# DOUZIEME COLLOQUE GRETSI - JUAN-LES-PINS 12 AU 16 JUIN 1989



# NUMBER OF FEATURES FOR HANDWRITTEN CHARACTER RECOGNITION

R.Krishnan Tampi<sup>1</sup> C.S.Sridhar<sup>2</sup>

<sup>1</sup> F.A.C.T.Ltd, Udyogamandal, Cochin, India <sup>2</sup> Cochin University of Science and Technology, Cochin - 682022, India

#### RESUME

### **SUMMARY**

Definition of features for general class of images is a difficult problem. For a particular class -- handwritten lower case English script -- it is shown that features can be defined on an arithmetic basis. This is achieved by dividing individual characters into fixed number of equal areas and using a finite set of primitives to define a best approximation to the contour existing in that area.

#### Introduction

Handwriting recognition requires the correct description of images and appropriate definition of features. Suen recommends a 30x40 pixel size for the recognition of binary images and alphanumeric characters [1]. However, feature ectraction remains a problem with no simple solution in sight. A recogniser which uses structural and global properties would be the most appropriate approach to the extraction of topological features. any case computation required will be enormous. A question arises whether for a small set of characters which have to conform to certain universal shapes such a complex approach is essential.

It is common knowledge that highly distorted, contour approximated letters in handwriting are easily recognised by Acknowledging the fact human beings. that humans have a large amount of learning the result remains unaltered that, in the recognition domain these characters are seperated by distances which are large and are invariant to the type of algorithm In this paper it is used by recognisers. shown that features can be defined on arithmetic basis for use in recognition algorithms. 36 features and two auxiliary are defined to assist recognition.



Feature definition on arithmetic basis

Consider the image of a single character. Let 144 pixels of its binary image arranged in 12 rows and 12 columns describe it completly. Let this be divided into 16 subregions called cells. Let each cell contain 9 pixels in 3 rows and 3 columns. This cell is used to define features.

The contour which passes through each pixel is binarised, its presence denoted by a 1 and absence by a 0. Thus the image is entirely represented by these 9 pixel square cells. This can be transformed to a 9 bit binary word with the central pixel bit as the MSB as shown in Fig.1. After analysing a large number of samples a few unique contours were selected as features for this class of images.

Since each cell has 9 elements there could exist 512 features provided they are unique. In natural writing all these do not exist and broken bits, tails and such artifacts cannot be treated as features. Looking at the binary words which represent the contours, we can find that not all of these are unique, being only rotations of other words. It can then be implied that such words which are mere rotations of others cannot be qualified as features.

Fig.2 shows two examples of describing the image of a character in the 16 cell paragraph format. 144 pixels have been divided into cells of 9 elements and also shown is the method of deriving the 9 bit binary word from the cells. Consider now a single cell. If it is rotated without altering its position in the paragraph, the central bit remains unshifted in position while the other 8 undergo shift. But the bit pattern remains unchanged. This is arithmetic operation clearly an on the

8 LSB's of the word. Each possible word which represents a contour can therefore have 7 rotations with reference to itself. Out of these eight only one can be unique. These unique words are to be considered as features. Conversely a feature can be defined as one which is not derived from another. Thus if a binary word representing a contour cannot be derived by lower order rotation from another such a word can be treated as a feature.

The number of features for this type of a description of an image can now be calculated by calculating the number of unique allowable 8 bit words. Such words are called Basic Words. A basic word must be odd decimal, and they can be listed out with mod 64 logic. It is also possible to arrive at a basic word from Consider the word say, any given word. 0 10 110 100. (For simplicity octal notation can be used). The word and its rotations are, 264, 151, 322, 245, 113, 226, 055 and 132, the MSB not considered since it does not undergo rotation. Amongst all these, 055 is alone unique and is the basic word. 055 is the smallest number and all the rest can be derived from it. A simple calculation shows that there are only 36 such basic numbers ( 8 bit words suitable for character cells) with the 9th bit 0 and an equal number with the 9th bit 1. A list of the basic words is given in Table-I. the remaining words from the 512 possible can be generated from these words.

For facilitating computer programme development, these basic words are given a 3 digit name. The first digit indicates the number of 1's in the word, the second indicates the family of the contour, but normally chosen so that the smallest basic



corresponds to the smallest 2nd digit of the name and the third digit representing the rotation with respect to the basic word. Some of the basic words have only four unique rotations while some others have none. Thus a total 512 are accounted for. For the type of image, cursive script handled here, three basic words are illegal as they represent broken contours. Considering all these factors 68 features can be defined.

Sufficiency of the features defined -

The 26 characters, the targets for classification in this work, show 0.6 as the cell occupancy and 0.25 as the pixel features occupancy rate. The above are sufficient for the classification of these images since the rotational forms can take care of any contour that may appear in handwriting. However, is a very large collection of words which will be difficult to implement on small machines. Hence data reduction is adopted. Data reduction, called filtering here, takes the form of replacing all words with more than 3 1's by words with 3 1's. words are retained, which when used for reconstruction of the character show remarkably good images. A further reduction was attempted, wherein only 4 words are These are symmetric horizontal, vertical, right and left slants. Fig.3 illustrates the results of filtering for some character images. The retention of 'letter' shape after the first filtering (8 words) is good. Also remarkable is the fact that when only four words are used the reconstructed images are highly distorted no doubt, but, the image of one character does in no way resemble that of another. This gives a good scope for classification.

'However, with only four features available syntactic methods for classification lead to a poor recognition score. be offset by a stricter replacement method during filtering. A preservation measure is defined to ensure stricter replacement. The trajectory and the continuity of the contour have to be used as measures for good shape retention. Therefore two secondary features are defined. These are called Tense and Tendency. They do not have the status of features. The extraction of these auxiliary features is fairly straight forward. The last digit of the word stands for the tendecncy, while the last two digits indicate the tense. Fig.3 illustrates these It is to enable the extraction aspects. of features like these that the word names for the basic words were appropriately With the paragraph of 16 cells chosen. and 4 features supported by auxiliary featu-

#### Conclusion -

macines.

It is shown that for lower case English (Roman), handwritten characters a small number of features can be defined on an arithmetic basis. Using auxiliary features recognition can be implemented based on four features. Though only handwritten characters were considered, this approach may hold good for similar character sets and images which have a universal quality.

res recognition is possible even using small

#### Reference -

[1] C.Y.Suen, "Character recognition by computer", Handbook of Pattern Recognition, Eds. K.S.Fu and T.Y.Young, Academic Press, N.Y., 1986.

'b'

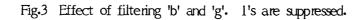
Filter I



## Number of Fc tures for Handwritten Chancter Recognition

Table I Basic words

Octal Word	Word Name	Word Name MSB	=0 Word Name MSB=1	Comment
000 001 003 005 007 011 013	000 001 200 210 300 220 310	000 001 200 210 300 220 310	000 001 307 380 490 390 4A0	Null Illegal Legal
	035, 037, 045 065, 067, 073	and similarl , 023, 025, 027, 03 5, 047, 053, 055, 05 8, 075, 077, 125, 12 7, 177 and 377	57, 063	
00000000000000000000000000110000000000	000 001 001 0 110 100 0 100 Example and the	3 3 4 4 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6	06 000 102 206 426 06 000 203 306 201 E4 304 382 3A0 000 paragraph descrpt 100 110 100 after one 100 2nd 100 rotation 000 ation Cell 306 and its rotati Names are 307 and 3	n 000 ion. 00.
			00000 000 000 00000 000 00 00000 00 0 00000 00 00000 00000 00000 00000 00000 00000 00000 00000 00000	0000 00000 0 000 00 0000 0 000 0000 0



Filter II

'g'

Filter I

Filter II