

A Natural Learning Algorithm based on Hough Transform for Text Lines Extraction in Handwritten Documents

Yao Pu, Zhixin Shi

*Center of Excellence for Document Analysis and Recognition
State University of New York at Buffalo, Buffalo, NY 14260, U.S.A.
E-mail: zshi@cedar.buffalo.edu*

A text lines extraction algorithm based on Hough Transform for unconstrained handwritten documents is proposed. A natural learning method which is similar to the learning procedure of human being is used in this algorithm. First some general informations have been obtained by applying Hough Transform. Then natural learning is used during the clustering. The algorithm gets more information step by step by natural learning until the final segmentation result is reached.

1 Introduction

Segmentation of a document image into blocks and lines is a necessary step in optical reading before applying symbol recognition. A text line has a linear structure which can be described as an ordered string of characters or words. The goal of line segmentation is to extract the lines from a document image with the linear structure preserved.

Many methods have been developed and successfully applied to the detection of text lines for machine printed documents. The most popular methods include projections, runlength smearing and Hough Transform[3]. These method could also be used for handwritten documents in the case of regular or simplified texts. Due to the nature of handwritten documents, methods developed for handwritten documents must adapt to the layout variability of each kind of the documents. For constrained documents, information including *a priori* given zones for the text to appear, positions and orientations of the text lines is used to drive the line extraction.

For unconstrained handwritten documents, the information such as the number of lines and the distance between lines is totally unknown. Unsuper-

vised learning must be applied first to get at least the initial information. Most existing algorithms are using only unsupervised learning. For human beings, however, after we have some basic clusters information, i.e. the classification of certain objects, the learning approach is actually the combination of unsupervised learning and supervised learning, called nature learning. When a new object appears, if it is very similar to an existing cluster, it will be assigned to this cluster. If it is different to all the existing clusters, it will be put aside for a while. After similar objects appear many times, a new cluster will be generated according to those relational information.

In this paper, a similar learning approach is proposed for extracting text lines in handwritten documents. The algorithm is based on Hough transform which has been successfully applied in detection of text lines in printed documents. For handwritten documents, different orientation of lines may appear in a page. The orientation of each line is fluctuating rather than of same angle, the spaces in between lines are irregular. Those characteristics limit the use of Hough transform. However, another observation is that the angle change of each line in a small area is small. The orientation of the line is of a fluctuation of main orientation of the line in the small area. Therefore Hough transform is more suitable to a small strip of text image rather than to a whole page.

The minima points of the connected components are used to represent the text lines. The natural learning algorithm based on Hough transform is first applied to the small strip of an image. The clusters computed from the small strip will then be taken as initial clusters to the whole page. The same learning approach will be used until the end of segmentation is reached.

2 Natural Learning Algorithm

The nature learning algorithm for text lines extraction can be described as the following:

1. Apply Hough transform to the minima points of the connected components in a small strip of the image to get the initial lines.
2. Cluster the minima points using moving window algorithm. Minima points will be put into the existing cluster, or put into a pool.
3. Check the minima points in the pool. Some of them may be put into the existing cluster by image information. To those which can not be put into any existing clusters, Hough transform will be applied again to generate new clusters. Then same algorithm in step 1 will be used to select lines. And then algorithm in step 2 will be applied to all minima points in this small strip of image.

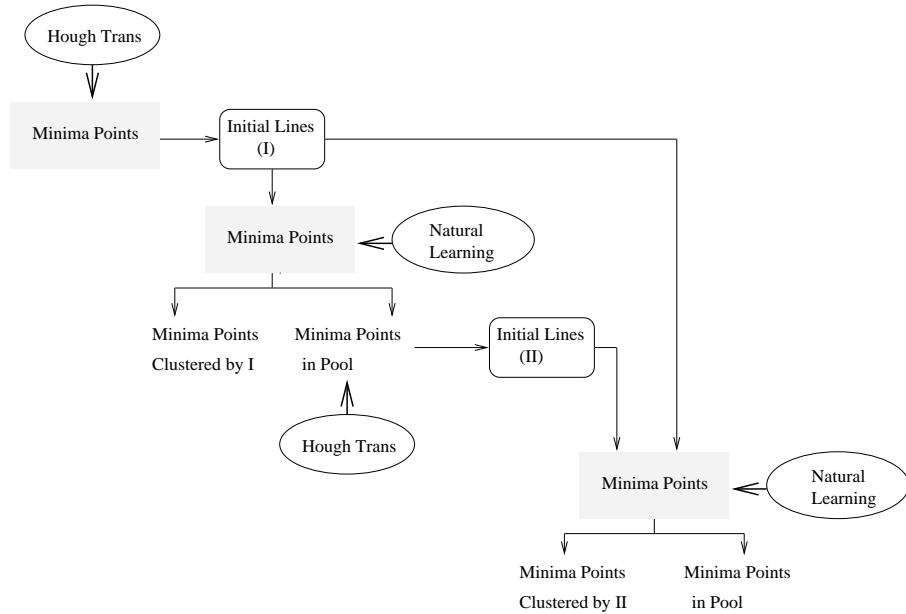


Figure 1: Natural Learning Algorithm.

4. Those clusters which have been generated by steps 2 and 3 will be regarded as the initial clusters to whole text image. Repeat steps 2 and 3 until the end of segmentation.

The key idea of natural learning algorithm is shown in Figure 1.

2.1 Lines Selection by Hough Transform

In first step, the left one fourth of the image is chosen to get the initial lines. Generally, we should choose a small strip of the image containing the most information. This can be achieved by choosing a strip which has highest black pixel density (if black pixels are used to carry the information).

Hough transform of a point (x, y) in the two dimensional plane is defined to be any point (r, θ) in the parametric plane by the following relation:

$$r = x * \cos \theta + y * \sin \theta \quad (1)$$

For any fixed point (x, y) , there are infinity many (r, θ) to be its Hough transform. Each (r, θ) pair corresponds to a straight line through the point (x, y) with equation in the form of (1). Points lying on a same straight line share a common Hough transform.

For the selected strip of the document image, we construct a discrete Hough domain by taking the resolution of θ to 1 degree, and the resolution of r to $R = 0.2 * H_w$ [3], where H_w is the average y-value difference between two consecutive maxima and minima of the connected components. The Hough transform domain for the minima points of our small strip is determined. An accumulation matrix is constructed as following:

For all minima points (x, y) in the small strip
 For $\theta=0$ to 180 degree
 $\{$
 $r = x*\cos\theta + y*\sin\theta;$
 $\text{accum}(r/H_w, \theta) = \text{accum}(r/H_w, \theta) + 1;$
 $\}$

Each entry of the matrix being an integer number corresponding to a pair of discrete value of (r, θ) is the number of possible points on the straight line determined by the pair.

We first compute the main orientation for the lines in the strip. As been discussed in [1] if a document is skewed along a particular orientation, then the value of θ when scanned for all possible value of r yields the maximum number of nonzero-zero-nonzero transition in the accumulation matrix. The nonzero values correspond to the extrema points which make the text lines, and the zero vales correspond to the white space between lines. To get the orientation, we consider the following function. Similar function was also given in [1],

$$R(\theta_i) = \sum_{j=2}^{N_r} [\text{accum}(r_i, \theta_i) - \text{accum}(r_{j-1}, \theta_i)]^2 \quad (2)$$

where $j = 1, 2, \dots, N_r$ are row index for the discrete r values and $i = 1, \dots, N_\theta$ are column index for the discrete θ values in the accumulation matrix. The main orientation as the particular angle is then calculated:

$$\theta_{main} = \max_i [R(\theta_i)].$$

After the main orientation of the strip image is computed, a part of accumulator matrix (corresponding to $\theta_{main} \pm 10$ degree) is extracted to be a new matrix A . We average each entry of A by its eight neighbors. Note this matrix A itself has a very good natural grouping characteristics.

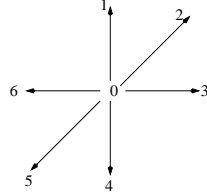


Figure 2: Directions for Line Select in Hough Domain.

Groups of non_zero elements are separated by groups of zero elements. They may overlap in horizontal direction. Each non_zero group corresponds to a line in the image domain.

Unlike other algorithms which used the information of images, such as the average height of components to estimate the cell size[3], an exhaust searching algorithm fully using the character of non_zero groups separated by zero groups is designed to catch those non_zero groups.

Each element in matrix A has eight neighbors in three rows. Due to the shape of each non_zero group we can percept, only six directions (Figure 2) are considered.

The search algorithm is described as following:

1. Generate a new all zero array B, with the same size of A to remember which element in A has already been searched. If the element has been searched in A, the same position element in B will be marked 1.
2. Find the biggest element (r_i, θ_i) in A which has not been marked in B.
3. Row search.

Search the nonzero elements in A in 3, 6 directions and mark them until two continuous zero is reached, or the edge of matrix A is reached, or a marked element is reached.

4. Up search.

From element (r_i, θ_i) , we consider the direction 2 first, if 2 is non_zero and not marked, apply row search to it. If it is zero or marked, we will consider direction 1. If both of them are zero, or have been marked, or edge of the matrix is reach, stop up search.

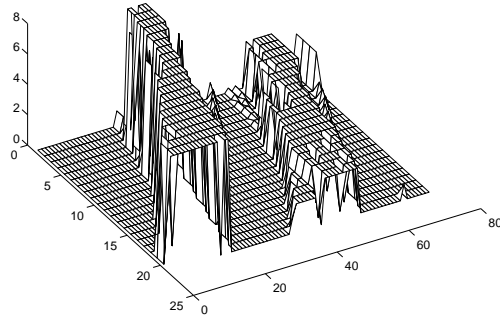


Figure 3: Two groups after exhaust search.

5. Down search.

Down search is symmetric to up search.

6. Find the r and θ to represent the line corresponds to a group.

All the elements connected by row search, up search, down search will be marked and be considered as a group. In order to tolerate the variation of line angle, The (r, θ) pair of the largest element in the group will be selected to represent this line. The number of minima points corresponding to the marked group is the size of this group, if it is bigger than a threshold, the line generated by this group is valid. Here we give a big threshold to minimize $P(\text{correct}|\text{error})$. The losing lines will be added after clustering discussed later. After the line has been selected, some image information will be calculated, such as the average height of the components in this group. Those information will be used later.

7. go to step 2 until all the elements in A have been searched.

Part of exhaust search result is show in Figure 3. There are two groups in this figure corresponding to two possible lines.

2.2 Moving Window Clustering

Natural learning algorithm will be applied based on the lines selected by Hough transform. It means that not all the minima points will be clustered to existing

lines determined by Hough Transform. A moving window is used to trace the angle change of each line.

The same number of clusters as the number of lines are generated. The first 5 points of each cluster are added according only to the line information (r, θ) . Then moving window clustering is applied, i.e. whether a new minima point belongs to this cluster or not will be decided by its 5 nearest points in this cluster (not by (r, θ)). If the distance (weighed x distance and y distance) is less than a threshold which is related to average height of components in this cluster, it will be assigned into this cluster. (This threshold is dynamically calculated based on the connected components in the cluster.) Otherwise it will be put into data pool.

2.3 Check the Minima Points in the Pool

The minima points in the pool will be checked again combining the information in the image domain. If minima points are on a component whose most minima points have already been clustered in a cluster, it will be assigned to this cluster. Otherwise, it is still in pool. By doing this, the minima points on the descender of character, such as g, q, will be clustered.

After the checking, Hough transform will be applied to the minima points to select lines in Hough domain. And then moving window clustering algorithm will be applied to those minima points. After that, there may still have some minima points in the pool. Those minima points may be a very small component such as the T hat or the point on i , or a new line begins at the edge of this small strip of image.

The result is shown in Figure 4. A line is added by checking the data in the pool. Minima points that marked by square and “+” are the points belong to this new line.

2.4 Deal with Whole Page

Then those clusters will be used as initial cluster to the whole page. The result is shown in Figure 5. After we have finished the moving window clustering algorithm. There are no points in the pool which can generate lines.

However, there may still be some points marked by “+” in the pool. They can generate new lines by using line selection algorithm. Two lines are added. The final result is shown in Figure 6.

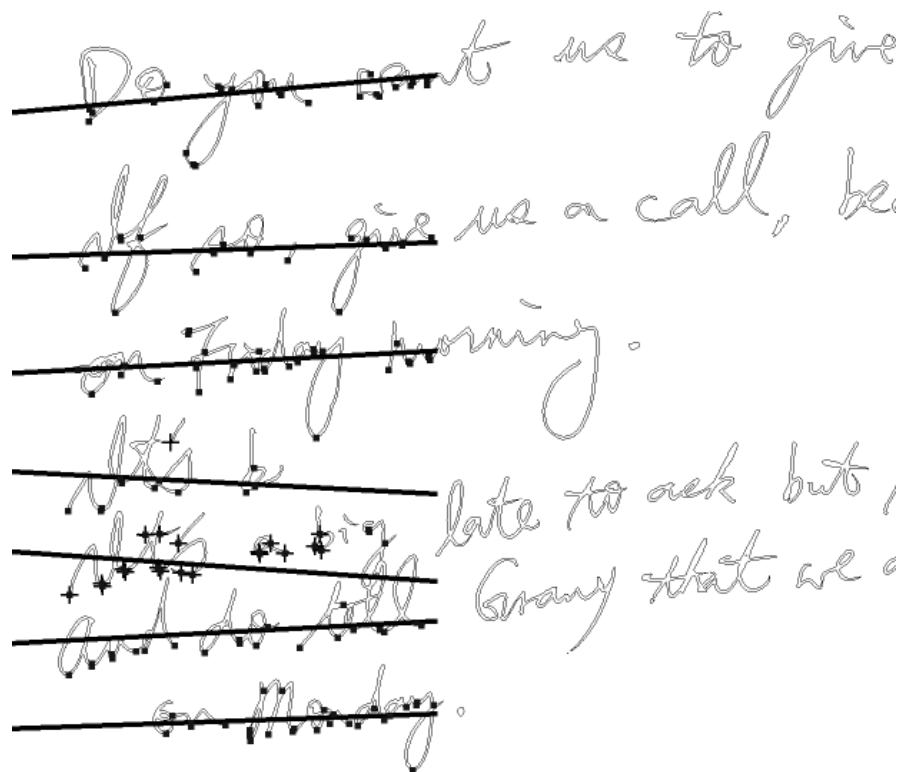
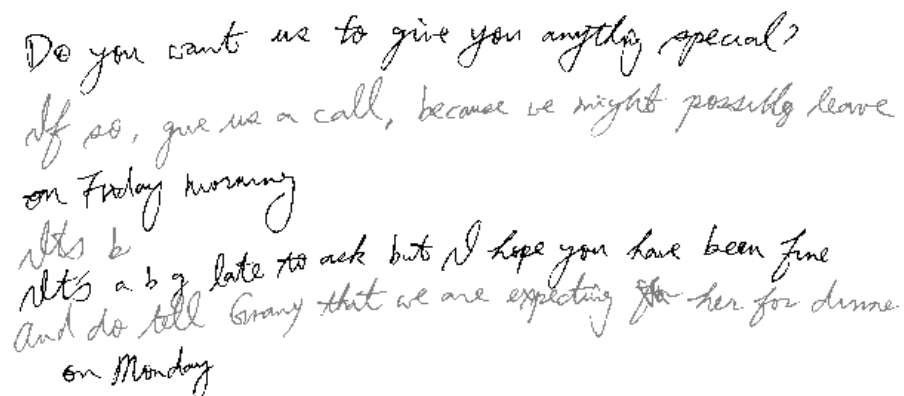


Figure 4: Result after moving window clustering check the minima points in pool.

3 Conclusion

The line selection algorithm described in this paper has been implemented. It has been tested on more than 100 handwritten text page images. Although most of the testing pages contain text lines with roughly horizontal orientation, some very difficult cases found with interfering lines, short lines starting in the middle of the pages and etc.

This algorithm tries to get the local information of each line and traces the variety of handwritten lines by natural learning method. The thresholds used in this algorithm are changing rather than constant. In this algorithm we used the first left strip to get the initial information. Further we should develop an



Do you want us to give you anything special?
If so, give us a call, because we might possibly leave
on Friday morning
It's a b g late to ask but I hope you have been fine
And do tell Grady that we are expecting her for dinner
on Monday

Figure 5: Final result 1.

algorithm to determine which part of the image we should start from to get most useful initial information.

References

1. Sargur N. Srihari and Venugopal Govindaraju, Analysis of Textual Images Using the Hough Transform, *Machine Vision and Applications*, pages 141–153. Springer-Verlag, New York, 1989.
2. Shapiro V., Gluhchev G., Sgurev V. Handwritten document image segmentation and analysis, *Pattern Recognition Letters*, No 14, pages 71–78, 1993.
3. Laurence Likform-sulem, Anahid Hanimyan, Claudie Faure, A Hough Based Algorithm for Extraction Text Lines in Handwritten Document, *Proceedings of the International Conference on Document Analysis and Recognition*. pages 774–777, 1995.
4. Icher L.A., Kasturi R. Text string segmentation from mixed text/graphics images, *IEEE PAMI*, Vol 10, No 3 pages 910–918, 1988.

~~Plot~~ ~~story~~
 "When Harry Met Sally" is an entertaining, romantic comedy.
 It takes place in New York City.
 The main characters are old acquaintances who meet again by chance.
 They find that they enjoy each other's company.
 At the same time, each has romantic relationships with other people.
 But the main characters ~~seemingly~~ ~~are~~ ~~happy~~ ~~to~~ ~~be~~ ~~together~~ ~~2 of them are~~
 together.
 They agree ~~on~~ ~~what~~ ~~the~~ ~~or~~ ~~not~~ ~~they~~ ~~are~~ ~~in~~ ~~love~~.
 Their friends think that they are perfect for each other.
 The ~~man~~ ~~and~~ ~~the~~ ~~woman~~ ~~are~~ ~~falling~~ ~~in~~ ~~love~~
 The reason ~~is~~ ~~that~~ ~~the~~ ~~man~~ ~~rejects~~ ~~the~~ ~~idea~~ ~~that~~ ~~is~~ ~~going~~
 to ~~marry~~ ~~her~~.
 Finally, the man apologizes and admits that he loves the woman.
 They ~~are~~ ~~in~~ ~~love~~ ~~and~~ ~~they~~ ~~are~~ ~~in~~ ~~love~~.
 The movie ends with an interesting question.
 It ~~asks~~ ~~us~~ ~~to~~ ~~think~~ ~~about~~ ~~whether~~ ~~it~~ ~~is~~ ~~possible~~ ~~for~~ ~~two~~ ~~people~~ ~~to~~ ~~be~~ ~~friends~~.
 The movie ~~ends~~ ~~with~~ ~~an~~ ~~answer~~.

Figure 6: Final result 2.