Learning Bayesian Networks by Genetic Algorithms. A case study in the prediction of survival in malignant skin melanoma

P. Larrañaga †, B. Sierra †, M. J. Gallego †, M. J. Michelena ‡, J. M. Picaza§

†Department of Computer Science and Artificial Intelligence, University of the Basque Country, Spain.
‡Oncological Institute of Gipuzkoa, Spain.

§Department of Computer Languages and Systems, University of the Basque Country, Spain.

e-mail: ccplamup@si.ehu.es

ABSTRACT In this work we introduce a methodology based on Genetic Algorithms for the automatic induction of Bayesian Networks from a file containing cases and variables related to the problem. The methodology is applied to the problem of predicting survival of people after one, three and five years of being diagnosed as having malignant skin melanoma. The accuracy of the obtained model, measured in terms of the percentage of well-classified subjects, is compared to that obtained by the called Naive-Bayes. In both cases, the estimation of the model accuracy is obtained from the 10-fold cross-validation method.

1. Introduction

Expert systems, one of the most developed areas in the field of Artificial Intelligence, are computer programs designed to help or replace humans beings in tasks in which the human experience and human knowledge are scarce and unreliable. Although, there are domains in which the tasks can be specifed by logic rules, other domains are characterized by an uncertainty inherent in them. Probability was not taken into account, for some time, as a reasoning method for expert systems trying to modelize uncertain domains, because the specifications and computer cost it requires are too expensive. At the end of the 80s, Lauritzen and Spiegelhalter [21] showed that these difficulties can be overcome by exploiting the modular character of the graphical models associated with the called probabilistic expert systems, that we call in this work Bayesian Networks.

Bayesian Networks (BNs) [13], [20], [22] constitute a probabilistic framework for reasoning under uncertainty. From an informal perspective, BNs are directed acyclic graphs (DAGs), where the nodes are random variables and the arcs specify the independence assumptions that must be held between the random variables. BNs are based upon the concept of conditional independence among variables. This concept makes possible a factorization of the probability distribution of the *n*-dimensional random variable $(X_1, ..., X_n)$ in the following way:

$$P(x_1, ..., x_n) = \prod_{i=1}^n P(x_i | pa(x_i))$$

where x_i represents the value of the random variable X_i , and $pa(x_i)$ represents the value of the random variables parents of X_i .

Thus, in order to specify the probability distribution of a BN, one must give prior probabilities for all root nodes (nodes with no predecessors) and conditional probabilities for all other nodes, given all possible combinations of their direct predecessors. These numbers in conjunction with the DAG, specify the BN completely. Once the network is constructed

⁰We thank Gregory F. Cooper for providing his simulation of the ALARM Network.

This work was supported by the Diputación Foral de Gipuzkoa, under grant OF 92/1996, and by the grant PI 95/52 from the Gobierno Vasco - Departamento de Educación, Universidades e Investigación.

it constitutes an efficient device to perform probabilistic inference. This probabilistic reasoning inside the net can be carried out by exact methods, as well as by approximated methods. Nevertheless, the problem of building such a network remains. The structure and conditional probabilities necessary for characterising the network can be either provided externally by experts or obtained from an algorithm which automatically induces them.

In this paper, a methodology for inducing automatically Bayesian Networks is introduced. This methodology is based on Genetic Algorithms and tries to obtain from the file of cases the most probable structure of the Bayesian Network. The work is organized as follows, in Section II some structure learning methods are reviewed, taking an special interest in the method proposed by Cooper and Herskovits [5]. Section III introduces Genetic Algorithms, while Section IV presents the structure learning methodology integrating both, the metric proposed by Cooper and Herskovits and the adaptative searching process characteristic of the Genetic Algorithms. In Section V we present the results obtained from applying the previous methodology to a file of cases, which contains information about 311 patients diagnosed as having malignant skin melanoma. The induced Bayesian network is used for classifying patients according to their prognosis of survival after one, three and five years of being diagnosed. These results are compared to those obtained by the called Naive-Bayes paradigm. Section VI gathers the conclusions.

2. Structure Learning in Bayesian Networks

2.1 Introduction

During the last five years a good number of algorithms whose aim is to induce the structure of the Bayesian Network that better represents the conditional independence relationships underlying in the file of cases have been developed. In our opinion, the main reason for continuing the research in the structure learning problem is that modelizing the expert knowledge has become an expensive, unreliable and time-consuming job.

The different approaches to the structure learning mentioned here are related with multiple connected networks, and have been grouped according to the necessity or not of imposing order on the variables. See Heckerman et al. [10] for a good review.

Assuming order among variables means that a variable X_i can have the variable X_j as parent only if, in the established order among the variables, X_j precedes X_i . With this restriction, the cardinality of the space that contains all the structures is given by $2^{\binom{n}{2}}$, where *n* is the number of variables in the system. Some methods under this restriction are those developed by Herskovits and Cooper [11], Cooper and Herskovits [5], and Bouckaert [4].

If we do not assume ordering between the nodes the cardinality of the search space is bigger, and it is given by the Robinson's formula [24]:

$$f(n) = \sum_{i=1}^{n} (-1)^{i+1} {n \choose i} 2^{i(n-i)} f(n-i); f(0) = 1; f(1) = 1.$$

Several authors have been working under these general assumptions. Among them, Bouckaert [3], Lam and Bacchus [14], and Provan and Singh [23].

2.2 The K2 algorithm

As it will be seen in Section IV, the three proposed approaches - based on Genetic Algorithms - use the CH metric proposed by Cooper and Herskovits [5] for evaluating the goodnes of a Bayesian Network structure, as well as the K2 algorithm developed by the previously mentioned authors. K2 is an algorithm that creates and evaluates a BN from a database of cases once an ordering between the system variables is given. The CH metric is used for the evaluation of the network that it constructs. K2 searches, given a database D for the BN structure B_{S^*} with maximal $P(B_S, D)$, where $P(B_S, D)$ is as described in the following theorem proved in [5].

Theorem Let Z be a set of n discrete variables, where a variable x_i in Z has r_i possible value assignments: $(v_{i1}, \ldots, v_{ir_i})$. Let D be a database of cases of m cases, where each case contains a value assignment for each variable in Z. Let B_S denote a BN structure containing just the variables in Z. Each variable x_i in B_S has a set of parents, which are represented with a list of variables π_i . Let w_{ij} denote the *j*th unique instantiation of π_i relative to D. Suppose there are q_i such unique instantiations of π_i . Define N_{ijk} to be the number of cases in D in which variable x_i has the value v_{ik} and π_i is instantiated as w_{ij} . Let $N_{ij} = \sum_{k=1}^{r_i} N_{ijk}$. If given a BN model, the cases occur independently, there are not cases that have variables with missing values and the density function $f(B_P|B_S)$ is uniform, then it follows that

$$P(B_S|D) = P(B_S) \prod_{i=1}^n g(i, \pi_i), \text{ where } g(i, \pi_i) = \prod_{j=1}^{q_i} \frac{(r_i - 1)!}{(N_{ij} + r_i - 1)!} \prod_{k=1}^{r_i} N_{ijk}!$$

The K2 algorithm assumes that an ordering on the variables is available and that, a priori, all structures are equally likely. It searches, for every node, the set of parent nodes that maximizes $g(i, \pi_i)$ - CH metric-. K2 is a greedy heuristic. It starts by assuming that a node does not have parents, after which in every step it adds incrementally that parent whose addition most increases the probability of the resulting structure. K2 stops adding parents to the nodes when the addition of a single parent cannot increase the probability. Obviously, this approach does not guarantee the selection of a structure with the highest probability.

3. Genetic Algorithms

The computing complexity inherent in a great number of real problems of combinatorial optimization has carried, as a consequence, the development of heuristic methods that try to tackle these problems successfully. An heuristic is a procedure which will give a good solution - not necessarily the optimal - to problems which can be catalogued as difficult, if you try to solve them obtaining the exact solution. Although there are heuristics developed for especific problems, in the past years there have been an explosion in the applications of what we could call metaheuristics, because its formulation is independent of the problem to solve. Among the most studied metaheuristics we quote Simulated Annealing, Tabu Search and Genetic Algorithms.

Genetic Algorithms [9] are adaptive methods that can be used for solving problems of search and optimization. They are based on the genetic process of living organisms. Through generations the populations evolve in nature according to the principles of natural selection and survival of the fittest postulated by Darwin [6]. Imitating this process, the Genetic Algorithms are capable of creating solutions for real world problems.

Genetic Algorithms use a direct analogy with the natural behaviour. They work with a population of individuals, each individual representing a feasible solution to a given

 begin AGA

 Make initial population at random

 WHILE NOT stop DO

 BEGIN

 Select parents from the population.

 Produce children from the selected parents.

 Mutate the individuals.

 Extend the population by adding the children to it.

 Reduce the extended population.

 END

 Output the best individual found.

 end AGA

FIGURE 1. The pseudo-code of the Abstract Genetic Algorithm.

problem. To each individual we assign a value or score according to the goodness of that solution. The better the adaptation of the individual to the problem, the more probable is that the individual will be selected for reproduction, crossing its genetic material with another individual selected in the same way. This cross will produce new individual - offspring of the previous - which share some of the features of their parents. In this way a new population of feasible solutions is produced, replacing the previous one and verifing the interesting property of having greater proportion of good features than the previous population. Thus, through generations good features are propagated through the population. Favouring the cross of the fittest individuals, the most promising areas of the search space are being explored. If the Genetic Algorithms have been well designed, the population will converge [7] to an optimal solution of the problem.

Figure 1 summarizes the pseudocode for the so-called Abstract Genetic Algorithm. In it the parent selection doesn't need to be made by asigning to each individual a value proportional to its objetive function, as is usual in the so-called Simple Genetic Algorithm. This selection can be carried out by any function that selects parents in a natural way. It is worth notice that descendants are not necessarily the next generation of individuals, but that this generation is made by the union of parents and descendents. That is why we need the operations of extension and reduction in the cycle.

4. Genetic Algorithms in the Induction of Bayesian Networks

4.1 Searching in the space of networks structures

In this approach, each individual in the Genetic Algorithm will be a Bayesian Network structure.

4.1.1 Notation and Representation

Denoting with D the set of BN structures for a fixed domain with n variables, and the alphabet S being $\{0,1\}$, a Bayesian Network structure can be represented by an $n \times n$ connectivity matrix C, where its elements, c_{ij} , verify:

$$c_{ij} = \begin{cases} 1 & \text{if } j \text{ is a parent of } i, \\ 0 & \text{otherwise.} \end{cases}$$

4.1.2 Assuming an ordering between the nodes

In this case, the connectivity matrices of the network structures are triangulated and

therefore the genetic operators are closed operators with respect to the DAG conditions. We represent an individual of the population by the string:

$$c_{21}c_{31}c_{41}\ldots c_{n1},\ldots c_{32}c_{42}\ldots c_{n2},\ldots c_{n-2n-1},c_{n-2n},c_{n-1n}$$

With this representation in mind, we show how the crossover and mutation operators work by using simple examples.



FIGURE 2. With order assumption: Crossing over two BN structures.

Example 1. Consider a domain of 3 variables on which the two BN structures of Figure 2(a) are defined. Using the above described representation, the networks are represented by the strings : 110 and 101. Suppose now that the two network structures are crossed over and that the crossover point is chosen between the second and the third bit. This gives the offspring strings 111 and 100. Hence, the created offspring structures are the ones presented in Figure 2(b).

Example 2. Consider the DAG of Figure 3(a). It is represented by the string 100. Suppose



FIGURE 3. With order assumption: Mutating a BN structure.

that the third bit is alterated by mutation. This gives the string 101, which corresponds with the graph of Figure 3(b).

4.1.3 Without asuming an ordering between the nodes

If no ordering assumption on the variables is made, we represent an individual of the population by the string:

$$c_{11}c_{21}\ldots c_{n1}c_{12}c_{22}\ldots c_{n2}\ldots c_{1n}c_{2n}\ldots c_{nn}.$$

As can be seen in the following examples, the genetic operators are not closed operators with respect to the DAG conditions.

Example 3. Consider a domain of 3 variables on which the two BN structures of Figure

4(a) are defined. Using the above described representation, the networks are represented by the strings : 001001000 and 000000110. Suppose now that the two network structures are crossed over and that the crossover point is chosen between the sixth and the seventh bit. This gives the offspring strings 001001110 and 000000000. Hence, the created offspring structures are the ones presented in Figure 4(b). We see that the first offspring structure is not a DAG.



FIGURE 4. Without order assumption: The crossover operator is not a closed operator.

Example 4. Consider the DAG of Figure 5(a). It is represented by the string 010001000. Suppose that the seventh bit is alterated by mutation. This gives the string 010001100, which corresponds with the cyclic graph of Figure 5(b).

To assure the closeness of the genetic operators we introduce a *repair operator*, which transforms the child structures that do not verify the DAG conditions into DAGs, by randomly eliminating the edges that invalidate the DAG conditions.

4.2 Searching for the best ordering



FIGURE 5. Without order assumption: The mutation operator is not a closed operator.

The individuals of the population are orderings whose fitness is computed by applying the formula of [5] to the structure that is induced by applying the K2 algorithm to it. Now, the cardinality of the search space is n!.

In this case, the problem of the structure learning can be modelled as a problem that resembles the intensively studied Traveling Salesman Problem (TSP). While the TSP problem is assumed to be symmetrical, in this approach to the structure learning of Bayesian Networks the problem is not a symmetrical one, and we search for acyclic orderings. See Larrañaga et al. [17] for an empirical evaluation of this approach.

4.3 Experiments

The two approaches in which the search has been done in the space of networks structures

have been evaluated empirically with a simulation of the ALARM network. For details see Larrañaga et al. [15], [16]. For these experiments we use a database containing 3000 cases which was the result of a simulation of the ALARM network [2]. This database has become a benchmark for evaluating the performance of newly proposed algorithms. The cardinality of the search spaces assuming an ordering between the variables or wihout assuming it, are respectively 3.061e200 and 3.008e237.

In figure 6 we can see that the Hamming distance between the ALARM network structure and the induced structure is one - the arc from node 12 to node 32 has been deleted -.

4.4 Genetic Algorithms in others combinatorial problems related with the Bayesian Network paradigm

Genetic Algorithms have been used as optimizers in several combinatorial problems that arise from the Bayesian Networks context. Thus, for example, Larrañaga et al. [18], obtain good decompositions of the moral graph associated with the propagation algorithm proposed by Lauritzen and Spiegelhalter [21]. Larrañaga et al. [19] also treat the problem of the fusion of Bayesian Networks coming from different authors, seeking for the consensual BN. Finally, Rojas-Guzmán and Kramer [25], and Gelsema [8] also treat the problem of determining the most probable global state of the system using GAs.

5 Predicting survival in malignant skin melanoma

5.1 The malignant skin melanoma

In spite of the advances achieved in last years in the treatment of cancer, the prognosis of patients having developed skin melanoma has changed very little. The incidence of the disease has grown without stopping in the last decade. Annual incidence has increased from 4% to 8%, and the progressive reduction of the ozone layer, if not stopped, will expand it even more.

Experimental data and the results of epidemiological studies suggest two main risk factors: sun exposure along with phenotype characteristics of the individual. Thus, for example, the continuous sun exposure represents an odds ratio of 9, while the acute intermitent exposition has got associated an odds ratio of 5.7.

Malignant skin melanoma is a rather uncommon tumour in our environment. It entails between the 8% and the 10% of the total malignant tumours that affect the skin. According to the Cancer Register of the Basque Country [12], in 1990 the rate of incidence was 2.2 for every 100000 people for males and 3 for every 100.000 for females.

The database contains 311 cases - diagnosed at the Oncological Institute of Gipuzkoa in the period between the first of January, 1988, and the 31 of December, 1995 - and for each case we have information about eight variables. The five predictor variables are: sex (2 categories), age (5 categories), stage (4 categories), thickness (4 categories) and number of positive nodes (2 categories). The variable to predict has two categories taking into account if the person survives or not one, three or five years after being diagnosed as malign skin melanoma.

5.2 The Models

Two models have been taken into account. First, we have induced a BN structure using GAs, as explained in Section IV. In order to get it, we have searched in the space of all structures without imposing any order restriction among the variables. Therefore we



FIGURE 6. (a) The ALARM network structure. (b) The Bayesian Network structure learnt by the Genetic Algorithm when assuming an ordering between the variables, from a database containing a simulation of the ALARM network with 3000 cases.

have tried to find, given a file with cases, the a posteriori most probable structure. The second model used is the called Naive-Bayes. This model assumes independence among predictor variables. In both models the estimations of the rate of well-classified individuals have been obtained using 10-fold cross-validation[26]. The propagation of the evidence has been done using the software HUGIN [1]



Model I. The a posteriori most probable structure. CH-GA. Figures 7, 8 and 9 show the

FIGURE 7. The a posteriori most probably structure for the one year case.



FIGURE 8. The a posteriori most probably structure for the thee year case.

structures of the Bayesian Networks induced by the Genetic Algorithm. They correspond to the predictions of survival after one, three and five years of being diagnosed. In table 1, estimations of the probability of succes in classification obtained by each of the previous models can be seen.

Model II. Naive - Bayes classifier. N-B. In spite of the strong assumptions of independence upon which the model is built, Naive - Bayes classifier has proved itself competitive against other more refined classifiers. It is assumed that all variables are conditionally independent given the value of the variable to predict. Therefore, the model ignores the



 $\ensuremath{\mathrm{FIGURE}}$ 9. The a posteriori most probably structure for the five year case.



 FIGURE 10. The Naive-Bayes classifier.

correlations among variables which can prejudice its predictive capacity. In Figure 10 it can be seen the structure of the Bayesian Network corresponding to the Naive-Bayes. This structure is common to the three classification problems. Table 1 shows that the estimations obtained by two of the Naive-Bayes models are inferior to those obtained by the previous approach.

Survival of Malignant Skin Melanoma			
	1 year	3 years	5 years
CH-GA	93.06	81.95	69.57
N-B	91.43	79.02	71.43

TABLE 1. Accuracy of the differents approaches for the prediction of survival one-year, thee-years and five-years after be diagnosed.

6. Conclusions and futher research

A method of induction of Bayesian Networks has been introduced. This method is based on intelligent search made by Genetic Algorithms. The method uses the CH metric and tries to find the a posteriori most probable Bayesian Network structure given the file of cases.

The Bayesian Network structures induced by this method have been empirically compared to the Naive-Bayes structures in one classification problem consisting on the prediction of survival of individuals after one, three or five years of being diagnosed as having malignant skin melanoma. Although in the inductive method does not exist an especial treatment for the variable to classify, the estimations of the 10-fold croosvalidation for the probability of survival are better in two of the three examples than those obtained by the Naive-Bayes paradigm.

In the future, we plan to combine the previous two techniques for finding the a posteriori most probable structure that contains, at least, the Naive-Bayes. Besides, it would be interesting to develop a method of intelligent search based on Genetic Algorithms, Simulated Annealing or Tabu Search that take into account the purpose the induced Bayesian Network will be use for - in our case supervised classification -.

0.1 References

- Andersen, S.K., Olesen, K.G., Jensen, F.V. and Jensen, F. (1989): "HUGIN a shell for building Bayesian belief universes for expert systems". In *Eleventh International Joint Conference on Artificial Intelligence*, vol. I, pp. 1128-1133.
- [2] Beinlinch, I.A., Suermondt, H.J., R.M. Chavez R. M. and Cooper G.F. (1989): "The ALARM monitoring system: A case study with two probabilistic inference techniques for belief networks". In Proceedings of the Second European Conference on Artificial Intelligence in Medicine, pp. 247-256.
- [3] Bouckaert, R.R. (1992): "Optimizing causal orderings for generating DAGs from data". In Uncertainty in Artificial Intelligence. Proceedings of the Eighth Conference, pp. 9-16.
- [4] Bouckaert, R.R. (1994): "Properties of Bayesian belief networks learning algorithms". In Uncertainty in Artificial Intelligence. Tenth Annual Conference, pp. 102-109.
- [5] Cooper, G.F., and Herskovits, E.A. (1992): "A Bayesian method for the induction of probabilistic networks from data". Machine Learning, vol. 9, no. 4, pp. 309-347.

- [6] Darwin, C. (1859): "On the Origin of the Species by Means of Natural Selection". Murray, London.
- [7] Eiben, A.E., Aarts, E.H.L., and van Hee, K.M. (1990): "Global convergence of genetic algorithms: An infinite Markov chain analysis". Computing Science Notes, Eindhoven University of Technology.
- [8] Gelsema, E.S. (1995): "Abductive reasoning in Bayesian belief networks using a genetic algorithm". In Preliminary Papers of the Fifth International Workshop on Artificial Intelligence and Statistics, pp. 245-251.
- [9] Goldberg, D.E. (1989): "Genetic Algorithms in Search, Optimization and Machine Learning". Addison-Wesley, Reading, MA.
- [10] Heckerman, D., Geiger, D. and Chickering, D.M. (1994): "Learning Bayesian networks: The combination of knowledge and statistical data". Technical Report MSR-TR-94-09, Microsoft.
- [11] Herskovits, E. and Cooper, G.F. (1990): "Kutató: An entropy-driven system for construction of probabilistic expert systems from databases". *Report KSL-90-22*, Knowledge Systems Laboratory, Medical Computer Science, Stanford University.
- [12] Izarzugaza, M.I. (1994): "Informe del registro de Cáncer de Euskadi 1990". Osasunkaria, pp. 8-11.
- [13] Jensen, F. V. (1996): "Introduction to Bayesian networks". University College of London.
- [14] Lam, W. and Bacchus, F. (1994): "Learning Bayesian belief networks. An approach based on the MDL principle". Computational Intelligence, vol. 10, no. 4.
- [15] Larrañaga, P., Poza, M., Yurramendi, Y., Murga, R., and Kuijpers, C. (1996): "Structure Learning of Bayesian Networks by Genetic Algorithms: A Performance Analysis of Control Parameters". *IEEE Trans*actions on Pattern Analysis and Machine Intelligence. In press.
- [16] Larrañaga, P., Murga, R., Poza, M., and Kuijpers, C. (1996): "Structure Learning of Bayesian Networks by Hybrid Genetic Algorithms". In Learning from Data: AI and Statistics V, Lecture Notes in Statistics 112. D. Fisher, H.-J. Lenz (eds.), New York, NY: Spriger-Verlag, pp. 165-174.
- [17] Larrañaga, P., Kuijpers, C., Murga, R., and Yurramendi, Y. (1996): "Learning Bayesian Network Structures by searching for the best ordering with genetic algorithms". *IEEE Transactions on System, Man and Cybernetics*. Vol. 26, 4, pp. 487-493.
- [18] Larrañaga, P., Kuijpers, C., Poza, M., and Murga, R. (1996) "Decomposing Bayesian Networks by Genetic Algorithms". Statistics and Computing. In press.
- [19] Larrañaga, P., Kuijpers, C., Murga, R., Yurramendi, Y., Graña, M., Lozano, J.A., Albizuri, X., D'Anjou, A., Torrealdea, F.J. (1996): "Genetic Algorithms applied to Bayesian Networks". In A. Gammerman (ed.) Computational Learning and Probabilistic Reasoning. John Wiley, pp. 211-234.
- [20] Lauritzen, S.L. (1996): "Graphical Models". Oxford University Press.
- [21] Lauritzen, S.L., and Spiegelhalter, D.J. (1988): "Local computations with probabilities on graphical structures and their application on expert systems". Journal Royal of Statistical Society B, vol. 50, no. 2, pp. 157-224.
- [22] Pearl, J. (1988): "Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference". Morgan Kaufmann, San Mateo.
- [23] Provan, G.M., and Singh, M. (1995): "Learning Bayesian Networks Using Feature Selection". In Learning from Data: AI and Statistics V, Lecture Notes in Statistics 112. D. Fisher, H.-J. Lenz (eds.), New York, NY: Spriger-Verlag, pp. 291-300.
- [24] Robinson, R. W. (1977): "Counting unlabeled acyclic digraphs". In C. H. C. Little (ed.) Lectures Notes in Mathematics 622: Combinatorial Mathematics V, Springer-Verlag, New York, pp. 28-43.
- [25] Rojas-Guzmán, C., and Kramer, M.A. (1993): "GALGO: A Genetic ALGOrithm decision support tool for complex uncertain systems modeled with Bayesian belief networks". In Uncertainty in Artificial Intelligence. Proceedings of the Ninth Conference, pp. 368-375.
- [26] Stone, M. (1974): "Cross-validation choice and assessment of statistical procedures". Journal Royal of Statistical Society, vol 36, pp. 111-147.