

**Substructural Surrogates for Learning
Decomposable Classification Problems:
Implementation and First Results**

**Albert Orriols-Puig, Kumara Sastry,
David E. Goldberg, and Ester Bernadó-Mansilla**
IlliGAL Report No. 2007010
March 2007

Illinois Genetic Algorithms Laboratory
University of Illinois at Urbana-Champaign
117 Transportation Building
104 S. Mathews Avenue Urbana, IL 61801
Office: (217) 333-2346
Fax: (217) 244-5705

Substructural Surrogates for Learning Decomposable Classification Problems: Implementation and First Results

Albert Orriols-Puig^{1,2}, Kumara Sastry¹,
David E. Goldberg¹ and Ester Bernadó-Mansilla²

¹Illinois Genetic Algorithm Laboratory (IlliGAL)
Department of Industrial and Enterprise Systems Engineering
University of Illinois at Urbana-Champaign
Urbana, IL 61801

²Group of Research in Intelligent Systems
Computer Engineering Department
Enginyeria i Arquitectura La Salle — Ramon Llull University
Quatre Camins, 2. 08022, Barcelona, Catalonia, Spain.

aorriols@salle.url.edu, deg@uiuc.edu, ksastry@uiuc.edu, esterb@salle.url.edu

March 26, 2007

Abstract

This paper presents a learning methodology based on a substructural classification model to solve decomposable classification problems. The proposed method consists of three important components: (1) a structural model that represents salient interactions between attributes for a given data, (2) a surrogate model which provides a functional approximation of the output as a function of attributes, and (3) a classification model which predicts the class for new inputs. The structural model is used to infer the functional form of the surrogate and its coefficients are estimated using linear regression methods. The classification model uses a maximally-accurate, least-complex surrogate to predict the output for given inputs. The structural model that yields an optimal classification model is searched using an iterative greedy search heuristic. Results show that the proposed method successfully detects key variable interactions in hierarchical problems, group them in linkages groups, and build maximally accurate classification models. The initial results on non-trivial hierarchical test problems indicate that the proposed method holds promise and have also shed light on several improvements to enhance the capabilities of the proposed method.

1 Introduction

Nearly decomposable functions play a central role in the design, analysis and modeling of complex engineering systems (25; 6; 7). A design decomposition principle has been proposed for the successful design of scalable genetic algorithms (GAs) (7; 16; 18), genetic programming (24), and learning classifier systems and genetics based machine learning (GBML) (3; 14). For example, in (3), estimation of distribution algorithms (EDAs) were applied over the rule-based knowledge

evolved by XCS (29; 30) to discover linkages between the input variables, permitting XCS to solve hierarchical problems which were intractable with first-generation XCS.

However, previous approaches used the probabilistic models built by EDAs—GAs that replace variation operators by building and sampling probabilistic models of promising solution—for recombination. However, the probabilistic models can also be used to induce the form of surrogates which can be used for efficiency enhancement of GAs (23; 17; 22) and GBML (15). In this paper we develop use the substructural surrogates for learning from decomposable problems with nominal attributes. Similar to Sastry, Lima, and Goldberg (22), we use the structural model of EDAs to induce the form of the surrogate and linear regression for estimating the coefficients of the surrogate. The surrogate is subsequently used to predict the class of unknown input instances.

In this paper, we discuss the critical components of the proposed methodology and outline several ways to implement it. We then propose a greedy search heuristic for discovering the structural model that minimizes the test error of the classification model constructed from it. We address this method as *greedy Extraction of the Structural Model for Classification* (gESMC). We artificially design a set of hierarchical problems by means of concatenating essential blocks which output, provided by a boolean function, serves as the input of another function that determines the global output of the example. Thus, these problems may be decomposed and essential blocks should be correctly processed to predict the correct output. gESMC is able to detect the interactions between variables and build accurate classification models. Moreover, the system is compared to C4.5 and SMO, which clearly shows the advantage of extracting the problem structure. Finally, we review the limitations of applying a greedy search to obtain the best structural model, show in which circumstances these limitations may appear, and propose approaches to overcome them.

The paper is organized as follows. Section 2 discusses the proposed methodology followed by a description of gESMC. The test problems designed and used in this study are discussed in Section 4. Section 5 compares the results of gESMC with C4.5 and SMO on the hierarchical problems. Section 6 discusses some enhancements that are yet to be investigated. Section 7 provides summary and conclusions.

2 Methodology for Learning χ -ary input Problems

In this section, we discuss a methodology for learning the structural and the classification model from a set of labeled examples. The methodology consists of three layers: (1) the *structural model layer*, (2) the *surrogate model layer*, and (3) the *classification model layer*. The *structural model layer* extracts the dependencies between the attributes of the examples in the dataset. These dependencies can be expressed in form of linkage groups (8; 9), matrices (32), or Bayesian networks (16). However the dependencies are represented, the key idea is that the salient interactions between attributes are used as a basis for determining the output. The *surrogate model layer* uses the *structural model* to infer the functional form of the surrogate and the coefficients of the surrogate are determined using linear regression methods. The resulting surrogate is a function that approximates the output of each input instance. Finally, the *classification model layer* uses the surrogate function to predict the class of new input instances.

In essence, we infer the structure of the surrogate from the structural models of attribute interactions and then use linear regression methods to estimate value of the coefficients (or the partial contributions of subsolutions to the output) of the resulting surrogate function. Finally the surrogate is used to predict the class of new input instances. Details of each of the three components is discussed in the following sections.

2.1 Structural Model Layer

The *structural model layer* is the responsible for identifying salient interactions between attributes (for example, linkage groups) which need to be processed together to determine their contribution to the output. For example, consider a problem with two binary attributes (x_1, x_2) and whose output is determined by the *x-or* boolean function. If we considered each of the attribute independently, we cannot evolve a function that computes the output accurately for all possible inputs. However, when we consider the two attributes together, we can easily create a function that can accurately predict the output for all possible input sequences.

A number of linkage-learning methods (7) can be used to implement the structural model layer. Here, we will use estimation of distribution algorithms (EDAs) (16; 18), which learn the salient interactions between decision variables by building probabilistic models of promising candidate solutions. In learning classifier systems (LCSs) or genetics-based machine learning (GBML) realm, EDAs have been successfully combined with LCSs to extract the linkages between classifiers' alleles (3; 14; 15). However, unlike previous studies which used the structural model as a replacement of recombination, in this study we integrate the structural model and learning with the use of substructural surrogates.

In order to achieve this integration, the first step is to find the structural model of the given data. This can be done in several ways. As with EDAs, given a class of permissible structural models, we can search for the best structural model. Prior and domain-specific knowledge can also be used to propose the structural model, and a search mechanism could be used to refine it (1). In this study we use a greedy search heuristic that searches for the model structure that results in the most accurate surrogate model, detail of which are given in Section 3.

2.2 Surrogate Model Layer

The *surrogate model layer* preprocesses the input examples according to the structural model and builds a regression model from these preprocessed examples, as described in (22). In this section we summarize the procedure of building such a surrogate. Consider a matrix D of dimension $n \times \ell$ that contains all the input examples (where n is the number of examples and ℓ the number of attributes). Once the structural model is built, every linkage group is treated as a *building block* (10). Then, we consider all possible input combinations within in each linkage group to process the input examples.

For example, consider the following structural model of a binary problem of 3 variables: $\{[x_1, x_3], [x_2]\}$. That is, there is salient interaction between variables x_1 and x_3 , which are independent from the variable x_2 . In this case, we consider the following schemata: $\{0*0, 0*1, 1*0, 1*1, *0*, *1*\}$. In general, given m linkage groups, the total number of schemata m_{sch} to be considered is given by:

$$m_{sch} = \sum_{i=1}^m \left[\prod_{j=1}^{k_i} \chi_{i,j} \right], \quad (1)$$

where $\chi_{i,j}$ is the alphabet cardinality of the j^{th} variable of the i^{th} linkage group, and k_i is the size of the i^{th} linkage group.

Then, each example in D is mapped to a vector of size m_{sch} , creating the matrix A of dimensions $n \times m_{sch}$:

$$\mathbf{A} = \begin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,m_{sch}} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,m_{sch}} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1} & a_{n,2} & \cdots & a_{n,m_{sch}} \end{pmatrix}, \quad (2)$$

where $a_{i,j}$ will have value '1' if the i th example belongs to the j th schemata, and '0' otherwise. Note that, given an example, only one of the schemata for each linkage group can have value '1'.

We map different labels or classes of the examples to numeric values. For example: $\{class_1, class_2, \dots, class_k\} \longrightarrow \{\mathbb{Z}_1, \mathbb{Z}_2, \dots, \mathbb{Z}_k\}$. The label or class c_i of each example is also kept in a matrix C of dimensions $n \times 1$:

$$\mathbf{C} = \begin{pmatrix} c_1 & c_2 & \vdots & c_n \end{pmatrix}^T. \quad (3)$$

Now, the task of designing the surrogate can be formulated into a linear system of equations and our objective is to compute the coefficients of the matrix \mathbf{x} of dimensions $m_{sch} \times 1$ that satisfy the following equality:

$$\mathbf{Ax} = \mathbf{C}. \quad (4)$$

In practice, we may not find an \mathbf{x} that satisfies this expression. For this reason, we use a multi-dimensional least squares fitting approach. That is, the problem is reformulated by estimating the vector of coefficients \mathbf{x} that minimize the square error function χ :

$$\chi^2 = (\mathbf{Ax} - \mathbf{C})^T \cdot (\mathbf{Ax} - \mathbf{C}). \quad (5)$$

The problem of least-squares fitting is well-known, and so we do not provide insight in the resolution methodology herein. The interested reader is referred to (5; 21). Here, we used the multi-dimensional least squares fitting routine available with GNU scientific library ¹ (GSL).

2.3 Classification Model Layer

Once we obtain the matrix x with the regression coefficients, the output for a new example is computed as follows. The example is mapped to a vector \vec{e} of size m_{sch} . The mapping procedure used is identical to that used to create matrix \mathbf{A} and as outlined in the previous section, the elements of \vec{e} will have a value '1' if the example belongs to the corresponding schemata and '0' otherwise. Then, the predicted output is given by:

$$output = \vec{e} \cdot x. \quad (6)$$

Note that the *output* is a continuous value, and has to be transformed to one of the possible class labels. Therefore, we convert the continuous output to the closer integer \mathbb{Z}_i in $\{\mathbb{Z}_1, \mathbb{Z}_2, \dots, \mathbb{Z}_k\}$, and then, return the class label that corresponds to \mathbb{Z}_i .

In essence, the proposed method relies on the structural and the surrogate models extracted from the data to build the classification model. Therefore, we note that if this structural model does not reflect the variable interactions accurately, the accuracy of the classification model will be limited. Therefore, a critical task for the success of the proposed methodology is our ability to find reasonably accurate structural models. In the next section we propose an implementation of the methodology that search iteratively for the best structural model, and uses the classification model to evaluate its quality. We call this implementation as the *greedy extraction of the structural model for classification* (gESMC).

3 Implementing the Methodology: gESMC

The pseudocode of the implementation of our proposed method is shown in Algorithm 1. In the initialization stage, the algorithm divides the data into training and test sets. We start with a

¹<http://www.gnu.org/software/gsl>

Algorithm 1: Building of structural and classification model via a greedy search.

Data: *dataset* is the set of labeled examples.

Result: *function* is the classification function and *bestModel* the structural model.

```
1 begin
2   [train, test]  $\leftarrow$  divideData ( data )
3   bestModel  $\leftarrow$  [1], [2], ..., [n]
4   function  $\leftarrow$  createSurrogateFunction ( bestModel, train[i] ) ▷ See Sect. 2.2
5   mdl  $\leftarrow$  evaluateModel ( bestModel, test[i] )
6   isImproving  $\leftarrow$  true
7   while mdl >  $\theta$  and isImproving do
8     for  $i \in \{1, \dots, \text{length}(\text{bestModel}) - 1\}$  do
9       for  $j \in \{i + 1, \dots, \text{length}(\text{bestModel}) - 1\}$  do
10        newModel[count]  $\leftarrow$  joinLinkages( bestModel, i, j)
11        newFunction[count]  $\leftarrow$  createSurrogateFunction ( newModel[count], train )
12        newMdl[count]  $\leftarrow$  evaluateModel ( newModel[count], test[count] )
13        count  $\leftarrow$  count + 1
14      end
15    end
16    best  $\rightarrow$  position min. mdl( newMdl ) ▷ Selects the best model
17    if newMdl[best] significantly improves mdl then
18      | bestModel = newModels[i]
19      | mdl = newMdl[i]
20    else
21      | isImproving=false
22    end
23  end
24 end
```

structural model where all variables are treated as independent and build a surrogate function via regression over the training set as explained in the previous section (see Section 2.2). The quality of the classification model is evaluated with the test set and stored in the variable *mdl*.

Similar to the extended compact genetic algorithm (eCGA) (9), in gESMC we use a greedy search heuristic to partition the set of attributes into non-overlapping clusters such that the classification error is (locally) minimized. That is, starting from a model where variables are treated as independent, we continue to merge substructures till either (1) the *mdl* measure becomes less than a user set threshold of θ , or (2) the search produces no improvement. In every iteration of the inner loop (lines 10 to 13), we merge two linkage groups from the current best model, create the surrogate and the classification model, and evaluate it. That is, $\binom{m}{2}$ new structural models are formed (where m is the number of substructures in the current best model), and their surrogate functions created and evaluated. Among the evaluated $\binom{m}{2}$ models, the one with the lowest classification error is chosen as the current best model for the next iteration if it *significantly improves* the current best model; otherwise, we terminate the search, and the current best surrogate and classification models are returned as the (locally) best models.

Three elements of the implementation need further explanation: (1) procedure to divide the data into training and test sets (line 2), (2) evaluation of the model (lines 5 and 12), and (3) procedure for comparing two models and choosing the better one (line 17). Each of the three

elements are discussed in the following paragraph.

Partition of the data. The procedure used to partition the data into training and test sets affects the estimation of classification error. A number of approaches such as holdout validation, k -fold cross validation, and leave-one-out cross-validation methods can be used. Here, we use a k -fold cross validation (26) with $k = 10$.

Evaluation of the model. The quality of a structural model depends on (1) the complexity of this model, and (2) the test error of the classification model created from it. Again a number of measures such as minimum description length metrics and multiobjective approaches could be used to measure the relative quality of a given surrogate and classification model. We use the k -fold cross validation which provides a measure of both the test error and the model complexity in terms of overfitting of training data. That is, if the structural model is more complex than necessary, the surrogate function will tend to overfit the training instances, and the test error increases.

Comparison of models. Given a current-best model, in gESMC we consider all pairwise merges of the substructures of the current-best model. We need to choose the best model among all the models created via the pairwise merges and compare it to the current-best model. best model have to be compared. Again, this could be done in a number of ways. For example, we could accept the new model if its classification error is lower than that of the current-best model. However, this might lead to spurious linkages and more complex model might be accepted, especially if the data set is noisy. Therefore to avoid getting unnecessarily complicated structural models, we say that a model m_1 is significantly better than a model m_2 if:

$$error_{m_1} < error_{m_2} - \delta, \tag{7}$$

where δ is a user-set threshold. Alternatively, we can use different statistical tests as well. In our implementation, we use a paired t-test to determine if a new, more complex, structural model is better than the current best model (26). We use a confidence interval of $\alpha = 0.01$.

Before proceeding with a description of the test functions, we note two important properties of gESMC. First, in the current implementation gESMC the structural model is a partition of the variables into non-overlapping groups. However, this limitation can easily be relaxed by using other structural models (32; 16). Second, because of the greedy procedure, we need some guidance from lower-order, sub-optimal structural models towards an optimal structural model. This limitation can be alleviated by replacing the greedy search heuristic with a global optimization method.

4 Test Problems

In this section, we present a general set of decomposable problems to investigate the capabilities of gESMC in correctly identifying salient substructures and building an accurate classification model. As illustrated in Figure 1, we propose a class of two-level hierarchical problems where the lower level consists of functions that operate on a set of binary-input blocks and the upper level consists of functions that operate on the lower-level function values to produce the output. The lower- and upper-level function used in the study are explained in Section 4.1 and 4.2, respectively.

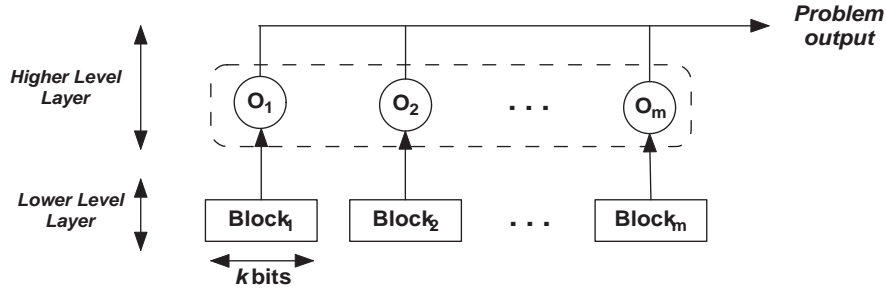


Figure 1: Example of the design of a two-level hierarchical problem. In the low level layer, m blocks of k bits are concatenated. Each block is evaluated resulting in the correspondent output. All the outputs are groups in an input string that is used to determine the global output.

4.1 Lower Level of the Hierarchy

At the lower level of the hierarchy, we considered the following two binary functions which operate independently on m blocks with k variables in each block. Moreover the variables within a block interact with each other and determine the output of the function.

The position problem. The *position* problem (2) is defined as follows. Given a binary input of length ℓ , the output is the position of the left-most one-valued bit. For example, $f(\underline{1}00)=3$, $f(0\underline{1}0)=2$, $f(00\underline{1})=1$, and $f(000)=0$. Note that every variable is linked to all the variables on its left. Therefore, all the variables need to be processed together to determine the output.

The parity problem. The *parity* problem (13) is a two-class binary problem defined as follows. Given a binary input of length ℓ , the output is the number of one-valued bits modulo two. For example, $f(110)=0$, $f(100)=1$, and $f(111)=1$. To predict the output accurately for the parity problem, all the variables have to be jointly processed. Additionally, for a k -bit parity problem, the structural model that represents all the variables to be independent yields a classification model of the same accuracy as the one that contains substructures of size $k - 1$ or less. That is till we get a structural model that groups all k variables together, the accuracy of the classification model does not increase.

4.2 Higher Level of the Hierarchy

At the higher level of the hierarchy, we use the following problems where each variable contributes independently to the output. That is, the structural information is contained in the lower-level of the hierarchy and the upper level function affects the salience of the substructures.

The decoder problem. The decoder problem (2) is a binary-input multi-class problem defined as follows. Given an input of length ℓ , the output is determined by the decimal value of the input. For example, $f(111) = 7$, $f(101) = 5$, and $f(000) = 0$. Note that each variable independently contributes to the output. That is, starting with class equal to zero, a '1' in i^{th} position adds 2^i to the output, irrespective of other variable values.

The count-ones problems. The *count-ones* is defined as follows. Given a binary input of size ℓ , the output is the number of one-valued bits. Again, the output of the count-ones problems can be predicted by treating the input variables independently.

As mentioned earlier, we concatenated m blocks of k bits of the two lower level problems with the two higher level problems to create four different hierarchical test problems. Specifically, we used the *position* at the lower level with the *decoder* (HPosDec) and the *count-ones* (HPosCount) in the higher level. Similarly, low order *parity* blocks were combined again with the *decoder* (HParDec) and the *count-ones* (HParCount). Additionally, we added some irrelevant bits, which do not contribute to the output, to see if our method was capable of ignoring them. Therefore, in our case, $\ell \geq m \cdot k$, where ℓ is the length of the input string.

With the above description of the test problems, we present the results of gESMC and compare with those of C4.5 and SMO in the following section.

5 Results

This section analyzes the behavior of gESMC for learning hierarchical problems, and compare the results to those obtained with two highly competitive learners in terms of performance and interpretability.

5.1 Experimental Methodology

We use the four hierarchical problems designed in the previous section to analyze the performance of gESMC. We start with concatenations of three minimum-order blocks in the lower level hierarchy (that is, $k=2$) and add 9 irrelevant bits to the input. Our aim is to analyze the capabilities of gESMC in (i) identifying salient substructures of interacting variables, and (ii) ignoring irrelevant variables, and not considering them in the classification model. Next, we increase the order of the lower-order blocks with a two-fold objective. For the problem with position blocks, we analyze if the system is able to identify and efficiently handle larger groups of linked variables. For the problems with parity blocks, we want to investigate the behavior of gESMC when there is a lack of guidance toward an accurate substructural model.

To illustrate the need for detecting linkage groups in classification tasks, we compare the results obtained with gESMC to those of two widely used learners: C4.5 (20), and SMO (19). C4.5 is a decision tree, derived from ID3, which has been widely used because of its ability to tackle a wider range of problems and because of the interpretability of the extracted knowledge. SMO is a *support vector machine* (27) that implements the *Sequential Minimal Optimization* algorithm. Although the interpretability is more difficult since it represents the knowledge as function weights, its competence has been demonstrated in different kind of problems. Both methods were run using WEKA (31). Unless otherwise noted, for C4.5 we used the default configuration, and for SMO, we used a polynomial kernel of order 1.

The three methods are compared in terms of performance (that is, test accuracy) and comprehensibility of the knowledge generated by the learner. As the datasets had a large number of instances, we used the *holdout* methodology² to estimate the test accuracy; that is, 70% of the instances were randomly selected and placed in the training set, and the rest formed the test set. To compare the performance of each pair of learners on a given problem, we applied a paired Student t-test (26) on the results. To study the interpretability of each method, we qualitatively compared the structural and the classification model evolved by gESMC to the the decision trees generated by C4.5, and the weights extracted by SMO.

²A holdout is the simplest cross-validation approach where the data is divided in two sets, the train and the test set.

Table 1: Test error and standard deviation obtained with gESMC, SMO and C4.5 on the problems HPosDec, HPosCount, HParDec, HParCount with $\ell=15$, $m = 3$, and $k = 2$. Results are averaged over ten runs with different holdouts and random seeds.

	gESMC	C4.5	SMO
<i>HPosDec</i>	0.00% \pm 0.00%	0.00% \pm 0.00%	0.00% \pm 0.00%
<i>HPosCount</i>	0.00% \pm 0.00%	0.00% \pm 0.00%	21.89% \pm 0.13%
<i>HParDec</i>	0.00% \pm 0.00%	3.32% \pm 2.90%	89.11% \pm 0.94%
<i>HParCount</i>	0.00% \pm 0.00%	5.15% \pm 4.43%	62.72% \pm 0.27%

5.2 Results with 2-bit Low Order Blocks

We first run gESMC, C4.5, and SMO on the problems HPosDec, HPosCount, HParDec, and HParCount with $\ell = 15$, $m = 3$ and $k = 2$. Therefore, the problems were formed by three lower level blocks of two bits, and there were 9 irrelevant bits at the end of the binary input. Next, we compare the results in terms of performance and interpretability.

5.2.1 Comparison of the Performance

Table 1 summarizes the test error resulting of applying gESMC, C4.5, and SMO on the four hierarchical problems. All the results were averaged over ten runs with different holdouts and random seeds.

The results show that gESMC obtained 0% test error for all the problems tested. This indicates that the method is able to process the variable linkages and build maximally accurate classification models. None of the other learners could achieve 0% error in all the problems. C4.5 achieved 0% test error for the problems HPosCount and HPosDec, the ones formed by position blocks. Nonetheless, on the problems that consist of parity blocks, C4.5 was significantly outperformed by gESMC according to a a paired t-test on a significance level of 0.99. Finally, SMO presents the worst behavior of the comparison. The learner could accurately generalize over the input data only on the HPosDec problem. For the problems HPosCount, HParDec, and HParCount, the results of SMO significantly degraded those obtained with gESMC and C4.5. Note the big difference in the test errors; for HParDec, SMO has 89.11% test error, C4.5 has 3.32%, and gESMC is maximally accurate. We repeated the experiments with a Gaussian kernel (11) to promote the discovering of the linkage groups, but no significant improvement was found.

These results highlight the importance of learning and incorporating the structural model into the classification model. gESMC found highly accurate classification models only after discovering the problem structure (examples of some structural and classification models are shown in the next section). However C4.5 and SMO failed since they were not able to identify this structure. Note that the problems formed by parity blocks resulted more problematic for both learners than the problems based on position blocks. This could be explained as follows. The variables linkages in the position are weaker than in the parity. That is, in the position problem every variable processed from left to right reduces the uncertainty of the output. In the parity, looking at a single variable does not reduce the uncertainty, and so, processing the linkages is crucial. We hypothesize that, for this reason, problems formed by parity are more difficult to learn for C4.5 and SMO.

Table 2: Structural models and surrogate functions build by gESMC for the problems HPosDec, HPosCount, HParDec, HParCount with $\ell=15$, $m = 3$, and $k = 2$.

HPosDec	<i>link. groups</i>	$[x_0x_1][x_2x_3][x_4][x_5][x_6][x_7][x_8][x_9][x_{10}][x_{11}][x_{12}][x_{13}][x_{14}]$
	<i>surr. func.</i>	$16.75 + 9(1 - \bar{x}_0x_1) - 6\bar{x}_2\bar{x}_3 - 3\bar{x}_2x_3 - 1.5x_4 + 0.5x_5$
HPosCount	<i>link. groups</i>	$[x_0x_1][x_2x_3][x_4][x_5][x_6][x_7][x_8][x_9][x_{10}][x_{11}][x_{12}][x_{13}][x_{14}]$
	<i>surr. func.</i>	$0.75 + \bar{x}_0x_1 + (1 - \bar{x}_2x_3) + 0.5\bar{x}_4 + 0.5x_5$
HParDec	<i>link. groups</i>	$[x_0x_1][x_2x_3][x_4x_5][x_6][x_7][x_8][x_9][x_{10}][x_{11}][x_{12}][x_{13}][x_{14}]$
	<i>surr. func.</i>	$2 + 4(\bar{x}_0x_1 + x_0\bar{x}_1) - 2(\bar{x}_2\bar{x}_3 + x_2x_3) - (\bar{x}_4x_5 + x_4\bar{x}_5)$
HParCount	<i>link. groups</i>	$[x_0x_1][x_2x_3][x_4x_5][x_6][x_7][x_8][x_9][x_{10}][x_{11}][x_{12}][x_{13}][x_{14}]$
	<i>surr. function</i>	$\bar{x}_0x_1 + x_0\bar{x}_1 + \bar{x}_2x_3 + x_2\bar{x}_3 + \bar{x}_4x_5 + x_4\bar{x}_5$

5.2.2 Comparison of the Interpretability

We now analyze the interpretability of the models created by gESMC, and qualitatively compare them to those obtained by C4.5 and SMO. Table 2 shows the structural models and the associated surrogate functions built for each problem. For HPosDec and HPosCount, gESMC correctly detects the linkages between the groups of variables $[x_0, x_1]$ and $[x_2, x_3]$; all the other variables are considered independent. The reason that gESMC incorrectly identifies that variables x_4 and x_5 are independent, is because it reaches the termination criteria of 0% test error. For the problems HParDec and HParCount, gESMC discovers the linkage groups $[x_0, x_1]$, $[x_2, x_3]$, and $[x_4, x_5]$. Differently from the position problem, now gESMC needs to discover all the existing parity groups to remove the uncertainty, and so, build the most accurate classification model.

The availability of the structural model with gESMC is another advantage over other conventional classification techniques in terms of interpretability. The structural model facilitates easy visualization of the salient variable interactions and better understand the resulting surrogate function. Note that for all the problems, gESMC built easily interpretable functions and also efficiently ignores irrelevant variables. quite comprehensible functions. For example, consider the problem HParCount, in which the output is the number of '1s' resulting from the evaluation of each low-order parity block. The function evolved clearly indicates that if any of the linkage groups has the schemata '01' or '10' (values from which the parity would result in '1'), the output is incremented by one.

Let us now compare this knowledge representation to those obtained with C4.5 and SMO. For this purpose, we consider the size of the trees built by C4.5, and the machines constructed by SMO. For HPosDec and HPosCount, C4.5 built a tree with 53 nodes, from which 27 were leaves. The resulting trees specified the output for each combination of the six relevant bits (see a portion of a tree for HPosDec in Fig. 2). Although these trees detail the output given the value of the first six variables, they do not show the variable interactions. For HParCount, C4.5 built trees that, on average, had 136 leaves and 270 nodes. For HParDec, the trees had 142 leaves and 283 nodes. These high number of nodes makes the interpretability of tree very hard. Additionally, all the trees had some irrelevant attributes in the decision nodes. That is, C4.5 was overfitting the train instances to reduce the train error, resulting in more complicated trees, further hindering the interpretability of the classification model.

In contrast to gESMC and C4.5, SMO presented the less interpretable results. In general, SMO creates a machine for each pair of classes, and adjust $\ell + 1$ weights for each machine (where ℓ is the number of attributes of the problem). For HPosDec, SMO built 351 machines with 16 weights ranging from 0 to 1. For HPoscount, HParDec, and HParCount, 6, 6, and 28 machines

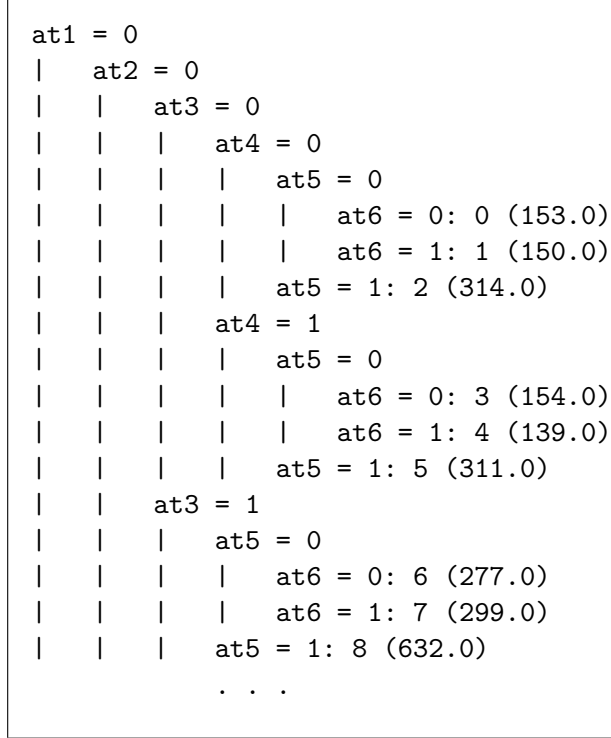


Figure 2: Portion of the tree built by C4.5 for the HPosDec problem.

were created, all them with 16 weights ranging from 0 to 1. Although some of these weights were zero, the machines evolved could not be interpreted at all. Thus, the human expert would not be able to extract any information from these knowledge models.

Although both SMO and gESMC represent the knowledge in terms functions that partition the search space, gESMC yields the structural model which permit easy visualization of salient variable interactions. Additionally, while SMO weights the input variables, while gESMC weights the different subsolution of the identified substructures. Classification models obtained via gESMC show the relative influence of subsolutions to the output and therefore is more easily interpretable than SMO.

5.3 Results Increasing the Low Level Block Size

We now increase the interaction order of the lower-level blocks to analyze the effect of order of interaction on the performance of gESMC. Additionally, for the test problems with parity blocks, we also want to investigate the effect of having no guidance from lower-order substructures towards obtaining accurate structural model on the accuracy of the classification model obtained via gESMC. Specifically, we use HPosDec, HPosCount, HParDec, and HParCount with $\ell = 15$, $m = 2$, and $k = 3$, and compared gESMC to C4.5 and SMO.

Table 3 shows the test errors for gESMC, C4.5 and SMO. For HPosDec and HPosCount, gESMC obtained 0% error test and for both problems, two different structural models were created during independent runs:

$$\begin{aligned}
Model_1 &: [x_0x_1][x_3x_4x_5][x_2][x_6][x_7][x_8][x_9][x_{10}][x_{11}][x_{12}][x_{13}][x_{14}], \\
Model_2 &: [x_0x_1x_2][x_3x_4][x_5][x_6][x_7][x_8][x_9][x_{10}][x_{11}][x_{12}][x_{13}][x_{14}].
\end{aligned}$$

Table 3: Test error and standard deviation obtained with gESMC, SMO and C4.5 on the problems HPosDec, HPosCount, HParDec, HParCount with $\ell=15$, $m = 3$, and $k = 3$. Results are averages over ten runs with different holdouts and random seeds.

	gESMC	C4.5	SMO
<i>HPosDec</i>	0.00% \pm 0.00%	0.00% \pm 0.00%	0.00% \pm 0.00%
<i>HPosCount</i>	0.00% \pm 0.00%	0.00% \pm 0.00%	14.24% \pm 1.03%
<i>HParDec</i>	24.00% \pm 25.50%	9.01% \pm 5.86%	76.91% \pm 2.01%
<i>HParCount</i>	49.99% \pm 0.00%	12.25% \pm 6.69%	49.94% \pm 0.22%

In both the above models, the variables of one of the lower-level blocks is correctly identified, and only two variables of the other lower-level blocks form a linkage group. Again as observed in the previous section, because gESMC meets the convergence criteria of 0% test errors for both the models even when one of the substructures is partially identified.

The surrogate functions evolved are qualitatively similar to those obtained in the previous section. In all cases, only the six relevant variables were taken in consideration, and specifically, some of their schemata. For example, one of the surrogate functions created for the HPosDec is

$$14.9\bar{x}_5 + 15x_5 - 2.5\bar{x}_3\bar{x}_4 + \bar{x}_3x_4 - 12\bar{x}_0\bar{x}_1\bar{x}_2 - 4\bar{x}_0x_1\bar{x}_2 - 8\bar{x}_0\bar{x}_1x_2 - 4\bar{x}_0x_1x_2, \quad (8)$$

which only contains the six relevant variables x_0, x_1, \dots, x_5 .

As expected, for HParDec and HParCount, gESMC yields poorer results. For HParDec, the average test error was 24% with a high standard deviation. This high deviation is because gESMC yielded a maximally accurate classification model for 50% of the runs. For the rest 50%, gESMC could not discover accurate structural model. Further investigation showed that this is due to the stochasticity of the holdout estimation. Since we randomly select 70% of the instances as the training set, the symmetry of the parity problem may be broken leading the greedy search heuristic to yield the accurate structural model. This indicates that introduction of stochasticity might break symmetry of parity-like functions and render the accurate structural model hill-climbable. However, the efficacy of adding exogenous noise to break symmetry needs to be further investigated.

For HPosDec, gESMC was not able to discover the accurate structural model in any of the cases, and therefore yields a test error of 50%. As mentioned earlier, the reason for this failure is due of the greedy search of the structural model. For the k -bit parity function, since all structural models with substructures of order $k - 1$ or lower yield classification models with the same error, the optimal structural model is not hill-climbable. Therefore, the greedy search heuristic fails to identify the accurate structural model and therefore yields inaccurate classification models. This limitation can easily be alleviated by a number of ways, some of which are outlined in the next section.

Finally, we compare the results of gESMC to those obtained with C4.5 and SMO. All three algorithms perform equally well in tackling HPosDec, and gESMC and C4.5 outperform SMO on HPosCount. However, on HParDec and HParCount, C4.5 outperforms both gESMC and SMO. However as with the 2-bit lower-order blocks, the trees of C4.5 had some irrelevant variables in the decision nodes indicating overfitting to the training data.

These results clearly show that gESMC can discover the accurate structural model given that it is hill-climbable from lower-order structural models. In the following section, we discuss some approaches to relax this limitation of gESMC. We also discuss ways to represent structural models with overlapping substructures.

6 Discussion

The results presented in the previous section highlighted both the strengths and limitations of gESMC. In this section we discuss some approaches to overcome the limitations of gESMC which have to be further investigated. We discuss approaches to discover accurate structural models even when there is a lack of guidance from lower-level structural models. We also discuss ways to represent structural models with overlapping substructures.

6.1 Lack of Guidance from Lower-Order Substructures

As mentioned earlier, the greedy search used in gESMC needs some guidance from lower-order substructural models towards the optimal structural model. That is, in order to discover a k -variable substructure the greedy search needs a classification model built with at least one of the substructures of order 2 to be more accurate than that with substructures of order 1, and the classification model built with at least one of the substructures of order 3 has to be more accurate than those with substructures of order 2 and so on. In the absence of such a guidance, the greedy search may stop because it cannot find any structural model that decreases the classification error. To alleviate this limitation, we propose the following two approaches:

Increase the order of substructural merges. We can increase the order of the linkages that the greedy search does if the test error is high and no better structural model is found. That is, at each iteration, instead of pairwise merges, we could permit higher-order merges if the pairwise merges yield no improvement.

We implemented this approach and tested gESMC on the four hierarchical problems. The results show that gESMC obtained 0% test error in all the four problems, and the structural models are correctly evolved. However, the limitation of this approach is the increase in the complexity and cost of the algorithm which is dictated by the maximum order of linkages permitted (ℓ_{max}):

$$Cost = \binom{\ell}{2} \cdot s + \binom{\ell}{3} \cdot s + \dots + \binom{\ell}{\ell_{max}} \cdot s, \quad (9)$$

where s is the cost of building a surrogate. Note that the cost of this approach increases with ℓ_{max} . For this reason, we do not consider this approach as a general solution, although it can be really useful in certain problem domains.

Select randomly one of the new structural models. If the test error is high, and the greedy search cannot find any structural model that significantly decreases this test error, a new structural model can be chosen randomly, or using a technique similar to *simulated annealing* (12). In this case, we would accept a structural model with a higher error in the hope of getting a better model in the subsequent iterations.

Preliminary results using this strategy indicates that gESMC yields maximally accurate classification models for problems consisting of lower-order parity blocks with $k > 2$. However, the structural models evolved are slightly more complicated and contains spurious interactions between variables. Nevertheless, these spurious linkages can be removed by analyzing the classification model and the relative contribution of different schemata to the output.

6.2 Creating Structural Models with Overlapping Substructures

Finally, we look at problems with overlapping linkages where some variables interact with different groups of variables depending on the input that has to be classified. A widely used test problem

with overlapping linkages is the *multiplexer* problem (4; 28), which is defined as follows. Given a bit string of length ℓ , where the first $\log_2 \ell$ bits are the *address bits* and the remaining bits are the *position bits*, the output is the value of the position bit referred by the decimal value of the address bits. For example, for the 6-bit multiplexer, $f(00\ 0101)=0$ and $f(10\ 1011)=1$. Thus, a surrogate with a group formed by all the address bits and the corresponding position bit as a basis accurately determines the output.

We tested gESMC on the 6-bit and 11-bit multiplexer problems. The structural models evolved contained all the address and the position bits in the same linkage group. For example, we obtained the following structural model for the 6-bit multiplexer:

$$[x_0x_1x_2x_3x_4x_5],$$

which resulted in a 0% test error. Since gESMC builds structural model with non-overlapping substructures, one way to handle overlapping substructures is by grouping the substructures together. However, such a merger is unnecessary and other methods which can build structural model with overlapping surrogates such as the design structure matrix genetic algorithm (DSMGA) (32; 15), can evolve a structural model such as:

$$[x_0x_1x_2][x_0x_1x_3][x_0x_1x_4][x_0x_1x_5],$$

The above structural model also yields a surrogate with 0% test error, and gives more information than the former one. Therefore, we will investigate the use of DSMGA and other similar methods that can discover structural models with overlapping variables in developing maximally accurate classification models.

7 Summary and Conclusions

In this paper, we proposed a methodology for learning by building a classification model that uses the structural and surrogate model of a data set. First, we discover the structural model of a set of examples, identifying salient groups of interacting variables to determine the output. Then, we use the structural model to infer the functional form of a surrogate function and the coefficients of the surrogate are estimated using linear regression. Finally, using the substructural surrogate, we build a classification model to predict the class of a given new set of inputs.

We presented gESMC, an implementation of the methodology which uses a greedy search heuristic to search for the structural, surrogate, and classification models that minimize the classification error. Without any problem knowledge, gESMC starts with a simplest model of independent variables and proceeds to explore more complex structural model till the classification error no longer improves or is below a user-defined threshold.

The results of using gESMC on four hierarchical test problems, and its comparison with C4.5 and SMO shows interesting results. Results show that gESMC significantly outperforms C4.5 and SMO in problems that consisted of 2-bit low order blocks in terms of learning accuracy and interpretability. The main difference between gESMC and other learners is that gESMC detects the structure of the data and uses it to predict the class of given inputs. In essence, gESMC not only yielded accurate classification models, but also the classification model is *interpretable*. That is, gESMC not only provides the classification model, but also the structure of the data, making it amenable to human interpretation.

Along with the strengths, the results also highlighted some limitations of gESMC. Specifically, the accuracy of the structural model to capture salient variable interactions depend on the guidance

from lower-order substructures. Therefore, the accuracy of the structural model and consequently the accuracy of the classification model suffers when there is no guidance from lower-order substructures. This limitation is expected given that we use a minimum description length style metric and also a greedy search heuristic which only considers pairwise merges of the substructures. Several approaches were outlined to overcome this limitation which need further investigation.

Acknowledgments

We thank the support of *Enginyeria i Arquitectura La Salle*, Ramon Llull University, *Ministerio de Ciencia y Tecnología* under project TIN2005-08386-C05-04, and *Generalitat de Catalunya* under Grants 2005FI-00252 and 2005SGR-00302.

This work was also sponsored by the Air Force Office of Scientific Research, Air Force Materiel Command, USAF, under grant FA9550-06-1-0096, the National Science Foundation under grant ITR grant DMR-03-25939 at the Materials Computation Center. The U.S. Government is authorized to reproduce and distribute reprints for government purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the Air Force Office of Scientific Research, the National Science Foundation, or the U.S. Government.

References

- [1] Shumeet Baluja. Incorporating a priori knowledge in probabilistic-model based optimization. In Martin Pelikan, Kumara Sastry, and Erick Cantú-Paz, editors, *Scalable Optimization via Probabilistic Modeling: From Algorithms to Applications*, chapter 9, pages 205–219. Springer, Berlin, 2006.
- [2] E. Bernadó-Mansilla and J.M. Garrell. Accuracy-Based Learning Classifier Systems: Models, Analysis and Applications to Classification Tasks. *Evolutionary Computation*, 11(3):209–238, 2003.
- [3] M.V. Butz, M. Pelikan, X. Llorà, and D.E. Goldberg. Automated Global Structure Extraction for Effective Local Building Block Processing in XCS. *Evolutionary Computation*, 14(3):345–380, 2006.
- [4] De Jong, K.A. and Spears, W.M. Learning Concept Classification Rules Using Genetic Algorithms. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pages 651–656. Sidney, Australia, 1991.
- [5] N.R. Drapper and H. Smith. *Applied Regression Analysis*. John Wiley & Sons, New York, USA, 1966.
- [6] J.J. Gibson. *The Ecological Approach to Visual Perception*. Mahwah, NJ: Lawrence Erlbaum Associates, 1979.
- [7] D.E. Goldberg. *The Design of Innovation: Lessons from and for Competent Genetic Algorithms*. Kluwer Academic Publishers, 1 edition, 2002.
- [8] G. Harik. Linkage Learning via Probabilistic Modeling in the ECGA. Technical report, (Ill-GAL Report No. 99010). Urbana, IL: University of Illinois at Urbana-Champaign, January 1999.

- [9] Georges R. Harik, Fernando G. Lobo, and Kumara Sastry. Linkage learning via probabilistic modeling in the ECGA. In Martin Pelikan, Kumara Sastry, and Erick Cantú-Paz, editors, *Scalable Optimization via Probabilistic Modeling: From Algorithms to Applications*, chapter 3, pages 39–61. Springer, Berlin, 2006. (Also IlliGAL Report No. 99010).
- [10] J.H. Holland. *Adaptation in Natural and Artificial Systems*. The University of Michigan Press, 1975.
- [11] S.S. Keerthi and C.J. Lin. Asymptotic Behaviors of Support Vector Machines with Gaussian Kernel. *Neural Computation*, 15(7):1667–1689, 2003.
- [12] J. Korst and E. Aarts. *Simulated Annealing and Boltzmann Machines*. Wiley-Interscience, New York, 1997.
- [13] T. Kovacs. Deletion Schemes for Classifier Systems. In *GECCO'99: Proceedings of the Genetic and Evolutionary Computation Conference*, pages 329–336. Morgan Kaufmann, 1999.
- [14] X. Llorà, K. Sastry, D. E. Goldberg, and L. de la Ossa. The χ -ary extended compact genetic algorithm: Linkage learning in pittsburgh lcs. In T. Kovacs, X. Llorà, and K. Takadama, editors, *Advances at the frontier of LCS*. Springer, Berlin, 2007. (Also IlliGAL Report No. 2006015).
- [15] X. Llorà, K. Sastry, T.-L. Yu, and D. E. Goldberg. Do not match, inherit: Fitness surrogates for genetics-based machine learning. *Proceedings of the 2007 Genetic and Evolutionary Computation Conference*, page Accepted, 2007.
- [16] M. Pelikan. *Hierarchical Bayesian Optimization Algorithm: Toward a new Generation of Evolutionary Algorithms*. Berlin: Springer Verlag, 2005.
- [17] M. Pelikan and K. Sastry. Fitness inheritance in the Bayesian optimization algorithm. *Proceedings of the 2004 Genetic and Evolutionary Computation Conference*, 2:48–59, 2004. (Also IlliGAL Report No. 2004009).
- [18] M. Pelikan, K. Sastry, and E. Cantú-Paz, editors. *Scalable Optimization via Probabilistic Modeling: From Algorithms to Applications*, volume 33 of *Studies in Computational Intelligence*. Springer, 2006.
- [19] J. Platt. Fast Training of Support Vector Machines using Sequential Minimal Optimization. In *Advances in Kernel Methods - Support Vector Learning*, pages 557–563. MIT Press, 1998.
- [20] J.R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers, San Mateo, California, 1995.
- [21] C.R. Rao and H. Toutenburg. *Linear Models: Least Squares and Alternatives*. Springer, Berlin, 1999.
- [22] K. Sastry, C.F. Lima, and D.E. Goldberg. Evaluation Relaxation Using Substructural Information and Linear Estimation. In *GECCO'06: Proceedings of the 8th annual Conference on Genetic and Evolutionary Computation*, pages 419–426, New York, NY, USA, 2006. ACM Press.
- [23] K. Sastry, M. Pelikan, and D. E. Goldberg. Efficiency enhancement of genetic algorithms via building-block-wise fitness estimation. *Proceedings of the IEEE International Conference on Evolutionary Computation*, pages 720–727, 2004. Also IlliGAL Report No. 2004010.
- [24] Kumara Sastry and David E. Goldberg. Probabilistic model building and competent genetic programming. In Rick L. Riolo and Bill Worzel, editors, *Genetic Programming Theory and Practise*, chapter 13, pages 205–220. Kluwer, 2003.

- [25] H.A. Simon. *Sciences of the Artificial*. Cambridge, MA: MIT Press, 1969.
- [26] T.G. Dietterich. Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms. *Neural Comp.*, 10(7):1895–1924, 1998.
- [27] V. Vapnik. *The Nature of Statistical Learning Theory*. Springer Verlag, New York, 1995.
- [28] S.W. Wilson. Quasi-Darwinian Learning in a Classifier System. In *4th IWML*, pages 59–65. Morgan kaufman, 1987.
- [29] S.W. Wilson. Classifier Fitness Based on Accuracy. *Evolutionary Computation*, 3(2):149–175, 1995.
- [30] S.W. Wilson. Generalization in the XCS Classifier System. In *3rd Annual Conf. on Genetic Programming*, pages 665–674. Morgan Kaufmann, 1998.
- [31] I.H Witten and E. Frank. *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann, San Francisco, 2nd edition, 2005.
- [32] T.-L. Yu. *A matrix approach for finding extrema: Problems with modularity, hierarchy, and overlap*. PhD thesis, University of Illinois at Urbana-Champaign, Urbana, IL, 2006.