

Active Character Recognition Using “A*-like” Algorithm

Jaehwa Park and Venu Govindaraju
Center of Excellence for Document Analysis and Recognition (CEDAR)
Department of Computer Science and Engineering
State University of New York at Buffalo
Amherst, NY 14260

Abstract

This paper describes an Active Character Recognition methodology, henceforth referred to as ACR. We present in this paper a method that uses an active heuristic function similar to the one used by A search algorithm that adaptively determines the length of the feature vector as well as the features themselves used to classify an input pattern.*

ACR adapts to factors such as the quality of the input pattern, its intrinsic similarities and differences from patterns of other classes it is being compared against and the processing time available. Furthermore, the finer resolution is accorded to only certain “zones” of the input pattern which are deemed important given the classes that are being discriminated.

Experimental results support the methodology presented. Recognition rate of ACR is about 96% on the NIST data sets and the speed is better than traditional classification methods.

1 Introduction

In this paper, we present an active character recognizer which uses a hierarchical feature space to achieve multi-resolution in an adaptive manner. The recognizer begins with features extracted in a coarse resolution and focuses on smaller sub-images of the character on each recursive pass thus effectively working with a finer resolution of a sub-image each time till the classification meets acceptance criteria. In other words, successive passes render the same features at a finer resolution.

The sub-images are selected by a function that actively guides the successive passes through a series of sub-images where it believes the features of interest are present. This is done in an adaptive and dynamic manner. Depending on the classes competing for the top spot different sub-images might be selected at any given stage. Figure 1 illustrates this advantage. If the contending classes are ‘3’ and ‘5’, the

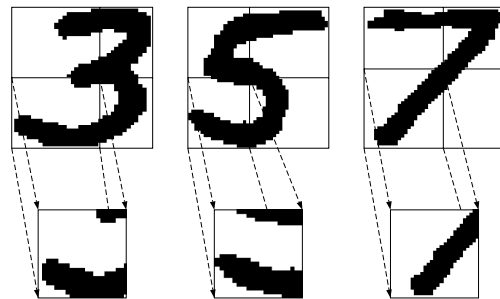


Figure 1. Discriminating characters by looking for relevant features in the sub-images

lower left sub-image does not carry any information that will help in discriminating the two classes. However, both ‘3’ and ‘5’ can be separated from a ‘7’ by that very sub-image.

1.1 Background

Most character recognition techniques described in the literature use a “one model fits all” approach, *i.e.*, a set of features and a classification method are developed and every test pattern is subjected to the same process, irrespective of the constraints present in the problem domain.

The digit recognition method described by Favata et al [2] is a classical method that typifies the passive character recognizers. It uses a 512 binary feature vector. Multi-resolution is achieved by choosing features of three different scales. The Gradient, Structural, Concavity, (GSC) features are the symbolic multi-resolutional features. Gradient of the image contour captures the local shape of a character. The Gradient features are extended to Structural features by encoding the relationships between strokes. Concavity features capture the global shape of characters. **G** features are the finest and the **C** features are the coarsest. Features at the 3 levels, **G**, **S**, **C** are combined in a k-nearest neigh-

bor classification method to produce a multi-resolution digit recognizer.

We note that the features and hence the feature vector is fixed and so is the classification method. There is no adaptation to the problem constraints. For example, it is possible that certain applications call for only the 2 class problem $\{0,1\}$ as would be in the case of dates¹. Passive methods, such as the GSC, handle such constraints only in the post-processing stage. After the digit is recognized as any of the 10 classes, the results are examined to see if '0' or '1' comes on the top. Clearly, comparing with all the classes increases the computational time unnecessarily. Further, features from certain zones of the images are not useful (Figure 1) in the discrimination of the 2 classes. In fact, considering such zones only adds to the confusion.

1.2 Objective

Our objective is to design a character recognizer that dynamically decides the following parameters based on the constraints presented by any instance of recognition.

1. Length of the feature vector
2. Subset of features chosen from a set of features
3. Resolution of features
4. Termination condition of acceptance

2 Methodology

A method that dynamically takes advantage of the constraints is called for. ACR is such a recognizer. The central idea underlying this methodology is the following. We will use a small set of features (8 in our implementation) that can be computed from a pattern image. These same features can be computed from different sub-images derived from the original image. Sub-images are defined by the quad tree rule [3].

The bounding box of a character image is divided into four rectangular regions (Figure 2c). The center of mass of the contour is first computed. A vertical and horizontal line through the center of mass determines the four regions. Subsequent layers of the quad tree are successively constructed by the same method.

The exact same features are computed for each node (sub-image) of the quad tree that is visited.

Let $I(d, l)$ represent the sub-image of the l^{th} node of a tree at depth d , where $d \in \{0 \dots N_D - 1\}$ and $l \in \{0 \dots 4^d - 1\}$. N_D is the maximum allowable depth of the quad tree, 4 in our implementation.

¹If the first digit in a date is recognized as a '3' then the second digit can only be a '0' or a '1'.

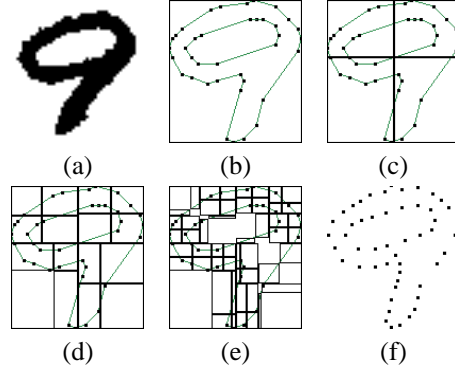


Figure 2. Character image and sub images of quad tree: (a) original image, (b) piece-wise linear contour, (c) layer 1, (d) layer 2, (e) layer 3, and (f) reconstruction of contour image from critical points of final layer

Let $\mathbf{x} = [x_0, x_1, \dots, x_{N_x-1}]$ be a binary feature vector with dimension N_x and

$C = \{c_0, c_1, \dots, c_{N_c-1}\}$ be a set of N_c classes, then by Bayes' rule the *a posteriori* probability of c_i given \mathbf{x} is given by

$$P(c_i|\mathbf{x}) = \frac{p(\mathbf{x}|c_i)P(c_i)}{P(\mathbf{x})} = \frac{p(\mathbf{x}|c_i)P(c_i)}{\sum_j p(\mathbf{x}|c_j)P(c_j)}, \quad (1)$$

where $P(c_i)$ is the *a priori* probability of class c_i and $p(\mathbf{x}|c_i)$ is the class conditional probability.

A vector code space $V = \{0, 1, 2, \dots, 2^{N_x} - 1\}$, can be defined where each point in V represents a different binary bit pattern of size N_x , thus mapping every binary feature vector, \mathbf{x} onto a unique point $v \in V$ in the vector code space. The probabilities $P(c_i|v)$, denoted as $P_i[v]$ can be stored in a look-up table of size $|V| \times N_c = 2^{N_x} \times N_c$ during training.

Let $v(d, l)$ represent the feature vector corresponding to sub-image $I(d, l)$. From the training image set, probabilities $P(c)$ and $P(v(d, l)|c)$ in each sub-image are obtained. Probabilities, $P(c_i|v(d, l))$ in every sub-image are calculated and stored in a look-up table in row v and column (d, l) (Table 1). This value is obtained from Equation 1, where \mathbf{x} is the feature vector computed from the sub-image, $I(d, l)$ and c_i is replaced by $c_i(d, l)$.

The sub-images are extracted dynamically, wherein the order in which the nodes of the quad tree are expanded is determined by an active function that is similar to A* heuristic function.

The following aspects of ACR are noteworthy. They establish our methodology as being ***active*** and meet all the objectives set forth in Section 1.2

v	$I(0,0)$	\dots	$I(N_D, 4^{N_D} - 1)$
0	$P_0[0], \dots, P_{N_c}[0]$	\dots	$P_0[0], \dots, P_{N_c}[0]$
1	$P_0[1], \dots, P_{N_c}[1]$	\dots	$P_0[1], \dots, P_{N_c}[1]$
2	$P_0[2], \dots, P_{N_c}[2]$	\dots	$P_0[2], \dots, P_{N_c}[2]$
3	$P_0[3], \dots, P_{N_c}[3]$	\dots	$P_0[3], \dots, P_{N_c}[3]$
\vdots	\vdots	\ddots	\vdots
v	$P_0[v], \dots, P_{N_c}[v]$	\dots	$P_0[v], \dots, P_{N_c}[v]$
\vdots	\vdots	\ddots	\vdots
$2^{N_x} - 1$			

Table 1. 3-dimensional look-up table: Each column represents a sub-image $I(d, l)$. Each entry in the table represents N_c probability values, one for each class. $N_c = 10$ in the digit recognition problem.

1. Since a different number of sub-images (nodes of the quad tree) are used at any given instance, the length of the feature vector is variable and given by $|N_x| \times$ (number of nodes visited).
2. The composition of the feature set is a variable because the order in which nodes are visited and the number of nodes visited varies.
3. The termination condition is dependent on the classes competing for the top spot at any given time.
4. Features extracted in earlier nodes (sub-images) correspond to coarse global features while those in lower nodes (sub-images) provide local features. The actual feature extraction process and the type of the features remains the same, however, the features are now computed from a smaller portion of the image.

2.1 Active Function (A*-like)

A confidence function Q_i^d for class c_i is defined as the product of the *a posteriori* probabilities corresponding to all sub-nodes belonging to layer d ,

$$Q_i^d = \prod_{l=0}^{4^d-1} P_i[v(d, l)]. \quad (2)$$

The confidence Q_i of a class c_i considering all the layers of the tree is defined as the product of all the *a posteriori* probabilities of all its sub-classes corresponding to all sub-nodes processed up to that time, *i.e.*,

$$Q_i = \prod_{d=0}^{N_D-1} Q_i^d = \prod_d \prod_l P_i[v(d, l)]. \quad (3)$$

If the classes and sub-classes are statistically independent then Q_i represents the class *a posteriori* probability. This is clearly not the case (since the sub-images are derived from the parent image). The value of Q_i lies in the interval $[0, 1]$.

Active Function is a means of expanding only some of the nodes in the tree. A node is expanded only if its *a posteriori* probabilities do not satisfy a termination criterion *e.g.*, the highest class probability is not sufficiently larger than the second highest class probability (Eq. 5). We can fix the number of nodes whose probabilities are used to update confidences on each pass and we limit the selection of new nodes to a subset of *potential nodes* used in the previous pass but has unprocessed child nodes or sub-nodes (Figure 3). The intermediate confidences can be computed recursively at pass t as follows: We use the notation $Q_i(t)$ to indicate the recursive confidence computations when all sub-nodes are not used.

$$Q_i(t) = Q_i(t-1) \cdot Q'_i(t-1), \quad (4)$$

where $Q'_i(t)$ is the confidence generated from the sub-nodes of a node selected at time t , where $t \in \{0, 1, \dots, \sum_{d=1}^{N_D-1} 4^{d-1} - 1\}$. This process is illustrated for a quad tree in Figure 3 where at depth 1, only one node is expanded.

Usually a classification result of a recognizer is accepted based on the confidence of the top choice and its separability from the second choice [2, 4]. Using the same criterion, normalized confidence, $\kappa_i(t)$, of a class c_i at recursion time t is defined by the ratio of the current (intermediate) confidence and the maximum possible confidence of the class regardless of the specific feature vector taken from the look-up tables of the all corresponding processed and potential nodes.

$$\kappa_i(t) = \frac{Q_i(t)}{R_i(t)} \quad (5)$$

$R_i(t)$ is the product of the maximum possible probability of class c_i in any vector code through the selected node sequence of processed and potential nodes.

Separability between top two classes is defined as,

$$\sigma(t) = \frac{Q_{top}(t)}{Q_{second}(t)} \quad (6)$$

where class c_{top} is the top choice and class c_{second} is the second choice in intermediate confidence result of time t .

Local separability $\sigma^{(d,l)}(t)$ of a potential node (d, l) is given by

$$\sigma^{(d,l)}(t) = \frac{P_{top}[v(d, l)](t)}{P_{second}[v(d, l)](t)} \quad (7)$$

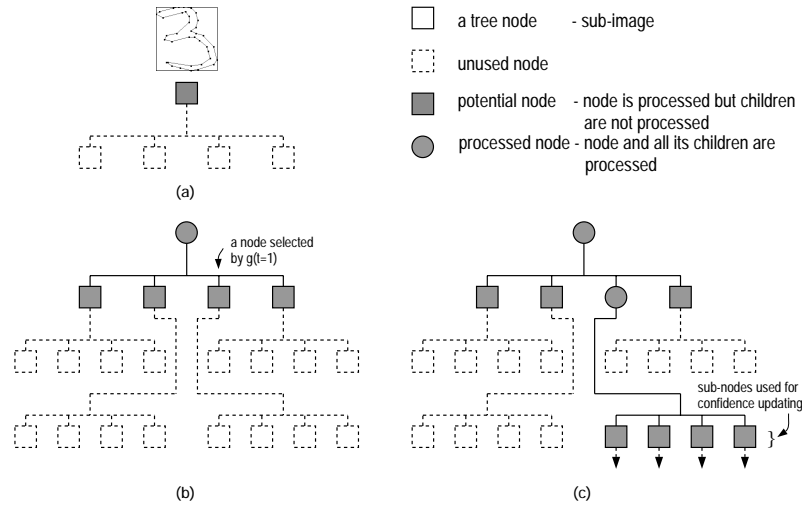


Figure 3. Confidence update recursion using a quad-tree; (a) starting from root, (b) after one recursive pass, (c) after two recursive passes.

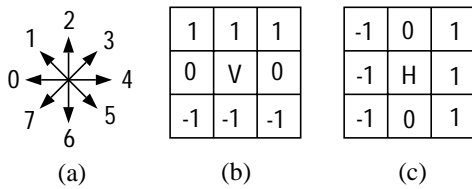


Figure 4. Contour representation; (a) contour slope division, (b) vertical weight mask and (c) horizontal weight mask

If the $\sigma(t)$ does not meet acceptance threshold recursive computations continue. The next node to be expanded can be determined by comparing local separability of the remaining nodes using the active function. Thus, on each pass the local separability $\sigma^{(d,l)}(t)$ serves as the A*-like heuristic in determining the new node to expand

2.2 Implementation

A binarized character image is represented by *piecewise linear contours* using chain code [5] (Figure 4a). A window of weights is convolved with a contour component and its neighbors to obtain the curvature at the contour point. The window size is determined by a function of the height of the bounding box of the character image and the density of black pixels within the bounding box.

At each contour point, horizontal and vertical slopes are found using the masks of Figure 4b and c and binarized using a threshold (0.65 in our implementation). Contour points which have high curvature (derivative of the bina-

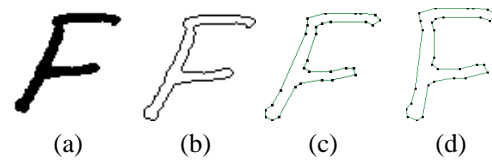


Figure 5. Contour representation and pre-processing; (a) an input image, (b) pixel based contour representation, (c) piecewise linearized contour representation and (d) piecewise linear contour representation after slant correction

riized slope) are identified as critical points. Contour segments between adjacent critical points are approximated to be linear. This representation adds to the overall processing speed by reducing the total data points and also by eliminating the need for smoothing and noise removal. Figures 5-(b) and (c) illustrate the procedure.

2.2.1 Slant correction

Continuous vertical line components (spanning several critical points) whose length is greater than a threshold are identified. Slant angle is estimated by the average of the slant angles (from the horizontal) of all vertical line components, where the angles are weighted in proportion to their length in pixels. Using the estimated slant angle, the critical points are mapped for slant correction. Figure 5d shows the translated points after slant correction.

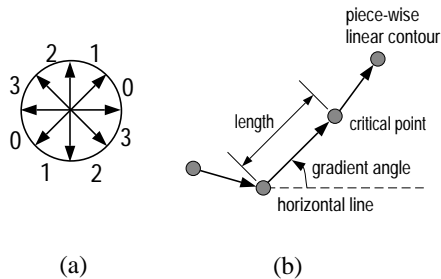


Figure 6. Gradient feature generation: (a) gradient division and (b) gradient angle measure

2.2.2 Features

Histogram of gradient and moment based projections are used as features. Feature vector for each sub-image is generated after binarizations of the feature measurements which are extracted from corresponding sub-regions. 8 bit binary feature vector is used.

Contour slope angles are quantized in pairs of $\pi/4$ slots that are opposite each other labeled 0 to 3 in Figure 6a. The angle between the horizontal line and the linearized contour piece between two adjacent critical points is assumed to be the slope angle of the contour segment and the length of a contour piece is taken as the Euclidean distance between adjacent critical points. (Figure 6b).

The ratio of the total contour length corresponding to the contour segments whose slope angle falls in division i and the total contour perimeter in that sub-image gives 4 gradient features. 4 moment features are obtained from the distribution of the contour within the sub-images. The ratio of the length of the sum of contour segments which are present within a sub-image to the total contour perimeter generate the moment features. In our implementation, all the thresholds were set as 0.25 uniformly for feature binarization.

3 Experiments

Experiments were performed using isolated digit and character images present in the NIST database 3 (SD 3) [6], and CEDAR database collected from handwritten postal addresses [7]. The images were divided into mutually exclusive training and testing sets for various experiments shown in Table 2. A 782 MIPS *SUN ULTRA 2* is the experimental platform.

The number of sub-images are 85 (1+4+16+64) yielding a maximum of 680 dimensional binary feature vector. The recursion of our recognition method can be terminated at any time by using various degrees of separability between top two choices as a terminal condition (Eq. 6)

Table 3 shows a comparison of different classification methods on data set 3. To keep the comparison meaningful, the same features have been used by all 3 classification methods: (i) ACR, (ii) a neural network method (NN), and (iii) a nearest neighbor method (k -NN). The neural network is a fully connected. It has three input layers of 680 units, one hidden layer of 353 units, and an output layer of 26 units. The k -nearest neighbor method is implemented using a Hamming distance metric and 6 nearest neighbors. Storage requirements for ACR method refers to the size of the lookup table. In case of the neural network it refers to the total number of weights and in case of the nearest neighbor classifier it refers to the prototypes in the training set. The relative merits of each of the three methods are highlighted in the table. The hierarchical OCR certainly has the speed advantage.

4 Discussion

Most character recognition techniques described in the literature use a “one model fits all” approach. We argue that such methods are passive. Most recognition methods that adhere to this paradigm have pre-determined dimension of feature vector and equal computational resources available to every test pattern. In very general terms, a large feature vector in multi-resolution provides higher accuracy at the cost of computational resources. However, approaches that mix features in high dimension have several disadvantages as well, such as correlation, common distance metric, computational complexity, etc.

We have shown a method of achieving the benefits of a multi-resolution feature space without missing features at different resolutions and also by keeping the feature vector small.

We have defined the confidence of a class based on the assumption that features in sub-images of an input pattern image are strongly correlated and the degree of correlation is uniquely characterized by the class. In other words, we expect the features to be highly correlated from the viewpoint of the true class and correlated to a lesser degree or not at all when considered from the viewpoint of the “non-true” classes. In equation 2, we keep updating the confidence values (Q_i) from successive recursive cycles. This would indeed lower the updated confidence of a class. However, based on the correlation assumption made, the confidences of the classes that are not true should decrease at a more rapid rate. Hence the separability between the true class (top choice) and the non-true class (second choice) will increase. This justifies the subsequent recursive cycles and using the product of confidences. The assumption of correlations is based on empirical observations. It is supported by the recognition accuracy of the method which is comparable to those described in [2, 8] that have trained and tested

Data set	No. of classes	Training images	Testing images	Image source(s)	Training set		Testing set	
					top 1	top 2	top 1	top 2
1	10	25656	12242	Postal+NIST	99.5	99.9	96.2	98.8
2	10	159228	53301	NIST digit	98.6	99.7	98.2	99.5
3	26	34961	11481	NIST alpha upper-case	99.2	99.9	96.4	98.8
4	26	29059	9623	NIST alpha lower-case	98.7	99.8	90.4	97.1
5	52	64020	21104	NIST alpha full-class	93.4	98.7	78.2	94.7
6	52	64020	21104	NIST alpha folding	97.5	99.5	92.4	97.1

Table 2. Experiments are conducted with different size of training and testing data sets

Classification Method	ACT	NN	6-NN
Top 1 correct rate (%)	95.7	96.4	95.7
Template size (Kbytes)	612	976	3,777
Recognition Time (msec/char)	1.45	11.5	384
Training Time (hour)	1	24	1

Table 3. Performance in various classification methods using data set 3. Proposed method operated on dynamic mode with $\log(\sigma) = 40$

on similar, if not same data sets.

Equations 6 and 7 do suggest that if an incorrect class is picked as the top choice in an early recursive cycle the error is going to be propagated in subsequent recursive cycles. Figure 7 illustrates that a if miss-classification is made in an early recursive cycle confidence value of the top choice (not the true class) rapidly drops and the separability becomes unstable and remains near zero even though the number of sub-images used by new recursive increases. On the other hand, when the top choice is the true class, the separability parameter holds steady.

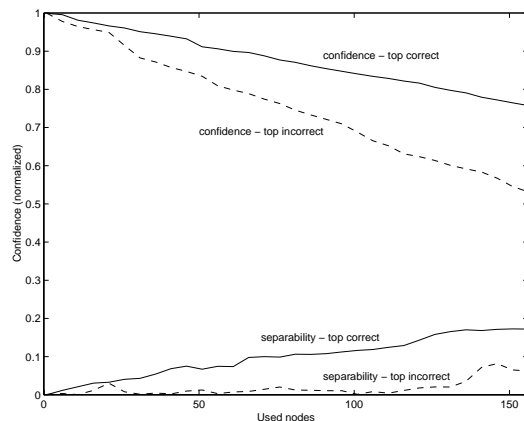


Figure 7. Confidence and separability changes - in data set 1. Solid line shows successful cases and dashed line shows error cases

References

- [1] C. Y. Suen. Computer Recognition of Unconstrained Handwritten Numerals. *Proceedings of the IEEE*, 80:1162–1180, 1992.
- [2] J. T. Favata and G. Srikantan. A Multiple Feature/Resolution Approach to Handprinted Digit and Character Recognition. *International Journal of Imaging Systems and Technology*, 7:304–311, 1996.
- [3] S. K. Chang. *Principles of Pictorial Information Systems Design*. Prentice Hall, 1989.
- [4] G. Srikantan S. W. Lam and S. N. Srihari. Gradient-Based Contour Encoding for Character Recognition. *Pattern Recognition*, 29(7):1147–1160, 1996.
- [5] H. Freeman. Computer Processing of Line-Drawing Images. *Computing Surveys*, 6(1):57–97, 1974.
- [6] M. D. Garris and R. A. Wilkinson. *NIST Special Database 3; Binary Images for Handwritten Segmented Characters (HWSC)*. John Wiley & Sons, Inc, New York, 1992.
- [7] J. J. Hull. A Database for Handwritten Text Recognition Research. *IEEE Transactions on Pattern Analysis and Machine Interlligence*, 16(5):550–554, 1994.
- [8] T. M. Ha and H. Bunke. Off-Line, Handwritten Numeral Recognition by Perturbation Method. *IEEE Transactions on Pattern Analysis and Machine Interlligence*, 19(5):535–539, 1997.