

# Text Search in Document Images Based on Hausdorff Distance Measures

Andrey Andreev, Nikolay Kirov

**Abstract:** *The Hausdorff type distances between the sets of points on the plane are the commonly used similarity measures for binary images. In this work we present several such measures in a unified manner and introduce a new, naturally arisen variant of Hausdorff distance. The matching performance of all similarity measures is compared by computer experiments, using real word images from a scanned book.*

**Key words:** *Binary Text Images, Hausdorff Distance, Similarity Measures, Word Searching*

## Introduction

Libraries contain huge amounts of historical documents which cannot be made available online because they do not have a searchable index. The wordspotting idea has been proposed as a solution for creating indexes for such documents by matching word images. Optical character recognition is the usual way of conducting text retrieval from scanned document images. Moreover recognizing full text in images is a wasteful task for information retrieval. The motivation of our work is to choose effective search in scanned documents by simply considering the image similarities. One of the most widespread ideas is to use Hausdorff type measures for word image similarity.

The classical Hausdorff distance (HD) between two point sets  $A$  and  $B$  is defined as

$$H(A, B) = \max\{h(A, B), h(B, A)\}, \quad (1)$$

where  $h(A, B)$  and  $h(B, A)$  are co-called directed distances between the sets. For original Hausdorff metrics

$$h(A, B) = \max_{a \in A} d(a, B), \quad \text{where} \quad d(a, B) = \min_{b \in B} \rho(a, b),$$

i.e.  $d(a, B)$  is the distance from a point  $a$  to the set  $B$ , and  $\rho(a, b)$  is a point distance.

Let  $A$  and  $B$  be two finite sets on the plane,  $|A| = N_A$  and  $|B| = N_B$  denote their number of points. The classical point distance in the plane is Euclidean distance

$$\rho(a, b) = \sqrt{(a_x - b_x)^2 + (a_y - b_y)^2}.$$

Here  $a_x$  and  $b_x$  are  $x$ -coordinates and  $a_y$  and  $b_y$  are  $y$ -coordinates of the points  $a$  and  $b$ . This distance is called Minkowski distance of order 2.

Manhattan distance (Minkowski distance of order 1) is often used  $\rho(a, b) = |a_x - b_x| + |a_y - b_y|$  as well as infinity norm distance  $\rho(a, b) = \max\{|a_x - b_x|, |a_y - b_y|\}$ . The last two variants are easy to be calculated, without multiplication and not using square root. We note that 0-1 distance

$$\rho(a, b) = \begin{cases} 0 & \text{if } a \equiv b \\ 1 & \text{otherwise} \end{cases} \quad (2)$$

---

This research has been partially supported by a Marie Curie Fellowship of the European Community program "Knowledge Transfer for Digitalization of Cultural and Scientific Heritage in Bulgaria" under contract number MTKD-CT-2004-509754.

defines also a metric in the plane.

Huttenlocher *at al.* [4] proposed the partial Hausdorff distance (PHD) for comparing images containing a lot of degradations or occlusions. For directed distance they take the  $K$ -th ranked point of  $A$  instead of the largest one

$$h_K(A, B) = K_{a \in A}^{th} d(a, B), \quad (3)$$

where  $K_{a \in A}^{th}$  denotes the  $K$ -th ranked value in the set of distances  $\{d(a, B) : a \in A\}$ , i.e. for each point of  $A$ , the distance to the closest point of  $B$  is computed, and then, the points of  $A$  are ranked by their respective values to this distance,

$$d(a_1, B) \geq d(a_2, B) \geq \dots \geq d(a_K, B) \geq \dots \geq d(a_{N_A}, B). \quad (4)$$

This HD measure requires one parameter, often represented by  $f = \frac{K}{N_A}$  ( $0 \leq f \leq 1$ ). Sim *at al.* [6] claim that a value in the interval  $[0.6, 0.8]$  gives good matching results. Note that this measure is not a metric because  $h_K(A, A) > 0$ !

The idea of José Paumard [5] is that we do not take into account the  $L$  closest neighbors of  $a \in A$  in  $B$ . So we can define the distance from a point  $a \in A$  to the set  $B$  as follows

$$d_L(a, B) = L_{b \in B}^{th} \rho(a, b),$$

where  $L_{b \in B}^{th}$  denotes the  $L$ -th ranked value in the set of distances  $\{\rho(a, b) : b \in B\}$  for a given point  $a$  of  $A$ . Now the directional Censored Hausdorff Distance (CHD) can be defined as

$$h_{K,L}(A, B) = K_{a \in A}^{th} d_L(a, B) = K_{a \in A}^{th} L_{b \in B}^{th} \rho(a, b). \quad (5)$$

Let us set two parameters  $\alpha = \frac{K}{N_A}$  and  $\beta = \frac{L}{N_B}$  which are relative values with respect to the number of points in the sets  $A$  and  $B$ . Then the recommended values in [5] for these parameters are  $\alpha = 0.1$  and  $\beta = 0.01$ .

For all three described measures (HD, PHD and CHD), the directed distance can be considered as a choice a representative pair of points  $(a_0, b_0)$ ,  $a_0 \in A$  and  $b_0 \in B$  such that the point distance between them  $\rho(a_0, b_0)$  is equal to the corresponding directed distance between the sets  $A$  and  $B$ . Another approach for measuring similarity between two finite sets in the plane is to calculate a sum of point distances.

Dubuisson and Jain [3] examined a number of distance measures of Hausdorff type for determination to what extend two point sets on the plane  $A$  and  $B$  differ. They introduced so-called Modified Hausdorff Distance (MHD) with the following distance measure

$$h_{\text{MHD}}(A, B) = \frac{1}{N_A} \sum_{a \in A} d(a, B) = \frac{1}{N_A} \sum_{a \in A} \min_{b \in B} \rho(a, b). \quad (6)$$

They claim than it suites in best way the problem for object matching supposing that  $\rho$  is the Euclidean metrics. We use infinity norm distance for our experiments (see [1], [2]) measuring the word similarities in binary text documents and conclude that this is one of the best measures for word matching. For comparison reason we try also MHD with 0-1 point distance (2), which is easier for calculation.

A bit better results were obtained in our examples omitting the coefficient  $\frac{1}{N_A}$  in front of the sum (6). We called this modification Sum Hausdorff Distance (SHD), [2]

$$h_{\text{SHD}}(A, B) = \sum_{a \in A} d(a, B) = \sum_{a \in A} \min_{b \in B} \rho(a, b). \quad (7)$$

In 1999 D.-G. Sim *at al.* [6] described two variants of MHD for elimination of outliers – usually the points of outer noise. Based on robust statistics M-estimation and least trimmed square they introduced M-HD and LTS distances.

The directed distance for M-HD is defined by

$$h_M(A, B) = \frac{1}{N_A} \sum_{a \in A} f(d(a, B)), \quad (8)$$

where the function  $f$  is convex and symmetric and has a unique minimum value at zero. One possible function is

$$f(x) = \begin{cases} |x| & \text{if } |x| \leq \tau \\ \tau & \text{if } |x| > \tau \end{cases}$$

This means that we sum the distances  $d(a, B)$  which are less than the constant  $\tau$  and add  $\tau$  when the distance is greater than  $\tau$ . The recommended interval of  $\tau$  is [3, 5]. Note that MHD with 0-1 point distance (2) is M-HD for  $\tau = 1$ .

The second measure proposed in [6] is called Least Trimmed Square HD (LTS-HD). The directed distance is

$$h_{LST}(A, B) = \frac{1}{N_A - K} \sum_{i=K}^{N_A} d(a_i, B), \quad (9)$$

where  $K \leq N_A$  and  $a_1, a_2, \dots, a_{N_A}$  are points of  $A$  for which (4) is valid. Parametrization of the method can be done by a parameter  $\alpha = \frac{K}{N_A}$ . For comparing noisy binary images the suggested value for this parameter is 0.2.

Following the definition of CHD (5), we introduce its analogical method based on the sum of point distances. The directed distance is

$$h_{NEW}(A, B) = \frac{1}{N_A - K} \sum_{i=K}^{N_A} d_L(a_i, B) = \frac{1}{N_A - K} \sum_{i=K}^{N_A} L_{b \in B}^{th} \rho(a, b). \quad (10)$$

We can set again the parameters  $\alpha = \frac{K}{N_A}$  and  $\beta = \frac{L}{N_B}$  which are relative values with respect to the number of points in the sets  $A$  and  $B$ .

## A new approach to similarity measures

We can consider a linear order of points of  $A$  and give a sequence representation:  $A = \{a_1, a_2, \dots, a_{N_A}\}$ . For every  $a_k \in A$  ( $k = 1, 2, 3, \dots, N_A$ ) we can calculate the distances (with respect to a metric  $\rho$  in  $R^2$ ) from  $a_k$  to all points in  $B$ , i.e.

$$d_k^1 = \min_{b \in B} \rho(a_k, b) = \rho(a_k, b_k^1), \quad d_k^2 = \min_{b \in B \setminus \{b_k^1\}} \rho(a_k, b) = \rho(a_k, b_k^2), \dots,$$

$$d_k^l = \min\{\rho(a_k, b) : b \in B \setminus \{b_k^1, b_k^2, \dots, b_k^{l-1}\}\} = \rho(a_k, b_k^l), \dots,$$

obtaining in such a way a nondecreasing sequence of numbers

$$d_k^1 \leq d_k^2 \leq \dots \leq d_k^l \leq \dots \leq d_k^{N_B}.$$

Carrying out these calculations for every point in  $A$ , we define a distance matrix  $D$

$$D = \begin{pmatrix} d_1^1 & d_1^2 & d_1^3 & \dots & d_1^l & \dots & d_1^{N_B} \\ d_2^1 & d_2^2 & d_2^3 & \dots & d_2^l & \dots & d_2^{N_B} \\ d_3^1 & d_3^2 & d_3^3 & \dots & d_3^l & \dots & d_3^{N_B} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ d_k^1 & d_k^2 & d_k^3 & \dots & d_k^l & \dots & d_k^{N_B} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ d_{N_A}^1 & d_{N_A}^2 & d_{N_A}^3 & \dots & d_{N_A}^l & \dots & d_{N_A}^{N_B} \end{pmatrix}$$

following arbitrary order of points in  $A$ . Later we will choose ordering of rows, corresponding to an order in a column. For definition of MHD (6) and M-HD (8) we do not need any order

$$h_{\text{MHD}}(A, B) = \frac{1}{N_A} \sum_{i=1}^{N_A} d_i^1, \quad \text{and} \quad h_{\text{M}}(A, B) = \frac{1}{N_A} \sum_{i=1}^{N_A} \min\{d_i^1, \tau\}.$$

For finding the Hausdorff distance (1) in the distance matrix  $D$ , we consider the following order (obtained by swapping the rows) with respect to the first column of  $D$

$$h(A, B) = d_1^1 \geq d_2^1 \geq \dots \geq d_k^1 \geq \dots \geq d_{N_A}^1.$$

The directed distance for partial Hausdorff distance (3) is  $h_K(A, B) = d_K^1$ . Even more, now we can calculate LTS-HD distance (9) summing the part of the first column elements

$$h_{\text{LST}}(A, B) = \frac{1}{N_A - K} \sum_{i=K}^{N_A} d_i^1.$$

Also, we can find CHD directed distance (5) as an element of matrix  $D$ . For this purpose we swap the rows of the matrix in such way that the  $L$ -th column is sorted, i.e.

$$d_1^L \geq d_2^L \geq \dots \geq d_k^L \geq \dots \geq d_{N_A}^L.$$

where  $L$  is the same number as in (5). Then  $h_{K,L}(A, B) = d_K^L$ . In addition, it is easy to find the value of the directed NEW distance (10), namely

$$h_{\text{NEW}} = \frac{1}{N_A - K} \sum_{i=K}^{N_A} d_i^L.$$

## Experiments

We carried out our experiments using an old book (1884) – Bulgarian Chrestomathy, created by famous Bulgarian writers Ivan Vasov and Konstantin Velichkov. The quality of scanned images are quite bad because this was one of the first books, processing in the digitization center and operators' qualification was not on appropriate level. Many pages have slopes in rows, there are significant variations in gray levels, etc.

There is no text version till now of this book, which may be produced using appropriate OCR software. The first reason is the quality of images. The second reason is the absence of OCR software because the text contains old and abandoned Bulgarian letters. Also spelling and grammar are quite different in modern Bulgarian language.

We used 200 pages from about 1000 book pages scanned at a resolution of 200 DPI as shown in Figures 1 and 2. The images are about  $2300 \times 3600$  pixels (8.28 MPixels), 14.8 x 23.3 cm, grayscale 256 (8 BitsPerPixel). We use preprocessing to convert the images to 1 bit per pixel, black and white, by the help of Image Magic software [7] with 60% threshold value.

The goal of our experiments is to compare practically the efficiency of described methods counting the number of correctly retrieved words in a sequence of words, sorted by their similarity measures with respect to the corresponding HD. For all experiments the same segmentation is used. We choose a pattern word and then measure similarities between it and the words with approximately same width.

Tables 1 and 2 contains numbers of correct words in an ordered sequence with the corresponding distance  $D$ . The numbers  $m$  and  $n$  in the ratio  $m/n$  in the tables denote:

- $m$ , the number of correct words with distance  $D$ ;
- $n$ , the number of all words with distance  $D$ .

поет, сатирик и публицист. Първо-то пѣшто, което е издалъ е книжка стихотворения „Васненикъ“ и по-послѣ „Смѣсна Книга“ (Букурештѣ 1852 г.), съ които той доби първа-та си извѣстность у насъ, като български писателъ. Отъ 1857 год. се почва него-ва та многополезна дѣятелность въ борба-та ни съ Гръци-тѣ за черковна независимость. Той дохожда въ Цариградъ и издава свои-тѣ „Смѣшни Календари“ сатирически книги, въ които бичува съ единъ искусенъ и ядовитъ сарказмъ пороци-тѣ и недостатки-тѣ на тогавашно-то българско общество, и гръцко-то високо духовенство (1857—1863). На 1863 год. той прѣдприе издаване-то на сатирически вѣстникъ „Гайда“, който не трая много врѣме. Доста хубави статии все въ полемическо-сатирически духъ, напечата той тамъ. Слѣдъ двѣ години Славейковъ прѣдприе издаване-то на политически вѣстникъ „Македония“ (1867—1870). Тамъ при разискване-то на разни въпроси отъ обществени и черковенъ интересъ-Славейковъ се стараше да разбуди народно-то чувство у Македонски-тѣ Българи, които душеше нетърпимо-то влияние на гръкоман, ство-то и фанариотство-то. Най-послѣ подиръ нѣколко врѣмenni спирация и конфискация на вѣстникъ-тѣ, правителство-то съвѣтъ го унищожилъ и запрѣтилъ на Славейкова да издава вече какъвъ-да-е вѣстникъ, а и него самаго турѣ въ тъмница, по обвинение, че въ послѣдни-тѣ броеве на „Македония“ явно проповѣдвалъ революционни идеи между Българе-тѣ.

Figure 1: A half page of the book, grayscale

поетъ, сатирикъ и публицистъ. Първо-то пѣшто, което е издалъ е книжка стихотворения „Васненикъ“ и по-послѣ „Смѣсна Книга“ (Букурештѣ 1852 г.), съ които той доби първа-та си извѣстность у насъ, като български писателъ. Отъ 1857 год. се почва него-ва та многополезна дѣятелность въ борба-та ни съ Гръци-тѣ за черковна независимость. Той дохожда въ Цариградъ и издава свои-тѣ „Смѣшни Календари“ сатирически книги, въ които бичува съ единъ искусенъ и ядовитъ сарказмъ пороци-тѣ и недостатки-тѣ на тогавашно-то българско общество, и гръцко-то високо духовенство (1857—1863). На 1863 год. той прѣдприе издаване-то на сатирически вѣстникъ „Гайда“, който не трая много врѣме. Доста хубави статии все въ полемическо-сатирически духъ, напечата той тамъ. Слѣдъ двѣ години Славейковъ прѣдприе издаване-то на политически вѣстникъ „Македония“ (1867—1870). Тамъ при разискване-то на разни въпроси отъ обществени и черковенъ интересъ-Славейковъ се стараше да разбуди народно-то чувство у Македонски-тѣ Българи, които душеше нетърпимо-то влияние на гръкоман, ство-то и фанариотство-то. Най-послѣ подиръ нѣколко врѣмenni спирация и конфискация на вѣстникъ-тѣ, правителство-то съвѣтъ го унищожилъ и запрѣтилъ на Славейкова да издава вече какъвъ-да-е вѣстникъ, а и него самаго турѣ въ тъмница, по обвинение, че въ послѣдни-тѣ броеве на „Македония“ явно проповѣдвалъ революционни идеи между Българе-тѣ.

Figure 2: A half page of the book, b/w

In Table 3 we count the number of correctly retrieved words among first 100, 200, . . . , 500 words with approximately same width. In Table 4 the ration  $m/n$  has the same meaning but the distances are in different scales.

We set the parameter  $f = 0.9$  for PHD measure (3) and  $\alpha = 0.2$  for LST-HD measure (9). For M-HD (8) we obtain results with  $\tau = 4$ .  $\alpha = 0.1$  and  $\beta = 0.01$  are parameters for CHD (5) and NEW (10).

There are two relative words (derivatives) of the pattern word **всички**, namely **всичка** and **всичко**. We count as correct words all three of them. This is very useful in practice and show another advantage of methods under discussion and our approach in search. Also, there are 5 similar words of the word **Русия**: **Руски**, **Руска**, **Руско**, **руски** and **руска**.

The best results are in bold in all tables.

## Discussion and Conclusion

In this article we do not discuss the quality of image preprocessing particularly the important step of segmentation. Also we have no data of number of searching words in the text, because this is tedious work which cannot be done by computer. It follows than we cannot produce the standard recall/precision retrieval estimation (see [2]). In addition, we cannot catch the words which are incorrect segmented as well as these which are break at the end of a line and remaining part is placed on the next line. Nevertheless we think that our comparison of similarity methods is significant for their implementations in software searching systems. In spite of low efficiency of these Hausdorff type methods (the searching takes a lot of time) we believe that the modern, high level personal computers could be able to solve the problem in reasonable time.

For word **всички**

$D =$	4	5	6	7	8
Method					
HD	16/16	44/44	115/120	168/217	177/500
PHD+3	<b>77/77</b>	206/254	209/500	—	—
CHD	19/19	<b>213/252</b>	<b>214/500</b>	—	—

Table 1: “Point-distances”

For word **Русия**

$D =$	4	5	6	7
Method				
HD+1	2/2	3/3	5/5	5/6
PHD+3	3/3	<b>11/15</b>	—	—
CHD		<b>8/8</b>	<b>13/24</b>	—

Table 2: “Point-distances”

For word **всички**

$n =$ Method	100	200	300	400	500
HD01	97	158	186	195	206
MHD	<b>100</b>	169	199	207	212
SHD	<b>100</b>	177	205	213	220
M-HD	<b>100</b>	173	202	214	218
LTS-HD	<b>100</b>	<b>185</b>	<b>215</b>	<b>221</b>	<b>224</b>
NEW	97	164	198	213	<b>224</b>

Table 3: “Sum-distances”

For word **Русия**

$n =$ Method			
HD01	4/4	9/18	10/23
MHD	10/10	<b>14/23</b>	<b>15/49</b>
SHD	<b>11/11</b>	14/24	–
M-HD	7/7	12/14	–
LTS-HD	10/10	<b>14/23</b>	–
NEW	7/7	12/15	14/26

Table 4: “Sum-distances”

The main conclusions that we derive from are:

1. “Sum-distances” (see Tables 3 and 4) outmatch “point-distances” (see Tables 1 and 2).
2. There are no significant differences between the methods that we call “sum-distances” ones.

## References

- [1] A. Andreev, N. Kirov, *Hausdorff Distance and Word Matching*, Proceedings of the International Workshop “Computer Science and Education”, June 3-5, 2005, Borovetz-Sofia, Bulgaria, 19-28.
- [2] A. Andreev, N. Kirov, *Word image matching in Bulgarian historical documents*, Review of the National Center for Digitalization, 8, (2006), 29-35.
- [3] M.-P. Dubuisson, A. Jain, *A Modified Hausdorff Distance for Object Matching*, In: Proc. 12th Int. Conf. Pattern Recognition, Jerusalem, Israel, 1994, pp. 566-568.
- [4] Daniel P. Huttenlocher, Gregory A. Klanderman, and William J. Rucklidge, *Comparing Images Using the Hausdorff Distance*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 15, (1993), No.9, 850-863.
- [5] José Paumard, *Robust comparison of binary images*, Pattern Recognition Letters 18 (1997), 1057-1063.
- [6] Dong-Gyu Sim, Oh-Kyu Kwon, and Rae-Hong Park, *Object Matching Algorithms Using Robust Hausdorff Distance Measures*, IEEE Trans. on Image Processing, 8, (1999), No.3, 425-429.
- [7] Image Magick, [www.imagemagick.org](http://www.imagemagick.org)

### ABOUT THE AUTHORS

Assoc.Prof. Andrey Andreev, PhD, Institute of Mathematics and Informatics, Bulgarian Academy of Sciences, Phone: (+359 2) 979 3874, E-mail: [aandreev@math.bas.bg](mailto:aandreev@math.bas.bg)

Assoc.Prof. Nikolay Kirov Kirov, PhD, Department of Informatics, New Bulgarian University, Phone: (+359 2) 811 0611, E-mail: [nkirov@nbu.bg](mailto:nkirov@nbu.bg) and Institute of Mathematics and Informatics, Bulgarian Academy of Sciences, Phone: (+359 2) 979 2850, E-mail: [nkirov@math.bas.bg](mailto:nkirov@math.bas.bg), WEB: <http://www.math.bas.bg/~nkirov>