



# Analysis of directional patterns of lung nodules in computerized tomography using Getis statistics and their accumulated forms as malignancy and benignity indicators

Stelmo Magalhães Barros Netto<sup>a</sup>, Aristófanés Corrêa Silva<sup>a,\*</sup>, Rodolfo Acatauassú Nunes<sup>b</sup>, Marcelo Gattass<sup>c</sup>

<sup>a</sup> Federal University of Maranhão, UFMA Applied Computing Group, NCA/UFMA Av. dos Portugueses, SN, Campus do Bacanga, Bacanga, 65085-580 São Luís, MA, Brazil

<sup>b</sup> State University of Rio de Janeiro, UERJ, São Francisco de Xavier, 524 Maracanã, 20550-900 Rio de Janeiro, RJ, Brazil

<sup>c</sup> Pontifical Catholic University of Rio de Janeiro, PUC-Rio R. São Vicente, 225 Gávea, 22453-900 Rio de Janeiro, RJ, Brazil

## ARTICLE INFO

### Article history:

Received 2 March 2011

Available online 26 May 2012

Communicated by C. Kambhamettu

### Keywords:

Medical image

Computer-aided diagnosis (CADx)

Lung nodules

Getis\* statistics

Image processing

## ABSTRACT

The large incidence of lung cancer in Brazil and around the world, in addition to its difficult diagnosis, especially in the initial stages, has been driving efforts to develop tools that support image-based diagnosis. The main objective is to avoid invasive procedures, which usually pose risks to patients. This work uses Getis spatial autocorrelation statistics, Getis\*, plus its accumulated forms to verify patterns occurring in geographic areas, aiming to indicate the nature of the lung nodule (benign or malignant). Nodule analysis is performed on its volume in a directional way, checking whether there are distances inside the nodule with large intensity variability of the voxels, for malignant and benign nodules. The classification is done by selecting the best four features from the 2400 generated features, for each of the Getis estimates. The Lung Image Database Consortium (LIDC) is used to verify the efficacy of the measures in the diagnosis. Results have shown that all of the Getis estimates succeeded in the discrimination of nodules in LIDC, with accuracy higher than 80% and confirmed by three different classifiers.

© 2012 Elsevier B.V. All rights reserved.

## 1. Introduction

Lung cancer is the commonest type of cancer in the world and one of the most serious public health problems in Europe and North America, according to data from Brazil's National Cancer Institute (NCI) (N. I. of Cancer (INCA) et al., 2010). In Brazil, lung cancer occupies first place in cancer deaths among men and fourth among women. According to INCA (N. I. of Cancer (INCA) et al., 2010), it is estimated that pulmonary neoplasm was responsible for 27,270 deaths (17,810 men and 9460 women) in Brazil in 2008. These figures correspond to an estimated risk of 19 new cases for each 100,000 men and 10 for each 100,000 women. This disease is currently one of the biggest human health concerns, and the use of tobacco is still the main risk factor. The average five-year survival rate ranges between 13% and 21% in developed countries and between 7% and 10% in emergent countries (N. I. of Cancer (INCA) et al., 2010). Much research has been done in order to better understand this disease, aiming to discover its origins, detect it in its early stages, develop methodologies to detect it, and propose more efficient treatments with fewer adverse effects for the pa-

tient. One of the causes of the low survival rate from lung cancer is related to its difficult early diagnosis due to the absence of symptoms, and to poor diagnosis in more advanced stages of the disease (Jamnik et al., 2002). Therefore, significant efforts are being made targeting the early diagnosis of lung cancer. The detection of lung cancer in an initial stage has been improved by a wider use of non-invasive imaging techniques, such as radiography and computerized chest tomography (CT). However, invasive techniques are still necessary for the diagnostic definition, which occurs through the cytological and histopathological study of materials obtained through suction puncture or biopsy. In this scenario, where the application of non-invasive techniques is gaining special relevance, a large number of computer tools have been employed, such as Computer-Aided Detection (CAD) and computer-aided diagnosis (CADx), developed with image processing and computer vision techniques. Using digital images generated in the CT acquisition process, it is possible to identify the lung nodule and perform a series of measurements over it, in order to find a correlation between these measurements and the malignancy or benignity diagnosis (da Silva et al., 2009). In image processing, many techniques have been developed for image-based diagnosis. The main emphasis is still the morphologic evaluation of nodules. Shah et al. (2005) investigated the use of a CADx system to distinguish malignant and benign nodules using volumetric data of the acquired nodules, before and after contrast injection. They classified the nodules

\* Corresponding author. Tel.: +55 98 33018243; fax: +55 98 33018841.

E-mail addresses: [nobnet2000@yahoo.com.br](mailto:nobnet2000@yahoo.com.br) (S.M. Barros Netto), [ari@dee.ufma.br](mailto:ari@dee.ufma.br) (A.C. Silva), [rodolfoacatauassu@yahoo.com.br](mailto:rodolfoacatauassu@yahoo.com.br) (R. Acatauassú Nunes), [mgattass@tecgraf.puc-rio.br](mailto:mgattass@tecgraf.puc-rio.br) (M. Gattass).

using logistic regression and discriminant analysis, with results based on the ROC curve of 0.92 and 0.69, respectively. Gimelfarb et al. (2005), who recently proposed a lung nodule diagnosis method based on the evaluation of the growth rate of the nodules in a certain period using 3D volume registration to quantize it, obtained an accuracy of 100% in a confidence interval of 95%. Sousa et al. (2007) evaluated a set of three geometrical features in order to distinguish between nodules and non-nodules using support vector machine as the classifier. The results were correct classifications in 100% of the cases. Lee et al. (2009) used a random forest classifier to distinguish between nodules and non-nodules and achieved sensitivity of 97.92% and specificity of 96.28%. Chen et al. (2010), using a CAD system based on a neural network ensemble with the LIDC database, achieved 0.79 of area under the ROC curve. Texture analysis has been applied on single lung nodules as a quantitative technique to provide differential information. As an example, (Silva et al., 2004) presented the possibility of representing texture features to distinguish lung nodules using geostatistical functions, such as semivariogram, covariogram, correlogram and semimadogram, and obtained results with 80% accuracy. Liang et al. (2006) proposed a CADx system using neural networks where image registration techniques were used to extract the regions of interest. After that, texture features were obtained using the gray-level spatial dependency method to train and test the neural network. For 20 patients who were analyzed, the results had accuracy of 100% for the 15 training sets and 100% for the 5 test sets. Petkovsa et al. (2006) presented an approach that combines image registration and the extraction of texture features from CT images with iodized contrast. Silva et al. (2007) investigated other geostatistical functions, such as Moran's index and Geary's coefficient, to discriminate lung nodules as either malignant or benign obtaining good results, with accuracy above 90% and sensitivity of 96.55%. In another work, (da Silva et al., 2008) evaluated the use of Ripley's K functions as discriminant nodule measurements and obtained 97.4% of accuracy and 90% of sensitivity. Way et al. (2006) designed an automatic system which segments the lung nodule from the parenchyma region and extracts texture features based on run-length statistics on the image resulting from the rubber-band straightening transform, achieving  $0.83 \pm 0.04$  of area under the ROC curve. This work presents a methodology for recognizing directional patterns of spatial distribution, using the computer as a diagnosis assistance tool. We apply the Getis statistic and its variations to 3D CT images of lung nodules. The main contribution and objective of this work consists in assessing the discriminatory power of this statistic for distinguishing between benign and malignant lung nodules. Getis statistics are local spatial association indexes widely used in spatial studies of geographic areas by means of point pattern analysis – as in Ecology – but which remain underexplored in other areas or with other image types, especially medical ones. This statistic is a measure of non-parametric spatial association, which is intended to measure spatial dependencies. In this work we intend to verify whether the quality and efficiency of these indexes can be applied to analyze lung nodule textures (tissues) and discriminate them as either malignant or benign. The reasons for choosing these indexes for this study are: (1) their structure is easy to adapt to analyze 3D textures; (2) they allow analyzing nodule textures in many directions; (3) they can be analyzed or applied to data whose distribution is non-normal; and (4) they allow measuring spatial clusters, i.e., realizing the same random variable in different places in space, whereas one limitation of distance-based measurements is that they are applicable only to positive observations. The paper is organized as follows. Section 2 presents the methodology, which includes the image acquisition process, the segmentation of the lung nodule, texture extraction by means of 3D analysis of the boundaries with Getis statistics, and the classifiers used to validate

the results. In Section 3 we present the results and evaluate the proposed methodology. Finally, in Section 4, we present final remarks.

## 2. Materials and methods

This section describes the procedures used in the application of the proposed method. We will describe the image acquisition, the background of Getis statistics, the metrics for evaluating the results, and an overview of the classifiers used to discriminate lung nodules as either malignant or benign.

### 2.1. Image acquisition

The database used in this work was created and made available by Brazil's National Cancer Institute (NCI), through a consortium of institutions to develop an agreement and standards for lung nodules obtained through computerized tomography (CT). The Lung Image Database Consortium (LIDC) (Armato et al., 2004) is a group that aims to establish standard lung image formats, management processes, technical reports and clinical data necessary for the development, training and evolution of algorithms to detect and diagnose lung cancer.

Database specification shows a contour for lung nodules larger than 3 mm, while for those smaller than 3 mm only the centroid is given. The contour of the nodule is defined along the slices, as well as its classification, by four specialists. From 84 exams available in this database, only 58 exams have contour information.

For the description of the nodules characteristics there is a scale that allows the specialist some subjectivity with respect to the nodule. According to this scale, in this work we considered as malignant nodules those cases that were described as “highly” and “moderately” suspect of malignancy, and as benign those described as “highly” and “moderately” benign. From the 58 exams that contain contour information, only 198 nodules (99 benign and 99 malignant) met such descriptions and were used to validate our methodology.

### 2.2. Getis statistics

The texture indexes studied in this work are based on Getis statistics, proposed by Zang (2008), which, in their expression for local use, can verify aggregation patterns or anomalies in a more refined manner and easy to interpret.

Expressions are the quotient of a simple sum of the values of the neighborhood of a considered region divided by the sum of the attributes under analysis, representing the degree of intensity of each attribute analyzed in the neighborhood of the examined region. For normalized variables, we can have Eqs. 1 and 2, considering an analysis region  $i$  and the  $n$  neighbor regions  $j$ :

$$G_i(d) = \frac{\sum_{j=1}^n W_{ij}(d)x_j}{\sum_{j=1}^n x_j} \quad (1)$$

$$G_i^*(d) = \frac{\sum_{j=1}^n W_{ij}^*(d)x_j}{\sum_{j=1}^n x_j} \quad (2)$$

where  $W_{ij}(d)$  and  $W_{ij}^*(d)$  are spatial proximity matrices, and their elements have value 1 when the distance between region  $i$  and region  $j$  is smaller than  $d$  (maximum distance), and have value 0 otherwise;  $W_{ij}^*(d)$  includes region  $i$  in the analysis; and  $x_i$  and  $x_j$  are the values of the attributes of regions  $i$  and  $j$ .

Eqs. 1 and 2 differ because they do not include the region or point analysis in the calculations, so the spatial neighborhood matrix is different for both equations.

The presence of a high value in a region under analysis signals the presence of high correlation in the neighborhood in both formulations, while the existence of a low value means low correlations in the neighborhood. This feature allows us to perform the significance test more easily, because both have distribution very close to normality (Zang, 2008). We verify this after normalizing the values of  $G_i(d)$  and  $G_i^*(d)$  with mean zero and unitary deviation, defined by:

$$Z_i = \frac{[X(d) - E(X(d))]}{(var(X(d)))^{1/2}} \quad (3)$$

$$E(X(d)) = \frac{\sum_{i=1}^n X(d)}{(n-1)} \quad (4)$$

$$varX(d) = \sum_{i=1}^n \left( \frac{X(d) - E(X(d))}{(n-1)} \right)^2 \quad (5)$$

where  $E(X(d))$  is the mean and  $var(X(d))$  is the variance of  $X$  that can assume  $G_i$  or  $G_i^*$ .

Visualizing Eqs. 1 and 2, we can derive two more equations, which are accumulative versions of each one along the distance, named  $AG_i$  and  $AG_i^*$ . That is,

$$AX(d) = \sum_{i=0}^d X(i) \quad (6)$$

where  $X(d)$  can assume  $G_i$  or  $G_i^*$ . These were distinguished in order to obtain more discriminating measures for this analysis.

### 2.2.1. Getis statistics applied to 3D analysis

Many spatial autocorrelation indexes used to construct patterns are based on three fundamental concepts: distance, direction and relative position (Taylor et al., 1977). Getis statistics are techniques that depend only on the distance between two regions (or points) of analysis, therefore their obtainment is not limited to objects found either on the same plane or in the same volume, or even in one n-dimensional space. Only a formal definition of distance in the dimensional visualization space is necessary. Direction is measured as an angle, and as such it requires a reference axis.

In order to calculate the Getis index in 3D, it is necessary to determine the vector, which is defined by direction and lag (or distance or  $d$  from Eqs. 1 and 2) in the 3D space. This vector, in a spherical coordinate system, has its direction defined by the azimuth (rotation with axis  $Z$  and axis  $X$  as reference) and elevation (rotation with axis  $X$  and axis  $Z$  as reference). Its distance can be associated to any physical quantity, but in this case it is determined by the separation distance between points or spatial arrangements (lag). Fig. 1(b) shows a vector in a 3D space defined by the azimuth and elevation parameters.

In this work the lags were determined by the maximum diameter found among all lung nodules belonging to our database, in such a way that the maximum lag was estimated. The maximum lag was then equally divided into unitary parts, forming an array with lags varying from 1 to the maximum lag.

In the analysis of the neighborhood, we saw that the formalization of the spatial proximity matrix makes the objects under study contiguous. For the 3D space, we used the vector definition above to determine whether the vector formed by the pair of points (or regions) under study is found in a predefined list of data vectors with certain values of azimuth, elevation and lag, thus constituting the neighborhood for analysis. To increase the possibility for a certain vector to be in this list, tolerances are defined for azimuth, elevation and lag. A 2D representation of these parameters, considering a generalization of azimuth and elevation as angular tolerance, is shown in Fig. 1(a). Thus, the Getis statistics are calculated for all vectors that satisfy a certain direction defined by azimuth and elevation, and a certain lag defined by the Euclidean distance. The distance and the lag are defined prior to the analysis.

### 2.3. Selection of discriminant features

When the number of features increases, it is reasonable to perform a selection of features in the data in order to improve the performance of the classification and rely on more stable data, removing irrelevant or redundant information (Webb et al., 2002).

The selection of features in this work was done using prTools (Duin et al., 2007), a Matlab toolbox which orders the features in

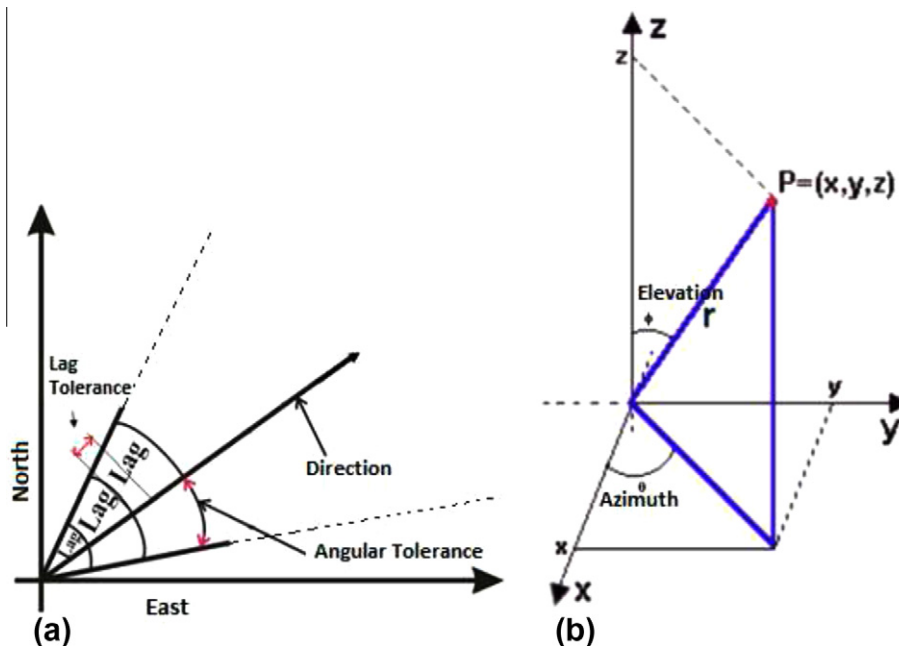


Fig. 1. Parameters used in the 3D neighborhood analysis. (a) 2D representation defined by direction, angular tolerance, lag and lag tolerance parameters. (b) 3D representation defined by azimuth and elevation parameters.

the vector by increasing the error of the classifier, or using another technique. The classifier used for selecting features was the Nearest Mean Classifier (NMC) (Webb et al., 2002), and the selected features were tested in all other classifiers used with Getis statistics.

#### 2.4. Classifiers used to validate the proposed method

The classification stage starts with the evaluation of types of features to be used in the recognition problem. An initial set of parameters is usually made available for this purpose. In the supervised hypothesis, this initial set, represented by  $nd$ -dimensional feature vectors, is used to develop the core, consisting of the training set.

System performance in pattern recognition is evaluated in terms of the error rate for all classes in the entire set. When the evaluation is done with the patterns in the training set, the average evaluation obtained is said to be optimal. In order to achieve better approximations with the pattern recognition system, we used another set of features that were not used during the training stage, to form the core. This new set was called test set, giving us the idea of the extent to which a classifier is able to generalize its capacity to recognize new patterns (Duda and Hart, 1973).

The classifiers used in this work were support vector machine (SVM) (Duda and Hart, 1973) with Gaussian kernel, nearest mean classifier (NMC) (Webb et al., 2002), and the linear classifier based on normal density (LDC) (Webb et al., 2002). Our motivation for choosing these classifiers was that each of them has specific features, and we wanted to observe the behavior of our method in relation to each one of them. The NMC is simple, fast, and the results are given in linear classification boundaries, but it is sensitive to scaling. The LDC is also simple and uses Bayes' rule for classification. Finally, the SVM was used because it is fast and has non-linear kernel, such as RBF, for example. The purpose was to analyze the results of each classifier with respect to our methodology.

#### 2.5. Evaluation of the classification methods

In order to evaluate the power to discriminate nodules through texture patterns, such as proposed by our methodology, we used sensitivity, specificity and accuracy as performance measurements for classifiers used in this study.

Sensitivity is given by  $TP/(TP + FN)$ , specificity is obtained by  $TN/(TN + FP)$ , and accuracy is given by  $(TP + TN)/(TP + TN + FP + FN)$ , where  $TP$  is true positive,  $TN$  is true negative,  $FP$  is false positive and  $FN$  is false negative (Duda and Hart, 1973). Correctly classified lung nodules are said to be true positive.

### 3. Results and discussion

Radiological features of benignity are well known and are based on calcification or fat texture patterns that shift the mean radiological density outside the range of soft tissues. Malignity does not

have such texture criteria, and its diagnosis is normally suggested by an irregular shape associated to certain clinical data, such as the load of tobacco (Swensen, 1997). Therefore, it is not easy to suggest a lung nodule diagnosis using texture features only. Nevertheless, this work investigates a new form of analyzing benignity and malignancy patterns in lung nodules. In this section we will show the application of Getis statistics as texture descriptors to suggest a benignity/malignancy diagnosis.

The methodology was evaluated using the image database described in Section 2.1. In the evaluation of the proposed methodology we used 198 nodules, being 99 malignant and 99 benign.

During the training and testing stages of the Getis statistics, we used the following input parameters based on the background presented in Section 2.2: maximum lag of 200 voxels (diameter of the largest nodule); azimuth angles of  $0^\circ$ ,  $45^\circ$ ,  $90^\circ$  and  $135^\circ$ ; elevation angles of  $0^\circ$ ,  $45^\circ$  and  $90^\circ$ ; and angular tolerance of  $\pm 22.45^\circ$ . The same parameters were used successfully in works like Silva et al. (2007), da Silva et al. (2008), da Silva et al. (2009) and da Silva Souza et al. (2007) to evaluate other geostatistical techniques applied to spatial correlation.

The nodules in our base have voxels with different dimensions, which can affect the performance of the method. So we forced the voxels of the base to have the same dimension, i.e., we made them isotropic. We did this by obtaining the smallest value for  $x$ ,  $y$  and  $z$  among all the voxels in the base and adopting it as the value for voxel isotropy.

The combination of all of these parameters resulted in a set of 2400 features ( $4 \text{ azimuths} \times 3 \text{ elevations} \times 200 \text{ lags}$ ) for each Getis statistic (Eqs. (1)–(6)). This number of features makes the classification computationally complex and costly, so we reduced the number of features using the NMC. The number of features selected for each Getis statistic varied from 4 to 9. However, we used only the first four features selected by the NMC. This strategy was used to make the classifications uniform and easier.

We executed each classifier ten times for different testing and training percentages, obtaining a mean value for sensitivity (Se), specificity (Sp) and accuracy (Ac). Tables 1–4 show the performance of each classifier for the Getis statistic applied, using training and testing percentages from 30% to 70%. The shaded lines in the tables indicate the best sensitivity and specificity percentages.

What we can notice in the results presented in Tables 1–4 is that even by increasing the number of training cases and decreasing the number of tests there was relative stability in the classification performance measurements. This demonstrates that with such features we can obtain accuracy around 80%, even with a reduced number of cases. We can also see in Tables 1–4 that the results given by the SVM classification presented low sensitivity but obtained high specificity, especially when analyzing Getis and Getis\* statistics, and its sensitivity improved with the accumulated forms.

These results show that the indexes  $G_i(d)$  and  $G_i^*(d)$  and their accumulated forms are good malignancy/benignity indicators with directional analysis.

**Table 1**

Results of the classification using Getis to distinguish malignant and benign nodules at various training and testing proportions.

Sample		Classifier								
Train	Test	LDC			NMC			SVM		
		Se	Sp	Ac	Se	Sp	Ac	Se	Sp	Ac
%		%			%			%		
30	70	78.7	76.38	77.54	74.64	82.9	78.77	74.71	77.43	76.07
40	60	78.7	76.38	77.54	74.64	82.9	78.77	70.17	82.17	76.17
50	50	81.03	79.23	80.13	76.92	82.82	79.87	77	81.6	79.3
60	40	75.86	76.9	76.38	70.69	83.45	77.07	74.5	80.25	77.38
70	30	75.86	76.9	76.38	70.69	83.45	77.07	76	81	78.5



**Table 2**

Results of the classification using Getis\* to distinguish malignant and benign nodules at various training and testing proportions.

Sample		Classifier								
Train	Test	LDC			NMC			SVM		
		Se	Sp	Ac	Se	Sp	Ac	Se	Sp	Ac
%		%			%			%		
30	70	73.91	83.91	78.91	69.71	86.67	78.19	78.14	77	77.57
40	60	78.81	77.8	78.31	71.36	84.07	77.71	74	83.67	78.83
50	50	78.57	81.63	80.1	72.24	86.94	79.59	71.4	83.2	77.3
60	40	79.49	80	79.74	71.79	86.92	79.36	71.25	81.5	76.38
70	30	75.17	80	77.59	68.97	85.86	77.41	77.33	83.33	80.33

**Table 3**

Results of the classification using Accumulated Getis statistics to distinguish malignant and benign nodules at various training and testing proportions.

Sample		Classifier								
Train	Test	LDC			NMC			SVM		
		Se	Sp	Ac	Se	Sp	Ac	Se	Sp	Ac
%		%			%			%		
30	70	81.59	76.67	79.13	82.46	79.13	80.8	83.14	73.57	78.36
40	60	82.2	75.93	79.07	83.39	78.81	81.1	83.17	79.33	81.25
50	50	81.84	79.39	80.61	82.24	80.41	81.33	80.2	80.4	80.3
60	40	80.77	79.74	80.26	81.54	80.51	81.03	82	79.25	80.63
70	30	77.93	80	78.97	77.93	81.72	79.83	78	77.33	77.67

**Table 4**

Results of the classification using Accumulated Getis\* statistics to distinguish malignant and benign nodules at various training and testing proportions.

Sample		Classifier								
Train	Test	LDC			NMC			SVM		
		Se	Sp	Ac	Se	Sp	Ac	Se	Sp	Ac
%		%			%			%		
30	70	83.33	74.93	79.13	83.91	78.12	81.01	80.29	77.29	78.79
40	60	80.51	80.68	80.59	81.69	80.85	81.27	82.33	74.67	78.5
50	50	81.43	78.16	79.8	83.27	79.39	81.33	81.40	77.40	79.4
60	40	82.82	80.77	81.79	85.13	79.49	82.31	82.5	76.5	79.5
70	30	82.82	80.77	81.79	85.13	79.49	82.31	82	77.33	79.67

**Table 5**

Comparison of literature on lung nodule diagnosis.

Works	Database	Se (%)	Sp (%)	Ac (%)	Az ROC
da Silva et al. (2008)	Proprietary	100	70	92.3	–
da Silva et al. (2009)	Proprietary	100	86.7	89.7	–
da Silva Sousa et al. (2007)	Proprietary	100	100	100	–
Lo et al. (2003)	Proprietary	–	–	–	0.887
Iwano et al. (2008)	Proprietary	76.9	80	–	–
Vittitoe et al. (1997)	Proprietary	–	–	–	0.82 ± 0.05
Way et al. (2006)	LIDC	–	–	–	0.83 ± 0.04
Chen et al. (2010)	LIDC	–	–	–	0.79
Lee et al. (2009)	LIDC	97.92	96.28	97.10	–
Kawata et al. (2003)	Proprietary	91.4	51.4	77.6	–
Yeh et al. (2008)	Proprietary	91.9	71.5	82.6	0.853
Shiraishi et al. (2006)	JSRT	87.7	66.7	–	0.778
Our method (Getis and SVM)	LIDC	82	79.25	80.63	–
Our method (Getis and NMC)	LIDC	81.54	80.51	81.03	–
Our method (Getis and LDC)	LIDC	82.82	80.77	81.79	–

It is difficult to make reliable comparisons with previously published CADx studies because of the different databases, different evaluation procedures and different optimization parameters used (Nishikawa et al., 1994; Kallergi et al., 1999). However, in terms of sensitivity, specificity and accuracy, the results presented in Tables 1 to 4 can be compared to those achieved by other texture techniques, such as those used by da Silva et al. (2008) and da Silva et al. (2009), and the geometric features used by da Silva Sousa et al. (2007), whose classification performances can be seen in Table 5. Suzuki et al. (2005) managed to discriminate 100% (76/76) of malignant nodules and 48% (200/413) of benign ones by training a neural network called MTANN with a group of voxels obtained by “windowing” pixel to pixel throughout all of the volume. Other results that execute the same task of diagnosing lung nodules can also be seen in Table 5. However, since the authors of these papers did not specify which subsets of the LIDC scans were used in their experiments, a direct comparison with our results is potentially misleading.

As can be seen in the results achieved, the proposed method did not present good performance compared to some of the studies mentioned. Using Getis statistics to analyze the texture of lung cancer images is an innovative proposal, and there are still many aspects to explore. Nonetheless, we have identified some

characteristics of our method that encourage us to continue investing in its development:

- The methodology was tested and evaluated with an image database, LIDC, whose images were acquired using varied protocols. This fact made the task of obtaining texture information that discriminates benign and malignant nodules more complex and challenging. Despite this adversity, Getis statistics proved to be robust, reaching a performance of around 80% in several tests. In much of the literature, a proprietary base is used with a single acquisition protocol, which facilitates the task of the texture analysis techniques.
- Differently from some of the cited works, there was no preprocessing to improve image quality. Our approach intended to assess the potential of the methodology considering the real conditions of the images. In fact, we wanted our method to analyze the images with the same difficulties faced by a specialized medical doctor. In this scenario, once again we observed that Getis statistics were robust, achieving satisfactory results. Nevertheless, we believe that the performance of the method will improve with the application of image enhancement filters.
- Another important point is that the proposed method analyzed only the texture of the images, disregarding a mandatory component in lung nodule diagnosis, namely its geometry. Specialized doctors rely on the texture and geometry of the nodule to complement diagnosis. In the works cited both components are used, so their performance is better than that of our method. Considering the robustness demonstrated by Getis statistics in this study, we believe that including geometry information (such as circularity, roughness, elongation, compactness, eccentricity, radial distance, etc.) will improve the method's performance.

Despite the small size of our sample, which does not allow us to draw generalized considerations, we notice the great potential of using Getis statistics. Our methodology was applied to a database with exams acquired from different CT scans and protocols, and nodules with varied shapes, sizes and tissue compositions. Even so, we achieved good sensitivity, specificity and accuracy rates overall. It is necessary to perform more tests and increase the sample in order to make a deeper and more conclusive assessment of the application of this methodology in the daily medical practice.

#### 4. Conclusion

This paper has shown that Getis statistics are promising indicators of malignancy and benignity in lung nodules when performing directional analysis. However, the validation of our results has shown that the maximum accuracy obtained was of approximately 80%. The number of nodules in our database is too small to draw definitive conclusions, but the preliminary results of this work are very encouraging and demonstrate that applying Getis statistics to three-dimensional data can help discriminate benign from malignant lung nodules in CT images. Based on the results, we believe that such measures provide significant support to a more detailed clinical investigation, and the outcomes were very encouraging when nodules were classified with support vector machine, nearest mean classifier and linear classifier based on normal density. Nevertheless, it is necessary to perform tests with a larger database and more complex cases in order to obtain a more precise behavior pattern. In addition, due to the relatively small size of existing CT lung nodule databases and the various CT imaging acquisition protocols, it is difficult to compare the diagnosis performance between the method presented here and others proposed in the literature, but we showed a comparison with some classic

works to position our method. In spite of the good results obtained only by analyzing the spatial association of textures using Getis statistics on three-dimensional data, further information could be obtained by analyzing point patterns using other rules to determine the plot volume. So, as a future development, we propose investigating these other methods in order to verify the possibility of a more precise and reliable diagnosis.

#### Acknowledgements

The authors acknowledge CAPES, CNPq and FAPEMA for financial support.

#### References

- Armato III, S.G., McLennan, G., McNitt-Gray, M.F., Meyer, C.R., Yankelevitz, D., Aberle, D.R., Henschke, C.I., Hoffman, E.A., Kazerooni, E.A., MacMahon, H., Reeves, A.P., Croft, B.Y., Clarke, L.P., 2004. Lung image database consortium: Developing a resource for the medical imaging research community 1. *Radiology* 232 (3), 739–748. <http://dx.doi.org/10.1148/radiol.2323032035>, <http://radiology.rsna.org/content/232/3/739.abstract>.
- Chen, H., Xu, Y., Ma, Y., Ma, B., 2010. Neural network ensemble-based computer-aided diagnosis for differentiation of lung nodules on CT images: Clinical evaluation. *Acad. Radiol.* 17 (5), 595–602. <http://dx.doi.org/10.1016/j.acra.2009.12.009>, <<http://www.sciencedirect.com/science/article/pii/S10766332090>>, ISSN 1076-6332..
- da Silva, E.C., Silva, A.C., de Paiva, A.C., Nunes, R.A., Gattass, M., 2008. Diagnosis of solitary lung nodules using the local form of Ripley's *K* function applied to three-dimensional CT data. *Comput. Methods Prog. Biomed.* 90 (3), 230–239. <http://dx.doi.org/10.1016/j.cmpb.2008.02.003>, ISSN 0169-2607.
- da Silva, C.A., Silva, A.C., Netto, S.M., de Paiva, B.A.C., Junior, G.B., Nunes, R.A., Lung nodules classification in CT images using simpson's index, Geometrical measures and one-class SVM. In: *MLDM*, 2009, pp. 810–822.
- da Silva Sousa, J., Silva, A., Cardoso de Paiva, A., 2007. Lung structure classification using 3D geometric measurements and SVM. In: *Rueda, L., Mery, D., Kittler, J. (Eds.), Progress in Pattern Recognition, Image Analysis and Applications*, vol. 4756. Springer, Berlin/Heidelberg, pp. 783–792. [http://dx.doi.org/10.1007/978-3-540-76725-1\\_81](http://dx.doi.org/10.1007/978-3-540-76725-1_81), ISBN 978-3-540-76724-4.
- Duda, R.O., Hart, P.E., 1973. *Pattern Classification and Scene Analysis*. Wiley-Interscience Publication, New York.
- Duin, R.P.W., Juszczak, P., Paclik, P., Pekalska, E., Ridder, D., Tax, D.M.J., Verzakov, S., 2007. *PRTools4.1*, a matlab toolbox for pattern recognition, Technical Report, Delft University of Technology.
- Gimelfarb, A.E.-B.G., El-Ghar, R.F.M.A., 2005. Computer aided characterization of the solitary pulmonary nodule using volumetric and contrast enhancement features. *Acad. Radiol.* 12, 1310–1319.
- Iwano, S., Nakamura, T., Kamioka, Y., Ikeda, M., Ishigaki, T., 2008. Computer-aided differentiation of malignant from benign solitary pulmonary nodules imaged by high-resolution CT. *Comput. Med. Imag. Graph.* 32 (5), 416–422. <http://dx.doi.org/10.1016/j.compmedimag.2008.04.001>, ISSN 0895-6111.
- Jamnik, S., Santoro, I.L., Uehara, C., 2002. Comparative study of prognostic factors among longer and shorter survival patients with bronchogenic carcinoma. *Pneumologia* 28 (5), 245–249. <http://dx.doi.org/10.1590/S0102-35862002000500002>, ISSN 0102-3586.
- Kallergi, M., Carney, G.M., Gaviria, J., 1999. Evaluating the performance of detection algorithms in digital mammography. *Med. Phys.* 26, 267–275.
- Kawata, Y., Niki, N., Ohmatsu, H., Moriyama, N., 2003. Example-based assisting approach for pulmonary nodule classification in three-dimensional thoracic computed tomography images. *Acad. Radiol.* 10 (12), 1402–1415. [http://dx.doi.org/10.1016/S1076-6332\(03\)00507-5](http://dx.doi.org/10.1016/S1076-6332(03)00507-5).
- Lee, S.L.A., Kouzani, A.Z., Nasierding, G., Hu, E.J., 2009. Pulmonary nodule classification aided by clustering. In: *Proc. 2009 IEEE Internat. Conf. on Systems, Man and Cybernetics, SMC'09*. IEEE Press, Piscataway, NJ, USA, pp. 906–911, ISBN 978-1-4244-2793-2, <http://dl.acm.org/citation.cfm?id=1732323.1732478>.
- Liang, T.K., Toshiyuki, T., Nakamura, H., Ishizaka, A., 2006. Automatic extraction and diagnosis of lung emphysema from lung CT image using artificial neural network. In: *SICE-ICASE Internat. Conf.*, pp. 2306–2311.
- Lo, S.-C.B., Hsu, L.-Y., Freedman, M.T., Lure, Y.M.F., Zhao, H., 2003. Classification of lung nodules in diagnostic CT: An approach based on 3D vascular features, nodule density distribution, and shape features. *SPIE* 5032, 183–189. <http://dx.doi.org/10.1117/12.481878>, <<http://link.aip.org/link/?PSI/5032/183/1>>..
- N. I. of Cancer (INCA), Estimativa 2010: incidência de câncer no Brasil, accessible in 15/10/2010, 2010.
- Nishikawa, R.M., Giger, M.L., Doi, K., Metz, C.E., Yin, F., Vyborny, C.J., Schmidt, R.A., 1994. Effect of case selection on the performance of computer-aided detection schemes. *Med. Phys.* 21, 265–269.
- Petkovska, I., McNitt-Gray, S.K.S.M.F., Goldin, J.G., Brown, M.S., Kim, H.J., Brown, K., Aberle, D.R., 2006. Pulmonary nodule characterization: A comparison of conventional with quantitative and visual semi-quantitative analyses using contrast enhancement maps. *Euro. J. Radiol.* 59, 244–252.

- Shah, S., McNitt-Gray, M., Rogers, S., Goldin, J., Suh, R., Sayre, J., Petkovska, I., Kim, H., Aberle, D., 2005. Computer aided characterization of the solitary pulmonary nodule using volumetric and contrast enhancement features. *Acad. Radiol.* 12 (10), 1310–1319.
- Shiraishi, J., Abe, H., Li, F., Engelmann, R., MacMahon, H., Doi, K., 2006. *Academic radiology* 13 (8), 995–1003. <http://dx.doi.org/10.1016/j.acra.2006.04.007>.
- Silva, A.C., Carvalho, P.C.P., Gattass, M., 2004. Analysis of spatial variability using geostatistical functions for diagnosis of lung nodule in computerized tomography images. *Pattern Anal. Appl.* 7 (3), 227–234.
- Silva, E.C., Silva, A.C., Paiva, A.C., Nunes, R.A., 2007. Diagnosis of lung nodule using Moran's index and Geary's coefficient in computerized tomography images. *Pattern Anal. Appl.* 11, 89–99.
- Sousa, J.R.F.S., Silva, A.C., Paiva, A.C., 2007. Lung structure classification using 3D geometric measurements and SVM. In: *Progress in Pattern Recognition, Image Analysis and Applications. Lecture Notes in Computer Science*, 4756. Springer, Berlin/Heidelberg, pp. 783–792.
- Suzuki, K., Li, F., Sone, S., Doi, K., 2005. Computer-aided diagnostic scheme for distinction between benign and malignant nodules in thoracic low-dose CT by use of massive training artificial neural network. *IEEE Trans. Med. Imag.* 24 (9), 1138–1150. <http://dx.doi.org/10.1109/TMI.2005.852048>, ISSN 0278-0062.
- Swensen, S.J., 1997. The probability of malignancy in solitary nodules. Application to small radiologically indeterminate nodules. *Arch. Internat. Med.* 8 (157), 849–855.
- Taylor, P.J., 1977. *Quantitative Methods in Geography: An Introduction to Spatial Analysis*. Waveland Press, Illinois.
- Vittitoe, N.F., Baker, J.A., Floyd, C.E., 1997. Fractal texture analysis in computer-aided diagnosis of solitary pulmonary nodules. *Acad. Radiol.* 4 (2), 96–101. [http://dx.doi.org/10.1016/S1076-6332\(97\)80005-0](http://dx.doi.org/10.1016/S1076-6332(97)80005-0), ISSN 1076-6332.
- Way, T.W., Hadjiiski, L.M., Sahiner, B., Chan, H.-P., Cascade, P.N., Kazerooni, E.A., Bogot, N., Zhou, C., 2006. Computer-aided diagnosis of pulmonary nodules on CT scans: Segmentation and classification using 3D active contours. *Med. Phys.* 33 (7), 2323–2337. <http://dx.doi.org/10.1118/1.2207129>, <<http://link.aip.org/link/?MPH/33/2323/1>>.
- Webb, A., 2002. *Statistical Pattern Recognition*. Wiley, Malvern, UK.
- Yeh, C., Lin, C.-L., Wu, M.-T., Yen, C.-W., Wang, J.-F., 2008. A neural network-based diagnostic method for solitary pulmonary nodules. *Neurocomputing* 72(1–3) ISSN 612–624, 0925–2312. <http://dx.doi.org/10.1016/j.neucom.2007.11.009>, *Machine Learning for Signal Processing (MLSP 2006)/ Life System Modelling, Simulation, and Bio-inspired Computing (LSMS 2007)*.
- Zhang, T., 2008. Limiting distribution of the G statistics. *Stat. Probab. Lett.* 78 (12), 1656–1661. <http://dx.doi.org/10.1016/j.spl.2008.01.023>, ISSN 0167-7152.