

---

# Globally Multimodal Problem Optimization Via an Estimation of Distribution Algorithm Based on Unsupervised Learning of Bayesian Networks

J. M. Peña

Computational Biology, Dept. of Physics and Measurement Technology,  
Linköping University, Sweden

jmp@ifm.liu.se

J. A. Lozano

P. Larrañaga

Intelligent Systems Group, Dept. of Computer Science and Artificial Intelligence,  
University of the Basque Country, Spain

ccploalj@si.ehu.es

ccplamup@si.ehu.es

---

## Abstract

Many optimization problems are what can be called globally multimodal, i.e., they present several global optima. Unfortunately, this is a major source of difficulties for most estimation of distribution algorithms, making their effectiveness and efficiency degrade, due to genetic drift. With the aim of overcoming these drawbacks for discrete globally multimodal problem optimization, this paper introduces and evaluates a new estimation of distribution algorithm based on unsupervised learning of Bayesian networks. We report the satisfactory results of our experiments with symmetrical binary optimization problems.

## Keywords

Estimation of distribution algorithms, Bayesian networks, unsupervised learning.

## 1 Introduction

Estimation of distribution algorithms (EDAs) (Larrañaga and (eds.), 2001; Mühlenbein and Paaß, 1996; Pelikan, 2002; Pelikan et al., 2000) are some relatively novel evolutionary algorithms (EAs) that are receiving increasing attention in the literature. Like any other class of EAs, EDAs solve a given optimization problem by evolving a population of individuals, i.e., a set of solutions to the optimization problem, towards promising zones of the search space. Such an evolution is mainly based on iterating between two steps: Selection of fit individuals from the current population, and combination of the selected individuals in order to create an offspring population and replace (partially) the current one. Unlike most EAs, EDAs do not make use of variation operators (e.g., crossover and/or mutation) in the combination step. Instead, EDAs generate the offspring population at each iteration by learning and subsequent simulation of a joint probability distribution for the individuals selected.

How the joint probability distribution is estimated from the individuals selected at each iteration as well as what assumptions are made for this process to be tractable is what distinguishes one EDA from another. However, existing EDAs typically disregard that many optimization problems are what can be called *globally multimodal*, i.e., they

present multiple global optima, and that often it is necessary or desirable to discover as many global optima as possible instead of only one of them. In this paper, we propose and evaluate a new EDA tailored to this scenario: The unsupervised estimation of Bayesian network algorithm (UEBNA). The only peculiarity of the UEBNA is the use of a Bayesian network for data clustering (Peña, 2001; Peña et al., 2001a; Peña et al., 1999; Peña et al., 2000; Peña et al., 2002) in order to factorize the joint probability distribution for the individuals selected at each iteration. This allows modelling simultaneously the basins of the different global optima represented by the selected individuals. Therefore, we conjecture that the UEBNA should be able to discover more global optima per run than existing EDAs while speeding up convergence in terms of number of evaluations of the objective function. The reason for this is because individuals belonging to different basins are not used together which usually results in poorly fit individuals and, thus, delays convergence.

The remainder of this paper is structured as follows. Section 2 reviews EDAs. Section 3 describes unsupervised learning of Bayesian networks and shows how this is incorporated into the EDA framework, resulting in the UEBNA. Section 4 evaluates the UEBNA on symmetrical binary optimization problems. Finally, Section 5 closes with some discussion and conclusions.

## 2 Estimation of Distribution Algorithms

Among stochastic heuristic search strategies for problem optimization, *evolutionary algorithms* (EAs) are well known for their good performance and wide applicability. Some classical EAs are genetic algorithms (Goldberg, 1989; Holland, 1975), evolutionary programming (Fogel, 1962; Fogel, 1964) and evolution strategies (Rechenberg, 1973; Schwefel, 1981). More recently, a novel class of EAs, known as *estimation of distribution algorithms* (EDAs) (Larrañaga and eds., 2001; Mühlenbein and Paaß, 1996; Pelikan, 2002; Pelikan et al., 2000), has been proposed and evaluated successfully in a wide variety of scenarios. This section first reviews EDAs and, then, discusses the difficulties that they encounter when optimizing globally multimodal problems. Prior to this, we introduce some terms used throughout the text.

The main feature shared by all the instances of the EA paradigm is being inspired by natural evolution of species. That is why much of the nomenclature of EAs is borrowed from the field of natural evolution. For instance, we talk about *populations* to refer to sets of solutions to an optimization problem, each solution is called an *individual*, and each basic component of an individual is called a *gene*. The main components of most EAs are: An initial population of individuals, a *selection method* over individuals, a set of *variation operators* over individuals, and a *replacement method* over individuals. Basically, all the EAs work in the same iterative way: At each iteration or *generation* some individuals of the current population are selected according to the selection method and modified by the variation operators in order to create new individuals and, consequently, a new population according to the replacement method. The objective of this iterative process is to evolve the population towards promising zones of the search space of the problem at hand.

### 2.1 Generic Estimation of Distribution Algorithm

The most distinctive characteristic of EDAs with respect to the rest of EAs is that EDAs replace the application of variation operators in order to generate the next population from the current one at each iteration by learning and subsequent simulation of a joint probability distribution for those individuals selected from the current population by

- 
1. Let  $\mathbf{po}_1$  be a population composed of  $Q$  uniformly generated individuals
  2. Evaluate the individuals in  $\mathbf{po}_1$
  3.  $u = 1$
  4. **while** the stopping condition is not met **do**
  5.   Let  $\mathbf{d}_u$  group  $N$  individuals selected from  $\mathbf{po}_u$  via the selection method
  6.   Let  $p_u(\mathbf{x})$  be the joint probability distribution for  $\mathbf{X}$  learnt from  $\mathbf{d}_u$
  7.   Let  $\mathbf{of}_u$  be the offspring population composed of  $M$  individuals sampled from  $p_u(\mathbf{x})$
  8.   Evaluate the individuals in  $\mathbf{of}_u$
  9.   Let  $\mathbf{po}_{u+1}$  be the population created from  $\mathbf{po}_u$  and  $\mathbf{of}_u$  via the replacement method
  10.    $u++$
  11. Return the best individuals found so far
- 

Figure 1: Pseudocode of the generic EDA.

means of the selection method. This results in two important advantages of EDAs over classical EAs: The sometimes necessary design of variation operators tailored to the particular optimization problem at hand is avoided, and the number of parameters to be assessed by the user is reduced. A further advantage of EDAs over classical EAs is that the relationships between the random variables that represent the genes of every individual selected can be explicitly expressed through the joint probability distribution learnt from them, instead of being implicitly kept by the individuals of successive populations as building blocks. In fact, it was already recognized in (Goldberg, 1989; Holland, 1975) that detecting interacting genes would be beneficial to genetic algorithms. This source of knowledge was called linkage information. This idea has been exploited by many researchers for the last few years in order to enhance the performance of genetic algorithms (Goldberg, 1989; Goldberg et al., 1993; Lobo et al., 1998). Finally, EDAs outperform classical EAs in deceptive optimization problems (Etxeberria and Larrañaga, 1999; Larrañaga and (eds.), 2001; Mühlenbein et al., 1999).

The generic EDA iterates between three main steps, after the individuals of the initial population  $\mathbf{po}_1$  have been generated, usually uniformly, and evaluated. The iterative process ends when the stopping criterion is met, e.g., performance of a maximum number of generations, uniformity in the current population, or no improvement with regard to the best individual of the previous generation. This causes the best solutions found so far being returned. The three main steps of the  $u$ -th iteration of the generic EDA are as follows for all  $u$ . First,  $N$  of the  $Q$  individuals of the current population  $\mathbf{po}_u$  are selected by means of the selection method. Then, these selected individuals are used to construct a learning database  $\mathbf{d}_u$  from which a joint probability distribution for  $\mathbf{X}$ ,  $p_u(\mathbf{x})$ , is induced.  $\mathbf{X} = (X_1, \dots, X_n)$  denotes an  $n$ -dimensional discrete random variable, where  $X_i$  is associated with the  $i$ -th gene of every individual in  $\mathbf{d}_u$ . Finally,  $M$  individuals are sampled from  $p_u(\mathbf{x})$  and evaluated in order to create the offspring population  $\mathbf{of}_u$  which, then, is used to generate the new population  $\mathbf{po}_{u+1}$  by replacing some individuals of  $\mathbf{po}_u$  via the replacement method. See Fig. 1 for a schematic of the generic EDA.

## 2.2 Families of Estimation of Distribution Algorithms

We have discussed above that replacing the application of variation operators by the learning and simulation of  $p_u(\mathbf{x})$  has immediate benefits. However, it carries some cost too because learning  $p_u(\mathbf{x})$  from  $\mathbf{d}_u$  is not a trivial task. As the computation of all the parameters needed to completely specify  $p_u(\mathbf{x})$  in the extensive representation is often prohibitive, several families of EDAs have arisen where this joint probability distribution is assumed to factorize according to a certain class of probabilistic models. The remainder of this section briefly reviews some of these families, according to an in-

creasing order of complexity. The interested reader may consult (Larrañaga and (eds.), 2001) for a more thorough review.

The simplest family of EDAs is based on the assumption that  $p_u(\mathbf{x})$  factorizes as a product of  $n$  univariate and mutually independent probability distributions, one for each  $X_i$ . Obviously, this is very far from what happens in difficult optimization problems, where relationships between the unidimensional random variables in  $\mathbf{X}$  usually exist. However, this assumption simplifies learning the probabilistic model for the factorization of  $p_u(\mathbf{x})$  from  $\mathbf{d}_u$ , as this process reduces to parameter fitting. Some examples of this approach are (Baluja, 1994; Kvasnicka et al., 1996; Mühlenbein, 1997; Santana and Ochoa, 1999).

A slightly more complex family of EDAs consists of those algorithms that take into account only bivariate dependencies between the unidimensional random variables in  $\mathbf{X}$  for the factorization of  $p_u(\mathbf{x})$ . Therefore, it is enough to use second order statistics to learn the probabilistic model for such a factorization. Some members of this family of EDAs are (Baluja and Davies, 1997; Baluja and Davies, 1998; De Bonet et al., 1997; Pelikan and Mühlenbein, 1999).

With the aim of improving performance, some researchers have proposed several instances of the generic EDA that involve statistics of order greater than two in the factorization of  $p_u(\mathbf{x})$ . See, for instance, (Harik, 1999; Mühlenbein et al., 1999; Soto et al., 1999). However, the most relevant research within this approach is based on Bayesian networks (Castillo et al., 1997; Cowell et al., 1999; Jensen, 2001; Lauritzen, 1996; Pearl, 1988), so that learning  $p_u(\mathbf{x})$  from  $\mathbf{d}_u$  reduces to learning a Bayesian network for  $\mathbf{X}$  from  $\mathbf{d}_u$ . As a result, the factorization of  $p_u(\mathbf{x})$  corresponds to the graphical factorization represented by the induced Bayesian network for  $\mathbf{X}$ . For a thorough discussion of these EDAs, the reader is referred to, for instance, (Etzeberria and Larrañaga, 1999; Larrañaga and (eds.), 2001; Larrañaga et al., 2000; Pelikan, 2002; Pelikan et al., 1999).

### 2.3 Globally Multimodal Problem Optimization

Many optimization problems are globally multimodal and, often, it is necessary or desirable to identify as many global optima as possible. In this scenario, classical EAs are ineffective, as they converge to at best a single global optima. The explanation is straightforward. When optimizing a globally multimodal problem, the basins of different global optima may be represented in the population. As there is no significant selective preference for one of the basins in the population over another, the stochastic variations due to the selection method make the population drift towards one of them and, thus, discover only one global optimum at most. Moreover, this global optimum is randomly chosen from the existing global optima. This phenomenon is known as *genetic drift* (De Jong, 1975; Goldberg, 1989; Goldberg and Segrest, 1987). In the absence of selective pressure, the stochastic nature of the selection method reduces population diversity.

Globally multimodal optimization problems are challenging for classical EAs not only in terms of effectiveness but also in terms of efficiency. The existence of several global peaks makes convergence speed slow down until the population drifts to one of the global peaks (Pelikan and Goldberg, 2000). Basically, the difficulties appear because combining good solutions coming from different parts of the search space or basins often results in poor solutions. The only mechanism that classical EAs have to make the population drift towards a single basin is genetic drift, but this phenomenon normally occurs very slowly. Therefore, there is not only quantitative, but also qualitative, interest in obtaining several global optima of globally multimodal optimization problems

by preventing genetic drift as much as possible is. Much attention has been devoted to the study of genetic drift as a cause of suboptimal convergence in classical EAs regarding, mainly, convergence time (Goldberg and Segrest, 1987), niching (Horn, 1993) and population sizing (Mahfoud, 1994). However, few works exist where some classical EAs have been modified in order to alleviate genetic drift as much as possible for globally multimodal problem optimization (Hocaoğlu and Sanderson, 1997).

EDAs should encounter exactly the same difficulties for globally multimodal problem optimization as classical EAs do, unless the class of probabilistic models that factorize  $p_u(\mathbf{x})$  is flexible enough to model simultaneously the different basins that may be represented in  $\mathbf{d}_u$ . One way to guarantee this is by using probabilistic models that are able to encode conditional dependencies between the unidimensional random variables in  $\mathbf{X}$ , i.e., conditional dependencies between the random variables corresponding to genes. Alternatively, those EDAs that do not model conditional dependencies between the unidimensional random variables in  $\mathbf{X}$  can perform well in globally multimodal problem optimization by incorporating *niching* (Goldberg, 1989; Goldberg and Richardson, 1987), i.e., the population is distributed in niches or subpopulations, in order to avoid combining solutions coming from different basins. Both approaches are suggested in (Pelikan and Goldberg, 2000), although the authors evaluate only the latter for symmetrical (globally multimodal) problem optimization. Basically, this paper implements niching based on partitional data clustering via the  $K$ -means algorithm (Anderberg, 1973; Hartigan, 1975) within one of the simplest EDAs, namely the univariate marginal distribution algorithm (Larrañaga and (eds.), 2001; Mühlenbein, 1997). In (Gallagher et al., 1999), the population-based incremental learning algorithm is extended to continuous problem optimization by learning and sampling a finite mixture model. This EDA can be seen as another example of niching, since each component in the mixture can represent a different niche. Unfortunately, the authors do not provide much evidence on the benefits of their algorithm for globally multimodal problems.

### 3 An Estimation of Distribution Algorithm Based on Unsupervised Learning of Bayesian Networks

The previous section has outlined two approaches to alleviate the poor performance of most EDAs for globally multimodal problem optimization: Either using probabilistic models that are able to encode conditional dependencies, or incorporating niching (e.g., based on partitional data clustering). Although these two approaches are apparently unrelated, they can be easily combined if data clustering is viewed from a model-based perspective and the class of models considered can encode conditional dependencies. Such a combined approach may benefit from the strengths of both original approaches and increase robustness and reliability on the problem optimization process. A sensible implementation of this consists in using unsupervised learning of Bayesian networks (Peña, 2001; Peña et al., 2001a; Peña et al., 1999; Peña et al., 2000; Peña et al., 2002). This section first introduces unsupervised learning of Bayesian networks and then shows how this can be incorporated into the EDA framework for effective and efficient globally multimodal problem optimization.

#### 3.1 Unsupervised Learning of Bayesian Networks

*Data clustering* is one of the main problems that arises in a great variety of fields (Anderberg, 1973; Duda and Hart, 1973; Hartigan, 1975; McLachlan and Basford, 1988; Peña, 2001). Given some data  $\mathbf{d}$  in the form of a set of instances with an underlying group-structure, data clustering may be roughly defined as the search for the best description

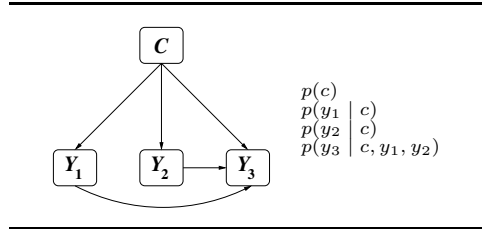


Figure 2: Structure and conditional probability distributions of a BN for data clustering for  $\mathbf{X} = (C, \mathbf{Y}) = (C, Y_1, Y_2, Y_3)$ .

of this group-structure, when the group membership of every instance is unobserved. Each of the groups in  $\mathbf{d}$  is called a *cluster*. The lack of knowledge of the cluster membership of every instance in  $\mathbf{d}$  makes data clustering also be referred to as *unsupervised learning*.

Most solutions to data clustering problems can be classified as being either *partitional*, *hierarchical* or *model-based*. Partitional and hierarchical approaches describe the group-structure underlying  $\mathbf{d}$  as a partition of  $\mathbf{d}$  or as a sequence of tree-like nested partitions of  $\mathbf{d}$ , respectively. On the other hand, model-based approaches describe the group-structure underlying  $\mathbf{d}$  through a probabilistic model induced from  $\mathbf{d}$ . In this paper, we take a model-based approach to data clustering. In particular, we assume that  $\mathbf{d}$  contains  $N$  instances or cases, i.e.,  $\mathbf{d} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ . The  $l$ -th case of  $\mathbf{d}$  is represented by an  $(n + 1)$ -dimensional discrete vector  $\mathbf{x}_l = (x_{l1}, \dots, x_{ln+1})$  partitioned as  $\mathbf{x}_l = (c_l, \mathbf{y}_l)$  for all  $l$ :  $c_l$  is the unobserved cluster membership, and  $\mathbf{y}_l = (y_{l1}, \dots, y_{ln})$  is the vector of observations or *predictive attributes*. We assume as well that the number of clusters underlying  $\mathbf{d}$ , denoted by  $K$ , is known. From a model-based perspective, every case in  $\mathbf{d}$  can be seen as a partial instance of an  $(n + 1)$ -dimensional discrete random variable  $\mathbf{X} = (X_1, \dots, X_{n+1})$  partitioned as  $\mathbf{X} = (C, \mathbf{Y})$ :  $C$  is a unidimensional discrete random variable representing the unobserved cluster membership, i.e., the *cluster random variable*, and  $\mathbf{Y} = (Y_1, \dots, Y_n)$  is an  $n$ -dimensional discrete random variable representing the set of predictive attributes, i.e., the *predictive random variable*. Therefore, model-based data clustering can be solved by learning a joint probability distribution for  $\mathbf{X}$  from  $\mathbf{d}$ . One of the paradigms especially well suited for such a purpose are Bayesian networks (Castillo et al., 1997; Cowell et al., 1999; Jensen, 2001; Lauritzen, 1996; Pearl, 1988).

Let  $\mathbf{X} = (C, \mathbf{Y})$  be a random variable as stated above. A *Bayesian network* (BN) for data clustering for  $\mathbf{X}$  is a pair  $(s, \theta)$ , where  $s$  is the *model structure* and  $\theta$  are the *model parameters* (Peña, 2001; Peña et al., 2001a; Peña et al., 1999; Peña et al., 2000; Peña et al., 2002). The model structure  $s$  is an acyclic directed graph whose nodes correspond to the unidimensional random variables in  $\mathbf{X}$ . Throughout the text, the terms node and random variable are used interchangeably. The model parameters  $\theta$  specify a conditional probability distribution for each node  $X_i$  in  $s$  given its parents  $\mathbf{Pa}_i$  in  $s$ ,  $p(x_i | \mathbf{pa}_i)$ . These conditional probability distributions are all typically multinomial.

A BN for data clustering  $(s, \theta)$  for  $\mathbf{X}$  represents a joint probability distribution for  $\mathbf{X}$ ,  $p(\mathbf{x})$ , through the following graphical factorization:

$$p(\mathbf{x}) = \prod_{i=1}^{n+1} p(x_i | \mathbf{pa}_i). \quad (1)$$

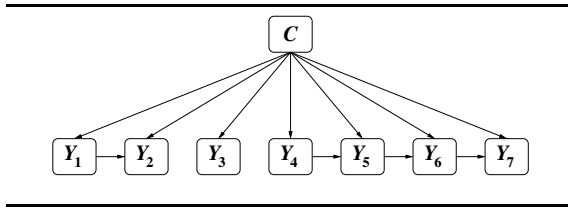


Figure 3: Structure of a TANB model for data clustering.

Therefore,  $s$  encodes a set of conditional (in)dependencies between the random variables in  $X$ . Moreover,  $s$  is usually constrained so that every  $Y_i$  is a child of  $C$ , i.e.,  $C \in \mathbf{Pa}_i$  for all  $i > 1$ . This restriction is imposed by the assumption that  $C$  has an impact on the joint probability distribution for  $Y$ . See Fig. 2 for an example of a BN for data clustering.

In this paper, we interpret unsupervised learning of BNs as an optimization problem. This is a challenging optimization problem in general. As a matter of fact, it has been proven in (Chickering, 1996) that the identification of the BN structure with the highest Bayesian Dirichlet equivalent score (Heckerman et al., 1995) among all the BN structures in which every node has no more than  $t$  parents is an NP-hard optimization problem for  $t > 1$ . It is usually assumed that this hardness holds for other common scores as well, though there is not yet a formal proof (Chickering, 2002). These results also apply to unsupervised learning of BNs.

As search space, we consider the space of structures of BNs for data clustering. This space can be restricted to the space of DAGs for  $Y$ , due to the fact that every  $Y_i$  is a child of  $C$ . Alternative search spaces include the space of equivalence classes of structures of BNs for data clustering (Chickering, 2002; Nielsen et al., 2003), and the space of ancestral orderings of structures of BNs for data clustering (Friedman and Koller, 2003; Larrañaga et al., 1996). Note that, as usual, model parameter fitting is considered a secondary optimization problem: Given a BN structure for data clustering, maximum likelihood (ML) or maximum a posteriori model parameter estimates can be effectively obtained via approximation techniques such as gradient descent methods (Binder et al., 1997), Gibbs sampling (Geman and Geman, 1984) or the EM algorithm (Dempster et al., 1977; McLachlan and Krishnan, 1997). In some cases, it may be desirable to restrict the attention to those BNs for data clustering that trade off expressivity for simplicity. This reduces the search space while the models considered are still expressive enough. One such example is the class of *tree augmented naive Bayes (TANB) models for data clustering* (Meilă, 1999; Peña, 2001; Peña et al., 2001a; Peña et al., 2000). TANB models for data clustering include those BNs for data clustering such that every  $Y_i$  has at most one other unidimensional predictive random variable in  $Y$  as a parent. See Fig. 3 for an example. TANB models have been also used in data classification (Friedman et al., 1997; Keogh and Pazzani, 1999).

As search strategy, we consider the *Bayesian structural EM (BSEM) algorithm* (Friedman, 1998). The BSEM algorithm relies on following the basic intuition of the iterations of the EM algorithm: Take advantage of the best estimate of the joint probability distribution found so far in order to compute quantities of interest that can not be directly obtained from the data at hand, and then to use effective and efficient model learning algorithms for complete data. Fig. 4 shows a pseudocode of the BSEM algorithm for unsupervised learning of BNs. The BSEM algorithm iterates between two main steps

- 
1. Let  $s_1$  be the initial model structure
  2. **for**  $u = 1, 2, \dots$  **do**
  3.   Run the EM algorithm in order to approximate the ML parameters  $\hat{\theta}_{s_u}$  for  $s_u$
  4.   Perform a greedy hill-climbing search over model structures, evaluating each one by  
 $Sc(s : s_u, d) = E[\log L(d | s) | d^Y, \hat{\theta}_{s_u}, s_u] = \sum_{d^C} \log L(d^C, d^Y | s) L(d^C | d^Y, \hat{\theta}_{s_u}, s_u)$
  5.   Let  $s_{u+1}$  be the model structure with the highest score among those visited in step 4
  6.   **if**  $Sc(s_{u+1} : s_u, d) = Sc(s_u : s_u, d)$  **then**
  7.     Return  $(s_u, \hat{\theta}_{s_u})$
- 

Figure 4: Pseudocode of the BSEM algorithm for unsupervised learning of BNs.

- 
1. Let  $po_1$  be a population composed of  $Q$  uniformly generated individuals
  2. Evaluate the individuals in  $po_1$
  3.  $u = 1$
  4. **while** the stopping condition is not met **do**
  5.   Let  $d_u$  group  $N$  individuals selected from  $po_u$  via the selection method
  6.   Let  $(s_u, \hat{\theta}_{s_u})$  be the BN for data clustering learnt from  $d_u$  via the BSEM algorithm
  7.   Let  $(s_u, \bar{\theta}_{s_u})$  be  $(s_u, \hat{\theta}_{s_u})$  with the exception that  $p(c)$  has been modified to be uniform
  8.   Let  $of_u$  be the offspring population composed of  $M$  individuals sampled from  $(s_u, \bar{\theta}_{s_u})$
  9.   Evaluate the individuals in  $of_u$
  10.   Let  $po_{u+1}$  be the population created from  $po_u$  and  $of_u$  via the replacement method
  11.    $u + 1$
  12. Return the best individuals found so far
- 

Figure 5: Pseudocode of the UEBNA.

that are as follows for the  $u$ -th iteration for all  $u$ . The first step (step 3 in Fig. 4) approximates the ML parameters  $\hat{\theta}_{s_u}$  for the current model structure  $s_u$  given the observed data  $d^Y$ , usually via the EM algorithm. On the other hand, the second step (step 4 in Fig. 4) searches for the highest scoring model structure in order to replace the current one. This latter step is usually solved through a greedy hill-climbing search considering all the possible additions, removals and reversals of a single arc at each point in the search. The score that guides the search is usually the expected  $\log L(d | s)$ , where  $s$  is the model structure being evaluated and the expectation is taken with respect to  $d^Y$ ,  $s_u$  and  $\hat{\theta}_{s_u}$ :

$$\begin{aligned} Sc(s : s_u, d) &= E[\log L(d | s) | d^Y, \hat{\theta}_{s_u}, s_u] \\ &= \sum_{d^C} \log L(d^C, d^Y | s) L(d^C | d^Y, \hat{\theta}_{s_u}, s_u) \end{aligned} \quad (2)$$

where  $d^C$  denotes a labelling or completion of  $d$ . Note that this score requires going through every possible completion  $d^C$  of  $d$ , which may be prohibitive. Instead, we approximate Eq. (2) by considering only the completion  $d^C$  that scores the highest  $L(d^C | d^Y, \hat{\theta}_{s_u}, s_u)$ . Therefore, the score for the structural search step of each iteration of the BSEM algorithm can be computed in factorable and closed form as reported in (Cooper and Herskovits, 1992; Heckerman et al., 1995).

### 3.2 Unsupervised Estimation of Bayesian Network Algorithm

This section describes the *unsupervised estimation of Bayesian network algorithm* (UEBNA), whose only peculiarity with respect to existing EDAs is being based on unsupervised learning of BNs. As discussed previously, incorporating unsupervised learning of BNs into the EDA framework seems a natural solution to alleviate the poor performance of most EDAs for globally multimodal problem optimization: It allows modelling simultaneously the different basins that may be represented by the individuals selected at each iteration, while preventing genetic drift.



As outlined in Fig. 5, the UEBNA consists in the iteration of the same three main steps as the generic EDA (see Fig. 1): Selection of promising individuals from the current population, probabilistic modelling of the selected individuals, and model sampling to create the new population. The singularities of the  $u$ -th iteration of the UEBNA are as follows for all  $u$ :

- The  $l$ -th case of  $\mathbf{d}_u$  is represented by an  $(n + 1)$ -dimensional discrete vector  $\mathbf{x}_l = (x_{l1}, \dots, x_{ln+1})$  partitioned as  $\mathbf{x}_l = (c_l, \mathbf{y}_l)$  for all  $l$ :  $c_l$  is the unobserved cluster membership, and  $\mathbf{y}_l = (y_{l1}, \dots, y_{ln})$  is the  $l$ -th selected individual. Therefore, every case in  $\mathbf{d}$  can be seen as a partial instance of an  $(n + 1)$ -dimensional discrete random variable  $\mathbf{X} = (X_1, \dots, X_{n+1})$  partitioned as  $\mathbf{X} = (C, \mathbf{Y})$ . This fits the discussion in Section 3.1.
- The BSEM algorithm should be provided with the number of clusters  $K$  underlying  $\mathbf{d}_u$ . In general, the higher the number of clusters, the higher the flexibility and expressivity but also the complexity of the model. Therefore, this parameter may have an impact on the performance of the UEBNA.
- When  $(s_u, \hat{\theta}_{s_u})$  is simulated to obtain  $\mathbf{of}_u$ , the number of individuals produced from each cluster is determined by the probability distribution for the cluster random variable  $p(c)$ . This implies that the number of individuals sampled from each cluster is proportional to its size. This sampling scheme favors large clusters, no matter their average fitness, and makes the population drift towards them. This seems unreasonable because it promotes genetic drift, and the UEBNA is aimed at preventing this phenomenon as much as possible. Yet, the explanation is straightforward. In the absence of selective pressure, i.e., when there is not significant selective preference for a cluster over another, the stochastic nature of the selection method makes some clusters have more representatives than others in the pool of selected individuals. Therefore, sampling more individuals from those clusters that have more members promotes genetic drift. On the other hand, sampling a number of individuals from each cluster that is proportional to its average fitness involves excessive selective pressure, i.e., those clusters with good average fitness are favored by the selection method but also by the sampling scheme, which is undesirable as well. With the purpose of avoiding promoting genetic drift and excessive selective pressure in the UEBNA, it seems justified to sample the same number of individuals from each of the clusters encoded by  $(s_u, \hat{\theta}_{s_u})$ . This is accomplished by first modifying  $p(c)$  in  $(s_u, \hat{\theta}_{s_u})$  to be a uniform distribution and, then, sampling the resulting model, here denoted by  $(s_u, \bar{\theta}_{s_u})$ . A similar discussion can be found in (Pelikan and Goldberg, 2000) regarding the number of individuals that should be sampled from each cluster in an attempt to implement niching.
- $\mathbf{of}_u$  is constructed by restricting the  $M$  instances sampled from  $(s_u, \bar{\theta}_{s_u})$  to their values for  $\mathbf{Y}$ .

## 4 Experimental Evaluation

This section evaluates the UEBNA for symmetrical (globally multimodal) problem optimization. Specifically, the evaluation involves optimization problems that show what is known as *symmetry on the alphabet* or *spin-flip symmetry* (Van Hoyweghen and Naudts, 2000): An optimization problem contains spin-flip symmetry when gene-complementary solutions score the same fitness. Therefore, spin-flip symmetrical optimization problems are globally multimodal. Some optimization problems that show

spin-flip symmetry and, thus, global multimodality are twomax problems, graph partitioning problems, random number partitioning problems and graph coloring problems. As reported in (Naudts and Naudts, 1998; Pelikan and Goldberg, 2000; Van Hoyweghen, 2001; Van Hoyweghen and Naudts, 2000), this class of globally multimodal optimization problems are challenging for most EAs, including EDAs. The evaluation of the UEBNA for some symmetrical optimization problems should provide us with sufficient insight to assess whether or not the UEBNA performs effectively and efficiently for globally multimodal problem optimization.

When one wants to identify several global optima of a spin-flip symmetrical optimization problem, it is tempting to consider searching for just one of them and, then, obtain another global optima by just changing the genes of the global optimum discovered to their complementary values. This approach is discarded in the evaluation of the UEBNA for the following reasons. Firstly, this shortcut is based on the knowledge that the search space contains spin-flip symmetry. However, a key feature of EAs in general and EDAs in particular is that they do not make assumptions about the search space of the optimization problem at hand. Secondly, working in this way only addresses the quantitative side (effectiveness) of the problem optimization process, i.e., the number of global peaks discovered, while ignoring the qualitative side (efficiency), i.e., the convergence speed (see Section 2.3). Moreover, this shortcut performs poorly even in terms of effectiveness if more than two global optima exist in the spin-flip symmetrical optimization problem at hand. Finally, this approach can not be applied to globally multimodal problem optimization in general. In other words, the knowledge of the search space being symmetrical is not used at all in the evaluation of the UEBNA.

This section first describes the evaluation setup and, then, presents the symmetrical optimization problems in the evaluation. Finally, the dynamics and performance of the UEBNA in these optimization problems are reported and discussed.

#### 4.1 Evaluation Setup

The BSEM algorithm run at each iteration of the UEBNA restricts the search to TANB models for data clustering. Furthermore, the BSEM algorithm should be provided with the number of clusters  $K$  underlying the set of selected individuals. As noted earlier, this can be seen as a parameter that sets the flexibility and expressive power of the models in the search space of the BSEM algorithm. In the evaluation, we consider different values of  $K$  in order to assess the impact of this parameter on the performance of the UEBNA. We start with  $K = 2$  and increase it until no further improvement is observed.

The convergence criterion for the EM algorithm run at each iteration of the BSEM algorithm is satisfied when either the relative difference between successive values for  $\log L(\mathbf{d} \mid \boldsymbol{\theta}_s, s)$  is less than 1 or 150 iterations are reached. Preliminary experiments with more demanding convergence criteria did not lead to significantly better results.

In the implementation of the UEBNA, we use HUGIN API version 3.1 (Jensen, 1997) wherever probabilistic inference or sampling of TANB models for data clustering is required (e.g., steps 3 and 4 in Fig. 4 and step 8 in Fig. 5). This means that probabilistic inference is done as indicated in (Jensen et al., 1990; Lauritzen and Spiegelhalter, 1988).

For comparison purposes, we benchmark the UEBNA against two well established EDAs, namely the *univariate marginal distribution algorithm* (UMDA) (Larrañaga and (eds.), 2001; Mühlenbein, 1997) and the *estimation of Bayesian network algorithm* (EBNA) (Etzeberria and Larrañaga, 1999; Larrañaga and (eds.), 2001; Larrañaga et al., 2000).

The UMDA is based on the assumption that  $p_u(\mathbf{x})$  factorizes as follows:

$$p_u(\mathbf{x}) = \prod_{i=1}^n p_u(x_i) \quad (3)$$

for all  $u$ . Moreover,  $p_u(x_i)$  is restricted to be a univariate multinomial distribution whose parameters are estimated from  $\mathbf{d}_u$  according to the ML criterion for all  $i$ . On the other hand, the EBNA reduces learning  $p_u(\mathbf{x})$  from  $\mathbf{d}_u$  to induction of a BN for  $\mathbf{X}$  from  $\mathbf{d}_u$  for all  $u$ . For this purpose, the EBNA runs a greedy hill-climbing search over BN structures for  $\mathbf{X}$  considering all the possible additions, removals and reversals of a single arc at each point in the search. The score that guides the search is the Bayesian information criterion (BIC) (Schwarz, 1978). The sample from the BN for  $\mathbf{X}$  learnt at each iteration of the EBNA is obtained by probabilistic logic sampling (Henrion, 1988).

The reasons for using the UMDA and the EBNA as benchmarks in the evaluation of the UEBNA are the following ones. First, both the UMDA and the EBNA have received much attention in the literature. Moreover, the EBNA is close in spirit to the UEBNA, as both are based on learning and simulation of BNs. Finally, the UMDA and the EBNA provide the opportunity to compare the performance of three different approaches for globally multimodal problem optimization: The UMDA is an EDA that neither encodes conditional dependencies nor implements niching, the EBNA is an EDA that can encode conditional dependencies but it does not use niching, and the UEBNA is an EDA that combines encoding of conditional dependencies with niching via unsupervised learning of BNs (see Section 2.3 and Section 3).

The three EDAs in the evaluation use *truncation selection* as the selection method, i.e., the most fit individuals in the current population are selected. Furthermore, the replacement method creates the new population by replacing the least fit individuals in the current one by all the offspring population. The algorithms stop when the relative difference between the sum of the objective function values of all the individuals of the population of two successive generations is 0 or 100 generations are reached. For the three EDAs in the evaluation, the population size, the number of selected individuals at each iteration, and the number of generated individuals at each iteration are 4000, 3000 and 3000, respectively. Preliminary experiments confirmed that these parameter values are well suited for the three EDAs in evaluation and that they do not favor any of the them over the rest. Having said this, using the same optimization schedule for all the EDAs in the evaluation eases comparison.

For each pair composed of one EDA and one globally multimodal optimization problem in the evaluation, the performance criteria measured are (i) the number of global optima discovered, (ii) the average deviation with respect to the expected number of individuals representing each global optima discovered, and (iii) the number of evaluations of the objective function and the runtime until convergence. The first two criteria reflect the effectiveness of the problem optimization process, while the last two criteria assesses its efficiency. The second performance criterion is calculated as follows:

$$\frac{1}{Op} \cdot \sum_{i=1}^{Op} \frac{|Q/Op - Q_i^*|}{Q/Op} \cdot 100 \quad (4)$$

where  $|\cdot|$  is the absolute value function,  $Op$  is the number of global optima captured,  $Q$  is the population size, and  $Q_i^*$  is the number of individuals in the population of the last generation representing the  $i$ -th global optima discovered. As there is not selective preference for a global peak over another, it is desirable that those global optima

identified are equally well represented in the population of the last generation, i.e.,  $Q_i^* \approx Q/Op$  for all  $i$ . Significant underrepresentation or overrepresentation of one or several of the global peaks discovered should be detected. Therefore, the closer the value of Eq. (4) to 0, the better. Likewise, the closer the value to 100, the worse.

## 4.2 Symmetrical Optimization Problems

The paragraphs below present the symmetrical optimization problems in the evaluation in terms of their search spaces and objective functions. These symmetrical optimization problems have been borrowed or adapted from (Pelikan and Goldberg, 2000; Pelikan et al., 2001).

### 4.2.1 Twomax Problem

The twomax problem is a simple symmetrical optimization problem whose search space is  $\{0, 1\}^n$ , i.e., the set of binary strings of length  $n$ , and whose objective function is as follows:

$$F_{twomax}(z) = F_{twomax}(z_1, \dots, z_n) = \left| \frac{n}{2} - \sum_{i=1}^n z_i \right|. \quad (5)$$

The objective is maximization and there are two global optima:  $z_1^* = (0, \dots, 0)$  and  $z_2^* = (1, \dots, 1)$  with fitness equal to  $\frac{n}{2}$ . In all the EDAs in the evaluation, every solution  $z$  is represented by an  $n$ -dimensional binary individual where the  $i$ -th gene coincides with  $z_i$  for all  $i$ . The evaluation involves two instances of the twomax problem with  $n = 50, 100$  and denoted by  $P_{twomax50}$  and  $P_{twomax100}$ , respectively.

### 4.2.2 Graph Bisection Problem

The graph bisection problem aims to split the set of nodes of a given graph into two equally sized subsets so that the number of edges between the two subsets is minimized. Consequently, the search space of the graph bisection problem is the set of all the partitions of the nodes of the given graph into two equally sized subsets. The fitness of a given solution is calculated as the number of nodes in the graph at hand minus the number of edges connecting the two subsets of nodes in the solution. Thus, the objective is maximization. In all the EDAs in the evaluation, every solution  $z$  is represented by an  $n$ -dimensional binary individual where the  $i$ -th gene corresponds to the  $i$ -th node of the graph for bisection for all  $i$ . Then, each gene of a given individual classifies one of the nodes of the graph into one of the two subsets. Under this codification, the search space of the graph bisection problem can be represented as the set  $\{(z_1, \dots, z_n) \mid (z_1, \dots, z_n) \in \{0, 1\}^n \text{ and } \sum_{i=1}^n z_i = \frac{n}{2}\}$ , where  $n$  is the number of nodes in the graph at hand. Note that only individuals with equal number of zeroes and ones represent feasible solutions. However, the generation of the offspring at each iteration of the EDAs in the evaluation is not a closed operation with respect to this feasibility condition. Thus, some individuals that may appear during the problem optimization process may need to be repaired. A simple randomized repair operator is used in the EDAs in the evaluation: An unfeasible solution is converted into a feasible one by, iteratively, picking at random a gene in the majority and changing it to its complementary value until a feasible solution is obtained.

The evaluation involves 10 instances of the graph bisection problem. The first three instances consist of three grid-like graphs, with  $n = 16, 36, 64$ , cut in halves and connected by two edges. There are two global optima with fitness equal to  $n - 2$ . In the following, these optimization problems are denoted by  $P_{grid16}$ ,  $P_{grid36}$  and  $P_{grid64}$ , respectively. The evaluation also involves three so-called caterpillar graphs, with sizes

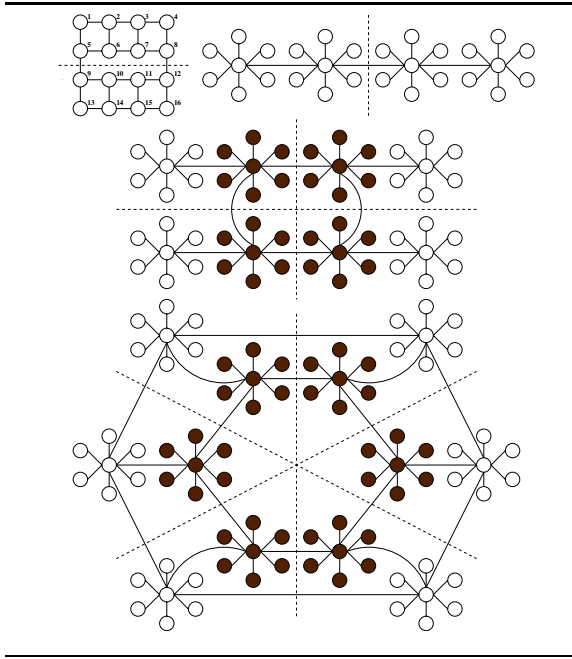


Figure 6: Graphs for  $P_{grid16}$  (top, left),  $P_{cat28}$  (top, right),  $P_{catring28}$  (middle, only dark nodes),  $P_{catring56}$  (middle, all the nodes),  $P_{catring42}$  (bottom, only dark nodes) and  $P_{catring84}$  (bottom, all the nodes). Dashed lines indicate optimal cuts.

$n = 28, 42, 56$ , composed of four, six and eight, respectively, seven node star-shaped graphs connected in a line. There are two global optima with fitness equal to  $n - 1$ . In the following, these optimization problems are referred to as  $P_{cat28}$ ,  $P_{cat42}$  and  $P_{cat56}$ , respectively. The last four instances of the graph bisection problem involve extensions of the caterpillar graphs so that there are more than two global optima.  $P_{catring28}$  and  $P_{catring56}$  involve graphs with  $n = 28, 56$ , respectively, and have four global optima with fitness equal to  $n - 2$ . On the other hand,  $P_{catring42}$  and  $P_{catring84}$  involve graphs with  $n = 42, 84$ , respectively, and have six global optima whose fitness is equal to  $n - 4$ . Fig. 6 shows most of the graphs for bisection. In addition to the difficulties derived from their symmetrical nature, these instances of the graph bisection problem present another source of difficulties: They are highly multimodal and present many local optima and only a few global optima (Pelikan and Goldberg, 2000; Pelikan et al., 2001; Schwarz and Ocenasek, 1999).

### 4.3 Results

Fig. 7 shows the dynamics of the UMDA, the EBNA and the UEBNA until convergence in one run for  $P_{twomax50}$ . The histograms summarize the number of solutions (vertical axis) in the population of different generations whose sum of genes is equal to the value of the horizontal axis. As previously stated, the two global optima of the optimization problem are gene-complementary and correspond to the left-most and right-most sides of the histograms. The histograms show that, as problem optimization progresses, the population drifts to one side in the case of the UMDA and to both sides in the case of the EBNA and the UEBNA. Moreover, the individuals of the population of the last gen-

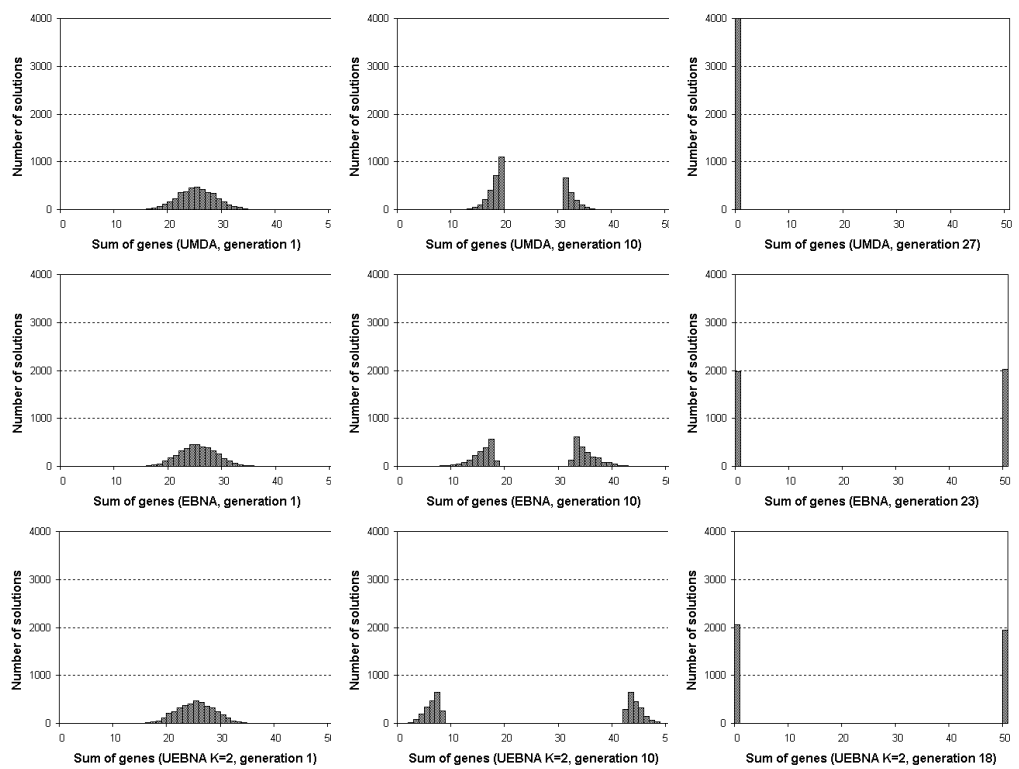


Figure 7: Dynamics of the UMDA (top row), the EBNA (middle row), and the UEBNA (bottom row) until convergence in one run for  $P_{twomax50}$ . The horizontal axis of each histogram represents the sum of the genes of a solution, whereas the vertical axis denotes the number of corresponding solutions in the population of different generations.

eration of the EBNA and the UEBNA are almost equally distributed among both global peaks. It can also be seen in the histograms that genetic drift occurs so slowly that the UMDA takes longer than the other two EDAs to converge. This clearly confirms what has been argued in Section 2.3 about the necessity of considering EDAs based on either encoding of conditional dependencies or niching, or both, for effective and efficient globally multimodal problem optimization. Finally, it should also be mentioned that, although the EBNA and the UEBNA perform equally well in terms of effectiveness, they differ in their efficiency: The UEBNA reaches convergence faster than the EBNA. This suggests that TANB models for data clustering are more appropriate than BNs for modelling the joint probability distribution for the individuals selected at each iteration. Note, however, that the EBNA relies on unrestricted BNs, which can potentially model more complex conditional dependencies than TANB models for data clustering. Therefore, this supports that combining model-based data clustering with the ability to model conditional dependencies is more robust and reliable against genetic drift when incorporated into the EDA framework than the ability to model conditional dependencies alone. As discussed in Section 3, this is the main motivation behind the development of the UEBNA.

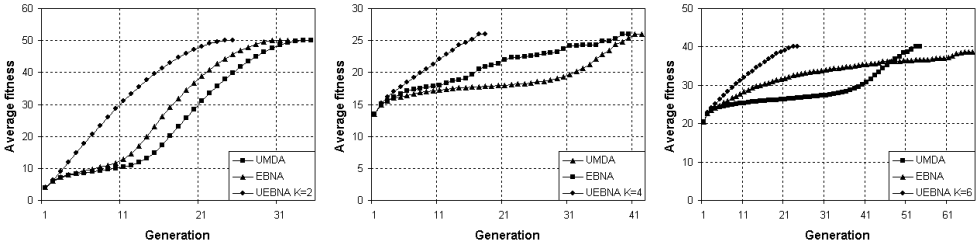


Figure 8: Average fitness of the individuals in the population as a function of the number of generations and until convergence for the UMDA, the EBNA and the UEBNA in one run for  $P_{twomax100}$  (left),  $P_{catring28}$  (middle) and  $P_{catring42}$  (right).

Fig. 8 provides the reader with additional evidence on the efficiency of the UEBNA for globally multimodal problem optimization. Specifically, the figure plots the average fitness of the individuals in the population as a function of the number of generations and until convergence for the UMDA, the EBNA and the UEBNA in one run for  $P_{twomax100}$ ,  $P_{catring28}$  and  $P_{catring42}$ . Note that the first optimization problem presents two global optima, the second four and the third six. These curves clearly show that the UEBNA speeds up converge without degrading the quality of the solutions obtained.

Figs. 7 and 8 illustrate the behavior of the UMDA, the EBNA and the UEBNA for a single run for  $P_{twomax50}$ ,  $P_{twomax100}$ ,  $P_{catring28}$  and  $P_{catring42}$ . The remainder of the 10 independent runs performed for these optimization problems leads to the same conclusions as those discussed above. Moreover, the observed patterns can be extended to the 10 independent runs performed for the rest of the symmetrical optimization problems in the evaluation. For the sake of brevity, figures are not reported. Instead, Tables 1 and 2 summarize the results that the UMDA, the EBNA and the UEBNA (with different values for  $K$ ) achieve for each of the 12 symmetrical optimization problems in the evaluation. For each combination of one EDA and one symmetrical optimization problem in the evaluation, the tables report performance in terms of average and standard deviation (i) over the 10 independent runs performed (*All runs*), and (ii) over successful runs (*Successful runs*), i.e., over those runs out of the 10 independent runs performed where at least one global peak of the symmetrical optimization problem at hand is identified. For *All runs*, the performance criteria measured are the number of global optima discovered (*Optima*), the number of evaluations of the objective function until convergence (*Eval*) and the runtime in seconds until convergence (*Time*).<sup>1</sup> For *Successful runs*, the average deviation with respect to the expected number of individuals representing each global optima discovered (*Deviation*) is calculated as indicated in Eq. (4) and reported, in addition to *Optima*, *Eval* and *Time*. *Optima* and *Deviation* relate to the effectiveness of the EDAs, while *Eval* and *Time* related to the efficiency. Finally, it should be mentioned that, when all the 10 independent runs performed for any of the symmetrical optimization problems in the evaluation are successful, the values for the performance criteria for *All runs* and *Successful runs* coincide. In this case, only the values for the performance criteria for *Successful runs* are reported, for the sake of readability.

Table 1 summarizes the effectiveness of the EDAs in the evaluation. The first conclusion that can be achieved from the results in the table is that the UEBNA enjoys

<sup>1</sup> All the experiments are run on a Pentium 900 MHz.

Table 1: Effectiveness of the UMDA, the EBNA and the UEBNA for the 12 symmetrical optimization problems in the evaluation. All the values are given in terms of average and standard deviation over 10 independent runs.

Problem	EDA	All runs	Successful runs	
		Optima $\pm$ sd	Optima $\pm$ sd	Deviation $\pm$ sd
$P_{twomax50}$ (2 global optima)	UMDA	— $\pm$ —	1.0 $\pm$ 0.0	0 $\pm$ 0
	EBNA	— $\pm$ —	1.5 $\pm$ 0.5	24 $\pm$ 37
	UEBNA $K=2$	— $\pm$ —	2.0 $\pm$ 0.0	1 $\pm$ 1
	UEBNA $K=4$	— $\pm$ —	2.0 $\pm$ 0.0	1 $\pm$ 1
$P_{twomax100}$ (2 global optima)	UMDA	— $\pm$ —	1.0 $\pm$ 0.0	0 $\pm$ 0
	EBNA	— $\pm$ —	1.0 $\pm$ 0.0	0 $\pm$ 0
	UEBNA $K=2$	— $\pm$ —	2.0 $\pm$ 0.0	1 $\pm$ 1
	UEBNA $K=4$	— $\pm$ —	2.0 $\pm$ 0.0	1 $\pm$ 1
$P_{grid16}$ (2 global optima)	UMDA	— $\pm$ —	1.0 $\pm$ 0.0	0 $\pm$ 0
	EBNA	— $\pm$ —	2.0 $\pm$ 0.0	34 $\pm$ 32
	UEBNA $K=2$	— $\pm$ —	2.0 $\pm$ 0.0	2 $\pm$ 2
	UEBNA $K=4$	— $\pm$ —	2.0 $\pm$ 0.0	2 $\pm$ 2
$P_{grid36}$ (2 global optima)	UMDA	0.7 $\pm$ 0.5	1.0 $\pm$ 0.0	13 $\pm$ 35
	EBNA	— $\pm$ —	1.8 $\pm$ 0.4	98 $\pm$ 4
	UEBNA $K=2$	— $\pm$ —	2.0 $\pm$ 0.0	2 $\pm$ 2
	UEBNA $K=4$	— $\pm$ —	2.0 $\pm$ 0.0	4 $\pm$ 6
$P_{grid64}$ (2 global optima)	UMDA	0.2 $\pm$ 0.4	1.0 $\pm$ 0.0	0 $\pm$ 0
	EBNA	0.9 $\pm$ 0.6	1.1 $\pm$ 0.4	47 $\pm$ 51
	UEBNA $K=2$	1.4 $\pm$ 1.0	2.0 $\pm$ 0.0	9 $\pm$ 18
	UEBNA $K=4$	— $\pm$ —	2.0 $\pm$ 0.0	7 $\pm$ 14
$P_{cat28}$ (2 global optima)	UMDA	0.9 $\pm$ 0.3	1.0 $\pm$ 0.0	0 $\pm$ 0
	EBNA	— $\pm$ —	2.0 $\pm$ 0.0	51 $\pm$ 30
	UEBNA $K=2$	— $\pm$ —	2.0 $\pm$ 0.0	3 $\pm$ 3
	UEBNA $K=4$	— $\pm$ —	2.0 $\pm$ 0.0	2 $\pm$ 1
$P_{cat42}$ (2 global optima)	UMDA	0.6 $\pm$ 0.5	1.0 $\pm$ 0.0	0 $\pm$ 0
	EBNA	1.2 $\pm$ 1.0	2.0 $\pm$ 0.0	85 $\pm$ 22
	UEBNA $K=2$	— $\pm$ —	2.0 $\pm$ 0.0	2 $\pm$ 1
	UEBNA $K=4$	— $\pm$ —	2.0 $\pm$ 0.0	1 $\pm$ 1
$P_{cat56}$ (2 global optima)	UMDA	0.5 $\pm$ 0.5	1.0 $\pm$ 0.0	0 $\pm$ 0
	EBNA	0.2 $\pm$ 0.6	2.0 $\pm$ 0.0	99 $\pm$ 0
	UEBNA $K=2$	— $\pm$ —	2.0 $\pm$ 0.0	9 $\pm$ 22
	UEBNA $K=4$	— $\pm$ —	2.0 $\pm$ 0.0	2 $\pm$ 2
$P_{catring28}$ (4 global optima)	UMDA	— $\pm$ —	1.0 $\pm$ 0.0	0 $\pm$ 0
	EBNA	2.8 $\pm$ 1.3	3.1 $\pm$ 0.9	81 $\pm$ 13
	UEBNA $K=2$	— $\pm$ —	4.0 $\pm$ 0.0	51 $\pm$ 9
	UEBNA $K=4$	— $\pm$ —	4.0 $\pm$ 0.0	28 $\pm$ 17
	UEBNA $K=6$	— $\pm$ —	4.0 $\pm$ 0.0	6 $\pm$ 12
$P_{catring56}$ (4 global optima)	UMDA	0.7 $\pm$ 0.5	1.0 $\pm$ 0.0	0 $\pm$ 0
	EBNA	1.0 $\pm$ 0.9	1.7 $\pm$ 0.5	94 $\pm$ 11
	UEBNA $K=2$	— $\pm$ —	2.9 $\pm$ 0.9	43 $\pm$ 39
	UEBNA $K=4$	— $\pm$ —	3.4 $\pm$ 0.7	6 $\pm$ 14
	UEBNA $K=6$	— $\pm$ —	3.7 $\pm$ 0.5	2 $\pm$ 1
	UEBNA $K=8$	— $\pm$ —	3.8 $\pm$ 0.4	2 $\pm$ 1
$P_{catring42}$ (6 global optima)	UMDA	— $\pm$ —	1.0 $\pm$ 0.0	0 $\pm$ 0
	EBNA	2.9 $\pm$ 1.7	3.2 $\pm$ 1.4	79 $\pm$ 17
	UEBNA $K=2$	— $\pm$ —	5.8 $\pm$ 0.6	51 $\pm$ 22
	UEBNA $K=4$	— $\pm$ —	5.6 $\pm$ 0.7	48 $\pm$ 12
	UEBNA $K=6$	— $\pm$ —	5.9 $\pm$ 0.3	25 $\pm$ 18
	UEBNA $K=8$	— $\pm$ —	5.8 $\pm$ 0.6	19 $\pm$ 22
$P_{catring84}$ (6 global optima)	UMDA	0.8 $\pm$ 0.4	1.0 $\pm$ 0.0	13 $\pm$ 35
	EBNA	0.2 $\pm$ 0.4	1.0 $\pm$ 0.0	50 $\pm$ 71
	UEBNA $K=2$	— $\pm$ —	2.2 $\pm$ 0.4	13 $\pm$ 24
	UEBNA $K=4$	— $\pm$ —	3.7 $\pm$ 0.7	10 $\pm$ 16
	UEBNA $K=6$	— $\pm$ —	4.3 $\pm$ 0.8	3 $\pm$ 1
	UEBNA $K=8$	— $\pm$ —	4.7 $\pm$ 1.0	5 $\pm$ 9
	UEBNA $K=10$	— $\pm$ —	4.8 $\pm$ 0.8	34 $\pm$ 20

higher rate of successful runs than the UMDA and the EBNA. In addition, the average number of global optima identified per run, i.e., *Optima* in *All runs*, indicates that the UEBNA outperforms by far both the UMDA and the EBNA in the 12 symmetrical optimization problems in the evaluation. The poor behavior of the UMDA and the EBNA illustrates that the globally multimodal optimization problems in the evaluation are challenging. The UEBNA behaves very effectively even in those optimization problems with four and six global peaks. Regarding effectiveness per successful run, i.e., *Optima* and *Deviation* in *Successful runs*, the results compiled in the table support what has been discussed in the paragraphs above. The UMDA is ineffective for glob-



Table 2: Efficiency of the UMDA, the EBNA and the UEBNA for the 12 symmetrical optimization problems in the evaluation. All the values are given in terms of average and standard deviation over 10 independent runs.

Problem	EDA	All runs		Successful runs	
		Eval $\pm$ sd	Time $\pm$ sd	Eval $\pm$ sd	Time $\pm$ sd
$P_{twomax50}$ (2 global optima)	UMDA	— $\pm$ —	— $\pm$ —	87100 $\pm$ 11229	157 $\pm$ 103
	EBNA	— $\pm$ —	— $\pm$ —	69700 $\pm$ 1703	416 $\pm$ 12
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	55000 $\pm$ 0	421 $\pm$ 33
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	56200 $\pm$ 2098	1011 $\pm$ 106
$P_{twomax100}$ (2 global optima)	UMDA	— $\pm$ —	— $\pm$ —	117400 $\pm$ 17596	215 $\pm$ 50
	EBNA	— $\pm$ —	— $\pm$ —	98800 $\pm$ 5138	3320 $\pm$ 279
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	76600 $\pm$ 1265	1976 $\pm$ 89
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	79900 $\pm$ 2470	5644 $\pm$ 463
$P_{grid16}$ (2 global optima)	UMDA	— $\pm$ —	— $\pm$ —	109300 $\pm$ 20418	201 $\pm$ 46
	EBNA	— $\pm$ —	— $\pm$ —	113200 $\pm$ 43317	215 $\pm$ 116
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	53500 $\pm$ 2916	156 $\pm$ 21
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	51400 $\pm$ 2366	210 $\pm$ 23
$P_{grid36}$ (2 global optima)	UMDA	217900 $\pm$ 57150	488 $\pm$ 162	200286 $\pm$ 53996	436 $\pm$ 144
	EBNA	— $\pm$ —	— $\pm$ —	244900 $\pm$ 71653	1056 $\pm$ 323
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	85600 $\pm$ 8462	620 $\pm$ 89
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	94000 $\pm$ 6782	909 $\pm$ 64
$P_{grid64}$ (2 global optima)	UMDA	299500 $\pm$ 12268	757 $\pm$ 57	281500 $\pm$ 23335	671 $\pm$ 67
	EBNA	249400 $\pm$ 61103	2801 $\pm$ 620	235750 $\pm$ 61120	2670 $\pm$ 630
	UEBNA $K=2$	128200 $\pm$ 12506	2424 $\pm$ 313	123143 $\pm$ 7690	2286 $\pm$ 226
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	124900 $\pm$ 3479	3809 $\pm$ 483
$P_{cat28}$ (2 global optima)	UMDA	128200 $\pm$ 24008	215 $\pm$ 57	124333 $\pm$ 21915	207 $\pm$ 55
	EBNA	— $\pm$ —	— $\pm$ —	138100 $\pm$ 67765	386 $\pm$ 215
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	57100 $\pm$ 2846	344 $\pm$ 35
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	60700 $\pm$ 949	435 $\pm$ 22
$P_{cat42}$ (2 global optima)	UMDA	175600 $\pm$ 26937	325 $\pm$ 59	166500 $\pm$ 21668	309 $\pm$ 51
	EBNA	238300 $\pm$ 74289	1196 $\pm$ 373	244000 $\pm$ 71875	1212 $\pm$ 344
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	73900 $\pm$ 1449	829 $\pm$ 65
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	76900 $\pm$ 1449	1064 $\pm$ 85
$P_{cat56}$ (2 global optima)	UMDA	209200 $\pm$ 22812	427 $\pm$ 59	197200 $\pm$ 7823	396 $\pm$ 20
	EBNA	277000 $\pm$ 52612	2305 $\pm$ 423	160000 $\pm$ 0	1357 $\pm$ 0
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	96700 $\pm$ 7675	1803 $\pm$ 242
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	94600 $\pm$ 2366	1956 $\pm$ 123
$P_{catring28}$ (4 global optima)	UMDA	— $\pm$ —	— $\pm$ —	127600 $\pm$ 16601	212 $\pm$ 35
	EBNA	203200 $\pm$ 90777	587 $\pm$ 290	192000 $\pm$ 88652	550 $\pm$ 281
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	54700 $\pm$ 949	347 $\pm$ 35
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	58000 $\pm$ 2000	495 $\pm$ 60
	UEBNA $K=6$	— $\pm$ —	— $\pm$ —	59800 $\pm$ 2530	514 $\pm$ 63
$P_{catring56}$ (4 global optima)	UMDA	218200 $\pm$ 38761	423 $\pm$ 79	200714 $\pm$ 28028	389 $\pm$ 63
	EBNA	238000 $\pm$ 50060	1992 $\pm$ 428	240000 $\pm$ 38683	2028 $\pm$ 338
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	97300 $\pm$ 4111	1761 $\pm$ 256
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	94600 $\pm$ 2757	1926 $\pm$ 177
	UEBNA $K=6$	— $\pm$ —	— $\pm$ —	96700 $\pm$ 3592	2349 $\pm$ 158
	UEBNA $K=8$	— $\pm$ —	— $\pm$ —	94600 $\pm$ 1897	2914 $\pm$ 216
$P_{catring42}$ (6 global optima)	UMDA	— $\pm$ —	— $\pm$ —	169000 $\pm$ 20000	313 $\pm$ 46
	EBNA	218800 $\pm$ 63815	1098 $\pm$ 323	221333 $\pm$ 67151	1111 $\pm$ 340
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	73000 $\pm$ 1414	853 $\pm$ 157
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	73000 $\pm$ 0	1038 $\pm$ 85
	UEBNA $K=6$	— $\pm$ —	— $\pm$ —	75700 $\pm$ 3302	1218 $\pm$ 87
	UEBNA $K=8$	— $\pm$ —	— $\pm$ —	75700 $\pm$ 2627	1398 $\pm$ 152
$P_{catring84}$ (6 global optima)	UMDA	253900 $\pm$ 33438	601 $\pm$ 99	248875 $\pm$ 32189	583 $\pm$ 89
	EBNA	277900 $\pm$ 38963	5609 $\pm$ 752	260500 $\pm$ 23335	5486 $\pm$ 481
	UEBNA $K=2$	— $\pm$ —	— $\pm$ —	123400 $\pm$ 4858	3944 $\pm$ 290
	UEBNA $K=4$	— $\pm$ —	— $\pm$ —	120100 $\pm$ 4483	4684 $\pm$ 221
	UEBNA $K=6$	— $\pm$ —	— $\pm$ —	124300 $\pm$ 3860	5567 $\pm$ 612
	UEBNA $K=8$	— $\pm$ —	— $\pm$ —	127900 $\pm$ 2470	6400 $\pm$ 515
	UEBNA $K=10$	— $\pm$ —	— $\pm$ —	121000 $\pm$ 3162	7535 $\pm$ 504

ally multimodal problem optimization, as at best a single global peak is identified per run. On the other hand, the EBNA and the UEBNA are able to discover several global optima per run. However, the results confirm the clear superiority of the UEBNA over the EBNA. On average, more global peaks are discovered per run and they are more equally represented in the population of the last generation. All this indicates that the UEBNA enjoys a robust and reliable behavior against genetic drift. Finally, it should also be observed that increasing the value of  $K$  for the UEBNA has a positive effect on the effectiveness, specially when optimizing the symmetrical problems with four and six global optima.

Table 2 summarizes the efficiency of the EDAs in the evaluation. Regarding the number of evaluations of the objective, i.e., *Eval*, the table shows that the UEBNA significantly speeds up convergence: The saving in number of evaluations that the UEBNA induces over the UMDA and the EBNA for any of the 12 symmetrical optimization problems is considerable. This proves that the UEBNA is able to alleviate genetic drift, accelerating convergence as a result. Furthermore, increasing the value of  $K$  for the UEBNA does not significantly increase *Eval* (it is even reduced in some cases) while, as observed above, effectiveness does improve. Unfortunately, one iteration of the UEBNA is much more time consuming than one iteration of the UMDA or the EBNA, as can be appreciated from the total runtime in *Time*, because it involves running the EM algorithm at least once. In any case, for 10 out of the 12 optimization problems in the evaluation, the UEBNA scores lower runtime than the EBNA for at least one of the values for  $K$  considered. This means that the UEBNA can identify more global optima than the EBNA with less evaluations of the objective function and in a shorter runtime.

It is worth mentioning that our current implementation of the UEBNA, being a proof of concept only, can be considerably improved in terms of runtime by accelerating the EM algorithm, which is the most time consuming part of the code. This implies that the runtime for the UEBNA reported in Table 2 should be read as an upper bound. Several accelerated versions of the EM algorithm have been proposed in the literature (Bauer et al., 1997; Fischer and Kersting, 2003; McLachlan and Krishnan, 1997). The results reported in these papers illustrate that these techniques can substantially reduce the runtime without degrading significantly the quality of the ML parameters. We could further accelerate the EM algorithm by implementing this simple observation: The ML parameters do not usually change substantially between consecutive generations of the UEBNA. Therefore, we could use the ML parameters obtained in one generation in order to initialize the EM algorithm in the next generation. This should reduce the number of iterations of the EM algorithm to converge. By implementing these improvements, the advantages of the UEBNA will be even more apparent. In this paper, we are primarily interested in evaluating the effectiveness of the UEBNA as a proof of concept, while we consider the runtime a secondary performance criterion because it depends very much on the implementation of the UEBNA. It is out of the scope of this paper to compare different implementations.

As summary, it can be said that the UEBNA behaves effectively as well as efficiently for symmetrical problem optimization. Specifically, the results discussed above confirm that the UEBNA is able to alleviate genetic drift with the help of unsupervised learning of TANB models. This means that the UEBNA reduces the likelihood of suboptimal convergence, obtains several global optima per run, and speeds up convergence.

## 5 Conclusions

The main contribution of this paper is the introduction and evaluation of a new estimation of distribution algorithm (EDA), called unsupervised estimation of Bayesian network algorithm (UEBNA), for effective and efficient globally multimodal problem optimization. The main steps of the UEBNA are the same as those of any other EDA: Selection of promising individuals, probabilistic modelling of the selected individuals, and model sampling in order to create the new population. The only peculiarity of the UEBNA with respect to existing EDAs is being based on unsupervised learning of Bayesian networks (BNs) in order to model the selected individuals at each iteration. This makes the UEBNA able to model simultaneously the different basins that may be

represented by the individuals selected at each iteration, whereas preventing genetic drift as much as possible, because this phenomenon is the main one responsible for the poor performance of most evolutionary algorithms (EAs), including EDAs, when optimizing globally multimodal problems.

We have evaluated the UEBNA for symmetrical (globally multimodal) problem optimization, which is known to be challenging for most EAs, including EDAs. We benchmarked the UEBNA against two well established EDAs, namely the univariate marginal distribution algorithm and the estimation of Bayesian network algorithm. The results obtained confirm the ability of the UEBNA to reduce the likelihood of sub-optimal convergence, obtain more global optima per run, and speed up convergence with respect to the two benchmarks. Thus, we can conclude that the UEBNA performs effectively and efficiently for symmetrical (globally multimodal) problem optimization.

In addition to BNs for data clustering, other classes of probabilistic graphical models for data clustering may be considered within the EDA framework, as illustrated in Fig. 5, for globally multimodal problem optimization. These are mixtures of Bayesian networks (Thiesson et al., 1998a; Thiesson et al., 1998b), and Bayesian multinets and recursive Bayesian multinets for data clustering (Peña, 2001; Peña et al., 2002). Despite having received little attention in the literature, these classes of probabilistic graphical models for data clustering offer greater flexibility and expressive power than BNs for data clustering: They can encode context-specific conditional (in)dependencies, whereas BNs for data clustering can encode only context-non-specific conditional (in)dependencies. These models can be particularly useful when the optimization problems at hand are known to be globally multimodal but not necessarily symmetrical. This is a line of research we are currently studying.

In this paper, we focus on discrete domains. The vast majority of the existing EDAs for discrete problem optimization have been already adapted to continuous problem optimization (see (Larrañaga and (eds.), 2001) for a revision). Likewise, we can extend the UEBNA to deal with continuous globally multimodal problem optimization by replacing unsupervised learning of BNs at each iteration by unsupervised learning of conditional Gaussian networks (Peña, 2001; Peña et al., 2001a; Peña et al., 2001b; Peña et al., 2001c). Like a BN for data clustering, a conditional Gaussian network for data clustering consists of a constrained acyclic directed graph and a set of conditional probability distributions. Unlike BNs for data clustering, the conditional probability density functions for the unidimensional predictive random variables are linear regression models conditioned on the value of the cluster random variable. As a result, the generalized joint probability distribution encoded is a conditional Gaussian distribution (Castillo et al., 1997; Cowell et al., 1999; Lauritzen, 1996). Currently, our main line of research concerns empirical evaluation of this extension of the UEBNA for continuous globally multimodal problem optimization. Preliminary experiments look promising.

## Acknowledgments

We thank the three anonymous referees for their useful suggestions. We also thank Marc Schoenauer for handling the review process.

## References

- Anderberg, M. R. (1973). *Cluster Analysis for Applications*. Academic Press.
- Baluja, S. (1994). Population-Based Incremental Learning: A Method for Integrating

Genetic Search Based Function Optimization and Competitive Learning. Technical Report CMU-CS-94-163, Carnegie Mellon University.

- Baluja, S. and Davies, S. (1997). Using Optimal Dependency-Trees for Combinatorial Optimization: Learning the Structure of the Search Space. In *Proceedings of the Fourteenth International Conference on Machine Learning*, pages 30–38. Morgan Kaufmann Publishers.
- Baluja, S. and Davies, S. (1998). Fast Probabilistic Modeling for Combinatorial Optimization. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence*, pages 469–476.
- Bauer, E., Koller, D., and Singer, Y. (1997). Update Rules for Parameter Estimation in Bayesian Networks. In *Proceedings of the Thirteenth Conference on Uncertainty on Artificial Intelligence*, pages 3–13. Morgan Kaufmann Publishers.
- Binder, J., Koller, D., Russell, S., and Kanazawa, K. (1997). Adaptive Probabilistic Networks with Hidden Variables. *Machine Learning*, 29:213–244.
- Castillo, E., Gutiérrez, J. M., and Hadi, A. S. (1997). *Expert Systems and Probabilistic Network Models*. Springer-Verlag.
- Chickering, D. M. (1996). Learning Bayesian Networks is NP-Complete. In *Proceedings of the Fifth International Workshop on Artificial Intelligence and Statistics*, pages 121–130. Springer-Verlag.
- Chickering, D. M. (2002). Learning Equivalence Classes of Bayesian-Network Structures. *Journal of Machine Learning Research*, 2:445–498.
- Cooper, G. F. and Herskovits, E. (1992). A Bayesian Method for the Induction of Probabilistic Networks from Data. *Machine Learning*, 9:309–347.
- Cowell, R. G., Dawid, A. P., Lauritzen, S. L., and Spiegelhalter, D. J. (1999). *Probabilistic Networks and Expert Systems*. Springer-Verlag.
- De Bonet, J. S., Isbell, C. L., and Viola, P. (1997). MIMIC: Finding Optima by Estimating Probability Densities. *Neural Information Processing Systems*, 9.
- De Jong, K. A. (1975). *An Analysis of the Behavior of a Class of Genetic Adaptive Systems*. PhD Thesis, University of Michigan.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum Likelihood from Incomplete Data Via the EM Algorithm. *Journal of the Royal Statistical Society B*, 39(1):1–38.
- Duda, R. O. and Hart, P. E. (1973). *Pattern Classification and Scene Analysis*. John Wiley and Sons.
- Etxeberria, R. and Larrañaga, P. (1999). Global Optimization Using Bayesian Networks. In *Proceedings of the Second Symposium on Artificial Intelligence*, pages 332–339.
- Fischer, J. and Kersting, K. (2003). Scaled CGEM: A Fast Accelerated EM. In *Proceedings of the Fourteenth European Conference on Machine Learning*, pages 133–144. Springer.
- Fogel, L. J. (1962). Autonomous Automata. *Industrial Research*, 4:14–19.

- Fogel, L. J. (1964). *On the Organization of Intellect*. PhD Thesis, University of California.
- Friedman, N. (1998). The Bayesian Structural EM algorithm. In *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence*, pages 129–138. Morgan Kaufmann Publishers.
- Friedman, N., Geiger, D., and Goldszmidt, M. (1997). Bayesian Network Classifiers. *Machine Learning*, 29:131–163.
- Friedman, N. and Koller, D. (2003). Being Bayesian About Network Structure. A Bayesian Approach to Structure Discovery in Bayesian Networks. *Machine Learning*, 50(1):95–125.
- Gallagher, M., Freat, M., and Downs, T. (1999). Real-Valued Evolutionary Optimization Using a Flexible Probability Density Estimator. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 840–846. Morgan Kaufmann Publishers.
- Geman, S. and Geman, D. (1984). Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:721–741.
- Goldberg, D. E. (1989). *Genetic Algorithms in Search, Optimization, and Machine Learning*. Addison-Wesley.
- Goldberg, D. E., Deb, K., Kargupta, H., and Harik, G. R. (1993). Rapid, Accurate Optimization of Difficult Problems Using Fast Messy Genetic Algorithms. In *Proceedings of the Fifth International Conference on Genetic Algorithms*, pages 56–64. Morgan Kaufmann Publishers.
- Goldberg, D. E. and Richardson, J. (1987). Genetic Algorithms with Sharing for Multimodal Function Optimization. In *Proceedings of the Second International Conference on Genetic Algorithms*, pages 41–49. Morgan Kaufmann Publishers.
- Goldberg, D. E. and Segrest, P. (1987). Finite Markov Chain Analysis of Genetic Algorithms. In *Proceedings of the Second International Conference on Genetic Algorithms*, pages 1–8. Morgan Kaufmann Publishers.
- Harik, G. R. (1999). Linkage Learning Via Probabilistic Modeling in the ECGA. Technical Report IlliGAL No. 1999010, University of Illinois at Urbana-Champaign.
- Hartigan, J. A. (1975). *Clustering Algorithms*. John Wiley and Sons.
- Heckerman, D., Geiger, D., and Chickering, D. M. (1995). Learning Bayesian Networks: The Combination of Knowledge and Statistical Data. *Machine Learning*, 20:197–243.
- Henrion, M. (1988). Propagation of Uncertainty by Probabilistic Logic Sampling in Bayes' Networks. In *Uncertainty in Artificial Intelligence 2*, pages 149–164.
- Hocaoğlu, C. and Sanderson, A. C. (1997). Multimodal Function Optimization Using Minimal Representation Size Clustering and Its Applications to Planning Multipaths. *Evolutionary Computation*, 5(1):81–104.
- Holland, J. H. (1975). *Adaptation in Natural and Artificial Systems*. The University of Michigan Press.

- Horn, J. (1993). Finite Markov Chain Analysis of Genetic Algorithms with Niching. In *Proceedings of the Fifth International Conference on Genetic Algorithms*, pages 110–117. Morgan Kaufmann Publishers.
- Jensen, F. (1997). *HUGIN API Reference Manual Version 3.1*.
- Jensen, F. V. (2001). *Bayesian Networks and Decision Graphs*. Springer-Verlag.
- Jensen, F. V., Lauritzen, S. L., and Olesen, K. G. (1990). Bayesian Updating in Causal Probabilistic Networks by Local Computations. *Computational Statistics Quarterly*, 5(4):269–282.
- Keogh, E. J. and Pazzani, M. J. (1999). Learning Augmented Bayesian Classifiers: A Comparison of Distribution-Based and Classification-Based Approaches. In *Proceedings of the Seventh International Workshop on Artificial Intelligence and Statistics*, pages 225–230. Morgan Kaufmann Publishers.
- Kvasnicka, V., Pelikan, M., and Pospichal, J. (1996). Hill Climbing with Learning (An Abstraction of Genetic Algorithms). *Neural Network World*, 6:773–796.
- Larrañaga, P. and (eds.), J. A. L. (2001). *Estimation of Distribution Algorithms. A New Tool for Evolutionary Computation*. Kluwer Academic Publishers.
- Larrañaga, P., Etxeberria, R., Lozano, J. A., and Peña, J. M. (2000). Combinatorial Optimization by Learning and Simulation of Bayesian Networks. In *Proceedings of the Sixteenth Conference on Uncertainty in Artificial Intelligence*, pages 343–352. Morgan Kaufmann Publishers.
- Larrañaga, P., Kuijpers, C. M. H., Murga, R. H., and Yurramendi, Y. (1996). Searching for the Best Ordering in the Structure Learning of Bayesian Networks. *IEEE Transactions on Systems, Man and Cybernetics*, 26(4):487–493.
- Lauritzen, S. L. (1996). *Graphical Models*. Clarendon Press.
- Lauritzen, S. L. and Spiegelhalter, D. J. (1988). Local Computations with Probabilities on Graphical Structures and Their Application to Expert Systems. *Journal of the Royal Statistical Society B*, 50(2):157–224.
- Lobo, F. G., Deb, K., Goldberg, D. E., Harik, G. R., and Wang, L. (1998). Compressed Introns in a Linkage Learning Genetic Algorithm. In *Proceedings of the Third Annual Conference on Genetic Programming*, pages 551–558. Morgan Kaufmann Publishers.
- Mahfoud, S. (1994). Population Sizing for Sharing Methods. In *Foundations of Genetic Algorithms 3*. Morgan Kaufmann Publishers.
- McLachlan, G. J. and Basford, K. E. (1988). *Mixture Models: Inference and Applications to Clustering*. Marcel Dekker.
- McLachlan, G. J. and Krishnan, T. (1997). *The EM Algorithm and Extensions*. John Wiley and Sons.
- Meilă, M. (1999). *Learning with Mixtures of Trees*. PhD Thesis, Massachusetts Institute of Technology.

- Mühlenbein, H. (1997). The Equation for Response to Selection and Its Use for Prediction. *Evolutionary Computation*, 5(3):303–346.
- Mühlenbein, H., Mahnig, T., and Ochoa, A. (1999). Schemata, Distributions and Graphical Models in Evolutionary Optimization. *Journal of Heuristics*, 5:215–247.
- Mühlenbein, H. and Paaß, G. (1996). From Recombination of Genes to the Estimation of Distributions I. Binary Parameters. In *Proceedings of Parallel Problem Solving from Nature IV*, pages 178–187.
- Naudts, B. and Naudts, J. (1998). The Effect of Spin-Flip Symmetry on the Performance of the Simple GA. In *Proceedings of Parallel Problem Solving from Nature V*, pages 67–76. Springer-Verlag.
- Nielsen, J. D., Kočka, T., and Peña, J. M. (2003). On Local Optima in Learning Bayesian Networks. In *Proceedings of the Nineteenth Conference on Uncertainty in Artificial Intelligence*.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann Publishers.
- Pelikan, M. (2002). *Bayesian Optimization Algorithm: From Single Level to Hierarchy*. PhD Thesis, University of Illinois at Urbana-Champaign.
- Pelikan, M. and Goldberg, D. E. (2000). Genetic Algorithms, Clustering, and the Breaking of Symmetry. In *Proceedings of Parallel Problem Solving from Nature VI*, pages 385–394. Springer-Verlag.
- Pelikan, M., Goldberg, D. E., and Cantú-Paz, E. (1999). BOA: The Bayesian Optimization Algorithm. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 525–532. Morgan Kaufmann Publishers.
- Pelikan, M., Goldberg, D. E., and Lobo, F. G. (2000). A Survey of Optimization by Building and Using Probabilistic Models. *Computational Optimization and Applications*, 21(1):5–20.
- Pelikan, M., Goldberg, D. E., and Sastry, K. (2001). Bayesian Optimization Algorithm, Decision Graphs, and Occam’s Razor. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 519–526. Morgan Kaufmann Publishers.
- Pelikan, M. and Mühlenbein, H. (1999). The Bivariate Marginal Distribution Algorithm. *Advances in Soft Computing-Engineering Design and Manufacturing*, pages 521–535.
- Peña, J. M. (2001). *On Unsupervised Learning of Bayesian Networks and Conditional Gaussian Networks*. PhD Thesis, University of the Basque Country.
- Peña, J. M., Izarzugaza, I., Lozano, J. A., Aldasoro, E., and Larrañaga, P. (2001a). Geographical Clustering of Cancer Incidence by Means of Bayesian Networks and Conditional Gaussian Networks. In *Proceedings of the Eighth International Workshop on Artificial Intelligence and Statistics*, pages 266–271. Morgan Kaufmann Publishers.
- Peña, J. M., Lozano, J. A., and Larrañaga, P. (1999). Learning Bayesian Networks for Clustering by Means of Constructive Induction. *Pattern Recognition Letters*, 20(11-13):1219–1230.

- Peña, J. M., Lozano, J. A., and Larrañaga, P. (2000). An Improved Bayesian Structural EM Algorithm for Learning Bayesian Networks for Clustering. *Pattern Recognition Letters*, 21(8):779–786.
- Peña, J. M., Lozano, J. A., and Larrañaga, P. (2001b). Performance Evaluation of Compromise Conditional Gaussian Networks for Data Clustering. *International Journal of Approximate Reasoning*, 28(1):23–50.
- Peña, J. M., Lozano, J. A., and Larrañaga, P. (2002). Learning Recursive Bayesian Multinets for Data Clustering by Means of Constructive Induction. *Machine Learning*, 47(1):63–89.
- Peña, J. M., Lozano, J. A., Larrañaga, P., and Inza, I. (2001c). Dimensionality Reduction in Unsupervised Learning of Conditional Gaussian Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(6):590–603.
- Rechenberg, I. (1973). *Evolutionstrategie: Optimierung Technischer Systeme Nach Prinzipien der Biologischen Evolution*. Fromman-Holzboog Verlag.
- Santana, R. and Ochoa, A. (1999). Dealing with Constraints with Estimation of Distribution Algorithms: The Univariate Case. In *Proceedings of the Second Symposium on Artificial Intelligence*, pages 378–384.
- Schwarz, G. (1978). Estimating the Dimension of a Model. *Annals of Statistics*, 6:461–464.
- Schwarz, J. and Ocenasek, J. (1999). Experimental Study: Hypergraph Partitioning Based on the Simple and Advanced Algorithms BMDA and BOA. In *Proceedings of the Fifth International Conference on Soft Computing*, pages 124–130.
- Schwefel, H. P. (1981). *Numerical Optimization of Computer Models*. John Wiley and Sons.
- Soto, M., Ochoa, A., Acid, S., and de Campos, L. M. (1999). Introducing the Polytree Approximation of Distribution Algorithm. In *Proceedings of the Second Symposium on Artificial Intelligence*, pages 360–367.
- Thiesson, B., Meek, C., Chickering, D. M., and Heckerman, D. (1998a). Learning Mixtures of Bayesian Networks. Technical Report MSR-TR-97-30, Microsoft Research.
- Thiesson, B., Meek, C., Chickering, D. M., and Heckerman, D. (1998b). Learning Mixtures of DAG Models. In *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence*, pages 504–513. Morgan Kaufmann Publishers.
- Van Hoyweghen, C. (2001). Detecting Spin-Flip Symmetry in Optimization Problems. In *Theoretical Aspects of Evolutionary Computing*, pages 175–206. Springer-Verlag.
- Van Hoyweghen, C. and Naudts, B. (2000). Symmetry in the Search Space. In *Proceedings of the Seventh IEEE International Conference on Evolutionary Computation*, pages 1072–1079. IEEE Press.