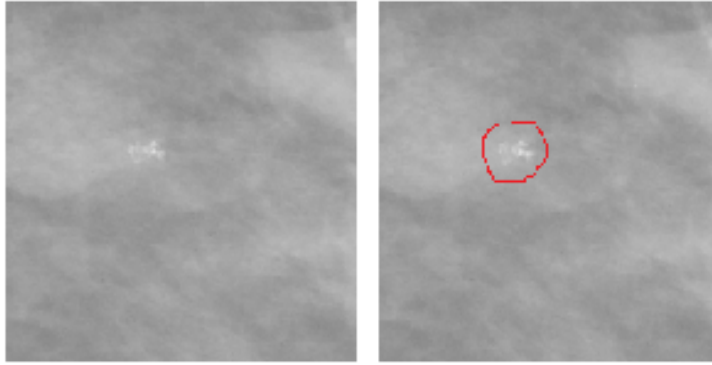


Figure 1.4: Microcalcifications in mammograms from mini-MIAS [5].

Table 1.1: Grading of imaging reports of microcalcifications according to risk of malignancy [1].

Grades	Degree of Suspicion	Mammographic Appearance
1	Normal	No abnormality seen
2	Consistent with a benign lesion	Popcorn, ring, micro cystic or diffuse bilateral
3	Indeterminate but probably benign	Localized cluster of round, punctuate
4	Suspicious of malignancy	Localized cluster of granular
5	Consistent with malignancy	Comedo calcification

### 1.2.2 Masses

Masses are areas that look abnormal in mammograms and they can be cysts, non-cancerous solid tumors or cancer. Since features of mass resembles to those of normal breast tissue, mass lesions are more difficult to detect in mammograms. Mass shape can be either round, oval, lobulated or irregular, and margins can be from circumscribed to spiculated, Figure 1.5.

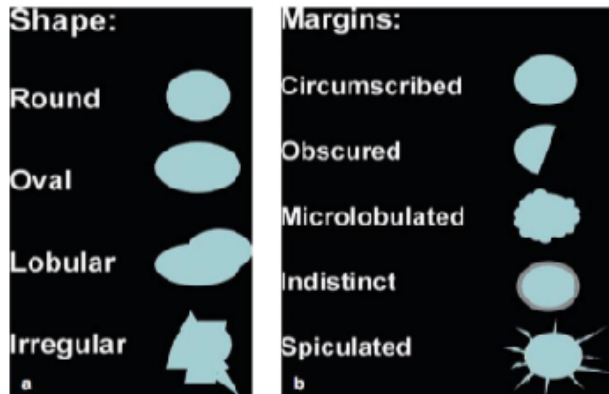


Figure 1.5: Mass Shapes in Mammograms [6].

After detection it is not easy to distinguish whether it is benign (B) or malignant (M). However, there are major differences in the shapes and textures [20]. Benign masses are typically smooth and distinct, and they are round shaped. On the other hand, malignant masses are irregular and their boundaries are usually blurry (Figure 1.6) [21]. A mass with irregular shape has a higher probability of being malignant; however a mass with an regular shape has a higher probability of being benign [22].

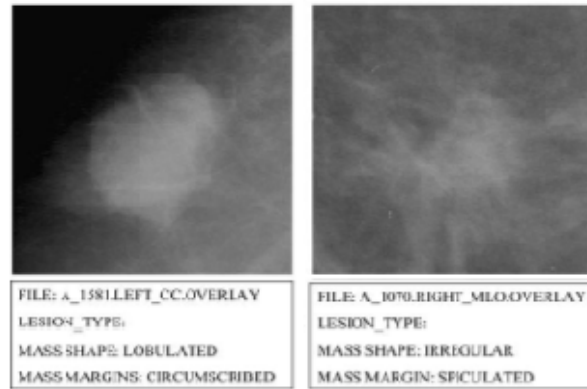


Figure 1.6: Two malignant mass ROIs cropped from full-sized mammograms with their pathological characteristics from DDSM [4].

### 1.2.3 Architectural distortions

Breast includes several linear structures that cause directionally oriented texture in mammograms and change of this textural orientation is an architectural distortion. Malignant architectural distortion includes cancer whereas benign includes scar and soft-tissue damage due to trauma. Due to its subtle appearance and variability in presentation, architectural distortion is the most commonly missed abnormality in false negative cases [23].

### 1.2.4 Bilateral asymmetry

Bilateral asymmetry of breast means a difference between corresponding regions in left and right breast and can be classified into global asymmetry and focal asymmetry [24]. The former, happens when a greater volume of fibroglandular tissue is present in one breast compared to another in the same region. The latter, corresponds to a circumscribed area of asymmetry seen on two views, and usually is an island of healthy fibroglandular tissue that is superimposed with surrounding fatty tissue [25].

### 1.3 CAD System for Mammograms

Computer aided detection (CAD) is an important application of image processing, pattern recognition, computer science and analysis techniques, aiming to assist doctors in making diagnostic decisions. If data is not easily interpretable, CAD systems may help doctors detecting subtle lesions and reduce the probability of failure. These computational systems are rising in detection of suspect cases [26]. Thus, in the past several years, CAD systems and related techniques have attracted attention of both researchers scientists and radiologists [27].

#### 1.3.1 CAD system evaluation

Result of a CAD system can be false positive (FP), true positive (TP), false negative (FN), true negative (TN) in terms of detecting the absence or presence of abnormality. Positive / negative refers the decision made by computer algorithm. False / true refers to the agreement between the decision and clinical state [28]. False positive cases result in critical operations such as biopsies.

Performance of a CAD system can be evaluated with two metrics. A CAD system's sensitivity (Equation 1.1) is the system success in noticing the abnormalities that really exist:

$$sensitivity = \frac{TP}{(TP + FN)} \quad (1.1)$$

Specificity (Equation 1.2), is a measure of how well the algorithm reports normal when there is no abnormality [29]:

$$specificity = \frac{TN}{(TN + FP)} \quad (1.2)$$

Possible trade offs between sensitivity and specificity are summarized in a receiver operating characteristic (ROC) curve, Figure 1.7. ROC curve is typically plotted with the TPF (True positive fraction = sensitivity) on Y axis and the FPF(False positive fraction = 1- specificity) on the X axis. Area under the ROC curve (Az) is expected to be 1 for an ideal detection algorithm.

In the literature, sensitivity concept is given as completeness (CP) in some papers . There is also correctness (1.3) metric that is related with the false positive rate (FPr) that  $FPr = 1 - CR$  [30]:

$$CR = \frac{TP}{(TP + FP)} \quad (1.3)$$

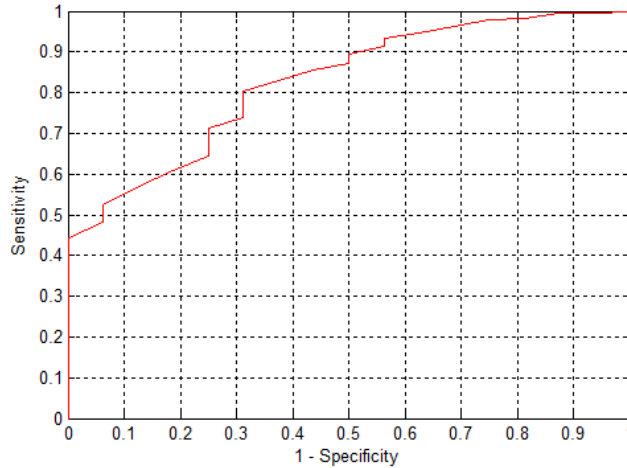


Figure 1.7: ROC curve.

### 1.3.2 Available databases

Available mammogram databases are given in [31]. DDSM, Digital database for screening mammography was created by Massachusetts General Hospital, the University of South Florida, and Sandia [4]. MIAS database was created by the Mammographic Image Analysis Society (MIAS).

MIAS database contains MLO view mammograms. Mammograms, in this database are examined by expert radiologists and for every mammogram, information of whether this mammogram has an abnormality is included. In addition, types of the lesions (microcalcification, circumscribed mass, spiculated mass, ill-defined masses, architectural distortions, bilateral asymmetry) and diagnosis results (B, M) are available. Moreover, the center, radius of the lesions, denseness (fatty, fatty-glandular, dense-glandular) of breast exist (Figure 1.8). Images are 1024 x 1024, each pixel is 200  $\mu m$  and 8 bit depth. Mammograms from mini-MIAS database is used in all phases of this thesis. 229 mammogram images from mini-MIAS database are selected in order to test our algorithms.

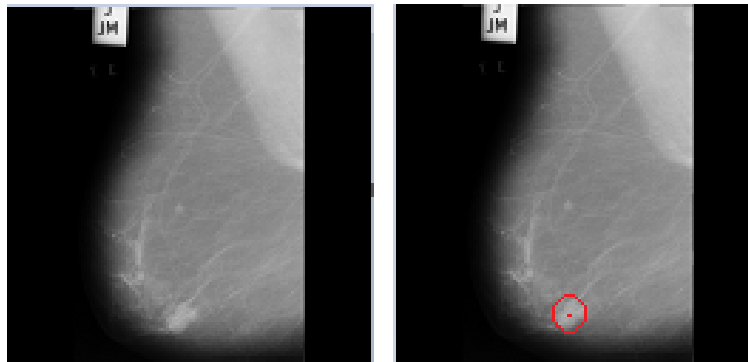


Figure 1.8: Benign circumscribed mass from mini-MIAS.

## 1.4 Literature Review

Generally, there are four basic steps of a conventional mass detection algorithm: pre-processing, segmentation, feature extraction, and classification as given in Figure 1.9. Background noise and mammogram labels are subtracted in the pre-processing step. Pectoral region segmentation and mass detection operations are implemented in segmentation step. Features of possible mass regions are determined and classification is made in feature extraction and classification steps. Following sections are the literature review of these steps. Literature review of pre-processing is not given since background is subtracted with simple image processing techniques in the literature. Segmentation includes pectoral muscle segmentation and mass detection steps for MLO view mammograms. Therefore, literature reviews of these two steps are given in stead of segmentation. Lastly, literature review of feature extraction and classification are given. Specifically, classification of textural features is concerned.

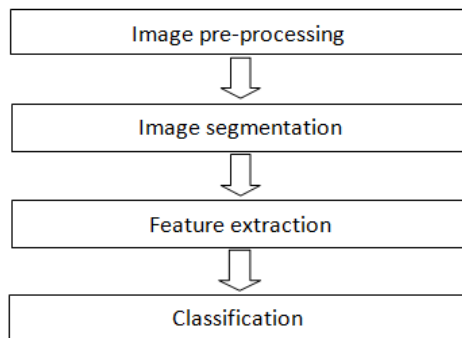


Figure 1.9: Conventional CAD mass detection algorithm [7].

### 1.4.1 Pectoral muscle segmentation

Pectoral muscle sometimes affects the performance of the algorithm due to its similar characteristic with abnormal tissues for MLO view mammograms. Mass detection should be done separately for pectoral muscle region. Therefore, pectoral muscle should be segmented before abnormal tissue search [32]. Several approaches were proposed on automatic identification of pectoral muscle in recent past. Some papers used straight line approximation [33]; Radon transform has been used to approximate the boundary with straight line in [34]. However, these methods fail when the boundary is curved. A curved edge is obtained by Hough transform in [35]. The problem is solved with Gabor filter bank in [36]. A discrete time Markov chain was applied in [37]. Two graph - based detection methods used in [38]. An iterative thresholding and gradient - based searching is applied in [39]. Average derivative and shape based features are used in [32] to obtain straight line. Then local maximum gradient search is made on this straight line to find pectoral muscle curve. The highest performance has been obtained in [32] in terms of FP pixel percentage and FN pixel percentage performance

metrics.

#### 1.4.2 Mass detection

Mass detection procedure can be carried out by implementation of one step or two distinct steps.

In the first approach regional feature extraction and classification are made. Since masses are small, these steps must be done in small regions to decide on whether there is a mass or not in a region. However, the breast region is very large when compared with mass region. Therefore, window with adequate size must be scrolled across the breast region. The drawbacks of this procedure is that it is time consuming.

In the second procedure, firstly suspicious regions are detected, then feature extraction and classification steps are applied to these suspicious regions. Detection stage aims to catch all masses in other words, it aims high sensitivity. False positives are acceptable because, false positives will be eliminated in the classification stage.

In this thesis, second procedure is implemented. Firstly, suspicious mass regions are detected with a gradient direction based adaptive filter, then classification is applied to these regions.

Low contrast characteristic of mammograms is the most important problem in the detection performance. An algorithm, resistant to low contrast is needed in order to detect all possible mass regions. Iris filter (Section 3.1) is suggested firstly in [40] and it depends only on the gradient directions of the image. Algorithm does not depend on gradient magnitudes on the image, so that it is resistant to low contrast property of mammogram images. Masses are mostly detected with Iris filter. It has the highest sensitivity performance when compared other techniques [30]. That is the main reason that we choose Iris filter in our possible mass region determination step. However, usually a large average FPPI performance is obtained when this filter is used.

First implementation of the Iris filter to detect masses in mammograms is done in [40]. After detection of suspicious regions, SNAKES algorithm is applied in order to obtain approximate boundary and 9 features' effect on the classification is discussed [40]. Mean shift algorithm and Iris filter detector is used in [41]. Mean shift segmentation is done before Iris filter implementation and the detection performance is discussed. A detailed study about the performance comparison of Iris filter in terms of implementation step size (grid size) and threshold values is made in [42]. Optimum threshold levels (highest sensitivity and lowest correctness) for different breast densities are explored in [30].

No classification step is performed in [30] [41] and [42] although large number of FPs are produced after mass detection.

### 1.4.3 Classification of textural features

Textural information such as edges, spots, lines, flat areas are important part of the visual world of animals and humans; they can successfully detect, discriminate, and segment textural characteristics using their visual systems [43]. It is not new to use textural information in order to reduce FPs of mass detection algorithms. Previous classification algorithms, using textural features of suspicious regions, are given in [2]. In addition, a new texture analysis technique using Local Binary Patterns (LBP) is offered in [2]. Classification performance with SVM is explored. The performance of the algorithm is compared with previous algorithms in Table 1.2.

Table1.2: Performance comparison of different textural approaches ([2]).

Work	Year	Features	Classifier	ROIs with mass	Normal with mass ROIs	Az
[44]	1996	Texture, morphologic	LDA, NNet	168	504	0,90
[45]	2001	Texture, shape	NNet	200	600	0,86
[46]	2005	Gray - level	NNet	681	984	0,84
[47]	2007	2DPCA	NN	256	1536	0,86
[2]	2009	LBP	SVM	512	512	0,94

Gabor filters have been used (e.g. see [48] and references therein) in order to detect breast cancer; however Gabor filters are applied on the whole image for extracting textural features, in these approaches.

Gabor filter bank is applied on different sub-regions of the ROIs extracted from mammograms and the moment based features from the magnitude Gabor responses are extracted in [49]. 256 normal and 256 mass regions are selected from DDSM database in [49]. SVM classifier is used and a performance of  $Az = 0,995$  is obtained. Such a high performance has not been obtained before so that we use the procedure applied in [49] in the feature extraction and classification steps.

## 1.5 Contributions

Steps of the proposed algorithm in this thesis is also parallel to the conventional algorithms. We propose an algorithm given in Figure 1.10. Firstly, we subtract the background region of mammogram in the pre-processing step. Secondly, segmentation of pectoral region is implemented. Thirdly, we determine suspicious mass regions with Iris filter which is an edge detection based adaptive filter. Fourthly, we extract textural features of these suspicious regions using Gabor filter bank. Lastly, we make classification with SVM classifier. Contributions are given in the following sections.



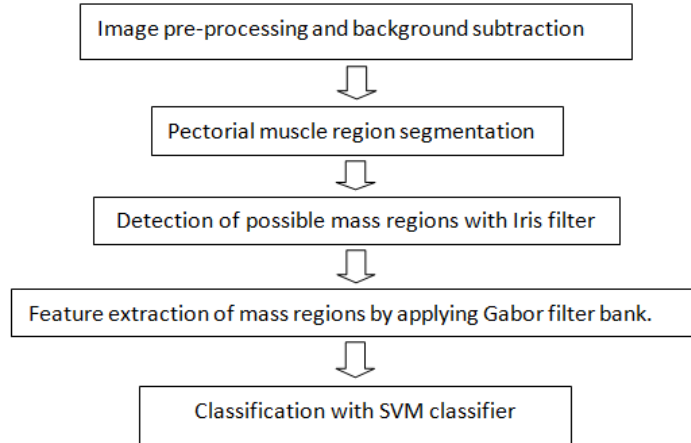


Figure 1.10: Proposed mass detection algorithm.

### 1.5.1 Pectorial muscle segmentation

- We use average derivative approach ([32]) on a pre-determined region of interest and calculate minimum derivative points for each row, firstly. Secondly, we obtain a guaranteed curve very close to real pectoral boundary for the top part of pectoral boundary. Instead of local search on a long straight line as implemented in [32], we make local search on short lines that are assumed to be the continuation of guaranteed curve on the region that boundary is not so clear. We assume that guaranteed curve continues as a line with the same slope of the last portion of curve. We call this line "line piece". Then, local minimum derivative is calculated for each row of the line piece. We fit optimum line piece for calculated minimum derivative points. Line fitting is made with least square solution [50]. We carry on this line piece fitting procedure until all boundary is found. We obtain the best performance in terms of FN pixel percentage and FP pixel percentage when compared to other algorithms that are used in the literature.

### 1.5.2 Mass detection

- We apply Iris filters with different support regions to the breast regions without pectoral muscle. We determine potential mass regions' sizes compatible with the support regions of the applied Iris filters. It is seen that the performance of classification algorithm, offered in [49] for regions with constant size, increases due to adaptive region size determination (Section 4.2.2).
- We add SVM classification step, based on features obtained by Gabor filter ([49]), to the mass detection step by Iris filter. FPPi ratio, obtained after mass detection step, is reduced 50 percent with a cost of missing 9 percent of the true mass regions.

## 1.6 Outline of the Thesis

The rest of the thesis is organized as follows.

Chapter 2 gives an overview of background subtraction and pectoral muscle segmentation. Algorithm developed for the pectoral muscle segmentation is explained.

Chapter 3 includes overview on Iris filter and its application on mammograms to detect possible mass regions.

Chapter 4 explains how textural features of suspicious regions are extracted with Gabor filter bank application. Moreover, SVM classification performance with respect to changing SVM parameters is discussed. Furthermore, classification of suspicious regions is given.

Chapter 5 discusses the results of the implemented classification in terms of sensitivity and specificity metrics.

Chapter 6 concludes the thesis.

## CHAPTER 2

### BACKGROUND SUBTRACTION AND PECTORAL MUSCLE REGION SEGMENTATION

Background subtraction and pectoral region segmentation algorithms are presented in this chapter.

#### 2.1 Image Processing Techniques Used in Background Subtraction

Binary erosion and dilation are the operators that use subtraction and addition of set elements. The dilation operation usually uses a structuring element for expanding the shapes contained in the input image whereas, erosion operation uses a structuring element for narrowing the shapes. They are morphologically dual of each other [51].

Erosion operation is defined in Equation 2.1 and illustrated on the left side of Figure 2.1. "A" is a binary image and pixels in blue region is 1. "B" is called structuring element with pixel values 1. Any point in region "B" is "b". While moving structuring element on binary image, erosion of A by B can be understood as the intersection of points reached by the center of B. For example, the erosion of a square of side 10 by a disc of radius 3 is a square of side 4. Erosion result of blue region is the green region.

$$A \ominus B = \bigcap_{b \in B} A_{-b} \quad (2.1)$$

Dilation operation is dual of erosion and defined in Equation 2.2, illustrated on the right side of Figure 2.1. While moving structuring element B on binary image "A", union set of "A" and "B" regions constitute dilation result. Dilation of a square of side 10, by a disk of radius 3 is a square of side 16 with rounded corners. Radius of the rounded corners is 3. Dilation result of blue region is the green region.

$$A \oplus B = \bigcup_{b \in B} A_b \quad (2.2)$$

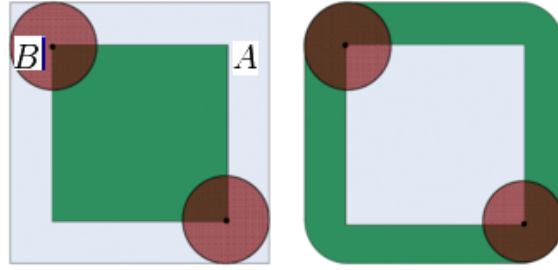


Figure 2.1: Left:Erosion operation and illustration, Right:Dilation operation and illustration [8].

Sometimes it is necessary to find of values connected pixels (1) are need to be found (Figure 2.2 left) on the binary image in image processing problems. They are found with 8 pixel neighbour search. If any pixel among 8 neighbour pixels is equal to pixel of interest this neighbor pixel and the pixel of interest are assumed connected pixels. Searching whole image, all connected components are found all labelled (Figure 2.2 right). In other words, a number is given for each component.



Figure 2.2: Left:Determination of connected components, Right:Labeling of connected components [9].

## 2.2 Background Subtraction

Mammogram images in mini-MIAS database are not obtained from a digital mammography but from a conventional one. They are digitized form of analog mammogram images and undesired parts such as label of the mammogram, noise, etc. exist on mammogram images. Therefore, background region must be discriminated from the breast region.

Global threshold is applied. Pixels below a threshold level are assigned 0 and over 1. This threshold level is chosen experimentally such that any breast region is not eliminated. After threshold operation, a binary image is obtained. Next operations are done on this binary image Figure 2.3 (Step 1). It is aimed to make all breast region pixels 1 and background pixels 0. Erosion operation is applied in order to

delete unnecessary parts, such as mammogram labels, over mammograms on binary image (Step 2). Mammogram label is deleted. However, while eroding binary image not only mammogram labels but also some part of breast region is also eroded. Thus, dilation operation is applied on image and the lost breast region is recovered (Step 3).

Operations are done on the binary image as mentioned. After the dilation operation, connected component labeling is required. Therefore, it is aimed to find all connected pixels for each region and each region is labelled (Step 4). Label of breast region is found and other regions' pixel values are set to 0. Consequently, only breast region with pixel values 1 and background region with pixel values 0, remains. This binary image is multiplied pixel by pixel with the original mammogram image and background is subtracted so that only breast region remains (Step 5).

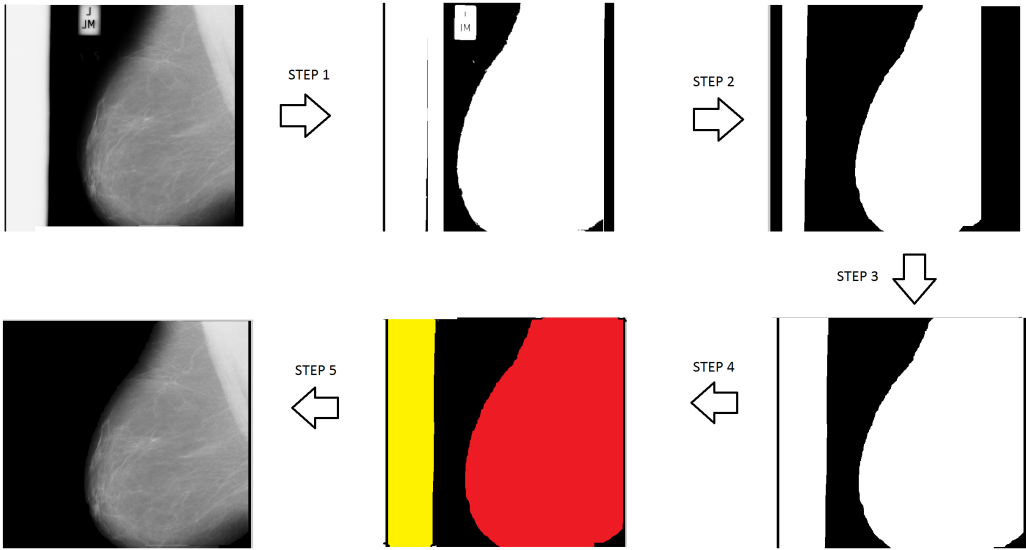


Figure 2.3: Background subtraction, Step1:Binary image obtainment after threshold, Step2:Erosion operation to eliminate mammogram label, Step3:Dilation operation to get back regions lost on mammograms after erosion step, Step4:Connected component labeling to discriminate breast region from other unnecessary parts.

### 2.3 Pectoral Muscle Segmentation

#### 2.3.1 Pectoral muscle segmentation algorithm

Figure 2.4 shows the algorithm applied in this thesis. Pectoral region is assumed to be always at the top left of the image. Therefore, all of the left turned breasts are rotated 180 degrees to satisfy this condition. 1024 x 1024 mammograms are resized to size of 512 x 512 in order to run the algorithms faster.

Firstly, ROI (Region of Interest) is determined. For this purpose edge point at the top

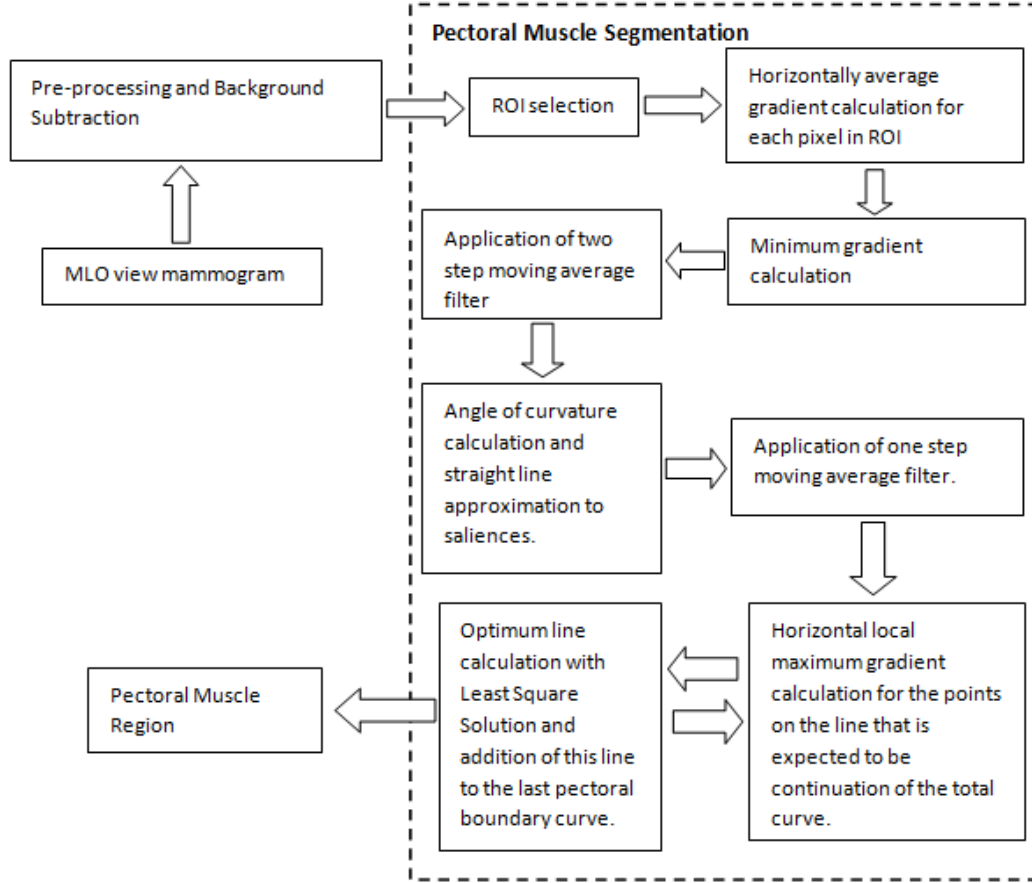


Figure 2.4: Pectoral muscle segmentation algorithm.

raw of the mammogram (A) and the bottom left point of the mammogram (B) are joined to designate hypotenuse edge of ROI as given in Figure 2.5.

Along the x axis in Figure 2.5 average derivative (Equation 2.3 [32]) is calculated for all  $x$  on  $y = 1$ . A typical average derivative change with respect to  $x$  is given in Figure 2.6. Use of average derivative reduces the effect of high intensity variation of noise spike and curvilinear structures [32]. The minimum gradient (actually maximum in magnitude but negative) point is assumed to be edge point for raw of interest. Therefore, point A is the pixel satisfying the minimum derivative condition for the first raw of the mammogram image.

$$Average\ derivative(x, y) = \frac{1}{N} \sum_{i=1}^N \frac{I(x+i, y) - I(x-i, y)}{2i} \quad (2.3)$$

where

$(x, y)$  is coordinate of the pixel where derivative is calculated.

$N$  is the number of pixel pairs used for average derivative computation.

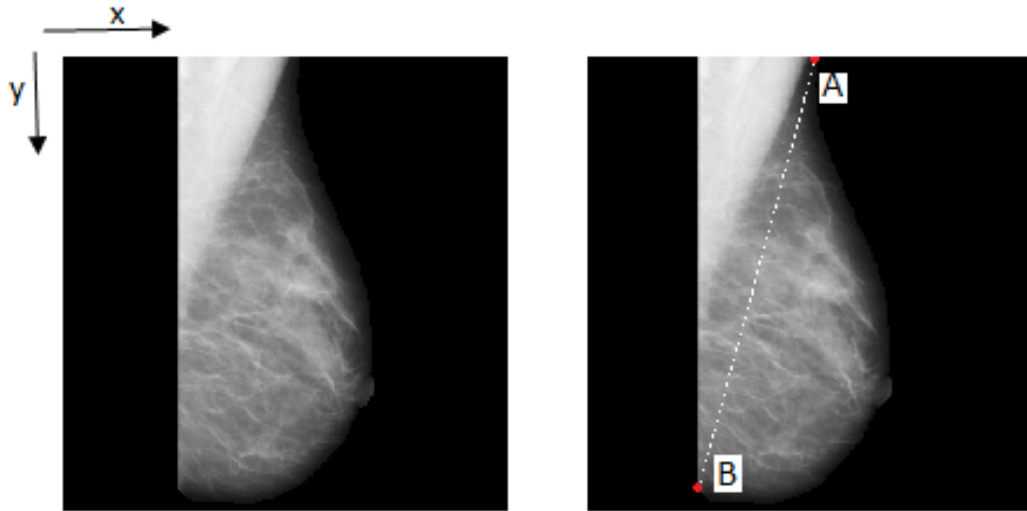


Figure 2.5: ROI determination.

$I(x, y)$  is the intensity at  $(x, y)$  position.

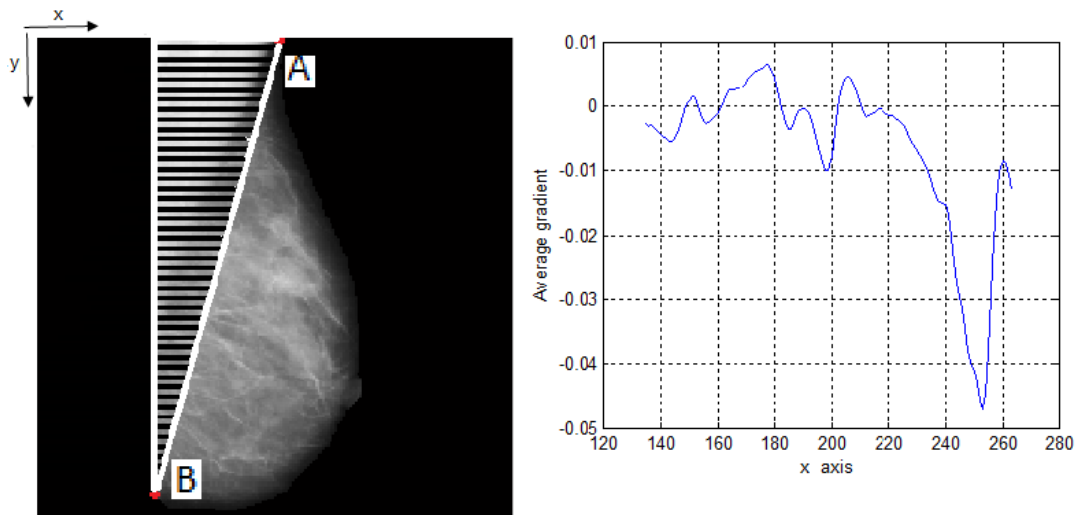


Figure 2.6: Average derivative calculation for each row and a typical average derivative plot

Pectoral muscle boundary points are aimed to be found for each  $y$  of the ROI. For this purpose, average derivative plot is obtained for each row and minimum average derivative value is determined for each plot (Figure 2.6).

For each  $y$  value in the ROI, an  $x$  value through which the pectoral muscle boundary is expected to pass, will be calculated. However, due to intensity variations apart from pectoral muscle boundary region in the ROI, the points determined will not always be on the pectoral muscle boundary (Figure 2.7). The points'  $x$  coordinate variation with respect to  $y$  is noisy.

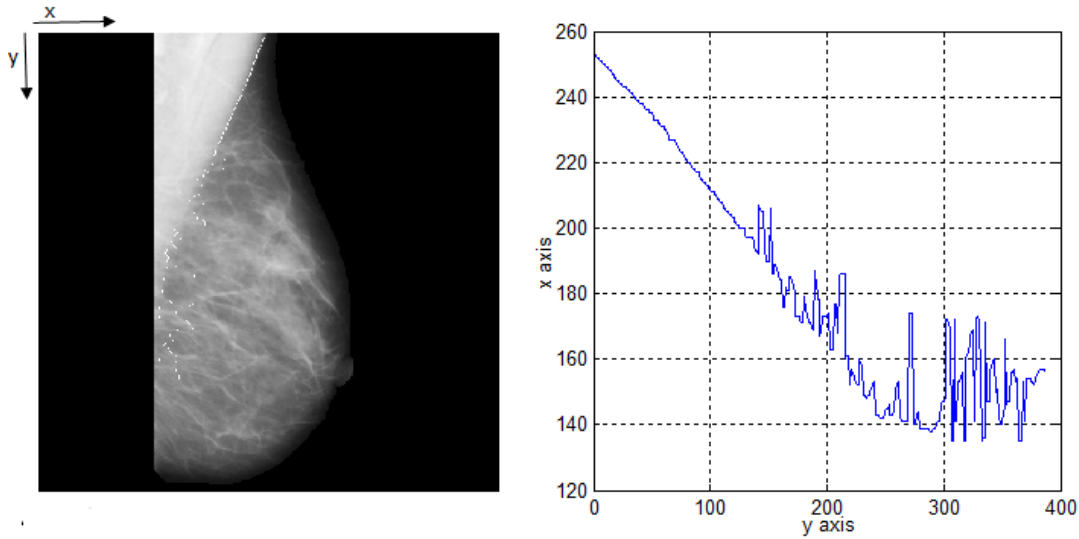


Figure 2.7: Minimum derivative points for each row in ROI.

Smoothing is required to clarify the pectoral muscle boundary; therefore two step moving average filter is applied (Figure 2.8, Figure 2.9).

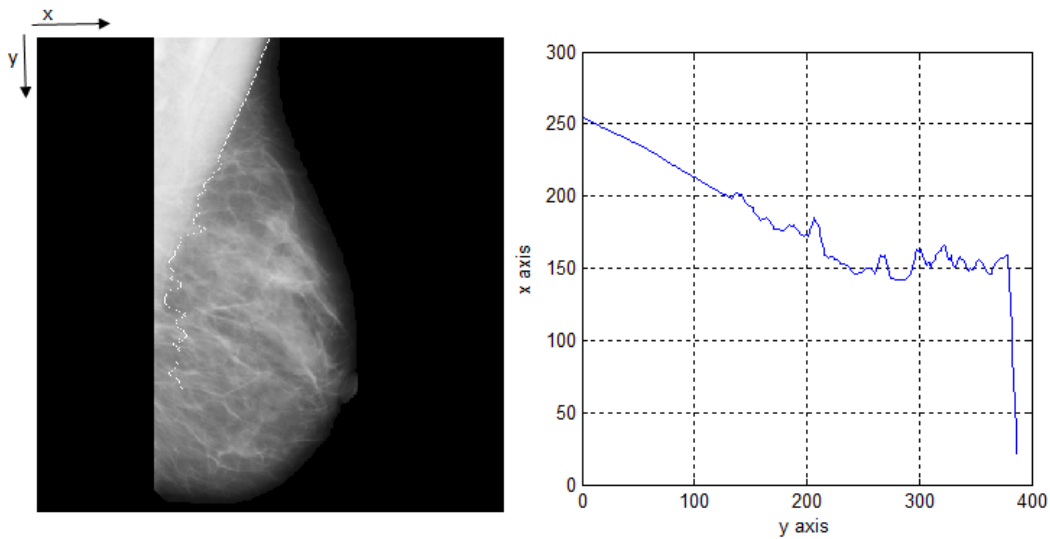


Figure 2.8: Minimum derivative points for each row in ROI after first step moving average operation.

It is observed that after some y value, determined curve behaves very different from pectoral muscle boundary. However, it is not expected from pectoral boundary to have abrupt saliencies. In fact, it is generally line or curved line. As a result of this fact, it is meaningful to obtain angle of curvature plot of each point (Figure 2.10). Angle of curvatures are calculated by Equation 2.4.



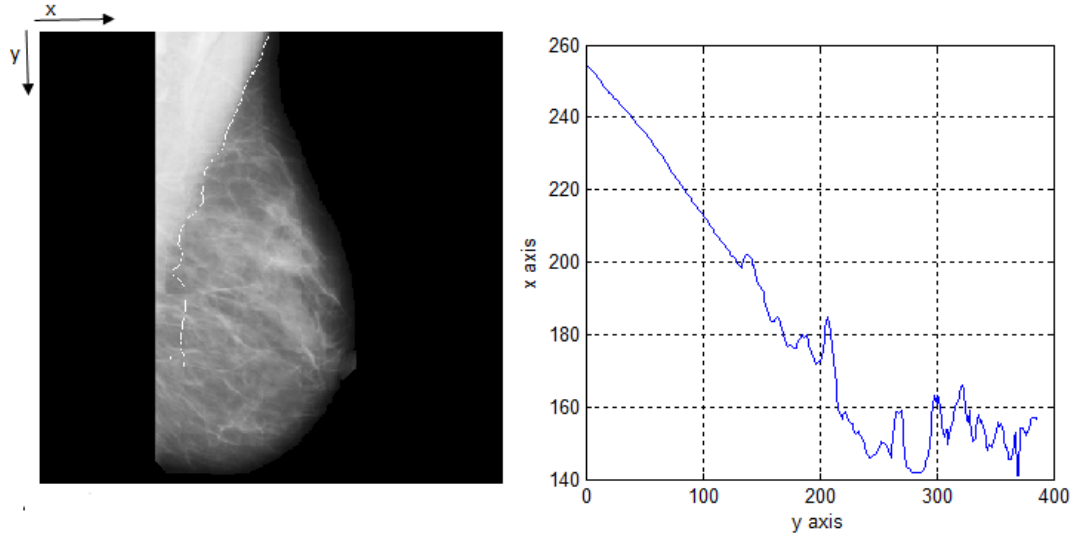


Figure 2.9: Minimum derivative points for each row in ROI after second step moving average operation.

$$m1 = \frac{x(y + \delta y) - x(y)}{\delta y} \quad (2.4a)$$

$$m2 = \frac{x(y) - x(y - \delta y)}{\delta y} \quad (2.4b)$$

$$\theta(y) = \text{atan} \left[ \frac{(m1 - m2)}{(1 + m1m2)} \right] \quad (2.4c)$$

where

$(x, y)$  is the coordinate of a current pixel  $(i, j)$  on the curve.

$\theta(y)$  is the angle of curvature at  $y$ .

$\delta y$  is the difference in  $y$  axis, between the pixel of which, angle of curvature will be determined and other two points (3 points needed to find the angle between two lines). It is chosen as 5 in this thesis.

If angle of curvature is above a positive threshold level or under a threshold level the incoming points after this point will be thrown away. In this thesis, this threshold levels are determined experimentally. Threshold level is 16 for positive curvature angles and -16 for negative curvature angles. In addition, if there are still saliencies a line will be fitted instead of this saliency. An acceptable curve close to the pectoral boundary has been obtained after this operation (Figure 2.11).

An acceptable curve, which is very close to real boundary, is obtained for one part of the boundary. In order to obtain the rest of the boundary curve where more uncertainty

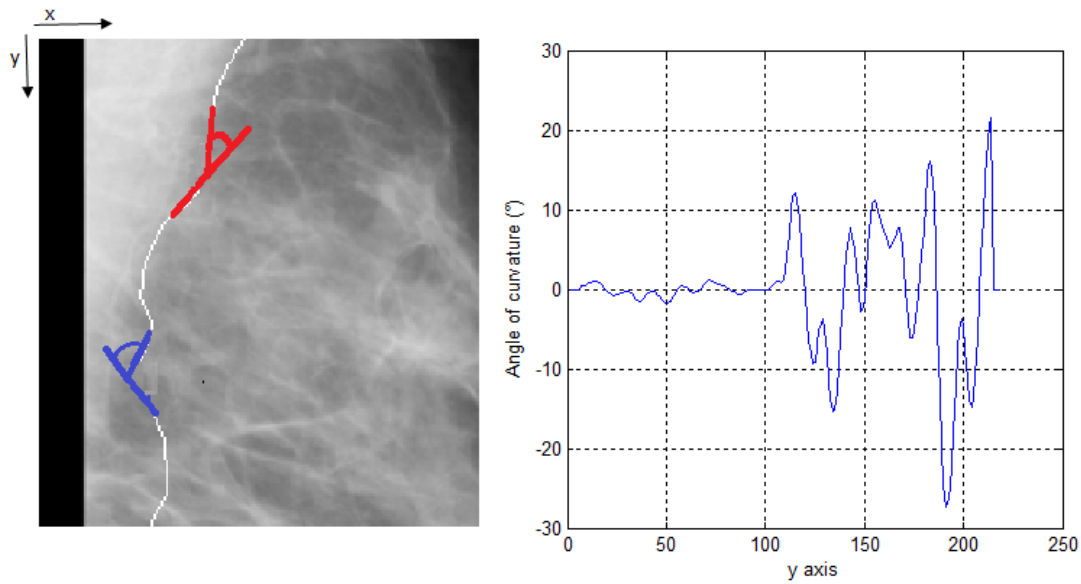


Figure 2.10: Angle of curvature values for the points, on pectoral muscle boundary.

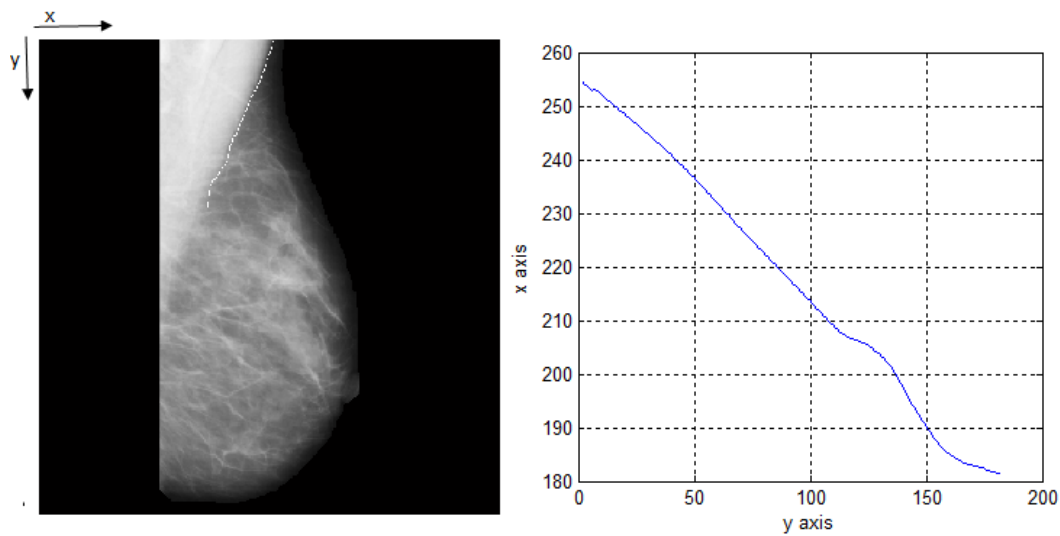


Figure 2.11: Pectoral muscle boundary points after the removal of saliencies.

exist in terms of boundary clearness, an algorithm based on local minimum derivative search and least square solution has been implemented.

Firstly, it is assumed that obtained curve continues with an imaginary line piece having the same slope with the last portion of the curve as shown in Figure 2.12. The size of the line piece is assumed to be the vertical size of the line piece.

At this point it should be mentioned that line piece size is not distinct and different line piece sizes result in different boundary curves. Radiologist or the expert will decide on which line piece size is producing the best looking boundary among different boundaries produced with different line piece sizes. Thus, offered algorithm is a semi-automatic

one. Line piece size is chosen 70 for this mammogram.



Figure 2.12: Imaginary line for horizontal derivative search muscle.

Derivative of each point on the horizontal line is drawn for every pixel on the line piece (Figure 2.13). The size of horizontal line is chosen as 21. 10 pixels, on the left side and 10 pixels, on the right side of the pixel of interest on line piece. There are as many horizontal line as the line piece size. Determined pixels, having minimum derivatives, are seen as black around the imaginary red line in Figure 2.13.

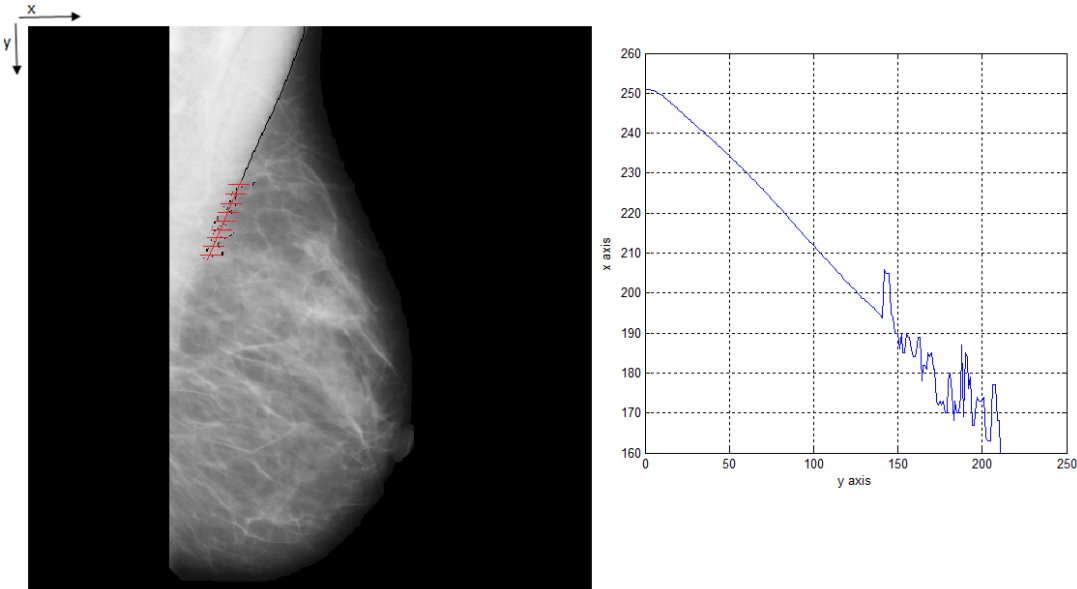


Figure 2.13: Minimum derivative calculation on the horizontal lines passing through imaginary line piece added.

Best line passing through these calculated pixels, is aimed to be estimated. An over determined system, as given in Equation 2.5, must be solved. Least square solution is offered in the literature to find unknowns for an overdetermined system [50]:

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \cdot \\ y_i \\ \cdot \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \\ \cdot & \cdot \\ x_i & 1 \\ \cdot & \cdot \\ x_n & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} \quad (2.5)$$

where

$(x_i, y_i)$  is the coordinate of a pixel after minimum derivative calculation.

$y = ax + b$  is the best line fitted to the pixels after minimum derivative calculation.

Least square solution is applied to this problem to estimate the optimum  $(a, b)$  pair as given in Equation 2.6:

$$\begin{bmatrix} a \\ b \end{bmatrix} = \left[ \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \\ \cdot & \cdot \\ x_i & 1 \\ \cdot & \cdot \\ x_n & 1 \end{bmatrix}^T \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \\ \cdot & \cdot \\ x_i & 1 \\ \cdot & \cdot \\ x_n & 1 \end{bmatrix} \right]^{-1} \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \\ \cdot & \cdot \\ x_i & 1 \\ \cdot & \cdot \\ x_n & 1 \end{bmatrix}^T \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \cdot \\ y_i \\ \cdot \\ y_n \end{bmatrix} \quad (2.6)$$

Estimated line will be added to the end of the last boundary curve obtained after removal of saliencies. The resultant boundary curve is given in Figure 2.14.

This procedure is repeated until the line piece approaches to the left boundary of breast. For this mammogram image at 3 steps the curve reaches left boundary of breast. The results after minimum derivative pixels and least square calculations for second and third steps are given in Figures 2.15, 2.16, 2.17, 2.18 respectively.

After reaching the left side of the breast region a moving average smoothing filter is applied on boundary curve and pectoral boundaries are obtained.

Different line piece sizes produce different pectoral boundaries as mentioned before. Pectoral boundaries, which are obtained for different line piece sizes, are given in Figure 2.19.

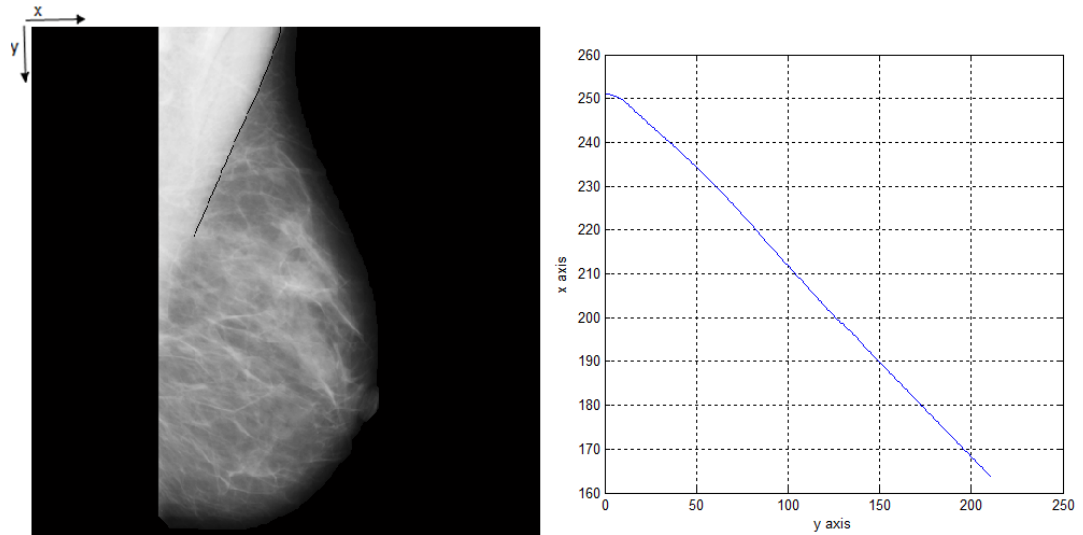


Figure 2.14: Least square solution applied to pixels.

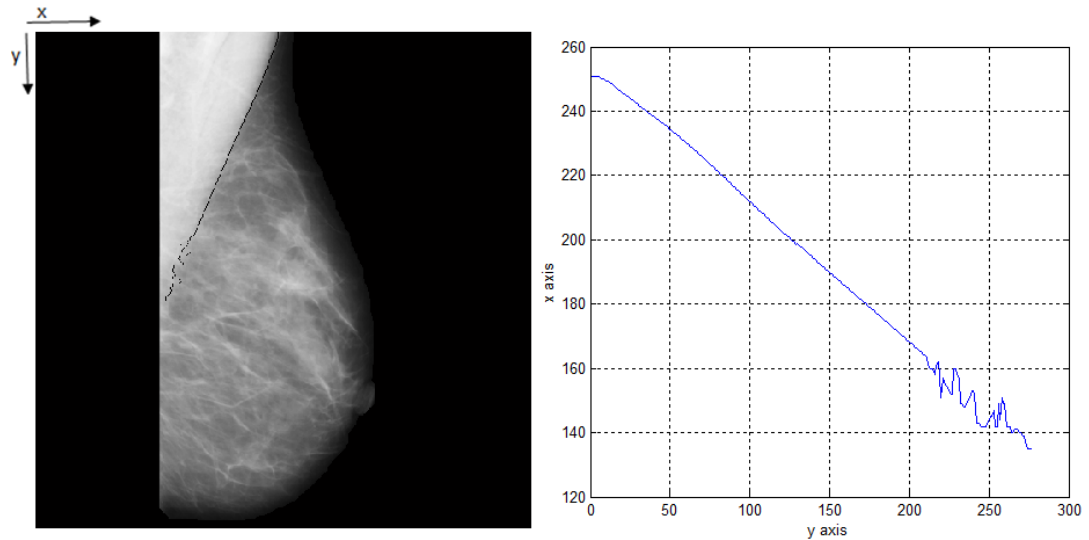


Figure 2.15: Minimum derivative calculation on the horizontal lines passing through imaginary line piece added (Second step).

Pectoral muscle boundary curves, calculated for some mammogram images, are given in Figure 2.20.

## 2.4 Evaluation of pectoral muscle segmentation algorithm

If the line piece is chosen to be constant and 30, we obtain satisfactory results for nearly 80 percent of pectoral boundaries. However, a semi-automatic algorithm, mentioned above, is offered for unsatisfactory boundaries and this may be the main drawback of proposed pectoral discrimination algorithm. In other words, someone has to compare

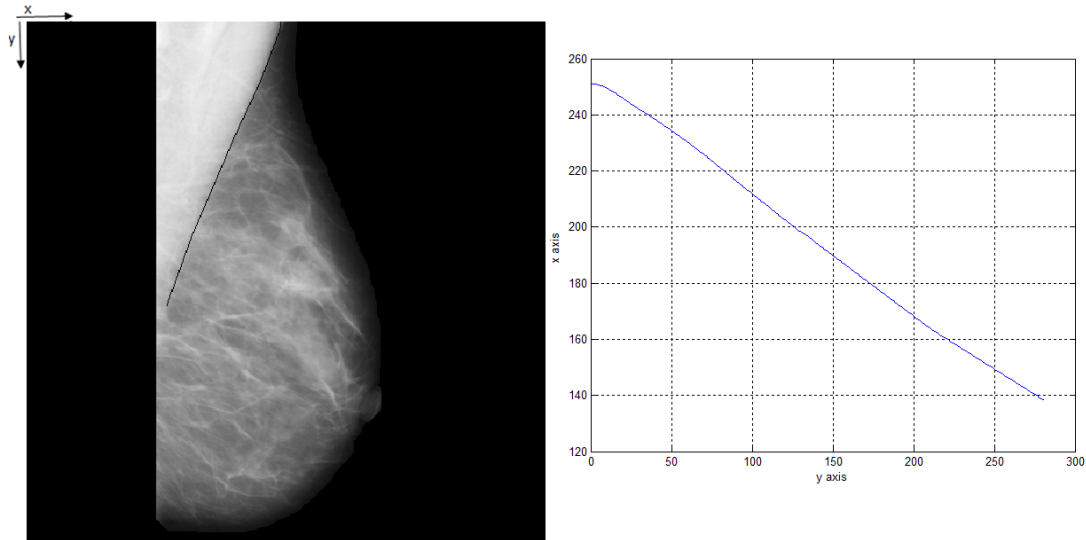


Figure 2.16: Least square solution applied to pixels (Second step).

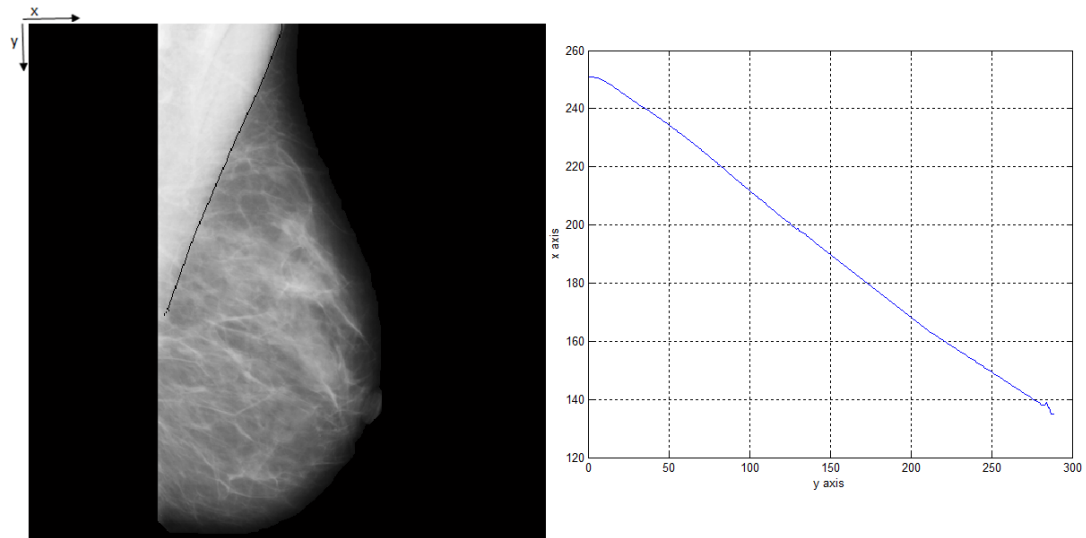


Figure 2.17: Minimum derivative calculation on the horizontal lines passing through imaginary line piece added (Third step).

different discrimination results, which are obtained with different line piece sizes, and choose the best line. On the other hand, pectoral boundaries, which are drawn with our algorithm, are better when compared to the other implemented algorithms in the literature.

In order to evaluate performance of the applied algorithm mammogram images are printed on A4 papers, firstly. Secondly, real boundary is drawn with a colored pen and photos of mammograms are taken with a professional camera (Nikon D7000). Thirdly, taken photos are transferred to computer. Green line on Figure 2.21 is true boundary whereas black line is drawn with offered algorithm.

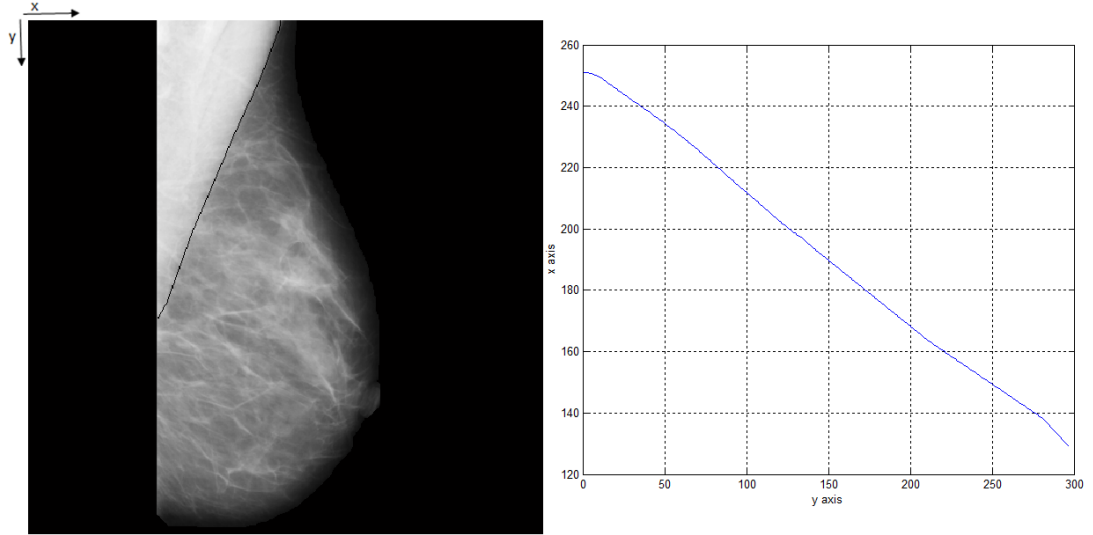


Figure 2.18: Least square solution applied to pixels (Third step).

Performance of algorithm is determined for each mammogram in terms of two metrics: FP pixel percentage and FN pixel percentage. FP and FN pixel percentages are mathematically defined in Equations 2.7 and 2.8, respectively:

$$FP \text{ pixel percentage} = \frac{|A \cap B| - |B|}{|B|} 100\% \quad (2.7)$$

$$FN \text{ pixel percentage} = \frac{|A \cap B| - |A|}{|B|} 100\% \quad (2.8)$$

where

$A$  is pectoral muscle region bounded with boundary obtained by offered algorithm

$B$  is true pectoral muscle region

We applied the proposed algorithm on 201 mammograms from mini-MIAS database. We calculated mean and standard deviation of FP pixel percentage, FN pixel percentage and compared with previous algorithms. 80 mammograms were used from mini-MIAS database in the previous works. Success of previously applied algorithms and offered algorithm is given in Table 2.1. It can be observed that offered method is more successful than applied methods in [32] and [36] in terms of both FP and FN pixel percentages. However, as mentioned earlier that it is a semi-automatic algorithm. It may be a disadvantage.

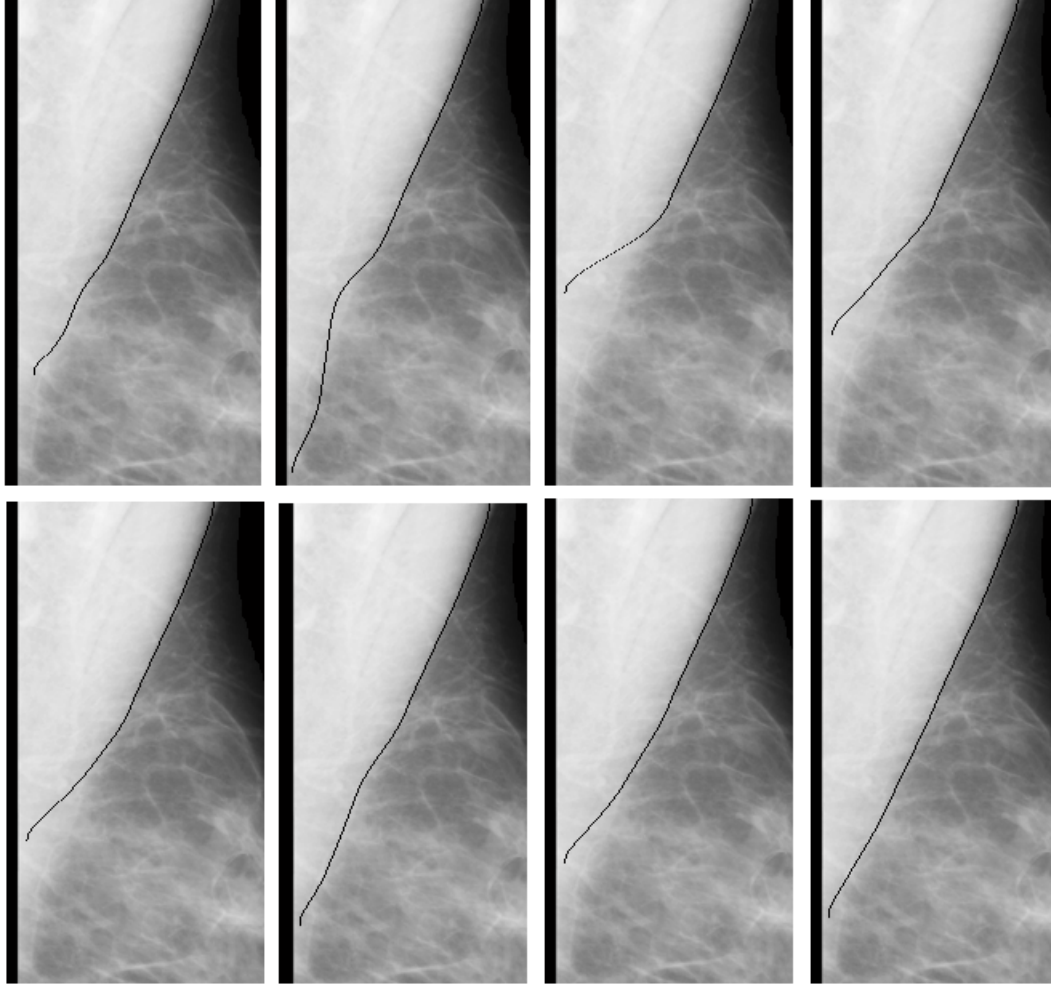


Figure 2.19: Calculated pectoral muscle boundaries with different line piece sizes for the same mammogram; First row:8 , 12, 16, 24, Second row:32, 36, 52, 60

## 2.5 Contrast Enhancement and Mass detection in the Pectoral Regions

Textural structures of pectoral region and breast region without pectoral muscle is very different. Therefore, image enhancement algorithms should be applied to these regions separately. For pectoral region, contrast enhancement is an important issue. In this thesis, conventional histogram equalization ([52]) is applied to the pectoral region, however satisfactory results are not obtained. It is thought that regional contrast enhancement would be a better solution for the contrast enhancement of pectoral regions. In addition, different mass detection and classification algorithms should be applied on pectoral regions. Since we do not have pectoral regions with mass, we could not applied any mass detection and classification algorithm for the pectoral regions. In this thesis, we handle only breast region without pectoral muscle for our mass detection and classification algorithms. Next chapter explains detection of possible mass regions



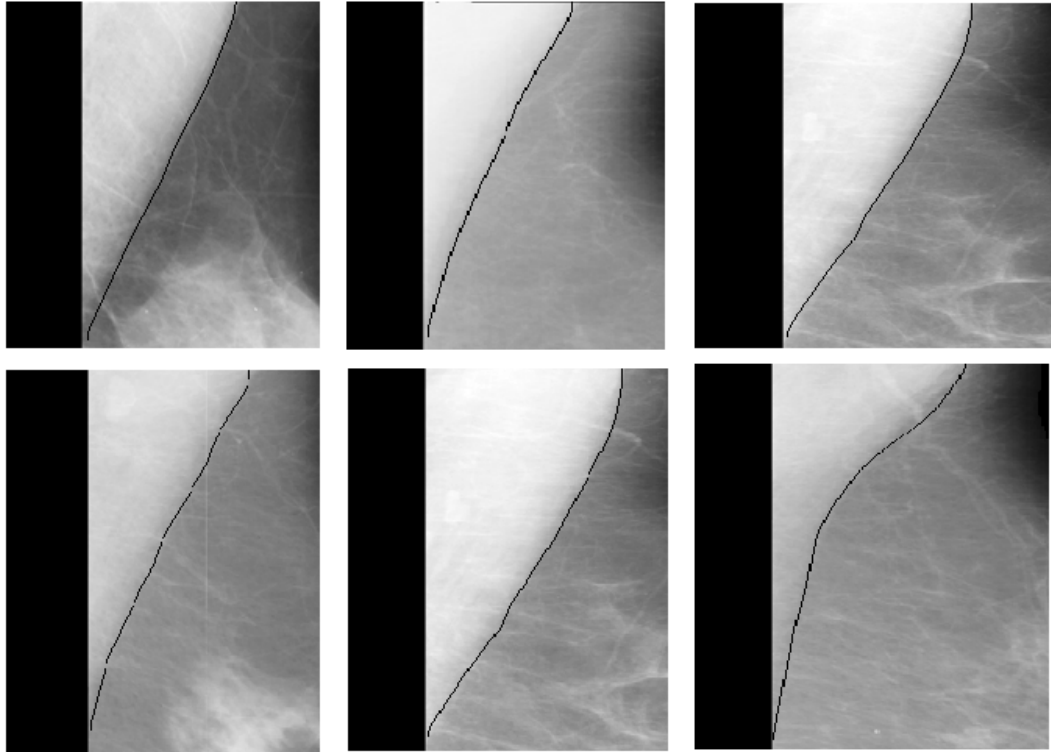


Figure 2.20: Calculated pectoral muscle boundaries for some mammograms, from mini-MIAS database.

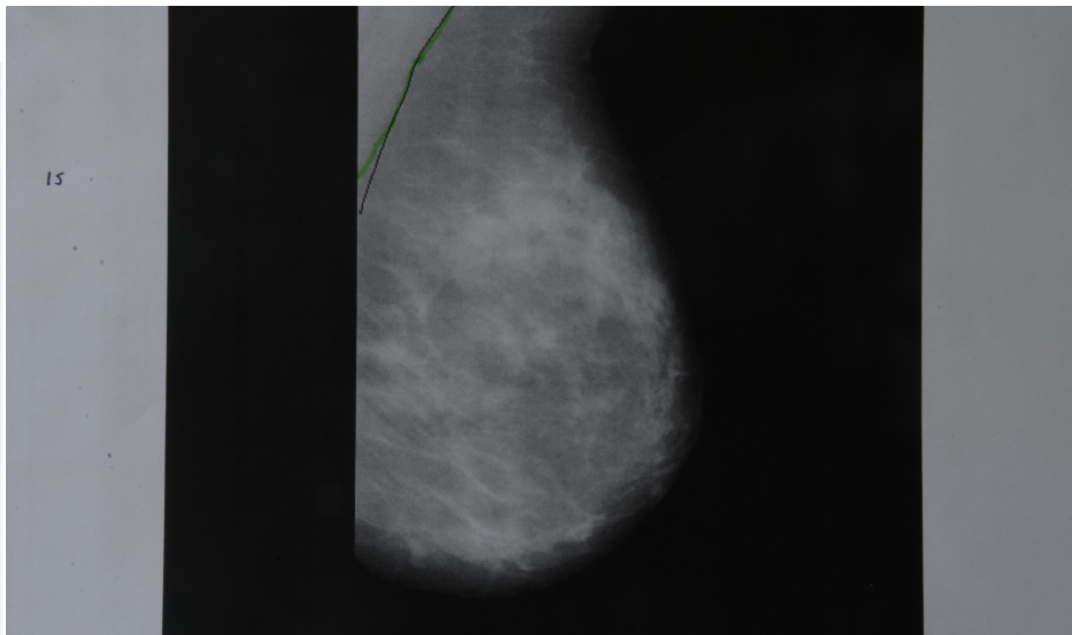


Figure 2.21: Mammogram image, taken with a camera.

in the breast region without pectoral region.

Table2.1: Mean and standard deviation of FP, FN pixel percentages for mammograms from mini-MIAS database.

Performance metric	Proposed method	[32]	[36]
FP pixel percentage	$1,30 \pm 2,56$	$1,84 \pm 2,83$	$4,64 \pm 5,03$
FN pixel percentage	$2,57 \pm 1,84$	$6,56 \pm 6,48$	$4,33 \pm 5,63$

## CHAPTER 3

### DETECTION OF POSSIBLE MASS REGIONS WITH IRIS FILTER

Mass detection with Iris filter in the breast region without pectoral muscle, is given in this chapter. Firstly, Iris filter is introduced. Secondly, Iris filter usage in the mass detection algorithm is explained. Thirdly, success of the algorithm is discussed in terms of sensitivity and FPPi performances.

#### 3.1 Iris Filter

Iris filter evaluates the degree of convergence of the gradient vectors within its region of support (support region = region between  $R_{min}$  ( $R_{in}$ ) and  $R_{max}$  ( $R_{out}$ ) circles) toward a pixel of interest [53] (Figure 3.1). Degree of convergence is related to the distribution of the directions of the gradient vectors and not to their magnitudes. Convergence index of a gradient vector at a given pixel is defined as the cosine of its orientation with respect to the line connecting the pixel and the pixel of interest. Equation 3.1 is the calculation of the convergence index of a gradient vector:

$$f(L_i) = \begin{cases} 0, & \text{if } |g| = 0. \\ \cos(\beta), & \text{otherwise.} \end{cases} \quad (3.1)$$

where

$L_i$  is a given pixel, Figure 3.1

$\beta$  is the angle between gradient vector at  $L_i$  and the vector from  $L_i$  to pixel of interest ( $P$ ), Figure 3.1.

$g$  is the gradient at pixel  $L_i$ .

$f(L_i)$  is the convergence index of gradient at point  $L_i$

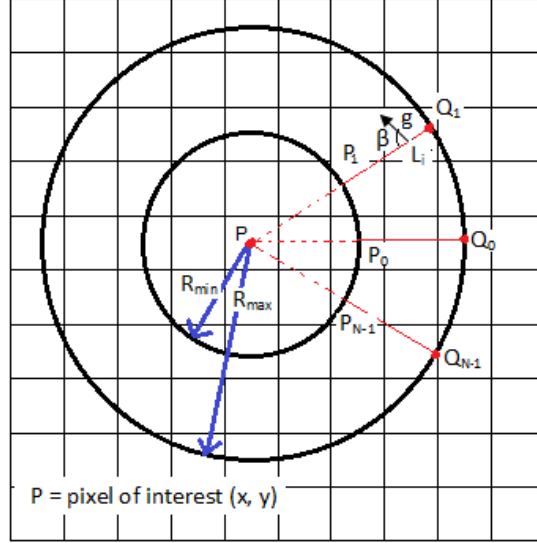


Figure 3.1: Iris filter definition.

The convergence degree of gradient vectors on the line of  $P_j L_i$ ,  $C_i$  can be defined as the average of convergence indexes over the length  $P_j L_i$  as given in Equation 3.2:

$$C_i = \frac{\int_{P_j}^{L_i} f(L_i) dl}{|P_j L_i|} \quad (3.2)$$

where

$P_j$  is the intersection pixel of  $j$  th line with the circle of radius  $R_{min}$ , Figure 3.1.

$l$  is length on line,  $P_j L_i$ .

$C_i$  is the convergence degree of gradient vectors on the line,  $P_j L_i$ .

$j$  is line number, defined in the interval  $[0, N - 1]$ ,  $N$  is the total number of lines.

The maximum convergence degree  $C_{i0}$  on the  $j$  th line is given in Equation 3.3:

$$C_{i0} = \max(C_i), L_i \in [P_j, Q_j] \quad (3.3)$$

where

$C_{i0}$  is the maximum convergence degree.

$Q_j$  is the intersection pixel of  $j$  th line with the circle of radius  $R_{max}$ , Figure 3.1.

Output of the Iris filter is the average of the convergence degree within its region of support and lies within the range  $[-1, 1]$ :

$$C(x, y) = \frac{1}{N} \sum_{j=0}^{N-1} C_{i0} \quad (3.4)$$

where

$C(x, y)$  is output of the Iris filter at the pixel of interest  $(x, y)$ .

### 3.2 Suspicious ROI Detection Algorithm

Generally, there is an increase in terms of intensity from border to center in a mass region. Sometimes this increase ends up in a local maximum point in the region as shown in Figure 3.1. All gradient vectors on the edge pixels are directed towards the mass region. Moreover, in mass region gradient vectors are directed towards a local maximum point. Iris filter application to mass detection arises from this idea indeed. If Iris filter's pixel of interest is on such a local maximum, the output of filter will be very close to 1. However, region of support of the filter is very critical. Gradient vectors in the region of support must be directed towards pixel of interest to obtain high output values. Otherwise, if the gradient vectors, directed toward pixel of interest, are out of this support region it is meaningless in terms of contribution to filter output.

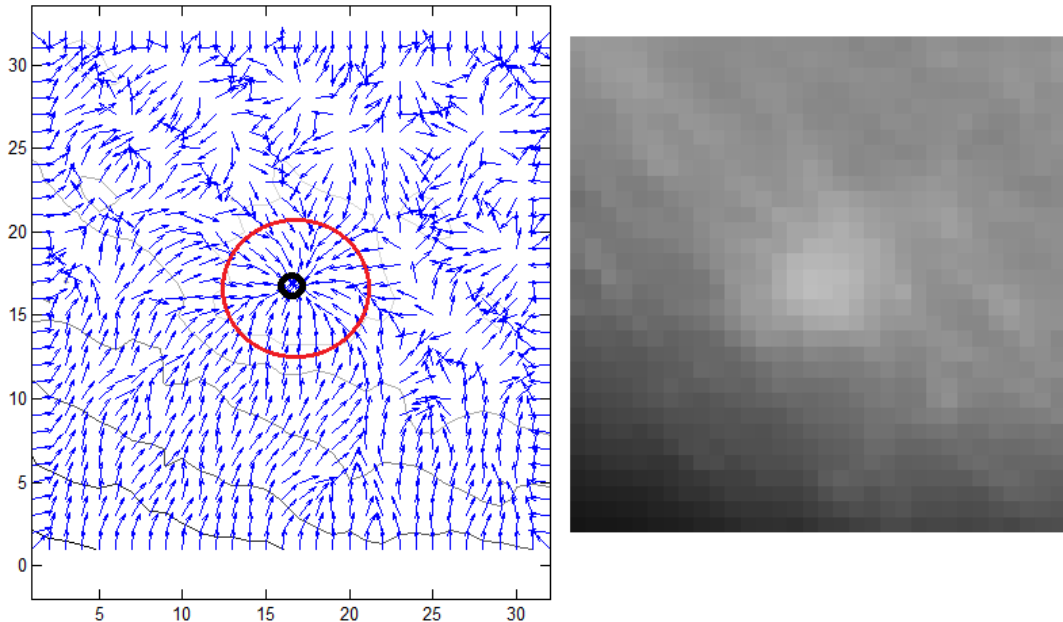


Figure 3.2: Left:Gradient map of the mass region with one local maximum, Right:Mass region.

Sometimes gradient vectors are directed towards more than one local maximum points as shown in Figure 3.3. Gradient vectors are directed towards the mass region on the edge pixels for this scenario, too. An Iris filter with a region of support including the edge pixels of the mass region will produce high output value for a pixel of interest in the mass region. However, if one of local maximums is the pixel of interest and we select the region of support such that it includes the gap between local maximums (Gradient vectors are not directed towards local maximum points in this region.) the Iris filter output will not be highest value.

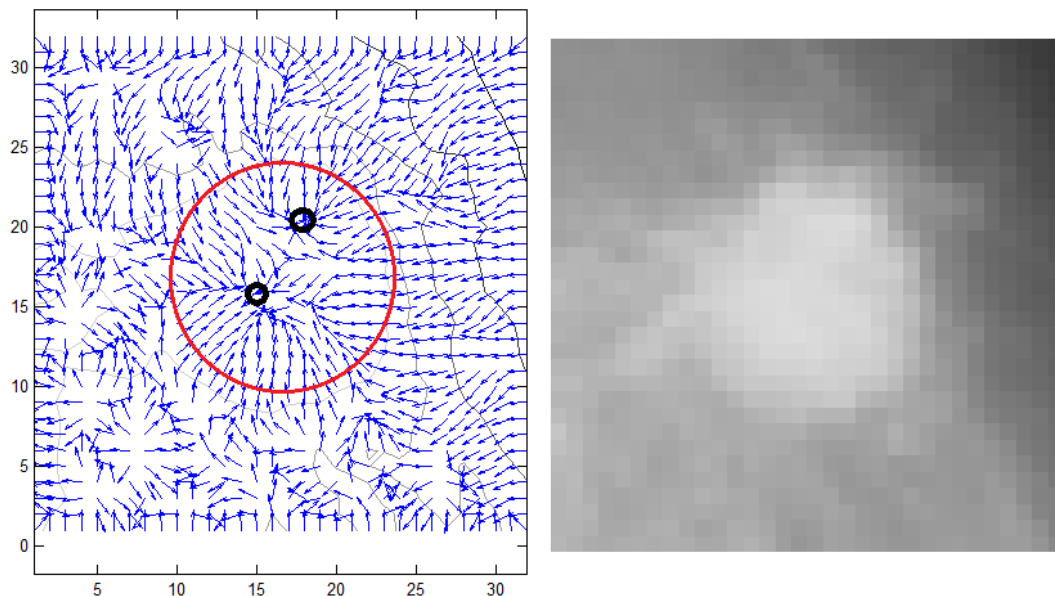


Figure 3.3: Left:Gradient map of the mass region with two local maximums, Right:Mass region.

Masses are detected by Iris filter generally when the edge region of the masses are in the region of support. When the region of support is chosen to be in the mass region, Iris filter may produce high values or not, it depends on the gradient vectors in the mass region. This situation is illustrated in Figure 3.4.

In order to understand this situation better convergence index maps are obtained for different support regions. Chosen pixel of interest and the convergence index map for a region of support, with  $R_{min} = 2$  and  $R_{max} = 62$  pixels, are given in Figure 3.5. Pixel of interest is chosen to be a point in a mass region to observe the change in convergence index map with respect to changing region of support. Different regions of support, that will be mentioned below, are obtained from this large region.

Figure 3.6 shows the convergence index map for a region of support with  $R_{min} = 2$  and  $R_{max} = 22$ , pixels. It should be noted that the high convergence index values do not start at pixels very close to  $R_{min}$ . Probably, these pixels are not on the edge region of the mass but in the mass region. Looking at the overall Iris filter algorithm this convergence index map may produce high filter output value for the pixel of interest,

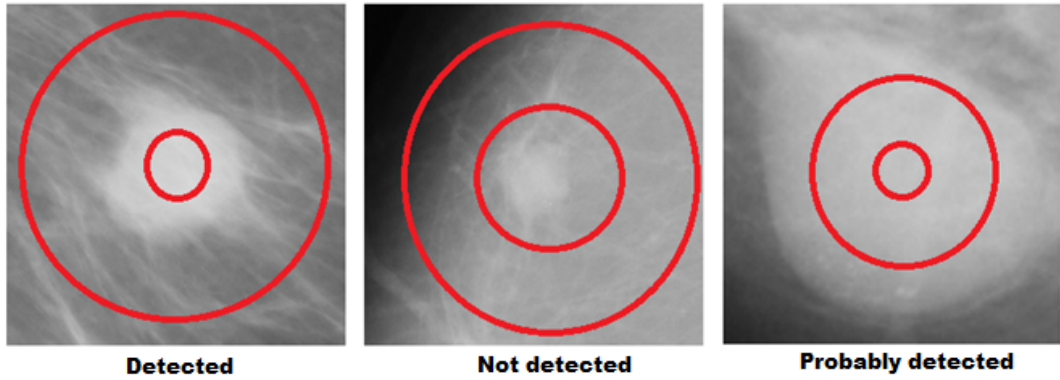


Figure 3.4: Iris filter’s support region and mass detection. Left:Mass edge falls between  $R_{min}$  circle and  $R_{max}$  circle. Middle:Mass is smaller than  $R_{min}$  circle. Right:Mass region covers  $R_{out}$  circle.

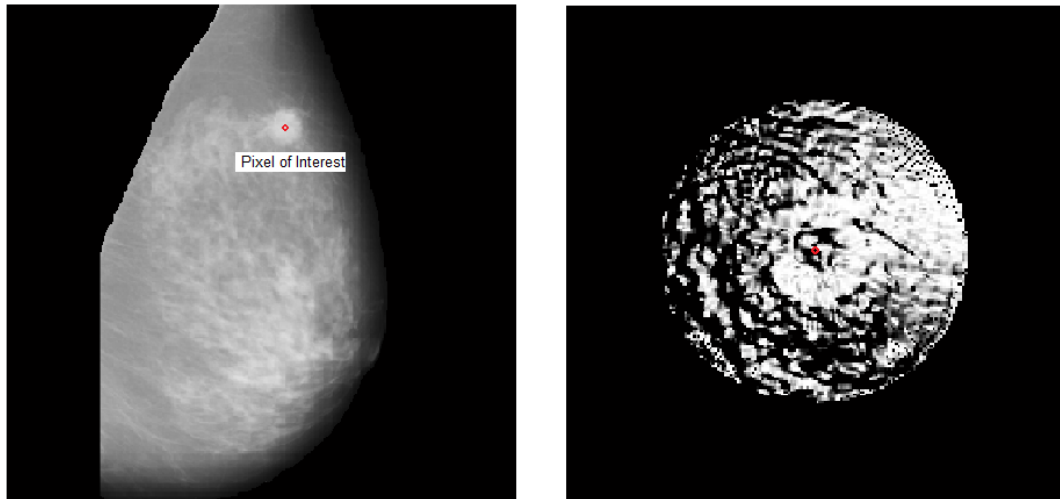


Figure 3.5: Left:Pixel of interest on the mammogram, Right:Convergence index map for each pixel in the region of support for chosen pixel of interest ( $R_{min} = 2$ ,  $R_{max} = 62$ ).

but not the highest one since pixels close to  $R_{min}$  have small convergence index values and the integral of convergence degree starts from  $R_{min}$ .

Figure 3.7 shows the convergence index map for a region of support with  $R_{min} = 12$  and  $R_{max} = 32$ . It should be noted that the high convergence index values start at pixels very close to  $R_{min}$ . Probably, these pixels are in the edge region of the mass on the yellow region drawn. Looking at the overall Iris filter algorithm this convergence index map may produce highest output value for the pixel of interest.

Figure 3.8 shows the convergence index map for a region of support with  $R_{min} = 22$  and  $R_{max} = 42$ . The highest values of convergence index map are on the top right part of the convergence index map. This region is breast boundary region drawn with

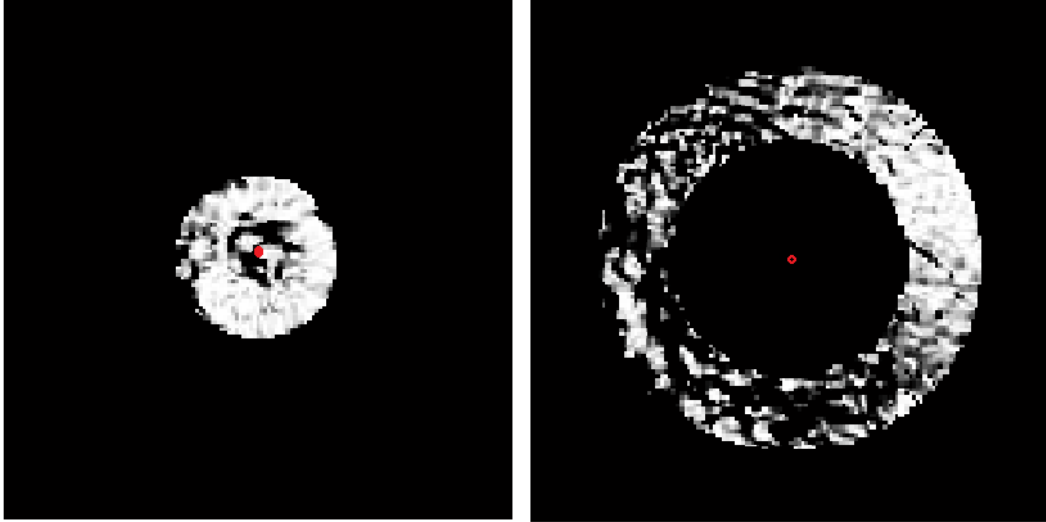


Figure 3.6: Left:Convergence index map for region of support, with  $R_{min} = 2$  and  $R_{max} = 22$ , Right:Convergence index map for region of support with  $R_{min} = 32$  and  $R_{max} = 52$ .

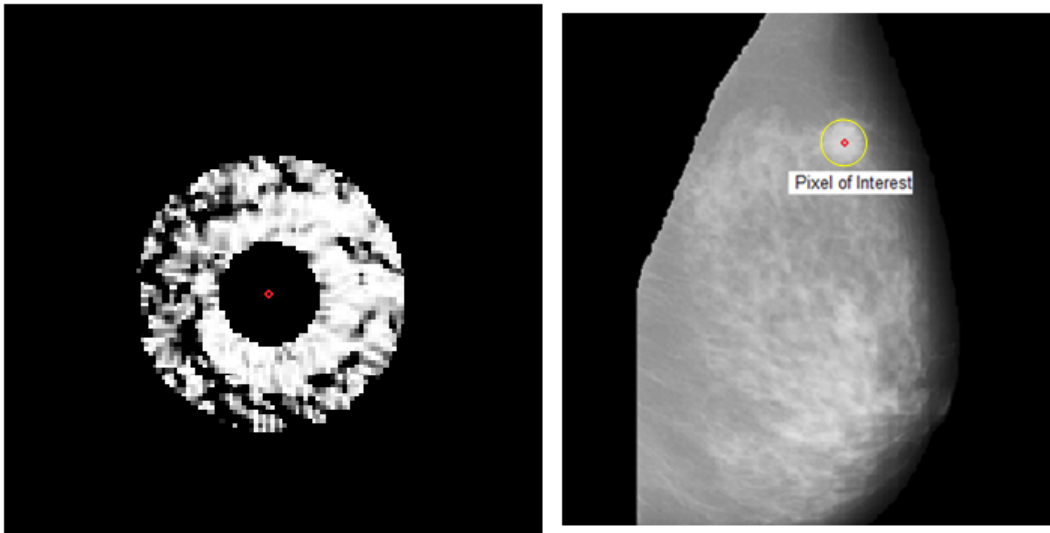


Figure 3.7: Left:Convergence index map for region of support, with  $R_{min} = 12$  and  $R_{max} = 32$ , Right:Pixel of interest and largest convergence index region.

yellow color in the figure. Breast boundary region is one of the source of false positives.

Figure 3.6 shows the convergence index map for a region of support with  $R_{min} = 32$  and  $R_{max} = 52$ . As in the previous map, the effect of breast boundary is also absolutely seen on this map.

Masses have different diameters from 3mm (15 pixels for an image in mini-MIAS) to 40 mm (200 pixels for an image in mini-MIAS) [30]. Hence, mass radius can occur



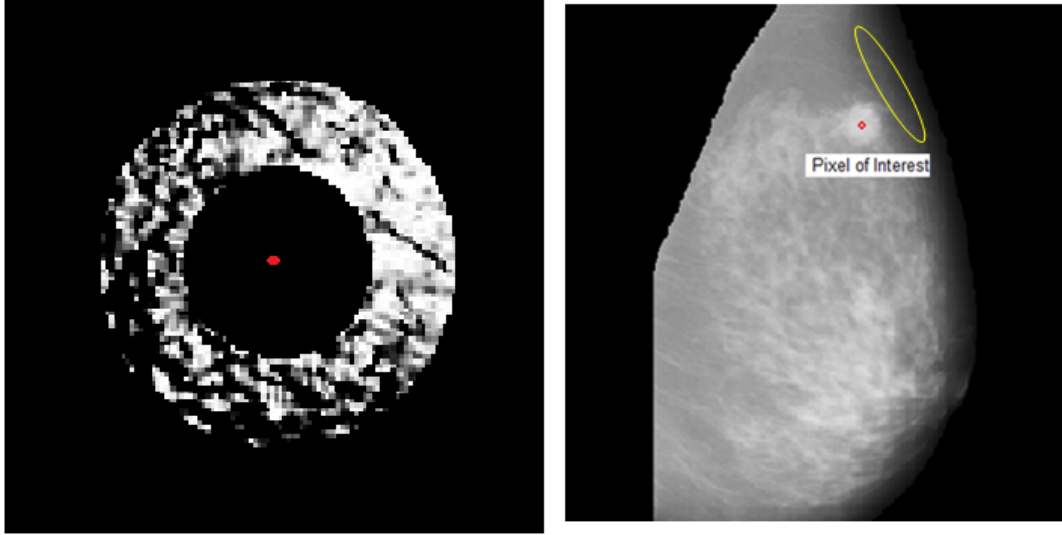


Figure 3.8: Left:Convergence index map for each pixel in the region of support, with  $R_{min} = 22$  and  $R_{max} = 42$ , for chosen pixel of interest, Right:Pixel of interest and largest convergence index region.

between 4 pixels and 50 pixels for mammograms with dimension 512 x 512. Mass size must be taken into account while running the Iris filter algorithm. In this thesis, five different Iris filters are applied on mammograms. The  $R_{min}$  and  $R_{max}$  values of these filters are given in Table 3.1. Regions are overlapping and wide enough not to miss masses.

Table3.1:  $R_{min}$ ,  $R_{max}$  values of Iris filters.

Filter	$R_{min}$	$R_{max}$
Iris filter 1	2	22
Iris filter 2	12	32
Iris filter 3	22	42
Iris filter 4	32	52
Iris filter 5	42	62

### 3.2.1 Iris filter implementation

Iris filter implementation algorithm is given in Figure 3.9. Firstly, Iris filters, with different support regions, are applied on mammogram. Therefore, 5 different Iris filter outputs are obtained. Simple threshold operation is implemented on these outputs; but the threshold levels may be different for each filter. Connected component analysis is made and the centers of the suspicious mass regions are calculated for each candidate region. Finally, regions of size 128 x 128 or 256 x 256 are obtained for each candidate mass center.

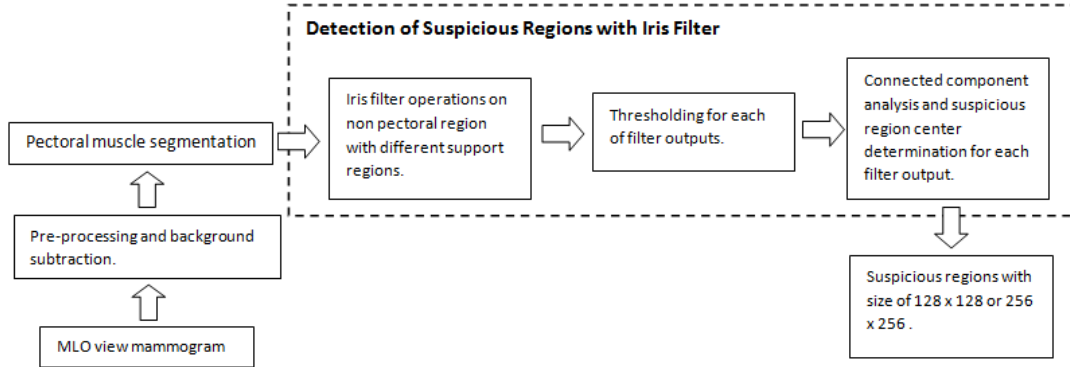


Figure 3.9: Iris filter implementation.

Iris filter may be time consuming in accordance with the chosen line number. Computation of convergence index and convergence degree for the pixels through which the lines pass, need process time. In [30] line number is determined as 20. In this thesis, detailed results of Iris filters are expected so that Iris filter with 360 lines are applied on mammograms.

Another critical thing that may be time consuming is the grid size. Grid size is the concept that the space difference between one pixel of interest and the other one, i.e if grid size is 1 the process will be done for all pixels. If grid size is different from one then for a horizontal or vertical line, the process will be done for one pixel among pixels as many as the grid size number. Disadvantage of this technique is that for large grid sizes the risk of missing masses is high since masses may occur in small sizes, too.

In this thesis, the process is done for selected pixels; it is not done randomly as grid size approach. In grid size approach the selected pixels may be a low intensity pixel (fat tissue etc.). However, it is known that masses are mostly high intensity looking parts of a mammogram. Therefore, pixel selection is done according to its density.

Firstly, conventional histogram equalization ([52]) is done on breast region without pectoral muscle. Histograms of an original breast region without pectoral muscle and histogram equalized image are given in Figure 3.10, 3.11, respectively.

Secondly, pixels close to intensity value 0,5 are found on histogram equalized image. Mean of these pixels' intensities on the original mammogram is calculated and this value is designated as a threshold value for Iris filter implementation. If the intensity value of the pixel of interest in the original image is greater than this value then Iris filter process is done for this pixel. If it is smaller, no process is done and 0 is assigned to this pixel on the Iris filter output. It is crucial to note that the task of histogram equalization is only determination of this intensity value. Iris filters are implemented on the breast region, without pectoral muscle, of original mammogram.

Separate runs are not made for different Iris filters, instead one run is made for a

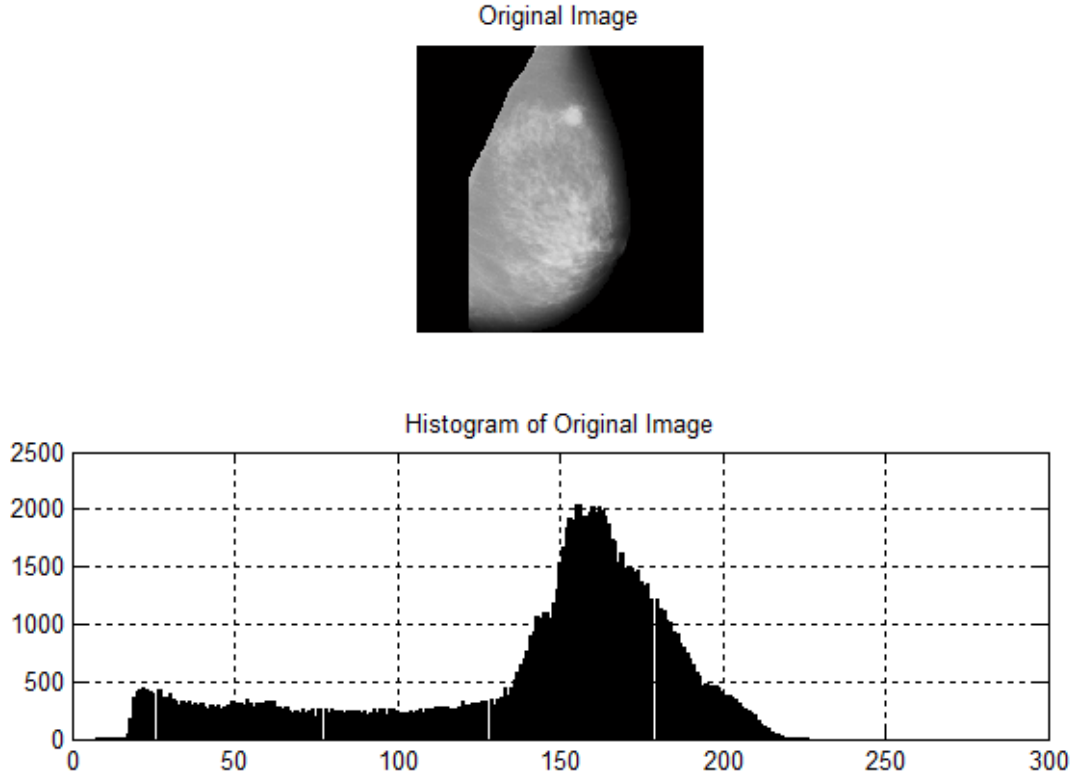


Figure 3.10: Original breast region and its histogram.

region of support including all other filter support regions. For instance, for Iris filters determined above only one convergence index map ( $R_{min} = 2$ ,  $R_{max} = 62$ ) is obtained and from this map, convergence degree and filter outputs are calculated to shorten run time.

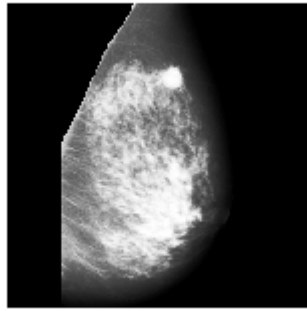
### 3.2.2 Outputs of Iris filters

Output of "Iris filter 1", applied on breast region, is given on the left side of Figure 3.12. Black regions on the breast texture are the pixels that no filter calculation is done and 0 is assigned to.

Pixels below the threshold level are assigned 0 and over the threshold level are assigned 1. Threshold applied result of Iris filter output, summed with original mammogram, is given on the right side of Figure 3.12. Although high threshold level is applied, great number of false positives come from this filter because the possibility of finding local maximum points with small region is high for a mammogram image.

Radius of the mass on top region of the mammogram given in right side of Figure 3.12 is 39. However, the suspicious points are found in the mass region. Actually, these points are local maximums on the mass region as mentioned before.

Histogram Equalized Image



Equalized Histogram of Image

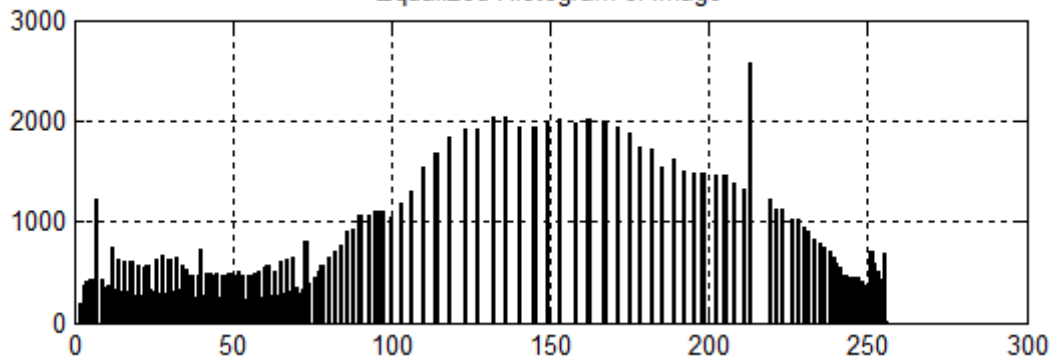


Figure 3.11: Histogram equalized breast region and its histogram.

Output of "Iris filter 2", applied on breast region is given on the left side of Figure 3.13.

Threshold applied result of Iris filter output, summed with original mammogram, is given on the right side of Figure 3.13. Mass edge falls in region of support for a pixel of interest on the mass region so that so many suspicious points are signed on the mass region. There are false positive pixels, having higher intensity values than the surrounding tissues.

Output of "Iris filter 3", applied on breast region, is given on the left side of Figure 3.14.

Threshold applied result of Iris filter output, summed with original mammogram, is given on the right side of Figure 3.14. Mass edge falls in region of support for a pixel of interest on the mass region so that many suspicious points are signed on the mass region. However, this time since radius of the mass falls between  $R_{min}$  and  $R_{max}$  values, the suspicious pixels clustered on the center of the mass region. It is also observed that some of the false positives signed on the left part of the breast are due to the breast boundary falling into the region of support.

Output of "Iris filter 4", applied on breast region, is given on the left side of Figure

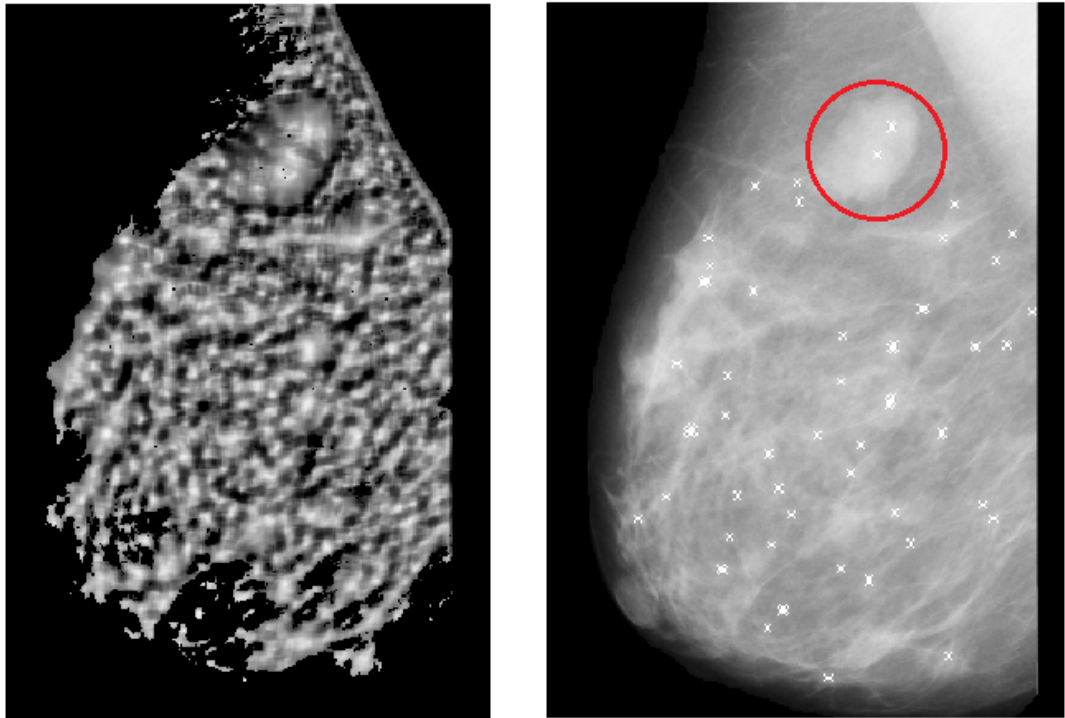


Figure 3.12: Left:Output of Iris filter 1, Right: Threshold level = 0,90 applied to filter output added to original mammogram.

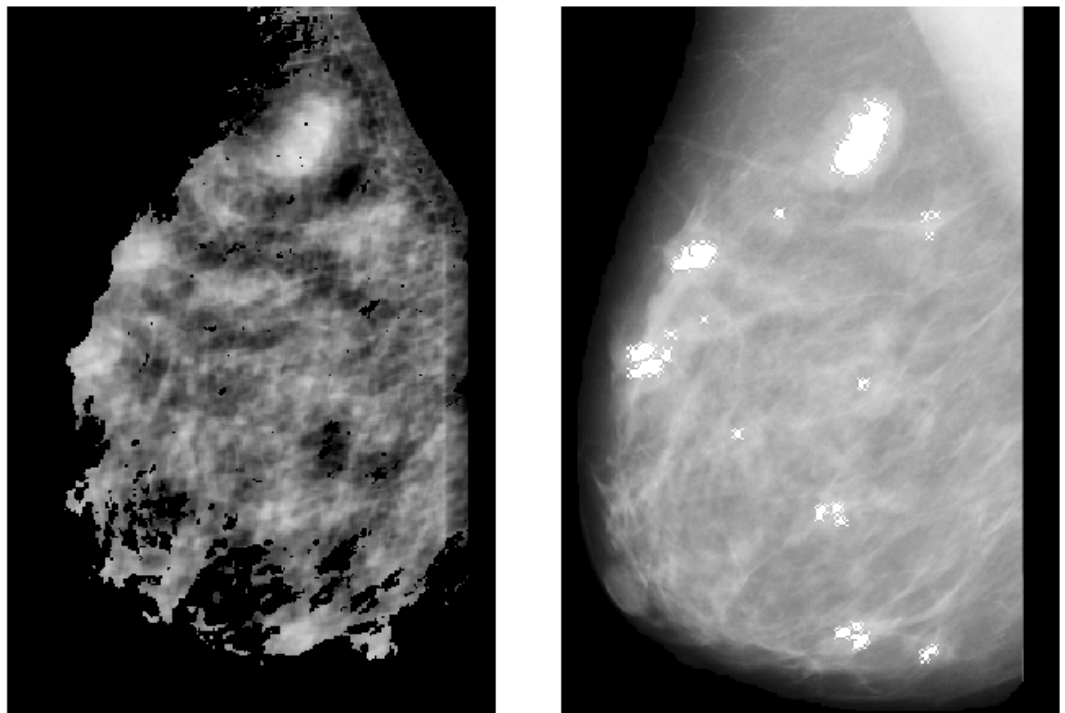


Figure 3.13: Left:Output of Iris filter 2, Right: Threshold level = 0,75 applied to filter output added to original mammogram.

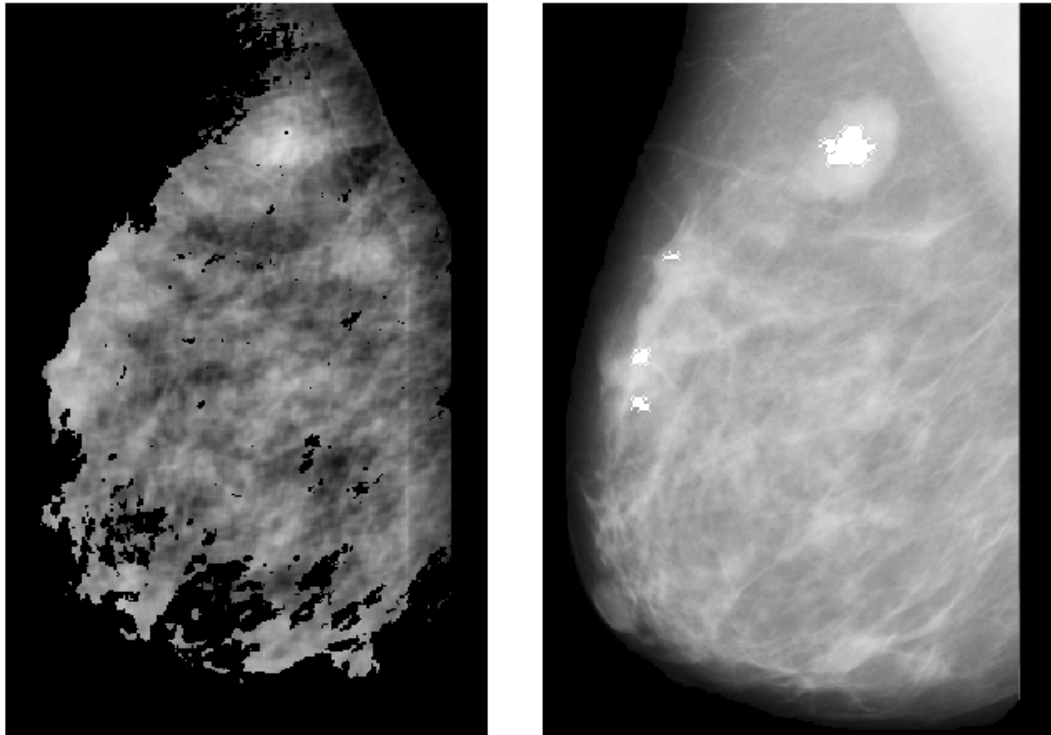


Figure 3.14: Left:Output of Iris filter 3, Right: Threshold level = 0,75 applied to filter output added to original mammogram.

3.15.

Threshold applied result of Iris filter output, summed with original mammogram, is given on the right side of Figure 3.15. If the suspicious pixels on the mass region are examined it is not difficult to see that the right part of the mass region falls in support region for these pixels and the left part of the region of support falls close to breast boundary. Since threshold is 0,75 it is possible for these pixels to exceed the threshold level and labeled as suspicious point.

Output of "Iris filter 5", applied on breast region, is given on left side of Figure 3.16.

Threshold applied result of Iris filter output, summed with original mammogram, is given on the right side of Figure 3.16. No pixel exceeds the threshold level for this filter output. If the mammogram is examined carefully correctness of this situation may be observed because there is not any region with an edge of this size.

### 3.2.3 Designation of suspicious regions

It is easily observed that the suspicious pixels are clustered and connected. Perhaps the choice of one region can contain all suspicious pixels. Therefore, a center may be designated for each pixel group then region may be determined. In order to group

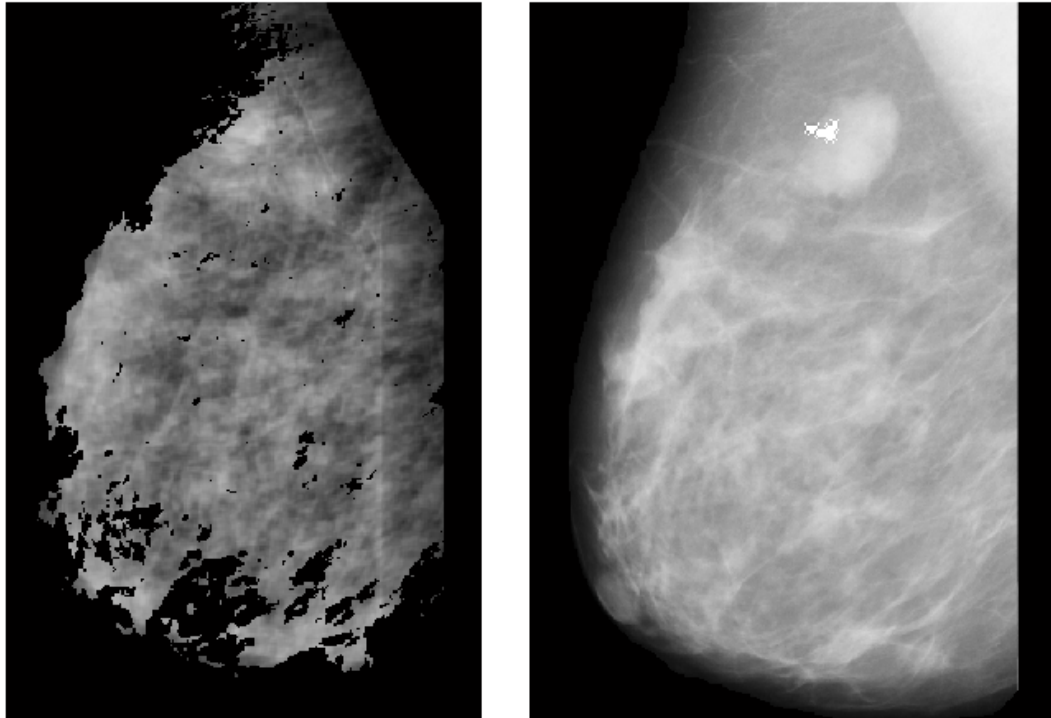


Figure 3.15: Left:Output of Iris filter 4, Right: Threshold level = 0,75 applied to filter output added to original mammogram.

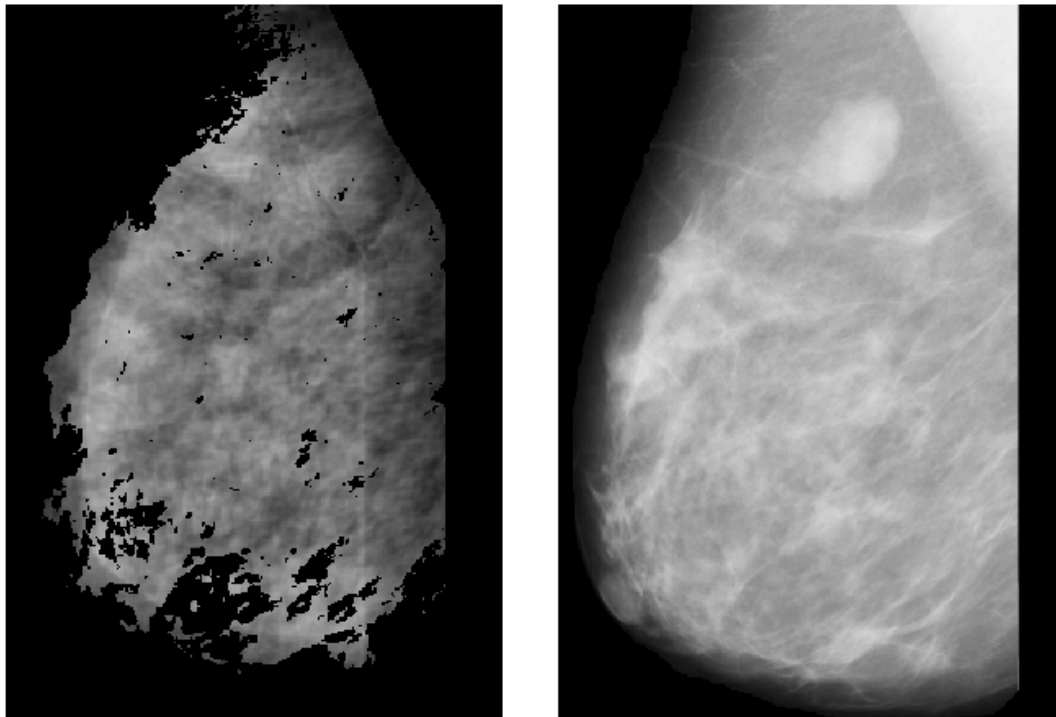


Figure 3.16: Left:Output of Iris filter 5, Right: Threshold level = 0,75 applied to filter output added to original mammogram.

suspicious pixels, connected component labeling is used.

After all groups are found the center of each group is calculated. The center is assumed to be simply the mean of the pixel coordinates in a group.

Center determination results of filtering and threshold operations for Iris filters 2 and 3 are given in Figure 3.17 and 3.18, respectively.

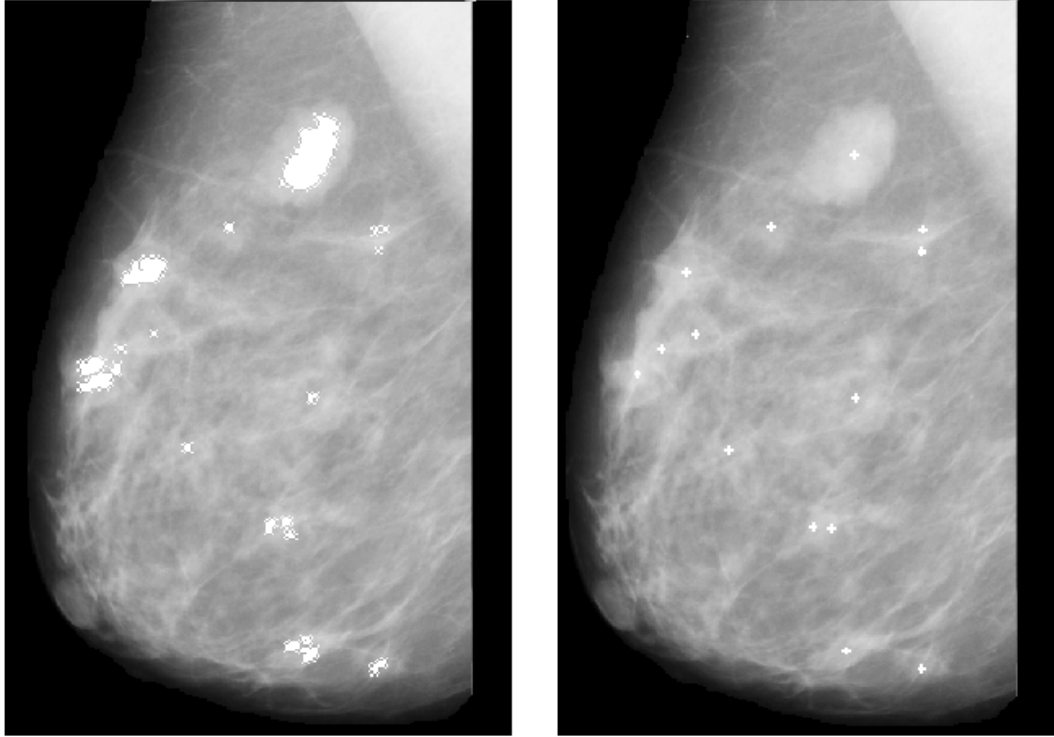


Figure 3.17: Left:Suspicious pixels obtained after filtering and threshold operation for Iris filter 2, Right:Connected component labeling and center determination result.

1024 x 1024 images from database were resized to dimension 512 x 512 before pectoral boundary calculations. The centers are calculated for images of size 512 x 512. Hereafter, the center pixel coordinates will be found on original images of dimension 1024 x 1024 and the regions will be selected from these images.

128 x 128 regions are obtained for center outputs of filtering with Iris filters 1 and 2. The center is assumed to be in the middle of the region (Figure 3.19). In addition, 256 x 256 regions are obtained for center outputs of filtering with Iris filters 3, 4 and 5 (Figure 3.20). Reason of region extraction with different dimensions is given in Section 4.2.2. This region dimension determination is critical in terms of classification performance. In other words, adaptive region size determination is one of the contributions of this thesis that affects the performance of the classification (Section 4.2.2).

In some cases while defining the region boundary for a center, second center may fall in the region as illustrated in Figure 3.21. Circle is drawn on the second center pixel



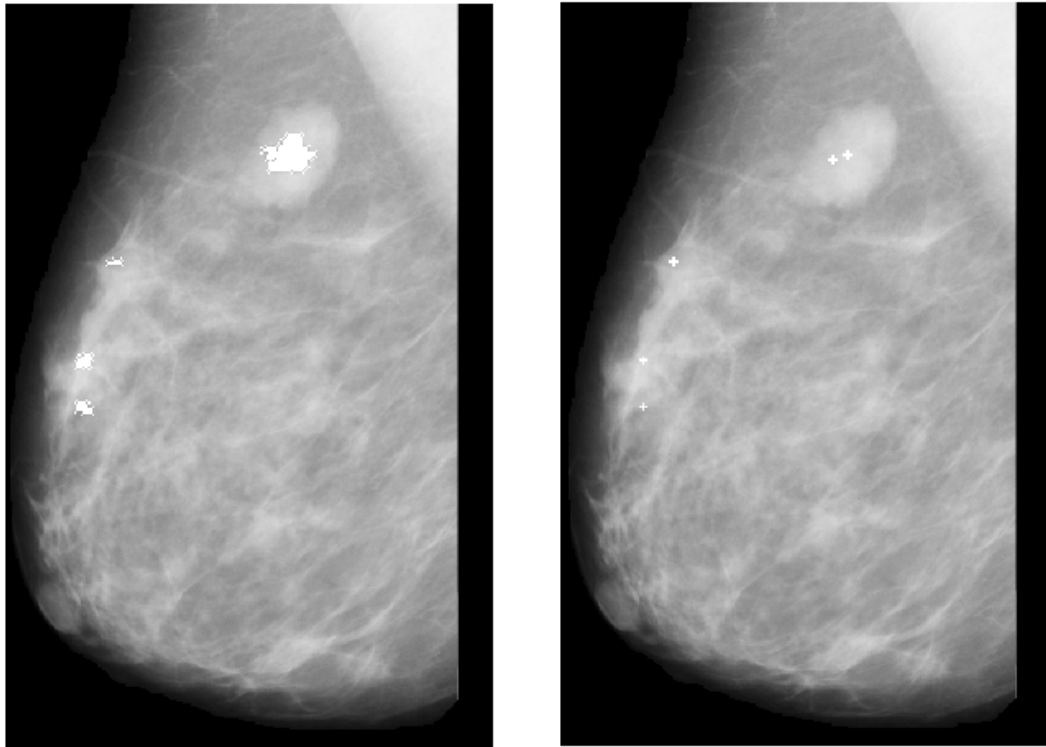


Figure 3.18: Left:Suspicious pixels obtained after filtering and threshold operation for Iris filter 3, Right:Connected component labeling and center determination result.

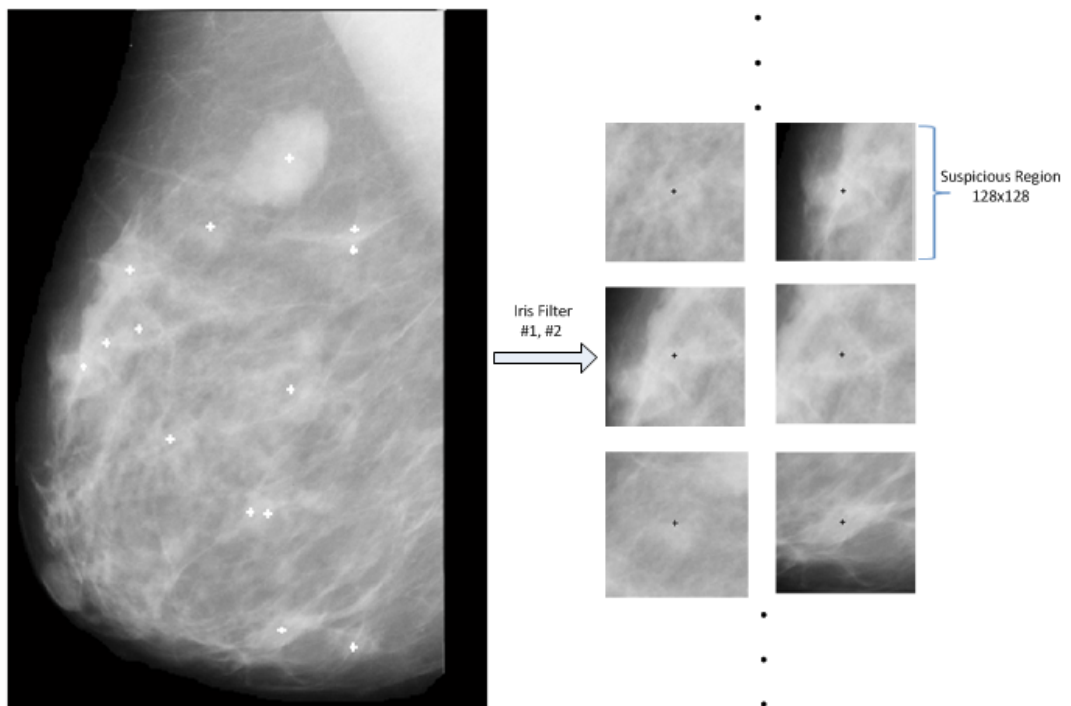


Figure 3.19: Region determination for centers obtained by Iris filters 1, 2 cases.

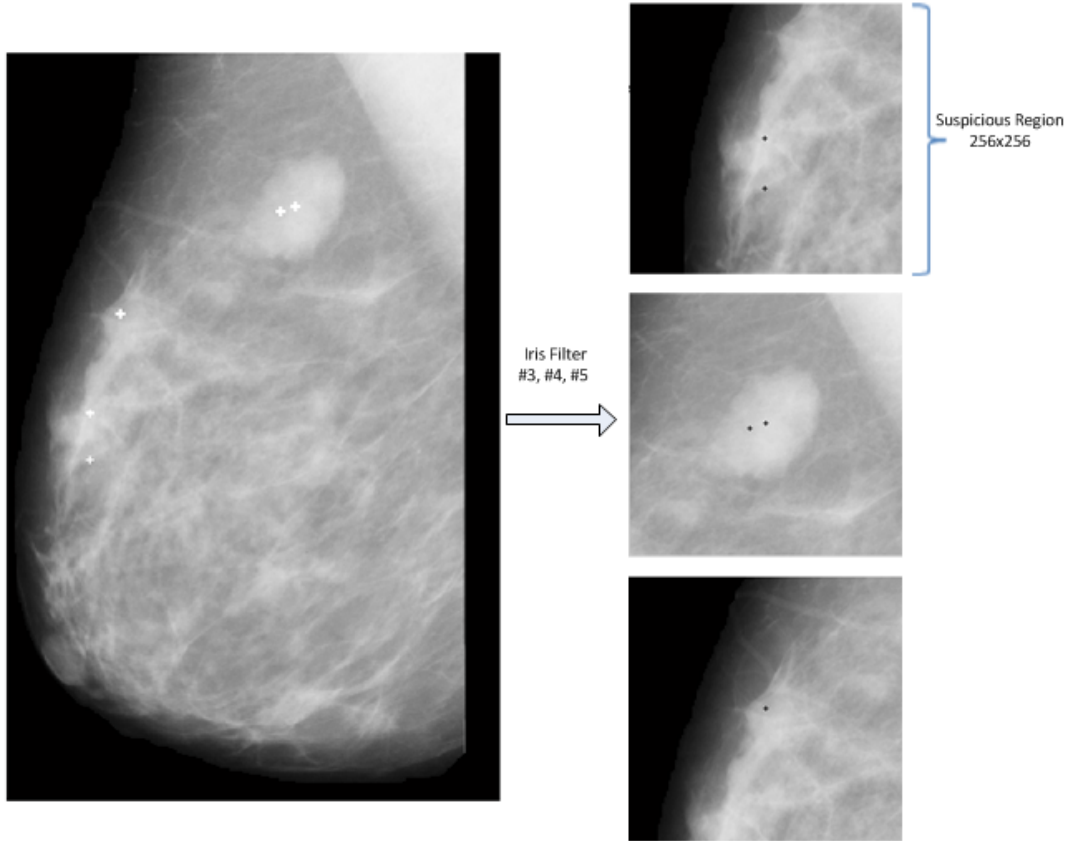


Figure 3.20: Region determination for centers obtained by Iris filters 3, 4, 5 cases.

with a radius  $R_{max}$  of the support region. If this circle is inside the region then no region determination is done for the second center and it is assumed to be in the region determined by the first center (Figure 3.21).

### 3.3 Results of Suspicious ROI Detection Algorithm

Iris filter approach is satisfactory in terms of sensitivity as mentioned before. Mass regions, in mammograms, are expected to be detected. All masses are detected except one mass, in this thesis. One spiculated malignant mass is not detected with Iris filter approach. Malignant mass and its gradient map are given in Figure 3.22. It is observed that gradients on the edge of mass are not directed towards a point in the mass region. Hence, Iris filters do not produce high outputs and pixels inside mass region stays below threshold levels for all filters.

Statistic of false regions detected by Iris filters is given in Table 3.2. It was previously mentioned that filter produces high FPpI rate. It is observed that different filters have different FPpI performance. Most of the FPs are gained from "Iris filter 1" although threshold level chosen is the highest one for this filter output. Moreover, in case of

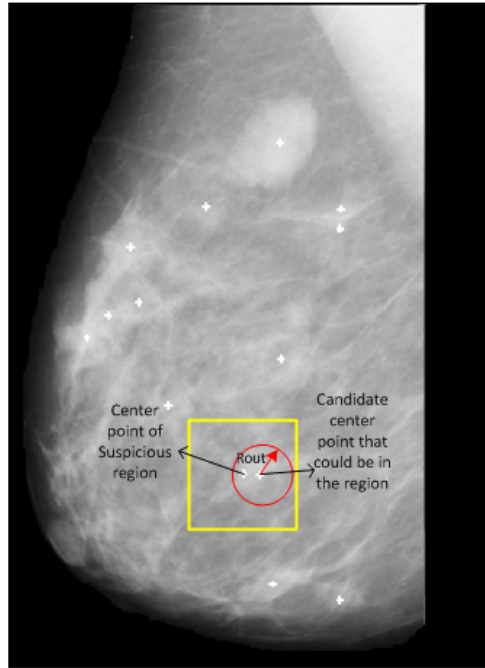


Figure 3.21: Region determination when centers are close to each other.

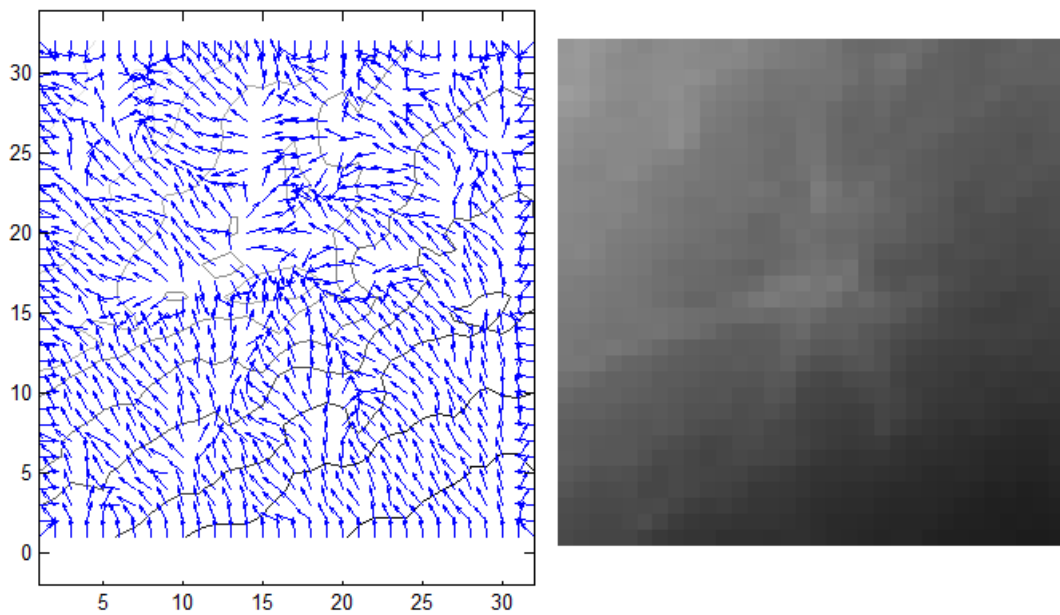


Figure 3.22: Left: Gradient map of spiculated malignant mass region, Right: Spiculated malignant mass region.

other filters in spite of the fact that threshold levels are chosen equal, different FPpI performance is obtained for each filter. Therefore, support region of a filter is critical in terms of FP performance. In this thesis, threshold levels are chosen experimentally and FP performance is not taken into consideration while determining threshold levels. A detailed analysis of threshold levels for each Iris filter can be made and optimum

thresholds can be determined in the future.

Table3.2: Mean FPpI for each Iris filter.

Iris filters	Mean FPpI
Iris filter 1	28,41
Iris filter 2	12,67
Iris filter 3	5,46
Iris filter 4	4,65
Iris filter 5	4,82

## CHAPTER 4

### SVM CLASSIFICATION USING FEATURES, EXTRACTED WITH GABOR FILTERS

After implementation of Iris filter, suspicious regions are determined. However, there are so many FPs as expected. In this chapter, FPs are tried to be reduced with a classification step. The features that are used in classification are obtained by a textural analysis, particularly Gabor filter bank application on sub-regions of possible mass regions.

#### 4.1 Feature Extraction with Gabor Filter Bank

All the stages of feature extraction with Gabor filter bank is given in Figure 4.1.

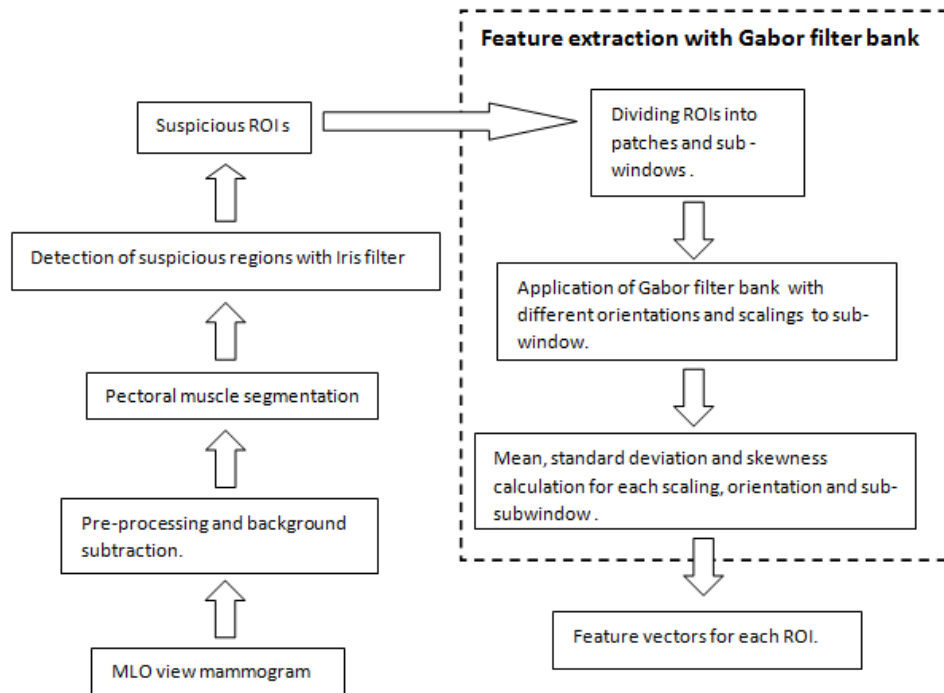


Figure 4.1: Feature extraction stages with Gabor filter bank.

### 4.1.1 Gabor filter bank

Gabor filters are kernels that are widely used in computer vision and image processing field applications e.g. face recognition ([54]), vehicle detection ([55]), etc. Gabor filters provide powerful statistics that could be used while extracting local spatial textural micro-patterns in mass detection problem since they have property to be tuned to different orientations and scales [56].

A two-dimensional Gabor filter is defined as a Gaussian kernel modulated by an oriented complex sinusoidal wave and can be described as in Equation 4.1:

$$g(x, y) = \frac{1}{2\pi\sigma_x\sigma_y} \exp^{-\frac{1}{2}\left(\frac{\tilde{x}^2}{\sigma_x^2} + \frac{\tilde{y}^2}{\sigma_y^2}\right)} \exp^{2\pi jW\tilde{x}} \quad (4.1a)$$

$$\tilde{x} = x\cos(\theta) + y\sin(\theta) \quad (4.1b)$$

$$\tilde{y} = -x\sin(\theta) + y\cos(\theta) \quad (4.1c)$$

where

$g(x, y)$  is the Gabor filter response of pixel of interest  $(x, y)$ .

$\sigma_x, \sigma_y$  are the scaling parameters of the filter and describe the neighborhood of a pixel where weighted summation takes place.

$W$  is the central frequency of the complex sinusoidal.

$\theta$  is the orientation of the normal to the parallel stripes of the Gabor function, is in the interval  $[0, \pi)$

A Gabor filter bank contain multiple individual Gabor filters that are adjusted with different parameters (scaling(S), orientation(O) and central frequency). 40 filters (5 S x 8 O) with initial maximum frequency equal to 0,2 and initial orientation set to 0 is used in this thesis. Selected frequencies are (0,2; 0,14; 0,1; 0,07; 0,05) and orientation angles are (0; 22,5; 45; 67,5; 90; 112,5; 135; 157,5 in degrees). Gabor filter kernels' dimension is 20 x 20. Space response of real part of the Gabor filter bank is given in Figure 4.2 for all Gabor filters. Total frequency response of the bank is also given in Figure 4.3.

The orientations and frequencies for a bank are calculated using Equation 4.2:

$$orientation(i) = \frac{(i-1)\pi}{m}; i = 1, 2, \dots, m \quad (4.2a)$$

$$frequency(i) = \frac{f_{max=2}}{(\sqrt{2})^{i-1}}; i = 1, 2, \dots, n \quad (4.2b)$$

where

$m$  is the total number of orientations.

$n$  is the total number of frequencies.

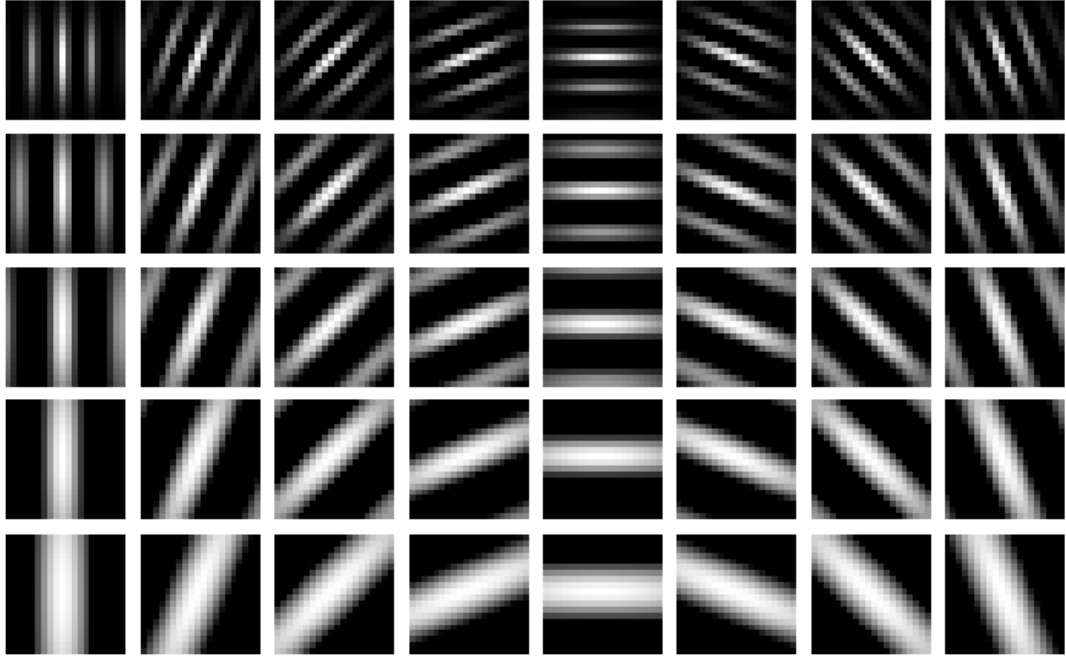


Figure 4.2: Gabor filter bank: filters in the same column have the same orientation, filters in the same row has the same frequency.

#### 4.1.2 Partition of ROIs

As mentioned before Gabor filter bank is applied on sub-windows but not the whole region. Suspicious ROIs extracted in the previous step are divided into patches and sub-windows as given in Figure 4.4.

Firstly, the whole ROI is divided into patches. There will be 16 patches for each ROI in this thesis. Patch number could be selected to be different ([49]). Size of a patch, in a region with a dimension 128 x 128, is 32 whereas patch size, in a region with a dimension 256 x 256, is 64 (Figure 4.4).

Secondly, sub-windows are determined. 4 neighbour patches creates one sub-window. For instance, patches 1, 2, 5, 6 creates first sub-window; 2, 3, 6, 7 creates the second one, etc. A total of 9 sub-windows will be created. Information of patch numbers included in sub-windows is given in Table 4.1.

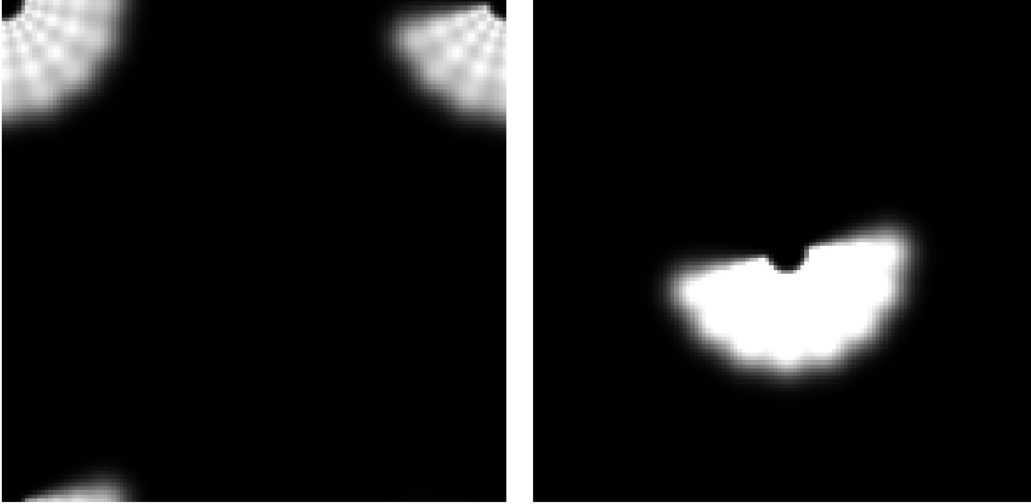


Figure 4.3: Left: Combined frequency response of Gabor filters, Right: Combined frequency response of Gabor filters without frequency shift.

Table 4.1: Patch numbers creating sub-windows.

Sub-window number	Patch numbers
1	1, 2, 5, 6
2	2, 3, 6, 7
3	3, 4, 7, 8
4	5, 6, 9, 10
5	6, 10, 7, 11
6	7, 8, 11, 12
7	9, 10, 13, 14
8	10, 11, 14, 15
9	11, 12, 15, 16

### 4.1.3 Feature extraction

Each sub-window is convolved with Gabor filter bank to extract the features. Magnitude response of each Gabor filter in the bank is collected from each sub-window and each one is represented by three moments: the mean ( $\mu_{l,m}$ ), the standard deviation ( $\sigma_{l,m}$ ) and the skewness ( $k_{l,m}$ ) (where  $l$  corresponds to the  $l$  th filter in the bank and  $m$  to the  $m$  th sub-window).

In this thesis, a Gabor filter bank of 40 filters (5 S x 8 O) is used. If this bank is applied on 9 sub-windows (Figure 4.4) of a single ROI, a feature vector of length 1080 will be obtained. One feature vector row obtained in this way and is given below:

$$\mu_{1,1}, \sigma_{1,1}, k_{1,1}, \dots, \mu_{40,1}, \sigma_{40,1}, k_{40,1}, \mu_{1,2}, \sigma_{1,2}, k_{1,2}, \dots, \mu_{40,9}, \sigma_{40,9}, k_{40,9}.$$



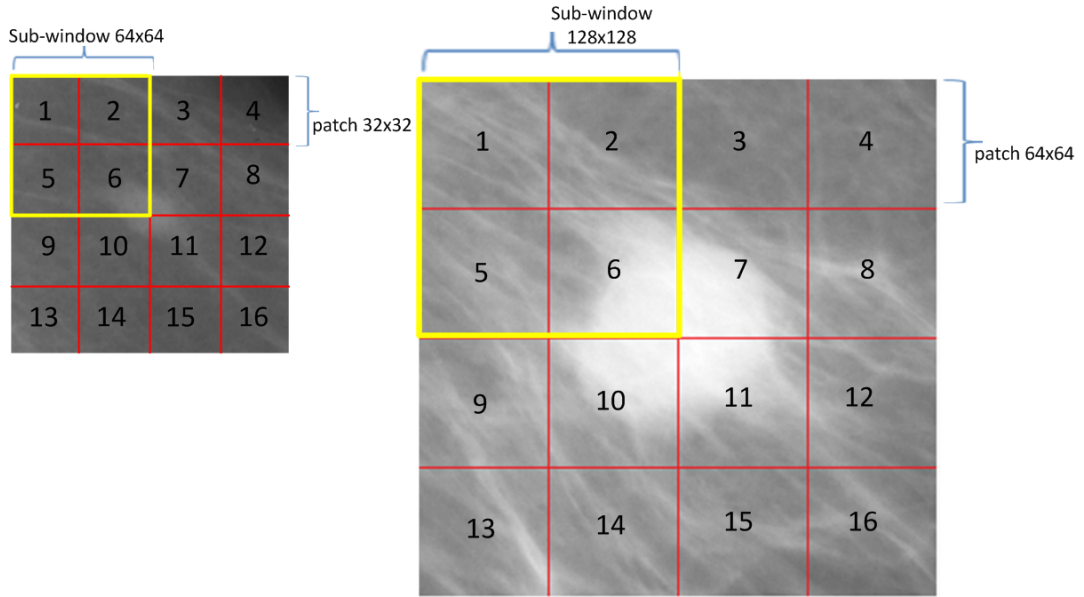


Figure 4.4: Left:Patches and sub-windows for ROIs with size 128 x 128, Right:Patches and sub-windows for ROIs with size 256 x 256.

## 4.2 Classification

A classification problem includes the assignment of an unseen pattern to a predefined class, with the characteristics of the pattern, presented in the form of a feature vector. However, a classifier is needed to be trained in order to perform this task firstly. Selected patterns of the same concept class are used to train the classifier. This set is called training set. Unseen patterns for which assignment process is done, form test set.

Numerous classification techniques exist. SVM (Support Vector Machine) is offered to solve binary classification problem ([56]). An optimal hyper-plane, separating the data belonging to different classes with large margin, is found by SVM [57]. This decision boundary is based on the most "informative" points of the training set. These informative points are called support vectors (Figure 4.5).

SVM is different from other other classifiers in terms of the way of handling risk concept. Although other classifiers deal with empirical risk that minimize error on training data, SVM deals with structural risk to maximize the the margin between samples for different classes.

Let  $Train = \{(x_i, y_i)\}_{i=1}^N$  be a training set.  $x_i$  is the  $i$ th training instance containing  $J$  features.  $y_i$  is the class label of  $x_i$ .  $y$  has two values  $+1, -1$ . A constrained optimization problem is solved using quadratic programming to find an optimal hyper-plane based on large margin framework (Equation 4.3):

$$f(x) = \sum_{i=1}^N a_i y_i k(x_i, x) + b \quad (4.3)$$

where

$a_i$  is the Lagrange multiplier.

$k(x_i, x)$  is the kernel function.

$f_x$ 's sign gives the membership class of  $x$ .

SVM kernel is simply dot product of the two given points in the input space if the problem is linearly separable. Otherwise, original input space is mapped to the higher dimensional space through a non-linear mapping function with suitable kernel. Misclassification penalty is controlled with regularization parameter ( $C$ ).

There are kernels such as linear, polynomial, sigmoid, radial-basis. Radial-basis function (rbf) kernel, which is widely used in the literature, is used in this work:

$$k(x_i, x) = \exp(-\gamma \|x_i - x\|) \quad (4.4)$$

where

$\gamma$  is the width of the kernel function and  $\gamma > 0$ .

Classification with SVM is illustrated in Figure 4.5 for a chosen  $C$  and  $\gamma$  couple. Test elements are represented as points and hyper-plane separates true (mass) and false (non-mass) points. Vector from test points to hyper-plane is the margin vector for this point. All points on separator line have 0 magnitude margin vector. Margin vectors with unit magnitude are called support vectors. Margin vectors with positive values are expected to be returned for all negative test element and negative values for all positive test elements in case of an ideal classification. However, this condition is not satisfied mostly, so that ROC curves are calculated to determine SVM performance. ROC curves are found with changing hyper-plane position, in other words adding margin offset to hyper-plane.

#### 4.2.1 Determination of SVM parameters

$\gamma$  and  $C$  are two parameters effecting the performance of SVM classifier in a classification process. These two parameters are chosen by user before classification of suspicious ROIs determined. User may specify these parameters according to his / her

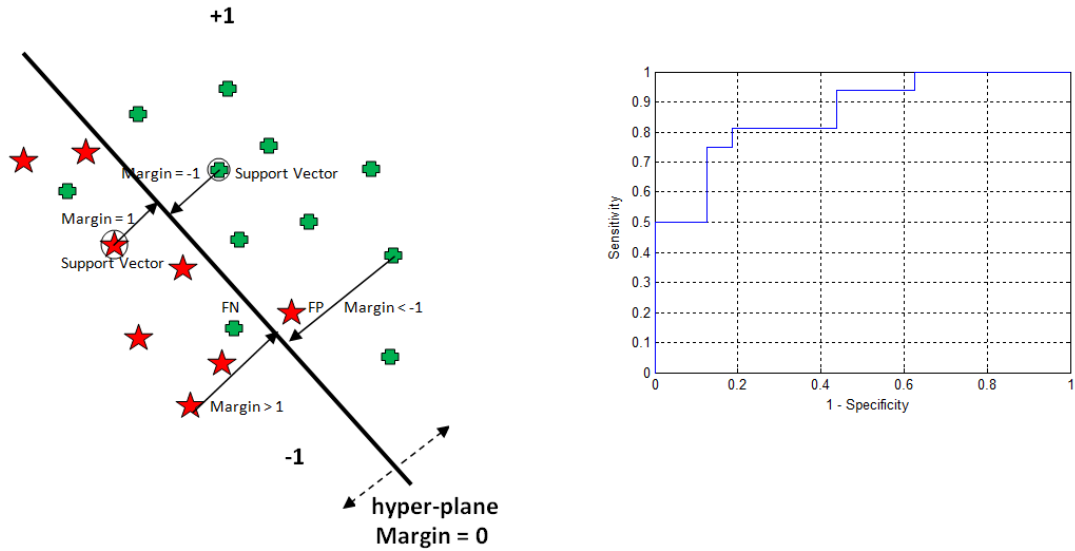


Figure 4.5: Left:SVM illustration, Right:ROC curve obtained by adding offsets to hyper-plane.

needs (high sensitivity, high truth rate, low FP, etc). Therefore, sensitivity, specificity truth rate analysis with changing  $\gamma$  and  $C$  parameters must be made with region sets. These region sets are selected from mini-MIAS. MIAS database includes information of whether there is a mass or not in a mammogram. Furthermore, if there is a mass, center pixel, radius and information of mass type is provided in database as mentioned before. Detailed information about quantity of mammograms, in terms of lesion existence and lesion type, is given in Table 4.2 for mini-MIAS database.

Table4.2: MIAS statistical data in terms of mass existence.

without mass	with mass	benign mass	malignant mass
179	50	37	13

Regions are selected with two different dimensions from database. If the radius of a mass is smaller than 55 pixels a region of dimension 128 x 128 is selected. If it is greater than 55 pixels a region of dimension 256 x 256 is selected. If there is no mass for a mammogram then a normal region is selected. Mass center is intersected with ROI center. Some of the selected 128 x 128 regions with mass and without mass are given in Figures 4.6, 4.7 respectively. In addition, selected 256 x 256 regions with mass and without mass are given in Figures 4.8, 4.9 respectively.

Statistic of selected ROIs is given in Table 4.3. Feature vectors are calculated for all selected ROIs to perform classification step.

Number of training and test sets affects the performance of SVM. For instance, when non-mass region number is higher in training set, classifier tends to produce negative (non-mass). However, if mass region number is higher, classifier is to produce posi-

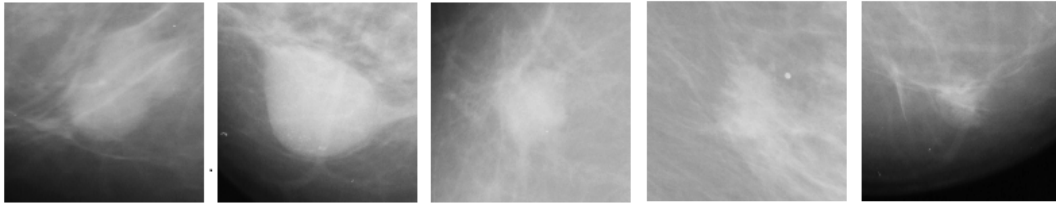


Figure 4.6: Selected regions of dimension 128 x 128 with mass.



Figure 4.7: Selected regions of dimension 128 x 128 without mass.

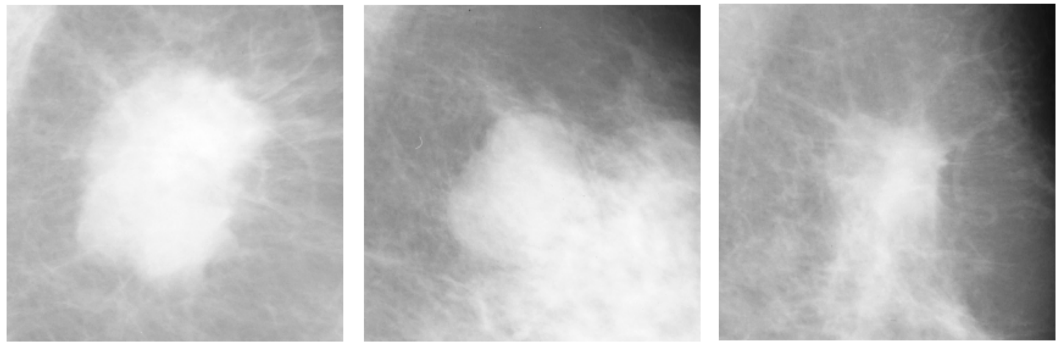


Figure 4.8: Selected regions of dimension 256 x 256 with mass.



Figure 4.9: Selected regions of dimension 256 x 256 without mass.

tive (mass) results. For this reason equal number of mass and non-mass regions are expected in training stage, to obtain balanced results. In addition, equal number of mass and non mass regions should be in test set to comment on performance metrics correctly.

Ideally, very large number mass and non-mass regions are expected to train SVM

Table4.3: ROI selection statistic.

Region size	with mass	without mass
128 x 128	34	178
256 x 256	16	178

in order to obtain exact and reliable results. However, as it is given in Table 4.3 number of regions with mass is very small (Only 34 for 128 x 128, 16 for 256 x 256) for this database. This is a limiting factor for the number of non-mass regions in the training set, too. Thus, classification is implemented with small number of regions. ROI numbers in training and test sets are given in Table 4.4. These sets are chosen randomly meaning that for each classification different training and test sets may be obtained and different performance could be acquired.

Table4.4: Statistic of training and test sets.

Region size	training sets with mass	training sets without mass	test sets with mass	test sets without mass
128 x 128	17	17	17	17
256 x 256	8	8	8	8

Classification is implemented for different  $C$  and  $\gamma$  values with randomly selected one training and one test set. These sets include as many regions as given in Table 4.4. Two separate classifications are made: one for 128 x 128 regions and the other for 256 x 256 ones. This is called one run. Classification performance, for  $C$  and  $\gamma$  values in a wide search range ( $-5 < \log_2(C) < 15$  and  $-2 < \log_2(\gamma) < 12$ ), is determined in one run. Performance metrics such as  $A_z$ , truth rate, sensitivity, (1 - specificity), obtained with respect to changing  $\gamma$  and  $C$  values for small and large selected regions, are given in Figures 4.10, 4.11, 4.12, 4.13, respectively.

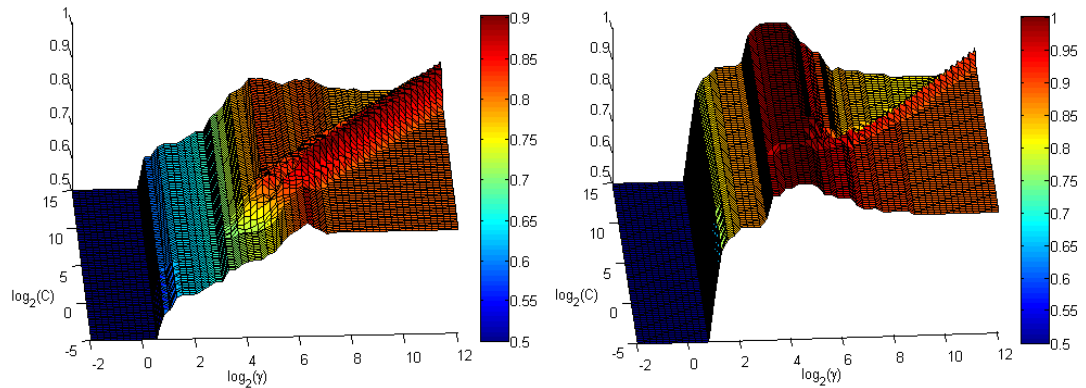


Figure 4.10: Left:  $A_z$  values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Right:  $A_z$  values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 1 run.

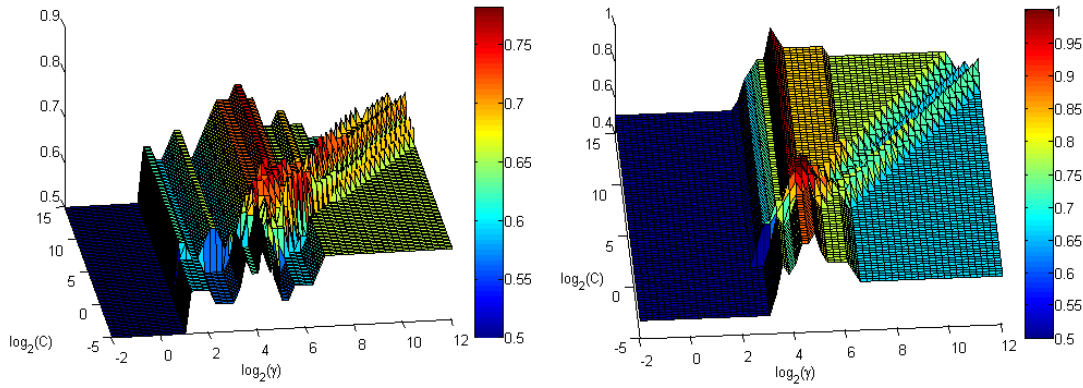


Figure 4.11: Left: Truth rate values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Right: Truth rate values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 1 run.

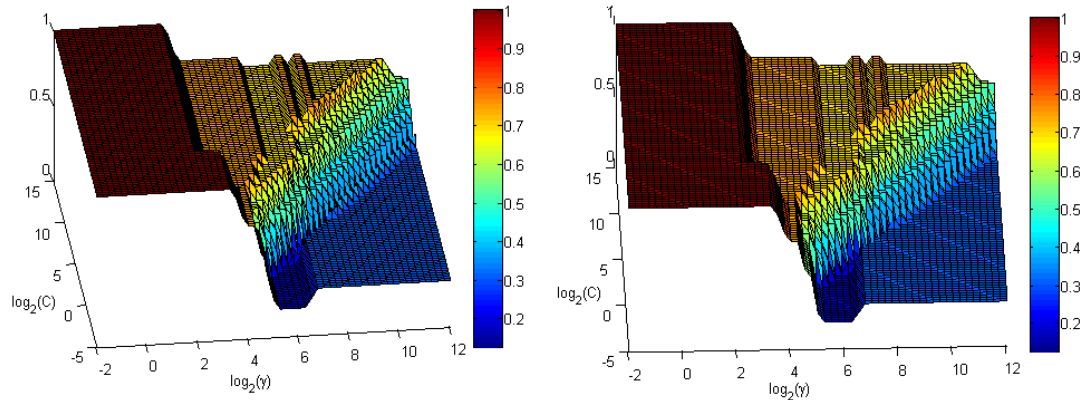


Figure 4.12: Left: Sensitivity values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Right: Sensitivity values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 1 run.

Truth rate is the rate of true decision for both mass and non-mass regions. Sensitivity and specificity concepts are explained previously. Specificity versus sensitivity performance may be seen on ROC curves as explained before. Here it is aimed to obtain performance points on such a plot to see (1 - specificity) versus sensitivity characteristic. This characteristic is illustrated in Figure 4.14. We name this illustration performance point plot and points on the plot, performance points.

It is crucial to note that a performance point in this graph, does not belong to only classification performance made with one  $C$  and  $\gamma$  couple. Same performance point could be obtained for different SVM parameters. In addition, there are performance points at only certain (1-specificity) and sensitivity values since limited number of mass and non-mass regions are tested. This quantization of the performance metrics' values is clearly seen for one run. Furthermore, there are more than one sensitivity value for a (1 - specificity) value meaning that sensitivity performance at one (1 - specificity)

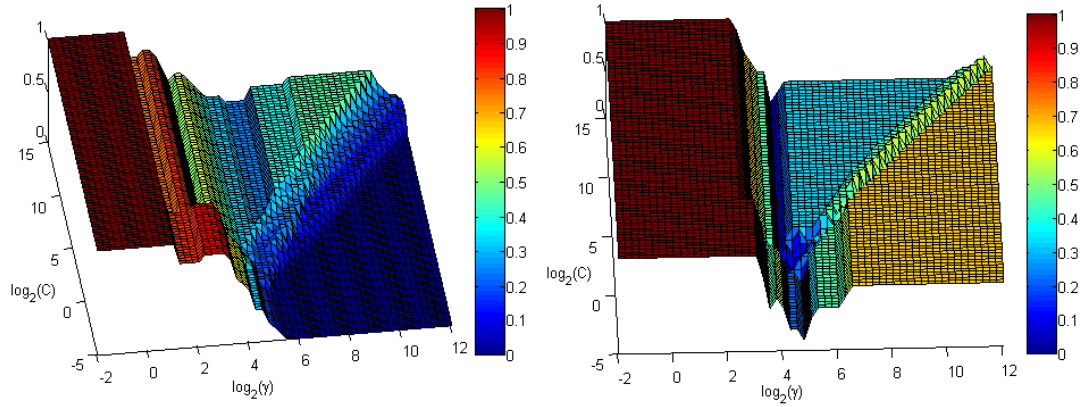


Figure 4.13: (1 - specificity) values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Right:(1 - specificity) values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 1 run.

value may change for different SVM parameters because limited number of regions are used for test. When more than one run is made in next parts this situation is partially eliminated. Performance point plots for the classifications of regions with dimension 128 x 128 are given here and in the rest of this thesis since there are more mass regions of size 128 x 128 than 256 x 256 so that it is more easy to comment on results.

One run is not enough to decide on, which of  $C$  and  $\gamma$  values will be used, for classification operation since there are not enough regions with mass, as mentioned before. So that 100 runs are made. Different training and test sets are chosen for each run. Mean and standard deviation of performance metrics are calculated. Means of  $A_z$ , truth rate, sensitivity and (1 - specificity) with respect to  $C$  and  $\gamma$  parameters are given in Figures 4.15, 4.16, 4.17, 4.18 respectively. Maximum standard deviation is around 0,1 level for each metric. This is reasonable so that detailed standard deviation maps are not given.

Maximum truth rate and  $A_z$  values for different region sizes are given in Table 4.5. Performance of classifier for 128 x 128 regions seems to be better since training is made better. However, truth rate and  $A_z$  values are worse than expected results given in [49]. It is thought that this is mainly due to training and selection of Gabor filter parameters. In [49] training and test are made with 512 regions with mass or without mass, so that a better training is made since number of training and test regions affect SVM classifier performance. Moreover, optimum Gabor filter bank parameter calculation is not made in this thesis. Gabor filter bank parameters, which are used, in [49] are set to extract features. However, these optimum parameters may change from database to database and from region size to region size. In future, these deficiencies may be resolved and a better performance could be obtained in terms of maximum truth rate and maximum  $A_z$  performance.

SVM performance points, which are obtained after 100 runs, are given in Figure 4.19.

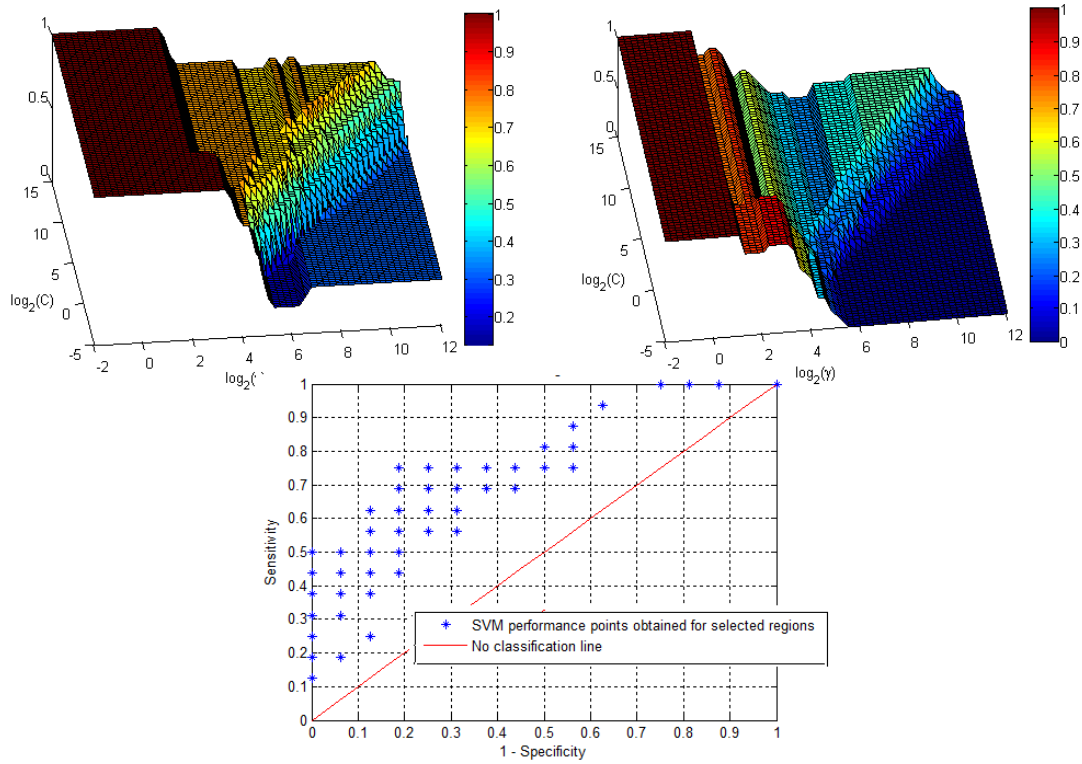


Figure 4.14: Top Left:Sensitivity values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Top Right:(1 - specificity) values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 1 run, Bottom:Performance points for different SVM parameters.

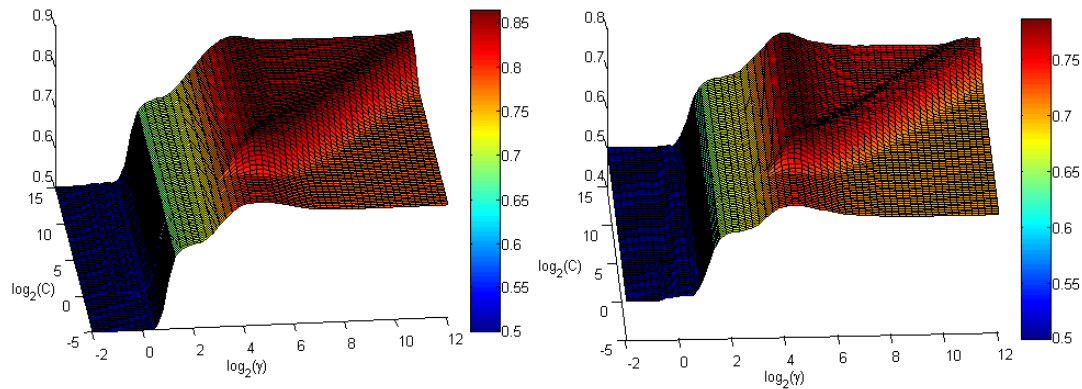


Figure 4.15: Left:Mean Az values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 100 runs, Right:Mean Az values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 100 runs.

It seems to be a ROC curve but is not exactly a ROC curve since for some values of (1 - specificity) and sensitivity no data exists. It is thought that if there were enough mass region, SVM parameters' resolution was better and SVM parameters' search ranges were wider this performance points would constitute a curve that is very close to a



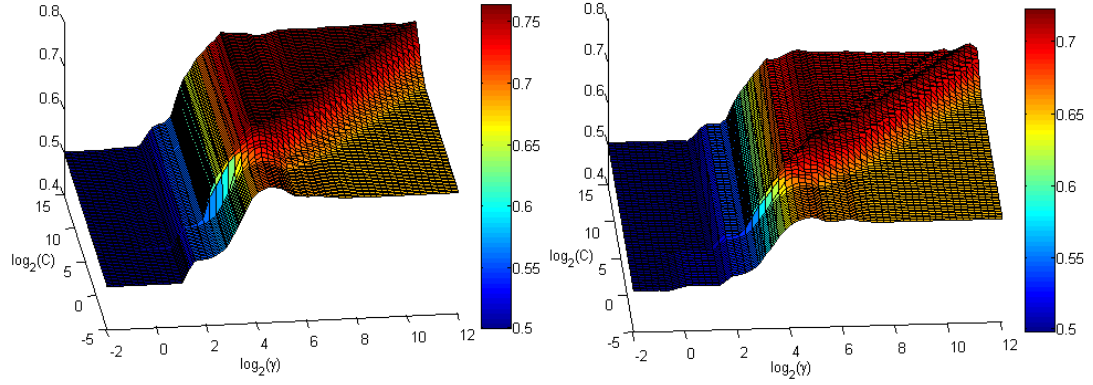


Figure 4.16: Left:Mean truth rate values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 100 runs, Right:Mean truth rate values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 100 runs.

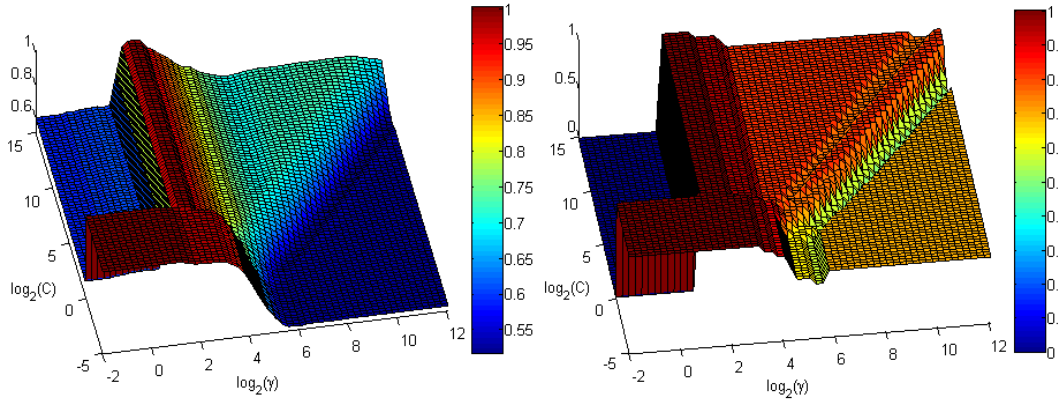


Figure 4.17: Left:Mean sensitivity values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after 100 runs, Right:Mean sensitivity values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after 100 runs.

ROC curve.

It is observed that there are performance points on "No classification line", meaning that classification is not performed for some SVM parameter couples. These parameter couples may be seen on Figures 4.15 and 4.16. SVM parameter couples with  $A_z = 0,5$  and truth rate = 0,5 are corresponded with points on "No classification line".

ROC curve determination with addition of margin to default hyper-plane was explained in the previous sections for SVM classifier. Therefore, average ROC curve of classification could be found for one  $C$  and  $\gamma$  couple from all runs. For instance, ROC curves with chosen  $C$  and  $\gamma$  parameters, satisfying maximum truth rate condition for 128 x 128 and 256 x 256 regions, are given in Figure 4.20.

It is important to note that user may choose operating point in two ways. Firstly,

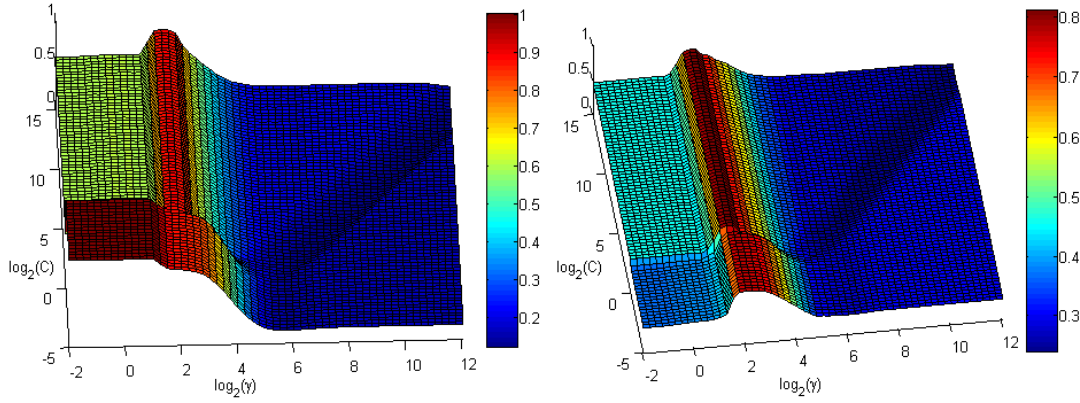


Figure 4.18: Mean (1 - specificity) values obtained with respect to SVM parameters for classification of selected regions of dimension 128 x 128 after one run, Right: Mean (1 - specificity) values obtained with respect to SVM parameters for classification of selected regions of dimension 256 x 256 after one run.

Table 4.5: Maximum truth rate and Az values for regions with dimensions.

Region size	Maximum truth rate	Maximum Az
128 x 128	0,76	0,86
256 x 256	0,72	0,78

$C$  and  $\gamma$  can be chosen satisfying desired sensitivity and (1 - specificity) condition on performance plot and next classifications can be made without adding offset to default hyper-plane obtained with these SVM parameters. Secondly,  $C$   $\gamma$  values can be chosen satisfying maximum truth rate or Az criteria and margin offset can be added on default hyper-plane to operate on the desired point of ROC curve.

#### 4.2.2 Classification performance comparison of regions with different dimensions including same masses

In this thesis, regions with different sizes were extracted from mammograms compatible with support region of Iris filter. 128 x 128 regions were extracted for "Iris filter 1", "Iris filter 2" and 256 x 256 regions were extracted for "Iris filter 3", "Iris filter 4", "Iris filter 5". Reason of this choice is explained in this section.

Selection of mass regions are explained in Section 4.2.1. If the radius of mass smaller than 55 pixels a region of size 128 x 128 is extracted for classification step. Classification performance comparison of these regions are made with selected regions of dimension 256 x 256, including the same masses. This time classifications are made with same number of training and test regions. Number of mass and non-mass regions are equal for both training and test sets. Mammogram numbers of selected mass and non-mass regions are also kept equal. In other words, only dimensions of regions are different. 50 runs are made. Mean Az and mean truth rate performance with respect

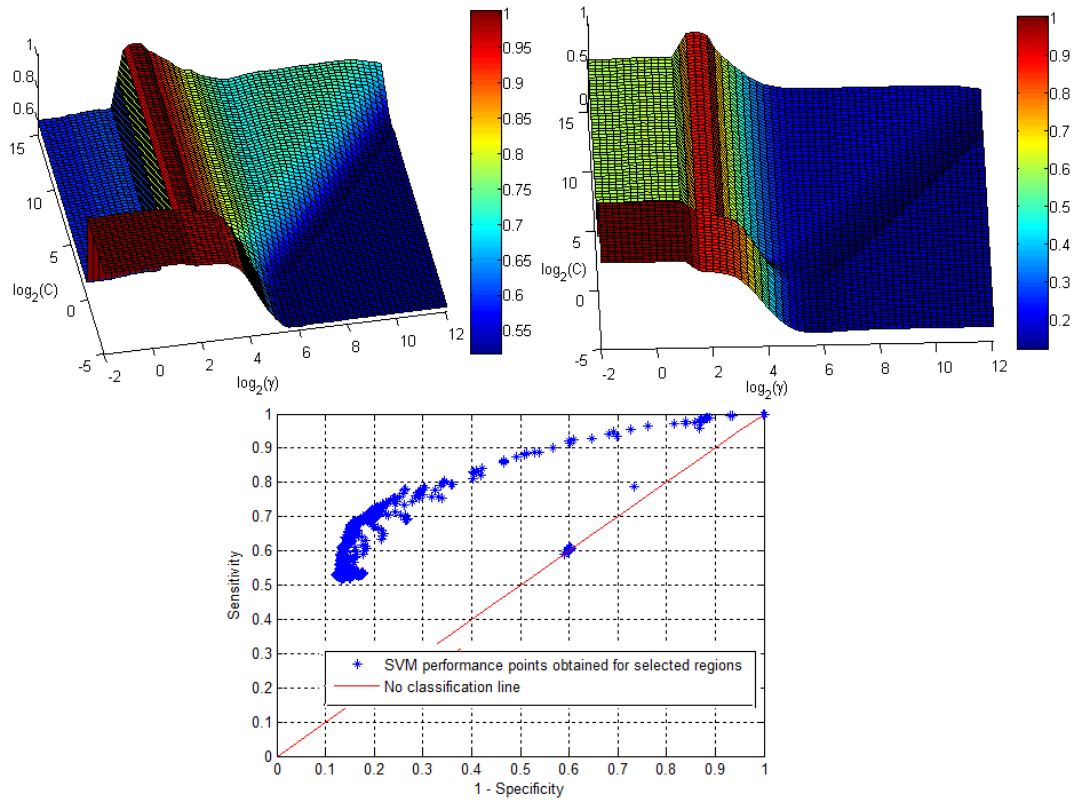


Figure 4.19: Top Left:Mean Sensitivity values obtained with respect to SVM parameters for selected regions of dimension 128 x 128 after 100 runs, Top Right:Mean (1 - specificity) values obtained with respect to SVM parameters for selected regions of dimension 128 x 128 after 100 runs, Bottom:Performance points for different SVM parameters after 100 runs.

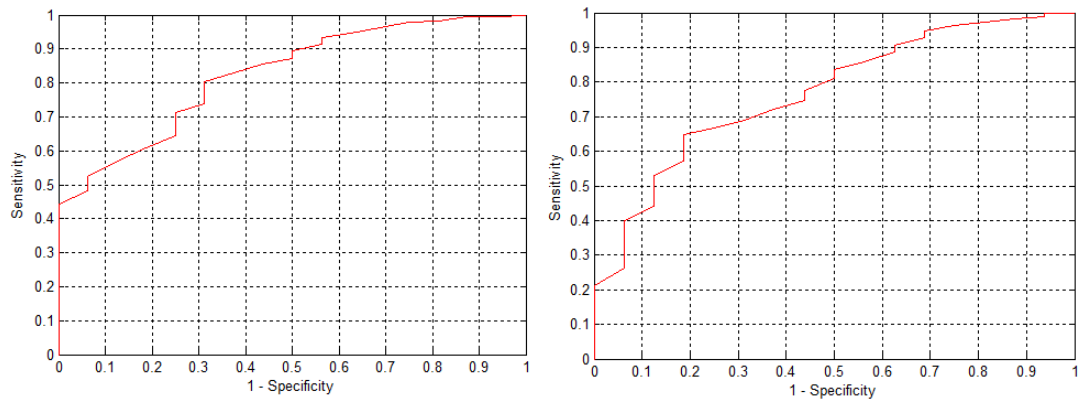


Figure 4.20: Left:ROC curve obtained with SVM parameters satisfying maximum truth rate condition for selected regions of dimension 128 x 128 after 100 runs, Right:ROC curve obtained with SVM parameters satisfying maximum truth rate condition for selected regions of dimension 256 x 256 after 100 runs.

to SVM parameters are determined. ROC curves are obtained with SVM parameters, satisfying maximum Az and maximum truth rate conditions for classification of regions of different sizes, are given in Figure 4.21.

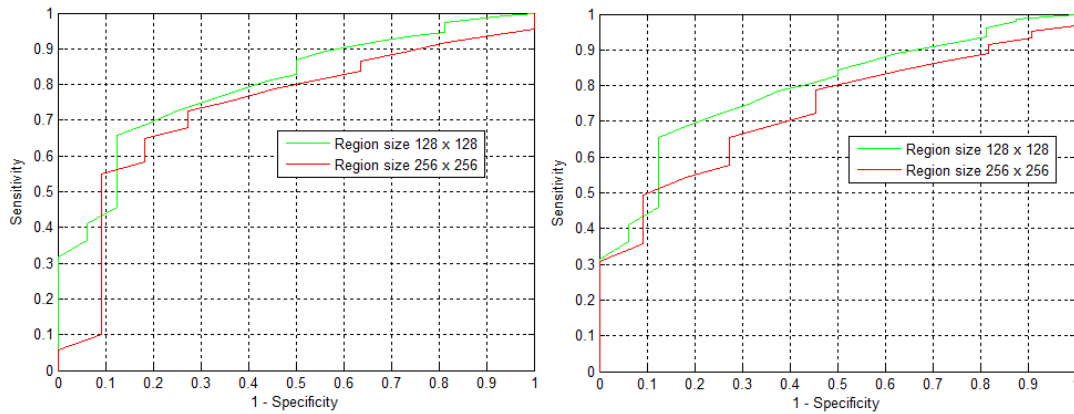


Figure 4.21: Left:ROC curve comparison of different selected region sizes with SVM parameters satisfying maximum Az condition, Right:ROC curve comparison of different selected region sizes with SVM parameters satisfying maximum truth rate condition.

Az performance is better for 128 x 128 regions than 256 x 256 ones in terms of both maximum Az and maximum truth rate conditions. Therefore, it is reasonable to choose extracted region sizes in parallel to mass sizes. Iris filter algorithm is very useful since size of detected mass is compatible with the applied Iris filter's region of support. In other words, size of regions, which are extracted, can be optimized using Iris filter's region of support information. In this thesis, only two dimensions are concerned due to limited number of mammograms including masses. However, number of dimensions could be increased if there were enough mammograms with masses and success of algorithm may be raised with this increment. For example, 5 dimensions of regions, matched with 5 Iris filters' output, may be extracted, etc.

### 4.2.3 Classification of suspicious regions

Classification of suspicious regions is performed with chosen SVM parameters. Overall block diagram for classification of suspicious regions is given in Figure 4.22. It should be noted that training set contains selected regions.

Suspicious regions, in one mammogram, are classified. Therefore, all selected regions of other mammograms may be in the training set since they are independent from test regions. Hence, number of training set could be increased. Number of regions in training set is given in Table 4.6. Equal number of non-mass regions to mass regions are chosen. There are more selected non-mass regions than selected mass region. As a result of this fact 50 runs are made for test regions. Different non-mass region sets are chosen for each run. Mean of all classifications determines overall performance of classifier.

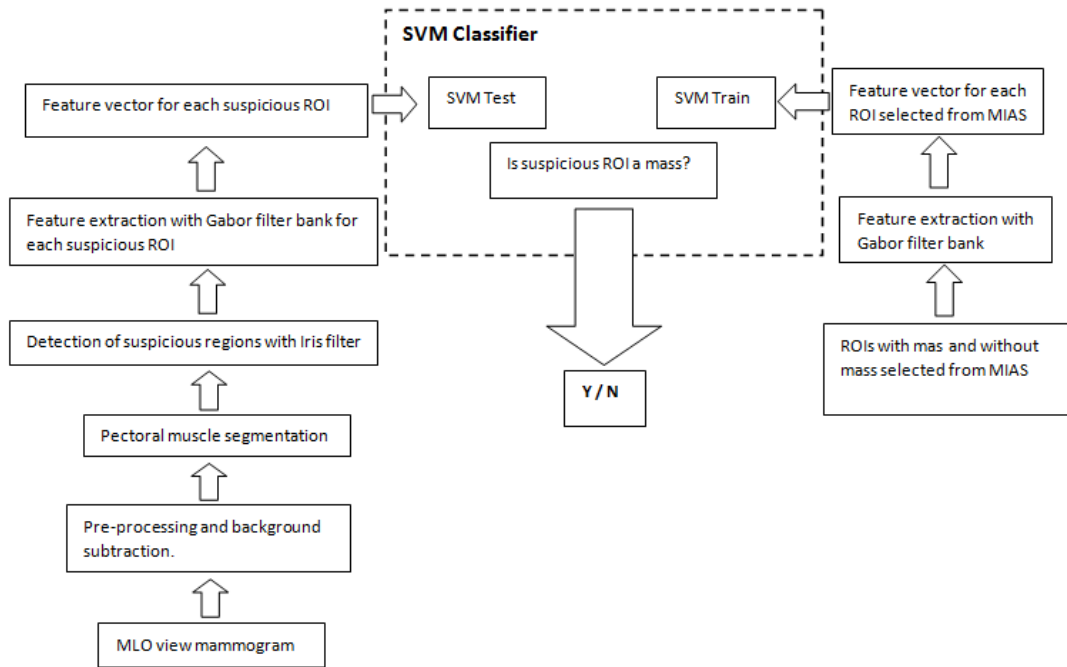


Figure 4.22: SVM implementation for classification of suspicious ROIs, training set is obtained from selected ROIs.

Table4.6: Number of regions in training set.

Region size	number of training regions with mass	number of training regions without mass
128 x 128	33	33
256 x 256	15	15



## CHAPTER 5

### RESULTS

Suspicious region classification results for chosen different SVM parameters are given in this section. Firstly, two SVM parameter couple are chosen according to mean performance metrics, suspicious regions' classifications are done with these parameters and classification results with chosen parameters are explored. Chosen  $C$ ,  $\gamma$  parameters and mean sensitivity, (1 - specificity), which are obtained in SVM determination step for these parameters, are given in Table 5.1. Parameter set 1 is chosen to keep sensitivity high and parameter set 2 is to keep specificity high.

Table5.1: Chosen SVM parameters.

Parameter sets	$\log_2(C)$	$\log_2(\gamma)$	Sensitivity	(1 - Specificity)
Parameter set 1 for 128 x 128 regions	0,25	3,00	0,85	0, 57
Parameter set 2 for 128 x 128 regions	11,50	10,25	0,62	0,19
Parameter set 1 for 256 x 256 regions	0,75	4,50	0,87	0,31
Parameter set 2 for 256 x 256 regions	6,75	7,25	0,62	0,24

Mean FPpI obtained after classification, which is implemented with chosen parameters for each Iris filter output, is given in Table 5.2. It is observed that FPs coming from Iris filters are reduced with classification step when Table 5.2 is compared with Table 3.2. FPpI reduction is more for parameter set 2 than parameter set 1 as expected.

(1 - specificity) values are greater than expected values since the ratio, number of test regions without mass over number of regions with mass, in test set is very larger than the ratio, regions without mass over regions with mass, in the training sets.

Results of mean FPpI can also be shown with respect to dimensions of regions as given in Table 5.3.

Statistic of true mass regions, which are detected after classification, is given in Table 5.4. Sensitivity is defined as the number of detected masses over the number of mam-

Table5.2: Mean FpPI after classification with respect to Iris filters.

Iris filters	Parameter set 1	Parameter set 2	(1 - Specificity) for set 1	(1 - Specificity) for set 2
Iris filter 1	11,62	8,16	0,40	0,29
Iris filter 2	8,14	6,32	0,64	0,50
Iris filter 3	4,36	4,28	0,79	0,78
Iris filter 4	4,15	4,06	0,89	0,87
Iris filter 5	4,35	4,13	0,90	0,86

Table5.3: Mean FpPI after classification with respect to region sizes.

Region size	(1 - Specificity) for set 1	(1 - Specificity) for set 2
128 x 128 regions	0,48	0,35
256 x 256 regions	0,86	0,84

mograms with mass. When compared with Table 5.1, sensitivity results are very close to, or better than expected values because training, which is applied to test suspicious regions, is better than training in SVM parameter determination step. In other words, more regions with mass and without mass are used to train SVM in suspicious regions' classification step (Table 4.6 and Table 4.4). Nearly 50 percent of the FPs are eliminated with a cost of missing 9 percent of true masses if parameter set 1 is chosen for 128 x 128 regions.

Table5.4: True mass regions detected after classification.

Parameter set	Region size	Mammograms with mass	Detected masses	Sensitivity
Parameter set 1	128 x 128	34	31	0,91
Parameter set 1	256 x 256	16	15	0,94
Parameter set 2	128 x 128	34	27	0,80
Parameter set 2	256 x 256	16	13	0,81

Secondly, this procedure is repeated for a wide range of SVM parameters and suspicious regions' classification results are explored for all of these SVM parameter couples. Sensitivity and (1 - specificity) performance of classification results with different SVM parameter couples are given in Figures 5.1 and 5.2 respectively.

Variation of sensitivity around 0,5 value is observed in SVM parameter region bounded with  $5 < \log_2(C) < 15$  and  $-2 < \log_2(\gamma) < 2$ . Looking at the mean Az and truth rate values on Figures 4.15 and 4.16 for this SVM parameter region, it is seen that Az and truth rate values are very close to 0,5 which means that no classification is performed. Moreover, limited number of suspicious mass regions, causing quantization as mentioned earlier, exist for test step. Therefore, variations are expected to occur



around 0,5 value for this region. Variation for regions of dimension 256 x 256 is greater than variation for regions of dimension 128 x 128 since there are less suspicious mas region to test. It is not offered to select SVM parameters in this region in classification as previously stated. There is variation for other regions also in terms of sensitivity performance, but it is acceptable since variation is small.

Variation of (1 - specificity) in SVM parameter region (bounded with  $5 < \log_2(C) < 15$  and  $-2 < \log_2(\gamma) < 2$ ) is smaller when compared with sensitivity because there are more non-mass suspicious regions to test than mass regions. Therefore, quantization does not effect performance so much in the same SVM parameter region.

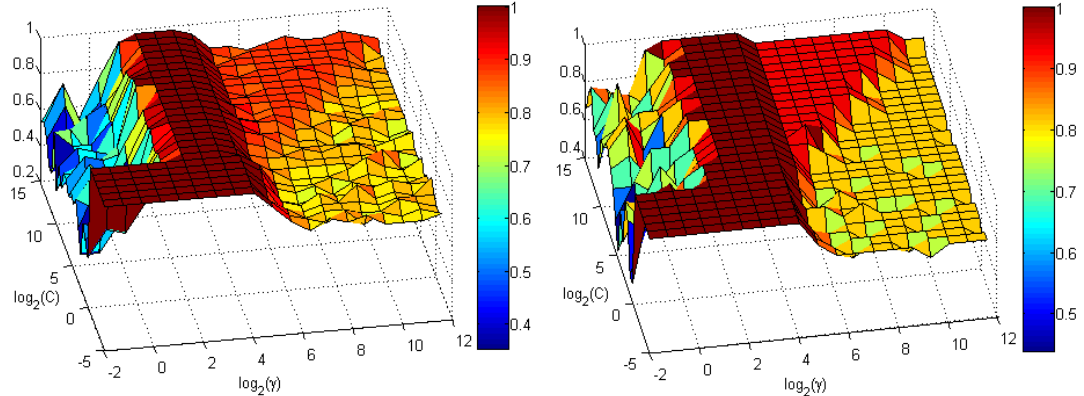


Figure 5.1: Left:Sensitivity values obtained with respect to SVM parameters after classification of suspicious regions with dimension 128 x 128, Right:Sensitivity values obtained with respect to SVM parameters after classification of suspicious regions with dimension 256 x 256.

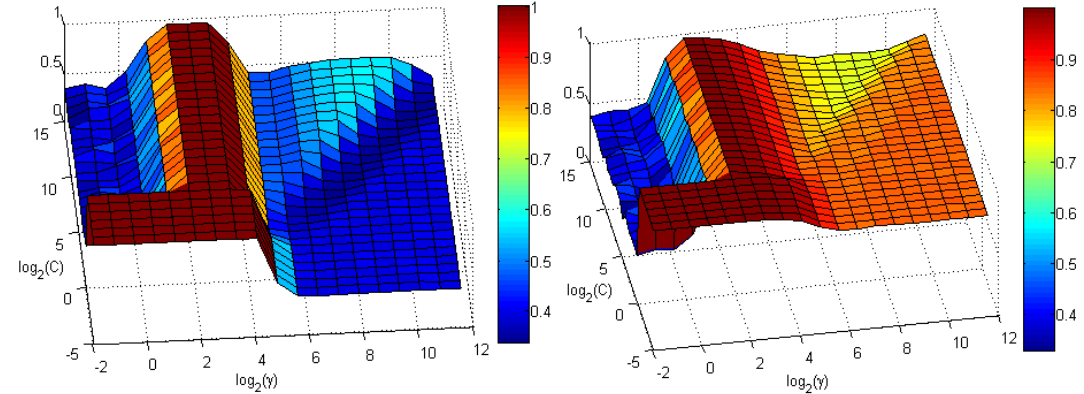


Figure 5.2: Left:(1 - specificity) values obtained with respect to SVM parameters after classification of suspicious regions with dimension 128 x 128, Right:(1 - specificity) values obtained with respect to SVM parameters after classification of suspicious regions with dimension 256 x 256.

Performance points obtained after classification of suspicious regions, with size 128 x 128, are shown in Figure 5.3 with red points. Variation in sensitivity and (1-specificity) values, mentioned previously, is also seen for the performance points obtained after

classification of the suspicious regions. In other words, there are performance points deviated from SVM performance points which are obtained for selected regions. Reasons of this circumstance explained above.

In addition to variation of sensitivity and (1-specificity) values it is noticed that minimum (1-specificity) values are greater than minimum (1-specificity) of SVM performance points in SVM parameter determination step. This is due to the fact that the ratio, suspicious regions' number with mass over suspicious regions' number without mass, is very smaller than the ratio, selected mass regions' number over selected non-mass regions' number in test steps of SVM parameter determination step. Mass over non-mass ratio is kept 0,5 for all training and test steps in SVM parameter determination process. However, training and test ratios are not equal for the classification of suspicious regions process. Classifier tends to produce more FPs since less information is trained about non-mass regions in spite of the fact that there are more non-mass suspicious regions to test.

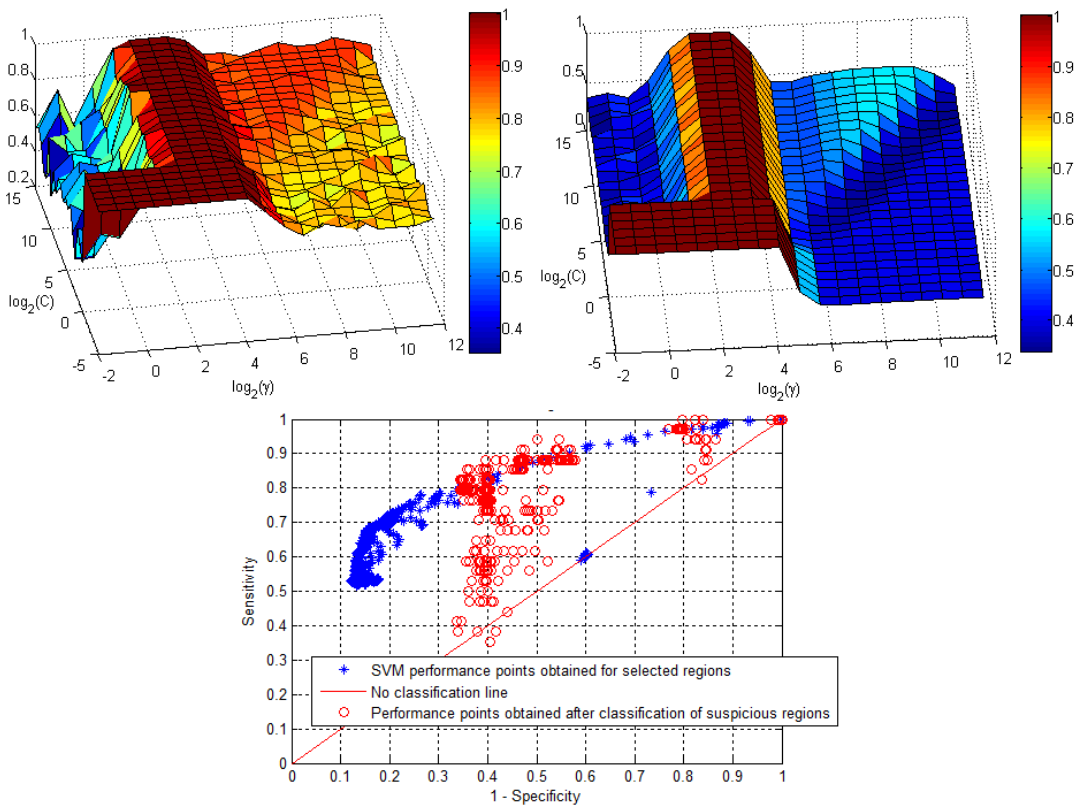


Figure 5.3: Top Left:Sensitivity values obtained with respect to SVM parameters after classification of suspicious regions with dimension 128 x 128, Top Right:(1 - specificity) values obtained with respect to SVM parameters for suspicious regions of dimension 128 x 128, Bottom:Performance points for different SVM parameters.

## CHAPTER 6

### CONCLUSION

To summarize, we have proposed a method for classification of mass regions in MLO view mammograms. We have determined the suspicious regions by Iris filters, with variable support regions, applied on the breast region without pectoral muscle. Suspicious regions are extracted compatible with the applied filters. Classification has been applied to the textural features are obtained using Gabor filters that are applied on these suspicious regions. False detection ratio has been reduced nearly 50 percent with a cost of missing 9 percent of true mass regions. In addition, we have proposed a novel algorithm, which is based on average derivative calculation and line fitting with least square solution, for pectoral muscle region determination. Our algorithm outperforms other algorithms given in the literature in terms of FP (False positive) pixel percentage and FN (False negative) pixel percentage metrics.

Novel pectoral muscle segmentation algorithm has been proposed based on two steps. Firstly, average derivative calculation applied for each raw of mammogram on a determined region of interest, minimum derivative points have been calculated and upper part of pectoral boundary has been obtained correctly after filtering operations. Secondly, line pieces, estimated with least square solution technique, have been fitted on calculated minimum derivative points for the lower part of mammogram. Obtained pectoral boundaries are evaluated with FP pixel percentage and FN pixel percentage metrics and compared with the performance of the previous algorithms. Average FP and FN pixel percentages are calculated as 1,30 and 2,57 respectively. Performance of proposed algorithm is better than the given algorithms in the literature. However, proposed algorithm is a man in the loop algorithm and needs someone to choose optimum line piece size to find the best pectoral boundary. In the future, selection of best pectoral boundary among boundaries that are obtained with different line pieces, can be implemented automatically.

Suspicious ROI detection has been performed with Iris filter algorithm. Iris filter is a gradient direction based filter. It produces high values for a pixel of interest when gradients in the support region of filter directed towards this pixel. Especially, when mass edge falls in the region of support filter it produces high output for pixel of interests inside mass region. Masses exist in different dimensions so that Iris filters

with different support regions should be used in order not to miss any mass. 5 different Iris filters have been applied covering all possible mass sizes. Simple threshold has been applied on each Iris filter and possible mass locations have been determined. Regions of dimension 128 x 128 have been extracted for outputs of "Iris filter 1", "Iris filter 2" and regions of 256 x 256 are extracted for outputs of "Iris filter 3", "Iris filter 4" and "Iris filter 5". Adaptive region size determination, compatible with the applied Iris filter, has increased classification performance.

All masses have been detected with Iris filters except one spiculated malignant mass. Gradients, in the support regions of all Iris filters are not directed towards mass region of this spiculated mass. Therefore, high values have not been produced for pixels of interest inside the mass region. After threshold no alarm has been given for this region. In short, Iris filter performance has not been good for this kind of situations. Otherwise, detection ratio is satisfactory.

Each Iris filter has produced different average FPPi. "Iris filter 1" has produced highest FPPi rate although threshold level applied is the highest one. In this thesis, threshold levels, for each filter, has been determined experimentally and FPPi performance has not been concerned since the main goal is to test classifier's success of FP elimination. Detection performance for different threshold values can be explored and an optimum threshold level may be found to obtain a better FPPi performance in mass detection step in the future.

Textural features have been extracted from sub-regions in feature extraction step. Gabor filter bank has been applied on sub-windows of suspicious regions. Mean, standard deviation and skewness for each filter and sub-region have been calculated so that a feature vector is obtained for each suspicious region. Gabor filter bank parameter values have been selected as the same values given in [49]. In this thesis, optimum Gabor filter bank parameter determination is not performed. In future, classification performance of features obtained for Gabor filter bank with different parameters may be searched and optimum parameters can be chosen to obtain a better classification performance.

Classification is implemented with SVM classifier that is used widely in the literature for binary classification problems. SVM classifier, with rbf kernel, has been used. SVM performance has been controlled with two parameters: regularization parameter ( $C$ ) and width of the kernel function ( $\gamma$ ). SVM performance, with different parameters, must be known in order to decide on which parameters will be chosen for classification of suspicious regions. Hence SVM classification performance has been explored with the classifications of selected mammogram regions from mini-MIAS database with respect to different SVM parameters. 100 runs are made due to lack of mass regions in the database. Different sets have been used for each run. Mass and non-mass regions' numbers have been kept equal for training and test sets of each run. Average of performance metrics such as Az, truth rate, sensitivity and (1 - specificity)) with

respect to SVM parameters have been determined. Maximum truth rate of 0,7634 has been obtained for 128 x 128 regions' classification and 0,7222 for 256 x 256 regions' classification. Maximum Az of 0,8633 is obtained for classification of 128 x 128 regions and 0,7854 for classification of 256 x 256 regions. Performance, in terms of Az and truth rate, is less than the performance given in [49]. It is thought that this is mainly due to the lack of selected mass regions in the database. Although 256 mass and 256 non-mass selected regions are used in [49], total 50 mass (128 x 128 and 256 x 256 regions) selected regions have been used in this thesis. Performance of classification for 128 x 128 regions is better than performance of classification for 256 x 256 regions since more mass and non-mass regions have been used in training steps. Moreover, performance points, obtained from mean sensitivity and (1 - specificity), have been shown on a plot on which user may choose SVM parameters to operate.

Suspicious regions, which are determined by Iris filters, have been classified with respect to different SVM parameters and classification performance has been explored for each case. Training of SVM has been done with selected mass and non-mass regions. FpPI ratio, obtained after mass detection step, has been reduced 50 percent with a cost of missing 9 percent of the true mass regions.

It is observed that for some SVM parameter regions, very close performance points to the performance points, which are determined in SVM parameter determination step, are obtained in terms of sensitivity and (1 - specificity) metrics. However, (1 - specificity) ratio is seemed to increase when compared with (1 - specificity) determined in SVM parameter determination step for the same SVM parameter couple. This is because of the fact that SVM has been trained and tested with equal number of selected mass and non-mass regions in SVM parameter determination step, however test set of one mammogram contains many suspicious regions without mass than suspicious regions with mass and training set contains equal number of selected mass and non-mass regions in suspicious regions' classification step.



## REFERENCES

- [1] Guidubaldo Querci della Rovere, Ruth Warren, and John R Benson. *Early breast cancer: from screening to multidisciplinary management*. Taylor & Francis, 2005.
- [2] Xavier Lladó, Arnau Oliver, Jordi Freixenet, Robert Martí, and Joan Martí. A textural approach for mass false positive reduction in mammography. *Computerized Medical Imaging and Graphics*, 33(6):415–422, 2009.
- [3] William E Erkonen and Wilbur L Smith. Radiología 101 las bases y fundamentos.
- [4] Michael Heath, Kevin Bowyer, Daniel Kopans, P Kegelmeyer Jr, Richard Moore, Kyong Chang, and S Munishkumaran. Current status of the digital database for screening mammography. In *Digital mammography*, pages 457–460. Springer, 1998.
- [5] The mini-MIAS database of mammograms. <http://peipa.essex.ac.uk/info/mias.html>, June 1994.
- [6] Rita Filipa dos Santos Teixeira. Computer analysis of mammography images to aid diagnosis. 2012.
- [7] HD CHENG, XJ SHI, R MIN, LM HU, XP CAI, and HN DU. Approaches for automated detection and classification of masses in mammograms. *Pattern recognition*, 39(4):646–668, 2006.
- [8] Xiangqian Jiang, Shan Lou, and Paul J Scott. Morphological method for surface metrology and dimensional metrology based on the alpha shape. *Measurement science and technology*, 23(1):015003, 2012.
- [9] MATLAB. *version 7.10.0 (R2010a)*. The MathWorks Inc., Natick, Massachusetts, 2010.
- [10] LAG Ries, D Melbert, M Krapcho, DG Stinchcomb, N Howlader, MJ Horner, A Mariotto, BA Miller, EJ Feuer, SF Altekruse, et al. Seer cancer statistics review, 1975-2005.
- [11] Arnau Oliver, Jordi Freixenet, Joan Marti, Elsa Pérez, Josep Pont, Erika RE Denton, and Reyer Zwiggelhaar. A review of automatic mass detection and segmentation in mammographic images. *Medical Image Analysis*, 14(2):87–110, 2010.
- [12] Jacques Estève, A Kricke, Jacques Ferlay, and D Maxwell Parkin. Facts and figures of cancer in the european community. 1993.
- [13] Mohamed Meselhy Eltoukhy, Ibrahima Faye, and Brahim Belhaouari Samir. Breast cancer diagnosis in digital mammogram using multiscale curvelet transform. *Computerized Medical Imaging and Graphics*, 34(4):269–276, 2010.

- [14] Radhika Sivaramakrishna and Richard Gordon. Detection of breast cancer at a smaller size can reduce the likelihood of metastatic spread: a quantitative analysis. *Academic radiology*, 4(1):8–12, 1997.
- [15] Catherine N Chinyama. *Benign breast diseases: radiology-pathology-risk assessment*. Springer Verlag, 2004.
- [16] Daniel B Kopans. *Breast imaging*. Wolters Kluwer Health, 2007.
- [17] Saskia van Engeland, Peter Snoeren, JHCL Hendriks, and Nico Karssemeijer. A comparison of methods for mammogram registration. *Medical Imaging, IEEE Transactions on*, 22(11):1436–1444, 2003.
- [18] K Kavitha and N Kumaravel. A comparative study of various microcalcification cluster detection methods in digitized mammograms. In *Systems, Signals and Image Processing, 2007 and 6th EURASIP Conference focused on Speech and Image Processing, Multimedia Communications and Services. 14th International Workshop on*, pages 405–409. IEEE, 2007.
- [19] RF Chang, SF Huang, LP Wang, DR Chen, and WK Moon. Microcalcification detection in 3-d breast ultrasound. In *Engineering in Medicine and Biology Society, 2005. IEEE-EMBS 2005. 27th Annual International Conference of the*, pages 6297–6300. IEEE, 2006.
- [20] Fatemeh Saki, Amir Tahmasbi, and Shahriar B Shokouhi. A novel opposition-based classifier for mass diagnosis in mammography images. In *Biomedical Engineering (ICBME), 2010 17th Iranian Conference of*, pages 1–4. IEEE, 2010.
- [21] Li Sun, Lihua Li, Weidong Xu, Wei Liu, Juan Zhang, and Guoliang Shao. A novel classification scheme for breast masses based on multi-view information fusion. In *Bioinformatics and Biomedical Engineering (iCBBE), 2010 4th International Conference on*, pages 1–4. IEEE, 2010.
- [22] Huanping Zhao, Weidong Xu, Lihua Li, and Juan Zhang. Classification of breast masses based on multi-view information fusion using multi-agent method. In *Bioinformatics and Biomedical Engineering, (iCBBE) 2011 5th International Conference on*, pages 1–4. IEEE, 2011.
- [23] Shantanu Banik, Rangaraj M Rangayyan, and JE Leo Desautels. Detection of architectural distortion in prior mammograms. *Medical Imaging, IEEE Transactions on*, 30(2):279–294, 2011.
- [24] Jelena Bozek, Emil Dumic, and Mislav Grgic. Bilateral asymmetry detection in digital mammography using b-spline interpolation. In *Systems, Signals and Image Processing, 2009. IWSSIP 2009. 16th International Conference on*, pages 1–4. IEEE, 2009.
- [25] Bao-Long Li, Xue-Qing Wang, and Zhi-Qing Fan. The bilateral information asymmetry on insurance market. In *Industrial Engineering and Engineering Management, 2009. IE&EM'09. 16th International Conference on*, pages 750–752. IEEE, 2009.



- [26] Claudio Marrocco, Mario Molinara, Ciro D’Elia, and Francesco Tortorella. A computer-aided detection system for clustered microcalcifications. *Artificial intelligence in medicine*, 50(1):23–32, 2010.
- [27] Jinshan Tang, Rangaraj M Rangayyan, Jun Xu, Issam El Naqa, and Yongyi Yang. Computer-aided detection and diagnosis of breast cancer with mammography: recent advances. *Information Technology in Biomedicine, IEEE Transactions on*, 13(2):236–251, 2009.
- [28] Bruno Boyer, Corinne Balleyguier, Olivier Granat, and Christian Pharaboz. Cad in questions/answers: review of the literature. *European journal of radiology*, 69(1):24–33, 2009.
- [29] Milan Sonka, J Michael Fitzpatrick, and Barry R Masters. Handbook of medical imaging, volume 2: Medical image processing and analysis. *Optics & Photonics News*, 13:50–51, 2002.
- [30] Arianna Mencattini, Giulia Rabottino, Marcello Salmeri, and Roberto Lojacono. Assessment of a breast mass identification procedure using an iris detector. *Instrumentation and Measurement, IEEE Transactions on*, 59(10):2505–2512, 2010.
- [31] Robert M Nishikawa. Mammographic databases. *Breast disease*, 10(3):137–150, 1998.
- [32] Jayasree Chakraborty, Sudipta Mukhopadhyay, Veenu Singla, Niranjana Khandelwal, and Pinakpani Bhattacharyya. Automatic detection of pectoral muscle using average gradient and shape based feature. *Journal of digital imaging*, 25(3):387–399, 2012.
- [33] Nico Karssemeijer. Automated classification of parenchymal patterns in mammograms. *Physics in medicine and biology*, 43(2):365, 1998.
- [34] Sérgio Koodi Kinoshita, Paulo M Azevedo-Marques, RR Pereira Jr, Jose Antônio Heisinger Rodrigues, and Rangaraj M Rangayyan. Radon-domain detection of the nipple and the pectoral muscle in mammograms. *Journal of Digital Imaging*, 21(1):37–49, 2008.
- [35] Margaret Yam, Michael Brady, Ralph Highnam, Christian Behrenbruch, Ruth English, and Yasuyo Kita. Three-dimensional reconstruction of microcalcification clusters from two mammographic views. *Medical Imaging, IEEE Transactions on*, 20(6):479–489, 2001.
- [36] Ricardo J Ferrari, Rangaraj M Rangayyan, JE Leo Desautels, RA Borges, and Annie France Frere. Automatic identification of the pectoral muscle in mammograms. *Medical Imaging, IEEE Transactions on*, 23(2):232–245, 2004.
- [37] Lei Wang, Miao-liang Zhu, Li-ping Deng, and Xin Yuan. Automatic pectoral muscle boundary detection in mammograms based on markov chain and active contour model. *Journal of Zhejiang University SCIENCE C*, 11(2):111–118, 2010.
- [38] Fei Ma, Mariusz Bajger, John P Slavotinek, and Murk J Bottema. Two graph theory based methods for identifying the pectoral muscle in mammograms. *Pattern Recognition*, 40(9):2592–2602, 2007.

- [39] Mariusz Bajger, Fei Ma, and Murk J Bottema. Minimum spanning trees and active contours for identification of the pectoral muscle in screening mammograms. In *Digital Image Computing: Techniques and Applications, 2005. DICTA'05. Proceedings 2005*, pages 47–47. IEEE, 2005.
- [40] Hidefumi Kobatake, Masayuki Murakami, Hideya Takeo, and Shigeru Nawano. Computerized detection of malignant tumors on digital mammograms. *Medical Imaging, IEEE Transactions on*, 18(5):369–378, 1999.
- [41] Toshihiko Terada, Yohei Fukumizu, Hironori Yamauchi, Hiroto Chou, and Yoshimasa Kurumi. Detecting mass and its region in mammograms using mean shift segmentation and iris filter. In *Communications and Information Technologies (ISCIT), 2010 International Symposium on*, pages 1176–1179. IEEE, 2010.
- [42] Arianna Mencattini, Giulia Rabottino, Marcello Salmeri, and Roberto Lojacono. An iris detector for tumoral masses identification in mammograms. In *Medical Measurements and Applications, 2009. MeMeA 2009. IEEE International Workshop on*, pages 215–218. IEEE, 2009.
- [43] Peter Kruizinga and Nikolay Petkov. Nonlinear operator for oriented texture. *Image Processing, IEEE Transactions on*, 8(10):1395–1407, 1999.
- [44] Berkman Sahiner, Heang-Ping Chan, Datong Wei, Nicholas Petrick, Mark A Helvie, Dorit D Adler, and Mitchell M Goodsitt. Image feature selection by a genetic algorithm: Application to classification of mass and normal breast tissue. *Medical Physics*, 23(10):1671–1684, 1996.
- [45] Wei Qian, Xuejun Sun, Dansheng Song, and Robert A Clark. Digital mammography: wavelet transform and kalman-filtering neural network in mass segmentation and detection. *Academic radiology*, 8(11):1074–1082, 2001.
- [46] Georgia D Tourassi, Nevine H Eltonsy, James H Graham, CE Floyd, and Adel S Elmaghraby. Feature and knowledge based analysis for reduction of false positives in the computerized detection of masses in screening mammography. In *Engineering in Medicine and Biology Society, 2005. IEEE-EMBS 2005. 27th Annual International Conference of the*, pages 6524–6527. IEEE, 2006.
- [47] Arnau Oliver, Xavier Lladó, Joan Martí, Robert Martí, and Jordi Freixenet. False positive reduction in breast mass detection using two-dimensional pca. In *Pattern Recognition and Image Analysis*, pages 154–161. Springer, 2007.
- [48] Yufeng Zheng. Breast cancer detection with gabor features from digital mammograms. *algorithms*, 3(1):44–62, 2010.
- [49] Muhammad Hussain, Salabat Khan, Ghulam Muhammad, Mohamed Berbar, and George Bebis. Mass detection in digital mammograms using gabor filter bank. In *Image Processing (IPR 2012), IET Conference on*, pages 1–5. IET, 2012.
- [50] Charles L Lawson and Richard J Hanson. *Solving least squares problems*, volume 161. SIAM, 1974.
- [51] F. Y. Shih. *Image Processing and Mathematical Morphology: Fundamentals and Applications*. CRC Press, 2012.

- [52] Hyunsup Yoon, Youngjoon Han, and Hernsoo Hahn. Image contrast enhancement based sub-histogram equalization technique without over-equalization noise. *World Academy of Science, Engineering and Technology*, 50:2009, 2009.
- [53] Hidefumi Kobatake and Shigeru Hashimoto. Convergence index filter for vector fields. *Image Processing, IEEE Transactions on*, 8(8):1029–1038, 1999.
- [54] Yu Su, Shiguang Shan, Xilin Chen, and Wen Gao. Hierarchical ensemble of global and local classifiers for face recognition. *Image Processing, IEEE Transactions on*, 18(8):1885–1896, 2009.
- [55] Zehang Sun, George Bebis, and Ronald Miller. Monocular precrash vehicle detection: features and classifiers. *Image Processing, IEEE Transactions on*, 15(7):2019–2034, 2006.
- [56] M Hussain and Naveed Khan. Automatic mass detection in mammograms using multiscale spatial weber local descriptor. In *Systems, Signals and Image Processing (IWSSIP), 2012 19th International Conference on*, pages 288–291. IEEE, 2012.
- [57] Bernhard E Boser, Isabelle M Guyon, and Vladimir N Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–152. ACM, 1992.