

A THREE-WAY CLASSIFICATION STRATEGY FOR REDUCING CLASS-ABUNDANCE: THE ZIP CODE RECOGNITION EXAMPLE

CHUN-HOUH CHEN AND KER-CHAU LI

ABSTRACT. For many real world problems in discriminant analysis, the abundance of classes is one major source of complexity. Let k be the number of classes. When k is large, it is difficult to discriminate all k classes at a time. To ignite the statistical community's interest in addressing this issue, we focus on the handwritten digit recognition problem using the Zip code data analyzed in the seminal article by LeCun *et al* (1989) wherein the neuron network with backpropagation was first applied successfully to a real world problem. We propose a three-way sub-classification strategy that first divides the original k -class problem into $C(k, 2) = k(k - 1)/2$ smaller problems. Each smaller problem concerns the discrimination between two classes from the initial k classes under the presence of a third class that consists of all cases from the rest of the $k - 2$ classes. The decisions from each sub-problem of three-way classification are then combined via a conditional error rate analysis that exploits the degree of consensus among various decisions. We report the existence of a large portion of high quality images found by our method; they are predicted with an error rate less than 1%.

1. INTRODUCTION

Classification problems concern the class membership prediction of an item from a population which is partitioned into k classes. One well-known example is the iris data studied by R. A. Fisher (1936). Three kinds of iris are under consideration and the problem is how to distinguish them based on data from four physical measurements. In search of a solution for this problem, Fisher introduced the celebrated linear discriminant analysis (LDA). Another problem which has recently received a great deal of attention is the recognition of handwritten digits such as zip codes from envelopes. The number of classes in this problem is $k = 10$ and each class represents one of the

KEY WORDS: Dimension reduction; Feature extraction; Hand-written digit recognition; Linear discriminant analysis.

ten digits, $0, 1, \dots, 9$. Classification problems may come from diverse scientific areas such as speech recognition, medical diagnosis, remote sensing, and genomics.

As expected, the task of classification becomes harder when the number of classes k increases. This is in large part due to the difficulty of finding main feature variables suitable for all k classes. In generic terms, features are just real-valued functions summarizing the useful aspects of the original variables in the training dataset. In complex applications, they are needed for reducing the data dimensionality. Ideally, feature variables must be constructed to capture distinctive properties between classes. When there are only two classes, feature extraction appears relatively simpler. In digit recognition for example, it is not hard to design some feature variables for distinguishing any pair of digits. However, the features designed for one specific pair may not be very helpful when applied to other digits.

Since binary classification is easier, an attractive strategy to solve the k -group problem is to consider each of the $C(k, 2) = k(k-1)/2$ two-class classification problems separately first. In this approach, only the items from the two classes in the learning sample are used to construct a classification rule for each two-class problem. After that, a final decision is made by combining together results from each problem. A simple way to do so is to use the majority rule which assigns the class membership of a new item in a way similar to a tournament with k players, each player representing one class and every player playing once against each of the other players. The player with the best winning record is awarded the class membership.

In this paper, a different approach is taken that also breaks the original k -class problems into $C(k, 2)$ smaller problems. However, instead of binary classification for each smaller problem, the new task is not only to discriminate the two selected classes, but also to tell if a case belongs to either class or not. To facilitate this, we form a three-way classification problem with a pair of classes A and B chosen from the original k classes and a third class whose members consist of all cases from the remaining $k-2$ classes. The third class will be referred to as "Other". Figure 1.1 shows a few alternatives in reducing complexity of class abundance. We have noticed that the majority of works in the literature have been on the all-at-one time approach. The intent of this article is to shift researchers' interests to other possibilities. Obviously,

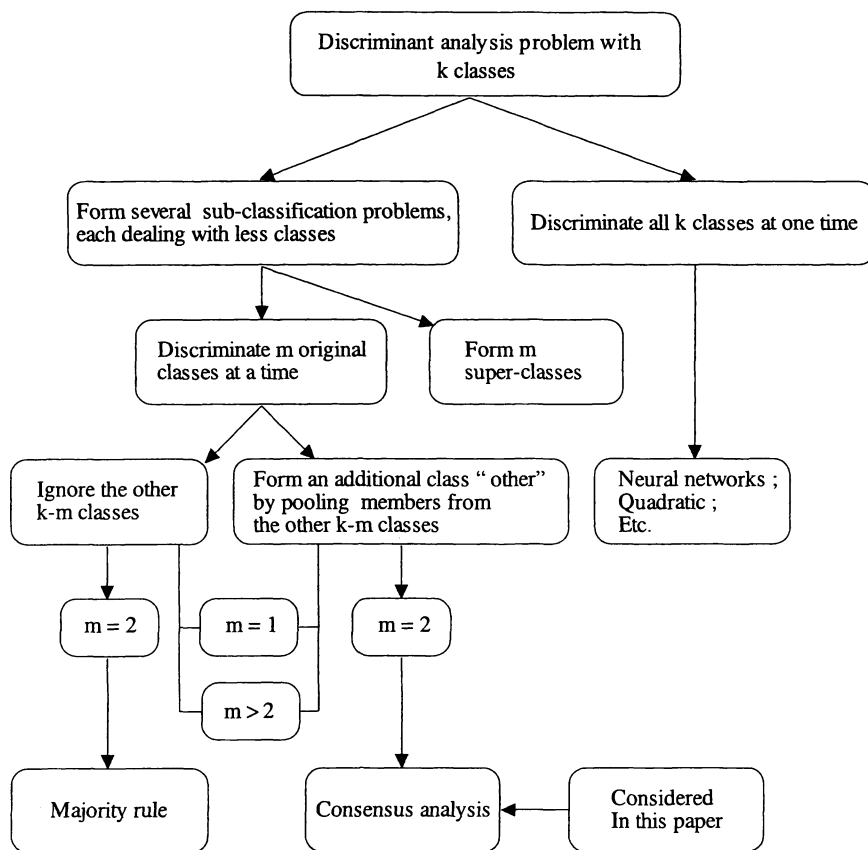


FIGURE 1.1. Alternative Strategies for Reducing Complexity of Class Abundance.

it is unlikely to find any any universal claim of superiority for one choice over another without looking at specific problems.

We shall use hand-written digit recognition as an example to illustrate the advantages of our approach. The data base consists of digitized images of zip codes on envelopes passing through Buffalo, NY. (LeCun *et al*, 1989). Using a conditional error rate analysis, we can identify a large portion of sample images in the test set that can be classified at an error rate of less than 1%.

In section 2, we first give a brief account of the zip code data. A centre-of-mass based method for feature extraction is introduced. This is followed by a preliminary analysis using linear discriminant rules. Section 3 describes details of the three-way

subclassification settings and discusses how the results can be combined using conditional error analysis. In Section 4, we compare our three-way classification with binary subclassification. We introduce another way of partitioning to extract feature variables from the zip code data. With this new feature space, we show how to combine the two-way and the three-way methods to get better results. Further discussion is given in Section 5.

2. DATA, FEATURES, AND LDA

In this section, we first describe the handwritten data set which will be used throughout this article. Then we introduce a method for extracting features that reduce the dimensionality to a level easier to manage. A preliminary analysis involving linear discriminant analysis (LDA) is provided.

2.1. The zip code data. Our data base comes from handwritten zip codes that appeared on some envelopes of U.S. mail passing through the Buffalo, NY post office. The digits were written by many different people with a great variety of writing styles and instruments. Each digit is converted into a 16 by 16 pixel image after some preprocessing as described in LeCun *et al* (1989). Figure 2.1 shows some of these normalized images.

The seminal work of LeCun *et al* (1989) uses a neural network with three hidden layers- 768, 192, and 30 hidden units, respectively, for each layer. A misclassification rate of 0.14% on the training data and 5.0% for the test set were reported. This remarkably low rate of error cannot be achieved without deliberate efforts on setting up the proper connection architecture and weight constraints. For example, the same authors reported that a fully connected network with one hidden layer of 40 nodes yields a worse result, 8.1% error rate on the test data. On page 546, lines 8-10, these authors further indicated the complexity involved: "The eight maps in H1 on which a map in H2 takes its inputs are chosen according to a scheme that will not be described here." Yet another factor affecting the error rate is the total number m of passes through the training set. The 5% rate for the test set was obtained when m is 23. Figure 2 of LeCun *et al* (1989) shows between 5 and 30, the error rate is between 5% and 6% - reaching about 5.5% at $m = 30$.

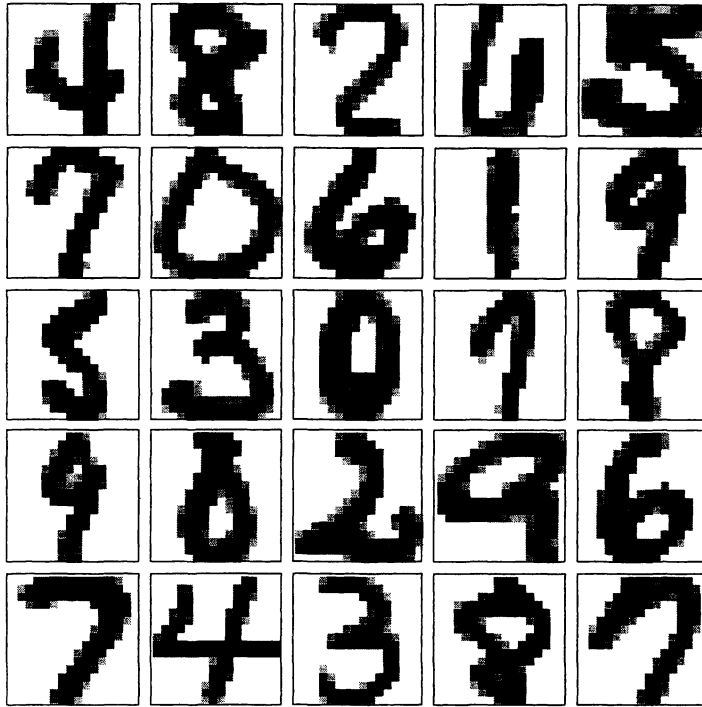


FIGURE 2.1. Sample Hand-Written Digits.

There are several attempts aiming directly at improving the error rates for this zip code data. The best result is 2.5% by Simard, LeCun and Denker (1993) who used a transformation-based nearest-neighbor method. The transformations incorporated various kinds of distortion due to factors such as shifting, scaling, rotating, and so on.

This dataset is also often used for testing classification methods which are not specially tailored for digit recognition. In such applications, the error rates are usually poorer than 5% or 6%. For example, using a penalized discriminant rule, Hastie, Buja, and Tibshirani (1995) obtained a rate of 8.2%. But this is already an improvement over the rate of about 10% obtained by the usual LDA.

2.2. Center-of-mass-based partition. Due to the spatial arrangement, the 256 variables representing the 16 by 16 image for each sample character are highly correlated.

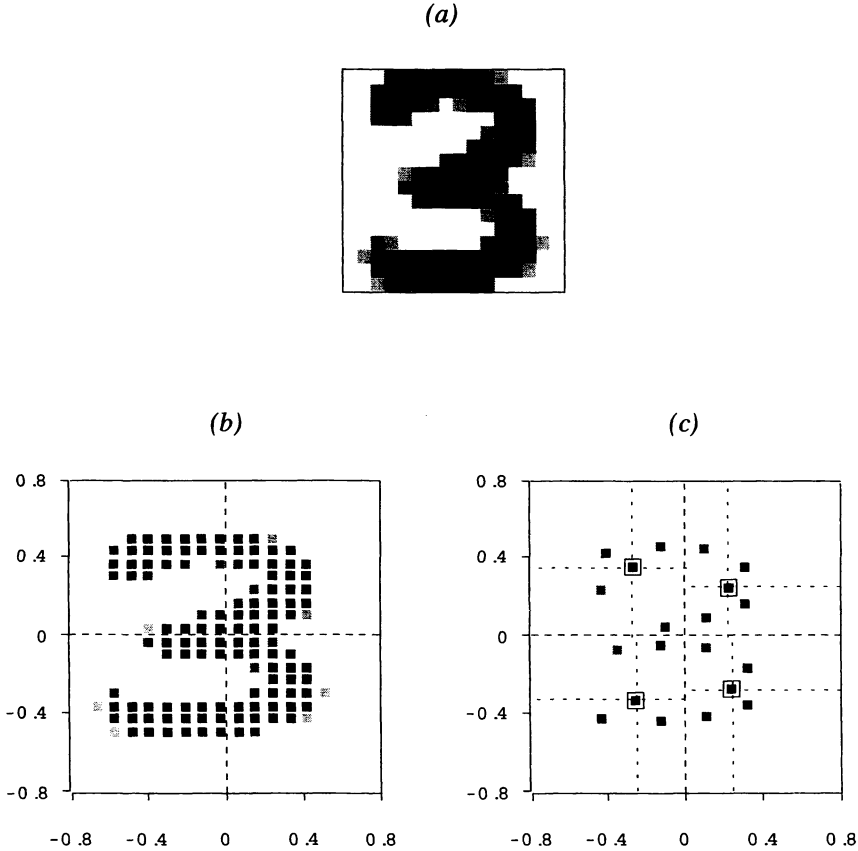


FIGURE 2.2. Locations of Mass Centers for Digit 3: (a). Original Image; (b). Transformed Digit; (c). Locations of Mass Centers with Weights.

Such redundancy among the input variables allows a certain degree of flexibility in designing strategies for feature extraction.

The method we use treats each character image as a piece of object with mass at each pixel (i, j) equals to its grey level w_{ij} . For each object, we first compute the centre of mass $(C_{01}, C_{02}) = (\sum w_{ij}^* i, \sum w_{ij}^* j)$ where $w_{ij}^* = w_{ij} / \sum w_{ij}$ is the weighted mass of pixel (i, j) . This mass center is then used to normalize the image by shifting (C_{01}, C_{02}) to the origin $(0, 0)$. The range of the two coordinate axes are scaled to take values within -1 and 1, using grey levels greater than a prespecified threshold value as a common denominator. Figure 2.2(b) gives a transformed digit 3 whose original

image is in Figure 2.2(a). The entire image is now partitioned into four pieces, one piece from each quadrant.

The next step is to calculate the mass center for the image in each quadrant. This generates 4 pairs of locational variables. After that, each quadrant is again partitioned into four quadrants, using the mass center calculated earlier as the origin. This yields a total of 16 rectangle regions and the mass center of the image in each rectangle is computed. The 20 locations of the mass centers for the digit 3 in Figure 2.2(b) are shown in Figure 2.2(c). Our feature space consists of these 40 locational variables and the 15 out of 16 weight variables for the total mass in each of the final 16 rectangles. We can not use all 16 weight variables because of collinearity. Thus the original 256 grey-level variables is reduced to 55 feature variables.

Quadrants are certainly not the only choices we can use in our center-of-mass-based partition. Later on a system of 8 radial lines as described in Section 4.2. will be used, which leads to better results.

In LeCun *et al* (1989), there are 7291 cases in the training set and 2007 cases in the test set. We use only 7188 and 1991 respectively as our training and test sets - this will be referred to as 7188/1991 data. Other cases are deleted for the moment because they have no mass in some of the final 16 rectangles. This problem can be corrected if we use the radial partition system - details to be discussed later.

It is not clear how the original 2007 test cases were selected in LeCun *et al* (1989). Typically, if the test set comes from the same population as the training set, then it should exhibit characteristics similar to any randomly selected subsample of equal size from the training set. However, this is not the case here. Several authors have reported an unexpected increase in error when applied to LeCun *et al*'s test set. For comparison, these authors often generate two randomly selected subsets of size 2000 each from the 7291 training set, one as the new training set and the other as the new test set. They find that the error rates are smaller for the new test set than the original size-2007 test set. We shall follow this practice and generate a 2000/2000 data - 2000 cases for the training set and 2000 cases for the test, both being randomly selected from the 7188 cases.

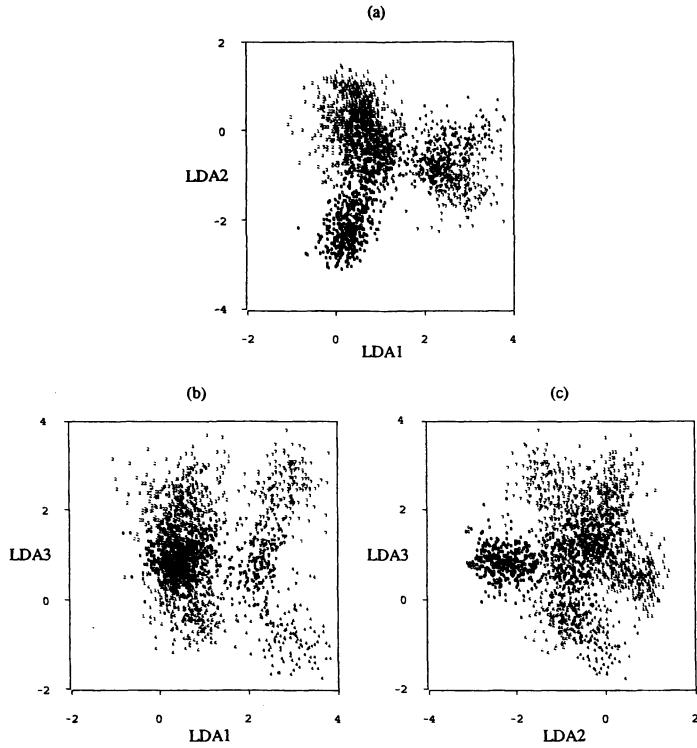


FIGURE 2.3. LDA Canonical Variates: (a). LDA1 versus LDA2; (b). LDA1 versus LDA3; (c). LDA2 versus LDA3.

2.3. A preliminary analysis. We first apply linear discriminant analysis (LDA) to the 2000/2000 data. The error rate is 5.75% for the training set and 8.3% for the test set. We then apply LDA to the 1991/7188 data. As expected, the result is even worse, 7.2% for the training set and 10.7% for the test set.

We take a closer look at the training set of 2000/2000 data. The first three canonical variates from LDA are shown in Figure 2.3. Noticeable clustering patterns can be found. For example, a good portion of digit 0 visually separable from other digits is found in the lower left corner of Figure 2.3(a). In the same figure, digit 1 appears to occupy another corner.

We are led to the suggestion of isolating these two clusters from the rest of data because they appear easier to classify than others. To do so, we can formulate a three-way classification problem by considering "0" as one class, "1" as another class, and

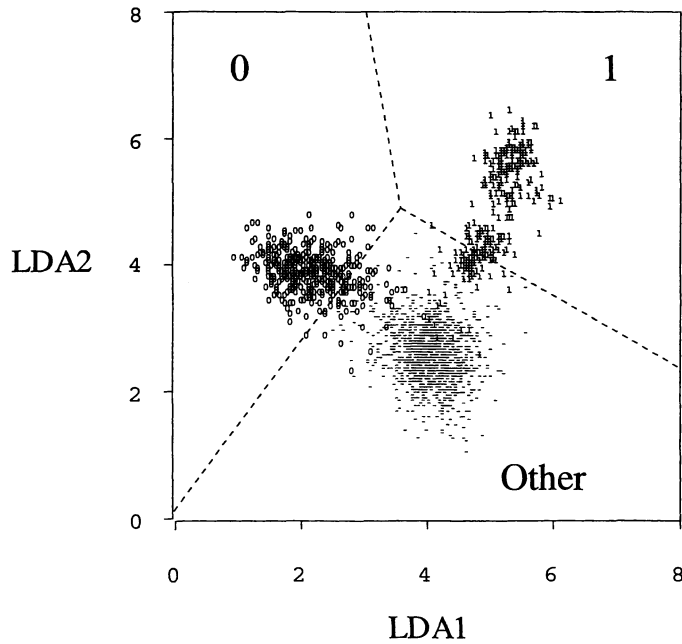


FIGURE 2.4. LDA Cannonocal Variates for Groups 0, 1 and Other(-).

pooling all other digits together as a third class called "Other". We first reduce the dimension of the feature space from 55 to 2 using the canonical variates from LDA. A discussion on the effectiveness of dimension reduction from the theory of sliced inverse regression (Li 1991) is offered in the Appendix. Figure 2.4 shows the scatterplot of the two canonical variates. The clustering pattern is more clear-cut than what is seen in Figure 2.3. Since we are interested in isolating cases that can be predicted with higher precision, the boundaries of discriminant regions for each class are pushed a bit toward the centers of 0 and 1 - this is done by setting the prior probability of 0's and 1's to be .0005 each. After we isolate these two clusters with mostly 0 and 1, we proceed by considering another three-way classification for the class "Other". This time we choose digits 6 and 7 as two classes and pooling all other digits together as "Other". We continue to partition the left-over "Other" until we obtain a three-way classification tree as shown in Figure 2.5.

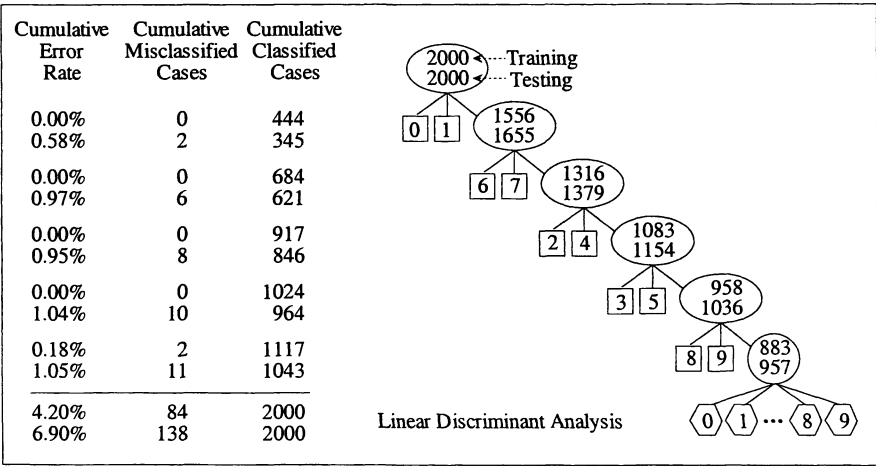


FIGURE 2.5. Three-Way Classification Tree for the Hand-Written Digit Recognition Problem.

From the tree obtained and the scatterplots generated in our analysis, we see that a good portion of the data form distinctive clusters which can be isolated in an iterative fashion. Out of the 2000 test cases, about a half of them can be classified with 1% of error. What should we do with the left-over cases? One simple suggestion is to apply the linear discriminant analysis. Another possibility is to repeat the same procedure again before linear discriminant analysis is applied. The result is given in Figure 2.6.

One problem with this three-way tree approach is how to decide the ordering of partition. Why are 0 and 1 chosen first? There are certainly many criteria that can be used. We simply choose the pair which gives the smallest error rate for the training data. Another problem is the choice of prior in determining how much the lines should be pushed away from the center of the "Other". This is another optimization problem, which needs to be solved.

While the above lines of thought may be worth pursuing further, in this article we shall alter our strategy a bit in order to bypass such difficult optimization decisions. There is no need to generate three-way trees. Instead, we shall synthesize results from all three-way classifications in a way to be discussed in the next section.

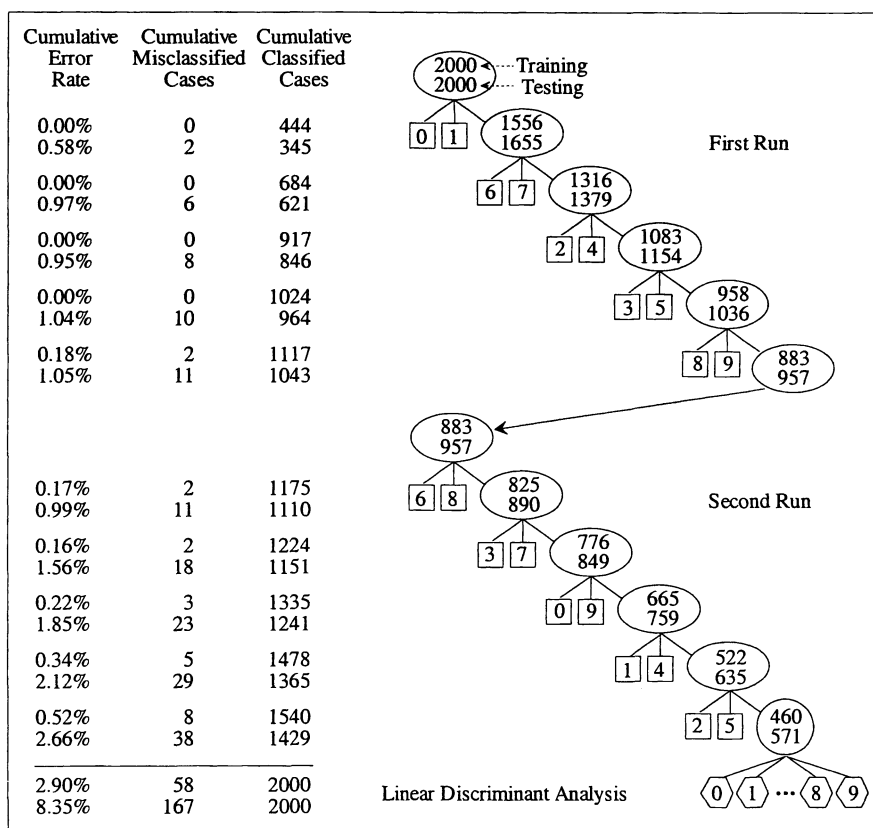


FIGURE 2.6. Repeated Three-Way Classification Tree for the Hand-Written Digit Recognition Problem.

3. THREE-WAY SUBCLASSIFICATION

We begin with a general description of our three-way subclassification method. Suppose there are k groups under consideration. For each pair of groups, i and j , we formulate a three-group problem - group i , group j , and a third group which is the union of all other $k - 2$ groups. Suppose from the training data, a classifier denoted by T_{ij} , is obtained. To a unit with feature \mathbf{x} , the classifier will assign the membership according to $T_{ij}(\mathbf{x}) = i, j$ or *Other* with *Other* standing for the third group. This is carried out for all of the $k(k - 1)/2$ three-group problems. Thus for each \mathbf{x} , we have $k(k - 1)/2$ values, $T_{ij}(\mathbf{x}), 1 \leq i < j \leq k$.

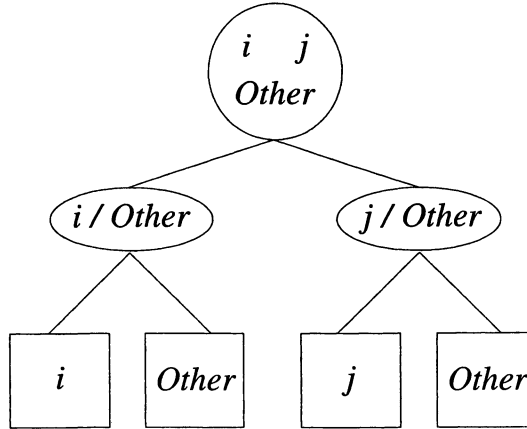


FIGURE 3.1. Tree-Structure for Three-Way Subclassification.

For the digit data, there are $k = 10$ classes, yielding a total of 45 subclassification problems. For each subclassification, we use a linear partition that has a simple tree structure (Figure 3.1). The tree consists of two levels of partitions - the first one is to divide the training set into two nodes $i/Other$ and $j/Other$ along the projection that best separates class i and class j . Node $i/Other$ is expected to contain most members from digit i and a good portion of members from digits other than i and j . This node is further divided into two children nodes, i and $Other$ again by a best linear partition rule. The other node $j/Other$ is similarly partitioned into two nodes, j and $Other$. This rule is chosen out of simplicity.

3.1. Majority rule. To combine the results from each smaller problem of classification, we first count the frequency that each class is assigned. For any fixed unit \mathbf{x} , we obtain a k -dimensional vector $(c_1(\mathbf{x}), \dots, c_k(\mathbf{x}))'$ where for each l between 1 and k , $c_l(\mathbf{x})$ equals the total number of times that $T_{ij}(\mathbf{x})$, $1 \leq i < j \leq k$, equals l . We can think of each classifier T_{ij} as assigning the winner for a match between players i and j , but with possibility that both lose and $T_{ij}(\mathbf{x})$ equals 'Other'. The vector $(c_1(\mathbf{x}), \dots, c_k(\mathbf{x}))'$ is just the score board showing the winning record for each player. Then it should be clear that the largest value that $c_l(\mathbf{x})$ can take is $k - 1$. We can rearrange the score board in a non-increasing order, $M_1(\mathbf{x}) \geq M_2(\mathbf{x}) \geq \dots \geq M_k(\mathbf{x})$.

Table 3.1. *Misclassification Matrices for Three-way Subclassification Using Majority Rule with the 55-Features-2000/2000 data: (a). Training Set; (b). Test Set.*

(a)

$y \backslash \hat{y}$	0	1	2	3	4	5	6	7	8	9	Other
0	0	0	0	1	0	0	2	0	0	1	6
1	0	0	2	0	0	0	0	0	0	1	6
2	0	1	0	1	1	1	3	0	2	1	2
3	0	0	3	0	0	0	0	0	0	1	4
4	0	1	0	0	0	0	1	0	0	11	0
5	2	0	0	1	0	0	0	0	1	0	0
6	3	1	0	0	0	1	0	0	0	0	2
7	0	0	1	0	0	0	0	0	0	3	0
8	2	0	0	0	0	0	0	0	0	1	4
9	0	0	0	0	0	1	0	1	1	0	5

(b)

$y \backslash \hat{y}$	0	1	2	3	4	5	6	7	8	9	Other
0	0	0	0	1	0	0	0	0	0	0	5
1	0	0	0	0	0	0	0	0	0	1	7
2	1	0	0	2	1	0	0	2	1	2	2
3	4	0	4	0	0	0	0	0	0	1	4
4	0	1	1	1	0	0	8	1	0	8	0
5	6	0	3	6	0	0	2	0	1	2	3
6	1	0	1	0	5	3	0	0	0	0	2
7	0	0	0	0	0	0	0	0	0	6	1
8	1	0	0	0	1	1	1	0	0	1	8
9	1	2	0	0	1	0	1	11	0	0	1

Consider the majority rule for the final class membership assignment. We assign class l to unit \mathbf{x} if $c_l(x) = M_1(\mathbf{x})$. Randomization is used as a tie-breaker.

We apply the majority rule to the 2000/2000 data and the 7188/1991 data. The results are summarized by two misclassification matrices in each case one for the training set and the other for the test set; Tables 3.1(a)-(b) and 3.2(a)-(b). A cell in each matrix shows the number of times a digit in the beginning of a row is misclassified as another digit on the top of the corresponding column. For example, take a look at the

Table 3.2. Misclassification Matrices for Three-way Subclassification Using Majority Rule with the 55-Features-7188/1991 data: (a). Training Set; (b). Test Set.

(a)

$y \backslash \hat{y}$	0	1	2	3	4	5	6	7	8	9	Other
0	0	0	0	2	0	1	5	0	0	2	16
1	0	0	3	0	0	0	0	0	0	2	28
2	2	3	0	13	1	3	5	2	6	6	16
3	4	0	5	0	0	6	0	1	2	2	20
4	0	7	2	2	0	0	24	2	3	34	0
5	13	0	4	6	0	0	5	0	2	2	15
6	8	3	0	0	2	7	0	0	1	0	10
7	1	0	2	0	0	0	0	0	1	20	3
8	4	0	0	2	2	6	1	0	0	4	24
9	0	2	0	1	1	1	1	14	1	0	12

(b)

$y \backslash \hat{y}$	0	1	2	3	4	5	6	7	8	9	Other
0	0	1	0	0	0	0	1	2	1	1	14
1	0	0	3	1	0	0	2	1	2	2	5
2	3	2	0	2	1	4	2	1	2	1	5
3	3	0	2	0	0	7	0	1	2	4	7
4	0	4	3	0	0	0	4	2	1	18	1
5	4	0	0	5	0	0	3	1	0	1	5
6	3	0	1	0	3	2	0	0	0	1	1
7	0	1	0	1	2	1	0	0	2	5	1
8	1	1	3	4	0	1	0	0	0	2	8
9	0	0	0	0	1	0	1	2	0	0	3

row with the digit 2 in the training set of 7188/1991 data. We see that 2 is misclassified as 0 two times, as 1 three times, and so on. We see that 2 is misclassified as 0 two times, as 1 three times, and so on. It is also unclassified 16 times as the last column “Other” indicates. For the 2000/2000 data, the overall error rate is 2.6% for the training and 4.85% for the test set; the unclassified rate is 1.45% for the training and 1.65% for the test set. For the 7188/1991 data, we have error rates 3.7% (training set) and 6.9% (test set) and unclassified rates 2.0% (training set) and 2.5% (test set). In general, these

Table 3.3. Conditional Error Matrices for Three-way Subclassification with the 55-Features-2000/2000 data (Number of Misclassification Cases / Number of Total Cases):
 (a). Training Set; (b). Test Set.

(a)

$M_1 \backslash M_2$	0	1	2	3	4	5	6	7	8
0	29/29								
1	0/8	2/2							
2	0/8	4/8	0/1						
3	3/7	0/5	3/4	0/1					
4	2/10	2/3	1/2	0/1	1/1				
5	0/16	0/3	1/3	0/1	2/3	0/1			
6	0/12	0/3	1/3	2/5	0/0	1/1	1/1		
7	0/19	3/10	1/2	0/0	0/1	0/0	2/4	0/1	
8	3/43	2/4	1/2	1/1	0/2	0/0	0/0	0/2	1/1
9	4/1604	0/64	0/20	1/15	0/13	0/12	0/6	0/8	0/3

(b)

$M_1 \backslash M_2$	0	1	2	3	4	5	6	7	8
0	33/33								
1	14/32	6/7							
2	3/19	4/5	1/1						
3	1/14	3/6	4/6	0/1					
4	1/9	1/3	0/2	1/2	1/2				
5	3/19	1/4	1/3	1/1	2/2	0/0			
6	2/21	1/9	1/2	2/3	0/1	2/3	0/0		
7	1/18	4/11	1/2	0/2	4/6	2/3	1/2	1/2	
8	2/29	2/12	1/6	0/1	3/7	0/3	0/2	1/4	0/2
9	10/1460	2/105	1/44	1/21	1/17	0/12	2/8	0/6	1/5

numbers are not impressive. There are better ways of using three-way subclassification than the straightforward application of the majority rule. This is to be discussed next.

3.2. Conditional error analysis and unanimity. Our conditional error analysis exploits the largest two values on the score board for each unit \mathbf{x} , $M_1(\mathbf{x})$ and $M_2(\mathbf{x})$.

The most interesting condition is when $M_1(\mathbf{x}) = k - 1$ and $M_2(\mathbf{x}) = 0$. Suppose there is an ideal unit \mathbf{x} from say class l , which is very easy to classify. Then for this unit the maximum score $M_1(\mathbf{x})$ is very likely to be $k - 1$ and it is expected to be achieved by class l , $c_l(\mathbf{x}) = M_1(\mathbf{x}) = k - 1$. This is because when class l competes with any other group ($k - 1$ times in total), the classifier should return l . On other hand, when the competition is between any two classes i, j other than class l , the classifier should return *Other* because class l is contained in “Other”. Thus $(M_1, M_2) = (k - 1, 0)$ represents the situation where an unanimous decision is reached. We anticipate the final classification to be most accurate under this condition.

For the 2000/2000 data, we find that there are 1604 cases in the training set with $(M_1, M_2) = (9, 0)$. From this subset, there are only 4 misclassification cases, representing an error rate of 0.25% which is much smaller than the overall error rate 2.6% given earlier. In the test set, there are 1460 cases with $(M_1, M_2) = (9, 0)$, and 10 out of them are misclassified. This amounts to 0.68% of conditional error, which are again much smaller than the overall error 4.85%. Substantial reduction in conditional error rate also occurs for the 7188/1991 data - 31 out of 5592 (=0.55%) for the train-

Table 3.4. Conditional Error Matrices for Three-way Subclassification with the 55-Features-7188/1991 data (Number of Misclassification Cases / Number of Total Cases): (a). Training Set; (b). Test Set.

(a)	$M_1 \backslash M_2$	0	1	2	3	4	5	6	7	8
	0	144/144								
	1	23/93	9/16							
	2	15/72	14/28	2/2						
	3	8/52	6/15	3/10	3/9					
	4	6/63	5/13	1/3	1/4	1/1				
	5	6/53	4/10	4/7	2/3	2/3	1/2			
	6	9/63	2/9	2/10	2/3	2/6	3/6	2/3		
	7	6/67	7/15	2/7	3/5	7/13	3/3	2/4	4/5	
	8	10/114	12/29	2/12	5/12	4/4	2/4	3/8	8/10	3/4
	9	31/5592	4/257	1/98	2/51	0/55	2/32	2/35	2/30	2/19

(b)

$M_1 \backslash M_2$	0	1	2	3	4	5	6	7	8
0	50/50								
1	14/42	5/10							
2	8/32	4/9	3/3						
3	1/17	0/1	4/6	2/3					
4	4/15	0/6	2/2	4/4	0/0				
5	1/18	2/2	2/2	2/2	1/2	1/1			
6	0/19	4/6	1/3	2/3	2/4	1/2	0/0		
7	2/20	1/4	1/4	0/1	1/2	1/1	0/0	2/2	
8	5/49	1/5	2/3	0/3	3/4	0/0	1/2	5/7	0/3
9	20/1407	5/100	5/33	2/19	1/11	1/19	2/10	3/8	4/10

ing set and 20 out of 1407 ($=1.4\%$) for the test set, as compared to the overall rates of 3.7% and 6.9% respectively.

For each fixed value of M_1 , we can also anticipate the quality of final classification to go down as M_2 increases because this reflects that the leader faces a stronger challenge from the runner-up and the degree of unanimity goes down. Similarly, for a fixed M_2 , the classification quality also degrades as M_1 decreases. Such trends can be found from Tables 3.3 (a)-(b) and 3.4 (a)-(b). For example, in Table 3.3 (a) and (b), combining numbers from 4 cells corresponding to $M_1 = 5, 6, 7, 8, M_2 = 0$ the error rates are 3/90 (training) and 8/87 (test) which are lower than the corresponding error rates for $M_1 = 5, 6, 7, 8, M_2 = 1$ - 5/20 (training) and 8/ 36 (test). For $M_1 = 9, M_2 = 1, \dots, 5$, the error rates are 1/124 (training), 5/ 199 (test) which are lower than the corresponding error rates for $M_1 = 8, M_2 = 1, \dots, 5$ - 4/9 (training) and 6/ 29 (test).

4. FURTHER CONSIDERATIONS

In this section, we discuss possible ways of enhancing the three-way subclassification approach.

4.1. Conditional error analysis for binary classification. Unlike three-way subclassification, binary subclassification does not have a straightforward apparatus for unanimity assessment. One possibility is to follow a similar conditional analysis as in the three-way approach. Let $M_1(\mathbf{x})$ and $M_2(\mathbf{x})$ be the highest two scores again, but

now obtained from the scoreboard by binary-subclassifications. The ideal condition for most accurate prediction requires $M_1(\mathbf{x})$ to be as large as possible and $M_2(\mathbf{x})$ be as small as possible. The larger the gap between them, the less competitive the runner up is, thus reflecting certain degree of unanimity.

For the digital problem, since there are 10 classes, it is easy to argue that if $M_1 = 9$, then M_2 cannot be smaller than 5. The condition $M_1 = 9, M_2 = 5$ represents the most favorable situation for better classification. We shall anticipate the error rate to increase as M_2 increases. This is indeed what we can find from Tables 4.1 (a)-(b) and 4.2 (a)-(b). For example, in the test set for 2000/2000 data, when fixing M_1 at 9, the error rates are seen to increase : 0/3, 0/96, 15/506 (=2.96%), 50/1337 (=3.7%), respectively for $M_2 = 5, 6, 7, 8$. However, an undesirable pattern is that the most favorable condition $M_2 = 5$ is only satisfied by three cases, while the least favorable condition $M_2 = 8$ has 1337 cases. Thus the conditional error analysis on binary classification does not lead to a useful way of finding a large portion of cases which can be classified with very high precision. The error rate has already reached $15/(3 + 96 + 506) = 2.5\%$ when conditioning on $M_1 = 9, 5 \leq M_2 \leq 7$ as compared to $10/1460 = 0.68\%$ for $M_1 = 9, M_2 = 0$ from three-way subclassification reported earlier.

However, a positive note for binary classification is that it can be used to improve the non-(9,0) group from three-way subclassification. The error rate is reduced from $120/540 = 22.22\%$ (among which $33/540 = 6.11\%$ were unclassified) to $94/540 = 17.41\%$.

Table 4.1. Conditional Error Matrices for Binary Classification with the 55-Features-2000/2000 data (Number of Misclassification Cases / Number of Total Cases): (a). Training Set; (b). Test Set.

(a)

$M_1 \backslash M_2$	5	6	7	8
5	0/0			
6	0/0	0/0		
7	0/0	0/0	0/0	
8	0/0	0/0	0/3	1/8
9	0/4	0/113	4/571	15/1301

(b)

$M_1 \backslash M_2$	5	6	7	8
5	0/0			
6	0/0	0/0		
7	0/0	1/1	4/5	
8	0/0	1/1	12/18	21/33
9	0/3	0/96	15/506	50/1337

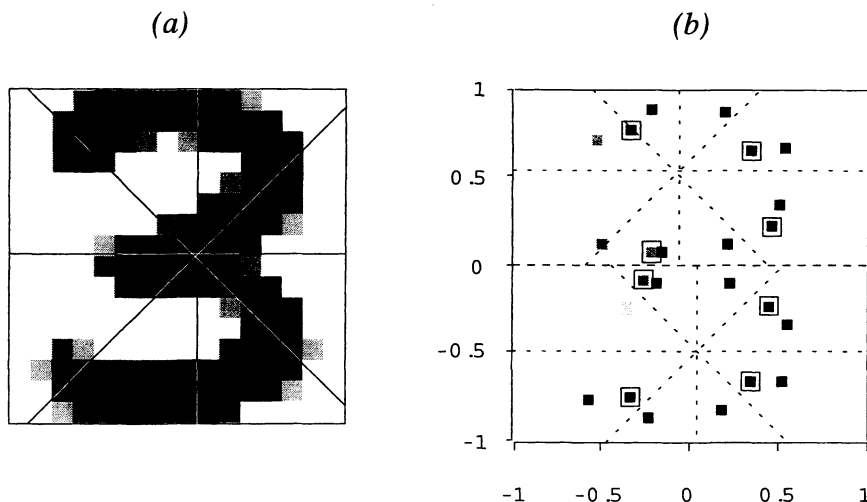


FIGURE 4.1. (a). First Level Radial Partition; (b). Locations of Mass Centers with Weights for three radial partitions

Table 4.2. Conditional Error Matrices for Binary Classification with the 55-Features-7188/1991 data (Number of Misclassification Cases / Number of Total Cases): (a). Training Set; (b). Test Set.

(a)

$M_1 \setminus M_2$	5	6	7	8
5	0/0			
6	0/0	0/0		
7	0/0	0/0	0/0	
8	0/0	0/0	12/20	41/62
9	0/4	1/117	15/1456	119/5529

(b)

$M_1 \setminus M_2$	5	6	7	8
5	0/0			
6	0/0	0/0		
7	0/0	1/1	2/2	
8	0/0	0/0	7/11	19/23
9	0/2	1/22	7/409	107/1521

4.2. Radial partition. Encouraged by the promising results from three-way analysis, we apply the same strategy to another set of feature variables. This new feature space has 69 variables. Just like the old feature space, they are also constructed using the centre of mass to guide the partition. The difference comes from the geometric configuration of the partitioning lines. We use an 8-region radial partition system.

We begin with the centre of mass of the whole digit. Instead of quadrants, 8 regions are obtained by further partitioning each quadrant diagonally into two equal

pieces. The location of the mass center for each of these 8 regions is computed. After that, the whole digit is horizontally divided into two halves - one half is above the x-axis and the other half is below the x-axis. For each half, we then compute the new centre of mass and apply the 8 region radial partition system to get another 8 locations of mass centers. Figure 4.1(a) shows how the same digit 3 is partitioned. A total of 24 mass centers are located in Figure 4.1(b). Our new feature space consists of these 48 locational variables and together with 21 weight variables. Each weight variable represents the total mass from one of the 24 regions obtained before. Note that due to colinearity among each of the three 8-region partitions, we cannot use all 24 weight variables.

Using this new set of feature variables, the error rate from LDA is $5.8\% = 425 / 7291$ for the training set and $9.9\% = 199 / 2007$ for the test set. They are not much different from the LDA results by 55 features. Can the three-way method help filter out a group of high quality cases?

The result is shown in Table 4.3. Here we find that for the unanimous winners, $M_1 = 9, M_2 = 0$, the error rate in the test set is reduced to below 1% (13 out of 1535 cases).

4.3. A combined use of different feature spaces. As pointed out before, three-way subclassification is especially effective in isolating high quality cases - cases that are easier to predict their membership. This is achieved by a conditional error analysis which exploits the degree of unanimity among different classification decisions from subclassification. After locating the very high quality cases, we can then focus on the rest of cases and try to find better classification, perhaps even with drastically different classification methods. For example, consider the nearest neighbor classifier used by LeCun *et al* (1989). As mentioned before, it has an error rate of 2.5% for the 2007 test cases, but is computationally very demanding. If we can use it only for the non- ($M_1 = 9, M_2 = 0$) cases, then the overall error rate would be still be at most around 3%. But in this way, we have allocated the total computation time more effectively without sacrificing much of the overall classification quality.

Perhaps an easier way of improvement is to combine the results from different feature spaces. We try the following path.

Table 4.3. Misclassification Matrix for Three-way Subclassification Using Majority Rule with the Feature 69-7291/2007 data.

$M_1 \backslash M_2$	0	1	2	3	4	5	6	7	8
0	35/35								
1	9/26	7/10							
2	3/13	6/10	2/3						
3	7/17	4/6	6/9	4/4					
4	2/14	1/2	0/2	0/3	1/1				
5	4/12	2/3	3/4	0/1	0/1	1/1			
6	1/11	1/4	2/3	2/2	2/2	1/2	0/0		
7	3/18	3/5	2/4	1/2	2/4	2/4	0/1	1/1	
8	3/28	2/11	4/7	0/3	1/2	3/5	2/2	3/6	1/1
9	13/1535	5/69	1/22	1/15	1/15	0/12	1/10	3/12	4/12

(1). Apply three-way subclassification trained by 69-features to all test cases and locate the $M_1 = 9, M_2 = 0$ group.

(2) Apply three-way subclassification trained by 55-features to the non-(9, 0) cases and locate the $M_1 = 9, M_2 = 0$ cases.

(3) Apply binary classification trained by 69-features to all other left-over cases.

Table 4.4. Breakdown of Error Rate for Combining Use of Different Feature Spaces.

Group	Cases	Misclassified Cases	Error Rate
(1). $M_1 = 9, M_2 = 0$ /3-way/69-features	1535	13	0.85%
(2). $M_1 = 9, M_2 = 0$ /3-way/55-features	160	15	9.38%
(3). Left-over/binary/69-features	312	87	27.88%
Total	2007	115	5.73%

A breakdown of the error rate is given in Table 4.4. The overall error rate is now about 5.7%, which is compatible with the result (between 5% and 6%) by the neural network approach. It is certainly a great improvement over the original LDA result which is about 10% for either feature. It is also significantly better than the 8.2% error rate obtained by Hastie, Buja, and Tibshirani (1995).

It is interesting to observe that there are a good number of tied scores in binary subclassification. In order to keep the procedure simple, we resolve these ties essentially by a random choice. Further investigation on how to handle these cases seems worthwhile.

5. DISCUSSION

Discriminant analysis is relatively easier when the number of classes is small. This prompts us to consider the subclassification strategy for alleviating the complexity of many classes. We propose a three-way subclassification method and show that it can be fruitfully applied to complement binary subclassification.

Three-way subclassification is designed to exploit the degree of unanimity among various decisions during subclassification. The information gathered from the conditional error rate analysis is used to tell if an incoming unit \mathbf{x} is easy to classify or not. If it falls into the unanimous decision group, $(M_1(\mathbf{x}), M_2(\mathbf{x})) = (k - 1, 0)$, then we are most confident about the classification accuracy. On the contrary, if it has small $M_1(\mathbf{x})$ or large $M_2(\mathbf{x})$, then this is a harder case and we may send it to other classifiers for a better result.

In many industrial applications, tolerable error rates are usually set up by concerns from the economic aspect. It is important to identify sub-populations that are easier to classify than others because it helps engineers to identify where the quality improvement may come from. Further theoretical investigation on discriminant analysis along this line of thoughts is worth pursuing.

Acknowledgments. Li's research is supported in part by NSF grants DMS-09803459, DMS-0104038 and DMS-0201005.

6. APPENDIX

First we carry out a first moment-based SIR (Li, 1991) analysis to reduce the dimension p of the feature variable from 55 to 2. In discriminant analysis, SIR can be carried out by finding the eigenvectors \hat{v}_i in the following eigenvalue decomposition:

$$\hat{\Sigma}_B \hat{v}_i = \hat{\lambda}_i \hat{\Sigma}_x \hat{v}_i$$

where $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$ are the eigenvalues, $\hat{\Sigma}_B$ denotes the covariance matrix between the class means, and $\hat{\Sigma}_x$ is the overall covariance matrix of \mathbf{x} in the training data. Let $\mathbf{x}_i, i = 1, \dots, n$ be the training data set, $\bar{\mathbf{x}} = n^{-1} \sum \mathbf{x}_i$ be the overall mean, and $\bar{\mathbf{x}}_l$ be the sample mean for the l -th class. We have $\hat{\Sigma}_x = n^{-1} \sum (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})'$ and $\hat{\Sigma}_B = n^{-1} \sum n_l (\bar{\mathbf{x}}_l - \bar{\mathbf{x}})(\bar{\mathbf{x}}_l - \bar{\mathbf{x}})'$. Since the rank of $\hat{\Sigma}_B$ is equal to 2 in our three-way subclassification, all but first two eigenvalues $\hat{\lambda}_i$ must be zero. We shall use only the first two directions \hat{v}_1, \hat{v}_2 to project the data.

Clearly this eigenvalue decomposition can be arranged in a form equivalent to Fisher-Rao's procedure for finding canonical variates :

$$(\hat{\Sigma}_W)^{-1} \Sigma_B \hat{v}_i = \hat{\theta}_i \hat{v}_i$$

where $\hat{\Sigma}_W = \hat{\Sigma}_x - \hat{\Sigma}_B$. Thus the general theory of SIR as described in Li(1991) can be applied to justify that the canonical variates $\hat{v}'_1 \mathbf{x}$ and $\hat{v}'_2 \mathbf{x}$ can be effective in reducing dimensionality even if the standard normal assumption used in justifying linear discriminant analysis is violated. As our figures have shown, the distribution of \mathbf{x} for each digit is not normal.

REFERENCES

- [1] Fisher, R.A. (1936). The use of Multiple Measurements in Taxonomic Problems, *Ann. Eugen.*, **7**, 179-188.
- [2] Hastie, T., Buja, A., and Tibshirani, R. (1995). Penalized discriminant analysis. *Ann. Stat.* **23**, 73-102.
- [3] LeCun, Y. Boser, B., Denker, J.S., Henderson, D, Howard, R.E., Hubbard, W., and Jackel, L.D. (1989). Backpropagation applied to handwritten zip code recognition. *Neural Computation.* **1**, 541-551.
- [4] Li, K.C. (1991) Sliced inverse regression for dimension reduction (with discussion). *J. Amer. Stat. Assoc.*, **86**, 316-342.
- [5] Simard, P.Y., LeCun, Y. and Denker, J. (1993). Efficient pattern recognition using a new transformation distance, in " *Advances in Neural Information Processing Systems*," Morgan Kaufman, San Mateo, CA, pp 50-58.

CHUN-HOUH CHEN
INSTITUTE OF STATISTICAL SCIENCE
ACADEMIA SINICA
TAIPEI 115, TAIWAN
cchen@stat.sinica.edu.tw

KER-CHAU LI
DEPARTMENT OF STATISTICS
8130 MATH SCIENCES BLDG.
BOX 951554
LOS ANGELES, CA 90095-1554
kcli@stat.ucla.edu