

An Image Analysis Algorithm for Soil Structure Identification

Małgorzata Charytanowicz^{1,2} and Piotr Kulczycki^{1,3}

¹ Polish Academy of Sciences, Systems Research Institute,
Centre of Information Technology for Data Analysis Methods,
Newelska 6, 01-447 Warsaw, Poland

² Catholic University of Lublin, Institute of Mathematics and Computer Science,
Konstantynow 1H, 20-708 Lublin, Poland

Malgorzata.Charytanowicz@ibspan.waw.pl

³ Cracow University of Technology, Department of Automatic Control and Information
Technology, Warszawska 24, 31-155 Cracow, Poland

Piotr.Kulczycki@ibspan.waw.pl

Abstract. Pore space study has been utilized as a general method for defining soil structures. This is because the characteristics particular to pore space impact the majority of physical and physicochemical soil parameters relevant due to plant growth. This paper presents an image segmentation approach for detecting the soil pore structures that have been studied by way of soil tomography sections. In so-doing, a research study was conducted using a density-based clustering method, and in turn, the nonparametric kernel estimation methodology. This overcomes the rigidity of arbitrary assumptions concerning the number or shape of clusters among data, and lets the researcher detect inherent data structures. After a short description of the method, the practical aspects and applications illustrated with a number of soil aggregates are presented.

Keywords: image processing and analysis, image segmentation, clustering, natural grouping, nonparametric estimation, kernel estimators, pore space, total porosity.

1 Introduction

Among all measurements characterizing the various aspects of a particular soil, the total porosity provides a more useful physical description that is relevant to plant growth. This measure is defined as the fraction of the total pore volume of soil material, including the solid and void components, that is taken up by the volume of void-space. Being simply a fraction of total volume, it can range between 0 and 1, typically falling between 0.3 and 0.7 for most soils. Moreover, a number of scientists have reported that studies of pore size distribution are useful as a general method for defining the soil structure [6, 17, 22]. Pore spaces location within the soil has different influences on fluid retention, conduction within the soil and the maximum space available for water, and, therefore, it is important to identify pore zones.

Advances in X-ray microtomography, together with image processing methods, provide a non-destructive alternative for detecting soil structures, especially pore space [15, 16]. However, the usefulness of computed tomography data for pore structure characterization depends on the accuracy of the grayscale images segmentation into binary pore and solid components. Different methods have been used to segment soil images, among these, simple binary and multiple thresholding [14], watershed, morphological, and normalized cut [3, 18, 21]. Furthermore a fuzzy approach successfully used in various data analysis problems [1, 10, 12, 13] is discovered for the segmentation methods [5]. These methods have given promising results, but they are very sensitive to image quality, low level of contrast and unintended noises. Based on these arguments, many studies now have begun to utilize the potential of image segmentation done by clustering methods [7, 11]. The objective of clustering process is to find pixel groups of a similar grey level intensity so as to organize them into more or less homogeneous groups and assign the same label to every pixel sharing certain visual characteristics. As a traditional clustering algorithm, *K*-means is popular for its simplicity in implementation, and it is commonly applied for grouping pixels in images. However, the quality of *K*-means suffers from being confined to being run with a fixed number of clusters. Therefore, many current research efforts have been focused on discovering and applying new approaches in segmenting by way of using various now available image processing techniques.

The main purpose of this investigation is to evolve a standard method of detecting pore space in the soil. A proposed methodology integrates image processing and clustering technique based on the Complete Gradient Clustering Algorithm [9, 10]. The principle of the proposed algorithm is based on the distribution of the data and the need to estimate its density. Within the algorithm, each cluster is identified by a local maximum of the kernel density estimator of the data distribution. As a result, regions of high densities of objects are recognized as clusters, while areas with sparse distributions of objects divide one group from another. The algorithm works in an iterative manner until a termination criterion has been satisfied. Data points are assigned to clusters by using an ascending gradient method, i.e. points moving to the same local maximum are put into the same cluster. It is worth underlining that the whole procedure does not need the application of any assumptions concerning the data distribution or fixed number of clusters. Rather, the parameter values are calculated using optimizing criteria, without any necessity of their arbitrary specification. However, by an appropriate change in values of these parameters, it is possible to influence the size of number of clusters, and also the proportion of their appearance in dense areas in relation to sparse regions of elements in a data set. As a result, the proposed procedure allows researchers to examine soil structure and extract pore space using the segmented images. Moreover, a comparison between the clustering results obtained from this method and the classical *K*-means clustering algorithm shows positive practical features of the Complete Gradient Clustering Algorithm.

2 Statistical Kernel Estimators

Let (Ω, Σ, P) be a probability space. Let also a real random variable $X : \Omega \rightarrow R$, whose distribution has the density function f , be given. The corresponding kernel estimator $\hat{f} : R \rightarrow [0, \infty)$, calculated using experimentally obtained values for the m -element random sample x_1, x_2, \dots, x_m , in its basic form is defined by

$$\hat{f}(x) = \frac{1}{mh} \sum_{i=1}^m K\left(\frac{x-x_i}{h}\right), \tag{1}$$

where $m \in N \setminus \{0\}$, the positive coefficient h is known as a smoothing parameter, whereas the measurable function $K : R \rightarrow [0, \infty)$ of unit integral, symmetrical with respect to zero and having a weak global maximum at this point, takes the name of a kernel. The influence of the smoothing parameter on particular kernels is the same for the basic definition of kernel estimator (1). Advantageous results are obtained thanks to the individualization of this effect, achieved through a so-called modification of the smoothing parameter. It relies on mapping the positive modifying parameters s_1, s_2, \dots, s_m on particular kernels, described as

$$s_i = \left(\frac{\hat{f}(x_i)}{\bar{s}}\right)^{-c}, \tag{2}$$

where $c \in [0, \infty)$, \hat{f} denotes the kernel estimator without modification, and \bar{s} is the geometrical mean of the numbers $\hat{f}(x_1), \hat{f}(x_2), \dots, \hat{f}(x_m)$. The parameter c stands for the intensity of the modification procedure and based on indications for the criterion of the integrated mean square error, the standard value $c = 0.5$ can be suggested. Finally, the kernel estimator with the smoothing parameter modification is defined in the following formula:

$$\hat{f}(x) = \frac{1}{mh} \sum_{i=1}^m \frac{1}{s_i} K\left(\frac{x-x_i}{hs_i}\right). \tag{3}$$

The choice of the kernel K form and the calculation of the smoothing parameter h is made most often with the criterion of the mean integrated square error. From a statistical point of view, the choice of the kernel form has no practical meaning and thanks to this, it becomes possible to take into account primarily properties of the estimator obtained or calculation aspects, advantageous from the viewpoint of the application problem under investigation. The standard normal kernel given by

$$K(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \tag{4}$$

is used most often. It is differentiable up to any order and assumes positive values in the whole domain.

The fixing of the smoothing parameter h has significant meaning for the quality of estimation. A smoothing parameter controls the tradeoff between bias and variance in the result. A large bandwidth leads to a very smooth density distribution, whereas a small bandwidth leads to an ragged density distribution. A frequently used bandwidth selection technique, called the “cross-validation method”, chooses h to minimize the function $g : R \rightarrow R$ defined as

$$g(h) = \frac{1}{m^2 h} \sum_{i=1}^m \sum_{j=1}^m \tilde{K}\left(\frac{x_j - x_i}{h}\right) + \frac{2}{mh} K(0) \quad , \quad (5)$$

where $\tilde{K}(x) = K^{*2}(x) - 2K(x)$, whilst K^{*2} denotes convolution function of K , i.e.

$$K^{*2}(x) = \int_R K(u)K(x-u) du \quad , \quad (6)$$

and for the standard normal kernel (4) is equal

$$K^{*2}(x) = \frac{1}{2\sqrt{\pi}} e^{-\frac{x^2}{4}} \quad . \quad (7)$$

The tasks concerning the choice of the kernel form, as well as additional procedures improving the quality of the estimator obtained, and all rules needed for calculating the smoothing parameter, are found in [8, 19, 20]. The utility of kernel estimation has been investigated in the context of the Complete Gradient Clustering Algorithm.

3 Complete Gradient Clustering Algorithm

Consider the data set containing m elements x_1, x_2, \dots, x_m , in n -dimensional space. Using the methodology introduced in Section 2, the kernel density estimator (3) may be constructed in n -dimensional space, i.e.:

$$\hat{f}(x) = \frac{1}{mh^n} \sum_{i=1}^m \frac{1}{s_i^n} K\left(\frac{x - x_i}{hs_i}\right) \quad , \quad (8)$$

where the kernel K is assumed to be radial or product [8, 19, 20]. The idea of the algorithm is based on the approach proposed by Fukunaga and Hostetler [4]. Thus given the start points:

$$x_j^0 = x_j \quad \text{for } j = 1, 2, \dots, m \quad (9)$$

each point is moved in an uphill gradient direction using the following iterative formula:

$$x_j^{k+1} = x_j^k + b \frac{\nabla \hat{f}(x_j^k)}{\hat{f}(x_j^k)} \quad \text{for } j = 1, 2, \dots, m \quad \text{and } k = 0, 1, \dots, k^*, \quad (10)$$

where $\nabla \hat{f}$ denotes the gradient of kernel estimator \hat{f} and the value of parameter b is proposed as $h^2 / (n + 2)$ while the coefficient h is the bandwidth of \hat{f} . The algorithm will be stopped when the following condition is fulfilled:

$$|D_k - D_{k-1}| \leq aD_0, \quad (11)$$

where D_0 and D_{k-1} , D_k denote sums of Euclidean distances between particular elements of the set x_1, x_2, \dots, x_m before starting the algorithm as well as after the $(k - 1)$ -th and k -th step, respectively. The positive parameter a is taken arbitrary and the value 0.001 is primarily recommended. This k -th step is the last one and will be denoted hereinafter by k^* where $k^* \in N \setminus \{0\}$.

Finally, after the k^* -th step of the algorithm (9)-(10) the set

$$x_1^{k^*}, x_2^{k^*}, \dots, x_m^{k^*}, \quad (12)$$

considered as the new representation of all points x_1, x_2, \dots, x_m , is obtained. Following this, the set of mutual Euclidean distances of the above elements:

$$\left\{ d(x_i^{k^*}, x_j^{k^*}) \right\}_{\substack{i=1, 2, \dots, m-1 \\ j=i+1, i+2, \dots, m}} \quad (13)$$

is defined. Using the methodology presented in Section 2, the auxiliary kernel estimator \hat{f}_d of the elements of the set (13), treated as a sample of a one-dimensional random variable, is created under the assumption of nonnegative support. Next, the first (i.e. obtained for the smallest value of an argument) local minimum of the function \hat{f}_d belonging to the interval $(0, D]$, where D means the maximum value of the set (13), is found. This local minimum will be denoted as x_d , and it can be interpreted as the half-distance between potential closest clusters. Finally, the clusters are created. First, the element of the set (13) is taken; it initially create a one-element cluster containing it. An element of the set (13) is added to the cluster if the distance between it and any element belonging to the cluster is less than x_d . Every added element is removed from the set (13). If there are no more elements belonging to the cluster, the new cluster is created. The procedure of assigning elements to clusters is repeated as long as the set (13) is not empty.

Procedures described above constitute the Complete Gradient Algorithm in its basic form. The values of the parameters used are calculated automatically. However, by an appropriate change in values of these parameters it is possible to influence the size of number of clusters, and also the proportion of their appearance in dense areas

in relation to sparse regions of elements in a data set. Too small a value of the smoothing parameter h results in the appearance of too many local extremes of the kernel estimator, and as a consequence, an increase in the number of clusters. On the other hand too great a value causes its excessive smoothing and an decrease in the number of clusters.

Next, the intensity of modification of the smoothing parameter is implied by the value of the parameter c . Its increase smoothes the kernel estimator in areas where elements of data set are sparse, and also it sharpens it in dense areas. In consequence, if the value of the parameter c is raised, then the number of clusters in sparse areas of data decreases, while at the same time, increases in dense regions. Inverse effects can be seen in the case of lowering this parameter value.

Detailed information on the CGCA procedures and their influences on the clustering results as well as applicational examples are described in the articles [2, 9, 10].

4 Methodology

The proposed methodology to elaborate an innovative image processing approach for detection pore space, based on computed tomography and the nonparametric kernel estimation methodology, is summarized as follows:

1. preparing the soil sample;
2. capturing the soil tomographic slices;
3. applying the contrast enhancement technique on the original soil images;
4. extracting the color components from the enhanced image;
5. applying the unsupervised segmentation technique that is based on the complete gradient clustering algorithm;
6. detecting the pore space from the segmented images.

4.1 Soil Classification

The investigated material was sampled from the cultivated soil layer, classified as silty loam (WRB Mollic Gleysols), explored at the Institute of Agrophysics, of the Polish Academy of Sciences in Lublin. The proportion of each particle size group in the soil was as follows: sand – 46%, silt – 28%, clay – 26%. Furthermore, the pH was: H₂O – 5.9, while KCl was 5.4.

On the experimental fields, a long-term fertilization trial had been executed. The adopted crop rotation from 1955 to 1989 was a cycle of potato – barley – rye, and from 1990 – a cycle of sugar beat – barley – rape – wheat. Three treatments concerning fertilization: control group – plant residues only, mineral fertilization – according to plant needs, and pig manure – 80 ton per ha, were studied. The aggregate soil organic matter was measured by the Multi N/C 3100 Autoanalyser (Analytic Jena, Germany).

Table 1. Aggregate soil organic matter measurements

Type of fertilization	Total organic carbon content [g/kg]	Total nitrogen content [g/kg]
Control (without fertilization)	13.54	1.35
Mineral fertilization	14.89	1.51
Pig manure fertilization	21.50	2.10

The total organic carbon shows the same tendency as total nitrogen, i.e. increasing in the same order: the lowest – control, middle – mineral fertilization, the highest – pig manure.

4.2 Soil Sample Preparation and Image Processing

The soil samples were air dried in room conditions, divided into smaller amounts, and gently sieved through 2 and 10 mm sieves. Soil aggregates remaining at 2 mm sieve and ranging from 2 to 10 mm, were then detected by means of X-ray computational tomography, using a GE Nanotom S device, with the voxel-resolution of 2.5 microns per volume pixel. Three 2D sections uniformly located within each aggregate were performed to characterize the aggregate structure. Next, tomography sections were processed using the Aphelion 4.0.1 package. In the initial step, the contrast enhancement technique was applied on the original soil images, and, subsequently, a rectangle ROI (region of interest) selection of the size of 128x128 pixels was performed upon the enhanced grayscale image. Thus, the ring artifacts of the original images were removed, and these ROI's were saved as a bitmap format. The color components data derived automatically from these images were then examined, as the Complete Gradient Clustering Algorithm (CGCA) allowed for soil image segmentation. In order to find a distribution density of the color components, the kernel estimators methodology, presented in Section 2 was used, with the application of the normal kernel, the cross-validation method, as well as the smoothing parameter modification procedure with standard intensity. Moreover, in the clustering algorithm, a modification of parameters values was employed to eliminate peripheral clusters. Finally, pore space detection was done automatically using the segmented images through choosing the cluster containing the lowest color components.

4.3 Image Segmentation Results

In order to assess the proposed segmentation method, three sections of each soil aggregate were captured from soil samples differing in term of fertilization. The color components that had been extracted from the grayscale images were subsequently fed as input to the Complete Gradient Clustering Algorithm for further segmentation processing. After that, the cluster of the lowest values corresponding to the pore space in

each sample was detected. What is more, all pixels of these values that had been distinguished in the original images in black, were captured. Fig. 1-3 show the 8 bit grayscale images of the captured soil samples and the corresponding resultant images with pore space shown in black. These images were subsequently composed in the table rows from the lowest to the highest sections, as cut with proportions 25%, 50%, and 75% of the aggregate height. After the pore space detection, a common quantitative analysis was conducted in order to assess the overall performance of the results obtained.

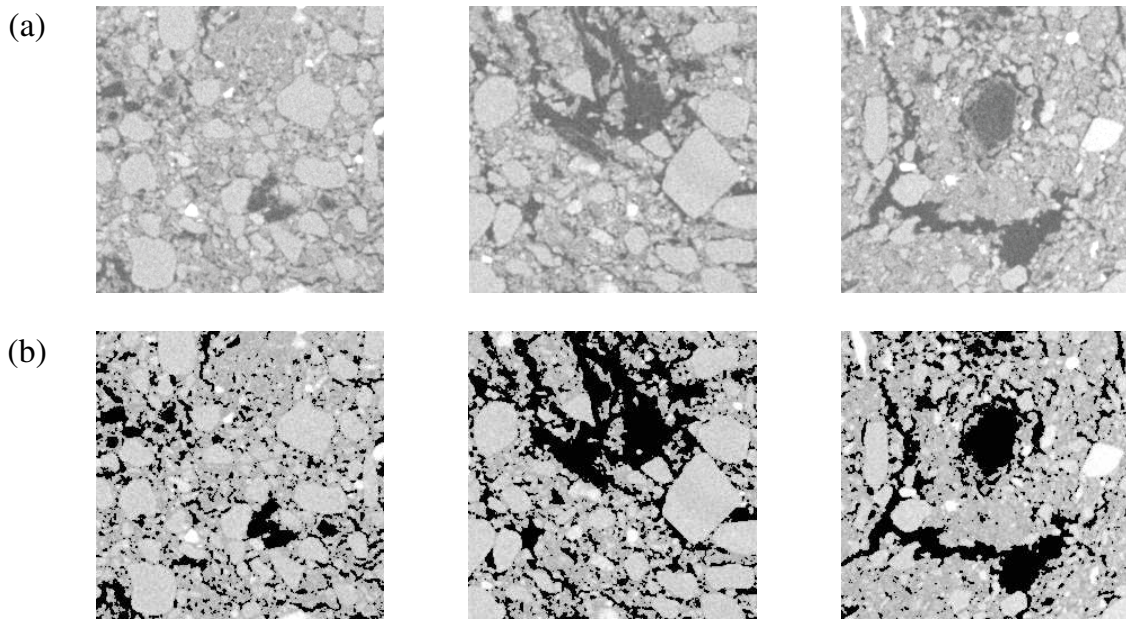


Fig. 1. The rectangle ROI selections of control group aggregates: original images (a), images with pore space in black detected by the CGCA (b)

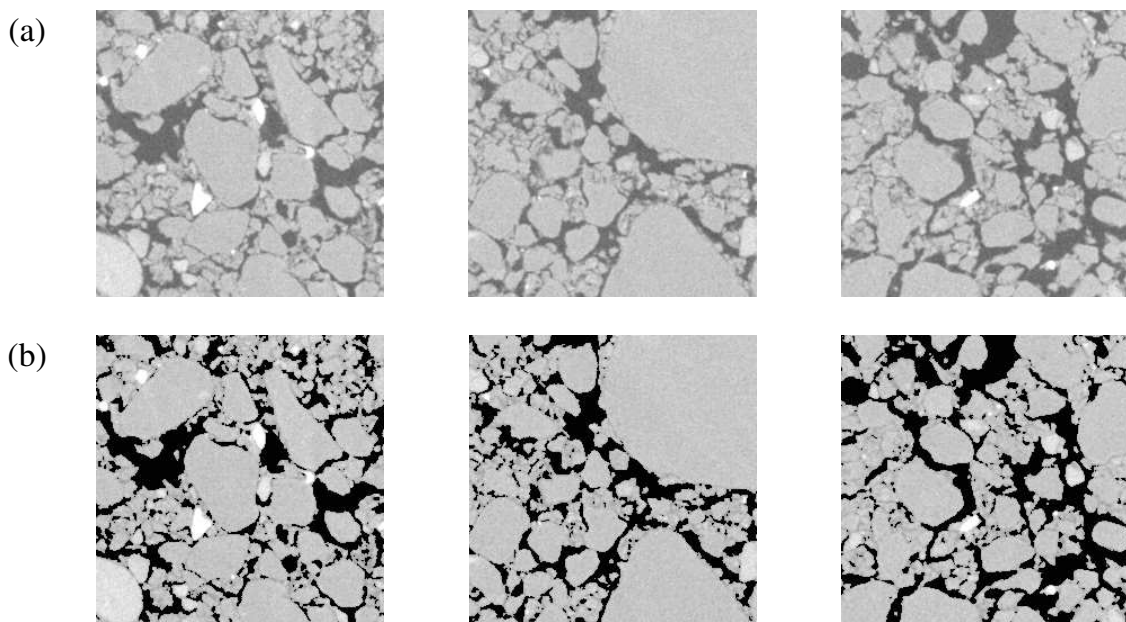


Fig. 2. The rectangle ROI selections of mineral fertilization aggregates: original images (a), images with pore space in black detected by the CGCA (b)

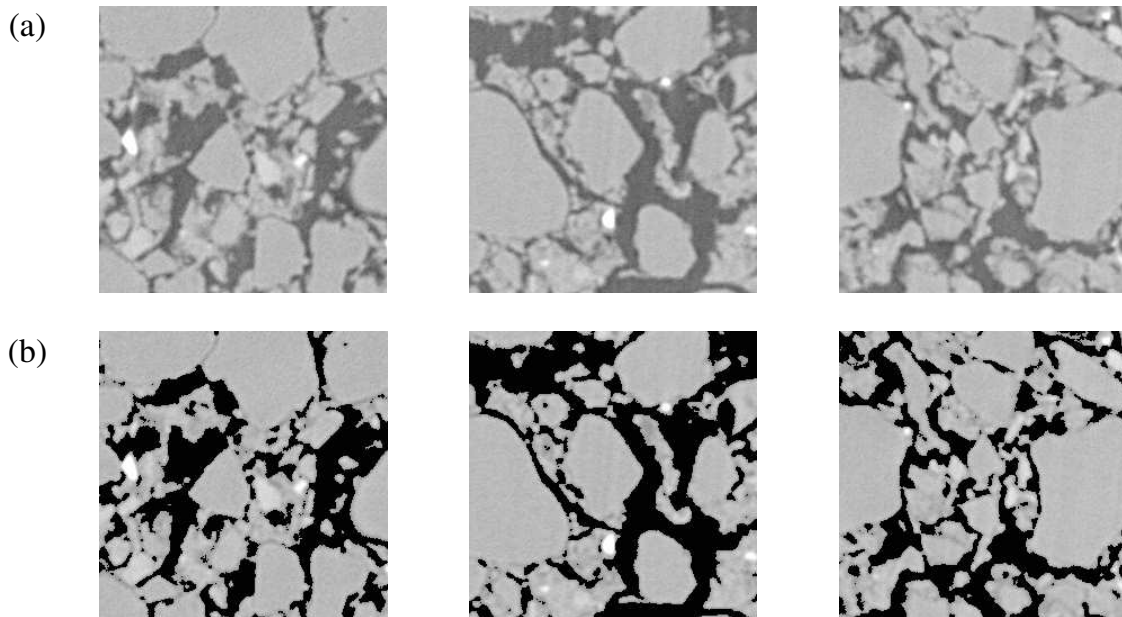


Fig. 3. The rectangle ROI selections of pig manure fertilization aggregates: original images (a), images with pore space in black detected by the CGCA (b)

Based on the observation of several soil images, it was found that the appearance of the pore space cells is similar within the aggregate of each treatment. Indeed, the upper limits of pixel values classified as pore space within the aggregate are almost the same and are nearly equal: 162, 165, and 163 for soil without fertilization (control group), 155, 156, and 155 for soil with mineral fertilization and 149, 148, and 151 for soil with pig manure fertilization. For each type of fertilization, values are ordered correspondingly to the image order as displayed on the Fig. 1-3.

As revealed, a diversity in pore space appearance is found between aggregates differing in term of fertilization. The largest upper limits of pixel values classified as pore space occurs in the soil without fertilization, and, despite this, the pore fraction, equaling 21% in an averaging of the three sections, is the smallest. The soil with pig manure fertilization incorporates the smallest values of the upper limits of these pixels, but despite this, has the largest fraction of pores, equaling 28% in an averaging of the three sections. Generally, the effect of fertilization is to increase the amounts of pores and its size in relation to the control group. The greater increase is for pig manure fertilization.

A comparable analysis of segmentation was obtained when the *K*-means algorithm with an arbitrary taken cluster number of two was used. Table 2 contains pixel value limits between the pore space and solid components, as calculated by the CGCA and the *K*-means algorithm for each type of fertilization.

This study has shown the adequacy of using nonparametric kernel estimation theory for determining soil structure. The limits obtained by the *K*-means algorithm are a bit greater than these obtained by the CGCA, and when used in the segmentation process, the *K*-means algorithm gives an overestimation of the pore space. Furthermore it is worth stressing, that this algorithm needs an a priori assumed correct number of clusters, which in many applications, may not be known. Indeed, even such a “correct” (from a theoretical point of view) number might not exist at all.

Table 2. The limits between the pore space and solid components calculated by the CGCA and *K*-means algorithm

Type of fertilization	The CGCA			The <i>K</i> -means algorithm		
	section 1	section 2	section 3	section 1	section 2	section 3
Control group	162	165	163	183	167	174
Mineral fertilization	155	156	155	165	165	161
Pig manure fertilization	149	148	151	153	151	156

The CGCA, instead, does not require strict assumptions regarding the desired number of cluster. This allows the number obtained to be better suited to the soil structure. Moreover, in its basic form, the values of the parameters may be calculated automatically, however, there exists the possibility of their optional change. A feature specific to it is the possibility that it can influence the proportion between the number of clusters in areas where data elements are dense, as opposed to their sparse regions. In addition, by the detection of peripheral clusters, the algorithm allows the identification of outliers. This enables their elimination or designation to more numerous clusters, thus increasing the homogeneity of the data set.

The segmentation of soil images using the proposed method has given promising results. The clustering algorithm enabled the detection and recognition of the soil features from which for our needs, the pore space was ascertained. However, its computation can be challenging even for recent computer hardware. The most significant trend towards facilitating this, is to increase the number of CPU cores and increase the CPU's ability to process more and more tasks in parallel. Even more important, and an integral part of this practice, is that it allows optimization, so that the complex algorithm could be performed in a reasonable time.

5 Summary

Recent advances in computed tomography and digital image processing provide non-destructive tools for studying the internal structures of soil aggregates. This seems very useful in characterizing the pore space and in quantifying the differences in pore structures of different types of soil. In so doing, a more detailed analysis may be obtained by quantifying the soil structure through using the proposed segmentation techniques based on the kernel density estimation.

In this paper, an alternative way of detecting pore space in computed tomography soil slices is proposed by way of using image processing and data clustering based on the kernel estimation methodology. This density-based clustering algorithm allows us to get a better comprehension and knowledge of data, with the objective of segmenting images into either pore space or into solid components that constitute homogeneous areas with respect to a property of interest.

The presented approach is more objective than classical parametric methods, and can be successfully applied for many tasks in data mining, particularly where arbitrary assumptions concerning the number or shape of clusters among data are not recommended. This approach is also motivated by the current rapid growth in computational power. Improved real-time data processing and algorithm efficiency have important add-on effects due to the concurrent increase in the quantity and complexity of the image data that are now being collected.

Acknowledgments. This work has been supported by the national grant Frame No NN 310 307 639 of the Polish Ministry of Science and Higher Education.

References

1. Bandyopadhyay, S., Saha, S., Pedrycz, W.: Use of a fuzzy granulation-degranulation criterion for assessing cluster validity. *Fuzzy Sets and Systems* 170, 22–42 (2011)
2. Charytanowicz, M., Niewczas, J., Kulczycki, P., Kowalski, P.A., Łukasik, S., Żak, S.: Complete Gradient Clustering Algorithm for Features Analysis of X-Ray Images. In: Piętka, E., Kawa, J. (eds.) *Information Technologies in Biomedicine*, vol. 2, pp. 15–24. Springer, Berlin (2010)
3. Das, D., Ghosh, M., Chakraborty, C., Maiti, A.K., Pal, M.: Probabilistic prediction of malaria using morphological and textural information. In: 2011 International Conference on Image Information Processing, Durgapur, India, November 3-5 (2011)
4. Fukunaga, K., Hostetler, L.D.: The estimation of the gradient of a density function, with applications in Pattern Recognition. *IEEE Transactions on Information Theory* 21, 32–40 (1975)
5. Ghosh, M., Das, D., Chakraborty, C., Ray, A.K.: Plasmodium vivax segmentation using modified fuzzy divergence. In: 2011 International Conference on Image Information Processing, Durgapur, India, November 3-5 (2011)
6. Hallett, P., Lichner, L., Czachor, H., Józefaciuk, G.: Pore shape and organic compounds drive major changes in the hydrological characteristics of agricultural soils. *European Journal of Soil Science* 64, 334–344 (2013)
7. Kaufman, L., Rousseeuw, P.J.: *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley, New York (1990)
8. Kulczycki, P.: *Estymatory jądrowe w analizie systemowej*. WNT, Warszawa (2005)
9. Kulczycki, P., Charytanowicz, M.: A Complete Gradient Clustering Algorithm Formed with Kernel Estimators. *International Journal of Applied Mathematics and Computer Science* 20, 123–134 (2010)
10. Kulczycki, P., Charytanowicz, M., Kowalski, P.A., Łukasik, S.: The Complete Gradient Clustering Algorithm: Properties in Practical Applications. *Journal of Applied Statistics* 39, 1211–1224 (2012)
11. Mirkin, B.: *Clustering for Data Mining: A Data Recovery Approach*. Chapman and Hall, London (2005)
12. Nowak, P., Romaniuk, M.: A fuzzy approach to option pricing in a Levy process setting. *International Journal of Applied Mathematics and Computer Science* 23, 613–622 (2013)
13. Nowak, P., Romaniuk, M.: Application of Levy processes and Esscher transformed martingale measures for option pricing in fuzzy framework. *Journal of Computational and Applied Mathematics* 263, 129–151 (2014)

14. Pal, N.R., Pal, S.K.: A review of image segmentation techniques. *Pattern Recognition* 29, 1277–1294 (1993)
15. Perret, J.S., Prasher, S.O., Kacimov, A.R.: Mass fractal dimension of soils macropores using computed tomography: from the box counting to the cube-counting algorithm. *Journal of Hydrology* 26, 285–297 (2003)
16. Peth, S., Nellesen, J., Fischer, G., Horn, R.: Non-invasive 3D analysis of local soil deformation under mechanical and hydraulic stresses by μ CT and digital image correlation. *Soil and Tillage Research* 111, 3–18 (2010)
17. Pires de Silva, A., Imhoff, S., Kay, B.: Plant response to mechanical resistance and air-filled porosity of soils under conventional and no-tillage system. *Scientia Agricola* 6, 451–456 (2004)
18. Ruberto, C.D., Dempster, A., Khan, S., Jarra, B.: Analysis of infected blood cell images using morphological operators. *Image and Vision Computing* 20, 133–146 (2002)
19. Silverman, B.W.: *Density Estimation for Statistics and Data Analysis*. Chapman and Hall, London (1986)
20. Wand, M.P., Jones, M.C.: *Kernel Smoothing*. Chapman and Hall, London (1994)
21. Wojnar, L., Majorek, M.: *Komputerowa analiza obrazu*. Computer Scanning System, Warszawa (1994)
22. Zdravkov, B., Cermak, J., Sefara, M., Janku, J.: Pore classification in the characterization of porous materials: A perspective. *Central European Journal of Chemistry* 5, 385–395 (2007)