

# IMPROVING CLASSIFIER ACCURACY BY SIMULATING FUZZY BOUNDARIES BETWEEN CLASSES

Venu Govindaraju, Ianakiev Krassimir and Sargur Srihari

Center of Excellence for Document Analysis and Recognition  
Department of Computer Science  
SUNY, Buffalo  
Amherst, NY 14260  
*email*:{govind}cedar.buffalo.edu  
*Fax*:716-645-6176

## 1. INTRODUCTION

The pattern classification problem can be defined as one of assigning a label to a pattern of unknown class based on labeled prototype patterns. Formally, given the pattern  $x$ , a set of  $c$  class labels ( $\omega_1 \dots \omega_c$ ) and a set of labeled prototypes, the classifier should assign a label to  $x$  such that the probability of misclassification is minimum.

A typical classifier  $C$  returns scores corresponding to each class which represents its "confidence" or "belief" in each class representing the unknown input pattern  $x$ . The scores may be represented by  $S^T$  which is a vector of  $c$  values where value  $s_i$  corresponds to class  $\omega_i$ . The scores reflect the degree to which the unknown pattern  $x$  represents the prototype(s) of the class. The performance of the classifier is measured by its correct classification probability, estimated by the relative frequency with which it assigns unknown patterns to their correct classes. Higher the accuracy, more desirable the classifier is in practical applications.

The method described in this paper is based on the following two ideas which appeal to our common sense [1].

1. When the correctness of a classifier on a pattern  $x$  is in question, it is best to consider the performance of the same classifier on the pat-

terns which are "similar" to  $x$ .

2. A classifier is usually accurate when the test pattern  $x$  falls close to the center of its class in feature space and prone to error when it falls near a class boundary.

## 2. BACKGROUND

Let us consider the problem of digit recognition where each of the digits  $0 \dots 9$  are the classes  $\omega_0 \dots \omega_9$ . Figure 1 illustrates the class boundaries and shows the prototypes in some feature space. We will assume for the purpose of this discussion that the classes are linearly separable. Samples that are "close" to the centroid of the classes are "stronger" representatives of the class while samples close to the boundaries are poorer representatives and at the same time prone to confusion with the neighboring classes.

The traditional classification techniques approach the problem as follows.

1. Given an input pattern  $x$ , a classifier typically computes the distance of the prototype from the centers of each of the classes in the feature space using some distance metric. The classes are then ranked in increasing order of the distances. In Figure 1,  $x$  is "closest" to

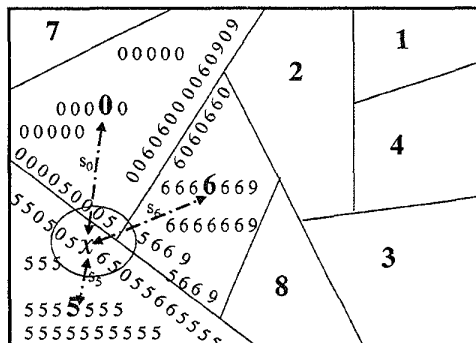


Figure 1: Feature space  $C_f$  of a typical digit recognizer. Samples that are close to the center of a class are strong representatives of the class and are usually correctly classified. Samples that are away from the center (close to the boundaries of classes) are weaker representatives and are prone to misclassification. The circle is essentially the sphere of the neighborhood of  $x$ .

the center of the class '5' with a distance of  $s_5$ . The 2<sup>nd</sup> and 3<sup>rd</sup> classes are '6' and '0' at distances  $s_6$  and  $s_0$  respectively.

It is possible that the test pattern  $x$  falls within the boundary of one class but its distance is closest to the center of a neighboring class. Since  $x$ , in the example, falls close to the boundary of classes  $\omega_5$ - $\omega_0$  and  $\omega_5$ - $\omega_6$  there is more than an average chance of misclassification of  $x$ .

2. A  $K$  Nearest Neighbor (KNN) classifier would label  $x$  as the majority class in the set of its  $k$  nearest neighbors in feature space [2]. This method has the advantage of not respecting fixed boundaries between classes. In the example of Figure 1, the KNN classifier would label  $x$  as '5' with a score of  $\frac{3}{6} = 0.5$ , as four of the samples in the  $K$ -neighborhood ( $k=6$ ) are labeled  $\omega_5$  by the recognizer (some erro-

neously).

### 3. FUZZY BOUNDARIES

Among the two traditional methods described above, the boundaries between classes are firm in the first method. The KNN method does not abide by firm boundaries. Each test pattern in the feature space is assigned a class based on the class of its neighbors.

Since potentially, a test pattern can fall anywhere in a feature space, we can precompute the class that will be assigned to each point in the feature space. If the KNN rule is used, the boundaries between classes will have a very complex shape. Further, the boundaries of classes will change based on the size of the neighborhood  $k$ .

Computation of such boundaries can be done ahead of time, so that when a test pattern falls at a particular point in the feature space a table lookup can determine its class. The table can prestore the class assignments for all points in feature space for various values of  $k$ . The distance of the test pattern from several neighbors has to be first computed and sorted. The top  $k$  form the  $k$  neighbors.

A typical binary vector of size 512 implies a feature space of size  $2^{512}$ . Hence the size of a table as described above will be prohibitive. Based on our observation stated earlier, that a classifier is prone to errors predominantly near the class boundaries, we can adopt the KNN rule only when a test pattern falls near a class boundary. For cases, where the test pattern is near the central core of a class it can be simply assigned to the same class. The cases when a test pattern is close to a boundary can be determined by simply observing the "confidence" score received by the pattern. Scores below a threshold value (each class can have a different threshold) indicate boundary cases. Such cases can be referred to the prestored table.

### 4. EXPERIMENTS

We have already described how the fuzzy boundaries change results only for those samples that fall in the proximity of class boundaries. Hence if the

the original recognizer has a high recognition accuracy, our method can only marginally affect the results. GSC digit recognizer is such a recognizer with very high accuracy. We do improve the overall accuracy, albeit, by just 0.05% (Table 1).

The performance improves in most classes but falls in some. Our future work will be to analyze the cause of drop in accuracy in some of the classes (Table 2).

GSC		Fuzzy	
reject %	error %	reject %	error %
0.0000	1.1763	0.0000	1.1763
0.3994	0.9956	0.3921	0.9923
0.7924	0.8234	0.8169	0.7577
1.3070	0.6704	1.3151	0.6374
1.8134	0.5491	1.8134	0.5075
2.6303	0.4195	2.6058	0.3019
2.8100	0.3866		
2.9897	0.3705	3.0142	0.2443
3.5125	0.3386		

Table 3: Improvement of the fuzzy method at selected accept rates

10E+R, a metric based on the notion that ever image classified in error is as expensive as rejecting (without a decision) 10 images. A plot of the Error versus Rejects is shown in Figure 2. Clearly, the fuzzy approach has a better score on this metric. Table 3 shows some selected points on the curve.

## 5. REFERENCES

- [1] D. Bouchaffra, V. Govindaraju, and S. Srihari, "A Methodology for determining probability of correctness of recognizer scores", *Accepted, IEEE Conference on Computer Vision and Pattern Recognition*, Santa Barbara, June 1998.
- [2] R. O. Duda and P. E. Hart, "Pattern Classification and Scene Analysis", Wiley, New York, 1973.

test set	Original		Improvement	
	correctly recognized	% of correctly recognized	correctly recognized	% of correctly recognized
30398	29532	97.15	29788	97.99

Table 1: Improvements to GSC recognizer using Fuzzy boundaries

class	size of the test set	GSC		Fuzzy	
		correctly recognized	% of correctly recognized	correctly recognized	% of correctly recognized
0	6986	6833	97.81	6833	97.81
1	6554	6496	99.12	6496	99.12
2	3954	3759	95.07	3845	97.24
3	2238	2170	96.96	2186	97.68
4	3191	3070	96.21	3023	97.87
5	1556	1497	96.21	1513	97.24
6	1474	1442	97.83	1449	98.30
7	1445	1391	96.26	1412	97.72
8	2049	1945	94.93	2001	97.61
9	951	929	97.69	930	97.93

Table 2: Improvements to GSC recognizer using Fuzzy approach

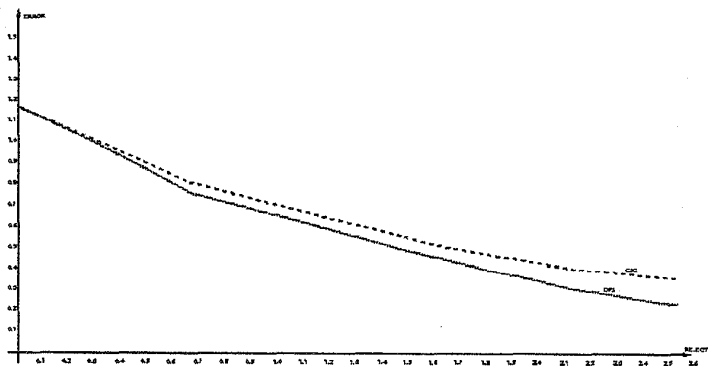


Figure 2: Error versus reject curve for GSC and Fuzzy method