

OCR IN A HIERARCHICAL FEATURE SPACE

Jaehwa Park, Venu Govindaraju, and Sargur Srihari

Center of Excellence for Document Analysis and Recognition
Department of Computer Science
SUNY, Buffalo
Amherst, NY 14260

ABSTRACT

This paper describes a methodology that allows fast and accurate character recognition while keeping the dimensionality of the feature space relatively small. Higher dimensionality can add to the discriminatory power of a recognizer but pays the price in an increase of computational time. We present a method that achieves high accuracy even with a low-dimensional feature space by simulating a multi-resolution feature space. Our approach is supported by promising experimental results. Recognition rate of 98% is achieved on a test set of about 16,000 handwritten numerals. Recognition rates on upper and lower case handprinted characters is about 95%.

1. INTRODUCTION

Recognizers using multi-resolution features have been found to be more accurate than conventional methods that work with features at a single scale [1]. The Gradient Structural Concavity (GSC) features described in [7] are 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. Features at the 3 levels, **G**, **S**, **C** are combined in a k-nearest neighbor classification method to produce a multi-resolution character recognizer. **G** features are the finest and the **C** features are the coarsest. The concept of a multiresolution feature space forms the

basis for the idea presented in this paper.

The computational complexity of GSC recognizer is increased because of the burden of additional features needed to accommodate the multi-resolution aspect of the feature space. Furthermore, the heterogeneous composition of the feature set makes it difficult to find an appropriate distance metric that can be used by the classifier in the feature space.

In this paper we present a hierarchical recognizer (Figure 1) that preserves the accuracy of multi-resolution methods while circumventing the disadvantages cited above. The recognizer begins with a coarse feature space and focuses on smaller subimages of the character on each pass, thus effectively working with a finer resolution of a subimage each time. The subimages are identified by an approach called *gaze planning*- it guides the successive passes through a series of subimages where the features of interest are present.

2. HIERARCHICAL CLASSIFICATION

Statistical pattern recognition techniques typically extract and evaluate features of a pattern image. Each feature is represented by one or more dimensions of the feature vector. A probability value along each feature dimension is computed based on the strength of the features extracted. If all the features are orthogonal, class probability of an input pattern is readily computed .

By dividing feature dimensions into non-over-

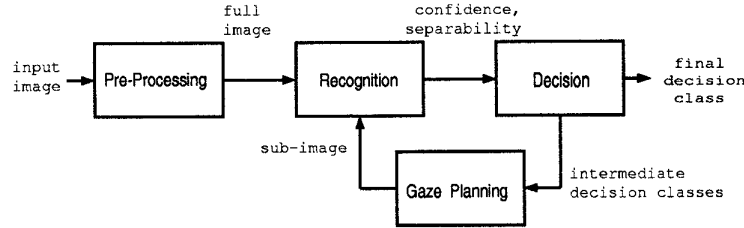


Figure 1: Hierarchical Classifier

lapping subsections, a quantized hyper feature space is established where the size of the hyper space is the product of all the feature sizes.

Let \mathbf{x} be a feature vector with dimension M and \mathbf{c} be a set of mutually exclusive classes of size N ,

$$\mathbf{x} = [x_1, x_2, \dots, x_M], \quad (1)$$

$$\mathbf{c} = [c_1, c_2, \dots, c_N]. \quad (2)$$

Posteriori probability of class c_i with feature vector \mathbf{x} of an input image is given by Bayes' Rule,

$$P(c_i/\mathbf{x}) = \frac{P(\mathbf{x}/c_i)P(c_i)}{P(\mathbf{x})} \quad (3)$$

where the probability of feature vector \mathbf{x} is

$$P(\mathbf{x}) = \sum_j P(\mathbf{x}/c_j)P(c_j). \quad (4)$$

2.1. Probability in Vector Quantized Hyper Feature Space

Given a threshold vector $\mathbf{T} = [T_1, T_2, \dots, T_M]$, we can obtain a binary feature vector,

$$\mathbf{b} = [b_1, b_2, \dots, b_M] \quad (5)$$

where each component is

$$b_{i|i=1}^M = \begin{cases} 0 & \text{if } x_i \leq T_i, \\ 1 & \text{otherwise.} \end{cases} \quad (6)$$

Then \mathbf{b} is a $M \times 1$ matrix and there 2^M bs such matrices possible.

By defining a vector code space,

$$v = \{0, 1, 2, \dots, 2^M - 1\} \quad (7)$$

every binary feature vector can be mapped onto the vector code space. In other words, every M dimensional hyper feature space can be then mapped to a unique code. The posteriori probability of class c_i with a vector code v is then given by

$$P_i = P(c_i/v) = \frac{P(v/c_i)P(c_i)}{P(v)} \quad (8)$$

Prior probabilities, $P(v)$ are obtained during training on preclassified image sets and stored as tables for future look-ups. The additional storage is computed as follows:

Given 9 features and 4 bytes per value, the table size for holding $P(v/c_i)$ is

$$2^{M=9} * 4\text{byte} = 2048\text{byte} \quad (9)$$

Accurate recognition of characters typically warrants many more features. We obtained 54% accuracy of top choice on a set of numerals which is not acceptable for practical applications. As mentioned earlier (9) feature space dimensionality must increase for accurate classification of noisy data. Common remedy therefore, as suggested by pattern recognition literature, would be to add new features. The GSC method introduced earlier uses a 512 binary feature vector ([7]). [4] uses 30 real coefficients to obtain satisfactory recognition performance. However, if we use 512 features, the size

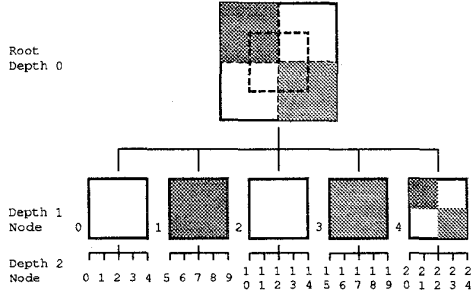


Figure 2: Quin Tree Structure

of the table will become prohibitively large ($2^5 12 \times 4$ bytes!).

The alternative is to keep the number of features the same but focus on subimages of the original image. We present a new architecture to support such a methodology. A process called gaze planning provides the subimage after each pass as long as we are not satisfied with results obtained. We will discuss this further in section 2.2.

Feature extracted in earlier passes provide the global level features while those in latter passes provide local features. The multi-resolutionality of features is achieved, although the actual feature extraction process and the nature of the features remains the same.

The subimages are obtained at any stage by using a division rule such as Quad Tree or Quin Tree (figure 2).

Let $P_i^{(d,l)}$ be a posteriori probability of class c_i and vector code v , within the sub-image of l^{th} node on tree depth d ,

$$P_i^{(d,l)} = P(c_i/v(d,l)) = \frac{P(v(d,l)/c_i)P(c_i)}{P(v(d,l))} \quad (10)$$

where

$$P(v(d,l)) = \sum_j P(v(d,l)/c_j)P(c_j) \quad (11)$$

$v(d,l)$ is a vector code that corresponds to feature vector from a sub-image bounded by node (d,l)

where,

$$d = 0 \sim D - 1, \quad (12)$$

$$l = 0 \sim L^d - 1, \quad (13)$$

D is the maximum depth of the tree and L is the number of branches per node.

If we assume that the class probabilities of a sub-image are independent of those of other sub-images, we can denote the probability P_i^d of class c_i of an entire layer at depth d by the products of probabilities of the same class within all sub-images.

$$P_i^d = P(c_i/v(d)) = \prod_l P(c_i/v(d,l)) \quad (14)$$

where $v(d)$ is a vector code matrix of $L^d \times 1$,

$$v(d) = [v(d,0), v(d,1), \dots, v(d, L^d - 1)]. \quad (15)$$

The overall probability of class c_i throughout the tree can be obtained from the probability values at each layer.

$$P_i = \prod_d P_i^d = \prod_d \prod_l P_i^{(d,l)} \quad (16)$$

$\sum (L^d - 1)$ multiplications are necessary to obtain the probability of a class, using equation (16). Applying the logarithm function to (16), we can convert to the equivalent add operations.

2.2. Gaze Planning

Since we assumed that the posteriori probabilities of a class within subimages are independent, we can update probabilities of a class c_i recursively. A subset of probabilities of a class is obtained by the products of those of sub-images chosen vertically as well as horizontally within a tree structure, starting from the root extending to the sub nodes. If we fix the number of classes whose probabilities get updated on each pass and only examine those which were updated in the previous pass (no new ones are initiated) the intermediate probabilities can be computed recursively at pass t ,

$$P_i(t) = P_i(t-1) \cdot P_i^{sub}(t) \quad (17)$$

where $P_i^{sub}(t)$ is products of posteriori probabilities of class c_i obtained from nodes which are the child nodes of a parent node that is selected at time t .

Let function $g(t)$ determine whether a node is to be examined in the next iteration. We can write $P_i^{sub}(t)$

$$P_i^{sub}(t) = P_i^{g(t-1)} = \prod_{k=sub(d,l)} P_i^{(d+1,k)}|_{(d,l)=g(t-1)} \quad (18)$$

where $sub(d,l)$ is a function which gives child node numbers on $d+1$ layer which belong to node (d,l) . If we choose $g(t)$ as a gazing function that finds the most informative node within the pool of examined nodes for updating class probabilities, probabilities of classes will rapidly approach their steady state.

The classification result of the recognizer is accepted based on confidence of the top choice and its separability (distance) from the second choice ([2] [7]).

We can define confidence and separability values as follows. Confidence $C_i(t)$ of a class c_i at recursion time t is defined by the ratio of the current (intermediate) probability and the maximum possible probability of the class,

$$C_i(t) = \frac{P_i(t)}{Q_i(t)} \quad (19)$$

$$Q_i(t) = \prod_{k=0}^{t-1} Q_i^{g(k)} \quad (20)$$

where $Q_i(t)$ is the products of possible maximum probabilities of class c_i with any vector code in a node sequence. The node sequence is accumulated outputs of a gazing function $g(k)$ till pass $t-1$.

In recursive representation,

$$Q_i(t) = Q_i(t-1) \cdot Q_i^{g(t-1)}, \quad (21)$$

where

$$Q_i^{g(t-1)} = \prod_{k=sub(d,l)} \max(P_i^{(d-1,k)})|_{(d,l)=g(t-1)} \quad (22)$$

Recursive equation of confidence for class c_i is given as

$$C_i(t) = C_i(t-1) \cdot C_i^{sub}(t) \quad (23)$$

where $C_i^{sub}(t)$ is the confidence of newly added branches of node (d,l) ,

$$C_i^{sub}(t) = C_i^{g(t-1)} = \frac{P_i^{g(t-1)}}{Q_i^{g(t-1)}} \quad (24)$$

From (18) and (23), the intermediate probabilities and confidences of each class are obtained at every pass. The classes are ordered using intermediate probabilities. Separability, $S_i(t)$, of the top two choices is given as

$$S_{i=top}(t) = \frac{P_{top}(t)}{P_{second}(t)} \quad (25)$$

when class c_i has the highest probability at time t .

If the result does not meet acceptable confidence and separability values, features within the same subspace is computed at different scale in the next pass.

For driving the gazing function, we define local separability at time t as a metric of measuring the support present for the top class choice. If class c_i is the top choice and class c_j is the second choice at time t , the local separability $S_i^{(d,l)}(t)$ of node (d,l) is

$$S_i^{(d,l)}(t) = \frac{P_i^{(d,l)}(t)}{P_j^{(d,l)}(t)} \quad (26)$$

We can now determine the next gazing node which is expected to be the most supportive of the current recognition results. If the top and second rated classes change, the local separability should be renewed at the nodes in which the updating has ended. Figure 3, summarizes the recognition algorithm.

To reduce computational complexity, if we use logarithm function, the equations of products for parameters are transformed to summations.

3. IMPLEMENTATION

1. Contour Representation:

A character image is scanned and binarized. Binarized character images are represented by

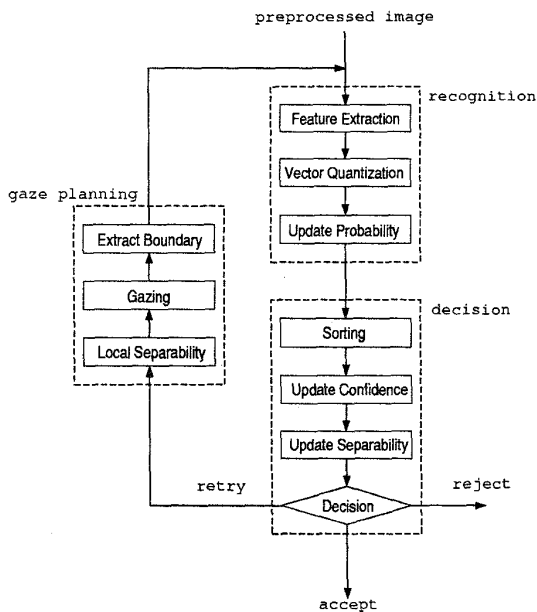


Figure 3: Recognition Algorithm

piecewise linear contours. High curvature points on contours are first identified. Contour segments between these points are assumed to be linear. This representation adds to overall speed by reducing the total data points and also by eliminating a the need for smoothing and noise removal. Figure 4 (a) shows the original binarized images and (b) shows slant and size normalized piecewise linear contours.

2. Grids

A variable size rectangular grid is used to define sub-regions for hierarchical feature extraction. Four exclusive rectangular sub-regions from the center of gravity of contours define the quad tree structure. The quin tree structure is same as the quad tree with an additional fifth sub-region formed by joining the centers of the other four sub-regions (Figure 2).

Quad and quin trees with a maximum depth

Source	Class	Training	Testing	Top%
Postal+NIST	10	25656	12242	96.3
NIST digit	10	159228	53301	98.2
NIST upper-case	26	34961	11481	96.4
NIST lower-case	26	29059	9623	90.4
NIST full	52	64020	21104	78.2

Table 1: Recognition Performance (Quin Tree)

4 were implemented. Figure 4-(c)s, (d)s and (e)s show sub-regions for each layer in the quin tree. Figure 4-(f)s show reconstructed contours using feature at the deepest level (4).

3. Features:

8 bit binary feature vector for the quad tree representation and 9 bit feature vector for the quin tree were used. The size of the feature vector is kept small for reasons described earlier. The direction of a contour ranges from 0 to π radians if we consider them to be unidirectional. 4 gradient feature bits are obtained from the number of pixels in each $\pi/4$ direction. Moment feature vectors are used to obtain the other 4 (or 5 in case of quin trees) bits.

4. EXPERIMENTS

The training and testing sets consist of isolated digits and characters collected from handwritten postal addresses and NIST digits sets.

1. The recognition results are shown in Table 1.
2. Table 2 shows the quin tree having superior recognition performance, but it increases the processing time.

5. CONCLUSION

Hierarchical classification allows high recognition rates with a small set of features. Finer resolution of features is simulated by sampling smaller subimages in each pass. We have assumed that there is no correlation between subimages. This assumption is not

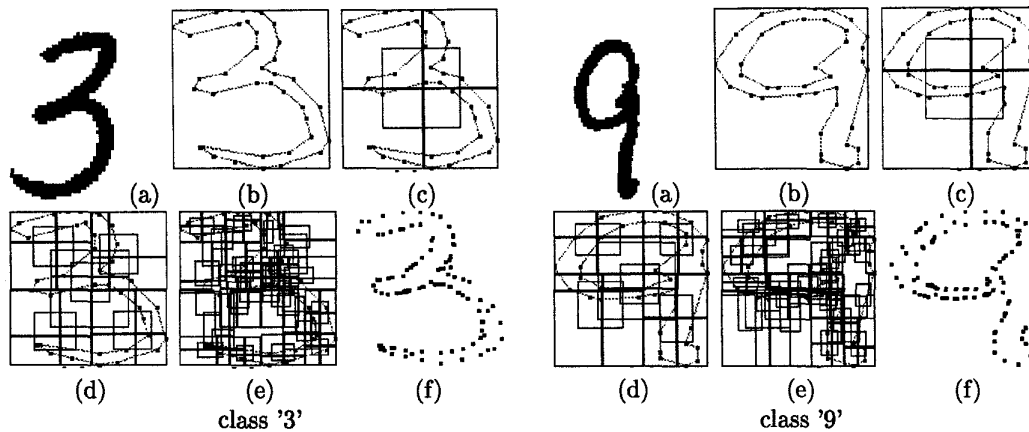


Figure 4: Character Image and Sub-Regions of Quin Tree: (a) original image, (b) piecewise linearized contour, sub-regions at (c) layer 1, (d) layer 2, (e) layer 3, and (f) pseudo reconstructed contour image from center points of final layer

Type	Quad Tree	Quin Tree
Total Number of Nodes	85	156
Top Correct Ratio (%)	94.0	96.3
2nd Correct Ratio (%)	97.9	98.5
Avg. Recog. Time (ms)	6.04	12.36

Table 2: Performance Comparing of Quad and Quin Tree Structure

Number of Classes	10	26	52
Recognition Time (ms)	12.36	13.05	13.97
Characters per sec	202	192	179

Table 3: Recognition Time

entirely true. Future work will address this issue to improve the recognition accuracy and speed further.

6. REFERENCES

- [1] A.K. Jain O.D.Trier and R.Taxt. Feature extraction methods for character recognition - a survey. *Pattern Recognition*, 29(4):641-662, 1996.
- [2] S.W.Lam G.Srikantan and S.N. Srihari. Gradient-based contour encoding for character recognition. *Pattern Recognition*, 29(7):1147-1160, 1996.
- [3] F.P.Kuhl. Elliptic fourier features of a closed contour. *Computer Graphics and Image Processing*, 18:236-258, 1982.
- [4] R.R. Bailey and M. Srinath. Orthogonal moment features for use with parametric and non-parametric classifiers. *IEEE PAMI*, 18(4), April 1996.
- [5] C. Y. Suen. Computer recognition of unconstrained handwritten numerals. *IEEE Proc.*, 80:1162-1180, 1992.
- [6] A. Crzyzak L. Xu and C. Y. Suen. Methods of combining multiple classifiers and their applications to handwriting recognition. *IEEE Trans. on System, Man and Cybernetics*, 22(3):418-435, 1992.
- [7] J.T.Favata and G. Srikantan. A multiple feature/resolution approach to handprinted digit and character recognition. *Intr. J. of Imaging Systems and Technology*, 7:304-311, 1996.