

Secondary Classification Using Key Features

V.Govindaraju, Z.Shi, A.Teredesai

Department of Computer Science and Engineering,
Center of Excellence for Document Analysis and Recognition,
State University of New York at Buffalo, Buffalo, NY 14260, U.S.A

ABSTRACT

A new multiple level classification method is introduced. With an available feature set, classification can be done in several steps. After first step of the classification using the full feature set, the high confidence recognition result will lead to an end of the recognition process. Otherwise a secondary classification designed using partial feature set and the information available from earlier classification step will help classify the input further. In comparison with the existing methods, our method is aimed for increasing recognition accuracy and reliability. A feature selection mechanism with help of genetic algorithms is employed to select important features that provide maximum separability between classes under consideration. These features are then used to get a sharper decision on fewer classes in the secondary classification. The full feature set is still used in earlier classification to retain complete information. There are no features dumped as they would be in feature selection methods described in most related publications.

Keywords: Classification, Feature Selection, K-nearest Neighbor, Handwritten character recognition

1. INTRODUCTION

We present in this paper a novel technique to improve the recognition accuracy of a digit recognizer without significantly increasing the computational time. It involves two levels of classification. The primary classification is simply the regular run of the recognizer. It is the secondary level of classification that is the focus of our methodology.

The approach we propose has some similarities with the method of serial combination of classifiers. Serial classifiers combination has emerged as a new direction in pattern recognition with the potential to improve classification performance by combining the strengths and overcoming the limitations of individual classifiers. When multiple classifiers or experts are available, recognition performance can be improved by combining classifiers at the decision level. The *parallel* combination of classifiers has received considerable attention in recent times. Voting schemes, rank-based schemes such as Borda Count, and theories of evidence combination such as Bayes, Logistic Regression, Fuzzy Integrals and Dempster-Shafer have been applied to combine classifiers at the abstract, rank and measurement levels.^{4,6}

In contrast, the *serial* combination of classifiers has not received much systematic study. In the typical serial combination of two classifiers, a fast classifier – the *reducer* – is used to reduce or filter the lexicon as a preprocessing step for a better and generally slower classifier – the *recognizer* (Figure 1). Serial combinations have been used in the past primarily for reasons of efficiency.

For example, the question could be as follows: Which of the 4 digits {1,4,7,9} is the image. What is called for here is not simply using the regular recognizer and looking for the ranks of the numerals of interest but building a specific classifier which focuses on only those features which are useful in telling the numerals in the subclass apart. The subclasses can vary dynamically in an application.

Genetic techniques like Genetic Algorithms and Genetic Programming have been used recently for feature selection. These techniques strive to provide a global optima in terms of the feature subset required for a particular task of classification. We employ genetic algorithms to provide us with secondary classifiers that are helpful in making a sharper decision when confusion between top choices from the first stage is very high.

Further author information: (Send correspondence to V.Govindaraju)

V.Govindaraju: E-mail: govind@cedar.buffalo.edu

Z.Shi: E-mail: zshi@cedar.buffalo.edu

A. Teredesai: E-mail : amt6@cedar.buffalo.edu

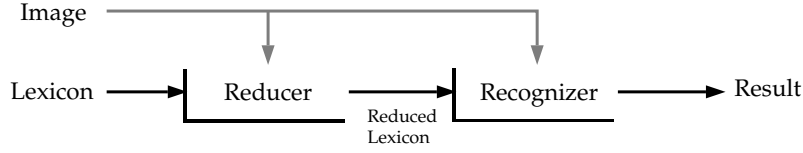


Figure 1. Serial combination of two classifiers

Using a set of isolated digit images (BHA set), we have built two data sets. The first set is to study the **correct rate** of the recognizer regardless of recognition confidence. Table 1 and table 2 show the confusion matrix.

Total character = 12888 The top 1 hit ratio is 0.9809

Table 1. Confusion matrix for digits 0-9.

	0	1	2	3	4	5	6	7	8	9
0	1640	4	5	2	1	0	6	3	1	0
1	1	3017	2	0	2	0	6	6	0	0
2	14	8	1998	6	4	2	2	24	9	3
3	1	1	6	774	1	4	0	3	3	2
4	3	13	0	0	2197	0	6	1	2	18
5	4	1	1	4	0	747	2	1	4	2
6	5	3	1	1	0	1	615	0	1	1
7	3	1	0	0	2	0	1	689	0	13
8	1	1	1	0	2	1	0	0	485	2
9	2	1	0	1	0	0	0	5	2	480

Table 2. Digit occurrence in top choice.

	top 1	top 2	top1 + top2	total
0	1640(0.987)	9(0.005)	1649(0.992)	1662
1	3017(0.994)	15(0.005)	3032(0.999)	3034
2	1998(0.965)	34(0.016)	2032(0.982)	2070
3	774(0.974)	12(0.015)	786(0.989)	795
4	2197(0.981)	31(0.014)	2228(0.995)	2240
5	747(0.975)	8(0.010)	755(0.986)	766
6	615(0.979)	4(0.006)	619(0.986)	628
7	689(0.972)	12(0.017)	701(0.989)	709
8	485(0.984)	3(0.006)	488(0.990)	493
9	480(0.978)	6(0.012)	486(0.990)	491

The (top 1 + top 2) hit ratio is 0.9913

On the second data set we use a threshold value of 0.8. The table shows the top choice correct rate after thresholding.

Based on the above analysis we want to build pairwise discriminators. The candidates for the pairs of digits of high interest are the ones with the high confusion rate. Since the top 2 choice correct rate for the GSC digit recognizer is very high we are justified in keeping this as a pairwise discrimination task and not expanding the subclass.

To build the pairwise discriminators we use the 512 GSC features and select a subset of features these features which are ‘key’ in telling the two digits in the pair apart. For example, when discriminating ‘7’ and ‘2’ the features localized at the bottom right corner will be ‘key’. The ‘key’ features are selected using a genetic algorithm. With

a generated set of prototypes, the genetic algorithm finds a subset of the 512 feature vector that will give a small distance among samples of the same class and large distance between samples of different classes.

A recent study of the GSC digit recognizer shows that the recognition rate is 96% after thresholding although the top one hit ratio is 98%. Based on this study, we are developed pairwise discriminators. Each of these pairwise discriminators works for a subclass having only two classes (i.e. 2 numerals). These subclasses are precisely the top two choices of the GSC.

A classifier using K-nearest neighbor algorithm combined with the key features has been implemented and tested. This classifier is based on the GSC classifier. Though the classifier gives high hit rate and meaningful confidence values improvements can be made by adding new prototypes for variations of handwritten character. However adding templates will cause increase in the size of the prototype database and thus hurt the speed of classification.

We experimented with the use of a neural network trained on the existing GSC prototype as well as a different prototype set of smaller size. Training was also done for further pairwise classification after the GSC stage as described above. For a pair of digits, 2 and 7, we build a feature set for three classes, each for the two digits and the third for the rest of the digits (may include anything other than the two digits). The feature set used is the key feature set developed with the help of GA. After training, the testing shows that the neural network is able to pick up the misclassified 2 or 7 to give a correct answer. If the input image is not of any the two digits, the answer will be the third class which is meant for a reject.

The GA employed for key feature selection is based on the following strategy: For each one of the 45 pairs of digit classes we want to find the positions in the feature vector which mark the distinguishing features between the pair. Using a set of 25656 prototypes from the GSC classifier, we have plotted the distribution of features for each class. These plots are studied to find the positions in the vector which contribute to the distinguishing characteristics. If the distance vector takes only those positions into consideration, the chances of telling the digits apart increase. The use of GA enables us to find those features that can tell the digits apart with maximal separability.

A set of 12242 handwritten digit images were used as testing data. 540 images are rejected when a thresholding value of 0.4 is used on the confidence value of top choices (the confidence value ranges from 0 to 0.6 for the GSC). Among the accepted (12242-540) images, 60 are incorrect.

$$\text{Accept rate} = (12242-540)/12242 = 95.6 \%$$

$$\text{Error rate} = 60/(12242-540) = 0.51\%$$

By using the key features on the 540 rejected cases, the following results were obtained as seen in table 3. The secondary classification returns results on 454 cases.

Table 3. Results after key features were used in recognition

method	correct	error	input
1.	281	74	454
2.	234	44	454
3.	138	14	299
4.	111	12	187
5.	66	8	117

Combining the two stages of classification using the 3rd level thresholding in the table above, the performance is as follows.

$$\text{Accept rate} = (12242-540 + 138)/12242 = 96.7\%$$

$$\text{Error rate} = (60+14)/(12242-540+138) = 0.63\%$$

So the accept rate increases by more than a point while the error rate increases by 0.1%.

Recognition of handwritten characters by a computer has been a topic of intense research for many years.⁵ Due to its pivotal role in many applications such as postal addresses interpretation, bank checks, tax forms and census forms reading, considerable research has been conducted in the area of handwritten character and numeral recognition.³

Due to the nature of handwritten characters, there is probably no one single method that can achieve high recognition and reliability rates for all applications. Thus in recent years, efforts have been directed toward more sophisticated considerations.²

2. METHODOLOGY

Feature extraction problem is defined by Devijver and Kittler as the problem of “extracting from the raw data the information which is most relevant for classification purposes, in the sense of minimizing the within-class pattern variability while enhancing the between-class pattern variability”.¹ Handwriting varies by writing styles, writing equipment, resolution and image quality. To design a good feature set for the purpose of machine classification with high recognition performance, many aspects should be put in consideration. This often leads to multi-resolution features for different scales of the character images and large size of feature sets being extracted. Classification methods designed for using these features will encounter difficulties in achieving high recognition rate. Moreover if the classification scheme is threshold based then choice of good thresholds will affect reliability of results.

Smaller number of classes of characters or numerals is easier to classify. This is so because the features required may be relatively fewer and the classification decisions are sharper. Then we can decouple the problem of achieving high recognition rate from the problem of achieving high reliability. The large feature set can be utilized to get high recognition rate (high number of hits without considering thresholding). Then the resulting outputs of this first step shrink the target down to a smaller range in which fewer number of classes, for example 2, are to be classified. For each possible group of smaller classes, a subset of features is selected for a secondary classification.

We don’t completely discard thresholding from the first level of classification. We rather use it effectively as follows. A strict thresholding is used at the first level to finish the process for those clearly right or wrong. The thresholding algorithm should also ensure that the smaller target set consisting of possible solutions should contain the correct solution. Only then will the secondary classification stage work correctly.

For each sub-group of classes of numerals, a feature selection scheme is designed to select a subset of the features. Then secondary classification is done on this subset of features to choose the best answer from among the down sized target class set.

2.1. Classification in the First Level

Over past several decades, many kinds of feature extraction schema have been designed for handwritten character recognition. Among all the available feature extraction implementations, we chose the GSC feature set designed at CEDAR for experiments of our proposed method. GSC features are multi-scale and multi-resolution.

The GSC classifier has been well characterized by extensive testing. In a typical run of the GSC classifier using a single thresholding value, the 10E+R value reaches its minimal value 6.9749 with the acceptance rate of 96.66% and correct rate of 99.64%. In table 4 we present the minimal 10E+R values for each of the 10 classes along with their performance when each class is thresholded independently.

From the table we can observe that the average reliable recognition accuracy rate after thresholding is about 96%. The top 2 hit (regardless of any thresholding) ratio is 99.56% for the ten classes and minimal top 2 hit ratio is 99.31% if the ten classes are recognized separately. These results suggest that the improvement of performance in digit recognizers is possible by classifier combination. Secondary classification can be done among the top 2 characters returned since the top 2 hit ratio is very high.

2.2. Feature Selection for the Second Level Classification

To further classify among only two possible classes for the input image rejected by the first level classification, the simplest way would have been to run the first step classification routine with the original number of features for the reduced target set of two instead of ten classes. This method usually does not work since the features are designed for classification of ten classes. Features that contribute to recognition in case of one class may contribute to confusion in case of another class or a smaller set of classes. This is actually a reason why the confidence value was not high enough and the input was rejected by the first level classification.

For each pair of possible classes, we want to get rid of the redundant features that confuse the recognizer and retain only those features which reveal maximal separability between the two classes. We build recognizers as pairwise discriminators with each giving decision for a subclass of only two classes(2 numerals). They will take the inputs rejected by GSC and further classify between the first two choices returned. The very high hit ratio for the first two returns of GSC ensures the right choice of a pairwise discriminator.

Table 4. Statistics shows the top hit ratio and top 2 hit ratio both are high.

class	top 1 hit(%)	top 2 hit(%)	image count	10E+R	accept (%)	correct (%)
0	99.08	99.31	1299	9.2916	96.2595	99.4449
1	99.85	100	1311	1.5985	99.1660	99.9235
2	98.69	99.38	1302	4.6808	97.6870	99.7632
3	98.16	99.54	1301	4.3068	97.2868	99.8406
4	97.93	99.38	1306	5.0361	97.3582	99.7606
5	98.87	99.76	1240	4.3102	96.5267	99.9163
6	98.78	99.54	1309	2.7742	98.0077	99.9218
7	99.13	99.62	1039	10.8758	94.1849	99.4939
8	98.57	99.52	1256	5.9833	96.4940	99.7523
9	98.75	99.54	879	10.4487	93.1614	99.6390
0-9	98.77	99.56	12242	6.9749	96.6590	99.6366

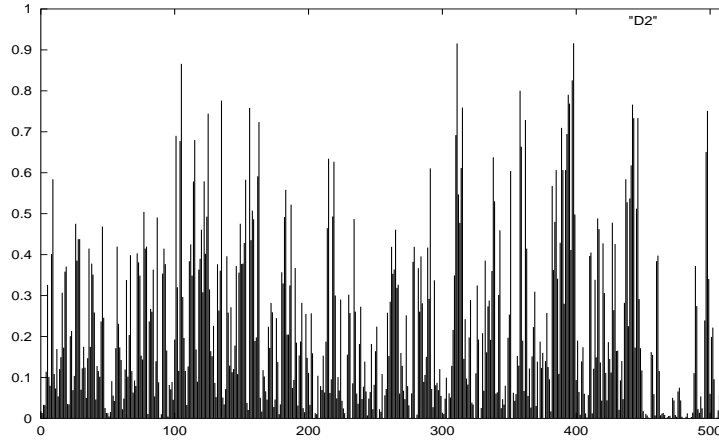


Figure 2. Feature distributions of numeral 2.

2.2.1. Distribution of GSC features

To build each pairwise discriminator, we select a subset of feature positions – key feature positions in the feature vectors of the two classes. The distributions of each of 512 binary features in the feature vectors are computed for each class. For example, Figure 2 is the distribution of features for numeral 2. For each pair of classes, the differences of the distributions at each feature positions are also computed, figure 3 gives the distributions of the differences for the features of numeral 2 and numeral 7.

The feature positions which show significant difference are selected to be key positions. For each pair, a threshold is computed and used on the distribution of difference. A feature position will be selected if the difference between the two class at the position is bigger than the threshold.

In some case, the distribution difference at a position is not bigger than the threshold, but the difference is significant comparing with the smaller distribution at the position. For example, the difference is bigger than a multiple of the smaller distribution. We also set a threshold for the multiple and select these positions based on the threshold.

The above discussion suggests that this problem of determining the key features can be rendered as an optimization problem in terms of selecting those positions in the feature vectors that present optimum separation between classes with respect to their distributions for those classes. Such a problem can be effectively handled with help of genetic algorithms. We encode the feature vectors as a bit string. For the problem of finding which feature set amongst the complete feature set provides optimal separation with different thresholds we perform a GA run with the feature vectors of the templates of these 2 classes as sample data. Each of the feature sets generated by GA is then

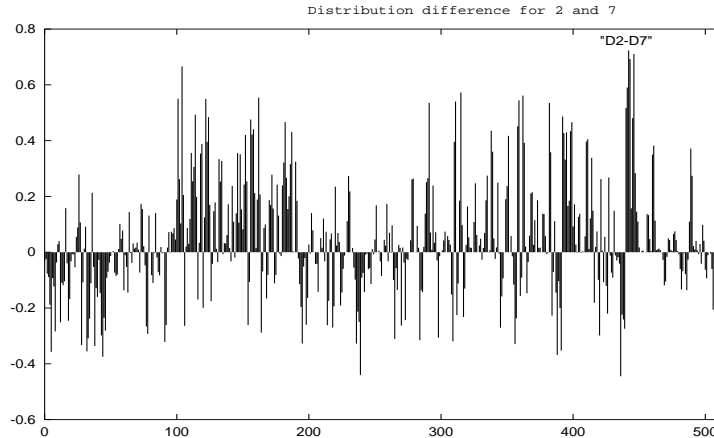


Figure 3. Differences in Feature distributions of numeral 2 and numeral 7.

evaluated to check if it expresses maximal separability. Maximal separability is quantified in terms of the threshold set previously.

The key feature positions are selected for all 100 pairs. For each pair of classes, a pattern vector is made with length of the original feature vector. The selected positions are denoted differently from the rest of positions.

2.2.2. Key feature Classification

A classifier using K-nearest neighbor algorithm combined with the subset feature pattern functionality is used for the secondary classification. Each input to this classifier is a feature vector rejected by the first level classifier. Since the first level classification behaves differently on different classes, a rejection scheme using different thresholds on different classes is used.

The design of the secondary classifier uses a set of neural networks each trained for a particular pair of digits. Such a neural network is combined with the K nearest neighbor classifiers. The prototype set for the K nearest neighbor is kept fixed and corresponding 2 class discriminators are trained.

3. EXPERIMENTS AND RESULTS

A combination of three classification routines using GSC features have been implemented. Each routine is based on K-nearest neighbor algorithm. The first two routines are full implementations of the K-nearest neighbor algorithm. The difference is that the first routine uses a selected set of prototypes which contains 1174 prototypes and the second routine uses a full set of 25656 prototypes. The smaller prototype set is selected to include only the hyper centers in the prototype space. About 80% of the input feature vectors can be classified to be accepted and are therefore decided by the first routine. We call the first routine *hyper center classifier*.

The rest of 20% feature vectors are presented to the second full K-nearest neighbor classification routine. From earlier statistical experiments, a threshold value is set for each class to determine the acceptance. This routine classifies another 16% of the input feature vectors.

After the first two classification routines, about 96% of the input images can be classified correctly. The small fraction of the total number of images rejected are sent to the second level classification routine – the secondary classification routine.

During the preprocessing stage of the experiments we identified the key feature sets used for pairwise discriminators at this secondary classification stage. This stage does not increase the time or cost of testing in any manner as key feature selection is done offline. Also the GA step is fairly easy to perform with any standard GA library. We used GAlib version 2.4.5 available at <http://lancet.mit.edu/ga/> for our experiments. GA routine is driven by the frequency distributions we obtained and presented in the previous section as a criterion function. Based on these

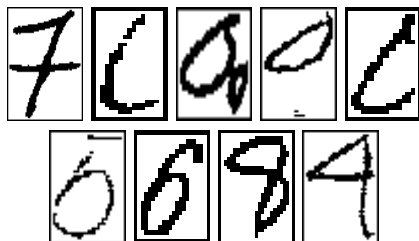


Figure 4. Images with truth not being in top 2 of the first level classification.

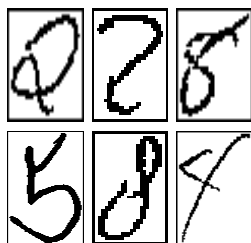


Figure 5. Images with error: rejected from the first level classifiers but incorrectly accepted by the secondary classifier.

and the input feature vectors of the templates the GA run produces feature subsets which are the key features for secondary classification of pairwise classifiers.

Several neural networks as pairwise classifiers were trained based on these key features. The output of each of these neural networks were the two classes on which the decision was to be made and a third class representing the fail to classify decision.

A set of 12242 handwritten digits images was used for our testing. After first two classification routines, 11848 were accepted among them and 50 were rejected completely. The rest of 394 feature vectors were sent to the secondary classifier. Using a preset threshold for the secondary classifier, 125 were accepted and among them 15 were misclassified. In these 15 incorrectly accepted images, 9 of them were not classified by the first level classifier into their top 2 choices.

For a view on such input images as those accepted incorrectly by the secondary classifier and which did not have the truth in the top 2 choices at the first stage of classification see figure 4 and figure 5. Based on these numbers the total number of images correctly recognized by the full classifier were 11908. The overall acceptance rate was increased from 96.78% to $(12242 - 394 + 125)/12242 = 97.80\%$ and the accuracy rate achieved was 99.45%.

This is a considerable improvement both in terms of accuracy and acceptance rates.

4. CONCLUSION

In this paper we proposed a multiple stage classification algorithm. Using the same set of features, the first level classification routine classifies simple feature vectors quickly and reliably. The left over feature vectors are difficult to classify and are hence presented to the second level classifier routine. The secondary classifier uses a sub-set of the original feature set called the key feature set. It helps determine the correct output among the set of possible targets which were the top 2 choices from the first stage. A small step was taken in employing genetic algorithms for feature selection in a handwritten character recognition problem domain. The testing experiments showed results with significant increase in performance with the high reliability preserved.

ACKNOWLEDGMENTS

We would like to thank Mr Xianzhong Chen for the excellent assistance.

REFERENCES

1. P. A. Devijver and J. Kittler. *Pattern Recognition: A statistical approach*. Prentice-Hall, London, 1982.
2. Paul D. Gader and Mohamed Ali Kjabou. Automatic feature generation for handwritten digit recognition. *IEEE transactions on pattern analysis and machine intelligence*, 18(12):1256–1261, December 1996.
3. U. Kressel and J. Schurmann. Pattern classification techniques based on function approximation. In H. Bunke and P. S. P. Wang, editors, *Handbook of Character Recognition and Document Image Analysis*, pages 49–78. World Scientific, 1997.
4. C.Y. Suen L. Xu, A. Krzyzak. Methods of combining multiple classifiers and their applications to handwriting recognition. *IEEE Systems, Man, and Cybernetics*, 22(3):418–435, 1992.
5. C. Y. Suen. Character recognition by computer and applications. In T. Y. Young and K.-S. Fu, editors, *Handbook of Pattern Recognition and Image Processing*, pages 569–589. Academic Press, 1986.
6. S.N. Srihari T.K. Ho, J.J. Hull. Decision combination in multiple classifier systems. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(1):66–75, 1994.