

Paweł FILIPCZUK<sup>1</sup>, Thomas FEVENS<sup>2</sup>, Adam KRZYŻAK<sup>2</sup>, Andrzej OBUCHOWICZ<sup>1</sup>

## GLCM AND GLRLM BASED TEXTURE FEATURES FOR COMPUTER-AIDED BREAST CANCER DIAGNOSIS

This paper presents 15 texture features based on GLCM (Gray-Level Co-occurrence Matrix) and GLRLM (Gray-Level Run-Length Matrix) to be used in an automatic computer system for breast cancer diagnosis. The task of the system is to distinguish benign from malignant tumors based on analysis of fine needle biopsy microscopic images. The features were tested whether they provide important diagnostic information. For this purpose the authors used a set of 550 real case medical images obtained from 50 patients of the Regional Hospital in Zielona Góra. The nuclei were isolated from other objects in the images using a hybrid segmentation method based on adaptive thresholding and k-means clustering. Described texture features were then extracted and used in the classification procedure. Classification was performed using KNN classifier. Obtained results reaching 90% show that presented features are important and may significantly improve computer-aided breast cancer detection based on FNB images.

### 1. INTRODUCTION

According to the International Agency for Research on Cancer and the National Cancer Registry in Poland, breast cancer is the most common cancer among women. Worldwide, in 2008, there were 1,384,155 diagnosed cases of breast cancer and 458,503 deaths caused by the disease [2, 3]. In 2009, there were 15,752 diagnosed cases in Polish women, 5,242 resulted in death. There has also been an increase in the number of breast cancer cases by 3-4% a year since the 1980s. The effectiveness of treatment largely depends on early detection of the disease. An important and often used diagnostic method is the so-called triple-test. It is based on 3 medical examinations and is used to achieve high confidence in the diagnosis. The triple-test includes self-examination (palpation), mammography or ultrasonography imaging, and FNB (Fine Needle Biopsy) [22]. FNB is an examination that consists in obtaining cytological material directly from the tumor. The collected material is then examined under a microscope to determine the prevalence of cancer cells. This approach requires deep knowledge and experience of the cytologist responsible for the diagnosis. Automatic morphometric diagnosis can so-called make the results objective and assist inexperienced specialists. It also allows screening on a large scale where only uncertain cases would require additional human attention. Along with the development of advanced vision systems and computer science, quantitative cytopathology has become a useful method for detection of diseases, infections as well as many other disorders [8, 20].

The paper presents some work in progress on a fully automatic breast cancer diagnostic system based on analysis of cytological images of FNB material. The task of the system is to mark a case as benign or malignant. In our previous work this was done using morphometric and topological features of nuclei like area or the distribution of the nuclei in the image [4, 5, 12, 13]. Recently we acquired a new database of cytological images of FNB from the Regional Hospital in Zielona Góra, Poland. The images were captured using entirely new technology. They offer much higher level of detail in comparison to the previous database. Pathologists from the hospital identified the distribution of chromatin inside nuclei, visible on the new images, as another important feature which was not included in our previous work. In cancer cells one might often notice distinct lumps of chromatin while in healthy ones the chromatin is usually distributed uniformly. This dependence can be represented by texture features. In this paper we propose 15 texture features based on GLCM (Gray-Level Co-occurrence Matrix) and GLRLM (Gray-

<sup>1</sup> Institute of Control and Computation Engineering, University of Zielona Góra, Zielona Góra, Poland.  
e-mail: {p.filipczuk, a.obuchowicz}@issi.uz.zgora.pl.

<sup>2</sup> Department of Computer Science and Software Engineering, Concordia University, Montreal, Canada.  
e-mail: {fevens, krzyzak}@cs.concordia.ca.

Level Run-Length Matrix) which might be used even on roughly segmented nuclei. In order to distinguish pixels representing the nuclei we used adaptive thresholding and k-means clustering. The entire automatic diagnostic procedure as well as the outcomes of the conducted experiments are described further in the paper. Achieved results are promising and allow looking optimistically to the future of the system.

The paper is divided into 6 sections. Section 1 introduces breast cancer diagnosis. Section 2 presents the process of acquisition of medical images used for testing. Section 3 shows nuclei segmentation. Section 4 describes in detail texture features used for the diagnosis. Section 5 delivers experimental results obtained using the proposed approach. The last part of the work includes conclusions and bibliography.

## 2. ACQUISITION OF MEDICAL IMAGES

The testing database contains 550 images of the cytological material obtained by FNB. The material was collected from 50 patients of the Regional Hospital in Zielona Góra, Poland. Biopsies without aspiration were performed under the control of ultrasonograph with a 0.5 mm diameter needle. Smears from the material were fixed in spray fixative (Cellfix by Shandon) and dyed with hematoxylin and eosin (h+e). The time between preparation of smears and their preservation in fixative never exceeded three seconds. Cytological preparations were then digitalized into virtual slides using Olympus VS120 Virtual Microscopy System. The system consists of a 2/3" CCD camera and 40× objective giving together 0.172 μm/pixel resolution. The average size of the slides is approximately 200,000×100,000 pixels. The scans were prepared using EFI (Extended Focal Imaging). EFI is scanning a preparation several times with the focus plane located at different places at the Z axis. Then the frames are put together in such a way as to keep only the sharp areas of each of them. This allows for extended focal depths impossible to obtain using only optics. Next on each slide a pathologist selected 11 areas which were converted to 8 bit/channel RGB TIFF files of size 1583×828 pixels compressed with lossless LZW algorithm. The number of areas per patient was recommended by the specialists from the hospital [16] and allows for correct diagnosis by a pathologist. The database contains 25 (275 images) benign and 25 (275 images) malignant cases. All cancers were histologically confirmed and all patients with benign disease were either biopsied or followed for a year.

## 3. NUCLEI SEGMENTATION

Classification of tumor malignancy is based on features extracted from nuclei. This requires isolating the nuclei from the background and other objects in the image (e.g., red blood cells). In literature many different approaches have been already proposed to extract cells from microscope images [1, 9, 10]. This task is usually done automatically, using one of the well-known image segmentation techniques [6, 18, 19]. However, reliable nuclei segmentation is a challenging task. FNB images are particularly difficult due to the way they are prepared. The material is taken by a needle and smeared on a slide. This may result in partial destruction of tissue structure, and sometimes even nuclei. The cells are usually not uniformly distributed on the preparation. They often form three-dimensional shapes, and they may be in contact with and/or occluded by other cells. In the presented approach we used automatic segmentation procedure that integrates results of image segmentation from two different methods. The algorithm uses adaptive thresholding segmentation to distinguish all dark objects (nuclei, red blood cells and others) from bright background. Next, nuclei are isolated from other objects using clustering algorithm.

The key idea of thresholding is to separate objects from the background based on pixel intensity fluctuations. Local threshold is calculated for each pixel using intensities of pixels from its neighborhood. This area was defined as a square window of size 75×75 pixels. The threshold is the mean intensity value of pixels inside the window.

The next step of image processing is distinguishing nuclei from the rest of the objects. This task is performed based on color information. It was decided to define three types of objects according to their color: nuclei, red blood cells and the background. The idea of image segmentation using clustering algorithms boils down to a search for clusters of pixels in color space. Derived clusters represent objects

that are characterized by a similar color. In the considered cases, the k-means algorithm was applied to calculate centers of 3 clusters and to determine pixel assignments.

The clustering procedure of k-means algorithm is based on minimizing the within-cluster sum of squared distances for  $K$  clusters:

$$J = \sum_{x=1}^X \sum_{y=1}^Y \sum_{k=1}^K \mu_{x,y,k} D_{x,y,k}^2, \quad (1)$$

where  $X$  and  $Y$  defines the size of the analyzed image,  $\mu_{x,y,k}$  is a function specifying whether  $(x, y)$ -th pixel belongs to the  $k$ -th cluster,  $D_{x,y,k}^2$  is squared Euclidean distance measure:

$$D_{x,y,k}^2 = (c_{x,y} - v_k)^T (c_{x,y} - v_k), \quad (2)$$

where  $c_{x,y} \in \mathfrak{R}^3$  is a vector of the coordinates of the  $(x, y)$ -th pixel in RGB space and  $v_k \in \mathfrak{R}^3$  is a vector of the coordinates of the  $k$ -th cluster center in RGB space. The k-means clustering procedure iteratively changes pixel assignments based on the distance to the nearest mean (cluster center) and updates the cluster centers to match the proper means of the data points they are responsible for. Detailed expressions for iterative updating cluster centers and pixel assignments can be found in the following papers [11, 14, 15].

The cluster representing nuclei is identified by comparing mean values of their coordinates. Since nuclei are the darkest objects the lowest mean value indicates cluster representing the nuclei. The obtained pixel assignment is then used to separate the nuclei from the rest of objects. A problem arises when the clustering algorithm is not able to generate a correct cluster due to relatively small number of pixels representing the nuclei. However, such cases are very rare and most images under consideration are correctly segmented using the described procedure. Finally, all objects smaller than  $2.32 \mu\text{m}^2$  (approximate smallest area of a nucleus) were removed. Sample segmentation results are presented in Fig. 1.

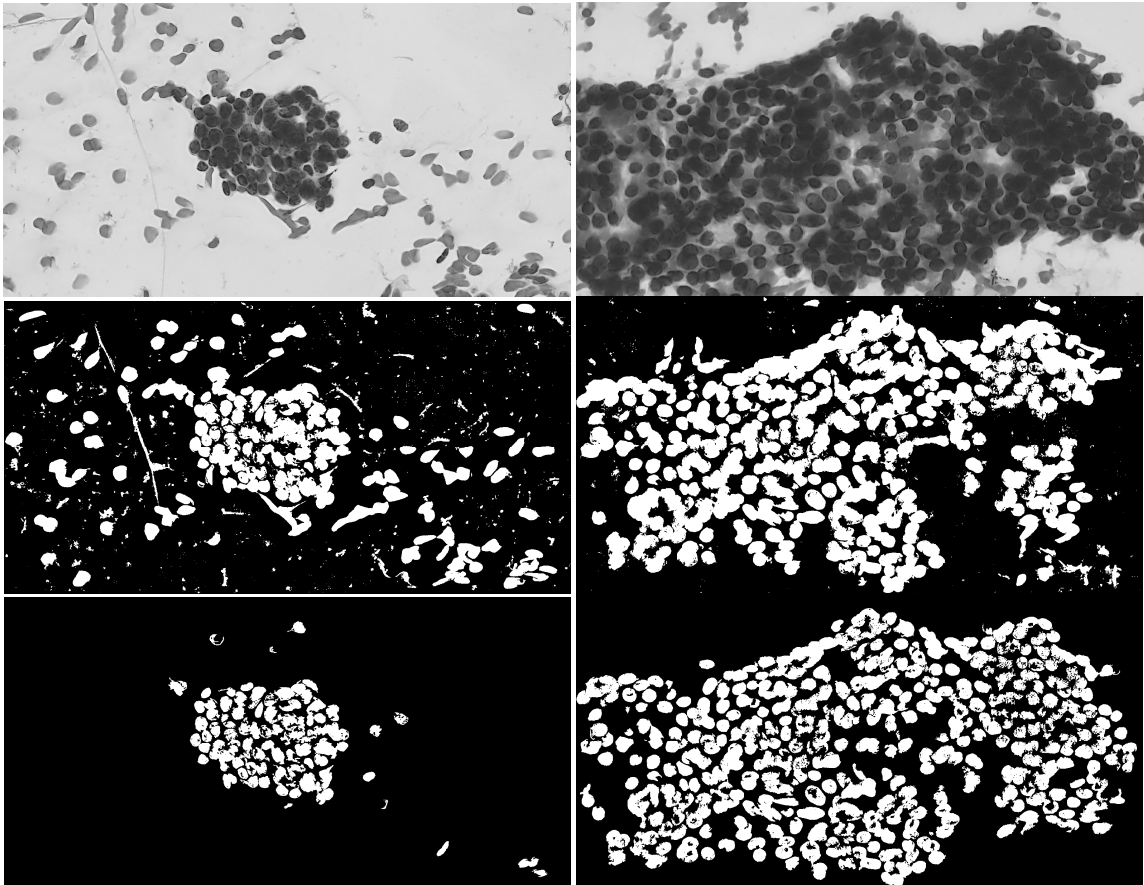


Fig. 1. Original images (top), adaptive thresholding (middle) and final segmentation result (bottom).

## 4. TEXTURE FEATURES

After isolation of nuclei from the images it is possible to extract texture features. For each image we extract the 15 features described below based on GLCM and GLRLM. At the end all the features were standardized.

### 4.1. GRAY-LEVEL CO-OCCURRENCE MATRIX FEATURES

The first four features are based on GLCM. The  $N \times N$  matrix  $P$ , where  $N$  is the number of gray levels, is defined over an image to be the distribution of co-occurring values of pixels at a given offset. In other words each element of  $P$  specifies the number of times a pixel with gray-level value  $i$  occurs shifted by a given distance to a pixel with the value  $j$  [7]. In our case we calculate the mean of four GLCMs determined for offsets corresponding to  $0^\circ$ ,  $45^\circ$ ,  $90^\circ$  and  $135^\circ$  using eight gray-levels.  $p$  is the normalized co-occurrence matrix:

- *contrast* - the intensity contrast between a pixel and its neighbor over the whole image:

$$contrast = \sum_{i,j=1}^N |i - j| p(i, j), \quad (3)$$

- *correlation* - the correlation of a pixel to its neighbor over the whole image:

$$correlation = \sum_{i,j=1}^N \frac{(i - \mu_i)(j - \mu_j) p(i, j)}{\sigma_i \sigma_j}, \quad (4)$$

- *energy* - also known as uniformity, the sum of squared elements in the GLCM:

$$energy = \sum_{i,j=1}^N p(i, j)^2, \quad (5)$$

- *homogeneity* - the closeness of the distribution of elements in the GLCM to the GLCM diagonal:

$$homogeneity = \sum_{i,j=1}^N \frac{p(i, j)}{1 + |i - j|}. \quad (6)$$

### 4.2. GRAY-LEVEL RUN-LENGTH MATRIX FEATURES

The remaining eleven texture features are based on GLRLM. The  $N \times M$  matrix  $p$ , where  $N$  is the number of gray levels and  $M$  is the maximum run length, is defined for a given image as the number of runs with pixels of gray level  $i$  and run length  $j$  [21]. Similarly to the GLCM, we compute run length matrices for  $0^\circ$ ,  $45^\circ$ ,  $90^\circ$  and  $135^\circ$  using eight gray-levels:

- *short run emphasis*:

$$SRE = \frac{1}{n_r} \sum_{i=1}^M \sum_{j=1}^N \frac{p(i, j)}{j^2}, \quad (7)$$

- *long run emphasis*:

$$LRE = \frac{1}{n_r} \sum_{i=1}^M \sum_{j=1}^N p(i, j) j^2, \quad (8)$$

- gray-level nonuniformity:

$$GLN = \frac{1}{n_r} \sum_{i=1}^M \left( \sum_{j=1}^N p(i, j) \right)^2, \quad (9)$$

- run length nonuniformity:

$$RLN = \frac{1}{n_r} \sum_{j=1}^N \left( \sum_{i=1}^M p(i, j) \right)^2, \quad (10)$$

- run percentage:

$$RP = \frac{n_r}{n_p}, \quad (11)$$

where  $n_r$  is the total number of runs and  $n_p$  is the number of pixels in the image,

- low gray-level run emphasis:

$$LGRE = \frac{1}{n_r} \sum_{i=1}^M \sum_{j=1}^N \frac{p(i, j)}{i^2}, \quad (12)$$

- high gray-level run emphasis:

$$HGRE = \frac{1}{n_r} \sum_{i=1}^M \sum_{j=1}^N p(i, j) i^2, \quad (13)$$

- short run low gray-level emphasis:

$$SRLGE = \frac{1}{n_r} \sum_{i=1}^M \sum_{j=1}^N \frac{p(i, j)}{i^2 j^2}, \quad (14)$$

- short run high gray-level emphasis:

$$SRHGE = \frac{1}{n_r} \sum_{i=1}^M \sum_{j=1}^N \frac{p(i, j) i^2}{j^2}, \quad (15)$$

- long run low gray-level emphasis:

$$LRLGE = \frac{1}{n_r} \sum_{i=1}^M \sum_{j=1}^N \frac{p(i, j) j^2}{i^2}, \quad (16)$$

- long run high gray-level emphasis:

$$LRHGE = \frac{1}{n_r} \sum_{i=1}^M \sum_{j=1}^N p(i, j) i^2 j^2. \quad (17)$$

## 5. EXPERIMENTAL INVESTIGATIONS

The features were tested for the classification efficiency, which is defined as a percentage of successfully recognized cases among all cases. For classification we used KNN classifier with value  $k=5$ .

There were 50 patients: 25 benign and 25 malignant. Each patient was represented by 11 images. The effectiveness was tested using the leave-one-out cross-validation technique [17], where a single case was a full set of 11 images representing 1 patient. This means the images belonging to the same patient were never at the same time in the training and testing set. The final diagnosis was obtained by a majority voting of the classification of individual images belonging to the patient (e.g. if 6 images were classified as benign and 5 as malignant then the final diagnosis for the patient would be benign).

The results show that there are four features providing important diagnostic information. They are *run length nonuniformity*, *high gray-level run emphasis*, *short run high gray-level emphasis* and *long run high gray-level emphasis* giving from 68% to 74% efficiency as an individual features (see Table 1). We also performed sequential forward selection to find optimal set of features to check the maximum efficiency using textural features. The best set is presented in Table 1 and gave very good result of 90%.

Table 1. The results of classification using individual texture features and an optimal set of features.

	Feature	Efficiency	Sensitivity	Specificity
GLCM	contrast	50 %	0.44	0.56
	correlation	44 %	0.44	0.44
	energy	56 %	0.60	0.52
	homogeneity	48 %	0.36	0.60
GLRLM	short run emphasis	46 %	0.52	0.40
	long run emphasis	34 %	0.32	0.36
	gray-level nonuniformity	62 %	0.56	0.68
	<b>run length nonuniformity</b>	<b>74 %</b>	<b>0.72</b>	<b>0.76</b>
	run percentage	42 %	0.32	0.52
	low gray-level run emphasis	44 %	0.40	0.48
	<b>high gray-level run emphasis</b>	<b>72 %</b>	<b>0.68</b>	<b>0.76</b>
	short run low gray-level emphasis	38 %	0.44	0.32
	<b>short run high gray-level emphasis</b>	<b>68 %</b>	<b>0.64</b>	<b>0.72</b>
	long run low gray-level emphasis	42 %	0.56	0.28
	<b>long run high gray-level emphasis</b>	<b>70 %</b>	<b>0.68</b>	<b>0.72</b>
Optimal set of features	<b>run length nonuniformity, high gray-level run emphasis, short run high gray-level emphasis</b>	<b>90 %</b>	<b>0.84</b>	<b>0.96</b>

## 6. CONCLUSIONS

The aim of the work was to test whether texture features might provide essential diagnostic information in automatic breast cancer diagnosis based on analysis of FNB images. To perform the experiment we used 550 cytological images from patients of the Regional Hospital in Zielona Góra. In order to segment nuclei we used a hybrid method based on adaptive thresholding and k-means clustering. We tested 15 GLCM and GLRLM texture features. The results of classification showed that not all of them are valuable in diagnostic process. However, some of them deliver important information giving up to 74% efficiency used individually. An optimal combination of features determined by sequential forward selection gave 90%, which is very promising result. This shows texture features are important in application to breast cancer detection and combined with morphometric and topological features may significantly improve computer-aided diagnosis.

## 7. ACKNOWLEDGMENT

One of the co-authors is a scholar within Sub-measure 8.2.2: Regional Innovation Strategies, Measure 8.2: Transfer of knowledge, Priority VIII: within the Regional human resources for the economy Human Capital Operational Programme co-financed by the European Social Fund and the state budget.



HUMAN CAPITAL  
HUMAN – BEST INVESTMENT!



Lubuskie  
Worth your while

EUROPEAN UNION  
EUROPEAN  
SOCIAL FUND



This research was partially supported by the National Science Centre in Poland and by the Natural Sciences and Engineering Research Council of Canada.

BIBLIOGRAPHY

- [1] AL-KOFAHI Y, LASSOUED W, LEE W, ROYSAM B., Improved Automatic Detection and Segmentation of Cell Nuclei in Histopathology Images, *IEEE Trans. on Biomedical Engineering*, 2010, Vol. 57, No. 4, pp. 841–852.
- [2] BRAY F., REN J., MASUYER E., FERLAY J., Estimates of global cancer prevalence for 27 sites in the adult population in 2008, *Int. J. Cancer*, DOI: 10.1002/ijc.27711, 2012.
- [3] FERLAY J., SHIN H., BRAY F., FORMAN D., MATHERS C., PARKIN D., Globocan 2008 v2.0, cancer incidence and mortality worldwide: Iarc cancerbase no. 10., International Agency for Research on Cancer, Lyon, France, 2010, online: <http://globocan.iarc.fr> (accessed on 30/08/2012).
- [4] FILIPCZUK P., KOWAL M., OBUCHOWICZ A., Automatic breast cancer diagnosis based on k-means clustering and adaptive thresholding hybrid segmentation, *Image processing and communications challenges 3, Advances in Intelligent and Soft Computing*, 2011, Vol. 102, pp. 295–303.
- [5] FILIPCZUK P., KOWAL M., OBUCHOWICZ A., fuzzy clustering and adaptive thresholding based segmentation method for breast cancer diagnosis, *Computer recognition systems 4, Advances in Intelligent and Soft Computing*, 2011, Vol. 95, pp. 613–622.
- [6] GONZALEZ R.C., WOODS R.E., *Digital Image Processing*, 3<sup>rd</sup> ed., Prentice Hall, New Jersey, 2008.
- [7] HARALICK R., SHANMUGAM K., DINSTEN I., Textural features for image classification, *IEEE Trans. on Systems, Man and Cybernetics*, 1973, Vol. 3, No. 6, pp. 610–621.
- [8] HASSAN M.R., HOSSAIN M.M., BEGG R.K., RAMAMOHANARAO K., MORSI Y., Breast-cancer identification using hmm-fuzzy approach, *Computers in Biology and Medicine*, 2010, Vol. 40, pp. 240–251.
- [9] HREBIEŃ M., STEĆ P., OBUCHOWICZ A., NIECZKOWSKI T., Segmentation of breast cancer fine needle biopsy cytological images, *Int. J. Appl. Math and Comp. Sci.*, 2008, Vol. 18, No. 2, pp. 159–170.
- [10] JELEŃ, Ł. FEVENS, T., KRZYŻAK A., Classification of breast cancer malignancy using cytological images of fine needle aspiration biopsies, *Int. J. Appl. Math and Comp. Sci.*, 2008, Vol. 18, No. 1, pp. 75–83.
- [11] KANUNGO T., MOUNT D.M., NETANYAHU N.S., PIATKO C.D., SILVERMAN R., WU A.Y., An efficient k-means clustering algorithm: Analysis and implementation, *IEEE Trans. Pattern Analysis and Machine Intelligence*, 2002, Vol. 24, pp. 881–892.
- [12] KOWAL M., FILIPCZUK P., OBUCHOWICZ A., KORBICZ J., Computer-aided diagnosis of breast cancer using Gaussian mixture cytological image segmentation, *Journal of Medical Informatics & Technologies*, 2011, Vol. 17, pp. 257–262.
- [13] KOWAL M., FILIPCZUK P., KORBICZ J., Hybrid cytological image segmentation method based on competitive neural network and adaptive thresholding, *Pomiary, Automatyka, Kontrola*, 2011, Vol. 57, No. 11, pp. 1448–1451.
- [14] LLOYD S.P., Least squares quantization in PCM, *IEEE Trans. Information Theory*, 1982, Vol. 28, No. 2, pp. 129–137.
- [15] MACKAY D., *Information Theory, Inference and Learning Algorithms*, Cambridge University Press, 2003.
- [16] MARCINIAK A., OBUCHOWICZ A., MONCZAK R., KOŁODZIŃSKI M., Cytomorphometry of Fine Needle Biopsy Material from the Breast Cancer, *Proc. 4<sup>th</sup> Int. Conf. on Computer Recognition Systems CORES' 05*, Springer, 2005, pp. 603–609.
- [17] MITCHELL T.M., *Machine Learning*, McGraw-Hill, 1997.
- [18] NAZ S., MAJEED H., IRSHAD H., Image segmentation using fuzzy clustering: A survey, *Proc. 6<sup>th</sup> Int. Conf. Emerging Technologies, ICET 2010*, 2010, pp. 181–186.
- [19] SURI J.S., SETAREHDAN K., SINGH S., *Advanced Algorithmic Approaches to Medical Image Segmentation*, Springer-Verlag, London, 2002.
- [20] ŚMIETANSKI J., TADEUSIEWICZ R., ŁUCZYŃSKA E., Texture Analysis in Perfusion Images of Prostate Cancer - a Case Study, *Int. J. Appl. Math and Comp. Sci.*, 2000, Vol. 20, No. 1, pp. 149–156.
- [21] TANG X., Texture information in run-length matrices, *IEEE Trans. On Image Processing*, 1998, Vol. 7, No. 11, pp. 1602–1609.
- [22] UNDERWOOD J.C.E., *Introduction to biopsy interpretation and surgical pathology*, Springer-Verlag, London, 1987.

