

Do not Match, Inherit: Fitness Surrogates for Genetics-Based Machine Learning Techniques

Xavier Llorà¹, Kumara Sastry², Tian-Li Yu³, and David E. Goldberg²

¹National Center for Super Computing Applications (NCSA)
University of Illinois at Urbana-Champaign, Urbana IL 61801

²Illinois Genetic Algorithms Laboratory (IlligAL), Dept. of Industrial and Enterprise Systems Eng.
University of Illinois at Urbana-Champaign, Urbana IL 61801

³Department of Electrical Engineering, National Taiwan University, Taipei 106, Taiwan.
xllora@uiuc.edu, ksastry@uiuc.edu, tianliyu@cc.ee.ntu.edu.tw, deg@uiuc.edu

ABSTRACT

A byproduct benefit of using probabilistic model-building genetic algorithms is the creation of cheap and accurate surrogate models. Learning classifier systems—and genetics-based machine learning in general—can greatly benefit from such surrogates which may replace the costly matching procedure of a rule against large data sets. In this paper we investigate the accuracy of such surrogate fitness functions when coupled with the probabilistic models evolved by the χ -ary extended compact classifier system (χ eCCS). To achieve such a goal, we show the need that the probabilistic models should be able to represent all the accurate basis functions required for creating an accurate surrogate. We also introduce a procedure to transform populations of rules based into dependency structure matrices (DSMs) which allows building accurate models of overlapping *building blocks*—a necessary condition to accurately estimate the fitness of the evolved rules.

Categories & Subject Descriptors

I.2.6 [Artificial Intelligence]: Learning—Concept Learning.

General Terms

Algorithms, Design, Theory.

Keywords

Learning Classifier Systems, Genetics-Based Machine Learning, fitness estimation, surrogate fitness, model-building GAs, EDAs, χ eCCS, DSMGA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

GECCO '07, July 7–11, 2007, London, England, United Kingdom.
Copyright 2007 ACM 978-1-59593-697-4/07/0007 ...\$5.00.

1. INTRODUCTION

A daunting challenge for learning classifier systems (LCS)—and genetics-based machine learning (GBML) in general—is the amount of time spent in the rule matching procedure required to estimate the fitness of a rule. Recently, Llorà and Sastry [20] proposed an efficient implementation using hardware accelerated vector operations. However, such efficient implementations do not eliminate the need to match candidate rules against the instance of a given data set. Recently, a renaissance of techniques for fitness inheritance has blossom in the estimation of distribution algorithms (EDAs) community. The goal is to build cheap surrogate fitness functions accurate enough to replace the calculation of expensive fitness functions. Sastry, Lima, and Goldberg [33] proposed a new surrogate fitness based on substructural information and linear estimation.

The introduction of the χ -ary extended compact classifier system (χ eCCS) by Llorà, Sastry, Goldbeg, and de la Ossa [21] showed that GBML approaches could greatly benefit from using competent genetic algorithms [11]. In this paper, we explore how surrogate fitness functions estimate the fitness of a rule to avoid matching it against the instances contained in a given data set. The fitness surrogate proposed by Sastry, Lima, and Goldberg [33] uses the substructural information obtaining from the probabilistic model-building done by eCGA [16], hence, it is a good candidate for χ eCCS.

We present how surrogate fitness functions can be built using substructural information evolved by χ eCCS. We explore how such surrogates behave on two test problems: the hidden XOR and the multiplexer problem. Both problems have different underlying properties. Using both problems we were able to identify the need for probabilistic models able to represent all the accurate basis functions required to create an accurate surrogate. In another words, the substructural model evolved by χ eCCS rely on non-overlapping *building blocks* (BBs) [21]. However, in a certain class of problems, such models may not align with the surrogate expression proposed by Sastry, Lima, and Goldberg [33]—which requires an overlapping model to solve, for instance, the multiplexer problem. For there reasons we will also adapt model-building techniques to evolve overlapping BBs—the dependency structure matrix genetic algorithm (DSMGA) by Yu, Goldberg, Yassine, and Chen [41].

The rest of this paper is structured as follows. Section 2 presents a general overview of the χ eCCS. Then, section 3 reviews some related work on fitness inheritance for EDAs, and section 4 describes the surrogate fitness used in this paper. We conducted initial experiments—as reported in section 5—that show the need, in certain class of problems, to use probabilistic models able to express overlapping BBs. We review one of such methods (DSMGA) in section 6, and what transformations may be required to model rule sets—section 7. Finally, we present some conclusions and further work in section 8.

2. THE χ -ary EXTENDED COMPACT CLASSIFIER SYSTEM

The χ -ary extended compact classifier system (χ eCCS) relies on a χ -ary extended compact genetic algorithm (χ eCGA) [7, 31] to identify *building blocks* among the rules. χ eCCS relies on a close-world assumption and uses rules defined over the ternary alphabet $\{0, 1, \#\}$. Another key element is the ability to provide proper niching capabilities—as already pointed out elsewhere by Bernadó-Mansilla et al. [2, 3].

The χ -ary extended compact genetic algorithm (χ eCGA) [7, 31], is an extension of Harik’s binary eCGA [16]. Unlike the original eCGA, χ eCGA can handle fixed-length chromosomes composed of genes with arbitrary cardinalities (denoted by χ). As in the original eCGA, χ eCGA assumes that a good probability distribution is equivalent to linkage learning. The measure of a good distribution is quantified based on minimum description length (MDL) models. The key concept behind MDL models is that given all things are equal, simpler distributions are better than the complex ones. The MDL restriction penalizes both inaccurate and complex models, thereby leading to an optimal probability distribution. The probability distribution used in eCGA is a class of probability models known as marginal product models (MPMs). MPMs are formed as a product of marginal distributions on a partition of the genes. MPMs also facilitate a direct linkage map with each partition separating tightly linked genes.

The χ eCGA can be algorithmically outlined as follows:

1. Initialize the population with random individuals.
2. Evaluate the fitness value of the individuals
3. Select good solutions by using s-wise tournament selection without replacement [13].
4. Build the probabilistic model: In χ eCGA, both the structure of the model as well as the parameters of the models are searched. A greedy search is used to search for the model of the selected individuals in the population.
5. Create new individuals by sampling the probabilistic model.
6. Evaluate the fitness value of all offspring
7. Replace the parental population (before selection) with the offspring population using restricted tournament replacement (RTR) [15]. We use RTR in order to maintaining multiple maximally general and maximally accurate rules as niches in the population.

8. Repeat steps 3–6 until the finalization criteria are met.

Three things need further explanation: (1) the fitness measure, (2) the identification of MPM using MDL, and (3) the creation of a new population based on MPM.

In order to promote maximally general and maximally accurate rules à la XCS [38], χ eCCS compute the *accuracy* (α) and the *error* (ε) of an individual [22]. In a Pittsburgh-style classifier, the accuracy may be computed as the proportion of overall examples correctly classified, and the error is the proportion of incorrect classifications issued. Let n_{t+} be the number of positive examples correctly classified, n_{t-} the number of negative examples correctly classified, n_m the number of times that a rule has been matched, and n_t the number of examples available. Using these values, the *accuracy* and *error* of a rule r can be computed as:

$$\alpha(r) = \frac{n_{t+}(r) + n_{t-}(r)}{n_t} \quad (1)$$

$$\varepsilon(r) = \frac{n_{t+}}{n_m} \quad (2)$$

We note that the error (equation 2) only takes into account the number of correct positive examples classified¹. This is due to the close-world assumption of the knowledge representation which follows from using a default rule. Once the *accuracy* and *error* of a rule are known, the fitness can be computed as follows.

$$f(r) = \alpha(r) \cdot \varepsilon(r) \quad (3)$$

The above fitness measure favors rules with a good classification accuracy and a low error, or maximally general and maximally accurate rules.

The identification of MPM in every generation is formulated as a constrained optimization problem. The goal is to minimize the *model complexity* and the *compress population complexity*. Further details may be found elsewhere [16, 21]. The greedy search heuristic used in χ -eCGA starts with a simplest model assuming all the variables to be independent and sequentially merges subsets until the MDL metric no longer improves. Once the model is built and the marginal probabilities are computed, a new population is generated based on the optimal MPM as follows, population of size $n(1 - p_c)$ where p_c is the crossover probability, is filled by the best individuals in the current population. The rest $n \cdot p_c$ individuals are generated by randomly choosing subsets from the current individuals according to the probabilities of the subsets as calculated in the model.

One of the critical parameters that determines the success of χ eCGA is the population size. Analytical models have been developed for predicting the population-sizing and the scalability of eCGA [32]. The models predict that the population size required to solve a problem with m building blocks of size k with a failure rate of $\alpha = 1/m$ is given by

$$n \propto \chi^k \left(\frac{\sigma_{BB}^2}{d^2} \right) m \log m, \quad (4)$$

where n is the population size, χ is the alphabet cardinality (here, $\chi = 3$), k is the building block size, $\frac{\sigma_{BB}^2}{d^2}$ is the noise-to-signal ratio [12], and m is the number of building blocks. For the experiments presented in this paper we used $k =$

¹The proportion of correctly classified examples is used instead to simplify the calculation of the final fitness

$|a|+1$ (where $|a|$ is the number of address inputs), $\frac{\sigma_{BB}^2}{d^2}=1.5$, and $m = \frac{\ell}{|I|}$ (where ℓ is the rule size).

Assembling a rule set that describes the concept requires maintaining multiple maximally accurate and maximally general rules. Thus, we need an efficient niching method, that does not adversely affect the quality of the probabilistic models. Therefore, following previous studies in EDAs [25], we use restricted tournament replacement (RTR) [15]. We note that a sub-structural niching method might be better than RTR in stably maintaining multiple niches [30], and it can be readily incorporated into the proposed algorithm. In RTR, each new offspring solution \mathbf{x} is incorporated into the original population using the following three steps: (1) select a random subset \mathbf{W} of size w (called window size) from the original population (before selection), (2) find the solution \mathbf{y} in \mathbf{W} that is most similar to \mathbf{x} (in terms of Euclidean distance), and (3) make a tournament between \mathbf{x} and \mathbf{y} where \mathbf{x} replaces \mathbf{y} if it is better than \mathbf{y} . The parameter w is called window size, and a good rule of thumb for setting this parameter is $w = \min\{\ell\}$, where ℓ is the problem size [25]. We note that the window size w affects the number of niches that can be maintained by RTR. That is increasing the window size can potentially increase the number of niches that can be maintained in the population and also increases the probability of maintaining the niches [15, 25].

We note that the population size n , affects the success probability of maintaining *all* maximally general, maximally accurate rules, γ . In essence, RTR requires larger population sizes to maintain the global optima for longer time. This is a well understood phenomena of niching methods and has been analyzed by Mahfoud for fitness sharing [24] and is applicable to RTR as well [30]. The minimum population size required by RTR for maintaining at least one copy of all but one maximally general maximally accurate rules in the population is given by [24, 30]

$$n \propto \frac{\log \left[\left(1 - \gamma^{1/t} \right) / n_{opt} \right]}{\log \left[(n_{opt} - 1) / n_{opt} \right]} \quad (5)$$

where t is the number of generations we need to maintain all the niches, n_{opt} is the total number of maximally general maximally accurate rules.

3. EVALUATION RELAXATION IN EDAS

In *evaluation relaxation*, an accurate, but computationally expensive fitness function—such as the matching procedure in the χ eCCS—is replaced by a less accurate, but inexpensive surrogate function, and thereby the total number of costly fitness evaluations are reduced [1, 14, 18, 26, 29, 34, 36]. The low-cost, less-accurate fitness estimate can either be (1) *exogenous*, as in the case of approximate fitness functions [1, 18, 23], where, external means are used to develop the fitness estimate, or (2) *endogenous*, as in the case of *fitness inheritance* [36] where, some of the offspring fitnesses are estimated based on fitness of parental solutions.

Sastry, Pelikan, and Goldberg [34] proposed a fitness inheritance method for EDAs, specifically for eCGA—a similar method was proposed for the Bayesian optimization algorithm (BOA) [25] by Pelikan and Sastry [26]. Similar to earlier fitness inheritance study [36], all the individuals in the initial population were evaluated using the expensive fitness function. Thereafter, an offspring was evaluated either using a surrogate with a user-specified inheritance probability

p_i , or using the expensive fitness function with a probability $1 - p_i$. However, the proposed method used the probabilistic models of eCGA to determine the structural form of the surrogate. That is, the MPM model used in eCGA, which partitions the variables of the underlying search problem into linkage groups, were used to determine the variable interactions used in the surrogate. Therefore, the process of learning a surrogate model was sub-divided into estimating the fitness contributions of all possible subsolutions in every partition according to the linkage map that is automatically and adaptively identified by the probabilistic model of eCGA. The authors used all evaluated parents and offspring in estimating the partial contributions of the subsolutions (or schemata) to the overall fitness of a candidate solution [34].

Specifically, they used schema theory basis for determining the relative and partial contribution of a schema to the overall fitness. That is, they defined fitness of a schema h as the difference between the average fitness of individuals that contain the schema and the average fitness of the population [34]:

$$\hat{f}_s(h) = \frac{1}{n_h} \sum_{\{i|x^{(i)} \supset h\}} f(x^{(i)}) - \frac{1}{M} \sum_{i=1}^M f(x^{(i)}), \quad (6)$$

where n_h is the total number of individuals that contain the schema h , $x^{(i)}$ is the i^{th} evaluated individual and $f(x^{(i)})$ its fitness, and M is the total number of individuals that were evaluated. If a particular schema is not present in the evaluated population, its fitness is arbitrarily set to zero.

4. FITNESS INHERITANCE IN EDAS USING LEAST SQUARES FITTING

To address the issues presented in the previous section, Sastry, Lima, and Goldberg [33] proposed a new surrogate fitness based on substructural information—as shown before—and linear estimation. We review in this section the main elements of such an approach since it is the basis of the fitness inheritance scheme used by χ eCCS. Similar to Sastry, Pelikan & Goldberg [34], individuals with exact fitness are used to estimate the sub-structural fitnesses of the remaining individuals. These sub-structures that are defined by the probabilistic model can be viewed and directly mapped into schemata. The fitness associated with the different schemas that match an individual is then combined to estimate his fitness. In this study, schema or building-block fitness is defined as the relative (to the average fitness of the population) fitness contribution to the overall fitness of an individual.

After the model is built the linkage groups are treated as building-blocks partitions, thus all possible schemata under this structure are considered. Considering a MPM example for a 4-bit problem, whose model is [1, 3] [2] [4], the schemata for which the fitness is predicted are $\{0^*0^*, 0^*1^*, 1^*0^*, 1^*1^*, *0^{**}, *1^{**}, ***0, ***1\}$. The total number of schemas is given by

$$N = \sum_{i=1}^m 2^{k_i}, \quad (7)$$

where m is the number of BBs and k_i is the size of the i^{th} BB (number of variables belonging to the BB).

The fitness values of the schemata are estimated as follows. Each individual used for learning is mapped into a binary vector of size N , where each variable of the vector uniquely identifies a given schema. That is, the vector is instantiated by the following delta function

$$\delta(x, h_j) = \begin{cases} 1, & \text{if } x \supset h_j \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

where x is the individual to be converted and h_j is the j^{th} schema. Basically, the vector will have value “1” for the schemas that contain individual x and “0” otherwise. After mapping M evaluated individuals using the above function, the following matrix with dimension $(M \times N)$ is obtained:

$$\mathbf{A} = \begin{pmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,N} \\ a_{2,1} & a_{2,2} & \dots & a_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ a_{M,1} & a_{M,2} & \dots & a_{M,N} \end{pmatrix}, \quad (9)$$

where $a_{i,j} = \delta(x^{(i)}, h_j)$. We note that $x^{(i)}$ denotes the i^{th} individual used for learning the surrogate fitness model. We note that the rank of matrix \mathbf{A} is $N - m + 1$.

Also, the relative (to the average) fitness of each evaluated individual is kept in a vector with dimension $(M \times 1)$ as

$$\mathbf{f} = \begin{pmatrix} f(x^{(1)}) - \bar{f} \\ f(x^{(2)}) - \bar{f} \\ \vdots \\ f(x^{(M)}) - \bar{f} \end{pmatrix}, \quad (10)$$

where $f(x^{(i)})$ is the evaluated fitness of the i^{th} individual used for learning and \bar{f} is the average fitness of all M evaluated individuals (both from parent and offspring population). The average fitness is then given by

$$\bar{f} = \frac{1}{M} \sum_{i=1}^M f(x^{(i)}). \quad (11)$$

Given that there are N different schema fitnesses to estimate, the fitness coefficients associated with the N binary variables can be displayed as vector of dimension $(N \times 1)$

$$\hat{\mathbf{f}}_s = \begin{pmatrix} \hat{f}_s(h_1) \\ \hat{f}_s(h_2) \\ \vdots \\ \hat{f}_s(h_N) \end{pmatrix}, \quad (12)$$

where $\hat{f}_s(h_j)$ is the fitness of schema h_j .

The task of estimating the relative fitness of each schema can be stated as finding a vector $\hat{\mathbf{f}}_s$ that satisfies the equality:

$$\mathbf{A}\hat{\mathbf{f}}_s = \mathbf{f}. \quad (13)$$

In practice, this equality might not be entirely satisfied and one must instead seek for minimizing the difference between left and right terms of Equation 13. For that, it is used a multi-dimensional least squares fitting approach. Thus, under the least squares fitting principle the problem of estimating the fitness of schemata can now be reformulated as finding the appropriate values for vector $\hat{\mathbf{f}}_s$ such that the following squared error function χ^2 is minimized:

$$\chi^2 = (\mathbf{A}\hat{\mathbf{f}}_s - \mathbf{f})^T (\mathbf{A}\hat{\mathbf{f}}_s - \mathbf{f}). \quad (14)$$

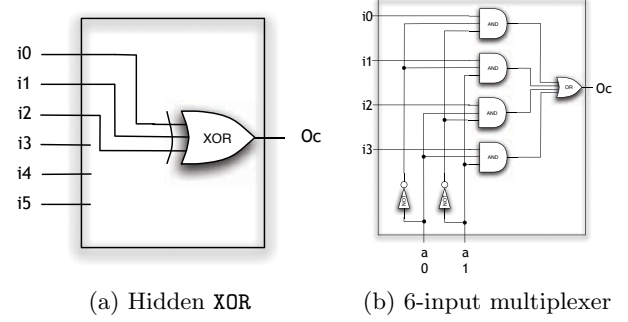


Figure 1: The two problems presented in this figure present different properties. XOR is a problem with no overlapping in the variable space. The 6-input multiplexer, on the other hand requires exhaustive reuse of variable to provide the proper output.

The solution to the above problem is a well-known result from literature, therefore details on the resolution are not provided and the interested reader should refer elsewhere [4, 8, 17, 19, 27]. The method used in this paper to perform multi-dimensional least squares fitting was provided by the R project for statistical computing².

After obtaining the estimates for schema fitnesses, the estimation of an individual’s fitness is a straightforward process that consists in summing the average fitness of the population to the fitness of each schema that contains the individual being considered. The estimated fitness of an individual x is then given by

$$f_{inh}(x) = \bar{f} + \sum_{j=1}^N \delta(x, h_j) \hat{f}_s(h_j), \quad (15)$$

where $\hat{f}_s(h_j)$ is given by the j^{th} element of vector $\hat{\mathbf{f}}_s$.

It can be easily seen that the surrogate obtained by using a structure inferred from a perfect model and the coefficients via least squares yields that is identical to Walsh transform [10] of the accurate fitness function. This clearly suggests that given an accurate probabilistic model, we can obtain a surrogate that accurately estimates the fitness of untested solutions.

5. FITNESS INHERITANCE IN χ ECCS: CHOOSING THE RIGHT BASIS FUNCTIONS

The surrogate fitness model proposed in the previous section has been successfully used in fitness inheritance schemes in the optimization domain [33]. To successfully apply such fitness surrogate in the learning domain of the χ ECCS [21] and other genetics-based machine learning, we conducted preliminary analysis. We used two well-know problems to the genetics-based machine learning community: (1) the hidden XOR problem, and (2) the 6-input multiplexer. We summarize in this section the insights obtain from this simple experimentation—see Figure 1.

The hidden XOR problem requires discovering a XOR function embedded in a set of inputs—most of them irrelevant to the XOR classification task. Our initial tests used a simple

²<http://www.r-project.org/>

3-bit XOR embedded in a 6 binary input space. χ eCCS easily solves the problem providing the following output rules, fitness, and the final MPMs model:

<i>Rules:</i>	111###	→ 1	$f(r) = 0.625$
	001###	→ 1	$f(r) = 0.625$
	010###	→ 1	$f(r) = 0.625$
	100###	→ 1	$f(r) = 0.625$
	<i>default</i>	→ 0	
<i>Model:</i>	[i0 i1 i2]	[i3]	[i4] [i5]

The rules provided by χ eCCS correctly describes the hidden XOR problem. Moreover, the model evolved by χ eCCS clearly groups the two variables [i0 i1 i2] involved in the hidden XOR. With this information we should be able to create the surrogate fitness following the steps described in the previous section. Once we obtained the surrogate, we generated all the possible rules (3^6) and computed the average error of the surrogate fitness. For illustrative purposes we present the average error of the surrogate for the evolved χ eCCS model and two other suboptimal models.

<i>Model</i>	$\epsilon(f_{inh}(x))$
[i0] [i1] [i2] [i3] [i4] [i5]	14.976%
[i0 i1] [i2] [i3] [i4] [i5]	14.952%
[i0 i1 i2] [i3] [i4] [i5]	0.890%

The results presented above show how when the correct model is provided, the surrogate model performs remarkably well. Encouraged by these results, we repeated the experiments for the 6-input multiplexer problem.

<i>Model</i>	$\epsilon(f_{inh}(x))$
[i0] [i1] [i2] [i3] [i4] [i5]	18.703%
[i0 i2] [i1 i5] [i3 i4]	18.665%

The results (see above) were not as compelling as the ones obtained on the hidden XOR problem. As stated by Sastry, Lima, and Goldberg [33] the proposed surrogate fitness based on substructural information and linear estimation should hold for any given accurate model. The reason for the poor performance of the surrogate fitness was not the method, but the model used to create such a surrogate. The model evolved by χ eCCS—the same as eCGA—is based on MPMs and, hence, by definition non-overlapping. This means that χ eCCS is able to solve the multiplexer problem using a non-overlapping and approximate model. This model does not prevent the χ eCCS to solve the multiplexer problem quickly, reliably, and accurately [21].

However, this approximated non-overlapping model does not produce an accurate surrogate fitness function. This is the result of choosing an inappropriate set of basis functions to fit the regression schema proposed by the surrogate [5]. In other words, to successfully use such surrogate for the multiplexer problem, we need to induce overlapping BBs. A validation of this intuition is presented below.

<i>Model</i>	$\epsilon(f_{inh}(x))$
[i0] [i1] [i2] [i3] [i4] [i5]	18.703%
[i0 i2] [i1 i5] [i3 i4]	18.665%
[i0 i1 i2] [i0 i1 i3]	
[i0 i1 i4] [i0 i1 i5]	0.733%

The last model shows how an overlapping model leads to a surrogate fitness for the multiplexer problem. This model easily follows the intuitive underlying model based on similar basis functions—as presented in Figure 1(b). Thus, if we want to use fitness inheritance we need to use a model builder that is able to produce such overlapping models. The following section introduces a competent GA that is able to evolve such models as the one required to solve the multiplexer problem.

6. OVERLAPPING BUILDING BLOCKS USING DSMGA

This section gives a brief introduction to the model-building process used in the dependency structure matrix genetic algorithm (DSMGA), which is later used in this paper as the basis of the proposed method. A detailed description of DSMGA is beyond the scope of this paper and can be found elsewhere [41]. DSMGA utilizes the dependency structure matrix (DSM) clustering techniques to extract overlapping BBs. In this section we introduce the concept of DSM and the DSM clustering problem. Then, we describe the metric to cluster DSMs and the algorithm used.

A dependency structure matrix is essentially an adjacency matrix representation of a graph where each entry d_{ij} represents the dependency between node i and node j [37, 39]. Entries d_{ij} can be real numbers or integers. The larger the d_{ij} is, the higher the interaction is between node i and node j . If we focus on the $[0, 1]$ domain, then $d_{ij} = 0$ means that node i and node j do not interact, and $d_{ij} = 1$ means that node i and node j interact with each other. The diagonal entries (d_{ii}) have no significance and are usually set to zero or blacked-out. For elaborate exposition of DSM, please see MIT DSM web site: <http://www.dsmweb.org/>.

The goal of DSM clustering is to find subsets of DSM elements (*i.e.*, clusters) so that nodes within a cluster are maximally interacting, and clusters are minimally interacting. In a typical DSM clustering problem, overlapping clusters (clusters that share same nodes) are permissible. The DSM model of linkage allows overlapping of variables, as opposed to the non-overlapping MPM proposed by eCGA [16]. However, rearranging a DSM to obtain the proper clusters requires a metric to compute the usefulness of the proposed clustering rearrangement [9, 35, 41].

DSMGA relies on a DSM clustering metric based on the minimal description length principle (MDL) [28]. Suppose that we have a model which describes a given data set, $DSM = [d_{ij}]$. Here, the model means a description that specifies which node belongs to which cluster. Usually, the model does not completely describe the given data; otherwise, the model would be too complex to use. Therefore, the description length that the model needs to describe the given data consists of two parts: the model description and the mismatched data description.

The minimum description length principle (MDL) [28] satisfies the needs for dealing with the above trade-off. The MDL can be interpreted as follows: among all possible models, choose the model that uses the minimal length for describing a given data set (that is, model description length plus mismatched data description length). There are two key points that should be noted when MDL is used: (1) the encoding should be uniquely decodable, and (2) the length of encoding should reflect the complexity.

Model Encoding. The description of each cluster starts with a number which is sequentially assigned to each cluster, and then this is followed by a sequence of nodes in the cluster. It is easily seen that the length of this model description is as follows:

$$\sum_{i=1}^{n_c} (\log_2 n_n + cl_i \cdot \log_2 n_n), \quad (16)$$

where n_c is the number of clusters in the model, n_n is the number of nodes, cl_i is the number of nodes in the

i -th cluster. If n_n and n_c are known, the above model description is uniquely decodable. When n_n is given, and by assuming $n_c \leq n_n$, then $\log n_n$ bits are needed to describe n_c . It is a constant for all models, and therefore they are omitted without loss of accuracy.

Mismatched Data Description. Based on the model, another DSM ($DSM' = [d'_{ij}]$) is constructed, where each entry d'_{ij} is 1 if and only if some cluster contains both node i and node j simultaneously. Then, DSM' is compared to the given DSM . For every mismatched entry, where $d'_{ij} \neq d_{ij}$, a description to indicate where the mismatch occurred (i and j) is needed and one additional bit to indicate whether the mismatch is zero-to-one or one-to-zero. Define a mismatch set $S = \{(i, j) | d'_{ij} \neq d_{ij}\}$. The mismatched data description length is given by:

$$\sum_{(i,j) \in S} (\log n_n + \log n_n + 1). \quad (17)$$

The first $\log n_n$ in the bracket indicates i , the second one indicates j , and the additional one bit indicates the type of mismatch.

The MDL clustering metric is given by the summation of the model description length and the mismatched data description. With some arithmetic manipulations, the metric can be expressed as follows:

$$f_{DSM}(M) = \log n_n \sum_{i=1}^{n_c} (cl_i + 1) + |S|(2 \log n_n + 1), \quad (18)$$

where n_c is the number of clusters, n_n is the number of nodes in the DSM, cl_i is the size of the i -th cluster, and S is a mismatch set.

With the above metric, the DSM clustering problem is converted into an optimization problem. Given a DSM, the objective is to find a DSM clustering arrangement (model, M) to minimize the above metric (f_{DSM}). A steepest descent algorithm was adopted by Yu, Goldberg, Yassine, and Chen [41] to optimize the DSM clustering problem. Based on the MDL metric, it add/remove one node to/from one cluster at each iteration. The steepest descent algorithm stops when no further improvement is possible. Further details can be found elsewhere [41].

7. FITNESS INHERITANCE, BASIS FUNCTIONS, AND OVERLAPPING BBS

As described in the previous section, DSM clustering is able to identify overlapping BBs. Using DSM clustering technique requires being able to express the interaction information among variables as a DSM. To achieve this goal we cannot directly use the transformation method proposed by DSMGA [40], since it deals with binary populations. However, if we can define a transformation from a rule set to a DSM that retains the interaction properties of the variables, then DSM cluster will provide us with the proper BB identification mechanism.

We can regard a rule r as a set of interactions among specific values. For instance, given the 6-input multiplexer rule $001### \rightarrow 1$, the first three positions contain specific values that need to interact with each other to properly assemble the rule. Hence, given a rule r which condition is defined among the set of possible variables X , we can define the interaction δ_s between the i th position and the j th position of

the rule r as:

$$\delta_s(r_i, r_j) = \begin{cases} 1, & \text{if } r_i \neq \# \wedge r_j \neq \# \wedge i \neq j \\ 0, & \text{otherwise} \end{cases} \quad (19)$$

That is, an interaction between position i and j exists if and only if both positions in the rule contain specific values. Thus, we should be able to define a matrix C where c_{ij} is the count of interactions among variables i and j in X for a rule set R . Such a matrix is defined as:

$$c_{ij} = \sum_{r \in R} \sum_{i \in X} \sum_{j \in X} \delta_s(r_i, r_j) \quad (20)$$

The count matrix C is the base of the DSM, where we define each of each entries d_{ij} as the normalized C where interactions belong to the $[0,1]$ domain. This can be simply achieved by defining d_{ij} as

$$d_{ij} = \frac{c_{ij}}{\max(C)} \quad (21)$$

Then, we have to decide which rules should be used when defining the DSM. The answer is easy, only the distinct accurate ones—as was already suggested elsewhere [6]. This can be easily achieved by filtering the rules of a population based on their computed error $\varepsilon(r)$. Only rules with no error should be used to build the DSM to cluster.

Another important consideration is if rules belonging to different classes should be mixed together in R . Due to the use of a default rule—close world assumption—when solving binary classification problems, only rules belonging to one class are evolved. For non-binary classification problems we should construct one DSM per class, avoiding the introduction of spurious interclass interactions.

Given the evolved rules for (1) the hidden XOR presented in section 5, and (2) the 6-input multiplexer presented below,

Rules:		$f(r)$
001###	$\rightarrow 1$	$f(r) = 0.625$
01#1##	$\rightarrow 1$	$f(r) = 0.625$
10##1#	$\rightarrow 1$	$f(r) = 0.625$
11###1	$\rightarrow 1$	$f(r) = 0.625$
0#11##	$\rightarrow 1$	$f(r) = 0.625$
#01#1#	$\rightarrow 1$	$f(r) = 0.625$
1###11	$\rightarrow 1$	$f(r) = 0.625$
#1#1#1	$\rightarrow 1$	$f(r) = 0.625$
default	$\rightarrow 0$	

the DSM matrices for both problem are

$$DSM_{XOR} = \begin{pmatrix} 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \quad (22)$$

$$DSM_{MUX} = \begin{pmatrix} 0 & 1 & 0.5 & 0.5 & 0.5 & 0.5 \\ 1 & 0 & 0.5 & 0.5 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0 & 0.25 & 0.25 & 0 \\ 0.5 & 0.5 & 0.25 & 0 & 0 & 0.25 \\ 0.5 & 0.5 & 0.25 & 0 & 0 & 0.25 \\ 0.5 & 0.5 & 0 & 0.25 & 0.25 & 0 \end{pmatrix} \quad (23)$$

Using this DSMs, DSM clustering returns [i0 i1 i2] [i3] [i4] [i5] as model for the hidden XOR problem. As we presented it in section 5, this is the right model the leads to the creation of an accurate fitness surrogate. On the other hand, the DSM clustering produces the model [i0 i1] <i2 i3 i4 i5> indicating that there is a building block

[i0 i1] that interacts with a bus of variables <i2 i3 i4 i5>. Hence, the bus can be expanded as [i0 i1 i2] [i0 i1 i3] [i0 i1 i4] [i0 i1 i5] giving the right overlapping model—see section 5—to build an accurate surrogate for the multiplexer problem.

8. CONCLUSIONS AND FURTHER WORK

We have shown how fitness inheritance for genetics-based machine learning techniques is possible. A surrogate fitness based on substructural information and least square fitting is able to accurately predict the fitness of the rules evolve by χ eCCS. Such a surrogate can replace the cost of computing the fitness of a rule against large data sets. We have also show how χ eCCS is able to solve hidden XOR and multiplexer problems quickly, reliably, and accurately by using approximating probabilistic models of the population of rules based on non-overlapping BBs.

However, such rough approximation models are not enough to build a proper surrogate fitness model. An accurate surrogate fitness function requires an accurate probabilistic model able to express overlapping BBs—as empirically shown for the multiplexer problem. The functional basis used to create the surrogate fitness need to be carefully chosen to allow an accurate regression of the fitness of the rules. In order to obtain such a surrogate we have defined a transformation which from a set of accurate rules creates a DSM that, when properly clustered using the DSM clustering method, provides accurate overlapping BBs. This overlapping BBs define the correct basis which allow the creation of accurate surrogates for certain class of Boolean problems such as the multiplexer. Future research should focus on introducing probabilistic model-building GAs able to express overlapping into the χ eCCS and bounding the speedup produce by the use of the surrogate.

9. ACKNOWLEDGMENTS

This work was sponsored by the Air Force Office of Scientific Research, Air Force Materiel Command, USAF, under grant FA9550-06-1-0370, the National Science Foundation under ITR grant DMR-03-25939 at Materials Computation Center and under grant ISS-02-09199 at the National Center for Supercomputing Applications, UIUC. 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.

10. REFERENCES

- [1] J.-F. M. Barthelemy and R. T. Haftka. Approximation concepts for optimum structural design—a review. *Structural Optimization*, 5:129–144, 1993.
- [2] E. Bernadó-Mansilla and J. M. Garrell-Guiu. MOLeCS: A MultiObjective Learning Classifier System. *Proceedings of the 2000 Conference on Genetic and Evolutionary Computation*, 1:390, 2000.
- [3] E. Bernadó-Mansilla, X. Llorà, and I. Traus. *MultiObjective Machine Learning*, chapter

- MultiObjective Learning Classifier System, pages 261–288. Springer, 2005.
- [4] Å. Björk. *Numerical method for least squares problems*. SIAM, Philadelphia, PA, 1996.
- [5] G. Box, J. Hunter, W.G. and Hunter, and W. Hunter. *Statistics for Experimenters: Design, Innovation, and Discovery*. Wiley, 2005.
- [6] M. V. Butz, M. Pelikan, X. Llorà, and D. E. Goldberg. Extracted global structure makes local building block processing effective in XCS. *Proceedings of the 2005 Conference on Genetic and Evolutionary Computation*, 1:655–662, 2005.
- [7] L. de la Ossa, K. Sastry, and F. G. Lobo. Extended compact genetic algorithm in C++: Version 1.1. IlliGAL Report No. 2006013, University of Illinois at Urbana-Champaign, Urbana, IL, March 2006.
- [8] N. R. Draper and H. Smith. *Applied Regression Analysis*. John Wiley & Sons, New York, USA, 1966.
- [9] C. Fernandez. *Integration Analysis of Product Architecture to Support Effective Team Co-location*. Master thesis, Massachusetts Institute of Technology, 1998.
- [10] D. E. Goldberg. Genetic algorithms and Walsh functions: Part I, a gentle introduction. *Complex Systems*, 3:129–152, 1989. (Also IlliGAL Report No. 88006).
- [11] D. E. Goldberg. *The Design of Innovation: Lessons from and for Competent Genetic Algorithms*. Kluwer Academic Publishers, Norwell, MA, 2002.
- [12] D. E. Goldberg, K. Deb, and J. H. Clark. Genetic algorithms, noise, and the sizing of populations. *Complex Systems*, 6:333–362, 1992. (Also IlliGAL Report No. 91010).
- [13] D. E. Goldberg, B. Korb, and K. Deb. Messy genetic algorithms: Motivation, analysis, and first results. *Complex Systems*, 3(5):493–530, 1989.
- [14] J. J. Grefenstette and J. M. Fitzpatrick. Genetic search with approximate function evaluations. *Proceedings of the International Conference on Genetic Algorithms and Their Applications*, pages 112–120, 1985.
- [15] G. R. Harik. Finding multimodal solutions using restricted tournament selection. *Proceedings of the Sixth International Conference on Genetic Algorithms*, pages 24–31, 1995. (Also IlliGAL Report No. 94002).
- [16] G. R. Harik, F. G. Lobo, and K. Sastry. Linkage learning via probabilistic modeling in the ECGA. In M. Pelikan, K. Sastry, and E. Cantú-Paz, editors, *Scalable Optimization via Probabilistic Modeling: From Algorithms to Applications*, chapter 3. Springer, Berlin, in press. (Also IlliGAL Report No. 99010).
- [17] S. Haykin. *Adaptive Filter Theory*. Prentice Hall, 1996.
- [18] Y. Jin. A comprehensive survey of fitness approximation in evolutionary computation. *Soft Computing Journal*, 9(1):3–12, 2005.
- [19] T. Kailath, A. H. Sayed, and B. Hassibi. *Linear estimation*. Prentice-Hall, Upper Saddle River, NJ, 2000.
- [20] X. Llorà and K. Sastry. Fast rule matching for learning classifier systems via vector instructions. In *Proceedings of the 2006 Genetic and Evolutionary*

- Computation Conference (GECCO 2006)*, page 15131520, 2006. (Also IlliGAL Report No. 2006001).
- [21] X. Llorà, K. Sastry, D. E. Goldberg, and L. de la Ossa. The χ -ary extended compact classifier system: Linkage learning in pittsburgh lcs. In *Proceedings of the 2006 Genetic and Evolutionary Computation Conference (GECCO 2006) Workshops: International Workshop on Learning Classifier Systems*, 2006. (Also IlliGAL Report No. 2006015).
- [22] X. Llorà, K. Sastry, D. E. Goldberg, A. Gupta, and L. Lakshmi. Combating user fatigue in iGAs: Partial ordering, support vector machines, and synthetic fitness. In *GECCO 2005: Proceedings of the 2005 conference on Genetic and evolutionary computation*, volume 2, pages 1363–1370, Washington DC, USA, 25–29 June 2005. ACM Press.
- [23] X. Llorà, K. Sastry, D. E. Goldberg, A. Gupta, and L. Lakshmi. Combating user fatigue in iGAs: Partial ordering, support vector machines, and synthetic fitness. *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 1363–1370, 2005. (Also IlliGAL Report No. 2005009).
- [24] S. W. Mahfoud. Population size and genetic drift in fitness sharing. *Foundations of Genetic Algorithms*, 3:185–224, 1994. (Also IlliGAL Report No. 94005).
- [25] M. Pelikan. *Hierarchical Bayesian Optimization Algorithm: Toward a New Generation of Evolutionary Algorithm*. Springer Verlag, Berlin, 2005.
- [26] M. Pelikan and K. Sastry. Fitness inheritance in the bayesian optimization algorithm. *Proceedings of the Genetic and Evolutionary Computation Conference*, 2:48–59, 2004. (Also IlliGAL Report No. 2004009).
- [27] C. R. Rao and H. Toutenburg. *Linear models: Least squares and alternatives*. Springer, Berlin, 1999.
- [28] J. Rissinen. Modeling by shortest data description. *Automatica*, 14:465–471, 1978.
- [29] K. Sastry. Evaluation-relaxation schemes for genetic and evolutionary algorithms. Master’s thesis, University of Illinois at Urbana-Champaign, General Engineering Department, Urbana, IL, 2001. (Also IlliGAL Report No. 2002004).
- [30] K. Sastry, H. A. Abbass, D. E. Goldberg, and D. D. Johnson. Sub-structural niching in estimation of distribution algorithms. *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 671–678, 2005. (Also IlliGAL Report No. 2005003).
- [31] K. Sastry and D. E. Goldberg. Probabilistic model building and competent genetic programming. In R. L. Riolo and B. Worzel, editors, *Genetic Programming Theory and Practise*, chapter 13, pages 205–220. Kluwer, 2003.
- [32] K. Sastry and D. E. Goldberg. Designing competent mutation operators via probabilistic model building of neighborhoods. *Proceedings of the Genetic and Evolutionary Computation Conference*, 2:114–125, 2004. Also IlliGAL Report No. 2004006.
- [33] K. Sastry, C. F. Lima, and D. E. Goldberg. Evaluation relaxation using substructural information and linear estimation. In *Proceedings of the 2006 Genetic and Evolutionary Computation Conference (GECCO 2006)*, pages 419–426, 2006. (Also IlliGAL Report No. 2006003).
- [34] 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.
- [35] D. Sharman, A. Yassine, and P. Carlile. Characterizing modular architectures. ASME 14th International Conference, pages DTM–34024, Sept. 2002.
- [36] R. Smith, B. Dike, and S. Stegmann. Fitness inheritance in genetic algorithms. In *Proceedings of the ACM Symposium on Applied Computing*, pages 345–350, New York, NY, USA, 1995. ACM.
- [37] D. V. Steward. The design structure system: A method for managing the design of complex systems. *IEEE Transactions on Engineering Management*, 28:77–74, 1981.
- [38] S. Wilson. Classifier fitness based on accuracy. *Evolutionary Computation*, 3(2):149–175, 1995.
- [39] A. Yassine, D. R. Falkenburg, and K. Chelst. Engineering design management: An informatoin structure approach. *International Journal of production research*, 37(13):2957–2975, 1999.
- [40] T.-L. Yu and D. E. Goldberg. Conquering hierarchical difficulty by explicit chunking: Substructural chromosome compression. In *Proceedings of the 2006 Genetic and Evolutionary Computation Conference (GECCO 2006) Workshops: International Workshop on Learning Classifier Systems*, pages 1385–1392, 2006. (Also IlliGAL Report No. 2006007).
- [41] T.-L. Yu, D. E. Goldberg, A. Yassine, and Y.-P. Chen. A genetic algorithm design inspired by organizational theory: Pilot study of a dependency structure matrix driven genetic algorithm. *Artificial Neural Networks in Engineering*, pages 327–332, 2003. (Also IlliGAL Report No. 2003007).