Active GP Classifiers for Handwritten Digit Recognition

Ankur M. Teredesai¹ and Venu Govindaraju²

- Department of Computer Science, Rochester Institute of Technology (RIT), Rochester, NY 14623 amt@cs.rit.edu
- ² 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 govind@cedar.buffalo.edu

Classification techniques strive to use features selectively to provide an optimal balance between classification accuracy and robustness. In traditional classification, all the classes in the n-class problem are equally important and the classification method tries to provide the best possible separation between all classes. Recent advances in classification termed *active* classification make selective decisions for specific subsets of classes within the multi-class problem. Typical examples of such scenarios are cases when there is a very high degree of confusion between two particular classes and a sharper solution is deemed useful. The design for active classifiers requires the underlying pattern recognition technique to blend feature discovery within the classifier training phase.

When features are derived from handwritten character images, different levels of informative detail can be present in different regions of the image. Certain features are obviously more distinguishing than others when it comes to separating various subclasses in a multi-class problem at least for this domain. We conjecture that this dual task of feature selection and classifier training can be combined effectively using the Genetic Programming (GP) paradigm as described in this chapter. This chapter includes a comprehensive analysis of these and similar issues. Comparative results on a variety of datasets in handwritten digit recognition using a variety of parameters are also reported.

1 Introduction

Pattern classification is a task of developing predictive models and can be viewed as a mathematical mapping from features to class values. The process of designing such mappings (also termed as classification models, classifiers or recognizers) is usually based on machine learning techniques. Traditional classifier development uses a wide range of techniques based on statistical pattern recognition and non-parametric methods like neural networks [7].

Traditional vs Active Classification

Let's consider the problem of developing a naïve digit classifier. The training dataset consists of labeled samples from ten classes i.e. digits 0 thru 9. In a traditional setting, the classifier is described by a model M such that:

$$M(X) \longmapsto \widehat{Y}$$
 (1)

where X is the input feature vector and \widehat{Y} is the corresponding estimated class-label that M suggests for X. We also note that Y is the true class-label. The number of input features X is fixed based on the feature extraction method. The classification method's objective is to find an optimal model \widehat{M} that separates the ten classes in a hyper-space formed by the feature vectors. The model \widehat{M} is trained to use all features and provide the best possible separability between all classes.

Let's now consider an application driven classification task posed as follows: Which of the 4 digits $\{1,4,7,9\}$ does the input match the best? Such a question is often posed by automated systems for zip-code recognition, check reading applications, etc. A traditional classifier will output the class and confidence values for the input without considering the limited set of class choices. What is called for here is to build a specific classifier which focuses on features useful in separating the numerals in the set instead of using a typical classifier. The classes in the set can vary dynamically. Such classifiers are termed 'active classifiers' because they focus on classification using a variable sub-set of features. As seen in Figure 1, a classifier can discriminate between digits '2' and '7' effectively by using an appropriate subset of features. The features corresponding to the lower portion of the image encode the differences of the two digit classes. Features corresponding to the top portion of the image (dark shaded) provide no useful information in this respect. In traditional handwritten character recognition based on 'passive' methods [19] [27] [10] the feature vector is fixed and all the features have a role in classification. One can design neural networks or nearest neighbor rules to work on a selective subset of the feature space but these methods are primarily based on heuristics [28] [13] [22]. The regions and corresponding set of features that provide maximal separability between classes are used in the design of active classifiers. The classifiers then use this information to provide an accurate decision. Formally, the classification model can now be described by:

$$M(X) = \bar{M}(X) \bigodot \bar{M}_{i,j}(x)$$
 (2)

where, $M_{i,j}$ is a sharper classifier utilized to separate specific subset of classes and $x \subset X$. Key feature discovery and dynamic selection of smaller feature

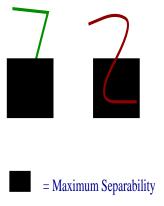


Fig. 1. The utility of discriminating regions in images can be observed. The shading highlights information rich zones which provide separability between classes. Lower portion is one such region in digits 2 and 7. GP based classifier identifies key regions of separability by using features corresponding to these regions as leaf-nodes in the GP-tree based classifiers.

subsets are important in the design of 'active' classifiers. Iterative gazing methods like those described by J. Park et.al. [21] can be used to convert the passive methods to active. The accuracy of classification then depends on the depth of iteration and the resources in terms of computation time available to make a decision. With the use of Genetic Programming (henceforth referred to as GP) [17] [18] [14] we try to eliminate the need for an iterative process and directly focus on the regions of interest. The GP solutions produced at the end of training use a small sub-set of the original feature set to provide a sharper decision in case of confusing classes.

Simplicity of solution, feature selection, search in program space, adaptability to dynamically changing classes, good resistance to noise and high accuracy are some of the advantages of the proposed paradigm [24]. In this paper we present the following: (i) issues in development of active classifiers using GP, (ii) automatic focus by GP techniques on image regions important from the viewpoint of the classification task and (iii) improved fitness functions for GP based classification.

The classification model of GP is a tree structure where the internal nodes are mathematical or logical operands and the terminal (leaf) nodes are features of the pattern. A detailed description of the GP technique can be found in [17]. The training process involves iteratively breeding solution trees and analyzing the suitability of each tree as a classifier. The suitability (fitness) of each tree is its ability to accurately classify the input image. A numeric value to quantify suitability (fitness) of the solution can be computed using a 'fitness function'.

In order to focus on regions of the image that provide high separability between digits, the GP search has to be guided using appropriate search mechanisms. There are various methods of generating new solution trees using genetic search operators like crossover, mutation and reproduction [3]. After the training phase, the classifier (solution tree) having the best performance is chosen for testing. Validation is done using a validation set of images drawn independent of the training and testing sets between training cycles to estimate the accuracy of potential classifiers in the GP-population and acts as a training termination criterion. GP being essentially a stochastic and nonparametric search method [17], it does not assume a given data distribution in its attempt to provide a globally optimal solution. The drawback is the tradeoff between training time and search for optimal solutions. We address several issues pertaining to development of GP-based handwritten digit classifiers in this paper and report results on the NIST dataset along with a comparative analysis with other handwritten digit classification methods described in the literature like Support Vector Machines, MLP-SVM combination, GSC and HOCR. The objective of this paper is to show that GP can provide performance comparable to other tradition pattern recognition techniques while using a smaller subset of features.

1.1 Paper Organization

Related research in the areas of handwritten character recognition and genetic programming is presented in section 2.1. Section 2.2 gives a generic framework for GP. Section 2.3 describes the active handwritten digit recognition approach using hierarchical multi-resolution feature space. Integrating pattern recognition with GP for the handwritten character recognition problem is discussed in section 3. Section 3.2 enumerates the new and improved fitness functions proposed. Effective modeling of GP for the digit recognition problem using different function sets, thresholds and breeding parameters is described in sections 3.3 and section 3.4. Experiments and analysis of results for each of the issues is presented in section 4.

2 Background

This section outlines related work in the areas of handwritten character recognition and GP.

2.1 Character Recognition

Recognition of handwritten characters by a computer has been a topic of extensive research for many years [23]. It plays a pivotal role in many applications such as postal addresses interpretation, bank checks, tax forms and

census forms reading. Due to the difficulty of the task, there is probably no one single method that can achieve high recognition and reliability rates for all applications. Thus in previous implementations, efforts have been directed towards sophisticated considerations like decision combination in multiple classifier systems [12] [26] and use of tolerant rough sets [15]. Recently we reported the use of genetic methods as a two-level classification model [11].

In the domain of handwritten character recognition there are two key strategies in current use. They can be broadly grouped as 'active' and 'passive' character recognition. Passive character recognizers use methods such as those described by Favata [8] and Suen [23]. These classical methods employ a One Model Fits all [22] paradigm. Some systems like adaptive resonance network learn as they perform [5] but lack feedback mechanisms to describe the learning strategy. Active character recognition as described by J.Park et.al. [21] [22] employs an active heuristic function that adaptively determines the length of the feature vector as well as the features themselves to classify an input pattern.

In terms of classifier modeling strategies for handwritten character recognition several machine learning techniques have been implemented. Multiple layer perceptron model has the number of neurons, the number of hidden layers, organization of the network, training parameters, etc. selected by an expert [7]. The architecture of the recognizer is expected to be more flexible and self organizing because learning by weight modification in a fixed network topology can limit the classifier to a locally-optimal solution. The complexity of the problem may not scale favorably with the complexity of the network selected. The classification standards might be met but there will always be a trade-off with learning time, network size, and other such issues. These issues and the trade-off in using GP have been discussed by Koza et.al. [18]. The solution representation of classifiers obtained from our approach are relatively simple compared to previously described methods. In terms of performance accuracy and generalization error, Support Vector Machines(SVM)[[25] have shown remarkable improvement in recent years. Ayat et.al [1] have reported a new kernel for digit recognition using SVM. A MLP-SVM hybrid recognizer was developed by Bellili [2]. We compare the performance of GP-based active classifiers with these and other active and passive methods in the experiments

Kuncheva [20] presents methods to make GA based classifiers less prone to over-training by including penalty terms in the fitness function accounting for the number of features used. It is difficult to quantify the quality of such fitness measures when features are hierarchical in nature. In a passive recognition scheme this type of penalty structure may provide good results where each feature is assumed independent. The relation between classification for data mining tasks and GP has been explored by Freitas [9]. Main emphasis is placed on the classification task but optimality of features sets is not explored. Bot [4] proposes a framework for automatic extraction of features with the aim of dimension reduction and improving the accuracy of a k-nearest neighbor

(k-nn) classifier. Method to change the dimensionality with respect to the number of classes is not presented.

The relation between pattern recognition and GP was recently reported by Kishore et.al [16]. Our work extends some of these methods and suggests improvements in terms of fitness function selection, data set organization and application to multi-class classification problems such as handwritten digit recognition. In the pattern recognition area our design of fitness functions is more appropriate than those currently in use [17]. Data set organization into class specific training sets as described in the paper is a new concept in the GP application domain. Many problems associated with stochastic learning concepts can be solved by gaining a better understanding of issues like penalized search functions, parameter tuning and data set reorganization.

2.2 Genetic Programming (GP)

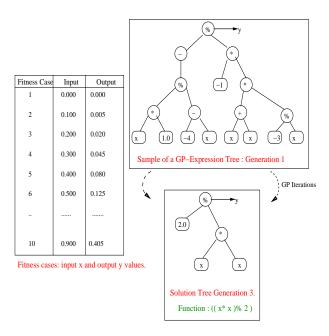


Fig. 2. Example showing modeling of a GP solution [3]. The inputs and expected outputs are given in the table and act as fitness cases. A sample tree randomly initialized in the first generation is shown. GP search for the best equation to represent the expected outputs (y) leads to the solution tree representing the equation $y=((x^*x)\%2)$ as shown in the second tree.

GP has been used over the past couple of years in several domains, primarily because of its ability to model problems effectively. The way trees are grown

and evolved can be observed by a simple example [3] given in Figure 2. GP consists of representing solutions as programs that can be evaluated. Consider a tree representation of any equation like (*(x+y)z). The internal nodes are the operators + and * and the leaf nodes are x, y, z. This is a typical solution tree in Genetic Programming. The problem at hand is easily represented using such expression trees. Inputs to the problem are represented as leaf nodes and are called terminals. Any operator used to operate on terminals is called a function. Functions and terminals together form a solution. Solutions are evaluated by instantiating the terminals and evaluating the functions. As solutions are represented with terminals and functions, the data types and length of these entities can vary from the simplest (e.g. integer) to more complex user defined types (e.g. Structures in C language). Thus solutions are logically referred to as programs and the search is based on genetic methods [17].

The search for the best program to solve a problem emulates the evolutionary process. An initial population of programs is randomly generated and evaluated. Best individual programs are selected for mating based on values obtained from fitness functions and new programs are created using genetic operators. Fitness functions essentially evaluate a program's ability to solve a problem. The new population is then evaluated and search for a better solution programs continues. This breeding process continues till an optimal solution is found or other user defined termination criterion are met.

2.3 Multi-Resolution Feature Space

Feature extraction is defined by Devijiver and Kittler [6] 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". Writing styles, writing equipment, resolution and image quality all play a major role in the variability of handwritten characters. To design a feature set for machine classification of handwritten characters one must consider features at different scales. However, one encounters difficulties when trying to achieve a high recognition rate and optimal resource utilization in tandem.

Hierarchical restructuring of the image domain using Quin and Quad trees produces multi-resolution features [21] [22]. Multi-resolution features accord a multi-level view of the image domain as shown in Figure 3. The figure shows how features extracted in this manner break up the image domain into a hierarchical feature space. This enables the classifier to iteratively focus on different regions of the image level by level. Features from different levels and regions combine to form information rich zones. For example, in Figure 3 the central region in multi-resolution space plus the start and end points of the digit image can combine to provide one information rich zone. The next stage in recognizer development is to implement a classifier such as a neural network or a K-NN rule. In order to optimize the resources in terms of training time

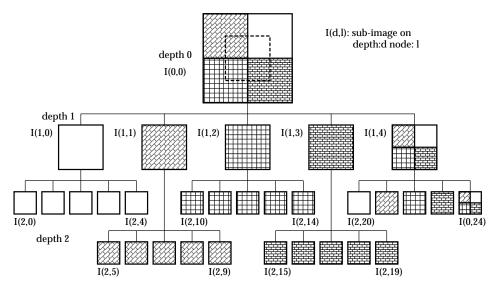


Fig. 3. Multi-resolution feature extraction using quin tree. The original image is divided into 5 parts at each level. The division is based on the gravity of pixels in the image. At depth 1 the image region I(1,0) is described by 9 features corresponding to the gradient and momentum of pixels in that region. The regions at greater depth provide greater detail. The level of resolution varies with depth.

and number of features used, it would be important to dynamically vary the solutions according to the scale and type of the inputs.

Previously active digit recognition using multi-resolution feature space made decision in an iterative matter by making decisions to use features at a greater resolution depending on the classification confidence at each stage. This made classification dependent on look-up tables and decision heuristics. Difference between iterative-active classification method and GP-based active classification method can be seen in Figure 4 and Figure 5. Each figure illustrates how a digit image is processed. In the iterative gazing procedure the decision to proceed to the next level is based on the confidence in classification at the current level. Compared to Figure 4, the nodes are not iteratively processed in Figure 5. Rather, the information about features which play an important role in classification is available directly at the end of the GP training as part of the classifier. In GP feature selection is a natural part of classifier evolution since solution trees are generated having a variety of feature combinations as terminal nodes. The best combination of features is automatically selected due to its recurring presence in good solutions over a number of generations. So far, we have discussed the motivation of developing a GP-based digit classifier. In the next section we outline the issues facing GP-based digit classification.

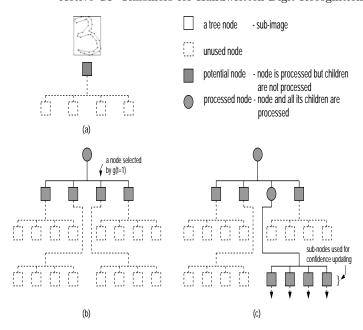


Fig. 4. Classification using an iterative classifier proceeds in levels. Based on the classifier confidence as a node is processed, the children are processed if required to improve the confidence and provide accurate classification. Classification process iterates over the multi-resolution space. Figure reproduced from J.Park and V.Govindaraju [21]

3 Active Classification Paradigm based on GP

In this section the strategy for implementing GP based classification is presented. This discussion is independent of the feature extraction methods and can be applied to both multi-level and single-level features. Once features are extracted the next task in classifier development is training. Classification of smaller number of digit classes may require fewer features. Subsets can be selected from the feature set to achieve high classification accuracy for a particular class. GP based binary classifiers (detectors) classify the input as belonging to either the class of interest C_i or belonging to the Φ class. The Φ class is defined as an aggregation of all classes except the class of interest C_i i.e. $\Phi = \{y - C_i | y = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}\}$. If the input feature vector belongs to class C_i the expected output is positive. Negative output corresponds to the class label Φ . A GP-based classification system is shown in Figure [6. A generic GP based training scheme has the following steps:

- 1. Decide the parameters for a GP run.
- 2. Feature vectors are computed from the training data and validation data and stored in flat files.

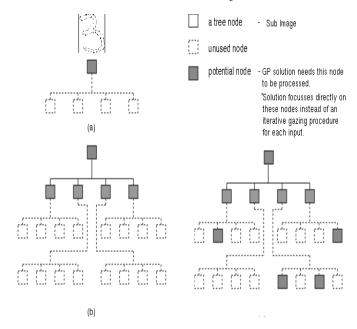


Fig. 5. GP based active classifier processes the root node level features. Makes a primary decision on the regions important for decision making and focuses directly on the regions of interest. No iterations are required since the classifier has been optimized for the level of detail required to make a decision during training.

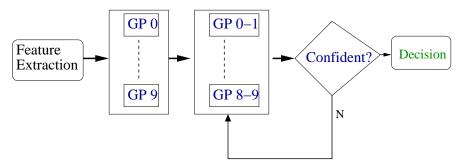


Fig. 6. Genetic Programming system for handwritten digit classification.

- 3. Since there are 10 classes in the digit datasets used in these experiments, we train 10 detectors. Each detector is trained to recognize one particular class of input.
- 4. Training for each detector proceeds as follows:
 - Each GP tree expression (acting as a detector) is evaluated against each feature vector corresponding to a sample in the training data.
 - The result of each evaluation is mapped to a real value that expresses a detection level for the sample. A detection level above +1 is positive and a level below -1 is negative. Any value in between is considered uncertain.
 - Fitness of each GP tree expression is computed based on the detection level for each input sample and the fitness function being used. Selection of fitness functions is discussed in section 3.2.
 - If a particular GP tree performs better than the best tree so far obtained on the training set, it undergoes validation using the validation set. Early termination criteria is then tested based on the performance of a particular detector on the validation set samples for that particular class.
 - If termination criteria is not met or maximum number of allowed generations is not over new trees are generated using the breeding parameters and search proceeds to the next generation.
- 5. The detectors are then combined in parallel to form an 11-class (0-9 and reject) classifier and are evaluated against the validation set once again. If the difference in accuracy of combined classifier on the validation set and individual detector performance is greater than 1% of accuracy for any class, the detector for that class undergoes further training.
- 6. As two classes might show high confusion at the combined classification stage, pair-wise classifiers for each pair of digits are also generated. These classifiers act as secondary classifiers.

Table 1. GP parameters that need to be tuned for the experiments.

| 0 Parameters | Values | | |
|--------------------------|---|--|--|
| 0 Objective: | Evolve function fitting classification of all digits 0-9 | | |
| Terminal Set: | Multi-Resolution Features | | |
| Function set: | +, -, *, DIV, SIN, COS, LOG, EXP and subsets of this set | | |
| 0 Population Sizes used: | 200, 400, 600 | | |
| Crossover probability: | 80 percent | | |
| Mutation probability: | 20 percent | | |
| Selection: | 'Fitness Proportionate' and 'Tournament Selection (size = 7 | | |
| Termination criterion: | Error rate less than 0.1% | | |
| Maximum Generations: | 2000 | | |
| Maximum Depth: | 17 | | |
| Initialization Method: | Half and Half | | |

Sample parameters used for development of GP classifiers are listed in Table 1. Modularization of the classification task by classifying two classes at a time brings up several issues such as data distribution, fitness evaluation, thresholding and use of breeding parameters. We address each of these issues subsequently.

3.1 Data Distribution

The digit distribution in training sets is an important issue. Since one classifier is supposed to handle only one class, care should be taken that this class is sufficiently represented in the training set. The training set used in this implementation needs to be biased. Around 35-50% of samples should belong to the class under consideration C_i with the rest of the classes making up the remaining. If a priori knowledge about confusion between classes is available, the training set should be accordingly distributed. For example, the digit 0 and digit 8 have a high confusion index, then while training for class 0 the training set should contain a sizeable number of samples from class 8, so the classifier can learn to separate between 0 and 8 effectively.

3.2 Fitness Functions

Fitness functions should reflect the overall goal of the GP run. Minimization of fitness over several iterations reflects generation of good solutions in the form of accurate classifiers. For GP based pattern recognition the fitness of individual solution trees should be proportional to the number of patterns classified correctly. Such fitness functions should levy suitable penalty on erroneous solutions and should reward good solutions. We have implemented several fitness assignment strategies (F_1, \dots, F_4) .

Fitness function F_1 is the generic function useful for developing traditional multi-class classifiers:

$$F_1 = N_b P_b + N_x P_x + P_c \tag{3}$$

where, N_b = number of training set samples on the boundary and thus rejected; N_x = number of training set samples belonging to class x incorrectly classified; P_b = penalty associated with boundary classification; P_x = penalty associated with incorrect class classification and P_c = penalty associated with non evolving solution criteria being met. Non-Evolving solution criteria occurs if the solution for multi-class classifier fails to classify at least one sample from each class under consideration.

For sharper binary classifiers the same fitness function takes the form:

$$F_2 = N_b P_b + N_{C_i} P_{C_i} + N_{\Phi} P_{\Phi} \tag{4}$$

where, N_{Φ} = number of training set samples belonging to class Φ incorrectly classified and P_{Φ} = penalty associated with incorrect Φ classification.

 F_3 is derived from F_2 but incorporates an adjustment to the raw fitness if the solution does not classify a certain minimum number of samples from the class C_i . This modification enables a GP run to climb out of a local-minima if the training set is not well balanced as discussed in the first issue.

$$F_3 = N_b P_b + N_{C_i} P_{C_i} + N_{\Phi} P_{\Phi} + P_i \tag{5}$$

where, J= the criterion deciding minimum number of samples of class C_i that must be targeted by every solution in order to maintain potency in the population; P_j = penalty associated with not meeting J. The value of J must be decided before we begin the run and should usually fall in the range of 1-5% of the number of samples belonging to class C_i in the training set.

F₄ incorporates adaptive parsimony pressure [29]. Controlling parsimony adaptively during evolution is a well founded method for improving the GP search. The motivation here is to emulate the way human visual system might try and distinguish images. Intuitively, one might try and learn higher level features that distinguish widely separable characters or numerals first, and later refine the classification abilities. Similarly a GP search can initially focus on finding solutions that provide maximal separability and subsequently refine the solution to achieve higher accuracy. The fitness function to implement such a strategy takes the form:

$$F_4 = F_3 + \alpha(g) * \Delta_i(g) \tag{6}$$

where,

$$\alpha(g) = \begin{cases} \frac{F3_{best}(g-1)}{N^2 * \Delta_{best}(g)} & \text{if } F3_{best}(g-1) > \epsilon \\ \frac{1}{N^2 * F3_{best}(g-1) * \Delta_{best}(g)} & \text{otherwise} \end{cases}$$

 $\Delta_i(g)$ = the complexity of the solution defined as the product of number of features used κ and the total depth of the solution tree λ ; N = number of training samples; g = current generation of the GP run and $F3_{best}(g)$ = Computed fitness F3 of the best individual in generation g.

The training process varies with each fitness strategy. Use of F3 results in more varied exploration of the solution space in our implementation. F4 generates more optimal solutions with fewer features expressed in the solution tree, but takes a longer time to converge. These fitness functions are more adaptive compared to the fitness proportionate heuristics based solely on the number of misclassified samples.

3.3 Function Sets

We mentioned previously that the terminal nodes of GP trees in our method are features extracted for each input image. Recall that mathematical operands are the function operators (internal nodes) in the GP tree. These functions can take terminals and output values of other functions as inputs. The choice of function sets defines the complexity of solutions produced. Our aim is to

develop simple solutions which can be interpreted easily to gather information regarding regions of the image important for classification. Assuming that suitable classification mechanisms applicable in handwritten digit recognition are non-linear expressions in the feature-space, we use function sets capable of producing non-linear tree classifiers.

The three operand sets used are:

- Arithmetic: +, -, DIV, *
- Trigo-Arithmetic: +, -, DIV, *, sin, cos
- Complex Arithmetic: +, -, DIV, *, sin, cos, exp, log

Function set consisting of +, -, *, DIV provided the best results in our experiments. Introduction of other mathematical operands like sin, exp, cos introduce fast convergence to local optima and lead to comparatively complex classifier trees. The arithmetic function set maintained the potential to evolve towards new solutions.

3.4 GP Parameters

GP based learning proceeds in populations. Breeding is a term used in GP for creation of new populations. Breeding parameters play an important role in controlling the population through a number of phases. Each phase has an operator and a rate suggesting how often that operator is used. Operators are similar to learning mechanisms like gradient descent. There are several breeding methods described in GP literature [17]. We tested two breeding strategies. The first one is based on selecting individuals using a tournament selection where a number of individuals are chosen at random with uniform probability with re-selection allowed. The best of the chosen individuals is selected to be in the new population. The second strategy has 4 breed phases with different crossover and mutation methods.

Initial attempts to converge prove unsuccessful when output values are +1 for the digit class x and -1 for the class Φ . Experimentation with values greater than equal to +1 for class x and less than equal to -1 for class Φ showed better results. The best separation between classes was accorded by values greater than or equal to +10 for class x and less than or equal to -10 for class Φ . Some input feature vectors were classified with output values between +10 to -10 and were placed in the reject class for this experiment. The width between class boundary seems to play an important role in classifier training suggesting that well separated boundaries play the same role as thresholds play in neural network classifiers.

4 EXPERIMENTS AND RESULTS

In this section we present an analysis of the various experiments conducted and results obtained. We have tested the GP based classification paradigm on various datasets. The NIST data set consists of 159,228 images in the training set and 53,300 images in the test set. The CEDAR-Digit data set consists of 18,000 images in the training set and 2,711 in the test set, obtained from handwritten zip-code images. Compared to the NIST dataset, the CEDAR-Digit data set is more noisy. Table[2 outlines the size of the training, validation and test sets used in our experiments.

Table 2. NIST and Cedar-Digit Datasets.

| Dataset | Training | Validation | Test |
|-------------|----------|------------|--------|
| NIST | 159,228 | 20,000 | 53,300 |
| CEDAR-Digit | 18,000 | 1,000 | 2,711 |

Table 3. Class-wise Results using GP Classification on NIST Test Data.

| Class | Images | Correct | Incorrect | Reject | Correct (%) |
|-------|--------|---------|-----------|--------|-------------|
| 0 | 5511 | 5373 | 138 | 20 | 97.5 |
| 1 | 5822 | 5665 | 156 | 30 | 97.3 |
| 2 | 5308 | 5143 | 165 | 12 | 96.9 |
| 3 | 5470 | 5311 | 159 | 17 | 97.1 |
| 4 | 5118 | 4949 | 169 | 15 | 96.7 |
| 5 | 4773 | 4639 | 134 | 12 | 97.2 |
| 6 | 5352 | 5207 | 145 | 22 | 97.3 |
| 7 | 5516 | 5317 | 199 | 28 | 96.4 |
| 8 | 5263 | 5047 | 216 | 11 | 95.9 |
| 9 | 5167 | 5043 | 124 | 34 | 97.6 |
| Total | 53300 | 51694 | 1605 | 201 | 97.1 |

Table 4. GP based classifier accuracy on NIST and Cedar-Digit Datasets.

| Dataset/Accuracy | Training | Validation | Test |
|------------------|----------|------------|------|
| NIST | 98.7 | 97.4 | 97.1 |
| CEDAR-Digit | 96.8 | 95.4 | 95.2 |

The results obtained by using GP classification are presented in Table 3. The column reject enumerates the number of images belonging to each class which were not accepted for classification by the classifier. We have included these numbers in the incorrect column i.e. the 1606 images incorrectly classifier include the 201 rejected images. The aggregated results on the two datasets described in Table 2 are presented in Table 4.

Examples of digit images correctly classified by GP based active classifiers are shown in Figure 7. These images are part of the noisy CEDAR-digit data set. Classifiers based on other techniques place these and similar images in the top 2 choice correct but fail to produce accurate decision. GP-based classifiers focus on the regions of maximal separability in each case and provide a correct decision thus increasing the acceptance rate. Sample classifier solutions for 2

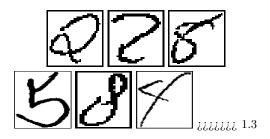


Fig. 7. Images with error: rejected by passive digit recognition methods but correctly accepted by the GP based active classifiers.

different digits are presented below: Some of the features used in the solution are F[30][1], F[8][5], F[16][0]. The reader should note the repetition of certain features over a number of generations in good (accurate) solutions.

These features can then be mapped back to the corresponding regions in the image domain to determine the information rich zones for a particular digit. In case of digit 3, it was found that the initial, middle and trailing strokes are such regions. 97.1% of the presented test patterns are recognized correctly using a subset of 100 (out of 9*31=279) features.

Solution for Digit 5

generation: 1192 nodes: 1145 depth: 7 hits: 4639/4773

accuracy:97.2
Part of the solution tree:

The classifier for digit 5 is also relatively simple in expression. An observation from this solution is that the solution space shows a lot of variation, suggesting that GP tries to focus on the key areas of separability and embeds them in individual solutions.

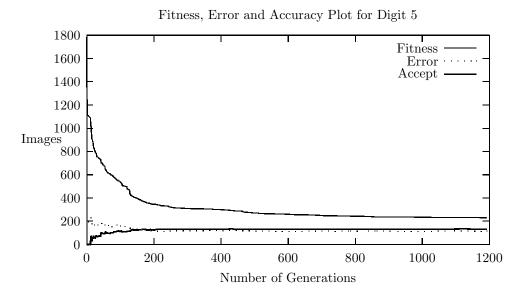


Fig. 8. Plot showing training error, fitness and accuracy over number of generations of a GP run for the digit 5 classifier. The plots shown are similar to other plots for other digits. The fitness decreases sharply during the initial runs and stabilizes over the number of generations with small decrements. Recall that lower fitness implies better solutions.

Consistent convergence towards a good solution can be observed from Figure 8. Over the number of generations in a given GP run this figure shows that fitness of individual solutions reduces as the accuracy of the classifier increases. The error rate decreases in proportion to fitness after the run has stabilized over initial generations. Fitness curves decrease over generations because fitness value is measured as a function of the number of test cases mis-classified (section 3.2). Hence lower fitness implies better classification. The spike in the error curve during the initial populations is suggestive of the amount of variability accorded by GP to discover accurate classifiers in initial generations.

Accuracy of the passive GSC recognizer [8] on the NIST test set 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) is 95.6%. In comparison GP based classifier has an accept rate of (53,300-201)/53,300=99.6% and accuracy of 97.1%. Hierarchical OCR which processes the input image using an active heuristic function called the gazing function uses a feature set similar to the GP based classifiers described here. The accuracy of HOCR[[21] on the NIST test set is reported at 98.2 which is higher than GP. GP outperforms HOCR in terms of the number of features used for classification and when pair-wise classifiers are required. The KMOD-kernel [1] based SVMs are comparable at 97.8%, and the best polynomial kernel with degree 4 provides 97.4% accuracy. Table 5 provides a comparative summary of accuracy and feature set sizes.

Table 5. Performance of Various Classification Methods Using NIST Data Set.

| Method | Correct | Feature Set Size | Decision Method |
|----------|---------|------------------|-----------------|
| GP-based | 97.1% | 100/279 | Active-Passive |
| GSC | 95.7% | 256 | Passive |
| HOCR | 98.2% | 279 | Active |
| KMOD-SVM | 97.8% | 277 | Passive |

4.1 Balancing datasets

We have discussed the need for organizing the datasets differently for the task of binary classifier development. A set of experiments was run in which the datasets were balanced for a class C_i . The balancing factor varied from 35 to 50% of the training samples. The plot in Figure 9 shows how data set re-organization plays a role in active classifier development. With normal distribution of samples in the training data, the error and accuracy curves stabilize in a local optimum quickly. After data re-organization, the two curves intersect and higher accuracy is achieved, with more variation in the solution space.

4.2 Fitness Functions

Several experiments were conducted to study the suitability of different fitness functions. During training, fitness functions play a key role in driving the search of GP classifiers. The rate of convergence and the quality of classifiers produced can vary widely. Fitness variation between fitness functions F2 and F3 (section 3.2) is shown in Figure 10. The y-axis represents $\frac{1}{Fitness}$.

GP runs converge quickly into a local optimum using fitness F2 which is similar to the defaults used in other GP implementations. Fitness function

Error with Normal Distribution Classified correctly Normal Distribution Error after Data Re-Organization Classified correctly after Data Re-organization · · · · Image\$00 No. of Generations

Error vs. Accuracy Plot for digit 3 to study Data Distribution Issues

Fig. 9. Plot showing error, fitness and accuracy over number of generations for digit 3 where data set is re-organized. Note the difference in the dotted curves suggesting the improved accuracy in classification. The convergence with normal distribution takes approximately 2000 generations but only 1200 generations are required with re-organized data distribution.

F3 and F4 perform better. Performance of F4 lies between F2 and F3 over most generations, since parsimony pressure is exerted to lower the solution complexity thereby compromising performance and resulting in longer times to converge to a good solution compared to F3. As can be observed from Figure [11] the number of features used by the best-in-the-generation solutions based on F3 is more than the features used by F4.

4.3 Parametric Variation

Each GP tree evaluates to a value. This value describes the phenotype of each solution. Analogous to pattern recognition techniques where the output of a classifier is a class label, the GP output value is an indication of the class assignment. There are various possible values a GP solution could evaluate to for a given problem. In our experiments we assign the value -1 and +1 to the two classes. A change to -10 and +10 shows a different error curve for the same dataset with all other parameters being the same. As seen in Figure 12 the general training curve is similar with both the thresholds with lower error rate for \pm 1 which is a tighter bound.

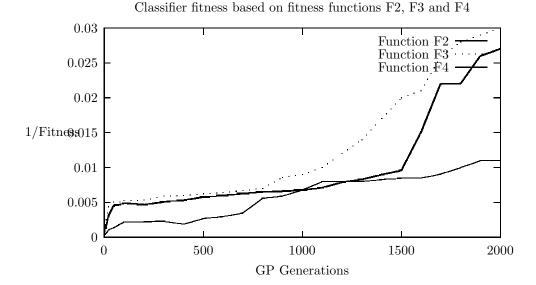


Fig. 10. Plot showing 1/fitness over number of generations for digit 1 for 2 different fitness functions. The increase in 1/fitness value describes better solutions. Function F3 results in better solutions compared to F2. The exploratory nature of F3 can be visualized from the fitness curve.

Figure 13 shows the number of images rejected along with the error. In this case the rejection rate stabilizes to a significant value with the threshold of ± 1 , but reduces if the threshold is well separated.

There were two sets of parameters used for experiments. Parameter Set 1 emulates the parameters described by Koza [17]. **Parameter Set 1:** breed phases = 2

breed[1].operator = crossover,keep trying=1, select=(tournament, size=7); breed[1].rate = 0.9

breed[2].operator = mutation,keep trying=1, select=(tournament, size=7); breed[2].rate = 0.1 The second set has function-point crossover 70% of the times and the simple any-point crossover 20% of the time. **Parameter Set** 2: probabilistic operators = off; breed phases = 4

 $\label{eq:breed} breed [1]. operator = crossover, select=fitness, internal=1.0; breed [1]. rate = 0.6 \\ breed [2]. operator = crossover, select=fitness, internal=0.0, external=0.0; breed [2]. rate = 0.2$

breed[3].operator = reproduction, select=fitness; breed[3].rate = 0.1

breed[4].operator = mutation, select=fitness, method=grow, depth=4; breed[4].rate = 0.1 Set 1 is useful when a quick convergence is required, specially for tasks like building pairwise classifiers. Set 2 is more exploratory in nature and is

Accuracy vs. Number of Features Utilized Number of Features Logarithmic Trend of Increase in Feature Utlization -Features Accuracy

Fig. 11. Accuracy increase in terms of increase in number of features used as terminal nodes in GP-trees is shown in the graph.

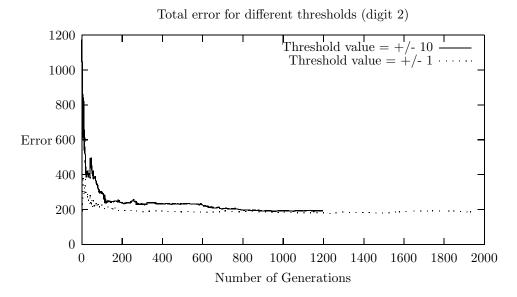


Fig. 12. Plot showing error variation using different thresholds for digit 2. In both cases the error stabilizes after 1200 generations. The variations in the curves suggest a difference in the way training proceeds due to different thresholds.

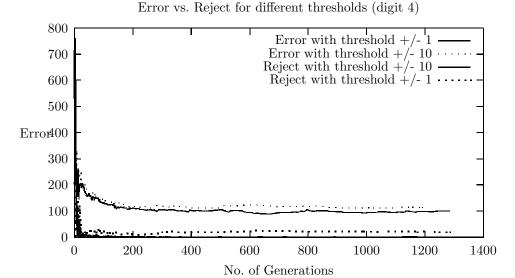


Fig. 13. Plot showing error variation using different thresholds for digit 4. The reject curve is plotted along with the error curve. The significant difference in reject curves for different thresholds supports the claim that expected values from the GP solution work as thresholding parameters in classification.

useful when we want to explore the solution space. The exploratory nature of the second set of parameters is seen in Figure 14.

5 Summary and Future Work

A new handwritten digit classification paradigm based on genetic programming techniques is presented. Classification accuracy of this technique is at par to the state-of-the-art methods. Classifiers that adapt dynamically to the scale of the classification task, are highly accurate, use a very small subset of features and conform to our intuitive understanding of the regions that provide maximum separability between classes can be developed using genetic programming. Different issues in implementation of this paradigm were presented in this paper.

The observation that in some cases of classification, the confusion introduced by certain regions of the image leads to a loss of classification accuracy is an important one. Traditional pattern recognition techniques start with a fixed feature set so it is difficult to dynamically adapt to the scale of the classification task. Future work on building better pair-wise classifiers is under consideration. The problems associated with noise in feature vectors and its

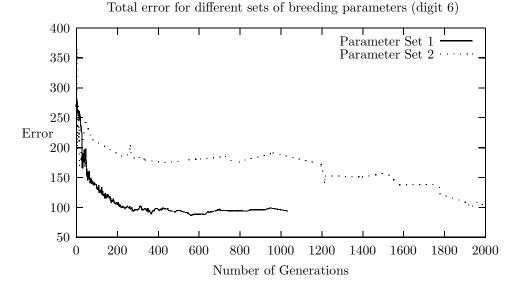


Fig. 14. Plot showing error variation using different breeding parameters. The parameter sets are the search mechanisms responsible for production of new improved solutions in GP based classification. The variation in error curves for two sets of parameters demonstrate the population variety which GP based paradigm provides.

implication on GP based training is a subject of further exploration. Since features are drawn from different levels correlations among them is often unavoidable. GP does try to select the optimal level but since no control structure is provided in generation of new GP-trees it is difficult to control feature selection. Future work in this direction will further our understanding of cognitively motivated pattern recognition approaches. It would be interesting to design genetic programming operators (functions sets) to compare different kind of search techniques. Our efforts are on to report the results of GP classifiers for "on-line" handwritten digit recognition which is another important application domain.

References

- N. Ayat, M. Cheriet, L. Remaki, and C. Suen. Kmod-a new support vector machine kernel with moderate decreasing for pattern recognition.application to digit image recognition. In K. Tombre, L. Spitz, and C. Fuh, editors, Proceedings of the 6th International Conference on Document Analysis and Recognition (ICDAR 02), pages 1215–1219, September 2001.
- 2. A. Bellili, M. Gilloux, and P. Gallinari. An hybrid mlp-svm handwritten digit recognizer. In K. Tombre, L. Spitz, and C. Fuh, editors, *Proceedings of the 6th*

- International Conference on Document Analysis and Recognition (ICDAR 02), pages 28–33, September 2001.
- W. Benzhaf, P. Nordin, R. Keller, and F. Francone. Introduction to Genetic Programming. Morgan Kaufmann Publishers, 1998.
- M.C.J. Bot. Feature extraction for the k-nearest neighbour classifier with genetic programming. In J.F. Miller, M. Tomassini, P.L. Lanzi, C. Ryan, A.G.B. Tettamanzi, and W.B. Langdon, editors, *Genetic Programming, Proceedings of EuroGP'2001*, volume 2038 of *LNCS*, pages 256–267, Lake Como, Italy, 18-20 April 2001. Springer-Verlag.
- G. Carpenter and S. Grossberg. A massively parallel architecture for a selforganizing neural pattern recognition machine. In Computer Vision, Graphics, and Image Processing, pages 35–115, 1987.
- P. Devijiver and J. Kittler. Pattern Recognition: A statistical approach. Prentice-Hall, London, 1982.
- R. Duda and P. Hart. Pattern Classification and Scene analysis. Wiley International, 1973.
- J. Favata and G. Srikantan. A multiple feature/resolution approach to handprinted digit and character recognition. *International Journal of Imageing Sys*tems and Technology, 7:304–311, 1996.
- A. Freitas. A genetic programming framework for two data mining tasks: Classification and generalized rule induction. In J. Koza, K. Deb, M. Dorigo, D. Fogel, M. Garzon, H. Iba, and R. Riolo, editors, Genetic Programming 1997: Proceedings of the Second Annual Conference, pages 96–101, Stanford University, CA, USA, 13-16 July 1997. Morgan Kaufmann.
- P. Gader and M. Kjabou. Automatic feature generation for handwritten digit recognition. *IEEE transactions on pattern analysis and machine intelligence*, 18(12):1256–1261, December 1996.
- 11. V. Govindaraju, Z. Shi, and A. Teredesai. Secondary classification using key features. In P. Kantor, D. Lopresti, and J. Zhou, editors, *Document Recognition and Retrieval VIII*, *Proceedings of the SPIE*, volume 4307, pages 272–279, 2001.
- T. Ho, J. Hull, and S. Srihari. Desicion combination in multiple classifier systems. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(1):66-75, 1994.
- A. Jain and B. Chandrasekaran. Dimensionality and sample size considerations in pattern recognition practice. In P. Krishnaiah and L. Kanal, editors, *Hand-book of statistics*, volume 2, pages 835–855. North-Holland Publishing Company, 1992.
- 14. K. Kinnear Jr. Advances in Genetic Programming. The MIT Press, 1994.
- 15. D. Kim and S. Bang. A handwritten numeral character classification using tolerant rough set. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(9):923–937, September 2000.
- J. Kishore, L. Patnaik, V. Mani, and V. Agrawal. Application of genetic programming for multicategory pattern classification. *IEEE Transactions on Evolutionary Computation*, 4(3):242–257, September 2000.
- J. Koza. Genetic Programming: On the Programming of Computers by Means of Natural Selection. MIT Press, 1992.
- J. Koza, F. Bennett, D. Andre, and M. Keane. Genetic Programming III. Morgan Kaufmann Publishers, 1999.

- 19. U. Kressel and J. Schurmann. Pattern classification techniques base on function approximation. In H. Bunke and P. Wang, editors, *Handbook of Character Recognition and Document Image Analysis*, pages 49–78. World Scientific, 1997.
- L.I. Kuncheva and L.C. Jain. Designing classifier fusion systems by genetic algorithms. IEEE Transactions on Evolutionary Computation, 4(4):327–336, November 2000.
- J. Park and V. Govindaraju. Active character recognition using a*-like algorithm. In Proceedings of Computer Vision and Pattern Recognition, 2000.
- 22. J. Park, V. Govindaraju, and S. Srihari. *OCR* in a hierarchical feature space. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(4):400–407, April 2000.
- C. Suen. Character recognition by computer and applications. In T.Y.Young and K.S. Fu, editors, *Handbook of Patter Recognition and Image Processing*, pages 569–589. Academic Press, 1986.
- A. Teredesai, J. Park, and V. Govindaraju. Active handwritten character recognition using genetic programming. In EuroGP2001 Proceedings, 2001.
- 25. V. Vapnik. Statistical learning theory, 1998.
- 26. L. Xu, A. Krzyzak, and C. Suen. Methods of combining multiple classifiers and their applications to handwritting recognition. *IEEE Systems, Man, and Cybernetics*, 22(3):418–435, 1992.
- 27. K. Yamamoto, H. Yamada, and T. Saito. Current state of recognition method for Japanese characters and database for research of handprinted character recognition. In S. Impedovo and J. Simon, editors, From pixels to features III: Frontiers in Handwriting Recognition, pages 250–260. Elsevier Science Publishers, 1992.
- J. Yang and V. Honavar. Feature subset selection using a genetic algorithm. IEEE Intelligent Systems, 13(2):44–49, 1998.
- 29. B.T. Zhang and H. Mühlenbein. Balancing accuracy and parsimony in genetic programming. *Evolutionary Computation*, 3(1):17–38, 1995.