

Distance between Histograms of Angular Measurements and its Application to Handwritten Character Similarity

Sung-Hyuk Cha and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition

State University of New York at Buffalo, U. S. A.

{scha,srihari}@cedar.buffalo.edu

Abstract

The similarity between two histograms has attracted many researchers in various fields. The type of histograms to be matched is often angular such as gradient directions in character images and hue values in color images. The distance between two angular type histograms differs from those of nominal or ordinal type histograms. Alas, conventional distance measures do not distinguish them. Thus we present a new algorithm to compute the newly proposed distance. We describe how to apply the distance to handwritten character similarity.

Key words: Character Similarity, Distance Measure, Histogram

1 Introduction

Finding similarity between two histograms is an important issue in pattern classification. Much of this study has been motivated and has found applications in various fields such as histogram based image indexing [7, 9, 10, 11] and pattern classification and cluster analysis [6]. A distance between histograms enables the histogram based indexing as one can retrieve similar distributions or enumerate them in order.

A histogram we consider is a one-dimensional array that represents the frequency of each discrete value of samples in finite range. Let X be a set of n samples $X = \{x_1, x_2, \dots, x_n\}$ and $H(X)$ be the histogram for the set X . $H_i(X)$ denotes the frequency of samples that fall into the i 'th bin of the histogram, thereby the value of samples in the i 'th bin is i . The total number of bins is b . For example of a set of samples, $X = \{1, 0, 7, 6, 0, 0, 2, 6, 6, 0\}$ where $n = 10$ and $b = 8$, the histogram $H(X) = \{4, 1, 1, 0, 0, 0, 3, 1\}$ and $H_0(X) = 4$. We will use these notations throughout the rest of this paper.

The $\Theta(b)$ algorithm for the distance between linear type histograms was introduced earlier to expedite the template matching problem [2, 3]. In this paper, we consider an angular type of histograms. Swain and Ballard used the *histogram intersection* measure as the similarity between color histograms to identify an object [11]. In their approach, histograms are treated as fixed-dimensional vectors, so called "*Nominal Type Histogram*". However, it is shown that conventional vector norm approach is inadequate when it is used

for certain types of measurements such as angular or ordinal. They do not consider the correlation between levels.

The subsequent sections are constructed as follows. In section 2, the distance between angular type histograms is defined and arrow representation of two histograms is introduced. In section 3, the $O(b^2)$ computational time algorithm to compute the distance between angular type histograms is described. Finally, section 4 discusses how to apply the distance to handwritten character similarity.

2 Angular Type Histograms

We define the distance as the minimum amount of movements of cells to transform one histogram $H(X)$ to the target histogram $H(Y)$. The definition for ordinal type histograms [2, 3] is as follows:

$$D[H(X), H(Y)] = \sum_{i=0}^{b-1} \left| \sum_{j=0}^i (H_j(X) - H_j(Y)) \right| \quad (1)$$

It is the sum of absolute values of prefix sum of difference for each level. Therefore, the algorithm for finding the minimum distance between two histograms consists of three steps. The first step is to obtain the differences for each level. The second step is to calculate the prefix sum of the differences for each level. Finally, the absolute values of the prefix sums are added.

However, this approach is not suitable for angular type histograms as the first bin and the last bin are neighbors as shown in Figure 1. Consider gradient directions. The values are angular degree where values are 0° to 359° degrees. The difference between levels in angular histograms is always less than or equal to π that is $b/2$. Let $d(x, y)$ be the difference between two levels or samples x and y . Unlike in linear histograms where $d(x, y) = |x - y|$, that in angular histograms is defined as follows:

$$d(x, y) = \begin{cases} |x - y| & : |x - y| \leq \frac{b}{2} \\ b - |x - y| & : \text{Otherwise} \end{cases}$$

For example, the difference between the degree 350° and 1° is 11° not 349° . This difference in measure necessitates a new algorithm to compute the distance in Angular type histograms.

Consider angular type histograms, $H(X)$ and $H(Y)$ as shown in Figure 1. Blocks or cells can move to clockwise or counter-clockwise directions. Each cell movement to the next level in either direction costs 1. The minimum cost required to build the target histogram from a given histogram is the distance. An intuitively appealing way to explain the distance is to use the arrow representation of two histograms as shown in 1. If one establishes an arbitrary one-to-one map-

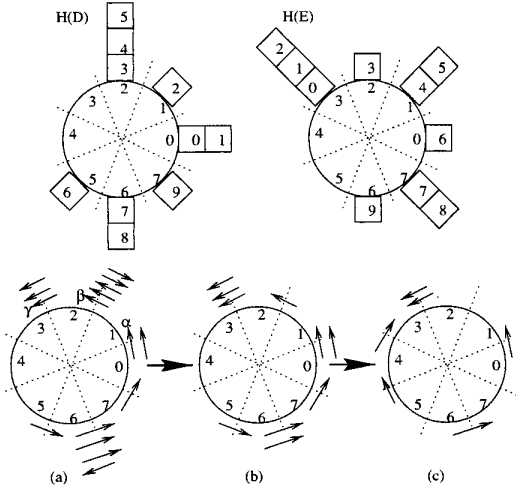


Figure 1. Angular Histograms

ping for the cells between two histograms, one can transform $H(X)$ into $H(Y)$ by moving cells in $H(X)$ to the corresponding position in $H(Y)$. For the example in Figure 1 (a), the arrows α, β and γ indicate the path from the cell 0 in $H(X)$ to the cell 0 in $H(Y)$. There are $n!$ number of ways to transform in this manner. Among these ways, there exists a minimum distance whose number of movements is the lowest. Some sample movements are illustrated as an arrow representation in Figure 1. The more similar two histograms are, the lower the new distance is.

One can observe some properties of an angular representation that gives the minimum configuration of $D[H(X), H(Y)]$. First, arrows must be one directional on each border line. Second, the number of border lines of one direction cannot exceed $b/2$. The configuration in Figure 1 (a) becomes one in Figure 1 (b) by cancelling out the opposite arrows on each border line. Note that these movements are redundant ones. Thus, there exists no border line of mixed directional arrows and each border line has either clockwise or counter-clockwise directional arrows. Suppose that the number of border lines of one direction is $b' > b/2$ and the first number of arrows were k , then after adding the circle, the number of arrows becomes $k + b$. By cancelling out the opposite directions on the same border line, we have $k + b - 2b' < k$. Therefore, if the number of border lines of one direction exceeds $b/2$, this is not a configuration of the

minimum distance.

Two important basic operations are derived. A complete circle in the chain of same directional arrows can be added in any direction and then the opposite arrows on the same line are cancelled out. After transforming, population of cells on each level of both transformed histogram and target histogram must be the same. Let A_l denote the number of counter-clockwise arrows from bin $l - 1$ to l . A_l is positive if arrows are counter-clockwise or it is negative, otherwise. It is a fact that $H_l(X) = H_l(Y) + A_l - A_{l+1}$. Since the equation is true for all levels, affecting cells in one bin means affecting all other bins as a chain reaction. Hence, there are only two possible operations to affect changes. The arrow representation of the minimum distance value is always constructible by the combination of these two basic operations. Consider a non-minimum distance arrow representation. By applying one of two operations of adding a clockwise or counter-clockwise circle, the lower number of arrows is achieved whereas it is unchanged or increased if the arrow representation has the minimum distance. In all, when the arrow representation is the minimum distance, no more reduction occurs when either basic operations are applied.

3 Algorithm

An algorithm to compute the distance between angular type histograms in $O(b^2)$ is presented. It gets an initial arrow representation from the Algorithm for the linear type histograms [2, 3] and then use two basic operations to derive the minimum distance arrow representation that guarantees all properties discuss in the previous section.

Algorithm 1 Distance-angular-histogram(int *x, int *y)

```

1 prefixsum[0] = x[0] - y[0]
2 for  $\forall i$  prefixsum[i] = prefixsum[i-1] + x[i] - y[i]
3 h_dist =  $\sum_{i=0}^{b-1} |\text{prefixsum}[i]|$ 
4 for(;;)
5     d = min(positive prefixsum[i])
6     for  $\forall i$  temp[i] = prefixsum[i] - d
7     h_dist2 =  $\sum_{i=0}^{b-1} |\text{temp}[i]|$ 
8     if h_dist2 < h_dist
9         h_dist = h_dist2
10        for  $\forall i$  prefixsum[i] = temp[i]
11    else break;
12 for(;;)
13     d = max(negative prefixsum[i])
14     for  $\forall i$  temp[i] = prefixsum[i] - d
15     h_dist2 =  $\sum_{i=0}^{b-1} |\text{temp}[i]|$ 
16     if h_dist2 < h_dist
17         h_dist = h_dist2
18        for  $\forall i$  prefixsum[i] = temp[i]
19    else break;
20 return(h_dist)

```

The algorithm is explained using the example shown in Figure 1 along with the following calculations:

$$\begin{array}{r}
2 \ 1 \ 3 \ 0 \ 0 \ 1 \ 2 \ 1 \ \dots \quad (1) \\
1 \ 2 \ 1 \ 3 \ 0 \ 0 \ 1 \ 2 \ \dots \quad (2) \\
\hline
1 \ -1 \ 2 \ -3 \ 0 \ 1 \ 1 \ -1 \ \dots \quad (3) \\
1 \ 0 \ 2 \ -1 \ -1 \ 0 \ 1 \ 0 \Rightarrow 6 \ \dots \quad (4) \\
1 \ 0 \ 2 \ -1 \ -1 \ 0 \ 1 \ 0 \Rightarrow 6 \ \dots \quad (5) \\
0 \ -1 \ 1 \ -2 \ -2 \ -1 \ 0 \ -1 \Rightarrow 8 \ \dots \quad (6)
\end{array}$$

The line (3) is the difference between (1) and (2) (step 1 and 2 in Algorithm 1). The line(4) is the initial arrow representation that is the prefix sum of the difference and the sum of the absolute value of these numbers (step 3). Note that step 1 through 3 are exactly the same as the Algorithm for the linear type histograms [2] that guarantees the property 1.

To ensure the properties 2 and 3, two basic operations (adding a complete cycle of arrows) are applied repeatedly. First, circles of clockwise arrows are added to the current arrow representation until there is no more reduction on the total number of arrows (step 4 through 11). The line (5) is the result of these steps. Next, circles of the counter-clockwise arrows are added in the similar manner (step 12 through 19). The line (6) is the result of adding a circle and the resulting value is greater than the previous one. Therefore, the distance is 6.

The worst case time complexity of the algorithm is $O(b^2)$ time and it can be shown by an example:

$$\begin{array}{l}
H(x) = 5 \ 5 \ 5 \ 5 \ 5 \ 1 \ 1 \ 1 \ 1 \ 1 \\
H(y) = 1 \ 1 \ 1 \ 1 \ 1 \ 5 \ 5 \ 5 \ 5 \ 5
\end{array}$$

The distance is 52 and it takes $\frac{b^2}{2}$. The space required for the algorithm is $O(b)$.

4 Similarity Based Writer Identification

Given a set of reference document and a query document, the *Handwriting Identification problem* is to find out the possible author of the query document among the authors of reference documents. Document examiners or handwriting analysis practitioners find important features to characterize individual handwriting as features are consistent with writers in normal undisguised handwriting [1]. Authorship may be determined by retrieving documents that are similar to the query document due to the following hypothesis.

We consider two modulo histogram features of a character as character level image signatures for identification. We show the experimental results of the earlier defined new measure on two modulo histograms: *Gradient directions* and *angles to the centroid*.

Modulo type histograms are obtained by extracting two angular type data from character images. The first angular data is called *Gradient directions*. They are computed by the *Sobel edge detection mask operators* [12].

$$\text{direction} = \tan^{-1} \frac{S_y(i, j)}{S_x(i, j)} \quad (2)$$

A sample of the gradient direction maps of a character image is shown in Figure 2

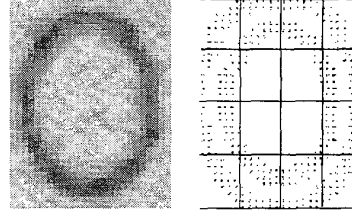


Figure 2. Gradient direction map

Another angular data is the angle to the centroid. First, the centroid of the character is computed as follows:

$$(\bar{i}, \bar{j}) = \left(\frac{\sum_i \sum_j i B(i, j)}{\sum_i \sum_j B(i, j)}, \frac{\sum_i \sum_j j B(i, j)}{\sum_i \sum_j B(i, j)} \right)$$

where $B(i, j) = 1$ if $I(i, j)$ is a black pixel and $B(i, j) = 0$ otherwise.

$$\text{angle to the centroid} = \tan^{-1} \frac{\bar{i} - i}{\bar{j} - j} \quad (3)$$

Some sample letters are shown in Figure 3 as examples (see [8] for description of this character set in detail). Two angular histograms of data (2) and (3) are built for every sample "A" character and the gradient direction histograms are shown in Figure 3.

Next, the distance between two modulo histograms is computed by the algorithm. The smaller the distance is, the more similar two characters are. Tab. 1 is the combined distance matrix for both the histogram for the gradient directions (2) and the angle to the centroid (3). Since the matrices for each feature are symmetry, the lower and upper triangles of the table are dedicated to the gradient directions and the angle to the centroid, respectively. Elements of the resulting matrix M have the following properties.

$$\begin{array}{ll}
M[x, y] = 0 \Leftrightarrow x = y & : \text{identity} \\
\geq 0 \Leftrightarrow x \neq y & : \text{non-negativity} \\
= M[y, x] & : \text{commutative law} \\
\leq M[x, z] + M[z, y] & : \text{triangle ineq.}
\end{array}$$

To show the compatibility of the new similarity measures with the human decision, we asked 25 persons at CEDAR to select similar pairs from the list in Figure 3. It is noted that similar letters in human eyes tend to have small distance values between them. The pairs $\{(3, 4), (11, 12), (8, 9)\}$ are chosen to be similar by more than 15 people and their distances are relatively small. By combining two features, the top three results are $\{(3, 4), (5, 6), (8, 9)\}$ with their corresponding average of two distance values $\{0.44, 0.75, 0.87\}$. Although the pair (11, 12) yields 1.15, when the character 11

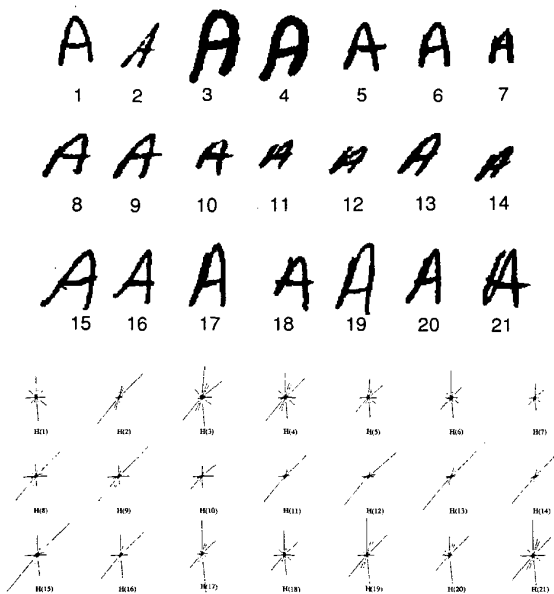


Figure 3. Sample A's and their Angular Representation of gradient direction histograms

is used as a query and one would like to find the closest one, the Tab. 1 suggests the character 12.

When two-dimensional information is represented in the one-dimensional histogram, certain information is lost. Therefore, while it is true that the two histograms from the similar character images tend to be similar, the reverse statement is not always true that two images with the similar histograms tend to be similar. For example, the pair (5,6) has small distances between histograms but they do not really look alike. Yet, the confusion of authorship based on this histogram distance is much smaller than those of slant, stroke width, etc. Together with other features, a writer verification system, called *dichotomizer* is developed with its 97% performance experimental result [5]. The complete comparison among features is given at [4]. Thus, histogram information are very helpful features to determine the similarity of two letters and we claim that distances in Tab. 1 tend to be small if they were written by a same author.

5 Concluding remarks

We have presented the edit distance between two angular type histograms. Conventional measures are not suitable to compute the angular type histograms. There are many direct and indirect applications such as color histogram based indexing, gradient histogram similarity in handwriting analysis, etc. These edit distances can be extended to any application which utilizes the distance between two distributions.

Table 1. Distances between Sample A's
lower triangle is the matrix for the gradient directions
upper triangle is the matrix for the angle to the centroid

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
1	0.0	3.8	1.2	1.0	1.3	0.8	1.8	2.7	2.6	2.9	2.7	2.8	2.8	3.3	2.0	1.9	1.4	1.7	1.6	1.7	1.3
2	4.7	0.0	3.2	3.1	4.0	3.9	3.6	3.7	2.7	3.8	3.1	3.2	2.2	1.9	3.2	2.6	3.7	4.2	3.6	3.3	3.8
3	2.3	2.5	0.0	0.4	1.4	1.2	1.2	2.3	1.9	2.0	2.3	2.4	2.2	2.8	1.5	1.7	1.3	1.7	1.1	1.4	1.3
4	2.6	2.1	0.3	0.0	1.4	1.2	1.5	2.4	1.9	2.0	2.2	2.4	2.2	2.8	1.5	1.5	1.3	1.7	1.1	1.3	1.4
5	1.4	3.3	0.9	1.2	0.0	0.7	1.1	1.7	1.9	2.0	2.1	2.2	2.5	3.2	1.4	2.1	2.1	1.4	2.1	2.0	1.2
6	1.5	3.4	1.0	1.3	0.8	0.0	1.1	2.2	2.3	2.3	2.4	2.4	2.6	3.3	1.7	1.9	1.6	1.4	1.6	1.7	0.9
7	2.1	3.4	1.5	1.8	1.3	1.3	0.0	1.6	1.8	2.0	2.1	2.1	2.4	3.0	1.4	2.1	2.2	1.7	2.1	1.9	1.5
8	4.1	1.4	2.1	1.8	2.8	2.9	2.8	0.0	1.1	1.5	1.5	1.3	2.1	2.7	1.2	2.6	3.4	2.3	3.3	3.0	2.6
9	4.5	0.9	2.3	1.9	3.1	3.2	3.2	0.6	0.0	1.9	1.2	1.1	1.3	1.8	0.9	1.9	2.9	2.6	2.8	2.4	2.4
10	3.4	2.3	1.7	1.4	2.2	2.4	2.7	1.2	1.6	0.0	1.8	1.8	2.3	3.2	1.8	2.6	3.2	1.7	2.8	3.1	2.6
11	5.1	1.1	2.8	2.5	3.7	3.8	3.8	1.3	0.9	2.3	0.0	1.3	1.7	1.8	1.3	2.4	3.2	2.9	3.1	2.9	2.7
12	5.6	1.9	3.6	3.2	4.3	4.4	4.4	1.7	1.5	2.6	1.0	0.0	1.7	1.8	1.5	2.8	3.1	2.9	3.0	2.9	2.6
13	5.7	1.4	3.5	3.1	4.3	4.3	4.5	1.9	1.3	2.7	0.7	1.1	0.0	1.2	1.5	1.8	2.7	3.0	2.4	2.5	2.6
14	5.4	1.3	3.1	2.8	4.0	4.1	3.8	1.6	1.5	2.8	0.7	1.3	1.1	0.0	2.4	2.6	3.2	4.0	3.2	3.0	3.3
15	4.0	1.2	1.7	1.4	2.6	2.6	2.7	0.9	0.9	1.2	1.3	2.0	1.8	1.9	0.0	1.5	2.5	2.0	2.5	2.0	1.8
16	3.8	1.6	1.5	1.3	2.4	2.3	2.5	1.7	1.5	1.4	2.0	2.7	2.0	2.6	0.9	0.0	1.9	1.9	1.9	1.6	1.7
17	2.7	2.2	1.2	1.6	1.4	1.2	1.7	3.1	3.0	2.6	3.5	4.3	3.7	3.8	2.4	1.9	0.0	2.2	0.7	0.9	1.1
18	1.7	3.1	0.9	1.1	0.5	0.8	1.4	2.6	2.9	2.1	3.5	4.1	4.1	3.7	2.4	2.2	1.6	0.0	1.1	2.3	1.4
19	3.1	2.6	1.2	1.4	1.7	1.7	2.4	3.0	2.9	2.4	3.3	3.9	3.0	3.7	2.1	1.3	1.0	1.8	0.0	1.3	
20	2.9	1.9	1.1	1.2	1.6	1.7	1.6	2.1	2.0	2.0	2.6	3.4	3.1	2.9	1.6	1.2	1.2	1.7	1.5	0.0	1.1
21	2.8	2.6	1.1	1.1	1.6	1.4	2.4	2.6	2.5	2.0	3.0	3.7	3.1	3.5	1.8	1.2	1.4	1.5	0.7	1.6	0.0

References

- [1] R. R. Bradford and R. B. Bradford. *Introduction to Handwriting Examination and Identification*. Nelson-Hall Publishers: Chicago, 1992.
- [2] S.-H. Cha. Fast image template and dictionary matching algorithms. In *Proceedings of ACCV'98 LNCS - Computer Vision*, volume 1351, pages 370–377. Springer-Verlag, Jan 1998.
- [3] S.-H. Cha. Efficient algorithms for image template and dictionary matching. *Journal of Mathematical Imaging and Vision*, 12(1):81–90, February 2000.
- [4] S.-H. Cha and S. N. Srihari. Multiple feature integration for writer identification. In *Proceedings of 7th IWFHR 2000*, pages –. Springer-Verlag, September 2000.
- [5] S.-H. Cha and S. N. Srihari. Writer identification: Statistical analysis and dichotomizer. In *Proceedings of SS&SPR 2000*, pages –. Springer-Verlag, September 2000.
- [6] R. O. Duda, D. G. Stork, and P. E. Hart. *Pattern Classification and Scene Analysis*. John Wiley & Sons, Inc., 2nd edition, 2000.
- [7] M. Flickner, H. Sawhney, W. Niblack, J. Ashley, Q. Huang, B. Dom, M. Gorkani, J. Hafner, D. Lee, D. Petkovic, D. Steele, and P. Yanker. Query by image and video content: the qbic. *Computer*, 28(9):23–32, Sept 1995.
- [8] J. J. Hull. A database for handwritten text recognition research. *Transactions on Pattern Analysis and Machine Intelligence*, pages 550–554, May 1994.
- [9] G. Pass, R. Zabih, and J. Miller. Comparing images using color coherence. In *ACM International Multimedia Conference*, pages 65–73. ACM, 1996.
- [10] H. S. Sawhney and J. L. Hafner. Efficient color histogram indexing. In *International Conference on Image Processing*, volume 1, pages 66–70. IEEE, 1994.
- [11] M. J. Swain and D. H. Ballard. Color indexing. *International Journal of Computer Vision*, 7(1):11–32, Nov 1991.
- [12] S. E. Umbaugh. *Computer Vision and Image Processing*. Prentice Hall PRT, 1998.