# A Novel Feature Extraction Method Based on Histogram and Mathematical Morphology for Isolated Handwritten Greek Characters Recognition

Badre-eddine El Kessab, Cherki Daoui and Belaid Bouikhalene

*Laboratory of Information Processing and Decision Support,
Faculty of Science and Technology, BP 523, Beni Mellal, Morocco
bade10@hotmail.fr*

## *Abstract*

*The isolated handwritten character recognition with multiple styles is a challenging research problem. In this paper, we propose a novel method of features extraction for character recognition based on the mathematical morphology and histogram techniques into vertical, horizontal, diagonal and anti-diagonal directions, knowing that the features extarction method is an important step in many image processing tasks. In this context, we present two comparisons in isolated handwritten Greek characters recognition, in fact the first comparison is between the hybrid methods exploited in features extraction which are the mathematical morphology combined with the histogram method; in contrast the second comparison is performed in order to deduce what is the most powerful between third genres of distances used in classification The Euclidean, Manhattan, and Minkowski distances. For this purpose, we have pre-processing each character image with different techniques. Furthermore, in the experiments results we provide extensive comparisons which demonstrate that our method outperforms for different characters' recognition, the results that we have obtained demonstrates really in one hand the performance of a novel method used in features extraction and the Euclidean distance in classification in the other hand.*

***Keywords:*** *Isolated handwritten Greek character, Feature extraction, histogram, Mathematical morphology, Classification, Euclidean distance*

## 1. Introduction

The Optical character recognition (OCR) is widely used in artificial intelligence and computer vision for handwriting character recognition, has been considered as a very dynamic field for research given that its applicability in many different domains such as passport documents, invoices, bank statements, computerized receipts, business cards, mail, The postal automation, automatic processing of administrative files *etc*.

Furthermore, here this work addresses to the problem of handwritten characters recognition is one of the most challenging research areas due to the large variations encountered in the writing styles of different writers. Our recognition system requires: different techniques in preprocessing. Further a novel structural method based on hybridization in features extraction then the classification with different distances by the k-nearest neighbours classifiers suggested.

Moreover, the several efficient techniques in each of the three principal phases forming our system of recognition are used in this work which is firstly the pre-processing used in order to render the character image in a best quality. In Second time, the features extraction phase exploited to extract the features from character and to convert it to a vector. Then, in last time the learning-classification phase employed for entraining all

character images of learning database and classifying those of test database. In this framework, several studies has been done for recognition of isolated handwritten Greek characters recognition while we extracted the features of each character by the mathematical morphology method in first time then with the histogram in second time then with finally the both methods are combined to create a powerful novel method in third time in one hand, or about the learning-classification phase for the recognition of each unknown character we have used the k-nearest neighbours (KNN) [15-21] with the different distances on the other hand. In this sense and in order to achieve this task we have pre-processed each character image by the median filter, the thresholding, the normalization, the rotation and the centering, techniques while we extracted the features of each word by the mathematical morphology and histogram methods. In fact, our targeted purpose is being able to compare between the precision of these both pervious methods of features extraction in one side and between the performances of third distances used in the identification: the Euclidean, Minkowski and Manhattan distances on the other side for the isolated handwritten Greek characters recognition.

This paper is organized in the following sections:

Section 2 presents a schema of the recognition system. Section 3, explains suggested preprocessing, and Section 4 is devoted to feature extraction. Section 5 describes the recognition with the different distances. Lastly in Section 6 the results are explained. The result section is then followed by conclusions.

## 2. Recognition System

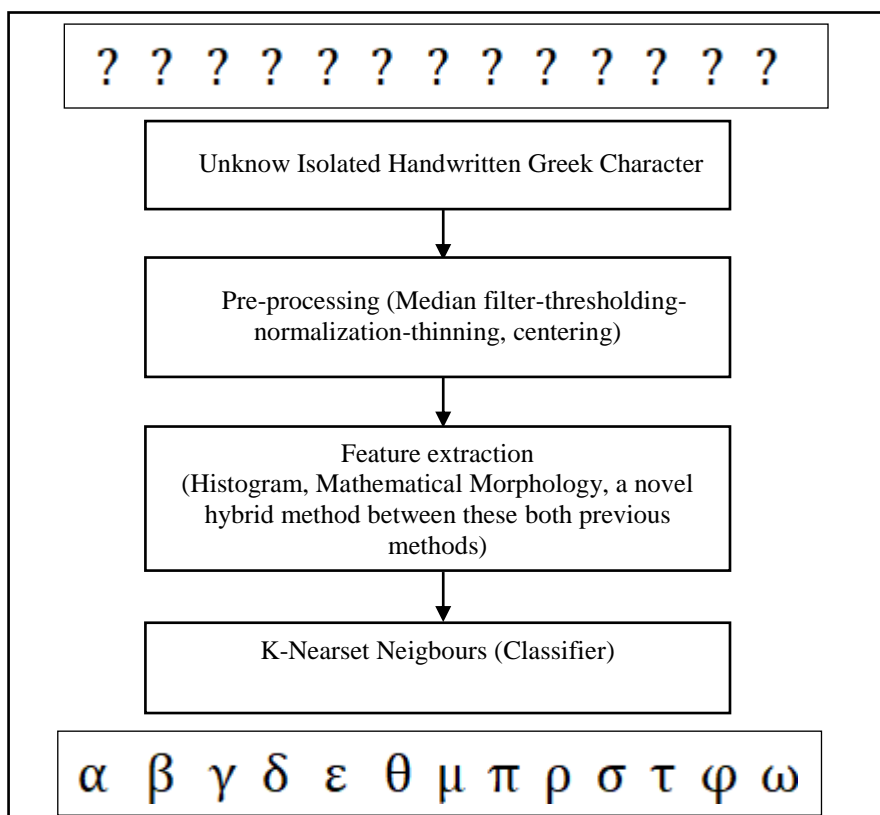The recognition system that we have opted in this study is presented in the following figure:



**Figure 1. The Proposed System for Isolated Handwritten Greek Character Recognition**

## 3. Database

The Greek alphabet is the character used to write in the Greek language and was the first alphabetic script to have distinct letters for vowels as well as consonants. The Greek alphabet today serves as a source of technical symbols and labels in several domains of mathematics and sciences. The alphabet has 24 letters as you see in Figure 2.
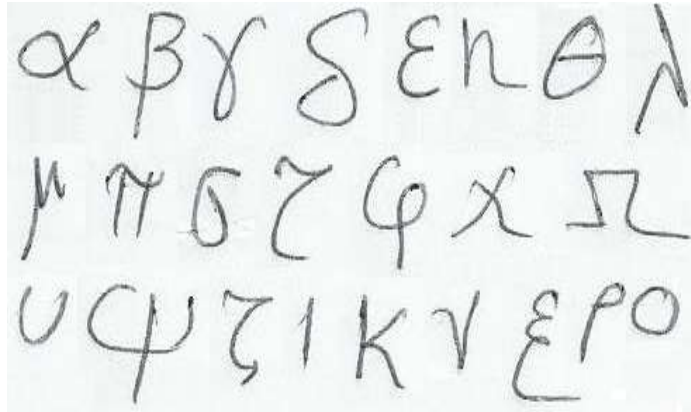


**Figure 2. Example of Isolated Handwritten Greek Characters from the Proposed Data Base**

## 4. Pre-processing

The pre-processing phase is a primary step in each OCR system is very important to creating a capable system to identify the Greek handwritten characters, we have collected the 24 characters in Greek language from α to ω written with 50 different students with different styles and a white back ground. On this work our problem reside in the identification of handwritten characters, these characters are written in different sizes, we need to all put it in standard size of 30x30 pixels, in fact the recognition for handwritten case is more complex due to varying writing styles from person to another due to different types of pens and papers *etc*.

The goal of this first phase in each OCR system is to remove each needless pixel including noise and redundant information in order to render in a best quality the character image so that it can be used in an efficient manner in the following phase which is the features extraction. Of this fact, to achieve this task, we have pre-processed in this research the images by the following techniques:

- The median filter applied for performing a filtration of image.

- The thresholding used to render each image contains only the black and white colors according a pre-selected threshold.

- The normalization of a character size to reduce the characters to be the same size.

- The thinning of a character to make the image one easier to process,

- The centering exploited for localizing the numeral justly in center of its image.

## 5. Features Extraction

The features extraction phase in each OCR system has as important role in the identification character by the great discrimination between all unknow characters entered to the system. More precisely, in this work we use both structural methods [1, 8] in feature extraction. The first method is based on mathematical morphology and the second

method is interested to histogram methods in order to converting the image of character to an extraction vector obtained by these previous methods can be used as an abstracted feature for identification of a character for the pattern showing in Figure 1. For Greek character recognition vector is created for all the characters which are converted into 114 elements after preprocessing. In this framework, we have chosen to use these methods which are:

### 5.1. Method based on Histogram

In the image processing, the histogram [9-14] is a method that associates each intensity value of the number of pixels taking this value. The determination of the histogram is carried out by counting the number of pixels for each image intensity. In this work we count the number of white pixels in each row, each column and each diagonal each anti diagonal as you see in the following figure.
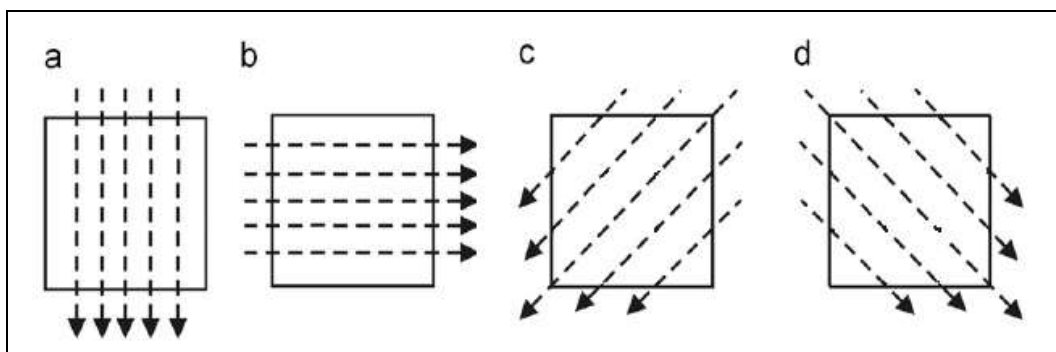


**Figure 3. Original Image**



**Figure 4. Process of Pixels Scanning by Columns, Rows, Anti-Diagonals and Diagonals**
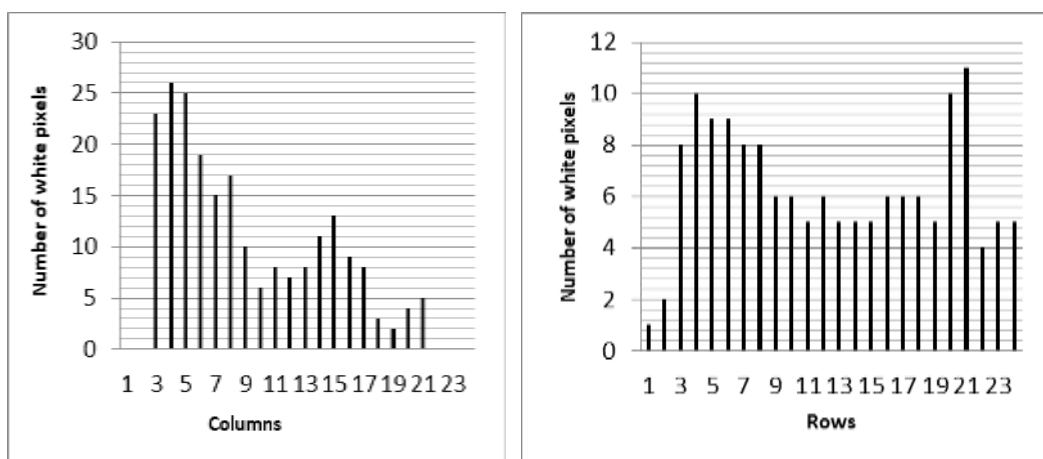


**Figure 5. Graphical Representation of Pixels Numbers in each Columns, Rows**
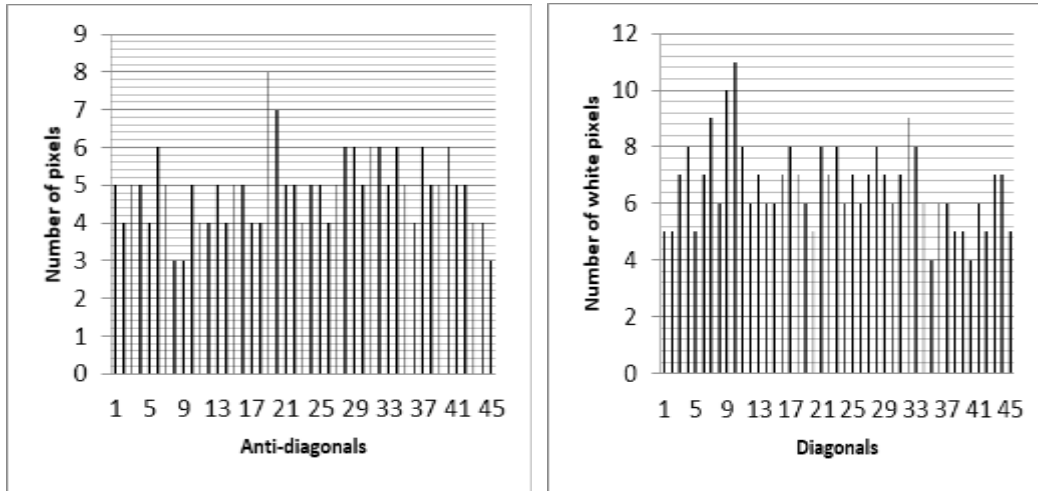
**Figure 6. Graphical representation of Pixels Numbers in each Diagonals and Anti-Diagonals**

We'll find four extraction vectors characterizing the image:

- Vector representing the number of white pixels in rows: $V_{Rows}$.
- Vector representing the number of white pixels in columns $V_{Columns}$.
- Vector representing the number of white pixels in diagonals $V_{Diagonals}$.
- Vector representing the number of white pixels in anti-diagonals $V_{AntiDiagonales}$.

The global extraction vector of parameters obtained by this method is given by:

$$V_{Global} = [\ V_{Rows},\ V_{Columns},\ V_{Diagonales},\ V_{AntiDiagonales}].$$

## 5.2. Method based on Mathematical Morphology

### 5.2.1. Dilatation

The Dilation is one of the basic operations in mathematical morphology [2-8]. Originally developed for binary images, it has been expanded first to grayscale images. The dilation operation usually uses a structuring element for probing and expanding the shapes contained in the input image is based on the change of pixels values (black or white) of the processed image. In this work we expand the processed image without the structuring element like you see in the followed algorithm:

➢ **Algorithm for Dilatation to the East:**

1. Begin
2. For each row
3. For each column
4. {
5. If ImageCharacter (row, column) = 1
6. ImageCharacter (row, column + 1) = 1
7. Increment
8. }
9. Find the Dilatation image to The East.
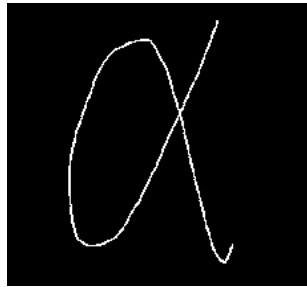10. End

> ➢ **Example of the dilatations of Alpha character.**



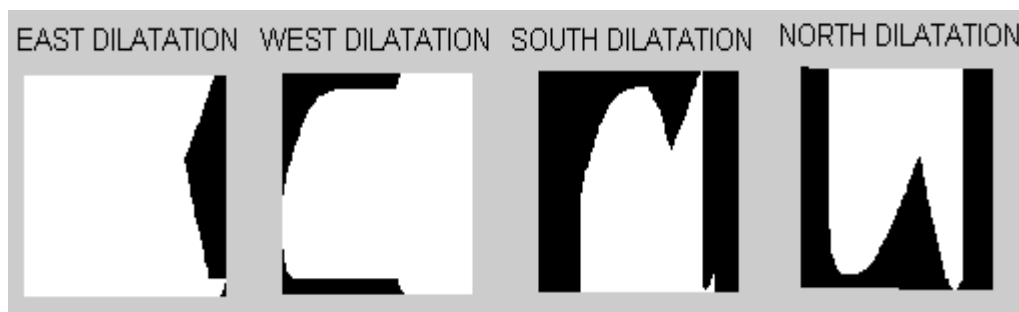**Figure 7. Original Image pf Alpha Character**



**Figure 8. Dilatation of Character Alpha**

## 5.3. The Characteristic Zones

The characteristic zones can be detected by the intersections of dilations found to the East, West, North and South. We define for each image five types of characteristic zones: East, West, North, South, and Central zone (this is just the image object or the white pixels in original image).

### 5.3.1. Extraction of Central Characteristic Zone:

A point of the original image (see Figure 7) belongs to the Central characteristic zone (Figure 9) if and only if:
- This point does not belong to the object (the white pixels in original image).
- From this point, moving in a straight line to the South, North, East and West on crosses the object. The result of the extraction is illustrated in (Figure 9). The same applies to the other zones.



**Figure 9. Original Image pf Alpha Character with Central Characteristics Zone**

## 5.4. Hybrid Method based on Histogram and Mathematical Morphology

This method is based on the combination between the two structural methods as you have seen above in order to increase the performance of the extraction phase and to

discriminate as much as possible the characters in order to facilitate its recognition in the classification phase, in order to convert the character image in a vector we count the number of white pixels in central characteristic zone (see Figures 11) then we count the number of white pixels in each row, each colum n, each diagonal and each anti diagonal and we put it in a extraction vector as you see in the following figures.
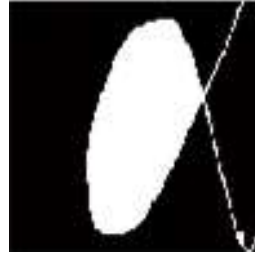


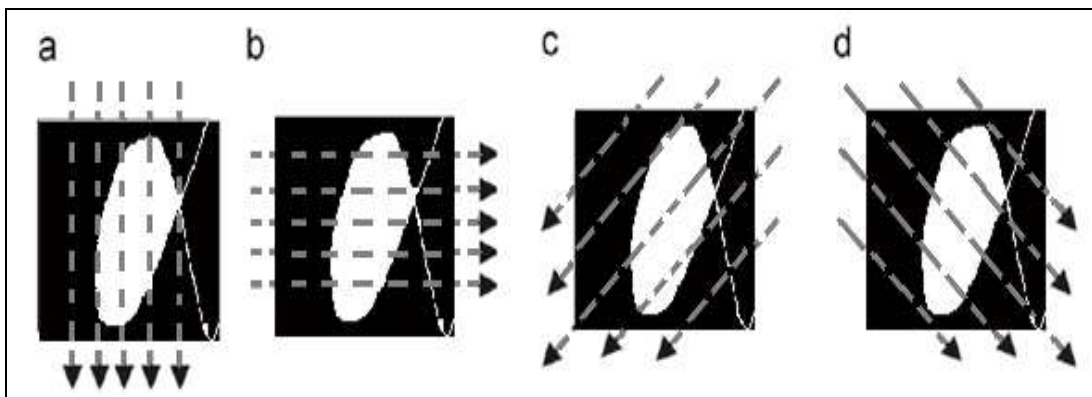**Figure 10. The Central Characteristics Zone of Alpha Character**



**Figure 11. Process of Pixels Scanning by Columns, Rows, Anti-Diagonals and Diagonals using Hybrid Method**
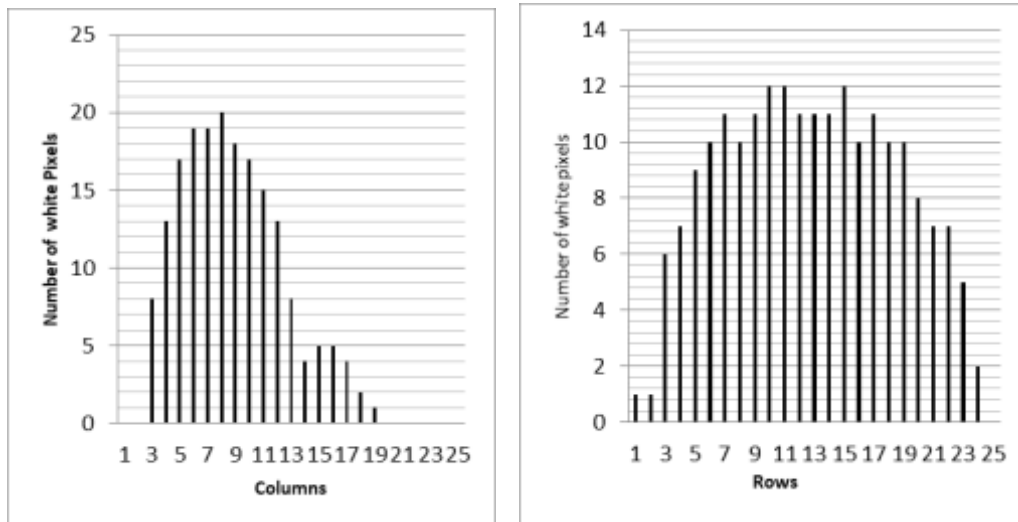


**Figure 12. Graphical Representation of Pixels Numbers in each Columns, and Rows**
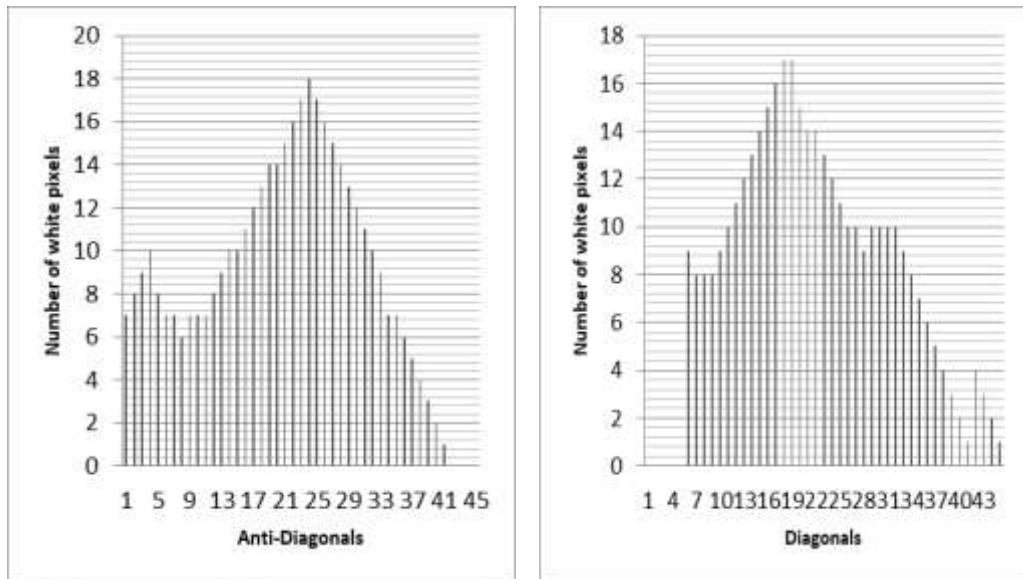
**Figure 13. Graphical Representation of Pixels Numbers in each Diagonals and Anti-Diagonals**

## 6. Recognition

The classification is the final procedure in character recognition for sorting the unknow character in our input system and assigning them to specific categories that is to say the correct identification. It is very important to use an efficient classifier in such a way that a system recognition becomes very performant. Here we have used the k-nearest neighbours classifier such as:

In practice, we have exploited the k nearest neighbors which is a method used for classifying pattern based on closest training examples in the feature space. The k-nearest neighbor algorithm [15-21] is considered as amongst and simplest of all machine learning algorithms: its principle of functioning consists to classify an unknown pattern (pattern of test) by a majority vote of its nearest neighbors in terms of distance, that is to say the unknown pattern is attributed to the class containing the greatest number of k nearest neighbors, it is noteworthy in this sense that the number k is a positive integer, typically not great. For well explaining how the classification by this method is carried, we present the following example:
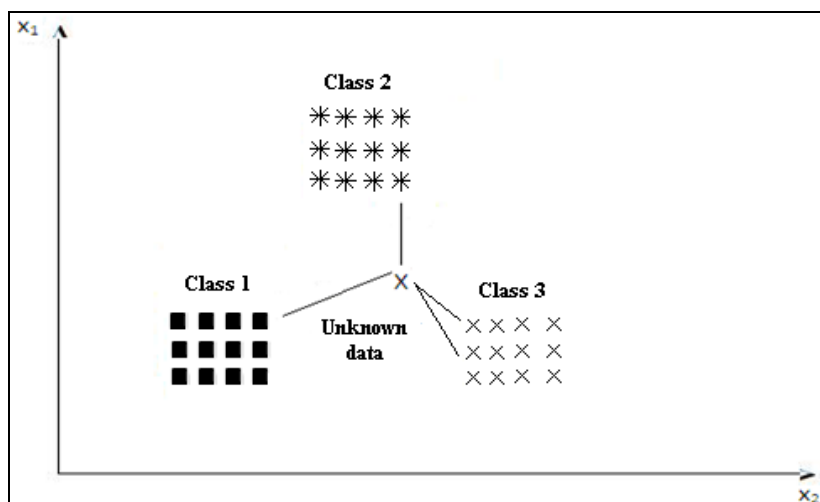


**Figure 14. Example of 4 Nearest Neighbors Classifier**

The Figure 4 shows that the unknown data is attributed to class 3.

Furthermore, hence several types of distance can be used for calculating the k nearest neighbors of unknown pattern, in this side there we have employed the following distances in order to compare between their performances in the isolated handwritten Arabic numerals recognition.

To this effect, given two vectors:

$$X = (x_1, x_2 \ldots\ldots x_N), Y = (y_1, y_2 \ldots\ldots y_N) \in IR^N \tag{1}$$

Between X and Y different distances d (X, Y) are defined by:

- **The Euclidean distance :**

$$d(X,Y) = \sqrt{\sum_{i=1}^{n}(x_i - y_i)^2} \tag{2}$$

- **The city bloc distance or Manhattan distance:**

$$d(X,Y) = \sum_{i=1}^{n} |(x_i - y_i)| \tag{3}$$

- **The Minkowski distance :**

$$d(X,Y) = \sqrt[p]{\sum_{i=1}^{n} |x_i - y_i|^p} \tag{4}$$

These distances can be briefly expressed by:

$$d_\lambda(X, Y) = \left[\sum_{i=1}^{n} |x_i - y_i|^\lambda\right]^{\frac{1}{\lambda}} \tag{5}$$

- $\lambda=1$ Manhattan distance.
- $\lambda=2$ Euclidean distance.
- $\lambda \rightarrow p$ Minkowski distance.

## 7. Experiments and Results

Firstly, we present the different parameters used in this work:
In order to achieve the desired comparison, we have used the following data:

- Each character image has a size equal to images 24x24 pixels.

- The number of all images of learning and of test that we have used is equal to 3000 images.

After several tries, we have chosen:

- The number k of nearest neighbors equals to 8

- The number p used in Minkowski distance equal to 4.

- The obtained extraction vector converting from the image of character contains 114 elements.

Therefore, we grouped the values that we obtained of the recognition rate $\tau_c$ of each character (given in %) and of the global rate $\tau_g$ of all characters (given in %) in the following table:

## Table 1. The Obtained Recognition Rates τc and τg by Different Four Distances

| Numeral | τc (Euclidean distance) | | τc (Manhattan distance) | | τc (Minkowski distance) | |
|---|---|---|---|---|---|---|
| | Histogram method | Hybrid method | Histogram method | Hybrid method | Histogram method | Hybrid method |
| α | 94,67 | 95,00 | 92,67 | 93,00 | 91,67 | 94,55 |
| β | 89,00 | 93,11 | 84,00 | 90,00 | 86,00 | 91,00 |
| γ | 86,30 | 92,00 | 85,67 | 90,67 | 88,33 | 90,67 |
| δ | 88,25 | 90,00 | 81,67 | 89,77 | 84,33 | 89,33 |
| ε | 87,00 | 91,21 | 88,33 | 88,30 | 85,00 | 90,24 |
| ζ | 77,17 | 78,12 | 68,32 | 76,32 | 70,60 | 71,35 |
| η | 72,21 | 77,14 | 66,22 | 74,34 | 66,76 | 70,26 |
| θ | 90,16 | 89,10 | 87,44 | 86,21 | 82,37 | 88,56 |
| ι | 91,00 | 94,00 | 88,00 | 90,22 | 88,00 | 90,34 |
| κ | 81,11 | 84,33 | 80,63 | 80,46 | 78,67 | 81,21 |
| λ | 74,13 | 75,00 | 73,27 | 76,94 | 70,43 | 75,64 |
| μ | 77,20 | 78,00 | 76,40 | 81,00 | 72,20 | 77,13 |
| ν | 63,54 | 70,21 | 62,00 | 78,34 | 62,70 | 63,50 |
| ξ | 80,12 | 81,00 | 76,46 | 75,61 | 70,32 | 77,09 |
| o | 86,69 | 91,31 | 84,24 | 83,78 | 81,45 | 80,21 |
| π | 66,62 | 68,20 | 64,37 | 94,68 | 89,83 | 74,25 |
| ρ | 74,70 | 76,12 | 68,80 | 72,00 | 69,55 | 72,18 |
| ς | 80,71 | 83,51 | 77,90 | 78,56 | 71,24 | 80,66 |
| σ | 67,50 | 70,67 | 66,71 | 66,60 | 62,31 | 74,46 |
| τ | 80,14 | 88,12 | 79,40 | 71,96 | 68,11 | 80,31 |
| φ | 75,41 | 80,71 | 73,00 | 78,00 | 71,00 | 80,00 |
| χ | 66,55 | 70,00 | 65,00 | 67,00 | 62,00 | 68,48 |
| ψ | 70,62 | 79,27 | 67,00 | 70,00 | 66,00 | 69,74 |
| ω | 74,00 | 86,72 | 70,00 | 75,00 | 68,00 | 74,71 |
| τg | **78,95** | **82,62** | **76,15** | **80,37** | **75,29** | **79,42** |

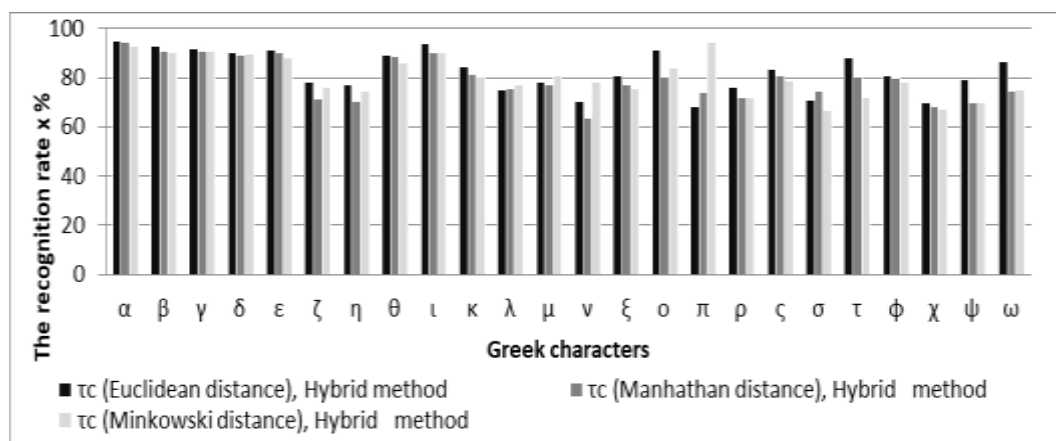The graphical representation to recognition rate of each numeral $\tau_c$ is:



## Figure 15. The Graphical Representation of Recognition Rate τc of Hybrid Method with Different Distances
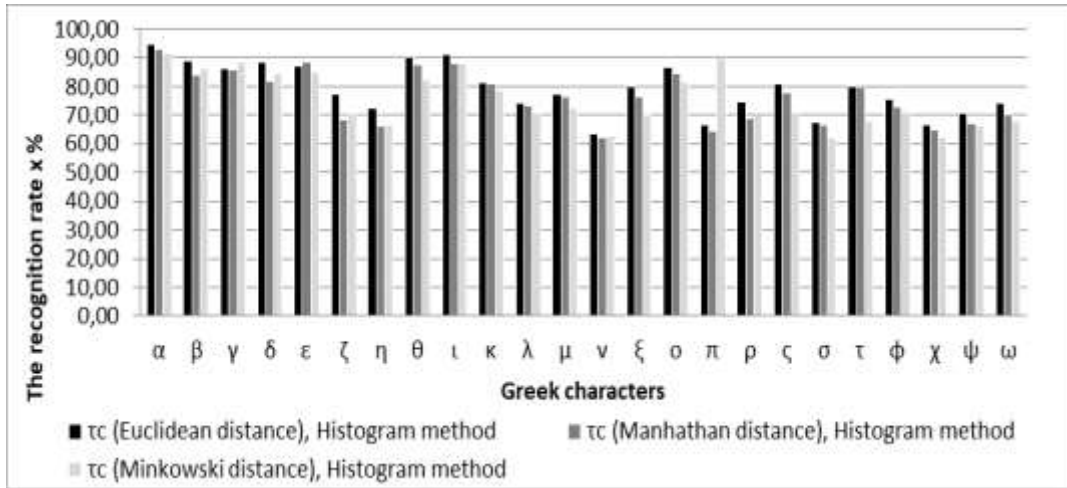
**Figure 16. The Graphical Representation of Recognition Rate τc of Histogram Method with Different Distances**

The graphical representation to recognition rate of all numerals $\tau_g$ is presented in the following figure:



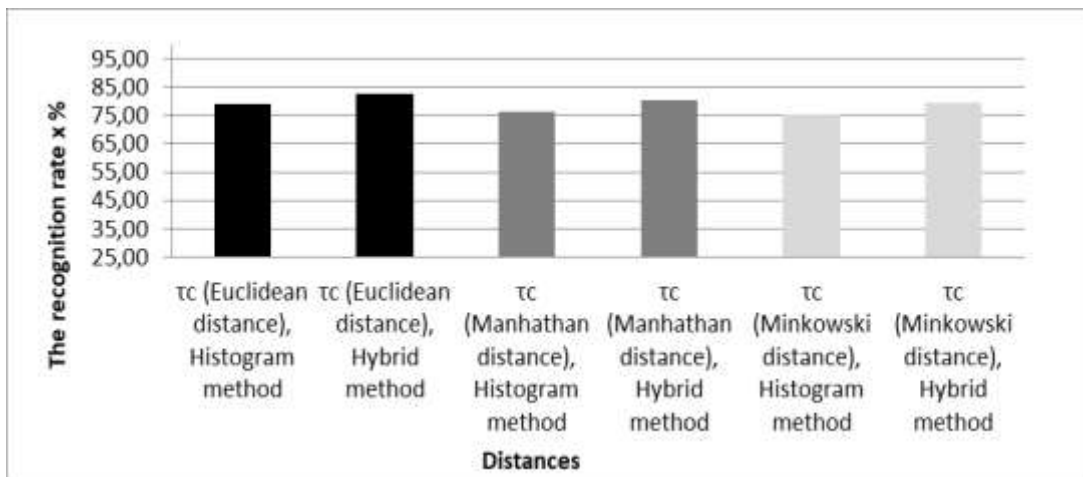**Figure 17. The Graphical Representation of Global Rate Recognition τg of each Method and of each Distance**

➤ *Analysis and comment :*

Taking into account all the results that we obtained, we really can to conclude that:
- The most performant extraction method is the mathematical morphology combined with the histogram.

- All the distances exploited in the k nearest neighbours classifier are generally almost equivalents in classification.

## 8. Conclusion

This paper presents a comparison between the performances of both structural methods which are the mathematical morphology combined with to histogram for recognition of handwritten Greek characters.

In this experimental study, we have verified that the recognition systems used in this approach which contains in the first phase the preprocessing by using the median filter, the thresholding, the normalization, the rotation and the centering and in second phase both structural methods in feature extraction which are the mathematical morphology and histogram in order to converting the all images of characters used in recognition to a extraction vector in reason the discriminate and facilitate its recognition, Finally the k-nearest neighbours with different distances are tested in the recognition phase. The obtained results really show that the most powerful recognition system is that contains the hybrid method which are the mathematical morphology combined with the histogram in features extraction phase and the all distances exploited in the k nearest neighbor's classifier are generally almost equivalents in the recognition phase.

The proposed system achieves a recognition accuracy of 82.62% using the hybrids methods. All these methods have given equal importance to different distances used in classification. So, our future work aims at finding other effective combination by introducing other method whatsoever statistical or structural in features extraction and other method the classification in recognition phase.

## References

[1] A. A. Desai, "Gujarati handwritten numeral optical character reorganization through neural network", Pattern Recognition, vol. 43, no. 7, **(2010)** July, pp. 2582-2589.

[2] Special Issue on Mathematical Morphology & Nonlinear Image Processing, Pattern Recognition, vol. 33, no. 6, **(2000)** June, pp. 875-876.

[3] B. El kessab, C. Daoui, B. Bouikhalene, M. Fakir and K. Moro, "Extraction Method of Handwritten Digit Recognition Tested on the MNIST Database", International Journal of Advanced Science and Technology, vol. 50, **(2013)** January.

[4] B. El Kessab, C. Daoui, B. Bouikhalene, M. Fakir and K. Moro, "Handwritten Tifinagh Text Recognition using Neural Networks and Hidden Markov Models", International Journal of Computer Applications (0975 – 8887), vol. 75, no. 18, **(2013)** August.

[5] B. El Kessab, C. Daoui, B. Bouikhalene and R. Salouan, "Some Comparative Studies for Cursive Handwritten Tifinagh Characters Recognition Systems", International Journal of Hybrid Information Technology, vol. 7, no. 6, **(2014)**, pp. 295-306.

[6] B. El Kessab, C. Daoui, B. Bouikhalene and R. Salouan, "A Comparative Study between the Support Vectors Machines and the K-Nearest Neighbors in the Handwritten Latin Numerals Recognition", International Journal of Signal Processing, Image Processing and Pattern Recognition, vol. 8, no. 2, **(2015)**, pp. 325-336.

[7] N. Batool and R. Chellappa, "Fast detection of facial wrinkles based on Gabor features using image morphology and geometric constraints", Pattern Recognition, vol. 48, no. 3, **(2015)** March, pp. 642-658.

[8] A. Căliman, M. Ivanovici and N. Richard, "Probabilistic pseudo-morphology for grayscale and color images", Pattern Recognition, vol. 47, no. 2, **(2014)** February, pp. 721-735.

[9] B. Yuan and M. Liu, "Power histogram for circle detection on images", Pattern Recognition, vol. 48, no. 10, **(2015)** October, pp. 3268-3280.

[10] B. Lei, E. Tan, S. Chen, D. Ni and T. Wang, "Saliency-driven image classification method based on histogram mining and image score", Pattern Recognition, vol. 48, no. 8, **(2015)** August, pp. 2567-2580.

[11] X. Zhao, Z. He, S. Zhang and D. Liang, "Robust pedestrian detection in thermal infrared imagery using a shape distribution histogram feature and modified sparse representation classification", Pattern Recognition, vol. 48, no. 6, **(2015)** June, pp. 1947-1960.

[12] A. Boulmerka, M. S. Allili and S. Ait-Aoudia, "A generalized multiclass histogram thresholding approach based on mixture modelling", Pattern Recognition, vol. 47, no. 3, **(2014)** March, pp. 1330-1348.

[13] A. Gordo, F. Perronnin and E. Valveny, "Large-scale document image retrieval and classification with runlength histograms and binary embeddings", Pattern Recognition, vol. 46, no. 7, **(2013)** July, pp. 1898-1905

[14] N. Lu, J. Silva, Y. Gu, S. Gerber, H. Wu, H. Gelbard, S. Dewhurst and H. Miao, "Directional histogram ratio at random probes: A local thresholding criterion for capillary images", Pattern Recognition, vol. 46, no. 7, **(2013)** July, pp. 1933-1948.

[15] T. Wakahara and Y. Yamashita, "K-NN classification of handwritten characters via accelerated GAT correlation", Pattern Recognition, vol. 47, no. 3, **(2014)** March, pp. 994-1001.

[16] Z. Liu, Q. Pan and J. Dezert, "A new belief-based K-nearest neighbor classification", Pattern Recognition, vol. 46, no. 3, **(2013)** March, pp. 834-844.

[17] L. Li, L. Zhang and J. SU, "Handwritten character recognition via direction sring and nearest neighbor matching", The Journal of China Universities of Posts and Telecommunications, vol. 19, Supplement 2, **(2012)** October, pp. 160-165,196.

[18] S. H. Rodríguez, J. F. M. Trinidad, J. Ariel and C. Ochoa, "Fast k most similar neighbor classifier for mixed data (tree k-MSN)", Pattern Recognition, vol. 43, no. 3, **(2010)** March, pp. 873-886.

[19] N. A. Samsudin and A. P. Bradley, "Nearest neighbor group-based classification", Pattern Recognition, vol. 43, no. 10, **(2010)** October, pp. 3458-3467.

[20] J. Yang and D. Zhang, "From classifiers to discriminators: A nearest neighbor rule induced discriminant analysis", Pattern Recognition, vol. 44, no. 7, **(2011)** July, pp. 1387-1402.

[21] J. Yang, L. Zhang, J. Yang and D. Zhang, "From classifiers to discriminators: A nearest neighbor rule induced discriminant analysis", Pattern Recognition, vol. 44, no. 7, **(2011)** July, pp. 1387-1402.

# Authors

**B. El Kessab**, received his Ph.D degree on informatics in 2014 and Master's degree in 2009 from Faculty of Sciences and Technology University Sultan Moulay Slimane Beni Mellal Morocco. The current research interests include pattern recognition, image analysis, document processing and automatic processing of natural languages.

**C. Daoui**, received his Ph.D degree on mathematics in 2002 from Mohamed V University Rabat Morocco. Currently is a professor in Faculty of Sciences and Technology, University Sultan Moulay Slimane Beni Mellal Morocco. His research topics are: the mathematics, operational research and pattern recognition.

**B. Bouikhalene**, received his Ph.D degree on mathematics in 2001 and Master's degree on Science of Computer and Telecommunications in 2007 from the University Ibn Tofel Kenitra. Currently is a professor in the Sultan Moulay Slimane University Beni Mellal Morocco. His research topics are: the pattern recognition, artificial intelligence and mathematics and its applications.