

# Preliminary Feature Extraction for Small Lesion Classification in Sonomammographic Images

Anna Pawlowska

*Department of Ultrasound  
Institute of Fundamental Technological  
Research, Polish Academy of Sciences  
Warsaw, Poland  
0000-0001-5070-1786*

Norbert Zolek

*Department of Ultrasound  
Institute of Fundamental Technological  
Research, Polish Academy of Sciences  
Warsaw, Poland  
0000-0002-2416-7783*

**Abstract**—In breast cancer diagnosis, early detection of tumors and accurate differentiation of malignant and benign breast lesions are key demands. Tumor size, as a measure of tumor progression, is related to recurrence rate and patient survival. This study aims to determine which sonographic features allow the differentiation of small breast lesions into benign and malignant. Inclusion criteria for the analysis were tumors with the longest diameter of less than or equal to 10 mm and tumors with confirmed classification by follow-up care or core needle biopsy result. Following the criteria, 1515 cases were analyzed, including 365 carcinomas and 1150 benign lesions. To quantitatively evaluate the images, 383 ultrasound parameters (BI-RADS features, morphological features, fractal features, histogram features, and texture parameters) were used. Univariate and multiple logistic regression analyses were used to assess the significance of various diagnostic features and their combinations. The combined classifier (based on 19 quantitative features) yields an area under the ROC curve of 0.91.

**Keywords**—Ultrasound imaging, Breast cancer, Small lesion, Feature extraction

## I. INTRODUCTION

Ultrasound is widely used to differentiate potentially malignant masses from benign breast lesions. An important diagnostic aspect is the degree of tumor progression. Tumor size correlates with the potential for metastasis, and thus with recurrence rates and patient survival [1]. Therefore, early detection and accurate classification are of great importance.

Differentiation of breast lesions is usually based on variable sonographic characteristics due to differences in histologic type, tumor grading and tissue components within the tumors. Smaller breast cancers tend to have a lower histological grade, fewer desmoplastic changes, less necrosis and less aggressive invasion into surrounding tissues [2]–[5]. Sonographic characteristics of breast lesions have been standardized using the BI-RADS lexicon [6], which provides a consistent breast imaging terminology, report structure and classification system.

Addressing the diagnostic importance of differentiating between benign and malignant tissue in small lesions, the studies [7], [8] were conducted to determine which individual and combined BI-RADS features have the highest diagnostic

accuracy (for 190 [7] and 1203 [8] cases included). Another study [9] identified differences in ultrasound features between small breast cancer (size  $\leq 5$  mm; 62 cases) and large breast cancer ( $> 5$  mm; 466 cases).

To date, no analysis of classification performance for small lesions ( $\leq 10$  mm in maximal diameter) using quantitative ultrasound parameters has been reported. Therefore, in this study, BI-RADS descriptors, morphological features, fractal features, histogram features, and texture parameters were considered, and classification evaluation was performed for each classifier separately and combined sonographic features.

## II. METHODS

### A. Data collection

The data collection protocol was approved by an institutional review board at the Lower Silesian Chamber of Medicine no. 2/BNR/2022. Inclusion criteria were tumors with the longest diameter of less than or equal to 10 mm and tumors with confirmed classification by follow-up care or core needle biopsy result. Applying the criteria to the database developed in [10], 1515 cases were analyzed, consisting of 365 carcinomas and 1150 benign lesions.

### B. BI-RADS descriptors

For each case, an experienced radiologist who performed the examination determined BI-RADS descriptors [6] and segmented the lesion. Among the BI-RADS descriptors [6], only those Bmode-related were included, i.e. shape (oval, round, irregular), orientation (parallel, not parallel), margin (circumscribed, not circumscribed), echogenicity (anechoic, hypoechoic, isoechoic, hyperechoic, complex), posterior acoustic features (enhancement, shadowing, combined pattern, no), hyperechoic halo (yes, no), calcifications (in a mass, outside of a mass, no), skin thickening (yes, no).

### C. Quantitative ultrasound parameters

To evaluate tumors quantitatively, 383 parameters from the BUSAT Toolbox [11] were used. From images and tumor masks, the following parameters were determined:

- 212 quantitative measures of BI-RADS descriptors for shape, orientation, margin, boundary, echogenicity, posterior behavior and spiculation;
- 13 histogram parameters;
- 11 fractal parameters (3 for contour and 8 for texture);
- 147 texture parameters (21 gray-level co-occurrence matrix features for 7 distances computed based on known image pixel size).

#### D. Statistical analysis

Univariate and multiple logistic regression analyses were used to assess the significance of various diagnostic features and their combinations. The chi-square test ( $p$ -value  $< 0.05$ ) and the backward elimination method were used to exclude parameters from the logistic regression model. To evaluate the prediction of the logistic regression model, ROC curves were used. All statistical analyses were carried out using R Statistical Software (version 4.3.1; R Core Team 2023).

### III. RESULTS AND DISCUSSION

#### A. Data characteristics

Of the 1,515 included lesions of 10 mm or less in size, 365 cases (24%) were carcinomas. For 98.26% of benign lesions and 94.79% of carcinomas, the examining physician found no signs. The patient's history reported no symptoms in 94.61% and 93.97% of cases for benign and malignant lesions, respectively. The average diameters for malignant and benign masses were 7.67 and 8.12 mm, respectively. Detailed characteristics of the data are shown in Table I.

TABLE I  
CLINICAL CHARACTERISTICS OF PATIENTS (AGE, SIGNS AND SYMPTOMS) AND TUMOR SIZE STATISTICS. SD - STANDARD DEVIATION.

	Benign (n=1150)	Malignant (n=365)	Total (n=1515)
Mean age (SD) [years]	46.9 (13.8)	60.4 (11.7)	50.6 (14.5)
Signs			
no	1130 (98.26%)	346 (94.79%)	1476 (97.43%)
palpable	20 (1.74%)	16 (4.38%)	36 (2.38%)
edema	0 (0.00%)	3 (0.82%)	3 (0.20%)
Symptoms			
no	1088 (94.61%)	343 (93.97%)	1431 (94.46%)
pain	2 (0.17%)	2 (0.55%)	4 (0.26%)
nipple discharge	24 (2.09%)	0 (0.00%)	24 (1.58%)
family history of cancer	34 (2.96%)	10 (2.74%)	44 (2.90%)
personal history of breast cancer	2 (0.17%)	2 (0.55%)	4 (0.26%)
mutation BRCA1(+)	0 (0.00%)	8 (2.19%)	8 (0.53%)
Tumor size [mm]			
Mean (SD)	7.67 (2.02)	8.12 (1.78)	7.18 (1.97)
Min-Max	2.26-9.99	2.86-10.00	2.26-10.00

#### B. BI-RADS descriptors

All BI-RADS descriptors were used as classifiers in univariate logistic regression. Only skin thickening was not an ultrasound feature significantly associated with malignancy. Then, a multivariate logistic regression analysis was performed with malignancy as the dependent variable and all BI-RADS descriptors as independent variables. The full model showed,

in addition to the skin thickening feature, posterior features as statistically insignificant factors. Next, using the backward elimination method, the final reduced model for predicting tumor malignancy was determined. The results of the logistic regression analyses are shown in Table II.

To evaluate the classification performance, ROC analysis was used. The following results were obtained:

- sensitivity: 0.89;
- specificity: 0.75;
- area under the ROC curve: 0.83.

The first study [12] focusing on small breast cancer features showed that small cancers were hypoechoic attenuating and not circumscribed. Other studies [7], [8] have included classifications of benign and malignant small lesions. In the study [7] based on 190 tumors, univariate analysis showed that shape, orientation, margin, hyperechoic halo and posterior features were factors significantly associated with malignancy. Shape and margin were significant variables in the multivariate model. The study [8] based on 135 small ( $\leq 1$  cm) tumors in the univariate analysis showed that margin, shape and echogenicity were significant factors. In contrast, in the multivariate analysis, the margin was the only significant differentiating factor. Differences in results between the cited studies and the present study (especially in multivariate analyses results, orientation, shape, margin, echogenicity, and halo were found to be significant in the present study) are likely due to different cohort sizes.

#### C. Quantitative ultrasound parameters

Out of 383 ultrasound parameters, 292 parameters were found to be statistically significant factors in univariate analysis, which were then reduced to 19 using the backward elimination method in the multivariate logistic regression model. The significant parameters from the final model are summarized in Table III. Based on the regression coefficients for these 19 parameters (five shape parameters, six echo pattern parameters, two histogram parameters, five texture parameters and one fractal texture parameter), the malignancy probabilities for all cases were determined (Fig. 1).

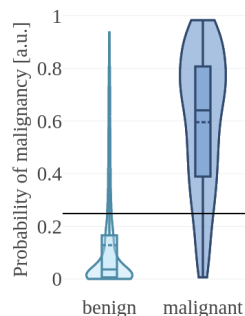


Fig. 1. Malignancy probability distributions for both tumor classes determined using the logistic regression model for quantitative parameters. The line indicates the cutoff value separating benign and malignant tumors.

TABLE II

UNIVARIATE AND MULTIVARIATE LOGISTIC REGRESSION ANALYSES FOR PREDICTING MALIGNANCY USING BI-RADS DESCRIPTORS. MULTIVARIATE REGRESSION ANALYSIS WAS PERFORMED FOR THE FULL AND REDUCED MODELS. THE RESULTS INCLUDED ODDS RATIO (OR), CONFIDENCE INTERVAL (CI) AND P-VALUE.

Descriptors	Descriptors' labels	Number of cases		Univariate mode OR (95% CI)	Multivariate full model Adjusted OR (95% CI)	Multivariate final model Adjusted OR (95% CI)
		benign	malignant			
Shape	irregular	379	326	-	-	-
	oval	671	27	0.05 (0.03-0.07, p<0.001)	0.39 (0.22-0.66, p=0.001)	0.38 (0.22-0.65, p<0.001)
	round	100	12	0.14 (0.07-0.25, p<0.001)	0.51 (0.22-1.10, p=0.097)	0.50 (0.22-1.07, p=0.086)
Orientation	not parallel	198	291	-	-	-
	parallel	952	74	0.05 (0.04-0.07, p<0.001)	0.45 (0.30-0.67, p<0.001)	0.43 (0.29-0.65, p<0.001)
Margin	not circumscribed	328	355	-	-	-
	circumscribed	822	10	0.01 (0.01-0.02, p<0.001)	0.06 (0.03-0.11, p<0.001)	0.06 (0.03-0.11, p<0.001)
Echogenicity	anechoic	140	9	-	-	-
	complex cystic/solid	79	1	0.20 (0.01-1.08, p=0.126)	0.12 (0.01-0.83, p=0.064)	0.08 (0.00-0.53, p=0.026)
	heterogeneous	223	29	2.02 (0.97-4.65, p=0.076)	0.49 (0.17-1.47, p=0.189)	0.47 (0.16-1.39, p=0.158)
	hyperechoic	41	0	0.00 (0.00-0.00, p=0.971)	0.00 (NA-4e+9, p=0.981)	0.00 (NA-1e+9, p=0.981)
	hypoechoic	544	269	7.69 (4.08-16.48, p<0.001)	0.93 (0.36-2.58, p=0.886)	0.95 (0.37-2.60, p=0.910)
	isoechoic	123	57	7.21 (3.59-16.14, p<0.001)	1.02 (0.37-3.00, p=0.975)	1.09 (0.40-3.20, p=0.865)
Posterior features	combined	7	10	-	-	-
	enhancement	122	9	0.05 (0.02-0.16, p<0.001)	0.29 (0.04-1.53, p=0.175)	-
	no shadowing	933	216	0.16 (0.06-0.43, p<0.001)	0.55 (0.09-2.30, p=0.453)	-
Halo	no	1100	168	-	-	-
	yes	50	197	25.80 (18.31-36.94, p<0.001)	5.34 (3.54-8.17, p<0.001)	6.11 (4.12-9.22, p<0.001)
Calcifications	in a mass	60	5	-	-	-
	no outside of a mass	1088	360	3.97 (1.75-11.43, p=0.003)	4.23 (1.57-13.67, p=0.008)	-
Skin thickening	no	2	0	0.00 (NA-8e+20, p=0.977)	0.00 (NA-4e+184, p=0.997)	-
	yes	1146	363	-	-	-
	yes	4	2	1.58 (0.22-8.12, p=0.599)	0.36 (0.04-2.61, p=0.314)	-

TABLE III

MULTIVARIATE LOGISTIC REGRESSION ANALYSIS FOR PREDICTING MALIGNANCY USING QUANTITATIVE ULTRASOUND PARAMETERS.

Parameter group	Parameter name	p-value	Regression coefficient
Shape parameters	sENC	p<0.001	-14.753
	sLS	p<0.001	3.597
	sAX_MX	p<0.001	-0.037
	sNRL_sd	p<0.001	59.149
	sNRL_ar	p<0.001	-81.246
Echo pattern parameters	eENTROr_D2_R8	p<0.001	-7.004
	eE5L5_mean	p<0.001	-1.761
	eS5L5_mean	p<0.001	4.307
	eR5E5_egy	p=0.002	-2.774
	eR5S5_egy	p=0.014	5.025
	eCORRm_D1_R4	p<0.001	30.521
Histogram parameters	hMSD	p<0.001	0.006
	hEgy	p=0.021	-30.052
Texture parameters	glcm_inf2h_D17_mean	p<0.001	-12.406
	glcm_dissi_D9_mean	p<0.001	1.359
	glcm_inf1h_D4_mean	p<0.001	-17.481
	glcm_indnc_D12_mean	p<0.001	202.768
	glcm_corr_m_D9_mean	p=0.003	2.842
Fractal texture parameter	FBM1	p<0.001	-50.987
	Intercept	p<0.001	-200.827

To evaluate the classification ability, the ROC curve was drawn (Fig. 2). Based on it, the cutoff point was chosen as the point closest to the point (0,1) on the ROC curve. The following results were obtained using the ROC curve analysis:

- sensitivity: 0.88;
- specificity: 0.83;
- area under the ROC curve: 0.91.

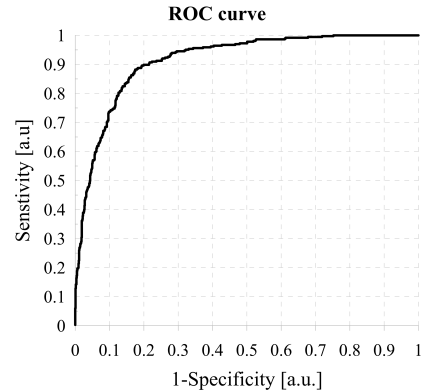


Fig. 2. ROC curve determined for the multivariate logistic regression model based on quantitative parameters.

Fig. 3 shows examples of images classified correctly and misclassified using the multivariate logistic regression model for quantitative parameters. Misclassified images, both false positive and false negative, are also misleading to the human observer.

To date, a study focused on classifying small lesion images using quantitative ultrasound parameters has not been published. However, studies using the same toolbox for parameter determination have been released. According to one of them [13], an accuracy of 0.88 was obtained using a classifier based on combined morphological features and 2054 images. Another study [14] developed a classifier using combined morphological and texture features and also yielded an accuracy of 0.88. Similar results (accuracy of 0.87-0.89) [15] were

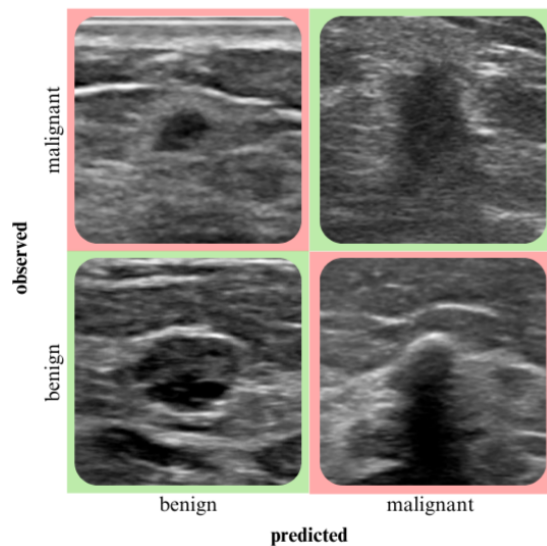


Fig. 3. Examples of images classified correctly and misclassified using the multivariate logistic regression model for quantitative parameters. The size of each image is 15x15mm.

achieved using a dataset of 2032 cases when different machine learning approaches were studied to classify breast lesions on ultrasound images.

#### IV. CONCLUSIONS

Ultrasound is useful in detecting and evaluating small tumors when BI-RADS guidelines are followed. Using logistic regression models with quantitative ultrasound parameters, diagnostic accuracy has improved (from 0.83 based on BI-RADS descriptors to 0.91) in differentiating between benign and malignant small lesions. However, due to the difference between the number of small cancers and benign lesions, data collection is required to be continued. Further work may also consider using complex models or neural networks for classification. Data-driven knowledge of ultrasound classification of small tumors may guide decisions on the biopsy of small lesions showing malignant features on ultrasound but yielding negative results on histopathology.

#### ACKNOWLEDGMENT

This work was supported by the Polish National Centre for Research and Development (INFOSTRATEG-I/0042/2021).

#### REFERENCES

- [1] "Effects of chemotherapy and hormonal therapy for early breast cancer on recurrence and 15-year survival: an overview of the randomised trials." *The Lancet*, vol. 365, no. 9472, pp. 1687–1717, May 2005, doi: 10.1016/S0140-6736(05)66544-0.
- [2] G. Rahbar et al., "Benign versus Malignant Solid Breast Masses: US Differentiation," *Radiology*, vol. 213, no. 3, pp. 889–894, Dec. 1999, doi: 10.1148/radiology.213.3.r99dc20889.
- [3] P. Skaane and K. Engedal, "Analysis of sonographic features in the differentiation of fibroadenoma and invasive ductal carcinoma.," *American Journal of Roentgenology*, vol. 170, no. 1, pp. 109–114, Jan. 1998, doi: 10.2214/ajr.170.1.9423610.

- [4] R. S. Butler, L. A. Venta, E. L. Wiley, R. L. Ellis, P. J. Dempsey, and E. Rubin, "Sonographic evaluation of infiltrating lobular carcinoma.," *American Journal of Roentgenology*, vol. 172, no. 2, pp. 325–330, Feb. 1999, doi: 10.2214/ajr.172.2.9930776.
- [5] P. M. Lamb, N. M. Perry, S. J. Vinnicombe, and C. A. Wells, "Correlation Between Ultrasound Characteristics, Mammographic Findings and Histological Grade in Patients with Invasive Ductal Carcinoma of the Breast," *Clinical Radiology*, vol. 55, no. 1, pp. 40–44, Jan. 2000, doi: 10.1053/crad.1999.0333.
- [6] American College of Radiology, C. J. D'Orsi, E. A. Sickles, E. B. Mendelson, and E. A. Morris, Eds., *ACR BI-RADS atlas: breast imaging reporting and data system ; mammography, ultrasound, magnetic resonance imaging, follow-up and outcome monitoring, data dictionary*, 5th edition. Reston, Va.: ACR, American College of Radiology, 2013.
- [7] P. Korpraphong, O. Tritanon, W. Tangcharoensathien, T. Angsuthin, and S. Chuthapisith, "Ultrasonographic Characteristics of Mammographically Occult Small Breast Cancer," *J Breast Cancer*, vol. 15, no. 3, p. 344, 2012, doi: 10.4048/jbc.2012.15.3.344.
- [8] S.-C. Chen, Y.-C. Cheung, C.-H. Su, M.-F. Chen, T.-L. Hwang, and S. Hsueh, "Analysis of sonographic features for the differentiation of benign and malignant breast tumors of different sizes," *Ultrasound in Obstet & Gyne*, vol. 23, no. 2, pp. 188–193, Feb. 2004, doi: 10.1002/uog.930.
- [9] H. Y. Kwon, E.-S. Cha, J. E. Lee, J. H. Kim, and J. Chung, "Small Breast Cancer ( $\leq 5$  mm): Ultrasonographic Features and Clinical and Pathological Characteristics," *J Korean Soc Radiol*, vol. 80, no. 4, p. 728, 2019, doi: 10.3348/jksr.2019.80.4.728.
- [10] A. Pawłowska et al., "Curated benchmark dataset for ultrasound based breast lesion analysis," *Sci Data*, vol. 11, no. 1, p. 148, Jan. 2024, doi: 10.1038/s41597-024-02984-z.
- [11] A. Rodríguez-Cristerna, W. Gómez-Flores, and W. C. De Albuquerque-Pereira, "BUSAT: A MATLAB Toolbox for Breast Ultrasound Image Analysis," in *Pattern Recognition*, vol. 10267, J. A. Carrasco-Ochoa, J. F. Martínez-Trinidad, and J. A. Olvera-López, Eds., in *Lecture Notes in Computer Science*, vol. 10267. Cham: Springer International Publishing, 2017, pp. 268–277. doi: 10.1007/978-3-319-59226-8\_26.
- [12] A. J. Potterton, D. J. Peakman, and J. R. Young, "Ultrasound demonstration of small breast cancers detected by mammographic screening," *Clinical Radiology*, vol. 49, no. 11, pp. 808–813, Nov. 1994, doi: 10.1016/S0009-9260(05)81973-7.
- [13] W. Gómez-Flores and J. Hernández-López, "Assessment of the invariance and discriminant power of morphological features under geometric transformations for breast tumor classification," *Computer Methods and Programs in Biomedicine*, vol. 185, p. 105173, Mar. 2020, doi: 10.1016/j.cmpb.2019.105173.
- [14] M. A. Perales-García and W. Gómez-Flores, "Genetic Programming for Feature Construction in Breast Ultrasound Tumor Classification," in *2024 Global Medical Engineering Physics Exchanges/Pan American Health Care Exchanges (GMEPE/PAHCE)*, Mexico City, Mexico: IEEE, Apr. 2024, pp. 1–6. doi: 10.1109/GMEPE/PAHCE62423.2024.10534649.
- [15] F. A. Gonzalez-Luna, J. Hernandez-Lopez, and W. Gomez-Flores, "A Performance Evaluation of Machine Learning Techniques for Breast Ultrasound Classification," in *2019 16th International Conference on Electrical Engineering, Computing Science and Automatic Control (CCE)*, Mexico City, Mexico: IEEE, Sep. 2019, pp. 1–5. doi: 10.1109/ICEEE.2019.8884547.