



Contents lists available at ScienceDirect

## Swarm and Evolutionary Computation

journal homepage: [www.elsevier.com/locate/swevo](http://www.elsevier.com/locate/swevo)

## Multidimensional genetic programming for multiclass classification

William La Cava<sup>a,\*</sup>, Sara Silva<sup>b,c,d</sup>, Kourosh Danai<sup>e</sup>, Lee Spector<sup>f</sup>, Leonardo Vanneschi<sup>c</sup>, Jason H. Moore<sup>a</sup><sup>a</sup> Institute for Biomedical Informatics, University of Pennsylvania, Philadelphia, PA, USA<sup>b</sup> BioISI – Biosystems & Integrative Sciences Institute, Faculty of Sciences, University of Lisbon, Lisbon, Portugal<sup>c</sup> NOVA IMS, Universidade Nova de Lisboa, 1070-312, Lisbon, Portugal<sup>d</sup> CISUC, Department of Informatics Engineering, University of Coimbra, Portugal<sup>e</sup> Dept. of Mechanical and Industrial Engineering, University of Massachusetts, Amherst, USA<sup>f</sup> School of Cognitive Science, Hampshire College, Amherst, USA

## ARTICLE INFO

## Keywords:

Genetic programming  
Multiclass classification  
Feature extraction  
Feature selection  
Feature synthesis  
Dimensionality reduction

## ABSTRACT

We describe a new multiclass classification method that learns multidimensional feature transformations using genetic programming. This method optimizes models by first performing a transformation of the feature space into a new space of potentially different dimensionality, and then performing classification using a distance function in the transformed space. We analyze a novel program representation for using genetic programming to represent multidimensional features and compare it to other approaches. Similarly, we analyze the use of a distance metric for classification in comparison to simpler techniques more commonly used when applying genetic programming to multiclass classification. Finally, we compare this method to several state-of-the-art classification techniques across a broad set of problems and show that this technique achieves competitive test accuracies while also producing concise models. We also quantify the scalability of the method on problems of varying dimensionality, sample size, and difficulty. The results suggest the proposed method scales well to large feature spaces.

## 1. Introduction

Feature selection and feature construction play fundamental roles in the application of machine learning (ML) to classification. Feature selection makes it possible, for example, to reduce high-dimensional datasets to a manageable size, and to refine experimental designs through measurement selection in some domains. The ML community has become increasingly aware of the need for automated and flexible feature engineering methods to complement the large set of classification methodologies that are now widely available in open-source packages such as Weka and Scikit-Learn [10,30]. Typical classification pipelines treat feature selection and feature construction as pre-processing steps, in which the attributes in the dataset are selected according to some heuristic [9] and then projected into more complex feature spaces using e.g. kernel functions [29]. In both cases the feature pre-processing is often conducted in a trial-and-error way rather than being automated or intrinsic to the learning method. The use of non-linear feature expansions can also lead to classifiers that are black-box, making it difficult for

researchers to gain insight into the modelled process by studying the model itself. In this paper we investigate a multiclass classification strategy designed to integrate feature selection, construction and model intelligibility goals into a distance-based classifier to improve its ability to build accurate and simple classifiers.

A well known learning method that implicitly conducts feature selection and construction is genetic programming (GP) [17], which has been proposed for classification [7,15]. GP incorporates feature selection and construction by optimizing a population of programs constructed from a set of instructions that operate on the dataset features to produce a model. Compared to traditional ML approaches such as logistic regression and decision tree classification, GP makes fewer *a priori* assumptions about the data [22] and allows for various program representations [26]. In addition, GP has well-established methods for optimizing the intelligibility of models [37]. There have been some promising real-world applications of GP to binary classification [43], but recent work has focused on extending GP to the multi-class classification problem [14,28], in which there are more than two outcomes

\* Corresponding author.

E-mail address: [lacava@upenn.edu](mailto:lacava@upenn.edu) (W. La Cava).<https://doi.org/10.1016/j.swevo.2018.03.015>

Received 15 November 2017; Received in revised form 26 February 2018; Accepted 30 March 2018

Available online XXX

2210-6502/© 2018 Published by Elsevier B.V.

to estimate. This previous work suggests that traditional GP fares worse in comparison to other classification methods in the multiclass setting. However, two recent GP-based methods, M2GP [14] and M3GP [28] were shown to perform on par with several other ML strategies in recent studies.

The performance improvements observed by M2GP and M3GP stemmed the incorporation of a distance-based classification strategy into a multi-output GP system. We recently proposed a new method [20] called M4GP, that, although inspired by M2GP and M3GP, significantly improves these two methods. In this paper we extend M4GP by introducing an archiving strategy and by comparing it to recently published methods on data challenges from two different domains. The contributions of this work are:

- M4GP uses a novel (stack-based) program representation, that simplifies the construction of multidimensional solutions compared to M2GP and M3GP (which, instead, used a tree-based representation). This makes the evolutionary process of M4GP more efficient and the final solutions more expressive, readable and easy to understand.
- M4GP incorporates a multiobjective parent selection and survival technique that allows it to clearly and consistently outperform M2GP and M3GP on a wide set of test problems. To the best of our knowledge, this technique had never been used for multi-class classification before.
- We introduce an archiving strategy that maintains a set of optimal trade-off solutions based on complexity and accuracy. The final model is selected from this archive using an internal validation set to reduce overfitting.

Thanks to these improvements, M4GP is able to improve the best known GP methods for multi-class classification, and finds results that are competitive with the state-of-the-art methods for the studied problems (a set of 26 classification problems, ranging in numbers of classes, attributes and samples). Furthermore, on a set of biomedical data sets with up to 5000 attributes, M4GP is shown to perform on par with state-of-the-art methods while producing smaller models in less time. All these features foster M4GP as the new state-of-the-art multi-class classification method with GP.

The paper is organized as follows: Section 2 presents M4GP. In Section 3 we discuss previous and related work, focusing on the similarities and differences between M2GP, M3GP and M4GP. Section 4 describes our experimental study, presenting the used test problems and the experimental settings. In Section 5, we discuss the obtained experimental results. Finally, Section 6 concludes the paper.

## 2. M4GP

In multiclass classification (classification into more than two classes), we wish to find a mapping  $\hat{y}(\mathbf{x}) : \mathbb{R}^p \rightarrow \mathcal{C}$  that associates the vector of attributes  $\mathbf{x} \in \mathbb{R}^p$  with  $K > 2$  class labels from the set  $\mathcal{C} = \{1 \dots K\}$  using  $n$  paired examples from the training set  $\mathcal{T} = \{(\mathbf{x}_i, y_i), i = 1 \dots n\}$ .

One way to conduct classification is to measure the similarity of each attribute to the bulk properties of the attributes within each class, and then assign the label corresponding to the most similar group. This strategy is embodied by the ‘nearest centroid’ classifier [8], which operates as follows. The attributes are partitioned into subsets  $\{X_1 \dots X_K\}$ , such that  $X_k$  is the subset of  $\mathbf{x}$  with class label  $k$ . To classify a new sample  $\mathbf{x}' \in \mathbb{R}^p$ , the distance of  $\mathbf{x}'$  to each subset  $\{X_1 \dots X_K\}$  is measured and the class label corresponding to the minimum distance is assigned as

$$\hat{y}(\mathbf{x}') = j, \text{ where } j = \arg \min_k D(\mathbf{x}', X_k), k = 1, \dots, K \quad (1)$$

In its simplest form,  $D$  can be the Euclidean distance between  $\mathbf{x}'$  and the centroid of  $X_k$ . However, previous work has indicated the appropriateness of the Mahalanobis distance (as used in linear discriminant analysis), clearly showing and justifying its advantages compared to other

types of metrics, like the Euclidean distance [14]. The Mahalanobis distance is defined as follows:

$$D_M(\mathbf{x}', X_k) = \sqrt{(\mathbf{x}' - \mu_k) \Sigma_k^{-1} (\mathbf{x}' - \mu_k)^T} \quad (2)$$

where  $\mu_k \in \mathbb{R}^p$  is the centroid of  $X_k$  and  $\Sigma_k \in \mathbb{R}^{p \times p}$  the within-class covariance matrix. Eqn. (2) renders  $D_M$  the equivalent Euclidean distance of  $\mathbf{x}'$  from  $X_k$  after scaling by the eigenvalues (variances) and rotating by the eigenvectors of  $\Sigma_k$  to account for the correlation between columns of  $X_k$ .

The nearest centroid classifier assumes the within-class samples to be normally distributed about their centroid and separated from the distributions of other classes. For this reason, the data in the right plot of Fig. 1 is easily classifiable using nearest centroids. Unfortunately, many real-world problems have data that violates these assumptions, as in the left plot of Fig. 1. Furthermore, for high dimensional data, calculating Eq. (2) can be impractical. However, by finding a set of transformations that project the data into a space that delineates class membership, the performance of a nearest centroid classifier can be improved and high-dimensional distance comparisons may be avoided. Furthermore, the relationship between the raw data and the transformed space can be made legible by using symbolic transformations, as in  $\phi_1$  and  $\phi_2$  in Fig. 1.

The goal of M4GP is to find a set of transformations  $\Phi(\mathbf{x}) : \mathbb{R}^p \rightarrow \mathbb{R}^d$  that projects  $\mathbf{x}$  into a space in which the samples are more accurately classified by the nearest centroid method. Classification is then conducted with centroids  $\mu_{\phi_k} \in \mathbb{R}^d$ , covariance matrix  $\Sigma_{\phi_k} \in \mathbb{R}^{d \times d}$ , and distances  $D_M(\Phi(\mathbf{x}), \Phi(X_k))$ . We use GP to find or approximate the optimal synthesized features  $\Phi^* = [\phi_1 \dots \phi_d]$  that maximize the number of correctly classified training samples, as:

$$\Phi^*(\mathbf{x}) = \arg \max_{\Phi \in \mathbb{S}} f(\Phi, \mathcal{T}) \quad (3)$$

$$f(\Phi, \mathcal{T}) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(\hat{y}(\Phi(\mathbf{x}_i)) = y_i) \quad (4)$$

where  $\mathbb{S}$  is the space of possible transformations  $\Phi$ ,  $f$  is the classification accuracy, and the indicator function  $\mathbb{1} = 1$  if  $\hat{y}(\Phi(\mathbf{x}_i)) = y_i$ , and 0 otherwise. GP simultaneously optimizes the form of  $\Phi(\mathbf{x})$ , the subset of  $\mathbf{x}$  used in  $\Phi(\mathbf{x})$ , and the transformation dimensionality ( $|\Phi| = d$ ). Therefore, M4GP can produce low-order models for high dimensional data sets, making Eq. (2) tractable, or produce high-order transformations to improve classification accuracy for highly non-linear problems.

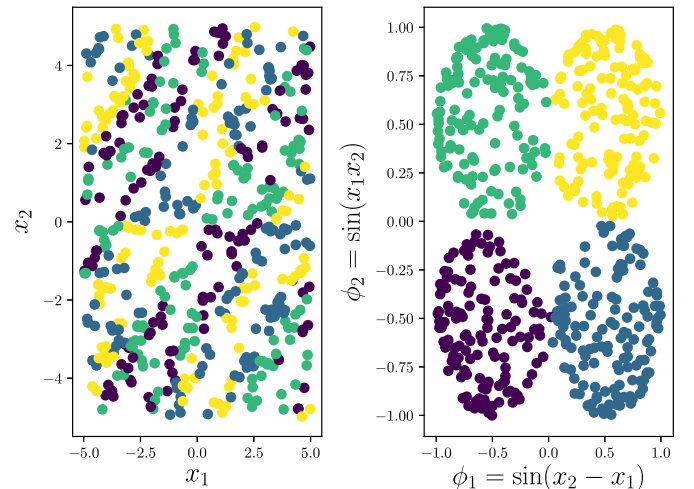


Fig. 1. Illustrative example of how data transformations can improve classification. On the left, the raw data for a synthetic 4-class problem for which there is significant overlap of data from different classes in the original space. On the right, the data is transformed into a new space in which the samples fall into neat clusters and class membership is easily discernible.

$$\mathbf{i} = \begin{bmatrix} x_1 & x_2 & x_1 & x_1 & * & x_2 & x_2 & * & x_1^2 & x_2 & * \\ \text{index: } & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \end{bmatrix}$$

program execution	stack
1. push ( $x_1$ ):	$[x_1]$
2. push ( $x_2$ ):	$[x_1 \ x_2]$
3. push ( $x_1$ ):	$[x_1 \ x_2 \ x_1]$
4. push ( $x_1$ ):	$[x_1 \ x_2 \ x_1 \ x_1]$
5. pull ( $x_1$ ), ( $x_1$ ); push ( $x_1 \cdot x_1$ )	$[x_1 \ x_2 \ x_1 x_1]$
6. push ( $x_2$ ):	$[x_1 \ x_2 \ x_1^2 \ x_2]$
7. push ( $x_2$ ):	$[x_1 \ x_2 \ x_1^2 \ x_2 \ x_2]$
8. pull ( $x_2$ ), ( $x_2$ ); push ( $x_2 \cdot x_2$ )	$[x_1 \ x_2 \ x_1^2 \ x_2^2]$
9. push ( $x_1$ ):	$[x_1 \ x_2 \ x_1^2 \ x_2^2 \ x_1]$
10. push ( $x_2$ ):	$[x_1 \ x_2 \ x_1^2 \ x_2^2 \ x_1 \ x_2]$
11. pull ( $x_1$ ), ( $x_2$ ); push ( $x_1 \cdot x_2$ )	$[x_1 \ x_2 \ x_1^2 \ x_2^2 \ x_1 x_2]$

$$\rightarrow \Phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

Fig. 2. Example of program representation of a multidimensional transformation. Arguments such as  $x_1$  are pushed to the stack, and operators such as ‘\*’ pull arguments from the stack and push the result.

### 2.1. Genetic programming

GP represents solutions as programs composed of basic building blocks and iteratively updates these programs based on their quality. In M4GP, each program consists of a set of equations that comprise the synthesized features  $\Phi(\mathbf{x})$  used to estimate  $\hat{y}$ . For example, a program  $\mathbf{i}$  might construct a polynomial expansion of two attributes as

$$\mathbf{i} \rightarrow \Phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2] \quad (5)$$

where  $\phi_1 = x_1$ ,  $\phi_2 = x_2$ ,  $\phi_5 = x_1 x_2$ , and  $|\Phi| = 5$ . Traditionally, a GP program consists of a single syntax tree with a single output generated at the root node [17]. For example,  $\phi_5$  above could be represented by a tree ( $(x_1) (x_2)$ ), where ‘\*’ is the root node and ( $x_1$ ) and ( $x_2$ ) are its leaves. In order to allow programs to represent multi-dimensional transformations, in M2GP and M3GP, program trees were modified with special nodes in order to allow for multiple outputs at the root [14,28]. This introduced unnecessary complexity to the representation. A contribution of this work is the introduction of a stack-based data flow to simplify the encoding of  $\Phi$ , presented in §2.1.1.

The GP population is optimized by probabilistically selecting programs based on their performance and recombining and mutating these programs to make steps in the search space. The selection mechanisms used in M4GP are described in §2.1.2.

#### 2.1.1. Representation

M4GP uses a stack-based representation [31] rather than the more traditional tree-based GP representations. Programs consist of post-fix notation equations, e.g.,  $\mathbf{i} = [x_1 \ x_2 \ +] \rightarrow \Phi = [x_1 + x_2]$ . This representation eliminates the need for specialized nodes used in M2GP and M3GP by supporting multiple outputs by default. Programs are evaluated by stack operations; for example the program in Eq. (5) can be constructed as

$$\mathbf{i} = [x_1 \ x_2 \ x_1 \ x_1 \ * \ x_2 \ x_2 \ * \ x_1 \ x_2 \ *]$$

The execution of program  $\mathbf{i}$  is illustrated in Fig. 2. In lieu of recursive tree evaluation, stack-based evaluation proceeds left to right, pushing and pulling instructions to and from a single stack. Arguments such as  $x_1$  are pushed to the stack, and operators such as ‘\*’ pull arguments from the stack and push the result. At the end of a program’s execution, the stack state is interpreted as the multi-dimensional transformation.

#### 2.1.2. Initialization, selection, and variation

M4GP initializes a population of programs of various sizes and dimensionality. Each equation in a program is initialized recursively in an analogous fashion to the grow method (see Ref. [32]) but limited by number of nodes rather than depth. Eq. (4) is used as program fitness for selection.

We test three population selection methods in this work: tournament selection [32], lexicase selection [11,39], and age-fitness Pareto selection [34]. Tournament selection proceeds by selecting a set of individuals (in this case, two) in the current population and choosing the one with better fitness as a parent for the next generation. We describe the other two methods below in more detail.

**Lexicase selection.** Lexicase selection is a parent selection technique that rewards individuals in the population for performing well on unique combinations of fitness cases, i.e. samples. Each parent selection begins with the entire population in the selection pool. The fitness cases are shuffled randomly and the first case is considered. Any programs in the pool that do not have *exactly* the best fitness on the first case are removed from the selection pool. If more than one individual remains in the pool, the removal process is repeated with the next case. This process continues with additional fitness cases until either 1) only one individual remains in the selection pool, in which case that individual is selected as a parent or 2) no more fitness cases are left, in which case a parent is selected from the remaining pool. This process is repeated for each parent selection.

In lexicase selection, test cases can be thought of as filters, and each parent selection represents a random path through these filters. The selected parents are Pareto-optimal with respect to the fitness cases [21]. Each fitness case has a filtering capacity that is directly proportional to its difficulty since it removes individuals from selection that perform worse than others on it. Selective pressure thereby shifts to cases that are difficult to solve. Because each parent selection uses a new ordering of fitness cases and the cases filter the population in proportion to their difficulty, parent selection favors individuals that perform well on unique combinations of test cases. Lexicase selection therefore promotes individuals with diverse performance observed in previous experiments [11,21].

**Age-fitness Pareto survival.** Each individual in Age-fitness Pareto survival is assigned an age equal to the number of generations since its oldest ancestor appeared. Each generation, one new individual is added to the population to provide a small amount of random restart. Selection for variation is random rather than guided, and during breeding a number of offspring equal to the overall population size is created. Afterwards, the set  $P$  consisting of the current population and the newly created individuals is culled to the original population size,  $N$ , using their accuracy and their age. In this case, a lower age is preferred, and the fitness is inverted so that a lower fitness is also better. Program survival is based on the environmental selection scheme of the Strength Pareto Evolutionary Algorithm 2 (SPEA2) [44]. This algorithm uses two measures to perform this reduction: 1) Pareto strength of an individual,  $S(\mathbf{i})$ , which is the number of individuals equal to or dominated by  $\mathbf{i}$  ( $i < j$ ),<sup>1</sup> divided by  $P + 1$ , and 2) a density estimate  $D(\mathbf{i}) < 1$ , based on the inverse of the distance to the  $k$ -th nearest neighbor [36] of  $\mathbf{i}$  in objective space (in this case the objectives are normalized between zero and one). These metrics are used to define a fitness value  $F(\mathbf{i})$  that combines the total strength of the individuals  $\mathbf{j} \in P$  that dominate  $\mathbf{i}$  with density estimate  $D(\mathbf{i})$ :

$$F(\mathbf{i}) = \sum_{\mathbf{j} \in P, \mathbf{j} < \mathbf{i}} S(\mathbf{j}) + D(\mathbf{i}) \quad (6)$$

Every nondominated solution is first copied to the new population. If the new population size is smaller than  $N$ , individuals are added in

<sup>1</sup> Individual  $\mathbf{i}_1$  dominates  $\mathbf{i}_2$ ; i.e.,  $(\mathbf{i}_1 < \mathbf{i}_2)$  if  $f_j(\mathbf{i}_1) \leq f_j(\mathbf{i}_2) \forall j$  and  $f_j(\mathbf{i}_1) < f_j(\mathbf{i}_2)$  for at least one  $j$ . Here a lower  $f_j$  is better.

order of lowest (i.e., best)  $F(i)$ . If the population is larger than  $N$ , signifying that there are more than  $N$  nondominated solutions, individuals are removed iteratively based on lowest  $D(i)$ , such that individuals in less dense locations in objective space are preferred. For the latter scenario, the use of  $D(i)$  for selection helps preserve spread of solutions along the Pareto front.

**Variation.** We use the mutation and crossover operators from M3GP [28] in order to simplify the points of comparison with M4GP. These search operators are biased to explore the dimensionality of  $\Phi$ . The search operators manipulate “sub-trees” of programs, which are equivalent to segments of interacting nodes in the stack-based representation. The mutation operator, with equal probability, chooses one of three actions: i) it replaces a sub-tree with a randomly generated sub-tree; ii) it adds a new randomly generated sub-tree to the end of the program, thereby increasing  $|\Phi|$  by one; iii) it deletes a sub-tree corresponding to a “root”, thereby reducing  $|\Phi|$  by one. The crossover operator similarly chooses between one of two equally probable actions: i) it performs standard sub-tree crossover of the parents, selecting non-“root” nodes; ii) it performs standard sub-tree crossover of “root” nodes. In this case, “roots” are those nodes in the program that produce a value in the final stack, and can be identified from stack-based programs in linear time.

### 2.1.3. Alternate classification strategies

In M4GP, Eqn. (1) is used to make classifications, in contrast to more traditional approaches in GP (see Ref. [7] for a review). More commonly, GP is used to evolve discriminant functions (mathematical expressions) that are combined with user-defined thresholds to classify samples. For example, for a binary classification problem, a mathematical expression is evolved to produce positive values for one class and negative values for another class - in other words the threshold is set to zero. When extending to multiclass problems, either the problem can be split into several one-versus-rest problems, as done with SVM, or several thresholds can be specified such that instances of are assigned class 0 for program outputs less than  $-1$ , class 1 for outputs between  $-1$  and  $1$ , and class 2 for outputs greater than  $1$ . Another approach is to evolve classification rules using boolean operators (e.g. AND, OR, NOT) and interpret the output as class label. The rule-based approach also requires the problem to be decomposed into binary classification problems to apply to multiclass problems.

We implement two simple approaches to classification using the stack-based, multi-output system described in §2.1 so that we can compare M4GP to more common GP implementations. The first, referred to as ‘bool’, includes boolean operators {AND, OR, NOT, <, >, ≤, ≥, =, IF-THEN, IF-THEN-ELSE} to evolve a set of classification rules, and interprets the output as a bitstring for classification. For a binary classification problem, the output of the boolean stack is interpreted as class 1 for output [0] and class 2 for output [1]. For a four class problem, the output of the boolean stack is interpreted as follows: [0, 0]:class 1; [0, 1]: class 2; [1, 0]: class 3; [1, 1]: class 4.

The second approach which we call ‘float’ evolves a set of mathematical expressions in equivalent fashion to the GP system in M4GP. However, classifications are made based on the index of the largest value in the stack output. For example, if the program in Fig. 2 produced the output [1.5, 3, 2.25, 9, 4.5], class 4 would be assigned to the evaluated sample.

In our experiments, we compare the centroid-based classification strategy used in M4GP to the ‘bool’ and ‘float’ strategies described above.

### 2.1.4. Archiving solutions

Typically in GP the individual with the highest training accuracy (the so-called ‘best-of-run’ individual) is chosen as the final model and evaluated on the test set. In order to compare to these previous meth-

ods, we do the same with M4GP on the benchmark problems. However, given the importance of interpretability, on subsequent problems we maintain an archive of solutions that represent the best trade-offs of training accuracy and complexity (i.e. the Pareto set [44]). To choose a final model, a small validation set is split from the training set before the run. At the end of the run, the archive is evaluated on the validation set, and the individual with the best score is chosen as the final model. This approach should guard against over-fitting and produce simpler models. The entire archive can also be studied by experts to get a better sense of the building blocks of good solutions to their problem, which can offer insight to the user.

## 3. Related work

GP has been used extensively for evolving classification functions  $\hat{y}(x)$  directly [7,15,26]. In application to multiple classes, the discriminant functions evolved by GP must be thresholded, or the problem must be split into several binary classification problems [7]. To overcome the need for arbitrary thresholds in multiclass problems, M2GP proposed a multi-output GP that evolved  $\Phi(x)$  and used the nearest centroid approach (Eq. (1)) [14]. M2GP demonstrated in particular that Mahalanobis distance outperformed Euclidean distance in this framework. M3GP extended M2GP to allow programs to change dimensionality during the run via specialized search operators that increased or decreased the dimensionality of a tree by modifying its root node [28]. An ensemble version of M3GP named eM3GP was also proposed that performed similarly to M3GP with smaller, more legible resultant programs [35]. In contrast to these methods, M4GP removes the need for explicit root nodes by using a stack-based data flow that also preserves multi-dimensionality and allows dimensionality to change flexibly. It also incorporates multi-objective selection strategies and an archive to promote concise solutions.

Our preliminary work with M4GP [20] explored the role of the centroid classifier and tested M4GP’s multi-dimensional GP strategy on a set of biomedical problems. A subset of the results suggested M4GP may be well-suited finding nonlinear interactions between genes in genome wide association studies, which motivated our detailed investigation in this paper. Here we expand upon our initial work in the following ways: 1) we rigorously compare M4GP to previous GP and ML methods on a set of benchmark problems; 2) we introduce an archiving strategy for model selection and interpretability; and 3) we conduct a detailed application study of M4GP to non-linear genetics problems with comparisons to recent literature.

In a broader sense, these GP-based methods highlight the unique challenge of feature construction and its role in learning systems [24]. A few recently developed ML methods have similarly merged GP with linear regression [1,13,27]. M4GP and its ancestors are more tailored to the multi-class case, which for these regression-based approaches requires thresholding of real-valued outputs. The stack-based approach used by M4GP is also novel compared to the GP representation used in these works. Another active area of research is the merger of GP with decision tree learners [18], with promising initial results.

Feature construction is considered an important aspect of image classification, and GP has been used in several methods in this domain. For example, GP can be used to learn image embeddings for ensemble methods [25], as an interactive learning tool for remote sensing [6], and for detection pulmonary nodes in medical imaging [5]. Liu et al. [25] also noted the GP’s potential as a dimensionality reduction technique for large-scale problems. M4GP differs from these approaches in two ways: first, it focuses on the capacity for low- and high-dimensionality feature extraction to flexibly suit the needs of the problem, and second, it applies to general multiclass classification problems. We show in application to a high-dimensional genetics dataset in Section 5.3 that M4GP can be useful for learning low-order representations.



**Table 1**

Data set properties used for benchmark comparisons and Opportunity Activity Recognition Challenge.

	Benchmark Comparisons								Activity Recognition	
	Heart	IM-3	IM-10	Mov1	Seg	Vowel	Wav	Yeast	Opp-S2	Opp-S3
Classes	2	3	10	15	7	11	3	10	4	4
Attributes	13	6	6	90	19	13	40	8	242	242
Samples	270	322	6798	360	2310	990	5000	7797	16667	15550

**Table 2**

M4GP settings. Changes to the settings for Opp-S2 and Opp-S3 are noted in parentheses.

Setting	Value
Population size	500 (1000)
Max Generations	100
Crossover/Mutation	50/50%
Ephemeral random constants [17] range	[0,1]
Program size limits by # nodes	[3, 100 (500)]
Initial dimensionality range ( $d$ )	[1,33]
Termination criterion	generations or perfect training accuracy
Trials	30 (10)

#### 4. Experimental analysis

The experimental analysis of M4GP is divided into three sections. First we conduct benchmark comparisons, comparing M4GP to alternative GP strategies and to results from related GP literature. In the subsequent two sections, we benchmark M4GP against published results from two different studies, one concerned with human activity recognition [3,33], and the other with disease prediction from genome-wide association studies [38]. Performance is quantified in a number of ways: 1) by classification accuracy on the test sets, 2) by run-time, and 3) by model interpretability according to dimensionality of the resultant classifiers as well their forms. The settings for M4GP are shown in Table 2, and whenever possible, match those used for M2GP, M3GP, and eM3GP. Programs are constrained to be between 3 and 100 nodes, and are initialized with a corresponding dimensionality between 1 and 33, constrained such that the minimum sub-program is at least 1 node.

##### 4.1. Benchmark comparisons

The 8 problems used for benchmark comparisons are shown in Table 1. Six of these problems are from the UCI data repository [23], and the two others, IM-3 and IM-10, are satellite data sets from a United States Geological Survey [42].

Our first experiment is to test whether incorporating centroid-based classification into GP improves its performance. To this end, we compare the ‘centroid’, ‘float’, and ‘bool’ configurations defined in §2.1.3 on the benchmark problems. Next, three versions of M4GP are tested: M4GP with lexibase selection (M4GP-lx), age-fitness Pareto survival (M4GP-ps), and tournament selection (M4GP-tn). On the first 8 problems, we benchmark M4GP against M2GP, M3GP, eM3GP, and several out-of-the-box classifiers from Weka [10]: random forests (RF), random subspace (RS), multi-layer perceptrons (MLP), and support vector machines (SVM). Each method is run for 30 trials, and for each trial the data is randomly partitioned into 70% training and 30% testing.

##### 4.2. Opportunity Activity Recognition Challenge

We compared M4GP’s performance to the results of the Activity Recognition Challenge at the 2011 IEEE International Conference on Systems, Man, and Cybernetics [3,33]. We used the Opportunity Activity Recognition data set (Opp-S2 and Opp-S3 in Table 1), which consists of 242 attributes and approximately 32,000 total samples. We bench-

mark M4GP’s ability to predict four classes of locomotion (stand, sit, walk, lie) from two test subjects (S2 and S3) against the entrants to the original challenge. In order to perform the comparison to published results, the weighted F-measure,  $F_1$ , is used as the fitness metric in place of Eq. (4).  $F_1$  measures classification performance as a function of precision ( $\frac{TP}{TP+FP}$ ) and recall ( $\frac{TP}{TP+FN}$ ) as:

$$F_1 = \sum_{\ell}^k 2w_{\ell} \frac{\text{precision}_{\ell} \cdot \text{recall}_{\ell}}{\text{precision}_{\ell} + \text{recall}_{\ell}} \quad (7)$$

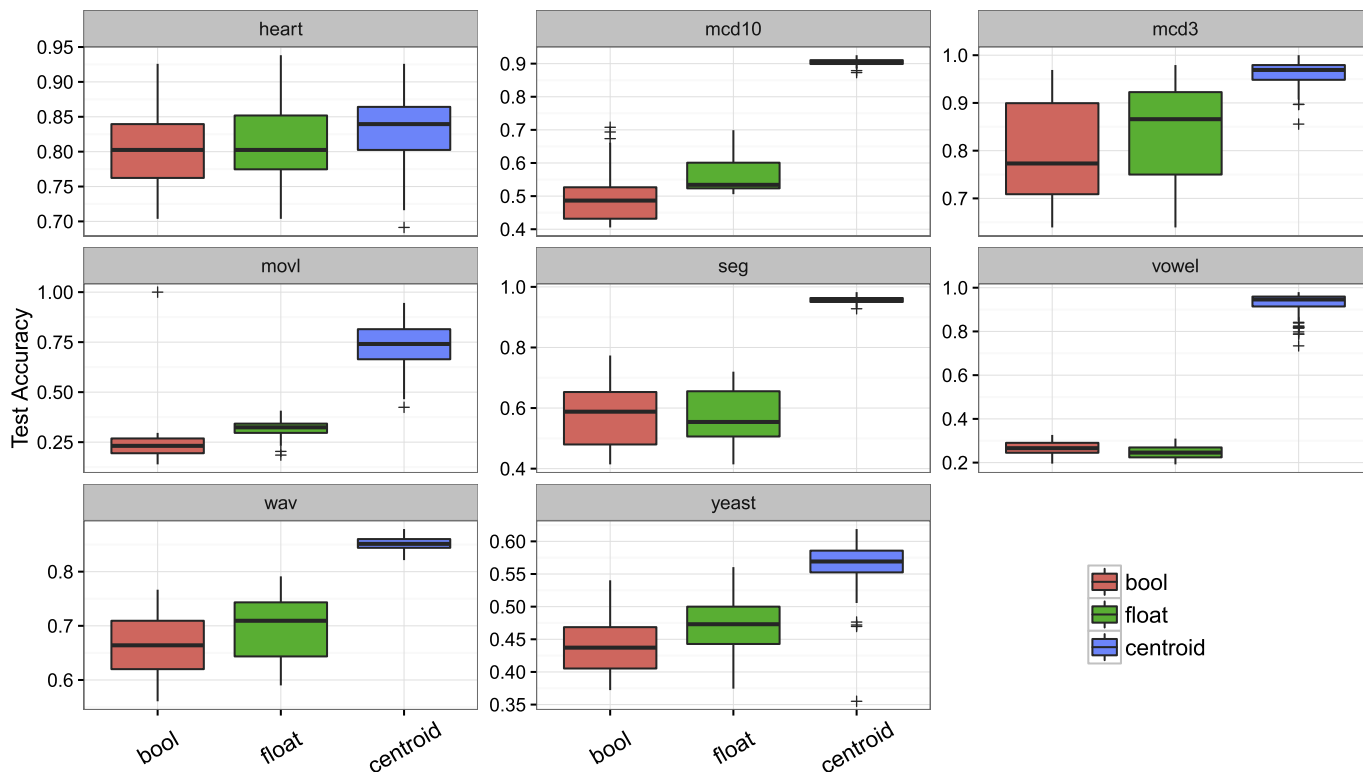
$TP$  is the number of true positives,  $FP$  is the number of false positives, and  $FN$  is the number of false negatives for class  $c_{\ell}$  that has  $k_{\ell}$  samples out of a total of  $n$ , yielding the proportion weight  $w_{\ell} = k_{\ell}/n$ . The raw time series data was downsampled as in Ref. [3] using a moving average with a window of 500 ms with 250 ms steps. Missing data is linearly interpolated, and all data was normalized to zero mean and unit variance.

##### 4.3. Biomedical application

Our preliminary work suggested that M4GP may be able to outperform traditional classification strategies in identifying nonlinear interactions in simulated genome-wide association studies (GWAS) [20]. Here, we compare M4GP to the state-of-the-art tools used for identifying epistatic interactions from noisy data sets. Recently, an adaptation of the tree-based pipeline optimization tool called TPOT-MDR [38] was shown to perform well on these types of problems in comparison to other methods. TPOT uses GP to optimize machine learning pipelines, searching over the space of Scikit-learn machine learning tools [30]. In Ref. [38], Sohn et al. coupled TPOT with multi-factor dimensionality reduction (MDR) and expert knowledge filter algorithms (EKF) to achieve good results. TPOT-MDR+EKF outperformed XGBoost [4], a state-of-the-art gradient-boosting tree algorithm, on these GWAS problems. Here, we compare M4GP to the results from that study to see if a lighter-weight algorithm (M4GP) with a similar ability to capture non-linearity in the original feature space, can perform as well as TPOT-MDR+EKF.

In addition to previously published methods, we compare to a deep neural network (DNN) approach implemented in PyTorch.<sup>2</sup> We construct a feed-forward neural network with 5 hidden layers of 10 nodes

<sup>2</sup> <http://pytorch.org/>.



**Fig. 3.** Accuracy of different GP classification methodologies on the test sets of the 8 benchmark problems. ‘bool’, ‘float’ and ‘centroid’ refer to different ways to interpreting the stack outputs of the GP programs, as described in §2.1.3. The ‘centroid’ method is the one used in M4GP.

**Table 3**

Comparison of best-of-run median test accuracy for the benchmark problems. The best result is highlighted in bold. Significant ( $p < 0.01$  according to a pairwise Wilcoxon rank-sum test with Holm correction) improvements with respect to each method is denoted by  $a - j$  according to the method labels.

Method	Heart	IM-3	IM-10	Movl	Seg	Vowel	Wav	Yeast
<sup>a</sup> RF	<sup>d</sup> 80.2	94.8	<sup>bcd</sup> <sub>fghij</sub> <b>96.9</b>	<sup>bcd</sup> <sub>fghij</sub> <b>89.4</b>	<sup>bcd</sup> <sub>fghij</sub> <b>97.3</b>	<sup>bcd</sup> <sub>g</sub> 89.4	81.5	<sup>d</sup> 57.5
<sup>b</sup> RS	<sup>d</sup> 81.5	92.8	<sup>cde</sup> <sub>fghij</sub> 93.9	<sup>d</sup> <sub>f</sub> 65.7	<sup>d</sup> <sub>g</sub> 96.0	<sub>g</sub> 82.8	82.2	<sup>d</sup> 56.6
<sup>c</sup> MLP	<sup>d</sup> 80.2	95.9	90.2	<sup>bd</sup> <sub>fghij</sub> 82.5	<sup>d</sup> <sub>g</sub> 96.3	<sub>g</sub> 82.5	<sup>ab</sup> <sub>g</sub> 83.3	<sup>d</sup> 58.0
<sup>d</sup> SVM	55.6	93.8	90.4	14.4	55.8	81.8	<sup>abc</sup> <sub>fgh</sub> 86.3	41.1
<sup>e</sup> M2GP	<sup>d</sup> 80.2	93.8	90.2	<sup>bcd</sup> <sub>fghij</sub> 85.9	<sup>d</sup> <sub>g</sub> 95.6	<sup>cd</sup> <sub>g</sub> 85.9	<sup>abc</sup> <sub>g</sub> 84.9	<sup>d</sup> 53.8
<sup>f</sup> M3GP	<sup>d</sup> 79.0	95.4	<sup>c</sup> <sub>ij</sub> 91.0	<sup>d</sup> 57.1	<sup>d</sup> 95.6	<sup>abcde</sup> <sub>g</sub> 93.8	<sup>ab</sup> <sub>g</sub> 84.3	<sup>d</sup> 56.2
<sup>g</sup> eM3GP	<sup>d</sup> 80.9	93.3	<sub>j</sub> 90.3	<sup>bd</sup> <sub>fh</sub> 78.6	<sup>d</sup> 94.7	78.6	81.2	<sup>d</sup> 56.2
<sup>h</sup> M4GP-lx	<sup>abcde</sup> <sub>fg</sub> 85.2	<sup>abcde</sup> <sub>fg</sub> <b>97.9</b>	<sub>ij</sub> 90.7	<sup>bd</sup> <sub>f</sub> 73.1	<sup>bde</sup> <sub>fgh</sub> 96.6	<sup>abcde</sup> <sub>fg</sub> 95.6	<sup>abc</sup> <sub>fg</sub> 85.3	<sup>d</sup> <sub>g</sub> <b>58.9</b>
<sup>i</sup> M4GP-ps	<sup>abcde</sup> <sub>fghj</sub> <b>90.1</b>	<sup>abcde</sup> <sub>fg</sub> 97.9	89.8	<sup>bd</sup> <sub>fh</sub> 80.1	<sup>d</sup> <sub>g</sub> 96.1	<sup>abcde</sup> <sub>fghj</sub> <b>97.5</b>	<sup>abc</sup> <sub>fghj</sub> <b>87.1</b>	<sup>d</sup> <sub>fg</sub> 58.9
<sup>j</sup> M4GP-tn	<sup>abcde</sup> <sub>fg</sub> 87.7	<sup>abcde</sup> <sub>fg</sub> 97.9	89.6	<sup>bd</sup> <sub>fh</sub> 76.9	<sup>d</sup> 95.1	<sup>abcde</sup> <sub>fg</sub> 96.0	<sup>abc</sup> <sub>fg</sub> 86.0	<sup>d</sup> 56.8

each. The DNN is trained via back propagation using stochastic gradient descent with a learning rate of 0.1 and a maximum of 1000 epochs. The non-linear activation functions in the hidden layers are defined as rectified linear units, and at the output, a softmax function is used. Dropout is used at the input layer for regularization [12].

We follow the experimental design in Sohn et al. by running 30 replication trials of 10-fold cross validation with balanced accuracy as the metric of performance. The simulated GWAS datasets used in the comparison are generated using GAMETES, which is a tool for embedding epistatic gene-gene interactions into noisy genetic datasets [41].

16 total problems are compared that vary according to two measures of difficulty: number of attributes (10, 100, 1000, 5000) and signal-to-noise ratio (0.05, 0.1, 0.2, 0.4), also known as heritability in the genetics community. For each problem, a two-way epistatic interaction is present that is predictive of the disease classification, and its presence is masked by the presence of confounding attributes and noisy disease classifications. We use multi-factor dimensionality reduction (MDR) in our analysis, which is a method developed specifically for detecting these interactions via exhaustive search of pairwise comparisons of attributes. Here, the comparison MDR-Pred is given the predictive fea-

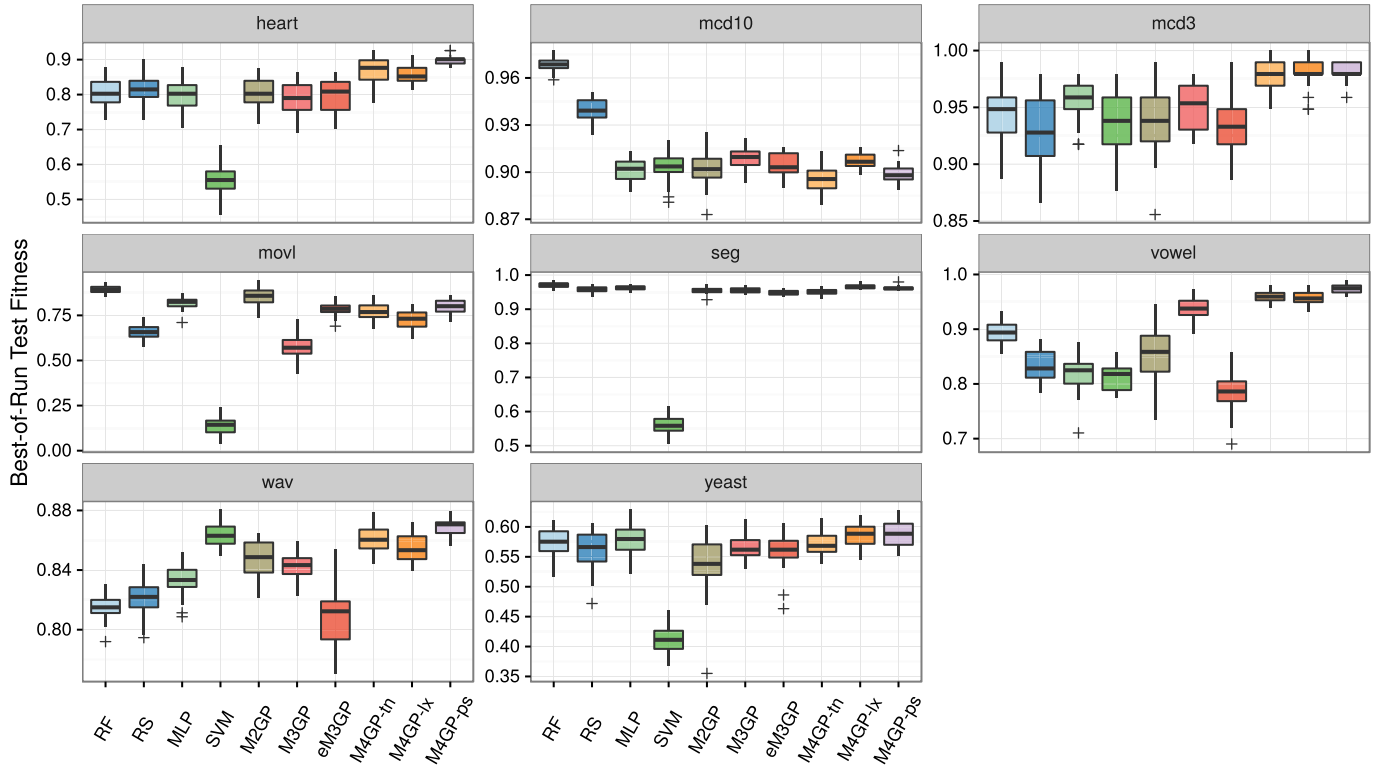


Fig. 4. Test set accuracy on the first eight benchmark problems for ten different classification methods.

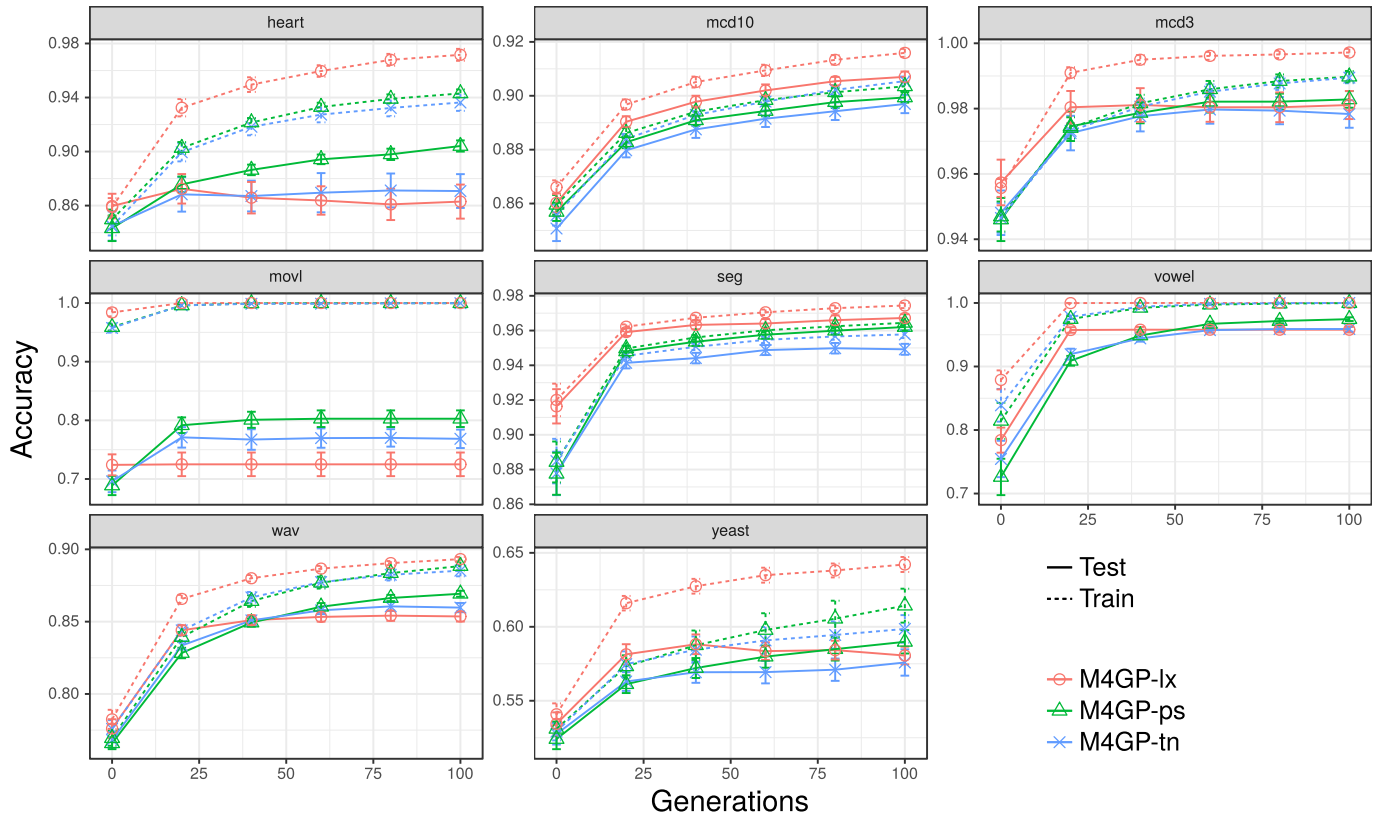


Fig. 5. Mean convergence characteristics of M4GP on the training set (dotted lines) and test set (solid lines) over 100 generations. Shapes indicate different selection methods. Error bars denote the confidence intervals.

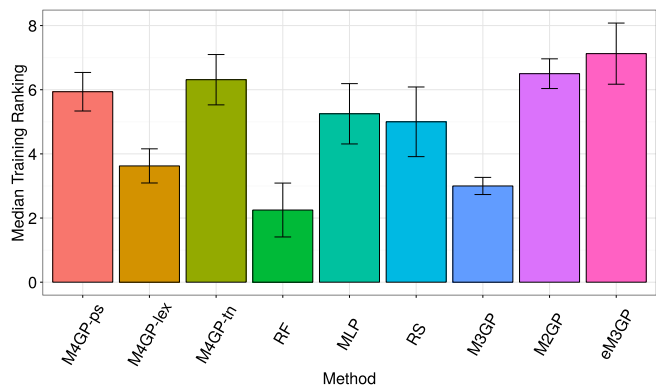


Fig. 6. Median training rankings of classifiers on the eight benchmark problems. Bars denote the 95% confidence interval.

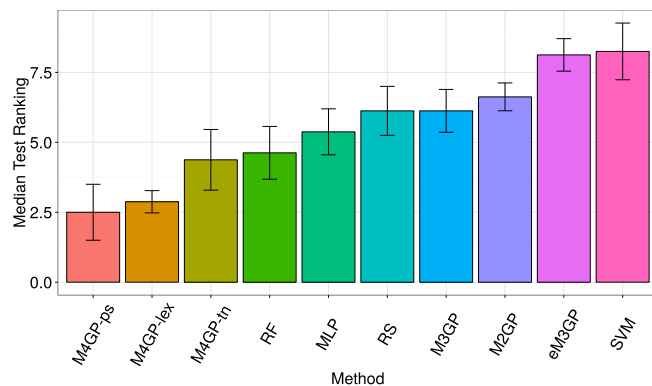


Fig. 7. Median test rankings of classifiers on the eight benchmark problems. Bars denote the 95% confidence interval.

tures beforehand, which is only possible in a simulated study such as this one. MDR-Pred provides a yardstick to indicate near-optimal performance on each dataset.

The EKF algorithms introduced in TPOT-MDR+EKF are based on the Relief family of feature selection algorithms [16]. They are pre-processing filters that are sensitive to interactions and are used in TPOT-MDR+EKF to reduce the attributes to a manageable size. We include a version of M4GP in our comparisons called M4GP+EKF that uses ReliefF [16] to perform feature selection before M4GP is run, reducing the dataset to the top 10 features found by the algorithm. This gives us a way to tease out the effect of feature selection versus the performance of the search algorithms themselves.

## 5. Results

The best classifiers generated by M4GP for each trial are compared first to benchmark methods in §5.1 and then to other published results in §5.2 and §5.3.

### 5.1. Benchmark comparisons

We first compare the centroid-based classification strategy to the ‘bool’ and ‘float’ strategies for classification, as described in §2.1.3. The test accuracies for the 8 benchmark problems are shown in Fig. 3. Across problems, the nearest centroid method for classification outperforms the other two strategies for interpreting stack outputs.

For the 8 benchmark problems, the median best fitness on the test sets for the first eight problems are shown in Table 3 with statistical comparisons. M4GP, across selection methods, produces the best classifiers in terms of test accuracy on five of the eight problems, and RF produces the best classifiers on the remaining three. The best-of-run classifier accuracies over all trials for all methods on the test sets are plotted in Fig. 4. Interestingly, the M4GP results are not typically best on the training sets but tend to generalize well. This generalization capacity is demonstrated by side-by-side comparisons of method rankings on training and test sets in Figs. 6 and 7, respectively.

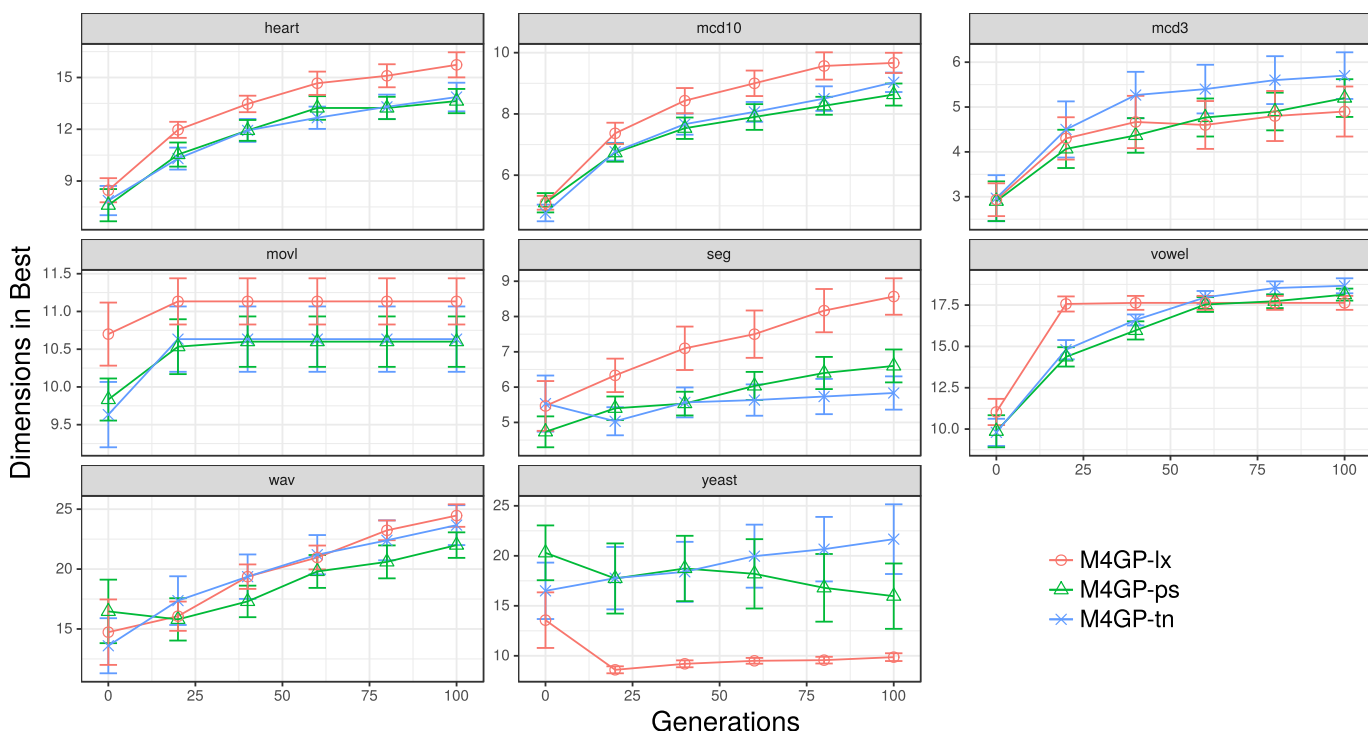


Fig. 8. Mean dimensionality of most accurate classifier (on training data) each generation. Error bars indicate the confidence interval.



**Table 4**

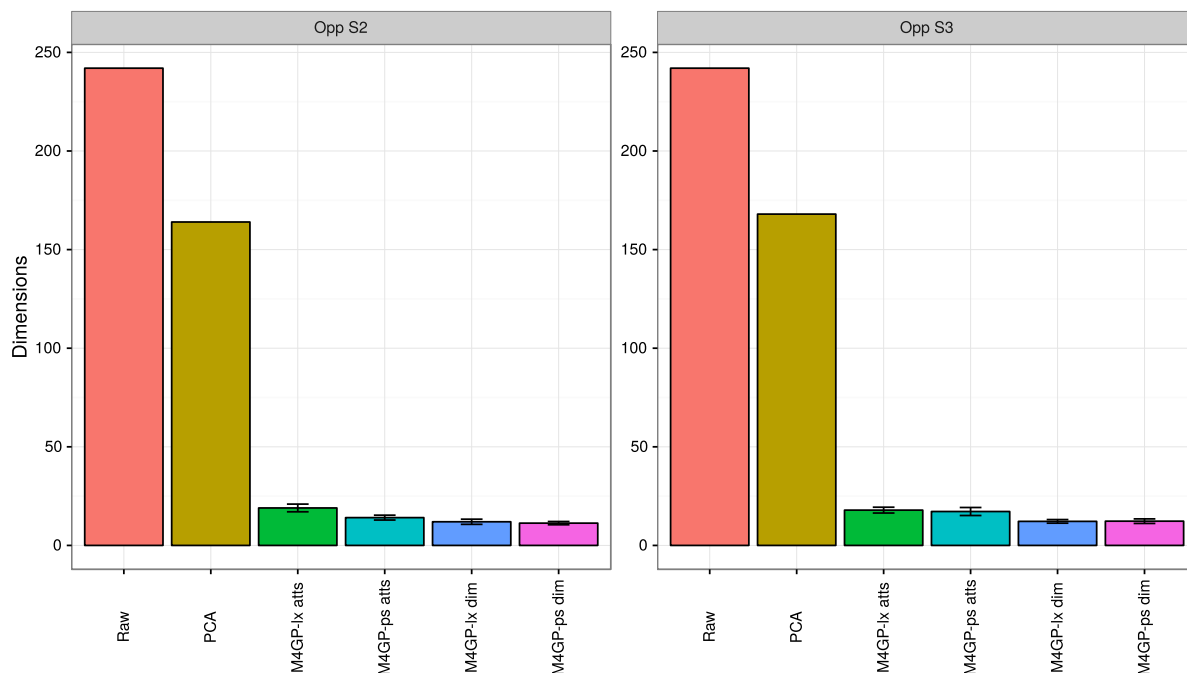
Comparison of F-measure on the Opportunity Activity Recognition data set (locomotion) for subjects 1 and 2 (S1 and S2). M4GP is compared to published results using one nearest neighbor (1-NN), SVM, SVM+1-NN, decision trees (C4.5), k-NN, decision tree (DT) grafting, and Adaboost. The team names are shown in parentheses. The best results in terms of F-measure are highlighted in bold. For M4GP results, the median best result of ten trials is shown with the 95% confidence interval.

Method	F-measure (no Null class)	
	S2	S3
1-NN (NStar)	0.88	0.85
SVM (SStar)	0.87	0.83
SVM+1-NN (CStar)	0.90	0.83
C4.5 (NU)	0.83	0.63
k-NN (MI)	0.87	0.86
DT grafting (MU)	0.86	0.87
Adaboost (UT)	0.74	0.72
M4GP-lx	<b>0.91 ± 0.00</b>	<b>0.89 ± 0.00</b>
M4GP-ps	0.90 ± 0.00	0.88 ± 0.00

The performance of the three M4GP methods on training and test sets over the duration of the run (generations) is shown in Fig. 5. It is worth noting that for the Movl problem, which is one of the three problems for which RF generalizes better, M4GP finds a perfect solution to the training data within the first few generations, and prematurely terminates. Past work also notes that this problem is an outlier in terms of GP's behavior [14,28]. On the Heart problem, M4GP-lx overfits to the training data, as evidenced by the slow decline in test fitness over the generations. Among selection methods, M4GP-ps ranks the best, followed by M4GP-lx and M4GP-tn. On individual problems, M4GP-tn did not outperform M4GP-ps or M4GP-lx; therefore, M4GP-tn was left out of further tests.

### 5.1.1. Dimensionality of solutions

Growth in the dimensionality of the solutions is a concern, and so we quantify how the best solutions change during training in Fig. 8.



**Fig. 9.** Dimensionality reductions afforded by M4GP on the Opportunity Activity Recognition Challenge. M4GP solution sizes are compared to the original number of attributes and a PCA dimensionality reduction preserving 98% variance. Here, 'atts' refers to the number of attributes used in the M4GP solutions, and 'dim' refers to the dimensions in the solution (i.e.  $|\Phi|$ ).

**Table 5**

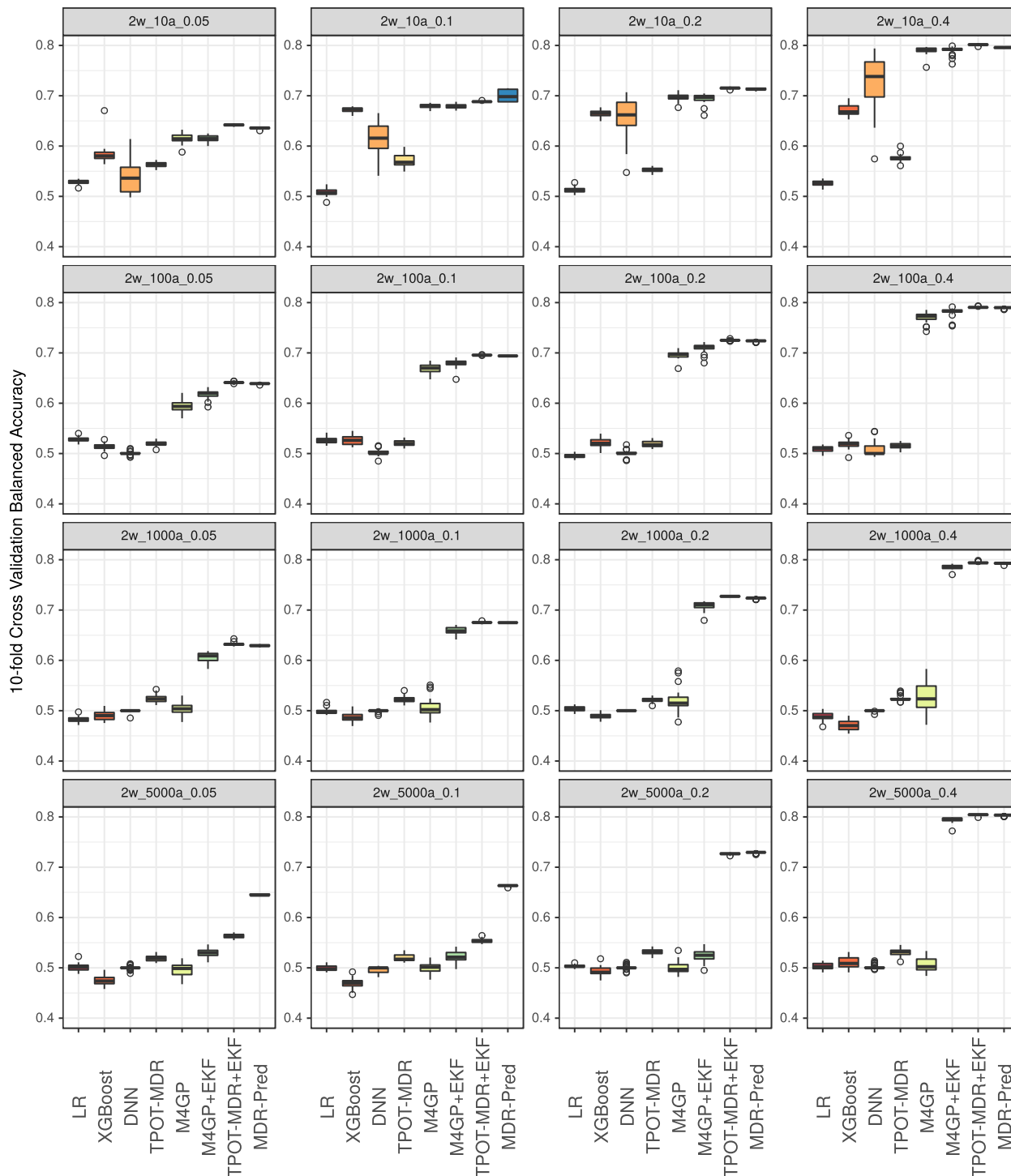
Run time for M4GP solutions.

Problem	Cores	Median Time (hr:min:s)		
		M4GP-lx	M4GP-ps	M4GP-tn
Heart	1	00:02:29	00:02:12	00:02:28
IM-3	1	00:02:30	00:02:49	00:02:14
IM-10	1	02:09:04	02:21:45	02:28:23
Movl	1	00:00:27	00:02:25	00:02:02
Seg	1	00:32:25	00:39:35	00:38:38
Vowel	1	00:04:21	00:20:37	00:15:57
Wav	1	01:36:44	01:30:17	01:35:12
Yeast	1	00:25:25	00:44:09	00:52:31
Opp S2	16	00:12:30	00:09:16	– (–)
Opp S3	16	00:11:51	00:08:47	– (–)

Differences between the selection methods are negligible except for on the Yeast problem, where lexicase selection produces smaller solutions. In general dimensionality growth is fairly flat aside for Wav, for which solutions tend to grow throughout training. The Pareto archive (see §2.1.4) used for the applications in the following sections helps address the overfitting issue that can result from solution growth.

### 5.1.2. Computational cost

As a population-based method, M4GP's main drawback in comparison to other methods is computational cost, shown in terms of wall-clock time in Table 5. Each trial of the initial eight problems from Table 1 were run on a single core, so the reported time is the time to evaluate population solutions in series. The times range from about 30 s for simple problems (e.g. Movl) to about two and a half hours for larger sample size problems (e.g. IM-10). Unfortunately run-times of the other methods were not available. The selection method used by M4GP (lx, ps, or tn) does not appear to affect the computation times in a significant way. For the Opportunity Activity Recognition problems, parallel processing was used, which reduced run-times, as discussed in the next section.



**Fig. 10.** Performance on the biomedical datasets. The subplot titles indicate the dataset; datasets are named by the convention 2w\_ $[\# \text{ attributes}]_a[\text{signal-to-noise ratio}]$ . From left to right, the signal-to-noise ratio increases; from top to bottom, the dimensionality of the dataset (and hence, difficulty) increases.

## 5.2. Opportunity Activity Recognition Challenge

The results of the Opportunity Activity Recognition problem are shown in Table 4 where the results of M4GP are compared to those from the original challenge in terms of F-measure (Eq. (7)). A number of methods are reported for this problem, including nearest neighbor (NN) classifiers, SVM, decision tree and ensemble versions thereof. For both subjects tested, M4GP-lx and M4GP-ps produce better classifiers than the competition, with lexicase selection performing slightly better than Pareto survival. Because only single best results are reported

in literature, we are unable to provide statistical tests for these results. However, we report the median F-measure on the test set for 10 trials of M4GP-lx and M4GP-ps along with confidence bounds, indicating that the median performance of M4GP exceeds the best reported result from the competition.

### 5.2.1. Dimensionality reduction

We characterized the interpretability of M4GP models for this problem in terms of 1) feature selection and 2) dimensionality reduction. To quantify 1), we compared the number of attributes in M4GP solutions

**Table 6**

Significance comparisons of methods on the GAMETES problems using a post-hoc Friedman's symmetry test. Results in bold indicate  $p < 0.01$ .

	LR	XGBoost	DNN	TPOT-MDR	M4GP	M4GP+EKF	TPOT-MDR+EKF
XGBoost	1	–	–	–	–	–	–
DNN	1	1	–	–	–	–	–
TPOT-MDR	0.5	0.8	0.6	–	–	–	–
M4GP	0.2	0.4	0.2	1	–	–	–
M4GP+EKF	<b>0.001</b>	<b>0.005</b>	<b>0.002</b>	0.4	0.8	–	–
TPOT-MDR+EKF	<b>1e-09</b>	<b>2e-08</b>	<b>4e-09</b>	<b>0.0002</b>	<b>0.003</b>	0.3	–
MDR-Pred	<b>2e-07</b>	<b>2e-06</b>	<b>2e-07</b>	<b>0.002</b>	0.02	0.6	1

to the problem dimensionality; for 2), we looked at the solution dimensions, i.e.  $|\Phi|$ , compared to a principal component analysis (PCA) preserving 98% variance. The results shown in Fig. 9 demonstrate that the M4GP solutions to Opp S2 and S3 were small in size and used very few features. The feature selection component is relevant to this particular problem since it may which be beneficial for sensor or measurement selection in future design of experiments.

### 5.2.2. Computational cost

The run-times of M4GP for Opp-S2 and Opp-S3 is shown in Table 5. These problems have close to three times as many training samples as IM-10 but were run on 16 core machines, which significantly reduced the run-time to around 10 min. In this case, the problem is run in parallel using an island model [2], in which the population is divided among the cores and shuffled at set intervals. The resulting computation times indicate the improvement afforded by parallel processing in GP methods.

### 5.3. Biomedical application

M4GP and M4GP+EKF are compared to LR, XGBoost, DNN, TPOT-MDR, and TPOT-MDR+EKF on 16 GAMETES datasets in Fig. 10. The results are compared in terms of 10-fold balanced accuracy over 30 trials of each algorithm. The problems are layed out in terms of difficulty: the dimensionality of the datasets grows along the y-axis, and the strength of the signal grows from left to right. Therefore, the bottom left corner represents the hardest problem, and the top-right represents the easiest problem. As mentioned earlier, MDR-Pred results indicate

approximately the best possible accuracy for each problem. The results show that M4GP and M4GP+EKF outperform LR on every problem, and outperform XGBoost and DNN on all but one problem (2w\_10a\_0.1). We also find that M4GP+EKF is able to perform near to optimal (within 5% of MDR-Pred's performance) on 13 of the 16 problems, whereas TPOT-MDR+EKF performs near to optimal on 14 out of 16 problems. The EKF feature selection filter clear improves the performance of both M4GP and TPOT-MDR. Without EKF, M4GP outperforms TPOT-MDR by greater than 10% accuracy on 8 of the problems, and performs similarly (within 10%) on the other 8 problems.

To test the significance of these results, we conduct a Friedman's test of multiple comparisons on the entire suite of problems, which indicates significant differences between the methods ( $p < 2.9e-15$ ). A post-hoc asymptotic symmetry test is then conducted, the results of which we present in Table 6. The test finds no significant differences between the three overall best methods: MDR-Pred, TPOT-MDR+EKF, and M4GP+EKF, which supports our observations of Fig. 10. These three methods all significantly outperform LR, XGBoost, and DNN across these problems. TPOT-MDR+EKF outperforms TPOT-MDR and M4GP significantly, whereas M4GP+EKF does not, indicating that TPOT-MDR+EKF produces more consistently strong results compared to M4GP+EKF.

### 5.3.1. Dimensionality of solutions

Fig. 12 demonstrates how the M4GP solutions appear in the Pareto archive as feature sets. The archived models are the best trade-offs between complexity and accuracy found during training. The red line

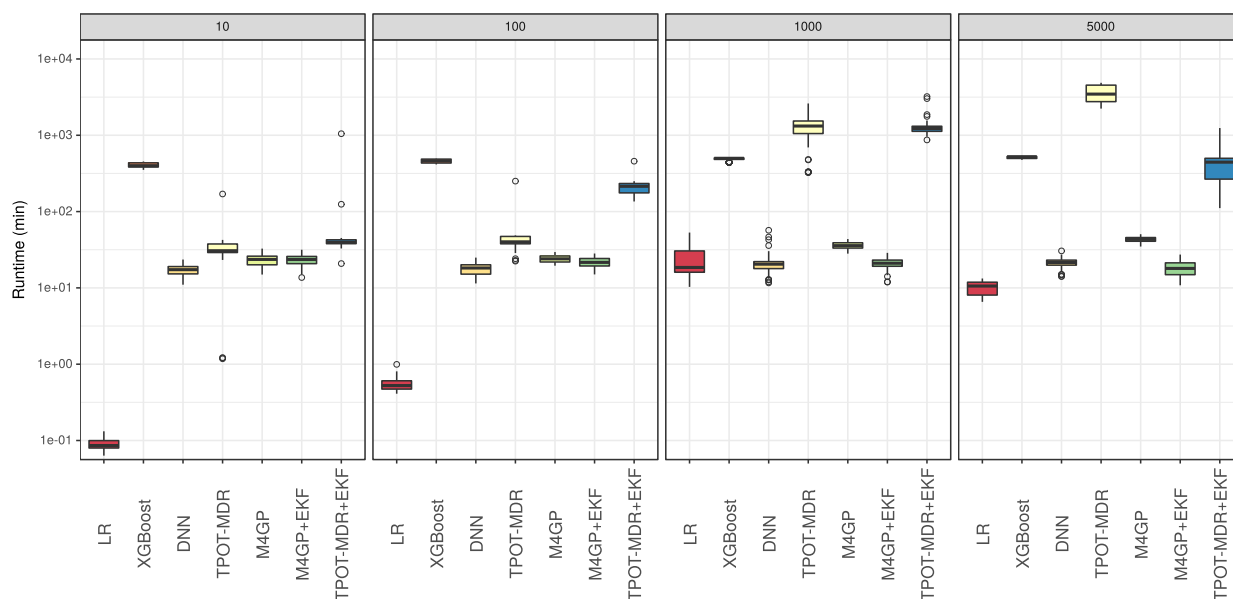


Fig. 11. Runtimes on the biomedical datasets. From left to right, the number of attributes in the datasets increase.

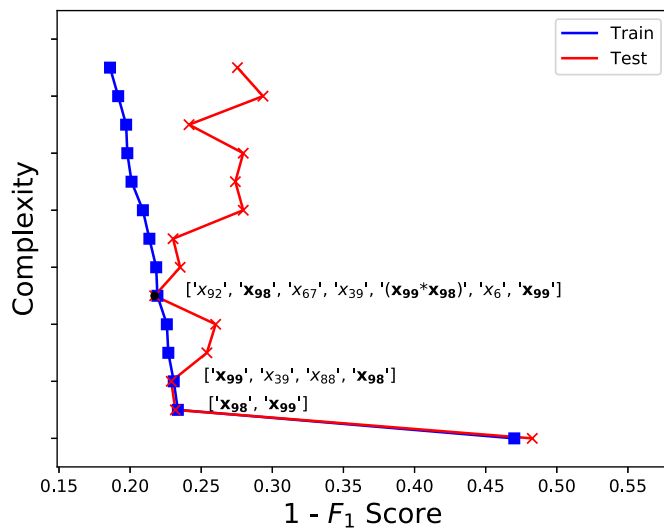


Fig. 12. Archive of solutions from M4GP for 2w\_100a\_0.4her. Blue dots are the training fitness of solutions in the archive, and red is the fitness of those solutions on the holdout test set. The top 3 M4GP solutions are shown with the interacting, predictive genes  $x_{98}$  and  $x_{99}$  in bold. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

represents the test fitness of these models. As is typical, more complex models often fail to generalize to the test set due to over-fitting. We print the solution form at the “elbow” of the Pareto archive, and for two models with the best generalization fitness. In both figures, the models produced by M4GP identify the correct underlying features in a simple form that is easily interpretable.

### 5.3.2. Computational cost and scaling

In general it appears that M4GP and M4GP+EKF achieve similar performance to TPOT-MDR and TPOT-MDR+EKF, respectively, on these problems. Notably, the M4GP methods produce these results with lower run-times, as shown in Fig. 11. Before we compare run-times, we should note that the cluster environment is subject to unequal computational loading between runs, which will affect the measurements; however, all comparisons were run on the same hardware. M4GP and M4GP+EKF take an average of 31 and 21 min to conduct 10-fold cross validation, respectively, whereas TPOT-MDR and TPOT-MDR+EKF take on average 1621 and 617 min, respectively. The run-time differences between M4GP and TPOT-MDR are expected: whereas M4GP operates over the space of mathematical operations and uses a fairly lightweight nearest centroid classifier, TPOT’s search space consists of entire machine learning pipelines, each of which must be trained each generation. The results here suggest the more lightweight approach is sufficient for capturing the interactions between features embedded in the GAMETES genetics data. Fig. 11 also shows that M4GP and M4GP+EKF scale well with increased numbers attributes in comparison to TPOT-MDR, TPOT-MDR+EKF, and LR. M4GP’s scalability to large feature spaces is due to its representation, which allows the dimensionality of solutions to adapt to the problem.

## 6. Discussion and conclusion

A new computational method for multi-class classification, based on GP, was studied in this paper. The new method is called M4GP, and it represents an improvement upon M2GP, M3GP and eM3GP, previous state-of-the-art techniques for multi-class classification with GP. It extends these methods by introducing a stack-based data flow, integrating advanced selection methods, and maintaining a Pareto archive that preserves concise models and integrates into the final model selec-

tion step. M4GP performs significantly better than these GP methods across all benchmark problems according to a Friedman test of rankings for test accuracy. We further demonstrated that advanced evolutionary selection methods, namely lexibase selection and age-fitness Pareto survival, perform better than tournament selection in producing accurate classifiers. Across the range of studied problems, the results indicate that M4GP is a competitive method for developing accurate classifiers with good generalization ability. In addition to being accurate, M4GP provides the flexibility to increase dimensionality for better classification of small dimension problems, and to decrease both the dimension of solutions and the number of necessary attributes for large dimension problems. For the Opp-S2 and Opp-S3 problems in particular, the dimensionality of the classifiers produced by M4GP constitute a major reduction from the original attribute space beyond that afforded by PCA reduction and motivate further applications of this method to large attribute problems. On a set of simulated genetics datasets from the biomedical community, M4GP was also able to produce simple classifiers that were as accurate as the best results in literature. The biomedical application demonstrates the usefulness of maintaining a Pareto archive of solutions when intelligibility is a key factor in modeling the process.

We have found that the computation time of M4GP is quite reasonable on a multi-core machine, with the largest sample size problems requiring approximately 10 min to run on a 16 core machine. On the biomedical datasets, M4GP scaled very well to large dimensions (5000 attributes) compared to other methods. Improvements to computation time could be made by further parallelizing the execution of individual models that can be run in parallel at the data (samples) and program (nodes) levels. Note that in this work we only parallelize the evaluation of the population of models.

There are several directions for future work that we briefly discuss here. For one, future work could investigate methods for mitigating the computational costs of M4GP. The largest computational cost stems from the inversion of the covariance matrix to compute the Mahalanobis distance (Eq. (2)). It may be possible to minimize this cost by only maintaining a single model  $\Phi$  and evolving a population that corresponds to the dimensions of  $\Phi$ , thereby minimizing the number of matrix inversions to once per generation. This approach is taken by ensemble-based GP methods like FEW [19] and EFS [1]. Another avenue of research could investigate other methods of pairing GP with classification strategies. M4GP implements a nearest centroid approach for this task, and previous work, summarized in §2.1.1, has explored other options. Nevertheless, it is not trivial to determine *a priori* which method of classification will pair best with GP for feature synthesis for a particular task. This suggests further research into classifier pairing optimization with GP feature selection and synthesis. Finally, more work is needed to fully examine the role of variation operators in this approach. Future work could analyze the role of mutation and crossover on multi-dimensional GP programs, with the potential to improve search by leveraging extra information regarding the performance of those sub-programs as features of the classifier. It may also be of interest to explore whether or not variation operators with semantic guarantees can be developed in this context.

## Acknowledgments

The authors would like to thank Mauro Castelli for his feedback as well as members of the Computational Intelligence Laboratory at Hampshire College. This work is partially supported by the National Science Foundation (NSF)-sponsored IGERT: Offshore Wind Energy Engineering, Environmental Science, and Policy (Grant Number 1068864), as well as Grant No. 1017817, 1129139, and 1331283. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the authors and do not necessarily reflect the views of the NSF. Biomedical results were supported by National Institute of Health grants AI116794 and ES013508. We acknowledge the finan-

cial support of PERSEIDS, PTDC/EMS-SIS/0642/2014, and BioISI R&D unit, UID/MULTI/04046/2013 funded by FCT/MCTES/PIDDAC, Portugal. Computing support was provided by Penn Medicine Academic Computing Services and the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by NSF grant number ACI-1053575 [40].

## References

- [1] I. Arnaldo, U.-M. O'Reilly, K. Veeramachaneni, Building predictive models via feature synthesis, in: Proceedings of the Genetic and Evolutionary Computation Conference, GECCO '15, ACM Press, 2015, pp. 983–990.
- [2] T.C. Belding, The distributed genetic algorithm revisited, in: ICGA-95, May 1995. arXiv: adap-org/9504007.
- [3] R. Chavarriga, H. Sagha, A. Calatroni, S.T. Digumarti, G. Tröster, J.d.R. Millán, D. Roggen, The Opportunity challenge: a benchmark database for on-body sensor-based activity recognition, *Pattern Recogn. Lett.* 34 (15) (2013) 2033–2042.
- [4] T. Chen, C. Guestrin, XGBoost: a scalable tree boosting system, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, ACM, New York, NY, USA, 2016, pp. 785–794.
- [5] W.-J. Choi, T.-S. Choi, Genetic programming-based feature transform and classification for the automatic detection of pulmonary nodules on computed tomography images, *Inf. Sci.* 212 (Dec. 2012) 57–78.
- [6] J.A. dos Santos, C.D. Ferreira, R. d. S. Torres, M.A. Gonçalves, R.A.C. Lamparelli, A relevance feedback method based on genetic programming for classification of remote sensing images, *Inf. Sci.* 181 (13) (July 2011) 2671–2684.
- [7] P.G. Espejo, S. Ventura, F. Herrera, A survey on the application of genetic programming to classification, *Syst. Man Cybern. Part C: Appl. Rev. IEEE Trans.* 40 (2) (2010) 121–144.
- [8] J. Friedman, T. Hastie, R. Tibshirani, *The Elements of Statistical Learning*, vol. 1, Springer series in statistics Springer, Berlin, 2001.
- [9] I. Guyon, A. Elisseeff, An introduction to variable and feature selection, *J. Mach. Learn. Res.* 3 (Mar. 2003) 1157–1182.
- [10] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, I.H. Witten, The WEKA data mining software: an update, *ACM SIGKDD Explor. Newsl.* 11 (1) (2009) 10–18.
- [11] T. Helmuth, L. Spector, J. Matheson, Solving uncompromising problems with lexicase selection, *IEEE Trans. Evol. Comput.* PP (99) (2014) 1.
- [12] G.E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, R.R. Salakhutdinov, Improving Neural Networks by Preventing Co-adaptation of Feature Detectors, arXiv:1207.0580 [cs], July 2012. arXiv: 1207.0580.
- [13] I. Icke, J.C. Bongard, Improving genetic programming based symbolic regression using deterministic machine learning, in: Evolutionary Computation (CEC), 2013 IEEE Congress on, IEEE, 2013, pp. 1763–1770.
- [14] V. Ingalalli, S. Silva, M. Castelli, L. Vanneschi, A multi-dimensional genetic programming approach for multi-class classification problems, in: Genetic Programming, Springer, 2014, pp. 48–60.
- [15] J.K. Kishore, L.M. Patnaik, V. Mani, V.K. Agrawal, Application of genetic programming for multicategory pattern classification, *Evolut. Comput. IEEE Trans.* 4 (3) (2000) 242–258.
- [16] I. Kononenko, Estimating attributes: analysis and extensions of RELIEF, in: European Conference on Machine Learning, Springer, 1994, pp. 171–182.
- [17] J.R. Koza, *Genetic Programming: on the Programming of Computers by Means of Natural Selection*, MIT Press, Cambridge, MA, USA, 1992.
- [18] K. Krawiec, Genetic programming-based construction of features for machine learning and knowledge discovery tasks, *Genet. Program. Evolvable Mach.* 3 (4) (2002) 329–343.
- [19] W. La Cava, J.H. Moore, Ensemble representation learning: an analysis of fitness and survival for wrapper-based genetic programming methods, in: GECCO '17: Proceedings of the 2017 Genetic and Evolutionary Computation Conference, ACM, Berlin, Germany, 2017, pp. 961–968.
- [20] W. La Cava, S. Silva, L. Vanneschi, L. Spector, J. Moore, Genetic programming representations for multi-dimensional feature learning in biomedical classification, in: Applications of Evolutionary Computation, Springer, Cham, Apr. 2017, pp. 158–173.
- [21] W. La Cava, L. Spector, K. Danai, Epsilon-lexicase selection for regression, in: Proceedings of the Genetic and Evolutionary Computation Conference 2016, GECCO '16, ACM, New York, NY, USA, 2016, pp. 741–748.
- [22] T. Li, C. Zhang, M. Ogihara, A comparative study of feature selection and multiclass classification methods for tissue classification based on gene expression, *Bioinformatics* 20 (15) (2004) 2429–2437.
- [23] M. Lichman, UCI Machine Learning Repository, University of California, Irvine, School of Information and Computer Sciences, 2013, <https://archive.ics.uci.edu/ml/index.php>.
- [24] H. Liu, L. Yu, Toward integrating feature selection algorithms for classification and clustering, *IEEE Trans. Knowl. Data Eng.* 17 (4) (Apr. 2005) 491–502.
- [25] L. Liu, L. Shao, X. Li, Evolutionary compact embedding for large-scale image classification, *Inf. Sci.* 316 (Sept. 2015) 567–581.
- [26] T. Loveard, V. Ciesielski, Representing classification problems in genetic programming, in: Evolutionary Computation, 2001, Proceedings of the 2001 Congress on, vol. 2, IEEE, 2001, pp. 1070–1077.
- [27] T. McConaghy, FFX: fast, scalable, deterministic symbolic regression technology, in: Genetic Programming Theory and Practice IX, Springer, 2011, pp. 235–260.
- [28] L. Muñoz, S. Silva, L. Trujillo, M3GP - multiclass classification with GP, in: Genetic Programming, Springer, 2015, pp. 78–91.
- [29] K.P. Murphy, *Machine Learning: a Probabilistic Perspective*, Adaptive Computation and Machine Learning Series, MIT Press, Cambridge, MA, 2012.
- [30] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al., Scikit-learn: machine learning in Python, *J. Mach. Learn. Res.* 12 (Oct) (2011) 2825–2830.
- [31] T. Perks, Stack-based genetic programming, in: Evolutionary Computation, 1994. IEEE World Congress on Computational Intelligence., Proceedings of the First IEEE Conference on, IEEE, 1994, pp. 148–153.
- [32] R. Poli, N.F. McPhee, J.R. Koza, *A Field Guide to Genetic Programming*, Lulu Press, 2008. lulu.com, [S.I.].
- [33] H. Sagha, S.T. Digumarti, J. d. R. Millán, R. Chavarriga, A. Calatroni, D. Roggen, G. Tröster, Benchmarking classification techniques using the Opportunity human activity dataset, in: Systems, Man, and Cybernetics (SMC), 2011 IEEE International Conference on, IEEE, 2011, pp. 36–40.
- [34] M. Schmidt, H. Lipson, Age-fitness pareto optimization, in: Genetic Programming Theory and Practice VIII, Springer, 2011, pp. 129–146.
- [35] S. Silva, L. Munoz, L. Trujillo, V. Ingalalli, M. Castelli, L. Vanneschi, Multiclass classification through multidimensional clustering, in: Genetic Programming Theory and Practice XIII, vol. 13, Springer, Ann Arbor, MI, May 2015.
- [36] B.W. Silverman, *Density Estimation for Statistics and Data Analysis*, vol. 26, CRC Press, 1986.
- [37] G.F. Smits, M. Kotanchek, Pareto-front exploitation in symbolic regression, in: Genetic Programming Theory and Practice II, Springer, 2005, pp. 283–299.
- [38] A. Sohn, R.S. Olson, J.H. Moore, Toward the Automated Analysis of Complex Diseases in Genome-wide Association Studies Using Genetic Programming, arXiv:1702.01780 [cs, q-bio, stat], Feb. 2017, arXiv: 1702.01780.
- [39] L. Spector, Assessment of problem modality by differential performance of lexicase selection in genetic programming: a preliminary report, in: Proceedings of the Fourteenth International Conference on Genetic and Evolutionary Computation Conference Companion, 2012, pp. 401–408.
- [40] J. Towns, T. Cockerill, M. Dahan, I. Foster, K. Gaither, A. Grimshaw, V. Hazlewood, S. Lathrop, D. Lifka, G.D. Peterson, R. Roskies, J.R. Scott, N. Wilkens-Diehr, XSEDE: accelerating scientific discovery, *Comput. Sci. Eng.* 16 (5) (2014) 62–74.
- [41] R.J. Urbanowicz, J. Kiralis, N.A. Sinnott-Armstrong, T. Heberling, J.M. Fisher, J.H. Moore, GAMEDES: a fast, direct algorithm for generating pure, strict, epistatic models with random architectures, *BioData Min.* 5 (1) (2012) 1.
- [42] United States Geological Survey (USGS), USGS Earth Resources Observation Systems (EROS) Data Center (EDC), 2012, <https://eros.usgs.gov/>.
- [43] L. Vanneschi, F. Archetti, M. Castelli, I. Giordani, Classification of oncologic data with genetic programming, *J. Artif. Evol. Appl.* 1–13 (2009) 2009.
- [44] E. Zitzler, M. Laumanns, L. Thiele, SPEA2: Improving the Strength Pareto Evolutionary Algorithm. Eidgenössische Technische Hochschule Zürich (ETH), Institut für Technische Informatik und Kommunikationsnetze (TIK), 2001.