

Genetic Programming for Evolving a Front of Interpretable Models for Data Visualisation

Andrew Lensen, *Member, IEEE*, Bing Xue, *Member, IEEE*, and Mengjie Zhang, *Fellow, IEEE*

Abstract—Data visualisation is a key tool in data mining for understanding big datasets. Many visualisation methods have been proposed, including the well-regarded state-of-the-art method t-Distributed Stochastic Neighbour Embedding. However, the most powerful visualisation methods have a significant limitation: the manner in which they create their visualisation from the original features of the dataset is completely opaque. Many domains require an understanding of the data in terms of the original features; there is hence a need for powerful visualisation methods which use understandable models. In this work, we propose a genetic programming approach named GP-tSNE for evolving interpretable mappings from a dataset to high-quality visualisations. A multi-objective approach is designed that produces a variety of visualisations in a single run which give different trade-offs between visual quality and model complexity. Testing against baseline methods on a variety of datasets shows the clear potential of GP-tSNE to allow deeper insight into data than that provided by existing visualisation methods. We further highlight the benefits of a multi-objective approach through an in-depth analysis of a candidate front, which shows how multiple models can be analysed jointly to give increased insight into a dataset.

Index Terms—Visualisation, t-Distributed Stochastic Neighbourhood Embedding, Genetic Programming, Feature Selection.

I. INTRODUCTION

VISUALISATION is a fundamental task in data mining which aims to represent data in a human-understandable graphical manner [1]. Visualisation is often touted as a useful tool for allowing practitioners to gain an understanding of the data being analysed, but state-of-the-art visualisation methods (e.g. t-Distributed Stochastic Neighbour Embedding (t-SNE) [2]) tend to be black boxes that give no insight into how the visualisation represents the original features of the data. Other methods such as autoencoders [3] and parametric t-SNE [4] produce a parametric mapping from the original dataset to the two-dimensional visualisation, but these use complex deep neural networks which are opaque to humans and provide little “intuitive” understanding. Even the most recent advances such as Uniform Manifold Approximation and Projection (UMAP) [5] acknowledge significant weaknesses with regards to interpretability. As McInnes et al. note (emphasis ours):

“For a number of uses [*sic*] cases the interpretability of the reduced dimension results is of critical importance. Similarly to most non-linear dimension reduction techniques (including t-SNE and Isomap), UMAP lacks the strong interpretability of Principal Component Analysis (PCA) and related techniques such a [*sic*] NonNegative Matrix Factorization (NMF). In particular the dimensions of the UMAP embedding space have no specific meaning, unlike PCA where the dimensions are the directions of greatest variance in the source data. Furthermore, since UMAP is based on the distance between observations rather than the source features, it does not have an equivalent of factor loadings that linear techniques such as PCA, or Factor Analysis can provide. If strong interpretability is critical we therefore recommend linear techniques such as PCA and NMF.” [5, p. 35]

While such linear techniques as PCA and NMF provide strong interpretability, they are inherently limited in the quality of their output dimensions by their use of linear weightings only. Tree-based genetic programming (GP) is often recognised for its ability to directly model functions by taking inputs as leaves and producing an output at the root [6] and is often used for dimensionality reduction in the form of feature construction (FC) [7]. Significant progress has been made on producing interpretable GP trees, which are simple enough to be understood by a human expert by using techniques such as parsimony pressure [6] and multi-objective optimisation [8], with multi-objective approaches showing particularly good results due to their ability to produce a range of solutions with different levels of complexity [9], [10]. Despite being clearly suitable for dimensionality reduction and evolving interpretable trees, GP has never been used to produce interpretable models that create understandable visualisations.

We recently investigated the use of GP to do manifold learning (GP-MaL) for dimensionality reduction [11]. We found that GP was able to significantly reduce the size of feature sets while still retaining high-dimensional structure by producing sophisticated constructed features which used a variety of non-linear and conditional transformations. Our findings also highlighted the considerable potential of GP-MaL for visualisation, but that the original fitness function needed substantial refinement and that the complexity of the evolved models would need to be significantly reduced.

In this paper we aim to significantly expand our previous work by using a multi-objective GP approach for visualisation which is expected to address the above limitations. To our knowledge, this is the first multi-objective GP method proposed which treats visualisation quality and model complexity as competing objectives. This paper will:

This work was supported in part by the Marsden Fund of New Zealand Government under Contracts VUW1509 and VUW1615, the Science for Technological Innovation Challenge (SfTI) fund under grant E3603/2903, and the University Research Fund at Victoria University of Wellington grant number 216378/3764 and 223805/3986.

The authors are with the Evolutionary Computation Research Group, Victoria University of Wellington, Wellington 6140, New Zealand (e-mail: andrew.lensen@ecs.vuw.ac.nz; bing.xue@ecs.vuw.ac.nz; mengjie.zhang@ecs.vuw.ac.nz).

- Propose a holistic set of functions and terminals for creating powerful and interpretable models for visualisation;
- Design a multi-objective approach to allow for the evolution of a solution front containing visualisations representing different levels of quality and complexity;
- Compare the quality and interpretability of the evolved visualisations with those produced by state-of-the-art visualisation methods on a range of datasets; and
- Perform an in-depth analysis of the trade-off between visualisation quality and tree complexity to demonstrate the unique advantages of the proposed approach.

This work is expected to provide a new and innovative approach which addresses the little-studied issue of model interpretability in visualisation. Of particular novelty is the in-depth analysis that will be performed that allows significant insight which is unattainable with existing single-solution black box approaches.

II. BACKGROUND

A. Dimensionality Reduction

Dimensionality reduction (DR), the process of reducing the number of features/attributes in a dataset to improve understanding and performance [12], has continued to grow in importance as datasets become increasingly big and deep neural networks become more and more un-interpretable. Common techniques to address this problem include feature selection (FS) and feature construction (FC), which reduce the feature space through removing unwanted features or creating fewer, more complex meta-features, respectively. Evolutionary Computation (EC) methods have been applied to these NP-hard problems with significant success [7], [13]. In particular, tree-based GP has proved to be a natural fit for FC problems due to its functional and interpretable nature; GP-based FC has been used in classification [14], image analysis [15], clustering [16] and other domains [17], [18].

Manifold learning (also known as non-linear dimensional reduction) [19] can be regarded as an unsupervised FC approach, where the task is to build a set of features which represent the non-linear manifold present in the high-dimensional feature space. One way of approaching this task is to attempt to build a *function*, which maps the high-dimensional space to the low-dimensional manifold; such an approach could provide an understandable mapping between the two. Until our recent work [11] there had been no investigation into whether GP could be used to evolve such a mapping.

B. Machine Learning for Performing Visualisation

The simplest commonly-used machine-learning visualisation methods are FS and Principal Component Analysis (PCA) [20]. By using FS to select only two (or perhaps three) features, one can visualise a dataset by plotting the feature values of each instance along the x- and y-axes. PCA uses a more sophisticated approach, where it creates a series of “principal components” which are linear combinations of the feature set, such that each successive component is orthogonal to all the preceding components. By plotting the first two

created components, a visualisation is created which is optimal when performing only linear transformations.

However, linear transformations fail to give clear visualisations of any underlying manifold/structure in a dataset beyond simple, low-dimensional data. Unfortunately, creating optimal non-linear transformations is an NP-hard problem, with many machine learning methods proposed as a result. The earliest methods include techniques such as Isomap [21] and Local Linear Embedding (LLE) [22], but t-SNE [2] is generally regarded as the mainstream state-of-the-art method.

t-SNE is an improvement to the previously proposed SNE method [23], which introduced a more nuanced probabilistic mapping from the high to low-dimensional spaces based on the similarity of neighbours in the high dimensional space. t-SNE improved upon SNE by introducing a cost function that could be better optimised by gradient-based optimisers; and by using a heavy-tailed t-distribution in place of a Gaussian to address the tendency of points in the low-dimensional space to be pushed together at the cost of separation between points, a characteristic known as the *crowding problem*. A range of improvements to t-SNE have since been proposed, including a more efficient tree-based nearest-neighbour search [24], and a parametric version [4], but the cost function, which is the key aspect in our work, has remained relatively unchanged.

C. Multi-objective Optimisation

Multi-objective optimisation (MO) is a technique used when a problem intrinsically has two (or more) *conflicting* objectives between which a trade-off must be made by any solution to the problem. The quality of a solution in this context is often relative to the objective function values of other solutions. For example, given two candidate solutions y, z to a problem with k objectives to be minimised, then the following test can be used to determine if y is strictly better than (or *dominates*) z :

$$\forall i : f_i(y) \leq f_i(z) \text{ and } \exists j : f_j(y) < f_j(z) \quad (1)$$

where $i, j \in \{1, 2, \dots, k\}$ for k objectives (in this paper, $k = 2$). A solution z for which this does not hold true for all other solutions y (i.e. it is *non-dominated*) is called a *Pareto-optimal* solution. The set of Pareto-optimal solutions is often called the *Pareto set*. The *Pareto front* is then the image of the Pareto set in the objective space. In practice, it is infeasible or very difficult to find the true Pareto front, and so MO algorithms hence search for the best *approximation front*.

EC methods are some of the most successful and widely used approaches to MO as their population-based structure naturally allows for multiple non-dominated solutions to be found in one run. Of the Evolutionary Multi-objective Optimisation (EMO) algorithms, the Multi-Objective Evolutionary Algorithm based on Decomposition (MOEA/D) [25] has seen widespread use in recent years. MOEA/D has been shown to produce approximation fronts with better spread and convergence than previous EMO methods [26]. Multi-objective genetic programming (MOGP) approaches have seen significant success in recent years [10], [27].

D. Related Work

1) *GP for Manifold Learning and Visualisation*: The only existing work we are aware of which used GP for MaL is our previous GP-MaL approach [11]. GP-MaL was proposed to tackle the general MaL problem, i.e. the reduction of a d -dimension high-dimensional space to a d' -dimension lower-dimensional space, where $d' \ll d$. A (single-objective) fitness function was proposed that measured how well manifold structure was preserved by measuring the preservation of n^{th} -neighbour orderings in the lower-dimensional space. GP-MaL was fundamentally designed for dimensionality reduction for data mining tasks (e.g. classification), and so $d' > 2$ was primarily used to retain as much structure of the data as possible. Some initial application of GP-MaL for visualisation was investigated with encouraging results, but it was found that a more specific fitness function would be needed to improve visualisation quality, and that tree size would need to be addressed for visualisations to be understandable.

Multi-objective GP approaches have been proposed to improve the visualisation quality of feature construction in a supervised learning context. The MOG3P method [28] aimed to produce solutions with a trade-off between three objectives: classifiability, visual interpretability, and semantic interpretability. Unlike in this work, MOG3P focused on performing *supervised learning*, using class labels to guide the learning process. The three objectives also do not appear to be strictly conflicting: on a naive classification dataset, it would be possible to achieve perfect classification performance with a simple hyperplane split, which could be represented by a simple GP tree which produces a visualisation with two very distinct classes. Even on more complex datasets, there is an inherent relationship between the visual interpretability of a dataset (i.e. how well classes are separated), and the classification performance: given well-separated classes, one would expect correct classification to result. A later approach focused on classifiability and visual interpretability only, but used several different measures for each of these objectives [10]. This work also is a supervised approach, and so tackles a significantly different problem to that of this paper. The use of GP for visualisation has also been applied to combinatorial optimisation problems, such as job shop scheduling [29].

2) *Other Approaches*: A variety of non-EC approaches have been proposed for MaL, some of which produce models that could theoretically be interpreted to understand the manifold in terms of the original features. Parametric t-SNE [4] is a variation of t-SNE that allows for re-use of learnt t-SNE representations on future samples. Parametric t-SNE constructs a mapping from the high- to low-dimensional space using restricted Boltzmann machines to construct a pre-trained feed-forward neural network model. The neural network used over 10,000 neurons on the largest dataset, heavily restricting the potential for interpretation of the network. Autoencoders [3] are another neural-network based approach, which attempt to compress the representation of the data into as narrow of a middle hidden layer as possible such that the original data can be re-created from the concise representation. To do so, autoencoders use a number of layers of varying sizes to

encode and decode the data. This provides a mapping to and from the learnt representation, but is unrealistic to interpret given the number of nodes and fully-connected topology. Attempts have been made to improve the interpretability of autoencoders [30], but success is fundamentally limited by the need for architectures which are differentiable. Self-organising maps (SOMs) [31], a variation of an unsupervised neural network, have been used for visualisation, but the number of weights scales with dimensionality and so their interpretability is limited on non-trivial datasets. The visualisation literature [32] discusses a few examples of mapping “synthesised dimensions” to original dimensions (in particular, page 4 of the supplementary material). These examples generally use a *post hoc* approach, by varying the model’s parameters or the instances [33] and analysing the effect, or using tools such as heatmaps. There is a distinct lack of literature that directly evolves simple models that use minimal unique features.

III. THE PROPOSED METHOD: GP-TSNE

Generally when applying GP to a problem, there are two main decisions to be made: what terminal and functional nodes are suitable, and how to formulate a fitness function to solve the problem. In addition, there are often other improvements made to the standard GP evolutionary process to further tailor it to the problem. Each of these three considerations are discussed in the following subsections. While each of these components builds on established work, the use of them together to produce a range of interpretable models producing high-quality visualisations is a new and substantial advance in this field. The full source code of GP-tSNE is available online¹.

A. GP Architecture

In order to project a high-dimensional space to a two-dimensional space for visualisation, it is important that a range of powerful and varied functions are available to the evolutionary process in order to produce compact but representative models. As exactly two dimensions are required for visualisation, we use a multi-tree approach where each individual contains two trees and each tree produces one dimension (i.e. the x- or y- axis) of the visualisation. A multi-tree representation is chosen rather than a co-operative co-evolution approach as the two trees must be tightly coupled (i.e. highly dependant) in order to give high-quality visualisations; co-operative co-evolution approaches tend to non-deterministically pair solutions from different sub-populations which greatly decreases the ability for such coupling to occur.

1) *Function Set*: Table I lists the nine functions and four terminals used in GP-tSNE. The four arithmetic functions are standard, with the exception of the $n+$ and $-$ function which are the only functions that can take the zero node as input. This allows a variable number of inputs to these functions, which allows the evolutionary process to easily perform mutation that effectively removes whole trees (hopefully introns: “useless” sub-trees whose output do not affect the tree output) in the

¹<https://github.com/AndLen/gptsne>

TABLE I

THE FUNCTION AND TERMINAL SETS OF THE PROPOSED METHOD.

Function No. of Inputs		Description
Arithmetic Functions		
$n+$	1-5	Flexible Addition
$-$	1-2	Std. Subtraction
\times	2	Std. Multiplication
\div	2	Protected Division
Non-Linear Functions		
Sigmoid	1	$\frac{1}{1+e^{-x}}$
ReLU	1	$\max(0, x)$
Conditional Functions		
Max	2	$\max(x, y)$
Min	2	$\min(x, y)$
If	3	if $(x < 0)$: y ; else z
Terminal Nodes		
F_i	0	i^{th} feature value
NF_i	0	Mean of i^{th} feature values from 3-nearest neighbours
Constant	0	From $U[-1, 1]$
Zero*	0	The number 0

pursuit of model simplification. The \div function performs protected division: if the denominator is 0, it returns 1.

The sigmoid and ReLU functions are unary operators that are included to allow easy transformation of (linear) inputs into a non-linear output space. These were chosen based on inspiration from auto-encoders and other neural network methods. The three conditional operators provide a different kind of non-linear transform which is quite unique to GP due to their non-differentiability, and they are expected to allow a single tree to exhibit varied behaviour depending on its inputs.

2) *Terminal Set*: The F_i terminal is commonly used in feature-based GP to use a feature as an input to a GP tree. Each feature in a dataset is assigned a distinct terminal which returns the values of the feature for a given instance. While this provides a great deal of flexibility to the EC search process, it does make the search space large for high numbers of features (m). To remedy this, we use PCA to produce the first PC on a dataset, and then select the j features from that PC which have the highest-magnitude weights. Features with higher-magnitude weights contribute more to a PC, and therefore choosing the j highest-magnitude features is a simple form of unsupervised feature selection. Each of these j features have an increased likelihood of being chosen from the terminal set, which helps the EC search process focus on a set of promising features, while still allowing it to select other features with a smaller likelihood. In this work, we set $j = \sqrt{m}$, and each of the j PCA-selected features have an *additional* j likelihood of being selected from the terminal set. For example, if $m = 64$ features, then $j = 8$: each of the 8 features chosen by PCA has $(1+8 = 9 \times$ the likelihood) of being chosen from the terminal set. $j = \sqrt{m}$ is chosen so that the number of selected features increases slowly with the total number of features. Note that while we use PCA to weight promising features more strongly, we do not use the actual transformed PCs in our terminal set, as this would make the trees very difficult to interpret.

The NF_i terminal serves two purposes: it provides a lower-noise version of a given feature, making trees less sensitive to

TABLE II

CALCULATION OF LIKELIHOOD (L.H.) FOR EACH TERMINAL TO BE SELECTED FROM THE TERMINAL SET.

Terminal	L.h.	L.h. for $m = 100$
“Normal” F_i	1	1
PCA-selected F_i	$1 + \sqrt{m}$	11
NF_i	1	1
Constant	$\lceil \frac{m}{10} \rceil$	10
Zero	$\lceil \frac{m}{10} \rceil$	10

noisy instances; and it allows a tree to encode local structure into the low-dimensional space more effectively by considering an instance’s neighbourhood. While individual features (F_i) are suitable for embedding global structure, they are not as effective for preserving local structure. The three nearest-neighbours for each instance are pre-computed by considering the ordering of the other instances by Euclidean distance.

The constant terminal is often used in GP, and is used here with a range of $[-1, +1]$ to allow different parts of the tree to have different impacts on the final output. The zero node is included solely for the $n+$ and $-$ function, and is not used by any other functions as it is either destructive or has no effect. The constant and zero nodes are chosen from the terminal set with a likelihood of $\lceil \frac{m}{10} \rceil$ as they are less likely to be useful to a function compared to a feature-based node. A summary of the terminal nodes’ weightings are shown in Table II. Note that the likelihood calculation is for a single instance in the case of a feature-based terminal, e.g. there are m different F_i terminals, each with a base likelihood of 1. All of the numerical values, including the “5” in the flexible addition, the use of 3-nearest neighbours, and the likelihood values, were empirically chosen by testing on a range of representative datasets.

B. Multi-Objective Approach

There is an intrinsic relationship in machine learning between the potential performance of a model and the complexity required to attain that performance. For example, the simplest model to separate two classes would be a decision boundary that simply thresholds at a certain point in space, whereas for three linearly-separable classes, at least two thresholds would be required (for some real α, β : $Class_A < \alpha < Class_B < \beta < Class_C$). The same is true in visualisation: the more granular (specific) a visualisation is, the more complex the function used to produce that visualisation must be. To recreate the high-dimensional structure of a complex dataset in two dimensions would require two very big and complex GP trees. Each node removed from a tree decreases the accuracy with which the tree can reproduce the high-dimensional probability distribution (in the case of t-SNE). As an analogy, consider evolving a very complex polynomial function with GP: the fewer components (nodes) in the evolved functions, the fewer inflection points available to approximate the function.

In this work, we employ a multi-objective approach to produce a set of solutions that allow for a trade-off between visualisation quality and model interpretability to be chosen. This is strictly a more difficult problem than optimising only visualisation quality (i.e. as in t-SNE) as a model must be found which maps the high-dimensional to the low-dimensional space. We choose to use MOEA/D [25] due to

the reasons discussed in Section II-C. The first objective uses t-SNE to measure the visualisation quality of a GP tree; the second objective uses tree size as a proxy measure for the interpretability of the tree. These will be discussed in turn.

C. Objective 1: Visualisation Quality

t-SNE uses conditional probabilities to represent the similarity between instances in a given dimensional space. Given two instances in the high-dimensional space, x_i and x_j , the conditional probability $p_{j|i}$ that x_j would be chosen as a neighbour of x_i is defined as [2]:

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma_i^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2/2\sigma_i^2)} \quad (2)$$

given a Gaussian with variance σ_i centred at x_i . t-SNE employs a symmetric approach where joint probability distributions are used; p_{ij} is computed as the symmetrised conditional probabilities which for n instances is defined as:

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n} \quad (3)$$

In order to remedy the crowding problem (see [2] for further details), t-SNE uses a slightly different approach in the low-dimensional space which is based on a Student t-distribution. The joint probabilities of two instances, y_i and y_j , in the low-dimensional space, called q_{ij} is computed as:

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}} \quad (4)$$

The cost function (C), which measures the extent to which the low-dimensional probability distribution does **not** match the high-dimensional probability distribution, is the difference between the two distributions as measured using the sum of the Kullback-Leibler (KL) divergences:

$$\text{Cost} = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad (5)$$

We use this cost function as the first objective, which should be **minimized**. The parameter σ is computed based on a perplexity of 40 using the approach outlined in [2]. We choose this cost function as it is well-tested and well-regarded in the MaL and visualisation literature.

D. Objective 2: Model Complexity

A common issue encountered in GP is the production of bloated trees, where a GP tree is significantly bigger than is necessary to achieve a given level of fitness. Traditionally, there is no evolutionary pressure to encourage compact trees, and so trees may contain unnecessarily complex sub-trees, or in the worst case introns. While being computationally inefficient, bloated trees are also much harder for humans to interpret and understand.

Parsimony pressure, which treats minimisation of model size as a (minor) objective in the optimisation process, is the most frequently method for controlling *bloat* in GP [34]. Weighted sum approaches, where a small component of overall fitness is based on tree size, has been often used for attempting to control bloat, but choosing the right weighting is difficult

and generally must be set empirically. More recently, multi-objective approaches have been proposed for addressing bloat, whereby tree size is used as a secondary objective to be minimised [9]. This approach allows a trade-off between fitness and model complexity to be found according to the approximation front produced by the GP process, while also producing a range of solutions of varying complexity in a single GP run which can be compared by the user for better insight into the problem being tackled.

We use a simple formula for complexity in this work which is based on the number of nodes in each of the two trees, T_a, T_b , in a GP individual I :

$$\text{Complexity}(I) = \sum_{T \in I} \sum_{N_i \in T} \begin{cases} 0, & \text{if } N_i = \text{ZeroNode} \\ 1, & \text{otherwise} \end{cases} \quad (6)$$

A given node N_i is counted towards the complexity unless it is a ZeroNode; these are not counted as they exist only to allow flexibility in the arity of the addition function, and can be removed from the tree structure when interpreting it.

E. Optimisation of Tree Constants

To further increase the visualisation quality without introducing additional model complexity, we use Particle Swarm Optimisation (PSO) [35] to fine-tune the Ephemeral random constants (ERCs) in each individual in the final front at the end of the evolutionary process. Standard GP has no ability to fine-tune its numerical parameters efficiently (as it randomly searches the parameter space); by employing PSO we can do so. Each ERC (from both trees) is allocated a dimension in the PSO representation, with the value of a given dimension representing how much the value of the given ERC varies from its original value. We use a very small range of initial position values ($[-0.15, 0.15]$) and low minimum and maximum velocities (-0.05 and 0.05) to focus the PSO search on fine-tuning the ERCs in the tree. One of the 30 particles is initialised with all values of 0 (i.e. the original ERC values) so that the PSO search is guaranteed not to produce inferior solutions. The PSO search is single-objective (as the tree structure is fixed), with the fitness function being the same cost function used as the first objective in the GP search (Eq. (5)). PSO is only used at the end of the GP evolutionary process both due to its computational cost and to prevent GP from falling into local minima that would result from fine-tuning during evolution.

Other techniques such as the Covariance Matrix Adaptation Evolution Strategy (CMA-ES) [36] have also been used for numerical optimisation, but PSO has been shown to give considerably better results on ill-conditioned functions² [37]. For GP trees of sufficient complexity, it is expected they have the characteristics of an ill-conditioned function as a change in a given ERC value is likely to significantly change the output of the tree due to the interactions present across different sub-trees. Preliminary testing confirmed that PSO could optimise the random constants of the final GP individuals more effectively and consistently than CMA-ES.

²An ill-conditioned function is one with a high condition number: its output is overly sensitive to changes in the values of its variables.

F. Other Considerations

We use a slight variation to the standard MOEA/D, whereby we prevent the breeding process from producing individuals which have already been created in earlier generations³. While this added a small amount of overhead to the evolutionary process, it gave a much more efficient and effective learning process which was ultimately found to improve the diversity and quality of the final approximated front significantly. The weight vectors in MOEA/D were scaled to encourage a uniformly distributed front. The cost objective was scaled to [0, 4] and complexity to [0, 4000].

The evolutionary process was sped up through the use of multi-threaded evaluation, caching of quality values for previously-seen trees, and the use of linear algebra libraries which allow significantly faster matrix operations through native code libraries. An overview of the GP-tSNE algorithm is shown in Algorithm 1.

Algorithm 1 Overall GP-tSNE Algorithm

```

Input: Dataset:  $X$ , maximum generations:  $G$ 
Output: Approximated Pareto front  $HOF$ 
1: Randomly initialise population  $P$ 
2:  $HOF \leftarrow \{\}$ 
   Evolutionary Loop:
3: for  $i = 1$  to  $G$  do
4:   for  $j = 1$  to  $|P|$  do
5:      $A, B \leftarrow \text{MOEAD}_{\text{Selection}}(P_j)$ 
6:      $Offspring \leftarrow \{\}$ 
7:     while  $|Offspring| < 2$  do
8:       if  $\text{Rand}() < CXPB$  then
9:          $Child_A, Child_B = \text{Crossover}(A, B)$ 
10:      else
11:         $Child_A, Child_B = \text{Mutate}(A), \text{Mutate}(B)$ 
12:      end if
13:      if  $\text{Unique}(Child_B)$  then
14:         $Offspring \leftarrow Offspring + Child_A$ 
15:      end if
16:      if  $\text{Unique}(Child_B)$  then
17:         $Offspring \leftarrow Offspring + Child_B$ 
18:      end if
19:    end while
20:     $Vis_A = \text{EvaluateInd}(Offspring_A, X)$ 
21:     $Vis_B = \text{EvaluateInd}(Offspring_B, X)$ 
22:    Evaluate  $\text{CachedQuality}(Vis)$  using Eq. (5)
    and  $\text{Complexity}(Offspring)$  using Eq. (6)
23:     $P \leftarrow \text{MOEAD}_{\text{UpdateFront}}(Offspring)$ 
24:  end for
25:   $HOF \leftarrow \text{NonDominated}(HOF, P)$ 
26: end for
27: for  $j = 1$  to  $|HOF|$  do
28:    $HOF_j \leftarrow \text{Optimise}_{CMA-ES}(HOF_j)$ 
29: end for
30: return  $HOF$ 

```

IV. EXPERIMENT SETUP

The proposed GP-tSNE method was applied to a range of representative datasets which contain well-formed classes that lend well to visualisation. The nine datasets are summarised in Table III, and have a range of numbers of instances, features, and classes. These datasets are from a number of different

³The breeding process tries up to 10 times to produce a non-duplicate individual — otherwise it returns one of the parents, chosen randomly.

TABLE III
CLASSIFICATION DATASETS USED FOR EXPERIMENTS.

Dataset	Instances	Features	Classes
Iris	150	3	3
Wine	178	13	3
Dermatology	358	34	6
Breast Cancer Wisconsin	683	9	2
COIL20	1440	1024	20
Isolet	1560	617	26
MFAT	2000	649	10
MNIST 2-class	2000	784	2
Image Segmentation	2310	19	7

TABLE IV
GP PARAMETER SETTINGS.

Parameter	Setting	Parameter	Setting
Generations	2500	Population Size	100
Mutation	20%	Crossover	80%
Min. Tree Depth	2	Max. Tree Depth	14
No. Trees	2	Pop. Initialisation	Half-and-half

domains including general classification, biology, and image analysis. Most of these datasets were sourced from the UCI repository [38]. The standard t-SNE implementation [2] was chosen as a baseline method as it has the same optimisation measure as GP-tSNE, and is regarded as the state-of-the-art in visualisation techniques. We used a standard perplexity value of 40 (the same as in GP-tSNE). We also compare to our previous GP-MaL method [11] which was the first GP method to perform MaL, using a single-objective approach. GP-MaL parameter settings were unchanged from the paper.

We found that the use of a multi-objective approach necessitated a large number of generations (2,500) in order to allow the biggest trees (with the best performance) to be trained sufficiently. However, only a small population size of 100 individuals was needed to achieve a reasonable cover of the front, especially due to the technique used to breed unique individuals. The remaining parameters (Table IV) are standard settings — tuning them further gave no change in performance.

Standard GP mutation is used, with one of the two trees in an individual being randomly chosen to be mutated. Crossover is performed by randomly selecting either the first or second tree to use, and then performing crossover between this tree in each parent — this crossover occurs between the same axis of the visualisation.

All three methods were run 30 times on each dataset. As our evaluation is primarily based on qualitative analysis (as is standard in visualisation [2], [5]), we chose the median result out of the 30 runs according to the objective function used by the method: t-SNE’s cost function for t-SNE and GP-tSNE, and the fitness function proposed in [11] for GP-MaL.

t-SNE runs (by default) for up to 1,000 iterations, corresponding to 1,000 evaluations of the t-SNE cost function. GP-tSNE runs for 2500 generations with a population size of 100, giving 250,000 evaluations. Clearly, GP-tSNE is computationally more expensive, although the use of fitness caching and multi-threading reduces this expense. Also, GP-tSNE produces many visualisations in a single run, whereas t-SNE produces only one. On the biggest dataset, Image Segmentation, GP-tSNE takes 30 hours, whereas t-SNE takes around two min-

utes. We intend to significantly reduce this difference in future work, through the use of surrogate techniques and tuning of the GP design, but note that GP-tSNE is highly parallelisable, and so can run much quicker in practice.

V. RESULTS AND DISCUSSION

The results across each of the nine datasets are shown in Figs. 1–9. Each figure contains eight plots, which we refer to as (a)–(h) reading left-to-right and top-to-bottom. Plot (a) shows the attained approximation front of the (median) GP result, with blue crosses showing each individual in the front, and the black line showing the shape of the front. The y and x-axes show each individual’s value for Objective 1 (cost/visualisation quality) and 2 (model complexity) respectively. The grey lines show the worst and best of the 30 GP-tSNE runs, which show the variance of the GP-tSNE algorithm. The green and yellow dotted lines represent the cost achieved by the baseline t-SNE and GP-MaL results respectively. Plot (b) shows the same front, but with the x-axis truncated at 200 complexity, so as to focus on the distribution of the simpler/more interpretable trees. Plots (c)–(f) show representative visualisations — coloured by class label — produced by individuals at different levels of model complexity (highlighted in red on the attained front): (c) is the simplest model found by GP-tSNE where each tree is a single feature (i.e. feature selection (FS)), (d) and (e) are two characteristic samples from along the front, and (f) is the most complex model, with the most accurate visualisation and lowest cost. A small amount of jitter is added to visualisations with low complexity (less than 20) so that overlapping points can be better distinguished. Plots (g) and (h) are the t-SNE and GP-MaL baseline results with median performance, where the cost written above the plots is calculated using t-SNE’s objective function (the first objective of GP-tSNE), to allow for sensible comparisons.

It is only possible to show a small sample of the evolved visualisations in this paper, but one of the significant advantages of GP-tSNE is its ability to produce a series of visualisations which progressively become more complex but more accurate representations of the high-dimensional feature space. To demonstrate this, we have included a video that shows each individual along the approximation front for each dataset, as the front is traversed from least to most complex. This is included with our submission, and is also available on YouTube⁴. Specific visualisations can be studied in more detail by lowering the playback speed (e.g. to 25% on YouTube).

A. Specific Analysis

On the simplest dataset, Iris (Fig. 1), all three methods provide a clear 2D visualisation, although t-SNE and GP-tSNE more distantly separate the green class. The Sample 1 visualisation for GP-tSNE is quite similar to that produced by GP-MaL when rotated, but with less continuous spacing of points, likely due to the trees being too small to convert the low-granularity feature space to a more complex output space. The front shows that diminishing returns are quickly

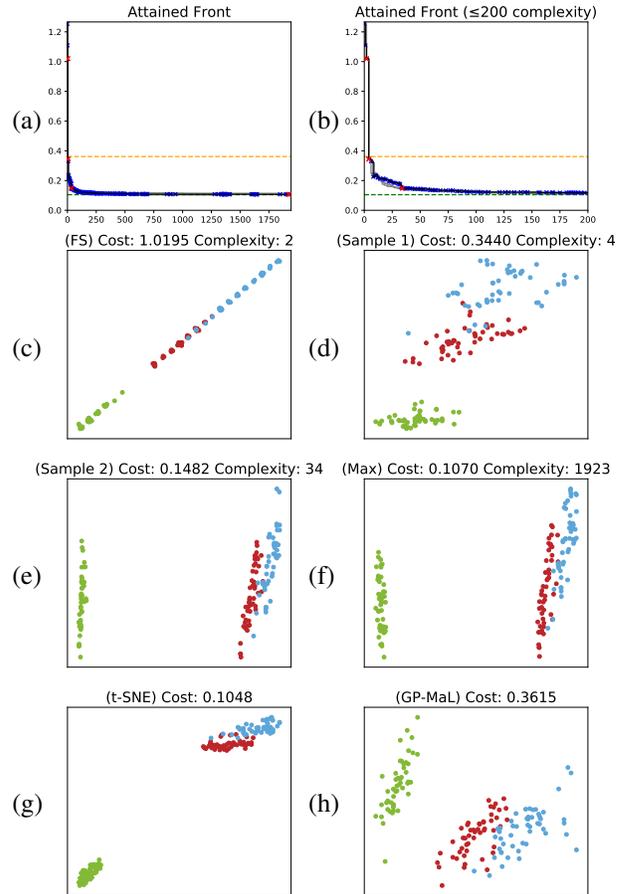


Fig. 1. Iris.

found on this dataset: after ~34 tree size (sample 2), quality only gradually improves as tree size increases. Indeed, the difference between the visualisations at complexity 34 and 1923 are very minor — clearly GP-tSNE can capture the bulk of the structure in the simple Iris dataset using a simple model.

Figure 2 shows the results on another simple dataset: Wine. All three methods produce clear visualisations, with the main difference being that t-SNE projects all 3 classes across the visualisation diagonally, whereas the two GP methods produce a curved pattern. Even at a low level of complexity, GP-tSNE clearly separates the three classes, and higher levels of complexity tend to mostly refine the local structure within each class. As on the Iris dataset, the cost achieved by GP-tSNE at its maximum complexity is reasonably close to t-SNE, despite being produced by a mapping rather than just an embedding.

Both t-SNE and GP-tSNE clearly separate the two classes on the Wisconsin Breast Cancer dataset (Fig. 3), whereas GP-MaL does not give as clear a separation boundary. GP-tSNE begins to show the division between the two classes as early as at a complexity of 33, with the red class becoming more tightly packed as the complexity increases further. Indeed, it appears that the GP models retain the same general “approach” in sample 1, 2 and at the maximum complexity, but become increasingly refined and accurate. In this way, it could be possible to use the simpler models to explore the more accurate patterns found by the bigger un-interpretable models by exploring how points move as complexity increases.

⁴<https://www.youtube.com/watch?v=K-z1jhBDHIY>

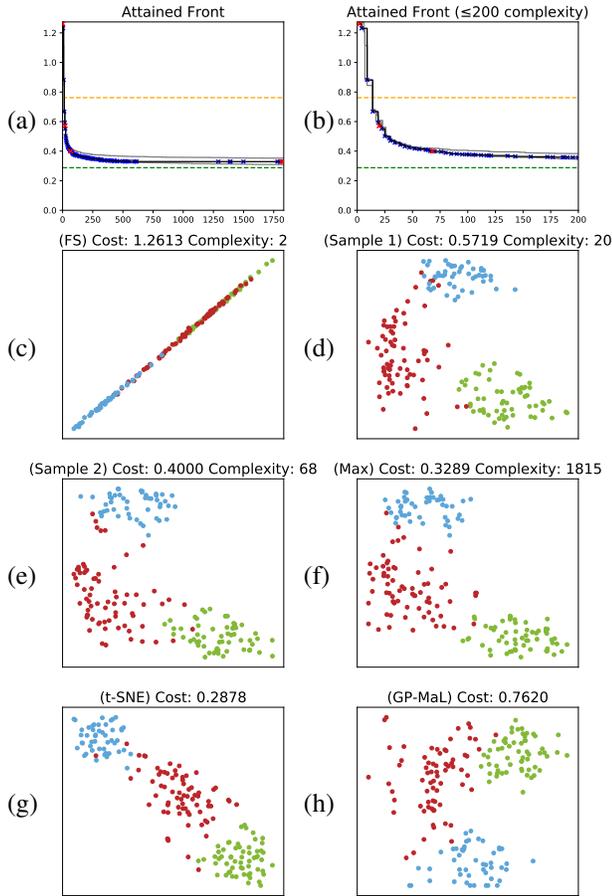


Fig. 2. Wine.

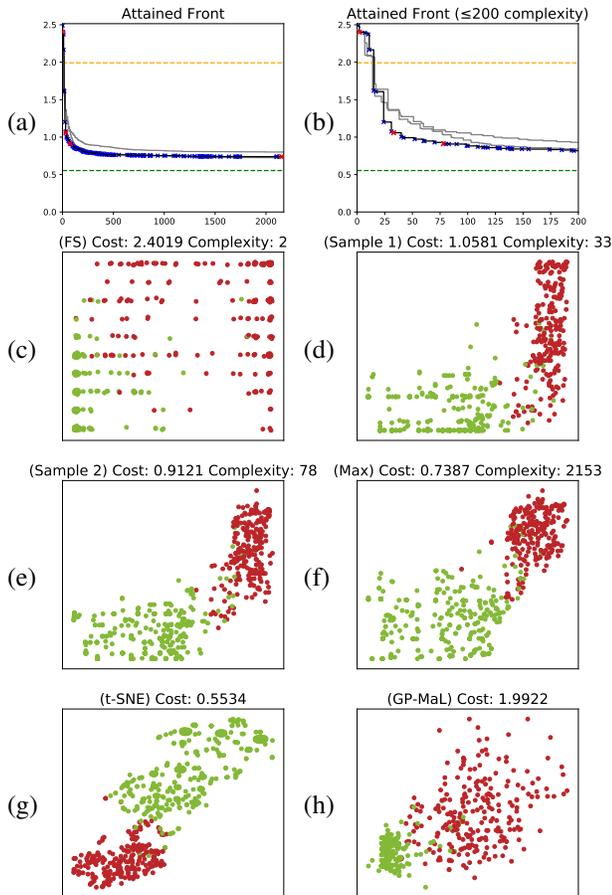


Fig. 3. Breast Cancer Wisconsin.

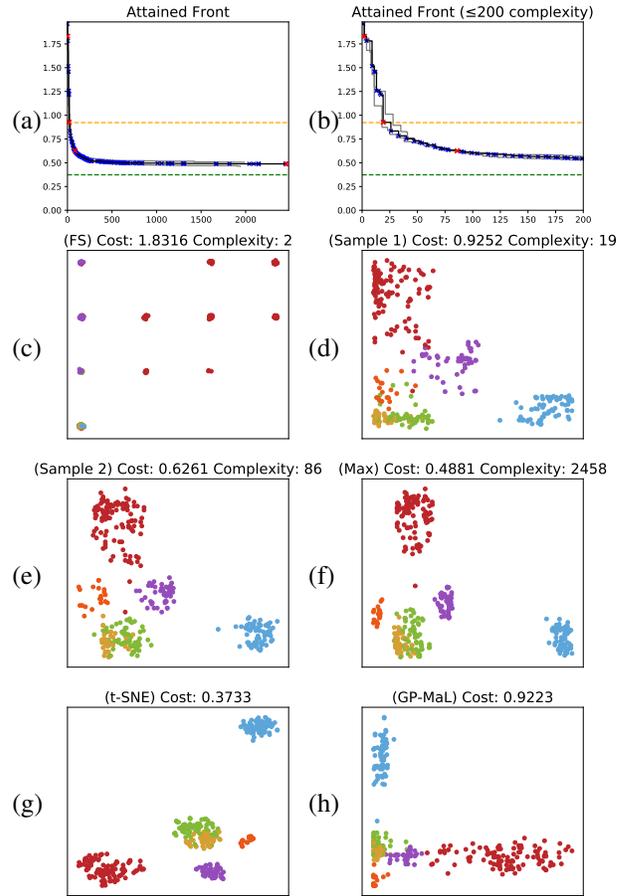


Fig. 4. Dermatology.

The Dermatology dataset (Fig. 4) further reinforces this pattern: as GP-tSNE produces increasingly complex models, the visualisation becomes of a “higher-resolution” and produces points that appear to be in a continuous space, rather than a discrete space. Sample 1 shows the blue class clearly being separated, and the red and purple classes beginning to separate from the other classes. Sample 2 shows a clear separation between all classes except the yellow and green ones, despite having a complexity of only 86. The maximum complexity separates the classes further, and compresses each class more tightly. The lowest cost for GP-tSNE gives a very similar visualisation to t-SNE, despite having a cost that is higher. GP-MaL is clearly worse than GP-tSNE, even when GP-tSNE has a tree size as low as 86.

The COIL20 dataset (Fig. 5) highlights t-SNE’s ability to freely move points throughout the 2D space, as it clearly produces the best visualisation. It is interesting to note however that GP-tSNE is able to separate some classes well, even at low model complexity. For example, in Sample 1 (complexity of 22), a number of classes (dark green, pink, blue) are already clearly visible — this suggests that these classes may be particularly well-defined, as they are easily separated. As the complexity increases, additional classes separate out, and the curve topology of the blue/pink/red classes seen in the t-SNE visualisation starts to form for GP-tSNE too.

Figure 6 shows how the Isolet dataset has many overlapping classes, with only a few classes distinctly separated by t-SNE. The GP methods generally overlap the same classes as

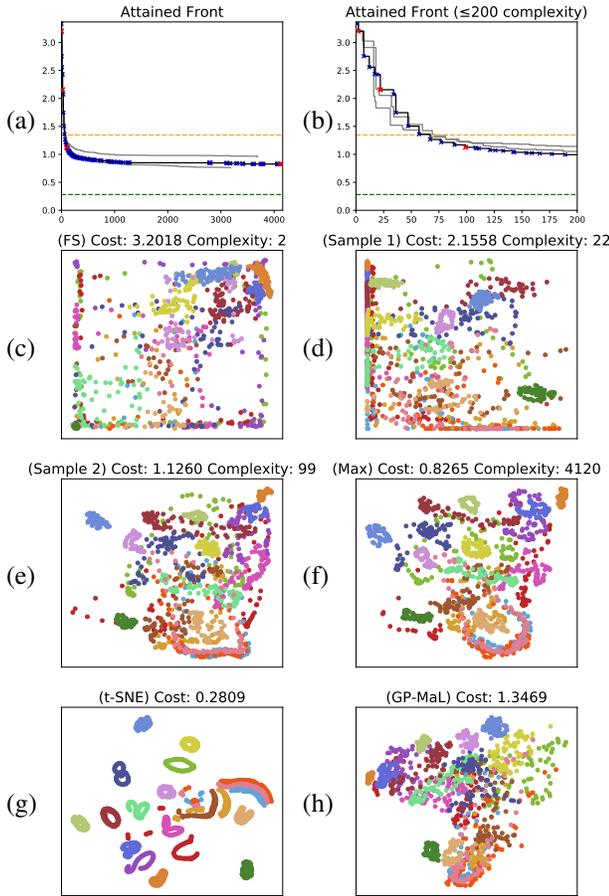


Fig. 5. COIL20.

t-SNE, but do not manage to separate the groups quite as successfully. Isolet has the highest number of classes (26) of all the datasets, which makes it especially challenging for GP-based methods, since evolving two functions that can provide sufficient granularity to separate 26 classes is clearly very challenging and requires sufficiently complex trees. Despite this, it may be possible to gain some insight into why given classes overlap, as the general overlapping patterns start to be visible even at lower complexities of 39 to 131.

On the MFAT dataset (Fig. 7), GP-tSNE is able to group each class effectively at the maximum complexity, but does not show the clear separation between classes that t-SNE provides. GP-tSNE clearly performs significantly better than GP-MaL, which overlaps the light green and purple classes onto other classes. The grouping of classes starts to appear at a complexity of 94, although the orange and purple classes overlap still. A complexity of 94 corresponds to two trees with 94 nodes in total — which is clearly quite restrictive given that the MFAT dataset has 649 features and 10 classes.

The visualisations in Fig. 8 (MNIST 2-class) provide an example of a large dataset with a small number of classes. All three methods are able to separate the two classes reasonably well, with GP-tSNE and t-SNE producing quite similar results. This separation is evident even at low model complexity in GP-tSNE – at a complexity of 51, it is possible to draw a vertical line through the middle of the visualisation that would separate the two classes with reasonable accuracy. At higher levels of complexity, GP-tSNE is able to push the two classes

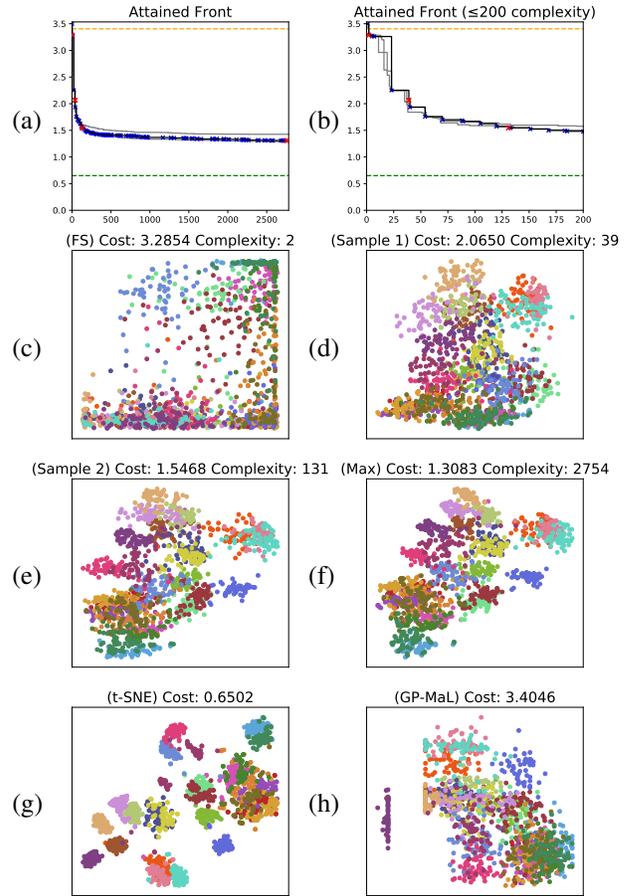


Fig. 6. Isolet.

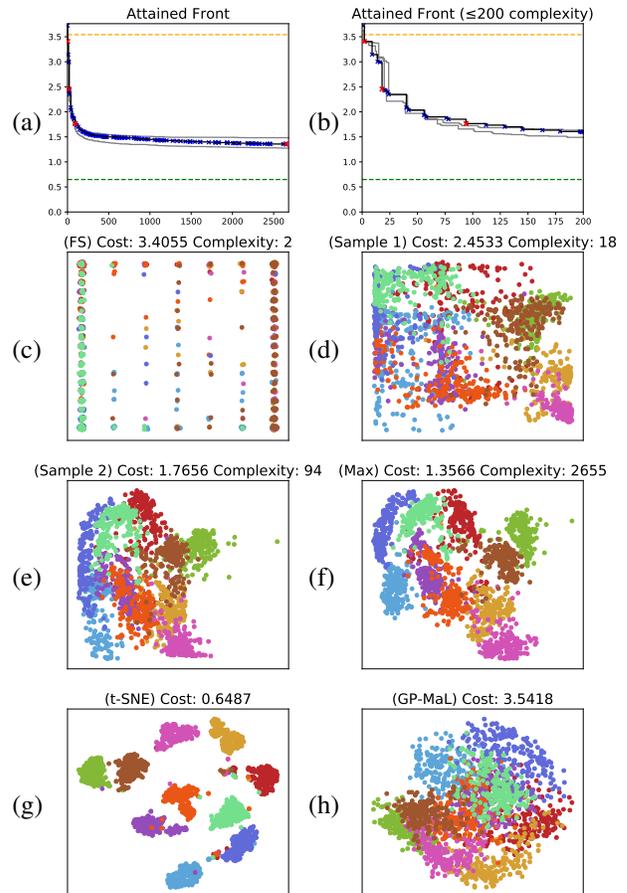


Fig. 7. MFAT.

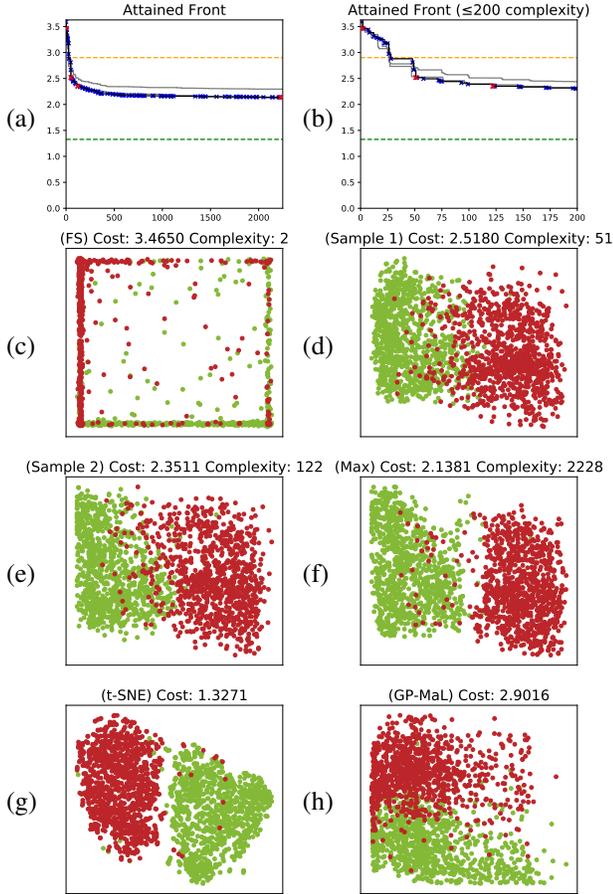


Fig. 8. MNIST 2-class.

apart to leave a clear gap between them. GP-MaL produces a less clear separation and struggles to distribute the points in the visualisation space well.

On the final dataset, Image Segmentation (Fig. 9), all three methods are able to separate the green and orange classes out, with t-SNE and GP-tSNE providing clearer separation margins. Both GP-tSNE and t-SNE also separate the red class to some extent, with t-SNE doing so slightly more effectively. An interesting observation is that the orange, and to a lesser extent the green classes, are separated at very low model complexity. The orange class is easily separated using only one feature, whereas the green class can be separated at a complexity of 31. This suggests that perhaps these are more natural/intrinsic classes in the dataset; perhaps they exhibit characteristics that make them particularly distinct from the other instances. Another interesting finding is that the red class is separated well by GP-tSNE at a complexity of 11 — but then overlaps with the purple class at higher complexity. This may be a limitation of the t-SNE objective function or may suggest that the red class is actually a *subclass* of the purple class and is incorrectly over-separated at a low complexity. Indeed, the red class corresponds to pixels that are part of a path in an image, and the purple class to cement pixels. Clearly, many different things can be constructed out of concrete, including paths, which explains the overlap that is evident at higher complexities. This phenomenon is another useful characteristic of GP-tSNE: it shows the different relationships between classes at different levels of model complexity.

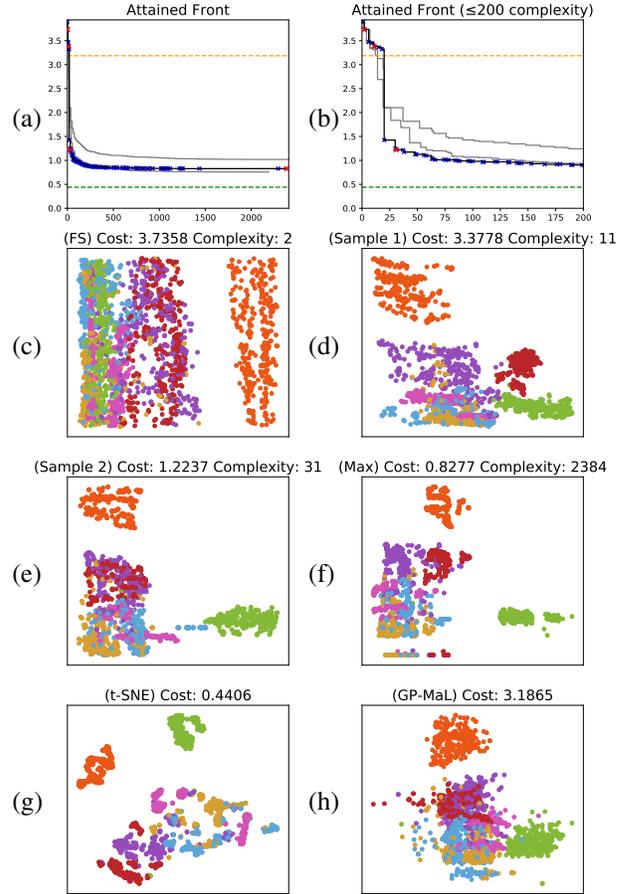


Fig. 9. Image Segmentation.

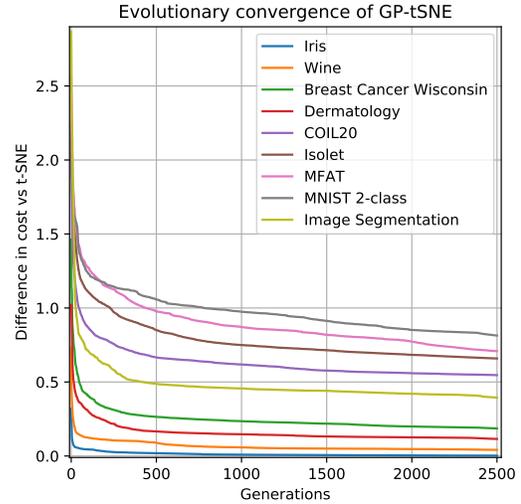


Fig. 10. Convergence of GP-tSNE over the evolutionary process: y-axis shows the difference in t-SNE cost between the best individual at a given generation and that attained by t-SNE.

B. General Findings

The visualisations produced by GP-tSNE were consistently superior to those of our previous GP method, GP-MaL. GP-MaL was developed for more general MaL (i.e. not just reducing to two dimensions); by using a visualisation-specific quality measure, GP-tSNE was able to clearly improve. As the number of instances increased, t-SNE began to produce clearer visualisations than GP-tSNE, although the same general patterns were shown by both methods. Given GP-tSNE is

in essence tackling a strictly more difficult task, of evolving a **functional model** to produce a visualisation, we see this as a success for the first such approach to this task. This pattern can be further seen in the convergence analysis shown in Figure 10. On the easier datasets (the first few in the legend), GP-tSNE converges before 1,000 generations. On the hardest datasets such as MNIST and MFAT, GP-tSNE is likely to benefit from additional generations. We are confident that future research in this new research direction will be able to further improve GP-tSNE on high-dimensional datasets with many classes.

One particularly interesting observation was how GP-tSNE produced similar overall visualisations at different model complexities, with the more complex models being more granular/refined than the simpler ones. This demonstrates the advantages of an evolutionary approach: good sub-trees of a complex tree can be used to improve simpler trees (and vice-versa) via crossover, allowing for more efficient search that produces models with different levels of trade-off. In the following section, we will explore this phenomenon in more depth to highlight the novel advantages of GP-tSNE.

VI. FURTHER ANALYSIS

The results of GP-tSNE on the Dermatology dataset (Fig. 4) were of particular interest as they showed similar results between GP-tSNE and t-SNE; showed the same general visualisation, but at different qualities across tree sizes; and achieved good visualisations even at reasonably small model complexities. To further demonstrate the value of GP-tSNE, this section analyses a selection of increasingly complex trees produced by the median GP run on this dataset.

A. Simple Models

The simplest model which gives a reasonable visualisation is shown in Fig. 11, which contains one tree with three nodes, and the other with 10, for a total complexity of 13. Even with such a low complexity, three of the classes start to appear, with the blue, purple, and red classes already showing a clear separation from the others. The top tree, which gives the x-axis of the visualisation, can be written as $x = f_5 + f_{11} + f_{32} + 2nf_{14} - f_{21} - nf_{21}$. The bottom tree, which gives the y-axis is simply $y = nf_{20} - f_{15}$. These two trees use a total of seven unique features out of the 34 in the original feature set, yet they are able to separate three classes of the dataset well. In particular, the blue, red and purple classes can be separated out by drawing two vertical lines through the plot. In other words, two thresholds can be applied to the output of the top tree (x-axis) in order to roughly separate the classes into three groups: red; orange and green; and purple and blue. Given that all features have the same range of $[0, 1]$ and the only feature subtracted in the tree is f_{21} , this suggests that the blue and purple classes have particularly small feature values for f_{21} (as they have high x values), and the red class has particularly large values (as it has low x values). f_{21} in the Dermatology dataset corresponds to the “thinning of the suprapapillary epidermis” feature, which is a feature commonly associated with the skin condition of psoriasis [39]. In the visualisation in Fig. 11, the red class corresponds to

a diagnosis (class label) of psoriasis — and indeed, the red instances appear on the left side of the plot, which is consistent with having a high value of f_{21} being subtracted from the rest of the tree. This sort of analysis can be used by a clinician to understand that this feature alone could be used to diagnose psoriasis with reasonable accuracy, and also provides greater confidence in the visualisation’s accuracy.

When the model complexity is increased to 19 (Fig. 12), the separation of points within classes starts to become better-defined, and the orange class starts to become more distinct from the yellow and green classes. The top tree actually uses fewer unique features (three) than at a complexity of 13, with significant weight put on f_{32} : $x = 4nf_{32} + f_{32} + 2nf_{14} + nf_3 + f_3$, whereas the bottom tree uses four unique features: $y = nf_{20} + f_{20} + nf_{19} + f_{19} + nf_8 - nf_{27}$, again for a total of seven unique features across both trees. The blue class (“lichen planus”) is clearly distinct from the other classes along the x-axis — given that the top tree weights f_{32} heavily, it seems likely that a high value of this feature is characteristic of the “lichen planus” diagnosis. Indeed, the dermatology literature commonly reports on this symptom being indicative of this diagnosis [39].

The trees shown in Fig. 13 appear very similar to the previous model, with the top tree varying only by the omission of one nf_{14} node, and the introduction of the subtraction of the nf_8 feature. This is, however, sufficient to cause some of the orange points to be separated from the yellow and green points along the x-axis — indicating that higher values of nf_8 may indicate a point belongs to the orange class rather than the yellow one. The major change in the bottom tree is the introduction of the nf_{21} and nf_6 features, which helps to compact the blue class and starts to separate the red class more strongly.

Figure 14 further pushes the blue class away from the other classes, at the expense of squashing the other classes together. This is achieved by making the top tree (x-axis) weight f_3 and f_{32} even more strongly. It is interesting to note that all the trees analysed so far have been strictly linear functions, utilising only arithmetic and subtraction. This is perhaps not unexpected: utilising more sophisticated functions is likely to require complex trees to fully take advantage of them — for example, taking the maximum of two simple sub-trees is unlikely to be more representative of the original dataset than simply adding together features. This also has the benefit of making the simpler trees easier to interpret and understand. While there are many methods for doing linear transformations for visualisation, such as PCA, these methods generally function by weighting **all** features to different extents; GP-tSNE only selects a subset of features to be used, and so is inherently more interpretable. The sequential analysis of increasingly more complex models clearly provides additional insight that would be lacking in a non-multi-objective approach. This analysis technique is also an exclusive benefit of GP-tSNE among other visualisation techniques which are primarily black boxes with one visualisation produced per run.

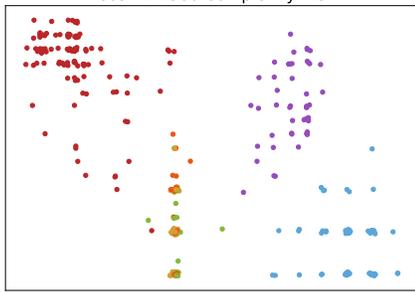
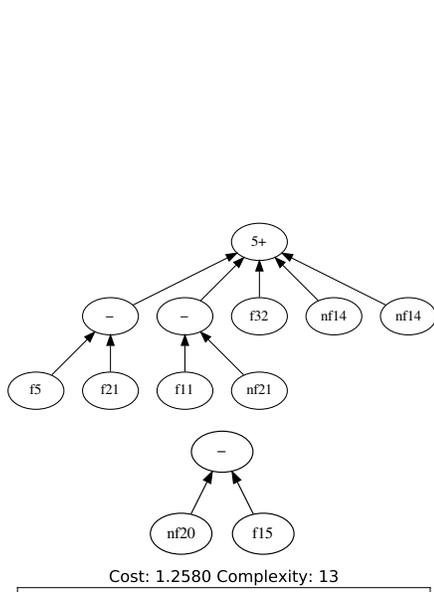


Fig. 11. Complexity of 13.

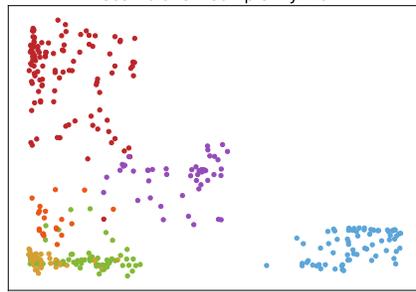
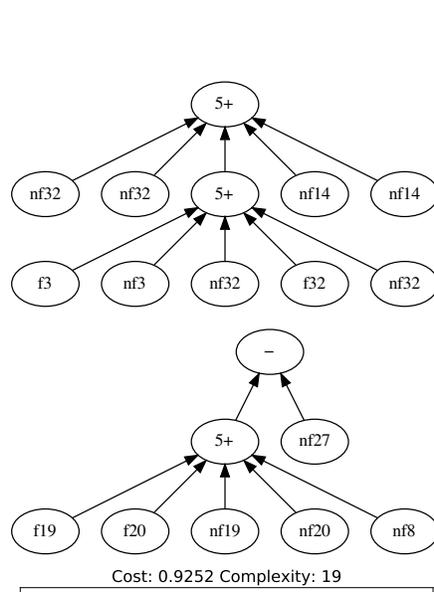


Fig. 12. Complexity of 19.

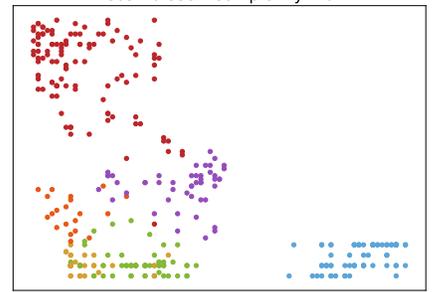
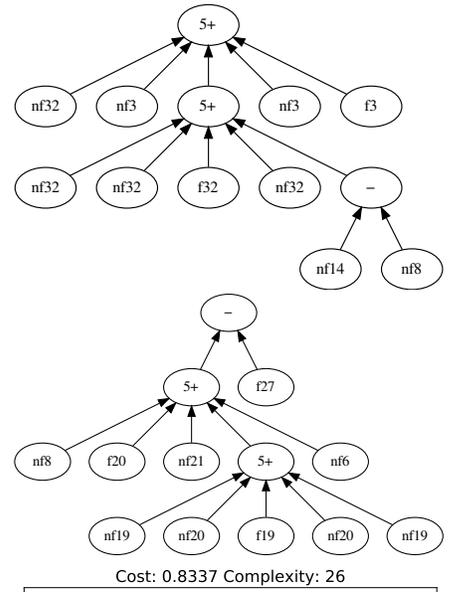


Fig. 13. Complexity of 26.

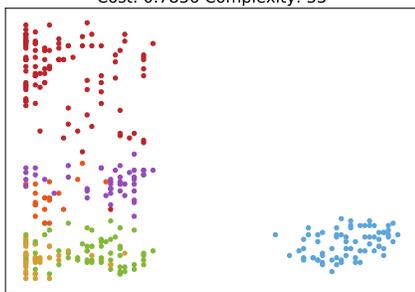
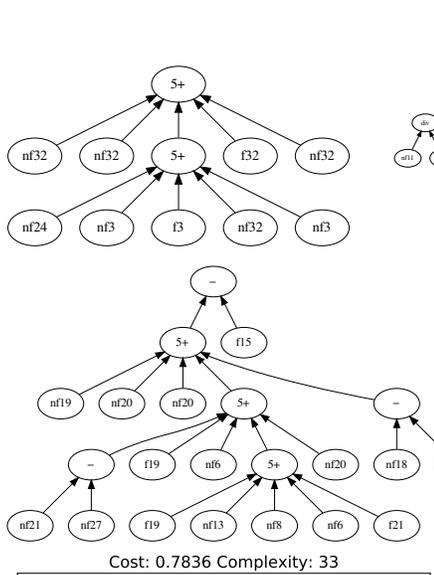


Fig. 14. Complexity of 33.

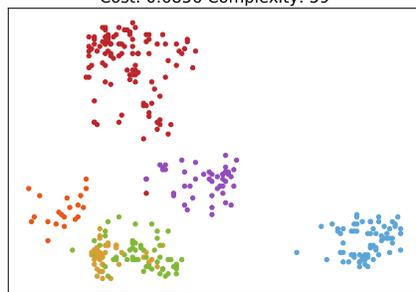
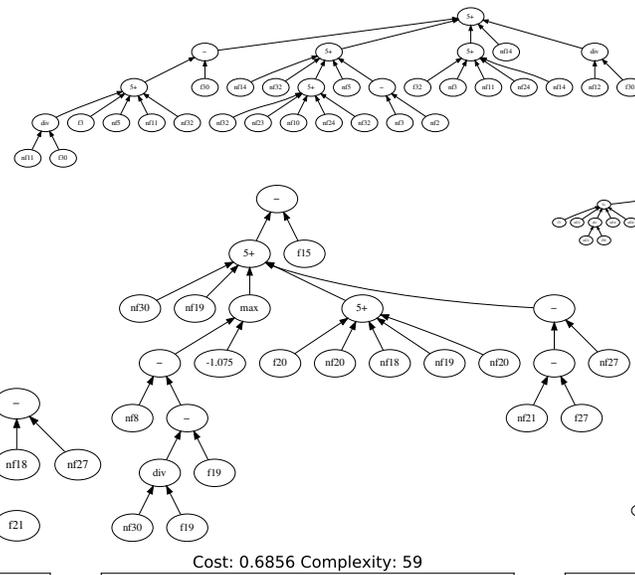


Fig. 15. Complexity of 59.

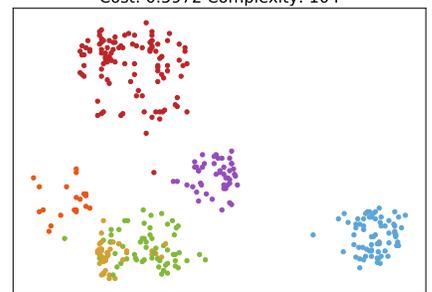
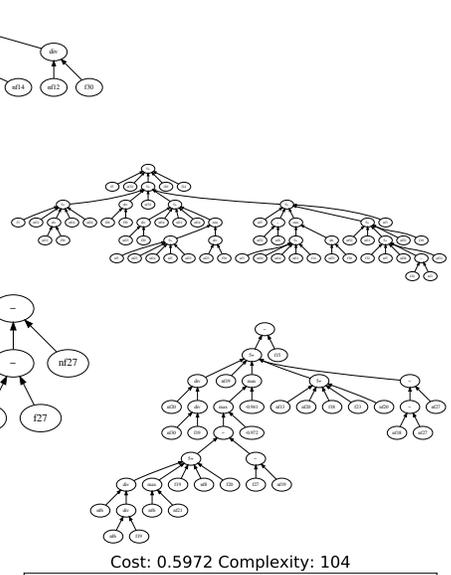


Fig. 16. Complexity of 104.

B. Complex Models

From here, the improvement in quality with the addition of more complexity begins to show clear diminishing returns, with the cost decreasing by only about 0.1 as the complexity is roughly doubled from 33 to 59, and from 59 to 104. Figure 15 gives a visualisation that is very similar to that of the maximum complexity (see Fig. 4), with the yellow and green classes overlapping but all other classes clearly distinctly separated. This is also the first stage at which constants and non-linear operators such as $-$ and \max are introduced. The trees used in previous models are still recognisable as sub-trees (or “building blocks”) in the more complex trees in Fig. 15 — it may be possible to analyse the more complex trees based on what has been added in addition to these sub-trees. Figure 16 does not change the class separations present in Fig. 15, but better separates within-class instances, as well as moving the purple class away from the other classes.

We do not attempt to analyse models with a complexity of over 100 in detail, as the trees become difficult to interpret easily, and eventually become as much of a black box as t-SNE itself is. The main aim of this paper is to produce good-quality (not state-of-the-art) visualisations that are *interpretable* in order to provide deeper insight into the relationships within a dataset in terms of the features it uses. If visualisation quality is the primary goal, then methods such as t-SNE are expected to be more appropriate, given they are not constrained to finding a functional mapping from the original feature space to the visualisation axes. While the very complex GP trees cannot be feasibly interpreted, they are still important to the GP-tSNE algorithm. During the evolutionary process, it is common to see complex trees with low cost being automatically simplified through the removal of superfluous sub-trees via crossover and mutation, allowing reductions in cost to “trickle down” to simpler models along the front. When we restricted the EMO search to only simple trees, we found that results were much poorer, with smaller trees having much higher cost. The complex trees are also useful outright as they can be directly applied for visualising future instances (e.g. streaming data), whereas t-SNE must re-optimize its embedding each time.

C. Summary

In this section, we showed how progressive examination of the approximation front from least to most complex models allowed insight into the structure of the data that is not easily achievable through traditional visualisation algorithms. Using the Dermatology dataset as a case study, we showed that even simple models could separate out some classes clearly, with the use of only a few features. As we increased the complexity, we found that the granularity/local structure within classes improved, but the overall patterns changed only slowly. In this way, it is possible to use simple, understandable models to provide insight into how more complex models are able to produce good-quality visualisations.

VII. CONCLUSION

This paper highlighted the need for research into machine learning methods that can not only produce clear visualisations, but do so through a model that can be interpreted

itself to give insight into how the visualisation arose from the features of the dataset. The use of GP was suggested due to its functional model structure, and a multi-objective approach was proposed that gave a staggered front of visualisations with different levels of trade-off between visual clarity and model complexity. Results on nine different datasets highlighted the promise in using a GP approach for this task, with visualisations produced that showed clear characteristics of datasets even at model complexities that could be low enough to understand. This was further showcased through an in-depth analysis of an evolved front on the Dermatology dataset, where concrete insight into how the dataset was structured was found by examining a range of models with different complexities.

As this is the first multi-objective GP approach to this problem, there is a clear need for continued future research. While the quality of visualisations were encouraging, on more complex datasets they fell short of those produced by t-SNE. This is not unexpected given t-SNE does not produce a functional mapping, but rather an embedding of the data. The measure of visualisation quality used in this work was based on the one used by t-SNE given it is the state-of-the-art measure. However, the design of t-SNE was constrained by the need for a differentiable cost function; as an EC-based method, GP-tSNE need not have this constraint on its objective function — there is likely to be better measures of visualisation quality that could be used in an EC context. It was also found that the more unique features used by a GP tree, the more difficult it was to understand; some sort of evolutionary pressure towards trees using a small set of (cohesive) features may improve this.

REFERENCES

- [1] U. Fayyad, G. Piatetsky-Shapiro, and P. Smyth, “From data mining to knowledge discovery in databases,” *AI Magazine*, vol. 17, no. 3, pp. 37–54, 1996.
- [2] L. van der Maaten and G. E. Hinton, “Visualizing high-dimensional data using t-SNE,” *Journal of Machine Learning Research*, vol. 9, pp. 2579–2605, 2008.
- [3] G. E. Hinton and R. R. Salakhutdinov, “Reducing the dimensionality of data with neural networks,” *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [4] L. van der Maaten, “Learning a parametric embedding by preserving local structure,” in *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics, AISTATS*, 2009, pp. 384–391.
- [5] L. McInnes and J. Healy, “UMAP: uniform manifold approximation and projection for dimension reduction,” *CoRR*, vol. abs/1802.03426, 2018. [Online]. Available: <http://arxiv.org/abs/1802.03426>
- [6] R. Poli, W. B. Langdon, and N. F. McPhee, *A Field Guide to Genetic Programming*. lulu.com, 2008, (Last Accessed: 27/09/19).
- [7] K. Neshatian, M. Zhang, and P. Andreae, “A filter approach to multiple feature construction for symbolic learning classifiers using genetic programming,” *IEEE Trans. Evolutionary Computation*, vol. 16, no. 5, pp. 645–661, 2012.
- [8] S. Bleuler, M. Brack, L. Thiele, and E. Zitzler, “Multiobjective genetic programming: Reducing bloat using SPEA2,” in *Proceedings of the Congress on Evolutionary Computation*. IEEE, 2001, pp. 536–543.
- [9] M. Wagner and F. Neumann, “Parsimony pressure versus multi-objective optimization for variable length representations,” in *Proceedings of the 12th International Conference on Parallel Problem Solving from Nature (PPSN) XII*, 2012, pp. 133–142.
- [10] A. Cano, S. Ventura, and K. J. Cios, “Multi-objective genetic programming for feature extraction and data visualization,” *Soft Comput.*, vol. 21, no. 8, pp. 2069–2089, 2017.
- [11] A. Lensen, B. Xue, and M. Zhang, “Can genetic programming do manifold learning too?” in *Proceedings of the European Conference on Genetic Programming (EuroGP)*, ser. Lecture Notes in Computer Science, vol. 11451. Springer, 2019, pp. 114–130.

- [12] H. Liu and H. Motoda, *Feature selection for knowledge discovery and data mining*. Springer Science & Business Media, 2012, vol. 454.
- [13] B. Xue, M. Zhang, W. N. Browne, and X. Yao, "A survey on evolutionary computation approaches to feature selection," *IEEE Trans. Evolutionary Computation*, vol. 20, no. 4, pp. 606–626, 2016.
- [14] B. Tran, B. Xue, and M. Zhang, "Genetic programming for feature construction and selection in classification on high-dimensional data," *Memetic Computing*, vol. 8, no. 1, pp. 3–15, 2016.
- [15] H. Al-Sahaf, A. Al-Sahaf, B. Xue, M. Johnston, and M. Zhang, "Automatically evolving rotation-invariant texture image descriptors by genetic programming," *IEEE Trans. Evolutionary Computation*, vol. 21, no. 1, pp. 83–101, 2017.
- [16] A. Lensen, B. Xue, and M. Zhang, "GPGC: genetic programming for automatic clustering using a flexible non-hyper-spherical graph-based approach," in *Proceedings of the Genetic and Evolutionary Computation Conference*. ACM, 2017, pp. 449–456.
- [17] S. Ahmed, M. Zhang, L. Peng, and B. Xue, "Multiple feature construction for effective biomarker identification and classification using genetic programming," in *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO '14*. ACM, 2014, pp. 249–256.
- [18] E. Hart, K. Sim, B. Gardiner, and K. Kamimura, "A hybrid method for feature construction and selection to improve wind-damage prediction in the forestry sector," in *Proceedings of the Genetic and Evolutionary Computation Conference*, 2017, pp. 1121–1128.
- [19] J. A. Lee and M. Verleysen, *Nonlinear dimensionality reduction*. Springer Science & Business Media, 2007.
- [20] I. T. Jolliffe, "Principal component analysis," in *International Encyclopedia of Statistical Science*. Springer, 2011, pp. 1094–1096.
- [21] J. B. Tenenbaum, V. d. Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, vol. 290, no. 5500, pp. 2319–2323, 2000.
- [22] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, vol. 290, no. 5500, pp. 2323–2326, 2000.
- [23] G. E. Hinton and S. T. Roweis, "Stochastic neighbor embedding," in *Advances in Neural Information Processing Systems 15 [Neural Information Processing Systems, NIPS 2002, December 9-14, 2002, Vancouver, British Columbia, Canada]*, 2002, pp. 833–840.
- [24] L. van der Maaten, "Accelerating t-SNE using tree-based algorithms," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 3221–3245, 2014.
- [25] Q. Zhang and H. Li, "MOEA/D: A multiobjective evolutionary algorithm based on decomposition," *IEEE Trans. Evolutionary Computation*, vol. 11, no. 6, pp. 712–731, 2007.
- [26] H. Li and Q. Zhang, "Multiobjective optimization problems with complicated pareto sets, MOEA/D and NSGA-II," *IEEE Trans. Evolutionary Computation*, vol. 13, no. 2, pp. 284–302, 2009.
- [27] U. Bhowan, M. Johnston, M. Zhang, and X. Yao, "Evolving diverse ensembles using genetic programming for classification with unbalanced data," *IEEE Trans. Evolutionary Computation*, vol. 17, no. 3, pp. 368–386, 2013.
- [28] I. Icke and A. Rosenberg, "Multi-objective genetic programming for visual analytics," in *Proceedings of the European Conference on Genetic Programming (EuroGP)*, 2011, pp. 322–334.
- [29] S. Nguyen, M. Zhang, D. Alahakoon, and K. C. Tan, "Visualizing the evolution of computer programs for genetic programming [research frontier]," *IEEE Computational Intelligence Magazine*, vol. 13, no. 4, pp. 77–94, 2018.
- [30] Y. Sun, G. G. Yen, and Z. Yi, "Evolving unsupervised deep neural networks for learning meaningful representations," *IEEE Trans. Evolutionary Computation*, vol. 23, no. 1, pp. 89–103, 2019.
- [31] T. Kohonen, "Self-organized formation of topologically correct feature maps," *Biological Cybernetics*, vol. 43, no. 1, pp. 59–69, Jan 1982.
- [32] L. G. Nonato and M. Aupetit, "Multidimensional projection for visual analytics: Linking techniques with distortions, tasks, and layout enrichment," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 8, pp. 2650–2673, Aug 2019.
- [33] R. Faust, D. Glickenstein, and C. Scheidegger, "Dimreader: Axis lines that explain non-linear projections," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 481–490, Jan 2019.
- [34] R. Poli and N. F. McPhee, "Parsimony pressure made easy: Solving the problem of bloat in GP," in *Theory and Principled Methods for the Design of Metaheuristics*, 2014, pp. 181–204.
- [35] J. Kennedy and R. C. Eberhart, "Particle swarm optimization," in *Proceedings of the IEEE International Conference on Neural Networks*, vol. 4. IEEE, 1995, pp. 1942–1948.
- [36] N. Hansen, S. D. Müller, and P. Koumoutsakos, "Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (CMA-ES)," *Evolutionary Computation*, vol. 11, no. 1, pp. 1–18, 2003.
- [37] N. Hansen, R. Ros, N. Mauny, M. Schoenauer, and A. Auger, "Impacts of invariance in search: When CMA-ES and PSO face ill-conditioned and non-separable problems," *Appl. Soft Comput.*, vol. 11, no. 8, pp. 5755–5769, 2011.
- [38] D. Dheeru and E. Karra Taniskidou, "UCI machine learning repository," 2017. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [39] D. E. Elder, *Lever's Histopathology of the Skin*, eleventh ed. Wolters Kluwer, 2015.



Andrew Lensen (M'17) received the B.Sc., B.Sc.(Hons 1st class), and Ph.D. degrees in computer science from Victoria University of Wellington, Wellington, New Zealand, in 2015, 2016, and 2019, respectively.

He is currently a Post-doctoral Research Fellow in the Evolutionary Computation Research Group within the School of Engineering and Computer Science at Victoria University of Wellington. His current research interests are mainly in the use of evolutionary computation for feature manipulation

in unsupervised learning, with a particular focus on the use of genetic programming for manifold learning, clustering, and feature synthesis.

Dr Lensen serves as a regular reviewer of several international conferences, including IEEE Congress on Evolutionary Computation, and international journals such as the IEEE Transactions on Evolutionary Computation, and the IEEE Transactions on Cybernetics.



Bing Xue (M10) received the B.Sc. degree from the Henan University of Economics and Law, Zhengzhou, China, in 2007, the M.Sc. degree in management from Shenzhen University, Shenzhen, China, in 2010, and the PhD degree in computer science in 2014 at Victoria University of Wellington, New Zealand.

She is currently an Associate Professor in the School of Engineering and Computer Science at Victoria University of Wellington. She has over 200 papers published in fully refereed international journals and conferences and her research focuses mainly on evolutionary computation, feature selection, feature construction, image analysis, and transfer learning.

Dr Xue is currently the Chair of IEEE Computational Intelligence Society (CIS) Data Mining and Big Data Analytics Technical Committee, Chair of the IEEE Task Force on Evolutionary Feature Selection and Construction, and Vice-Chair of IEEE CIS Task Force on Transfer Learning & Transfer Optimization, and Vice-Chair of IEEE CIS Task Force on Evolutionary Deep Learning and Applications.



Mengjie Zhang (M04SM10-F'19) received the B.E. and M.E. degrees from the Artificial Intelligence Research Center, Agricultural University of Hebei, Hebei, China, and the Ph.D. degree in computer science from RMIT University, Melbourne, Australia, in 1989, 1992, and 2000, respectively.

He is currently Professor of Computer Science, Head of the Evolutionary Computation Research Group, and the Associate Dean (Research and Innovation) in the Faculty of Engineering. His current research interests include evolutionary computation,

particularly genetic programming, particle swarm optimization, and learning classifier systems with application areas of image analysis, multi-objective optimization, feature selection and reduction, job shop scheduling, and transfer learning. He has published over 500 research papers in refereed international journals and conferences.

Prof. Zhang is a Fellow of the Royal Society of New Zealand and have been a Panel member of the Marsden Fund (New Zealand Government Funding), a Fellow of IEEE, and a member of ACM. He was the chair of the IEEE CIS Intelligent Systems and Applications Technical Committee, and chair for the IEEE CIS Emergent Technologies Technical Committee and the Evolutionary Computation Technical Committee, and a member of the IEEE CIS Award Committee. He is a vice-chair of the IEEE CIS Task Force on Evolutionary Feature Selection and Construction, a vice-chair of the Task Force on Evolutionary Computer Vision and Image Processing, and the founding chair of the IEEE Computational Intelligence Chapter in New Zealand. He is also a committee member of the IEEE NZ Central Section.