

# Automated feature-based machine learning classification models for breast lesions

Areeq I. Hasan, Sarina M. Hasan

**Abstract**—We report the effectiveness of the apparent lesion-size discrepancy determined by semi-automatic segmentation algorithms between the processed ultrasound B-Mode and elasticity image of *in vivo* patient data calculated via a novel two-step adaptive-stretching strain estimation algorithm as a classification criterion between benign fibroadenomatoid and adenocarcinomatoid breast tissue. Furthermore, we find both support-vector machine classifiers accepting RF echosignal sequences and trained on the basis of lesion-size discrepancy as determined by aforementioned methods assuming an identical I/O schema while trained on the basis of lesion-size discrepancy, intensity contrast discrepancy, intensity contrast in the elasticity image, and intensity contrast in the B-Mode image as determined by all four processing-segmentation algorithm permutations to be effective classification models by which to differentiate between the RF echosignals of fibroadenoma and adenocarcinoma.

**Index Terms**—lesion-size discrepancy, adaptive-stretching strain estimation, semi-automatic segmentation, support-vector machine classifier,

## I. INTRODUCTION

THE rate of breast cancer vulnerability in women typically increases until they are about 74 years old, and then begins to steadily decline [1]. As of 2017, there were 250,520 cases reported for females with breast cancer[1]. While 128.5 out of 100,000 new women are diagnosed with this cancer each year, 20.3 of these women die[2]. An estimated 12.9% of females are diagnosed with breast cancer at least once in their life; there were about 3,577,264 women diagnosed and living with breast cancer in 2017 in the United States[2]. Stage 1 breast cancer means that the cancer cells that have emerged, reside in the same area in which they emerged[3]. If breast cancer is detected in this early state, the survival rate is 100% because it is still confined[3]. On the other hand, stage 4 of breast cancer has much more perilous consequences[4]. Stage 4 breast cancer is when the cancer has spread to other parts of your body, such as your brain and lungs[4]. At this stage, the cancer is highly invasive; while women can get treatments to live longer, this stage of cancer is incurable[4]. The work we have completed is geared towards malignant lesions early on, thereby having the ability to save lives. Being that so many women in the United States suffer from this genetic disease, it is imperative to develop methods to diagnose this illness before it becomes carcinogenic. Breast cancer continues to pose a prevalent concern for many females, affecting their health and relationships, and it is a prevalent issue that must be addressed.

In this paper, we discuss our developed method for the classification of breast lesions. Benign lesions are typically not harmful for ones health. However, it is important to detect

malignant tumors, as they can be dangerous for ones well-being. We developed a semi-automatic method to differentiate between benign and malignant tumors. Our algorithm evaluated the size of a lesion from its B-mode and strain images. The size differences we calculated in the strain and B-mode images of breast lesions led to accurate classification performances. Garra et al (1997) first described that invasive ductal carcinoma (IDCA) appears significantly larger in strain images than B-mode images [5].

Garra et al (1997) explored the size of numerous breast lesions and how they can lead to the analysis of its classification in regards to being a benign or malignant mass[5]. Following the biopsies of the lesions, information about the sizes of the lesions in their elastograms and sonograms were documented and compared to the biopsies[5]. The results in their paper showed that softer regions on elastograms appeared more luminous, whereas firmer tissues, including masses, appeared to be darkened[5]. Along with being darker on imaging, the cancers were also substantially larger on elastograms over sonograms[5]. They were able to conclude that elastography can help classify masses through means of size and contrast[5]. Their work established that size discrepancies play an accurate role in diagnosing lesions[5]. Moreover, this paper also established that a mass's lightness or darkness can contribute to its classification which is what led us to explore the idea of contrast briefly in our work[5]. Garra et al's information was useful for developing our algorithm and served as the backbone for the validity of our reasoning. Barr et al (2012) performed a larger study (528 female patients with 635 lesions) from female patients[6]. It commented on many of the findings from Garra et al (1997), including how size discrepancy is related to the classification of tumors. This work supports the earlier findings with a much larger number of patients; these criteria have been used by physicians in clinics all over the world. In this study, Elasticity and B-mode imaging were provided to the researchers after the participants had received elastograms[6]. A ratio for lesion sizes on elasticity imaging and B-mode was established as 1.0[6]. If the ratio was higher than 1.0, the lesion was considered malignant[6]. If the ratio was lower than 1.0, the lesion was considered to be benign[6]. Their findings were then proven with actual biopsies of each lesion to support their hypothesis. Of the 635 total lesions that were imaged and biopsied, there were 222 malignant lesions and 413 benign ones[6]. 219 out of the 222 malignant lesions actually had a ratio that was at least 1.0[6]. 361 out of the 413 benign lesions turned out to have which was less than 1.0[6]. Through these findings, his study concluded that there was a relationship between lesion-size

ratios and classification. Both works used expert radiologists for reading and classification. The performance is highly dependent on clinician expertise. We developed a computer algorithm that would be objective and operator independent. It should be able to help a less experienced clinician with accurately diagnosing breast lesions. Our algorithm was based on these prior work and differentiated between benign and malignant tumors reliably based on its lesion-size ratio.

First, elastograms were computed using a novel two-step adaptive-stretching strain estimation algorithm. After, we used the Hilbert transform to compute our B-mode images. Using GrowCut algorithms the lesions were marked off. We then estimated the size discrepancies of the lesions and computed an intensity contrast of B-mode and strain images.

This algorithm can be useful to future clinicians by helping them determine the classification of lesions that they face. While a lot of the research community has been focused on inductive research, we created an algorithm that specifically focuses on analyzing data that classifies lesions. The following sections of this paper will discuss the specific work that was completed, in particular, the methods used to conduct our research and the results that followed our methods in addition to future works to improve our results.

AIH

July 4, 2020

## II. METHOD

By means of semi-automatic segmentation and feature-detection algorithms analyzing the ratios of binarized lesion image-matrix sums and comparative intensity contrasts, we use the apparent size and intensity discrepancies of lesions in their respective elasticity and ultrasound images to first collect feature data with which to train classification mechanisms for breast lesions, and later as an integrated constituent of the culminating classifiers. The algorithm is factored into four sub-mechanisms being (1) strain sequence computation and analysis, (2) strain and B-Mode image computation, (2) strain and B-Mode image processing, (3) the segmentation of strain and B-Mode images, and (4) lesion area ratio and contrast computation.

### A. The Dataset: UVM Sonix-500RP Dataset

In this analysis, we use *in vivo* breast lesion patient data taken using the Ultrasonix Sonix-500RP (Ultrasonix Medical Corporation, Richmond, BC, Canada) commercial ultrasound research interface with a L145/38 probe operating at 10 MHz, the manufacturer-provided nominal value, at the University of Vermont Medical Center (UVM). Each patient gave informed consent for data collection, and the Institutional Review Board (IRB) approved the study. For each lesion, compression was applied manually (freehand) for the RF Data (\*.rf) and real-time strain sequences (\*.avi) compression. In addition, histological information regarding the lesion including position information, size, and classification/pathological nature was also stored in the data repository. These features provided the ground truth for this analysis. Of the 46 lesions recorded in the database, eleven were classified as adenocarcinoma and eleven as fibroadenoma. These lesions were extensively analyzed and

together comprised the training dataset used for the final SVM classifier function.

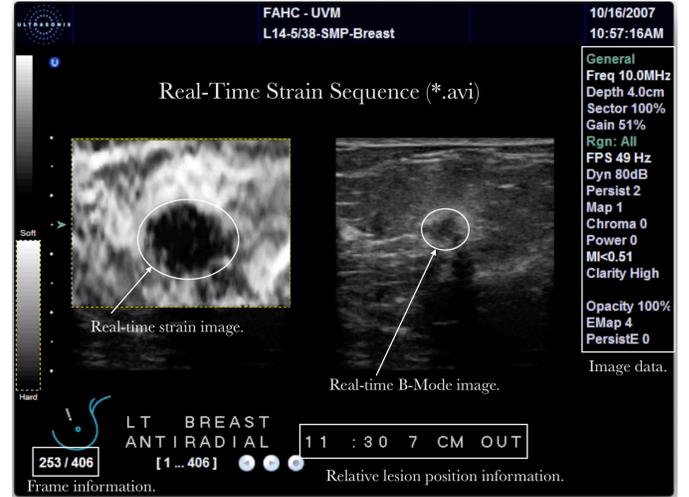


Fig. 1: Example of a real-time strain sequence (\*.avi) of a lesion from the UVM data repository.

### B. Strain Sequence Computation and Analysis

Adaptive-stretching strain estimation, an algorithm that maximizes the correlation between the pre-compression and post-compression RF echo-signals by iteratively stretching the latter, was used in order to compute the elasticity image sequences for each lesion in the data repository. From these sequences, one strain image was selected from the corresponding RF data to calculate the area of the lesion in the elastogram. The strain images were computed in a two-step approach where the mean, minimum, and maximum strain values estimated in the first iteration of the algorithm were used as the applied strain, minimum strain, and maximum strain processing parameters, respectively, for the computation of the final strain map in the second iteration. This novel two-step approach to adaptive stretching has never previously been attempted, and we have observed it to produce significantly more robust strain maps than with the single-step approach.

In order to determine the nominal frame pairs for the given RF data and ensure that the region depicted in a given strain image was the ROI (region of interest), a function iterated through the frames in each RF creating a sequence of strain images with the  $i^{th}$  and the  $i + f_{skip}^{th}$  frame pairs. An  $f_{skip}$  of 10 was found to be nominal and, thus, was used as the  $f_{skip}$  in processing all RF data. These strain sequences were created for each lesion in the dataset, and each frame in the sequence was then manually compared with the corresponding \*.avi to determine which frame pair produced a strain image in which the lesion was best articulated. These nominal frame pairs were later used to compute the strain image in the algorithm comparing the areas between the lesion as shown in the strain image and the lesion as shown in the B-Mode image.

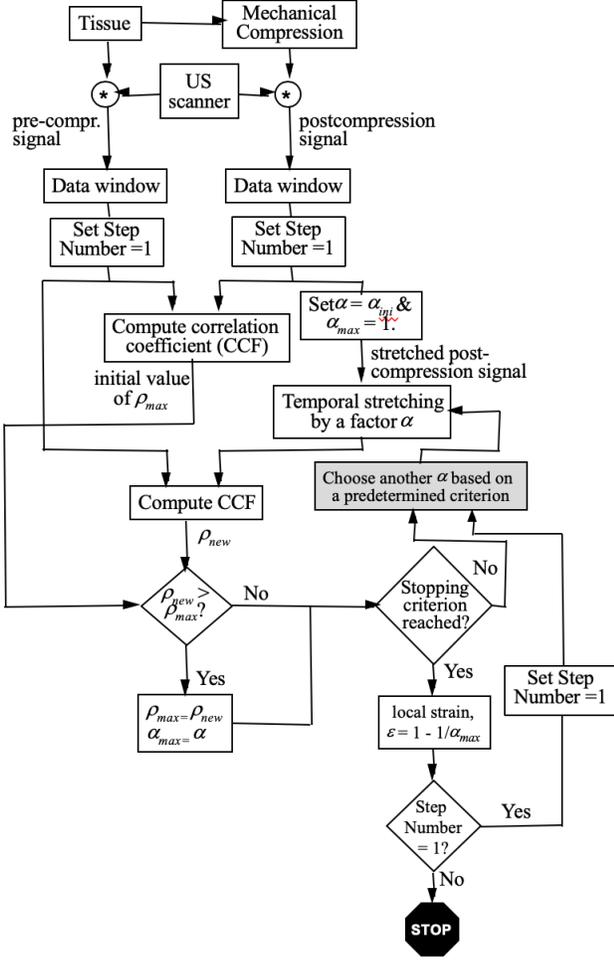


Fig. 2: Flowchart illustration of the novel two-step adaptive stretching strain estimation algorithm.

### C. Strain and B-Mode Image Computation and Processing

Once elasticity image sequences had been manually analyzed and nominal pre-compression and post-compression RF echo-signals frame pairs had been selected, the RF source data and frame pairs were passed to the feature classification algorithm for strain/B-Mode computation. The elastogram computations were done via the two-step adaptive-stretching strain estimation, and the B-Mode images were computed applying the conventional Hilbert transform method on RF frame data. The images were then processed to better articulate the lesion and improve the accuracy of the segmentation algorithm. The strain images were upsampled via bicubic interpolation in order to smoothen the lesion boundary from its pixelated appearance as well as match matrix dimensions with the B-Mode image necessary for the area comparison algorithm. The images were further histogram equalized and intensified to enhance the contrast between the lesion and its surrounding tissue. 2D external balloon force vector fields were then calculated from the strain image to process the image for energy minimization via the deformable splines in Active Contour segmentation. In addition, a phase-in-

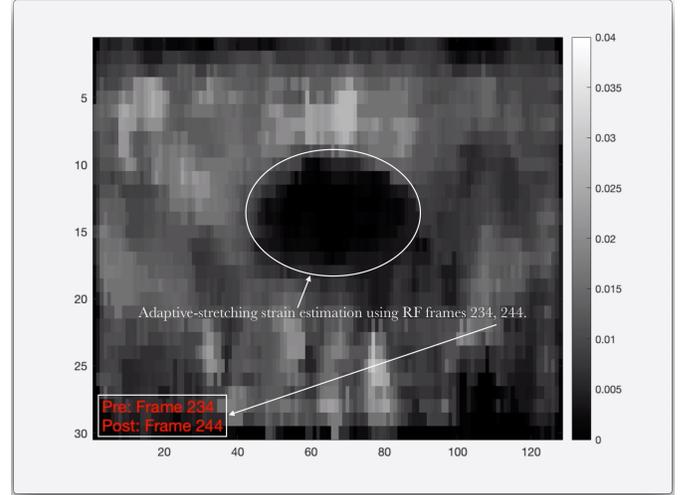


Fig. 3: Example of a frame from a strain sequence depicting an adenocarcinoma. Adaptive-stretching strain estimation with pre-compression and post-compression RF echo-signal frames 234 and 244 was used to calculate this frame.

maximum-orientation (PMO) image of the strain image processed with  $3 \times$  median filtering, brightening, and intensification was also computed. The matrix product of the vector field and the PMO image were added to the pre-processed strain image in order to suppress the effect of external forces on the strain computation. These image matrix sums were then logarithmically compressed to reduce the impact of a few high amplitude points overshadowing the rest of the image by nonlinearly mapping the amplitude values and adjusting the dynamic range. The final log compressed strain images were then intensified once more before the segmentation label regions are given to the user to be delineated.

The B-Mode segmentation pre-processing algorithm, while similar to strain processing, involved more processes given the lesions were more difficult to discern in the B-Mode image. The image was first logarithmically compressed prior to being processed with a speckle reducing anisotropic diffusion algorithm and a level of histogram equalization. Following a level of intensification, a 2D external balloon force vector field was calculated from the image as well as a PMO image of the strain image processed with  $3 \times$  median filtering, brightening, and intensification, the matrix product of which were added to the pre-processed B-Mode image. The resulting image was logarithmically compressed and intensified once more before the segmentation label regions were given to the user to be delineated.

### D. Segmentation of Strain and B-Mode Images

In order to determine the most accurate means by which to segment any such lesion, two segmentation mechanisms were tested. The first was GrowCut, an interactive multi-label N-dimensional image segmentation algorithm. The second was Active Contour, an energy minimizing, deformable-spline-based image segmentation algorithm which necessarily was used in addition to GrowCut serving as an empirical

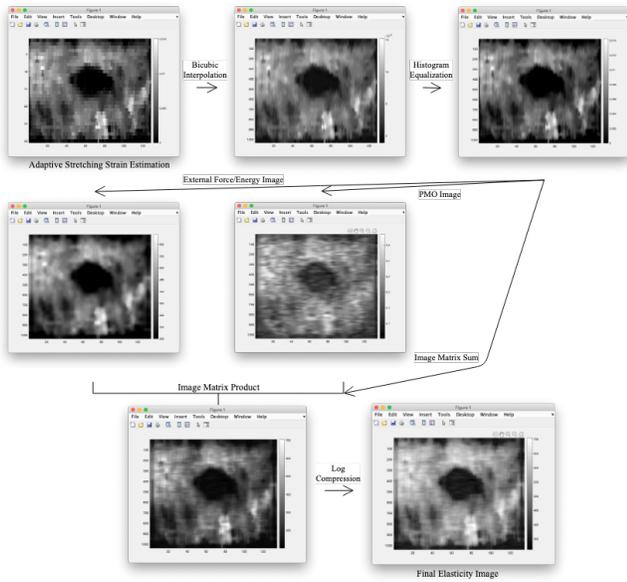


Fig. 4: Flowchart illustration of the elasticity image processing procedure.

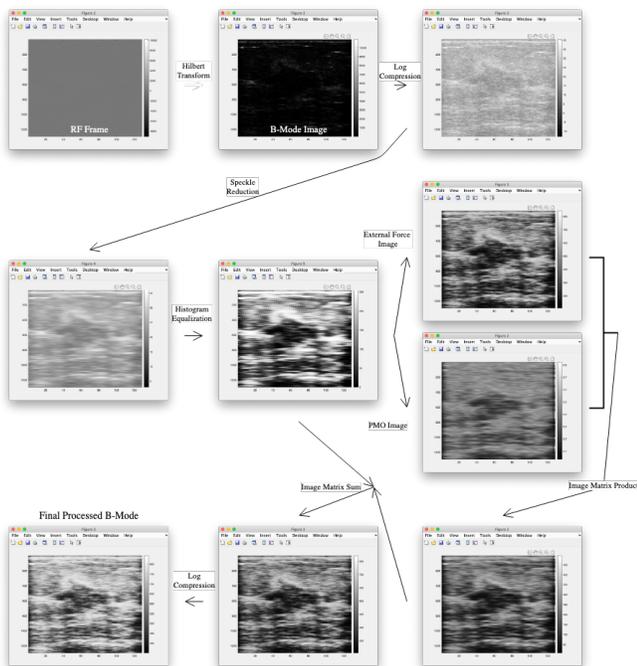


Fig. 5: Flowchart illustration of the B-Mode image processing procedure.

segmentation algorithm. GrowCut segmentation requires labels indicating two regions within the input image: one region outside of which no part of the ROI (the lesion) lies and one region inside of which lies only constituents of the ROI (the lesion). The GrowCut algorithm then takes the inner region and grows it outward no farther than the boundary of the outer region via cellular automata evolution. Via the inter-actability of figure manipulation mechanisms, the inner and outer regions are manually delineated, first on the strain

image and then on the B-Mode image. The markers are then casted as logical matrices, the manner in which label data is passed to the GrowCut algorithm. The empirical segmentation of the final binarized output of GrowCut in addition to the original strain and B-Mode images were passed as parameters to Active Contour. The segmentation algorithm used the point distribution model via the energy-minimization of deformable splines to produce the final segmented lesion images.

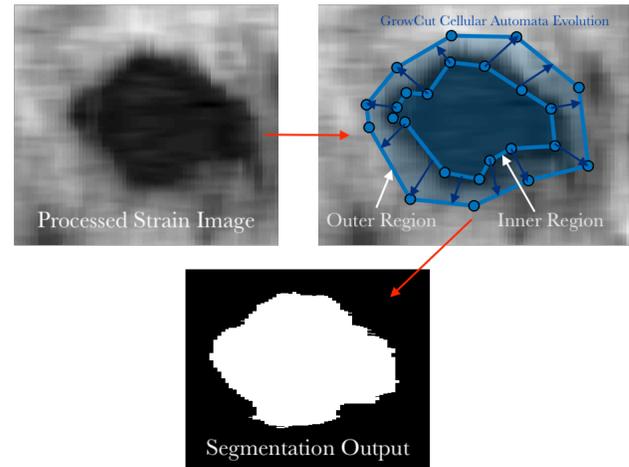


Fig. 6: Visual illustration of elasticity image segmentation via the GrowCut algorithm.

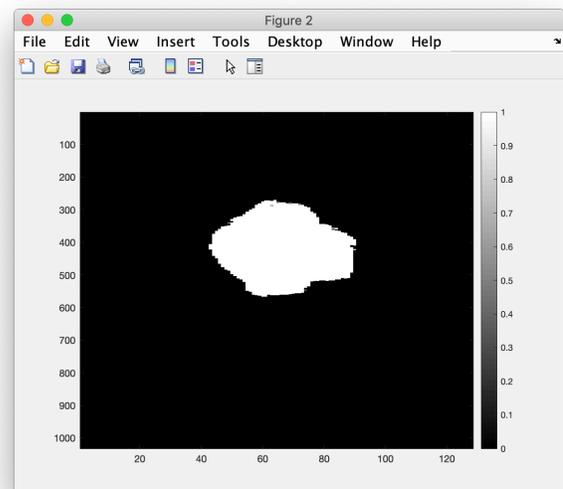


Fig. 7: The final segmented output of an elasticity image processed via the segmentation algorithm described in this paper. The segmentation above was

#### E. Lesion Area Ratio and Contrast Computation

Once the final segmented lesion elasticity and B-Mode images have been processed via Active Contour, the areas and contrasts of the lesion in both images are computed. Both

images are binarized logical matrices where the region in which the lesion lies is represented with a 1 and the region in which the lesion does not is represented with a 0 and are of matching dimensions. Thus, the elements of each matrix can be summed in order to determine a value representing the size of the lesion in the elasticity and B-Mode images. Furthermore, the area ratio can be calculated by taking the area of the lesion in the strain image and dividing it by that of the lesion in the B-Mode image (the order, as in which area is divided by which, is arbitrary as long as it remains consistent among computations). The contrast for both the strain and B-Mode images were computed by averaging the acoustic impedance/strain values inside the lesion and those outside the lesion (as defined by the final segmented image), taking the absolute value of their difference and dividing it by the sum of the squares of the acoustic impedance/strain variances inside the lesion and outside the lesion.

$$Area = \sum_{\langle i,j \rangle} I_{ij}$$

$$Contrast = \frac{|\bar{Z}_{Lesion} - \bar{Z}_{Surrounding}|}{\sigma_{Lesion}^2 + \sigma_{Surrounding}^2}$$

In addition to computing the lesion contrast in both the elasticity and B-Mode images, the ratio between the two (given by the contrast in the strain image divided by that of the B-Mode image) were also computed and analyzed for trends in this study.

### III. RESULTS

For each RF file, label regions defining the location of the ROI in both elasticity and B-Mode images were given to the feature detection image processing algorithm only once to minimize the variability due to the sensitivity of the output segmentation to small differences in the provided regions. The algorithm determined the values of four features: (1) the area ratio between the strain and B-Mode images, (2) (3) the contrast between the lesion and its surroundings in both strain and B-Mode images, and (4) the ratio between the strain and B-Mode contrasts. These features were computed via four different procedures: the GrowCut segmentation of an RF echosignal as processed by the algorithm detailed in II.C (PG), the GrowCut segmentation of an unprocessed RF echosignal (NG), the Active Contour segmentation of a processed RF echosignal as empirically segmented via GrowCut (PA), and the Active Contour segmentation of an unprocessed RF echosignal as empirically segmented via GrowCut (NA). This yielded a total of sixteen values describing feature information being recorded for each patient.

#### A. Analysis of Segmentation Procedures

In order to determine which of the four feature detection algorithms (PA, PG, NA, NG) provided the most significant distinction between adenocarcinoma and fibroadenoma for each feature, and, thus, could effectively be used as a classification mechanism for said criterion, the p-value percentages from two-tailed t-tests between the two classification sets were computed as follows:

Feature	Processing	Segmentation	P-value %
Area Ratio	Processed	GrowCut	0.0008
	Processed	Active Contour	0.0022
	Unprocessed	GrowCut	0.0684
	Unprocessed	Active Contour	0.3900
Strain Contrast	Processed	GrowCut	0.0781
	Processed	Active Contour	0.2053
	Unprocessed	GrowCut	0.0188
	Unprocessed	Active Contour	0.0290
B-Mode Contrast	Processed	GrowCut	63.04
	Processed	Active Contour	94.38
	Unprocessed	GrowCut	70.67
	Unprocessed	Active Contour	82.72
Contrast Ratio	Processed	GrowCut	79.42
	Processed	Active Contour	74.25
	Unprocessed	GrowCut	6.588
	Unprocessed	Active Contour	14.71

**Table 1.** Feature t-tests analyzing the significance of the distinction between fibroadenoma and adenocarcinoma with various processing-segmentation algorithm permutations.

As demonstrated in the table above, the GrowCut segmentation of **processed** elasticity and B-Mode images best differentiated between fibroadenoma and adenocarcinoma on the basis of *area ratio* and *B-Mode contrast*. Furthermore, the GrowCut segmentation of **unprocessed** strain and B-Mode images best differentiated between the two on the basis of *strain contrast* and *contrast ratio*. Not for any for the four features was an Active Contour segmentation procedure most effective in differentiating between lesion classified as fibroadenoma and those as adenocarcinoma. Moreover, the additional use of the processing algorithm as detailed in II.C produced a difference  $85.5\times$  more distinct between adenocarcinoma and fibroadenoma on the basis of area ratio when segmented with GrowCut. This is indicated in the p-value percentage reduction from 0.0684%  $\rightarrow$  0.0008% as well as a difference  $1.12106\times$  more distinct between the two on the basis of B-Mode contrast when segmented with GrowCut given the reduction from 70.6687%  $\rightarrow$  63.0371%. Finally all four processing-segmentation permutations provided a significant distinction ( $p\text{-value} < 0.05 = 5\%$ ) between adenocarcinoma and fibroadenoma on the basis of area ratio. Only the segmentation of unprocessed elasticity and B-Mode images via both GrowCut and Active Contour did so on the basis of strain contrast. None of the four algorithm permutations provided a significant distinction between adenocarcinoma and fibroadenoma on the basis of B-Mode contrast and contrast ratio. The lowest p-value percentages for each of the four features and the corresponding processing-segmentation algorithm permutation that produced it are shown below.

In order to quantify the distinction among the four tested processing-segmentation algorithm permutations, double-tailed t-tests between the output feature determination data of pairs of algorithms were computed between those pairs in which one of either the processing or segmentation mechanism remained constant.

As can be seen from the algorithmic comparison data, with

Feature	Processing	Segmentation	P-value %
Area Ratio	Processed	GrowCut	8.43e-6
Strain Contrast	Unprocessed	GrowCut	1.8824
Contrast Ratio	Unprocessed	GrowCut	6.5879
B-Mode Contrast	Processed	GrowCut	63.0371

**Table 2.** Feature t-tests analyzing the significance of the distinction between fibroadenoma and adenocarcinoma with the processing-segmentation algorithm permutations producing only the lowest p-value percentages.

Feature	Algorithm I	Algorithm II	P-value %
Area Ratio	PG	PA	8.43e-6
	NG	NA	2.18e-5
	PG	NG	0.000684
	PA	NA	0.00391
Strain Contrast	PG	PA	0.0781
	NG	NA	0.205
	PG	NG	0.0188
	PA	NA	0.0290
B-Mode Contrast	PG	PA	0.630
	NG	NA	0.944
	PG	NG	0.707
	PA	NA	0.827
Contrast Ratio	PG	PA	0.794
	NG	NA	0.743
	PG	NG	0.066
	PA	NA	0.147

**Table 3.** Algorithmic comparison t-tests analyzing the significance of the distinction between processing and not processing as well as segmenting via GrowCut as opposed to via ActiveContour.

regards to the determination of area ratio, the GrowCut segmentation of processed RF echosignals produced significantly different feature data from the Active Contour segmentation of processed RF echosignals. The GrowCut segmentation of unprocessed images, however, did not produce significantly different data from the Active Contour segmentation of unprocessed images. Given that the GrowCut segmentation of processed RF echosignals produced a stronger distinction in area ratio between fibroadenoma and adenocarcinoma than the Active-Contour segmentation of these echosignals, we conclude the processing algorithm is better optimized for GrowCut segmentation and intensifies the distinction between GrowCut and Active Contour segmentation with reference to area ratio. Furthermore, processing does not appear to significantly change the area ratio produced by either segmentation algorithm given the p-value percentages of 10.63% and 10.19% between the data produced by processed and unprocessed images segmented by GrowCut and that by processed and unprocessed images segmented by Active Contour, respectively. This insignificant distinction, however, rather than necessarily being a reflection on the processing algorithm could arise from the insensitivity of the segmentation algorithms to changes in contrast and compression.

With regards to the determination of strain contrast, there does not seem to be any significant distinction between the data produced by GrowCut and Active Contour segmenta-

tion for processed RF echosignals and even more so for unprocessed echosignals. There does, however, appear to be a significant distinction between the strain contrast data produced by processed and unprocessed images segmented by GrowCut and that by processed and unprocessed images segmented by Active Contour. Given that unprocessed images produced strain contrast data more strongly differentiating between adenocarcinoma and fibroadenoma, we can conclude that processing reduces the distinction between adenocarcinoma and fibroadenoma with reference to strain contrast.

With regards to the determination of B-Mode contrast, there appears to be a significant distinction between the data produced by GrowCut and Active Contour segmentation for unprocessed RF echosignals and even more so for processed echosignals. Given that the GrowCut segmentation of processed and unprocessed RF echosignals produced a stronger distinction in area ratio between fibroadenoma and adenocarcinoma than the ActiveContour segmentation of these echosignals, we conclude the processing algorithm is better optimized for GrowCut segmentation and intensifies the distinction between GrowCut and Active Contour segmentation with reference to area ratio. Furthermore, there appears to be a significant distinction between the B-Mode contrast data produced by processed RF echosignals when segmented by either GrowCut or ActiveContour. However, the fact that this distinction increases the distinction between fibroadenoma and adenocarcinoma for GrowCut segmentation and reduces the distinction for ActiveContour segmentation further supports the conclusion that the processing algorithm is better optimized for GrowCut segmentation.

With regards to the determination of contrast ratio, there does not seem to be any significant distinction between the data produced by GrowCut and Active Contour segmentation for processed RF echosignals and even more so for unprocessed echosignals. This manifests in a manner similar to the determination of strain contrast indicating that processing increases the distinction between the contrast ratio data produced by GrowCut and ActiveContour. Furthermore, there does not appear to be any significant distinction between the contrast ratio data produced by processed and unprocessed RF echosignals segmented by GrowCut or ActiveContour.

### B. Analysis of Lesion Features

A p-value of less than 0.05 ( $p < 5\%$ ) is considered significant. Two of the four parameters perform well, one perform almost adequately, and one is useless in discriminating between fibroadenoma and IDCA. Clearly, the ratio between lesion sizes in the strain image (computed using the algorithm in II.C) and the B-mode image performs the best job (of the four tested) ( $p < 0.0008\%$ , which is much smaller than 5%). Strain contrast (between the lesion and the background/surrounding tissue) also performs really well ( $p < 1.8824\%$ ). The ratio between strain and B-mode contrasts is not quite significant ( $p < 6.6\%$ ). B-mode contrast cannot distinguish between fibroadenoma and IDCA at all ( $p < 70.7\%$ , which is much larger than 5%).

[Given the corresponding p-value percentage of 0.0008%  $\ll$  5%, the ratio between the size of a lesion in

the B-Mode of an RF echo signal and the elasticity image calculated via two-step adaptive-stretching strain estimation from the same RF echosignal and the RF echosignal 10 frames ahead of it can serve as an effective criterion—the most of the four tested— by which to classify lesions in RF echosignals as benign fibroadenoma or cancerous adenocarcinoma when processed via the algorithm detailed in II.C and segmented via GrowCut cellular automata evolution. While certainly to less of an extent than area ratio with the corresponding p-value percentage of  $1.8824\% < 5\%$ , the elasticity intensity contrast between the ROI and surrounding tissue can further serve as a mostly effective classification criterion. On the other hand, the ratio between the intensity contrast in the strain and B-Mode images, while close to the  $p\text{-value} = 0.05 = 5\%$  significance threshold, cannot quite produce distinct enough of a differentiation between fibroadenoma and adenocarcinoma, and the B-Mode contrast offers little insight into the classification of a lesion in an RF echosignal with a p-value percentage  $12.6074\%$  greater than the significance threshold.]

### C. Classification Algorithms

The four features shown in Table 1 were computed using four processing permutations. We used two machine learning algorithms to classify the breast masses into adenocarcinoma or fibroadenoma: (1) a support vector machine (SVM) classifier. We used biopsy results as ground truth.

[Using the feature data collected from all four processing-segmentation algorithm permutations and the ground truth classifications provided by the dataset, we used a support vector machine (SVM) classifier machine-learning classification mechanism in order to determine whether a given RF echosignal depicted an adenocarcinoma or fibroadenoma]

1) *Support Vector Machine (SVM) Classification:* In training four support vector machine classification mechanisms with the feature data collected from the processing-segmentation algorithm permutations that produced the lowest p-values for each feature and the corresponding ground truth classifications, classification thresholds for each feature are determined to differentiate between RF echosignals depicting fibroadenoma and adenocarcinoma. Following the Z-Score normalization of the feature data, the leave-one-out iterative validation schema where each point serves as the validation set and the validation accuracy is calculated as an average over the dataset, divides the dataset into the training and validation sets. Gradient Ascent, a first-order iterative maximization algorithm of the concave objective function as defined by the training set, is applied with a linear kernel— given the features appear to indicate a certain classification by means of a constant threshold value. Since the training algorithm depends solely on the training data as represented as dot products in  $H$ , i.e. on functions of the form  $\Phi(x_i) \cdot \Phi(x_j)$ , if there were a kernel function  $K$  such that  $K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$ , only  $K$  would be a necessary parameter for the training algorithm and  $\Phi$  would never need to explicitly be determined. The optimal regularization parameter as determined by the highest f-measure yield is then used to determine the final model which

produces a constant line, the value of which is the classification threshold. The threshold values and the corresponding leave-one-out validation accuracy for each feature can be seen below.

Feature	Threshold	Accuracy (%)
Area Ratio	2.768960	100%
Strain Contrast	$1.17 \times 10^8$	81.82%
Contrast Ratio	$2.30 \times 10^{18}$	81.82%
B-Mode Contrast	$1.2 \times 10^{-5}$	54.55%

**Table 3.** The SVM-determined thresholds and their corresponding validation accuracies for the four tested features.

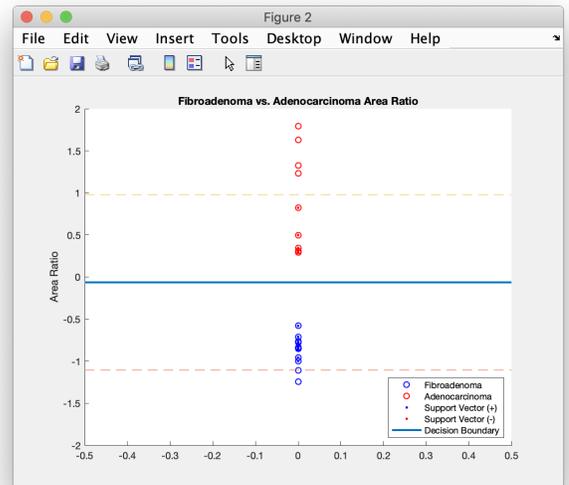


Fig. 8: The final output of the trained support-vector machine threshold determination algorithm. The fibroadenoma are demarcated with blue dots, the adenocarcinoma with red dots, and the threshold with a solid blue line.

Size discrepancy between the B-mode and strain images yielded a validation accuracy of 100% when classified with an SVM classifier.

[As demonstrated by the leave-one-out validation accuracy of 100%, the size discrepancy between a lesion in the B-Mode of an RF echosignal and the elasticity image calculated via the two-step adaptive-stretching strain estimation of that echosignal and one 10 RF frames ahead can serve as an effective classification criterion when implemented with a support-vector machine classifier.]

## IV. DISCUSSION

Early detection of breast cancers can save lives of those affected. In the previous sections, we have discussed a semi-automated method to discriminate between benign and malignant lesions in order to assist clinicians with the diagnosis of breast lesions. Our algorithm utilizes feature-based decision tree classification models to analyze the ratios of estimated lesion sizes in B-mode and strain images later in the classification stage. This can serve as an effective classification



expertise in elastography. Follow-up studies could introduce a fully automatic segmentation algorithm. We have already developed a novel algorithm that can automatically select a seed point (a point inside the lesion). We are perfecting our algorithm so that it can identify the lesion contour correctly for different scenarios.

Additionally, implementing a hyperparameter optimization mechanism for segmentation pre-processing minimizing double-tailed t-test outputs comparing the features extracted from fibroadenoma and adenocarcinoma RF echosignal segmentation would be an effective means to improve the PG-NG/PA-NA p-values [6]. A fully-automated algorithm could be beneficial for diagnostic purposes in the absence of a radiologist. Some important foundational work has begun, though much remains to be studied [5] [7]. The opportunity is great; as there have been multiple foundational studies to advance this topic of urgency. In fact, we are currently working on a fully automated segmentation method that serves as an effective computer-aided diagnostic (CAD) mechanism for classifying adenocarcinoma and fibroadenoma.

The major impact of the algorithms we designed to compute the size difference in strain and B-mode images of breast lesions is to assist in improving the modern methods radiologists currently utilize to observe lesions in strain images for adenocarcinoma lesions and expedite the development of fully-automated algorithms as mentioned in the previous paragraph [7]. Our hope is that our work will support our radiologists in their attempt to diminish breast cancer mortality, which is the most common cancer in American women. [1] Our study unlocks many future possibilities for fully automating the size difference in strain and B-mode images of breast lesions, in addition to providing considerable contributions to the field of region growing ultrasound. Our study and algorithms present effective classification criterion between fibroadenoma and adenocarcinoma in the case of both machine learning classifiers which are discussed in the next section.

## V. CONCLUSION

Using a semi-automatic feature-based support-vector machine, we analyze the ratios of binarized lesion image matrix sums and comparative intensity contrasts. Consequently, we find that the apparent size discrepancy as determined by GrowCut cellular automata evolution between a lesion in the processed B-Mode of an RF echosignal and the processed elasticity image calculated via the two-step adaptive-stretching strain estimation of that echosignal and one fixed quantity of frames ahead can serve as an effective classification criterion between fibroadenoma and adenocarcinoma in the case of both machine learning classifiers. However, this cannot be attributed to the intensity contrast between a lesion and its surrounding in neither the B-Mode nor elasticity image of the corresponding RF echosignal regardless of processing by means of the algorithm detailed in II.C and the applied segmentation algorithm between GrowCut and Active Contour and the ratio between the intensity contrasts as described above for the B-Mode images and strain images of an RF echosignal. Furthermore, we find both support-vector machine

classifiers accepting RF echosignal sequences as input and returning lesion classifications as output trained on the basis of the size discrepancy, contrast ratio, elasticity contrast, and B-Mode contrast to be effective classification models by which to differentiate between the RF echosignals of fibroadenoma and those of adenocarcinoma.

## ACKNOWLEDGMENT

The authors would like to thank The University of Vermont for providing access to the breast lesion data repository used for the training and validation of the classification mechanisms detailed here...

## REFERENCES

- [Dis17] Center for Disease Control {and} Prevention. *USCS Data Visualizations*. United States Cancer Statistics: Data Visualizations. Library Catalog: gis.cdc.gov. 2017. URL: <https://gis.cdc.gov/grasp/USCS/DataViz.html> (visited on 07/04/2020).
- [NIH16] NIH National Cancer Institute. *Cancer of the Breast (Female) - Cancer Stat Facts*. SEER. Library Catalog: seer.cancer.gov. 2016. URL: <https://seer.cancer.gov/statfacts/html/breast.html> (visited on 07/04/2020).
- [Nata] National Breast Cancer Foundation. *Breast Cancer Stage 0 & Stage 1*. National Breast Cancer Foundation. Library Catalog: www.nationalbreastcancer.org. URL: <https://www.nationalbreastcancer.org/breast-cancer-stage-0-and-stage-1/> (visited on 07/04/2020).
- [Natb] National Breast Cancer Foundation. *Breast Cancer Stage 4*. National Breast Cancer Foundation. Library Catalog: www.nationalbreastcancer.org. URL: <https://www.nationalbreastcancer.org/breast-cancer-stage-4/> (visited on 07/08/2020).
- [GC97] Brian S. Garra and E. Ignacio Cespedes. "Elastography of Breast Lesions: Initial Clinical Results". In: *Radiology* 202.1 (1997), pp. 79–86. (Visited on 07/04/2020).
- [Bar+12] Richard G. Barr, Stamatia Destounis, Logan B. Lackey, et al. "Evaluation of Breast Lesions Using Sonographic Elasticity Imaging: A Multicenter Trial". In: *Journal of Ultrasound in Medicine* 31.2 (Feb. 2012), pp. 281–287. ISSN: 02784297. DOI: 10.7863/jum.2012.31.2.281. URL: <http://doi.wiley.com/10.7863/jum.2012.31.2.281> (visited on 07/04/2020).
- [Muk+16] Rashid Al Mukaddim, Juan Shan, Irteza Enan Kabir, et al. "A novel and robust automatic seed point selection method for breast ultrasound images". In: *2016 International Conference on Medical Engineering, Health Informatics and Technology (MediTec)*. 2016 International Conference on Medical Engineering, Health Informatics and Technology (MediTec). Dhaka, Bangladesh: IEEE, Dec. 2016, pp. 1–5. ISBN: 978-1-5090-5421-3. DOI: 10.1109/MEDITEC.2016.7835370.

URL: <http://ieeexplore.ieee.org/document/7835370/> (visited on 07/04/2020).