



NLR-TP-98484

On the suitability of genetic-based algorithms for data mining

S. Choenni



NLR-TP-98484

On the suitability of genetic-based algorithms for data mining

S. Choenni

This investigation has been carried out under a contract awarded by the Ministry of Defence, contract number 16-98-3831-01.

The Ministry of Defence has granted NLR permission to publish this report.

This report is based on an article to be published in the Proc. Intern. Workshop on Advances in Database Technologies, by Springer Verlag (serie LNCS) 1998.

The contents of this report may be cited on condition that full credit is given to NLR and the author.

Division:	Informatics
Issued:	November 1998
Classification of title:	unclassified

Summary

Data mining has as goal to extract knowledge from large databases. To extract this knowledge, a database may be considered as a large search space, and a mining algorithm as a search strategy. In general, a search space consists of an enormous number of elements, making an exhaustive search infeasible. Therefore, efficient search strategies are of vital importance. Search strategies based on genetic-based algorithms have been applied successfully in a wide range of applications. In this paper, we discuss the suitability of genetic-based algorithms for data mining. We discuss the design and implementation of a genetic-based algorithm for data mining and illustrate its potentials.



Contents

1	Introduction	5
2	Preliminaries & problem limitations	8
3	Data Mining with Genetic algorithms	9
3.1	Representation	9
3.2	Fitness function	10
3.3	Manipulation operators	12
3.3.1	Mutation	12
3.3.2	Cross-over	13
4	Optimization rules	15
5	Algorithm	17
6	Implementation and preliminary results	20
7	Conclusions & further research	24

1 Table

4 Figures

(26 pages in total)

1 Introduction

The amount and kind of data that are stored in databases are still growing. These data may contain implicit knowledge that can improve the quality of decisions taken in the present or future. Due to the huge amount and various kind of data that are stored in databases, conventional data analysis tools are inadequate to extract knowledge from these databases.

A field that deals with extracting knowledge from databases, without putting restrictions on the amount or types of data in a database, is data mining. Nowadays, data mining is achieving more and more interest from academy as well as from industry [Ref 11]. The interest from academy stems from the many challenging research topics entailed by data mining. Interest from industry is due to the fact that data mining may result into knowledge that could be of vital importance for a company.

Research and development in data mining evolves in several directions, such as association rules, time series, and classification. The direction of association rules is focussing on the development of algorithms to find frequently occurring patterns in a database, see among others [Ref 2, 3, 16, 18, 21]. In time series databases, one tries to find all common patterns embedded in a database of sequences of events [Ref 4]. The classification of tuples in a number of groups on the basis of common characteristics and the derivation of rules from a group is another direction in data mining [Ref 1, 13, 15, 5].

Especially the latter field has our attention. We have developed an algorithm, in cooperation with the Centre for Mathematics and Computer Science in the Netherlands, to classify tuples in groups and to derive rules from these groups. In our view, a user formulates a mining question and the algorithm selects the group(s) that satisfy this question. For example, in an insurance environment, a question may be to identify persons with (more than average) chances of causing an accident. Then, the algorithm searches for the profiles of these persons.

In general, the search spaces that should be inspected in order to find answers on mining questions are very large, making exhaustive search infeasible. So, heuristic search strategies are of vital importance to data mining. Genetic algorithms, which are heuristic search strategies, have been successfully used in a wide range of applications, such as air traffic management [Ref 17], computational fluid dynamics, query optimization, etc. A genetic algorithm is capable of exploring different parts of a search space [Ref 19].

In this paper, we discuss the applicability of a genetic-based algorithm to the search process in



data mining. We show how a genetic algorithm can be suited for data mining problems. In our approach, a search space consists of expressions. An expression is a conjunction of predicates and each predicate is defined on a database attribute. Initially, a random number of expressions, called initial population, is selected. Then, the initial population is manipulated by applying a number of operations. The best individuals are selected to form the next generation and the manipulation process is repeated until no significant improvement of the population can be observed.

In general, data mining algorithms require a technique that partitions the domain values of an attribute in a limited set of ranges, simply because considering all possible ranges of domain values is infeasible. In our approach, this boils down on the selection of proper ranges in expressions. Suppose that we have an attribute *age* which has a domain between 18 to 65, and an expression of the form *age* **in** $[v_i, v_k]$, in which v_i and v_k are values from the domain of *age*, defining a range of values. The problem is how to choose the values for v_i and v_k . As illustrated in [Ref 21], this is in general a tough problem. Our solution to this problem is based on a suitable choice of the mutation operator (see Section 3.3). Furthermore, we have chosen a representation for individuals that fits in the field of databases. The same holds for the manipulation operators and the function to rank individuals (fitness function). The fitness function discussed in this paper is close to our intuition and gives rise to a speed up of the optimization process. Based on our approach, we have implemented a (prototype) tool for data mining, and have performed a preliminary evaluation. The results will be presented in this paper.

A genetic approach has been proposed in [Ref 5] to learn first order logic rules and in [Ref 12] a framework is proposed for data mining based on genetic programming. However, the authors neither come up with a implementation nor with experiments. The effort in [Ref 5] is focussed towards machine learning, and the important data mining issue of integration with databases is superficially discussed. The effort in [Ref 12] describes a framework for data mining based on genetic programming, and stresses on the integration of genetic programming and databases. However, an elaborated approach to implement and evaluate the framework is not presented. Other related research has been reported in [Ref 13, 15]. In these efforts, the authors use variants of a hill climber to identify the group(s) of tuples satisfying a mining question. However, the problem of partitioning attribute values has not been discussed in these efforts. We note that a genetic-based algorithm has, by nature, a better chance to escape from a local optimum than a hill climber.

The remainder of this paper is organized as follows. In Chapter 2, we outline some preliminaries and problem limitations. In Chapter 3, we identify the issues that play a role in genetic-based algorithms and how to solve them in a data mining context. In Chapter 4, we point out a number



of rules that may speed up the search process of a genetic-based algorithm. Chapter 5 is devoted to an overall algorithm for data mining. In Chapter 6, we discuss the implementation of the algorithm and some preliminary results. Finally, Chapter 7 contains conclusions and further work.

2 Preliminaries & problem limitations

In the following, a database consists of a universal relation. The relation is defined over some independent single valued attributes, such as $att_1, att_2, \dots, att_n$, and is a subset of the Cartesian product $\text{dom}(att_1) \times \text{dom}(att_2) \times \dots \times \text{dom}(att_n)$, in which $\text{dom}(att_j)$ is the set of values that can be assumed by attribute att_j . A tuple is an ordered list of attribute values to which a unique identifier (tid) is associated. So, we do not allow missing attribute values. Furthermore, we assume that the content of the database remains the same during the mining process.

An expression is used to derive a relation, and is defined as a conjunction of predicates over some attributes. The length of an expression is the number of attributes involved in the expression. An example of an expression of length 2 is ($age \text{ in } [19, 24] \wedge gender \text{ is 'male'}$), representing the males who are older than 18 and younger than 25. An expression with length 1 is called an *elementary* expression. In this paper, we deal with search spaces that contain expressions.

3 Data Mining with Genetic algorithms

Genetic algorithms are heuristic search strategies inspired by natural genetics and evolutionary principles [Ref 14]. They have been proven to be promising in a wide variety of applications. A genetic algorithm [Ref 19] randomly generates an initial population. Traditionally, individuals in the population are represented as bit strings. The quality of each individual, i.e., its fitness, is computed. On the basis of these qualities, a selection of individuals is made (an individual may be chosen more than once). Some of the selected individuals undergo a minor modification, called mutation. For some pairs of selected individuals a random point is selected, and the substrings behind this random point are exchanged; this process is called cross-over. The selected individuals, modified or not, form a new population and the same procedure is applied to this generation until some predefined criteria are met.

In this chapter, we discuss the issues that play a role in tailoring a genetic algorithm for data mining. Section 3.1 is devoted to the representation of individuals in a population. We propose a representation that seamlessly fits to relational databases. Then, in Section 3.2, we discuss a fitness function, which computes the quality of an individual. Finally, in Section 3.3, we discuss the two operators that are used to manipulate an individual.

3.1 Representation

As stated in Chapter 2, an expression is a conjunction of predicates and a predicate is defined on a database attribute. An individual can be regarded as an expression to which some restrictions are imposed with regard to the notation of elementary expressions and the number of times that an attribute may be involved in the expression.

The notation of an elementary expression depends on the domain type of the involved attribute. If there exists no ordering relationship between the attribute values of an attribute att , we represent an elementary expression as follows: $expression := att \text{ is } (v_1, v_2, \dots, v_n)$, in which $v_i \in \text{dom}(att)$, $1 \leq i \leq n$. In this way, we express that an attribute att assumes one of the values in the set $\{v_1, v_2, \dots, v_n\}$. If an ordering relationship exists between the domain values of an attribute, an elementary expression is denoted as $expression := att \text{ in } [v_i, v_k]$, $i \leq k$, in which $[v_i, v_k]$ represents the values within the range of v_i and v_k .

An attribute is allowed to participate at most once in an individual. This restriction is imposed to prevent the exploration of expressions to which no tuples satisfy. In the following, an expression to which no tuples qualify will be called an *empty* expression. Consider a database in which,



- $p_1 = \text{gender is ('male')} \wedge \text{age in [19,34]}$
- $p_2 = \text{age in [29,44]} \wedge \text{town is ('Almere', 'Amsterdam', 'Weesp')} \wedge \text{category is ('lease')}$
- $p_3 = \text{gender is ('male')} \wedge \text{age in [29,34]} \wedge \text{category is ('lease')}$
- $p_4 = \text{gender is ('female')} \wedge \text{age in [29,40]} \wedge \text{category is ('lease')} \wedge \text{price in [50K, 100K]}$
- $p_5 = \text{gender is ('male')} \wedge \text{price in [20K,45K]}$

Fig. 1 Example of a population

among others, the age of persons is recorded. Then, the expression $\text{age in [19,34]} \wedge \text{age in [39,44]}$ represents the class of persons whose age is between 19 and 34 as well as between 39 and 44. It should be clear that no persons will satisfy this expression, since *age* is a single-valued attribute.

In the following, a population is defined as a set of individuals. As a running example, we use a database that keeps record of cars and their owners. This artificial database consists of the following universal relation¹: $Ex(\text{gender}, \text{age}, \text{town}, \text{category}, \text{price}, \text{damage})$, in which the attributes *gender*, *age*, and *town* refer to the owner and the remainder of the attributes refer to the car. Attribute *category* records whether a car is leased or not and *damage* records whether a car has been involved in an accident or not. An example of a population consisting of 5 individuals is given in Figure 1.

3.2 Fitness function

A central instrument in a genetic algorithm is the fitness function. Since a genetic algorithm is aimed to the optimization of such a function, this function is one of the keys to success. Consequently, a fitness function should represent all issues that play a role in the optimization of a specific problem. Before enumerating the issues that play a role in the context of data mining, we introduce the notions of cover and success.

Definition 1: Let D be a database and p an individual defined on D . Then, the number of tuples that satisfies the expression corresponding to p is called the cover of p , and is denoted as $\|\sigma_p(D)\|$. The set of tuples satisfying p is denoted as $\sigma_p(D)$.

Note that p can be regarded as a description of a class in D and $\sigma_p(D)$ summarizes the tuples satisfying p . Within a class we can define subclasses. In the following, we regard classification problem as follows: *Given a target class t , search interesting subclasses, i.e., individuals, within class t .* We note that the target class is the class of tuples in which interesting knowledge should

¹We note that a universal relation can be obtained by performing a number of joins between the relations involved in a database.

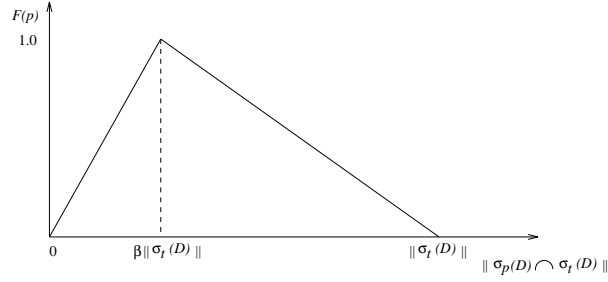


Fig. 2 Shape of the fitness function

be searched for. Suppose we want to expose the profiles of risky drivers, i.e., the class of persons with (more than average) chances of causing an accident, from the database $Ex(\text{gender}, \text{age}, \text{town}, \text{category}, \text{price}, \text{damage})$. Then, these profiles should be searched for in a class that records the characteristics of drivers that caused accidents. Such a class may be described as $\text{damage} = \text{'yes'}$.

We feel that the following issues play a role in classification problems.

- The cover of the target class. Since results from data mining are used for informed decision making, knowledge extracted from databases should be supported by a significant part of the database. This increases the reliability of the results. So, a fitness function should take into account that small covers are undesired.
- The ratio of the cover of an individual p to the cover of the target class t , i.e., $\frac{\|\sigma_p(D) \cap \sigma_t(D)\|}{\|\sigma_t(D)\|}$. If the ratio is close to 0, this means that only a few tuples of the target class satisfy individual p . This is undesired for the same reason as a small cover for a target class. If the ratio is close to 1, almost all tuples of the target class satisfy p . This is also undesired because this will result in knowledge that is often known. A fitness function should take these properties into account.

Taking into account above-mentioned issues, we have defined the following fitness function:

$$F(p) = \begin{cases} \frac{\|\sigma_p(D) \cap \sigma_t(D)\|}{\beta \|\sigma_t(D)\|} C(t) & \text{if } \|\sigma_p(D) \cap \sigma_t(D)\| \leq \beta \|\sigma_t(D)\| \\ \frac{\|\sigma_p(D) \cap \sigma_t(D)\| - \|\sigma_t(D)\|}{\|\sigma_t(D)\|(\beta - 1)} C(t) & \text{otherwise} \end{cases}$$

in which $0 < \beta \leq 1$, and

$$C(t) = \begin{cases} 0 & \text{if } \frac{\|\sigma_t(D)\|}{\|\sigma(D)\|} \leq \alpha \\ 1 & \text{otherwise} \end{cases}$$

and $0 \leq \alpha \leq 1$.



We note that the values for α and β should be defined by the user and will vary for different applications. The value of α defines the fraction of tuples that a target class should contain in order to be a candidate for further exploration. The value β defines the fraction of tuples that an individual should represent within a target class in order to obtain the maximal fitness. In Figure 2, the shape of the fitness function is presented.

The fitness grows linearly with the number of tuples satisfying the description of an individual p as well as satisfying a target class t , i.e., $\|\sigma_p(D) \cap \sigma_t(D)\|$, above a user-defined value α , and decreases linearly with $\|\sigma_p(D) \cap \sigma_t(D)\|$ after reaching the value $\beta\|\sigma_t(D)\|$.

It should be clear that our goal is to search for those individuals that approximate a fitness of $\beta\|\sigma_t(D)\|$. Consider the target class *damage = yes* that consists of 100.000 tuples. Assume that a profile is considered risky if about 30.000 out of 100.000 persons satisfy this profile. This means that $\beta \approx 0.3$. Assuming that 33.000 of the persons caused an accident are young males, the algorithm should find individuals like (*gender is ('male') ∧ age in [19,28]*).

3.3 Manipulation operators

In this section, we discuss the effect of the mutation and cross-over operator on an individual.

3.3.1 Mutation

As stated in the introduction of this section, a mutation modifies an individual. In defining the mutation operator, we take into account the domain type of an attribute. If there exists no ordering relationship between the domain values, then we select randomly an attribute value and replace it by another value, which can be a NULL value as well, in an expression that contains this attribute. For example, a mutation on attribute *town* of individual p_2 (see Figure 1) may result into $p'_2 = \text{age in [29,44]} \wedge \text{town is ('Almere', 'Den Haag', 'Weesp')} \wedge \text{category is ('lease')}$.

If there exists a relationship between the domain values of an attribute, the mutation operator acts as follows in the case that a single value is associated with this attribute in an expression, i.e., the expression looks as *att is (v_c)*. Let $[v_b, v_e]$ be the domain of attribute *att*. In order to mutate v_c , we choose randomly a value $\delta_v \in [0, (v_e - v_b)\mu]$, in which $0 \leq \mu \leq 1$. The parameter μ is used to control the maximal increase or decrease of an attribute value. The mutated value v'_c is defined as $v'_c = v_c + \delta_v$ or $v'_c = v_c - \delta_v$ as long as $v'_c \in [v_b, v_e]$. To handle overflow, i.e., if $v'_c \notin [v_b, v_e]$, we assume that the successor of v_e is v_b , and, consequently the predecessor of v_b is v_e . To compute a mutated value v'_c appropriately, we distinguish between whether v_c will be increased or decreased, which is randomly determined.

In the case that v_c is increased

$$v'_c = \begin{cases} v_c + \delta_v & \text{if } v_c + \delta_v \in [v_b, v_e] \\ v_b + \delta_v - (v_e - v_c) & \text{otherwise} \end{cases}$$

and in the case v_c is decreased

$$v'_c = \begin{cases} v_c - \delta_v & \text{if } v_c - \delta_v \in [v_b, v_e] \\ v_e - \delta_v + (v_c - v_b) & \text{otherwise} \end{cases}$$

Let us consider the situation in which more than one value is associated with an attribute *att* in an expression. If a list of non successive (enumerable) values is associated with *att*, we select one of the values and compute the new value according to one of the above-mentioned formulas. If a range of successive values, i.e., an interval, is associated with *att*, we select either the lower or upper bound value and mutate it. A potential disadvantage of this strategy for intervals is that an interval may be significantly enlarged, if the mutated value crosses a domain boundary. Suppose that the domain of *age* is [18,60], and we mutate the upper bound value of the expression *age in* [55, 59], i.e., the value 59. Assume that the value 59 is increased by 6, then 59 is mutated in the value 23. The new expression becomes *age in* [23, 55].

We note that the partitioning of attribute values, i.e., the selection of proper intervals in an expression, is simply adjusted by the mutation operator. As noted before, partitioning of attribute values is in general a tough problem [Ref 21].

3.3.2 Cross-over

The idea behind a crossover operation is as follows; it takes as input 2 expressions, selects a random point, and exchanges the subexpressions behind this point. To illustrate this idea, we consider a relation $R(att_1, att_2, \dots, att_n)$ and two expressions, e^i and e^j , in which *all* attributes are involved. Let e^i be defined as follows: $(e_1^i \wedge e_2^i \wedge e_3^i \wedge \dots \wedge e_{k-1}^i \wedge e_k^i \wedge e_{k+1}^i \wedge \dots \wedge e_n^i)$, in which e_l^i represents an elementary expression in which attribute *att* is involved. And, let e^j be defined as $(e_1^j \wedge e_2^j \wedge e_3^j \wedge \dots \wedge e_{k-1}^j \wedge e_k^j \wedge e_{k+1}^j \wedge \dots \wedge e_n^j)$. Then, a cross-over between e^i and e^j at point k may result into the following two expressions, namely $e^{i'} = (e_1^i \wedge e_2^i \wedge e_3^i \wedge \dots \wedge e_{k-1}^i \wedge e_k^j \wedge e_{k+1}^j \wedge \dots \wedge e_n^j)$ and $e^{j'} = (e_1^j \wedge e_2^j \wedge e_3^j \wedge \dots \wedge e_{k-1}^j \wedge e_k^i \wedge e_{k+1}^i \wedge \dots \wedge e_n^i)$.

In general, not all attributes will be involved in an expression. This may have some undesired effects for a cross-over. First, a cross-over may produce individuals in which an attribute is involved more than once. For example, a cross-over between the individuals $p_1 = \textit{gender is}$ ('male') $\wedge \textit{age in}$ [19,34] and $p_2 = \textit{age in}$ [29,44] $\wedge \textit{town is}$ ('Almere', 'Amsterdam', 'Weesp') $\wedge \textit{category is}$ ('lease') after the first attribute results into the following individuals: $p_1' = \textit{gender is}$ ('male') $\wedge \textit{town is}$ ('Almere', 'Amsterdam', 'Weesp') $\wedge \textit{category is}$ ('lease') and $p_2' = \textit{age in}$ [19,34] $\wedge \textit{age}$



in [29,44]. As we can see, the attribute *age* appears twice in p'_2 .

Second, a cross-over may result in an offspring that is exactly the same as the parents with probability 1.0. For example, a cross-over between p_1 and p_2 that occurs on a point that is beyond the last elementary expression of p_2 , i.e., *category is* ('lease'), will result into the equal new individuals.

To prevent the above-mentioned effects, we apply the following technique to perform cross-overs. Consider two individuals p_i and p_j that have been selected for crossover. Let A_i be the set of attributes that is not involved in p_i but is involved in p_j and A_j the set of attributes that is not involved in p_j but is involved in p_i . Then, for each attribute in A_i , we generate an empty elementary expression using this attribute and add it to p_i . The same procedure is applied to the attributes of A_j . We note that this procedure has as effect that the length of p_i and p_j become equal. Finally, we regard an individual as a sequence of elementary expressions, and order these in p_i and p_j according to the rule that elementary expressions having the same attribute will appear at the same position in p_i and p_j . Then, the cross-over between p_i and p_j can be performed. We note that the cross-over point should be chosen between the first and final position of p_i or p_j . The following example illustrates this technique.

Example 1 Consider the individuals $p_1 = \textit{gender is} ('male') \wedge *age in* [19,34] and $p_2 = \textit{age in}$ [29,44] \wedge *town is* ('Almere', 'Amsterdam', 'Weesp') \wedge *category is* ('lease') again. Then, $A_1 = \{\textit{town, category}\}$ and $A_2 = \{\textit{gender}\}$. So, we extend p_1 with the following expression *town is* (") \wedge *category is* (") and p_2 is extended with *gender is* (").$

After ordering the elementary expressions p_1 and p_2 look as follows:

$$\begin{aligned} p_1 &= \textit{gender is} \text{ ('male')} \wedge \textit{age in} [19,34] \wedge \textit{town is} \text{ (')} \wedge \textit{category is} \text{ (')} \\ p_2 &= \textit{gender is} \text{ (')} \wedge \textit{age in} [29,44] \wedge \textit{town is} \text{ ('Almere', 'Amsterdam', 'Weesp')} \wedge \\ &\quad \textit{category is} \text{ ('lease')} \end{aligned}$$

Now, a cross-over at position 2 results into

$$\begin{aligned} p'_1 &= \textit{gender is} \text{ ('male')} \wedge \textit{age in} [19,34] \wedge \textit{town is} \text{ ('Almere', 'Amsterdam', 'Weesp')} \wedge \\ &\quad \textit{category is} \text{ ('lease')} \\ p'_2 &= \textit{gender is} \text{ (')} \wedge \textit{age in} [29,44] \wedge \textit{town is} \text{ (')} \wedge \textit{category is} \text{ (')} \end{aligned}$$

Note that p'_2 is equal to *age in* [29,44]. \square

In the next chapter, we introduce a number of rules that may prevent the exploration of unpromising individuals.

4 Optimization rules

In this chapter, we discuss two propositions that may be used to prevent the exploration of unprofitable individuals. These propositions are derived from the shape of the fitness function. The complexity of a genetic-based algorithm is determined by the evaluation of the fitness function, since this is computationally the most expensive operation [Ref 20]. Before presenting these propositions, we introduce the notion of a similar of an individual.

Definition 2: Let $\text{length}(p)$ be the number of elementary expressions involved in p . An individual p_{sim} is a *similar* of p if each elementary expression of p_{sim} is contained in p or p_{sim} contains each elementary expression of p and $\text{length}(p_{sim}) \neq \text{length}(p)$.

As stated in the foregoing, we search for individuals with high values for the fitness function F . Recall that this function looks as follows:

$$F(p) = \begin{cases} \frac{\|\sigma_p(D) \cap \sigma_t(D)\|}{\beta \|\sigma_t(D)\|} C(t) & \text{if } \|\sigma_p(D) \cap \sigma_t(D)\| \leq \beta \|\sigma_t(D)\| \\ \frac{\|\sigma_p(D) \cap \sigma_t(D)\| - \|\sigma_t(D)\|}{\|\sigma_t(D)\|(\beta-1)} C(t) & \text{otherwise} \end{cases}$$

in which $\sigma_t(D)$, β , and $C(t)$ are constant during a mining session.

Note that the computation of $F(p)$ requires the number of tuples that satisfy individual p . So, these tuples should be searched for and retrieved from the database, which is a costly operation [Ref 10]. Although several techniques may be used to minimize the number of retrievals from a database [Ref 8, 9], still large amounts of tuples have to be retrieved from the database in mining applications. The techniques to minimize retrievals are mainly based on storing frequently used tuples in an efficient way in main memory. In this way, disk accesses are reduced.

In the following, two propositions will be presented that may be used to avoid the computation of fitness values of unprofitable individuals. These propositions decide if the fitness value of a similar of an individual p is worse than the fitness of p . If this is the case, this similar can be excluded from the search process.

Proposition 1: Let p_{sim} be a similar of p . If $\|\sigma_p(D) \cap \sigma_t(D)\| \leq \beta \|\sigma_t(D)\|$ and $\text{length}(p_{sim}) > \text{length}(p)$ then $F(p_{sim}) \leq F(p)$.

Proof. From $\text{length}(p_{sim}) > \text{length}(p)$ follows that $\sigma_{p_{sim}}(D) \subseteq \sigma_p(D)$. As a consequence, $\|\sigma_{p_{sim}}(D) \cap \sigma_t(D)\| \leq \|\sigma_p(D) \cap \sigma_t(D)\|$. Since $\|\sigma_p(D) \cap \sigma_t(D)\| \leq \beta \|\sigma_t(D)\|$, it follows $F(p_{sim}) \leq F(p)$. \square



Proposition 2: Let p_{sim} be a similar of p . If $\|\sigma_p(D) \cap \sigma_t(D)\| \geq \beta\|\sigma_t(D)\|$ and $\text{length}(p_{sim}) < \text{length}(p)$ then $F(p_{sim}) \leq F(p)$.

Proof. Similar to the proof of Proposition 1. \square

Note that the propositions do not require additional retrievals from a database to decide if $F(p_{sim}) \leq F(p)$.

We discuss two alternatives how these propositions may contribute in optimizing the search process. The first alternative is an application of the propositions at cross-over level, while the second is an application at population level.

As stated in the foregoing, a cross-over is applied on a mating pair and results into two offsprings. Suppose that a mutation is performed after a cross-over, and the parent and the offspring with the highest fitness values are eligible to be mutated (see Chapter 5). Consider an offspring p resulted from a cross-over, and let p_o be a similar of p , one of its parents. If we can decide that $F(p_o) \leq F(p)$, then it is efficient to mutate p_o . The reason is that computation on an unmutated p_o will be a wasting of effort.

Now, we discuss how above-mentioned propositions may be applied at population level. Before computing the fitness values of the individuals of generation P_{n+1} , we compare each individual with the individuals in generation P_n . For each individual p_{n+1} in P_{n+1} , we check if p_{n+1} has similars in the generation P_n . If this is the case and we can decide that the fitness of a similar is worse or equal than the fitness of p_{n+1} , we replace p_{n+1} by another individual (which is not a similar of p_{n+1}). In this way, we save computational efforts.

In the next chapter, we propose an overall algorithm, in which we apply the two propositions at cross-over level.

5 Algorithm

The previous chapter was devoted to the major issues that play a role in designing a genetic-based algorithm for data mining. In this chapter, we describe the overall algorithm. Before starting this description, we discuss a mechanism to select an individual for a next generation.

The mechanism to select individuals for a new generation is based on the technique of elitist recombination [Ref 22]. According to this technique, the individuals in a population are randomly shuffled. Then, the cross-over operation is applied on each mating pair, resulting into two offsprings. The parent and the offspring with the highest fitness value are selected for the next generation. In this way, there is a direct competition between the offsprings and their own parents. Note, the offspring provides the possibility to explore different parts of a search space.

The elitist recombination technique has been chosen for two reasons. First, there is no need to specify a particular cross-over probability, since each individual is involved in exactly one cross-over. Second, there is no need for intermediate populations in order to generate a new population as is the case in a traditional genetic algorithm. These properties simplify the implementation of a genetic algorithm. Let us outline the overall algorithm.

$P(t)$	=	population at time t
$F'(P(t))$	=	$\sum_{p \in P(t)} F(p)$
p_j	=	j -th individual in population $P(t)$
ϵ	=	threshold value
$\text{cross-over}(p_1, p_2, o_1, o_2)$	=	procedure that takes two individuals p_1 and p_2 as input, applies a cross-over, and produces two offsprings o_1 and o_2 .
$\text{mutate}(p, c, p')$	=	procedure that mutates an individual p with probability c into p' .
$\text{similar}(p_1, p_2)$	=	boolean function that decides whether two individuals are similar or not.

Table 1 Meaning of symbols and procedures used in Figure 3

The algorithm starts with the initialization of a population consisting of an even number of individuals, called $P(t)$. The individuals in this population are shuffled. Then, the cross-over operation is applied on two successive individuals. We note that the cross-over operator is applied exactly once for an individual. After completion of a cross-over, the fitness values of the parents are compared¹; the parent with the highest value is selected and it may be mutated with a probability c . This parent, p'_{sel} , is added to the next generation, and in case it is mutated its fitness value is computed.

¹These values are already computed and stored by the algorithm.



```

program Genetic Algorithm;
initialize( $P(t)$ );
FOR  $p \in P(t)$  DO  $F(p)$  OD;
 $F'(P(t+1)) := \epsilon + 1$ ;  $F'(P(t)) := 0$ ;
WHILE  $F'(P(t+1)) - F'(P(t)) \geq \epsilon$  DO
     $F'(P(t)) := F'(P(t+1)) := 0$ ;
     $j := 1$ ;
    shuffle( $P(t)$ );
    WHILE  $j < population\_size$  DO
        cross-over( $p_j, p_{j+1}, o_1, o_2$ );
        IF  $F(p_j) > F(p_{j+1})$  THEN  $p_{sel} := p_j$ 
            ELSE  $p_{sel} := p_{j+1}$ 
        FI;
        mutate( $p_{sel}, c, p'_{sel}$ );
        IF  $p'_{sel} \neq p_{sel}$  THEN  $F(p'_{sel})$  FI;
         $p'_{sel} \rightarrow P(t+1)$ ;
        FOR  $k = 1, 2$  DO
            IF (similar( $p'_{sel}, o_k$ ) AND  $F(o_k) \leq F(p'_{sel})$ )
                THEN mutate( $o_k, 1.0, o'_k$ );
                ELSE mutate( $o_k, c, o'_k$ );
            FI;
        OD;
         $F(o'_1); F(o'_2)$ ;
        IF  $F(o'_1) > F(o'_2)$  THEN  $o'_1 \rightarrow P(t+1)$ 
            ELSE  $o'_2 \rightarrow P(t+1)$ ;
        FI;
         $j := j + 2$ ;
    OD;
    For  $p \in P(t)$  DO  $F'(P(t)) := F'(P(t)) + F(p)$  OD;
    For  $p \in P(t+1)$  DO  $F'(P(t+1)) := F'(P(t+1)) + F(p)$  OD;
     $P(t+1) := P(t)$ ;
OD;
END.

```

Fig. 3 Sketch of the genetic-based algorithm

Then, for each offspring, p_o , we test if this offspring is a similar of p'_{sel} and if its fitness value is worse or equal than p'_{sel} . If this is the case, p_o is an unpromising individual, and, therefore, we always mutate p_o . Otherwise, we mutate p_o with probability c . Note, to compare the fitness value between p'_{sel} and p_o , the propositions of the previous chapter are used. So, no additional fitness

values are computed for this comparison. After possible mutation of the offsprings, their fitness values are computed, and the most fittest offspring is added to the new generation. This process is repeated for all individuals in a generation.

Once the new population has been built up, the total fitness of the existing as well as of the new population is computed, and compared. The algorithm terminates if the total fitness of the new population does not significantly improve compared with the total fitness of the existing population, i.e., that the improvement of the total fitness of the new population is less than a threshold value ϵ . In Figure 3, the algorithm is sketched, and in Table 1 the meaning of some symbols and procedures used in Figure 3 are listed. For a complexity analysis of the algorithm, we refer to [Ref 20].

In the next chapter, we discuss the implementation of the algorithm and some preliminary results.

6 Implementation and preliminary results

Based on the algorithm described in Chapter 5, we have built a re-targetable prototype of a data mining tool, which means that the tool can be coupled to different databases. The main goal of the current effort is to determine if the genetic-based algorithm is able to find hidden knowledge. Let us continue with the description of the tool.

The tool takes as input a mining question and produces a set of answers for the question. The prototype is running in a Microsoft Access 97 environment that uses the Microsoft Jet Engine¹. The genetic-based algorithm is implemented in Visual Basic. We have chosen this environment for two reasons. First, this environment is available at our laboratory. Second, the database that we want to mine is a Microsoft database.

In Figure 4, the architecture of our tool is represented. Once the tool receives a mining question,

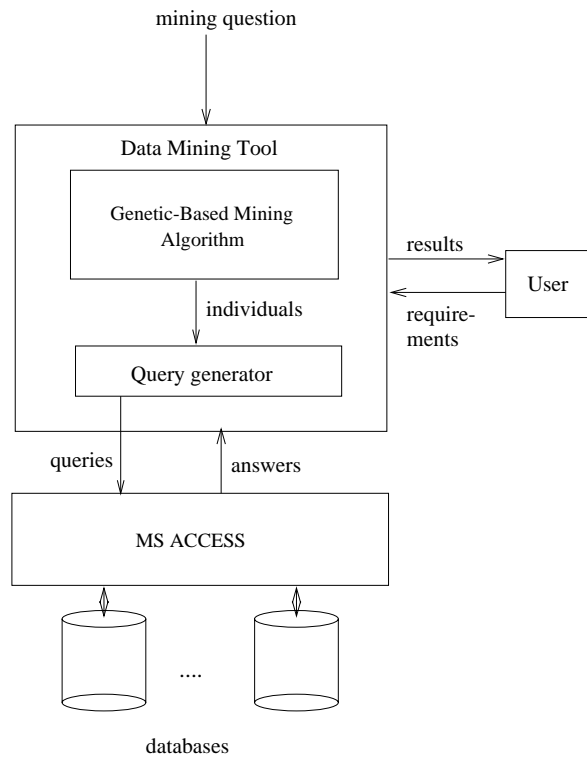


Fig. 4 Architecture of a data mining tool

it runs the genetic-based mining algorithm, which generates, among others, a population. The individuals of the population are passed to the Query Generator. For each individual a correspond-

¹Within the scope of a feasibility study, a previous version of the algorithm was implemented in C and connected to the Monet Database Server [Ref 6, 7].

ing SQL query is generated. These queries are passed to the MS ACCESS DBMS, which on its turn processes the queries and passes the results to the data mining tool. These results are used to compute the fitness of each individual. Upon request of the user individuals and their associated fitness values can be shown. The user has the possibility to modify the initial mining question or to specify additional requirements with regard to the question. We note that this is a very useful feature of a mining tool, since, in practice, a user starts a mining session with a rough idea of what information might be interesting, and during the mining session (with help of the results provided by the system) the user specifies more precisely what information should be searched for.

The generation of a query corresponding to an individual p is straightforward. Recall that an individual is a conjunction of predicates. So, this forms the major part of the WHERE clause of a query corresponding to p . Since we require the number of tuples satisfying p and a target class t , the query corresponding to p is: *select count(*) from database where $p \wedge t$.*

We have applied the tool on an artificial database, and are currently applying it on a real-life database. In the following, we describe the databases and the results obtained so far.

Artificial database

This database consists of the relation *Ex(gender, age, town, category, price, damage)*. For this database 100.000 tuples have been generated, of which 50% have a value "yes" for *damage*, i.e., 50% of the tuples relate to an accident. Furthermore, the fact that young men in lease cars have more than average chances to cause an accident was hidden in the database. The goal of mining this database was to determine whether the tool is capable to find the hidden fact. Therefore, we have set the target class as *damage = 'yes'*, and we searched for the profile of risky drivers. We note that the expression for the hidden profile is: *age in [19,24] \wedge category is ('lease') \wedge gender is ('male')*.

We have mined the database with varying initial populations, consisting of 36 individuals. The following classes of initial population were distinguished.

- random: This population contained a few individuals that could set the algorithm quickly on a promising route.
- modified random: Individuals that could apparently set the algorithm on a promising route were replaced by other (not promising) individuals.
- bad converged: This population contained individuals with low fitness values. The algorithm has to manipulate the initial population significantly in order to find the hidden profile.

We have observed that the algorithm usually finds near optimal solutions, i.e. profiles that look



like the hidden one, in less than 1000 fitness evaluations. The differences between the hidden profile and profiles found by the algorithm (for different initial populations) were mainly caused by variations in the range of attribute *age*.

The type of the initial population plays a role in the number of fitness evaluations required to find a near optimal solution. Running the algorithm on the database with random initial populations, the algorithm was able to find a near optimal expression quite rapidly, i.e., in about 100 fitness evaluations. Starting from modified random initial populations, 300 to 400 fitness evaluations were required for a near optimal solution. Starting from bad converged initial population, 900 to 1000 fitness evaluations were required.

With regard to the settings of the parameters α and β , we note that appropriate values could be easily selected, since the content of the database is precisely known. For example, α should be less than 0.5. The mutation probability was set on 0.20 and ϵ was set on 0.01. These values were simply chosen on the basis of our intuition. The impact of the parameters on the number of fitness evaluations and the solutions generated by the algorithm is a topic for further research. \square

Real-life database

Currently, we are mining a real-life database, the so-called FAA incident database, which is available at our aerospace laboratory, NLR. This database contains aircraft incident data that are recorded from 1978 to 1995. Incidents are potentially hazardous events that do not meet the aircraft damage or personal injury thresholds as defined by American National Transportation Safety Board (NTSB). For example, the database contains reports of collisions between aircraft and birds while on approach to or departure from an airport. While such a collision may not have resulted in sufficient aircraft damage to reach the damage threshold of an NTSB accident, the fact that the collision occurred is valuable information. The FAA database consists of more than 70 attributes and about 80.000 tuples.

The initial mining task on the FAA database was: search for the class of flights with (more than average) chances of causing an incident, i.e., profiles of risky flights. This search resulted in (valid) profiles but which could be easily declared. An example of such a profile is that aircraft with 1 or 2 engines are more often involved in incidents. The explanation for this profile is that these types of aircraft perform more flights.

During the mining process the mining question was refined in the following more specific questions:

- Given the fact that an incident was due to operational defects not inflicted by the pilot, what is the profile of this type of incident?
- Given the fact that an incident was due to mistakes of the pilot, what is the profile of this type of incident?
- Given the fact that an incident was due to improper maintenance, what is the profile of this type of incident?

We have proposed these questions to our tool with the following values for the parameters, $\alpha = 0$, $\beta = 0.25$, mutation probability (c) = 0.3, and $\epsilon = 0.1$. Furthermore, the population size was set on 50. On the first glance, the results of the tool appear to be promising. Safety experts at our laboratory are analysing the results. On the basis of their analysis, we will set up a plan to mine the database more systematically, and to study the impact of different parameter values on the results provided by the tool. The goal of the latter study is to formulate some guidelines for selecting parameter values for similar type of databases, such as an aircraft accident database. \square

Although our evaluation is not completed yet and a significant amount of research has to be done, e.g., on performance issues, in order to build an adequate genetic-based data mining tool, the preliminary results are promising. A second observation is that the range interval of an attribute in an expression may significantly enlarged, if a mutation occurs on a domain boundary. An interval that consists of (almost) the whole the domain slows down the search process. In a next version of the tool, we will enhance the mutation operator. An alternative is to clip on boundary values in cases of overflow.



7 Conclusions & further research

In order to answer mining questions, very large search spaces should be inspected, making an exhaustive search infeasible. So, heuristic search strategies are of vital importance in searching such spaces. In this paper, we have discussed a genetic-based algorithm that may be used for data mining. Contrary to the conventional bit string representation in genetic algorithms, we have chosen a representation that fits better in the field of databases. Furthermore, we have chosen a fitness function that is close to our intuition to rank individuals. The fitness function gives rise to an optimization of the search process.

A genetic-based algorithm for data mining has two major advantages. First, the problem of partitioning attribute values in proper ranges, which is in general a tough problem [Ref 21], could be solved by choosing a suitable mutation operator. Second, a genetic-based algorithm is able to escape a local optimum and does not pose any restrictions on the structure of a search space.

By means of a (prototype) implementation and a preliminary evaluation, we have shown the potentials of a genetic-based data mining tool. Since the preliminary results of the tool appear to be promising, we are setting up a research plan to evaluate this tool thoroughly. The outcome of the evaluation will determine our future research activities in this field.

Acknowledgements The author is grateful to Wim Pelt from the Royal Netherlands Navy, who made this research possible. Hein Veenhof and Egbert Boers from NLR are thanked for their valuable comments on earlier drafts of this paper. Finally, the students Leo de Penning and Martijn Suurd from the University of Twente are thanked for their implementation efforts.

References

1. Agrawal, R., Ghosh, S., Imielinski, T., Iyer, B., Swami, A., An Interval Classifier for Database Mining Applications, in Proc. of the 18th Very Large Data Base, 1992, pp. 560-573.
2. Agrawal, R., Imielinski, T., Swami, A., Mining Association Rules between Sets of Items in Large Databases, in Proc. ACM SIGMOD '93 Int. Conf. on Management of Data, 1993, pp. 207-216.
3. Agrawal, R., Srikant, R., Fast Algorithms for Mining Association Rules, in Proc. Int. Conf. on Very Large Databases, 1994, pp 487-499.
4. Agrawal, R., Srikant, R., Mining Sequential Patterns, in Proc. 11th Int. Conf. on Data Engineering, 1995, pp. 3-14.
5. Augier, S., Venturini, G., Kodratoff, Y., Learning First Order Logic Rules with a Genetic Algorithm, in Proc. 1st Int. Conf. on Knowledge Discovery and Data Mining, 1995, pp. 21-26.
6. Boncz, P., Kersten, M., Monet: an impressionist sketch of an advanced database system, in Proc. Basque Int. Workshop on Information Technology, 1995.
7. Boncz, P., Wilschut, A., Kersten, M., Flattening an Object Algebra to Provide Performance, in Proc. 14th Int. Conf. on Data Engineering, 1998, pp. 568-577.
8. Choenni, R., Siebes, A., Query Optimization to Support Data Mining, in Proc. DEXA '97 8th Int. Conference and Workshop on Database and Expert Systems Applications, 1997, pp. 658-663.
9. Choenni, R., Kersten, M., Saad, A., Akker van den, J., A Framework for Multi-Query Optimization, in Proc. COMAD '97 8th Int. Conference on Management of Data, 1997, pp. 165-182.
10. Elmasri, R., Navathe, S., Fundamentals of Database Systems, The Benjamin/Cummings Publishing Company, 1989.
11. Fayyad, U.M., Piatetsky-Shapiro, G., Smyth, P., Uthurusamy, R., (Eds), Advances in Knowledge Discovery and Data Mining, AAAI/The MIT Press, 1996.
12. Freitas, A., A Genetic Programming Framework for two Data Mining Tasks: Classification and Generalized Rule Induction, in Proc. Int. Conf on Genetic Programming, 1997, pp. 96-101.
13. Han, J., Cai, Y., Cerone, N., Knowledge Discovery in Databases: An Attribute-Oriented Approach, in Proc. of the 18th Very Large Data Base, 1992, pp. 547-559.
14. Holland, J.H., Adaption in Natural and Artificial Systems, Ann Arbor: University of Michigan Press, USA, 1975.



15. Holsheimer, M., Kersten, M.L., Architectural Support for Data Mining, in Proc. AAAI-94 Workshop on Knowledge Discovery, 1994, pp. 217-228.
16. Houtsma, M., Swami, A., Set-Oriented Mining for Association Rules in Relational Databases, in Proc. 11th Int. Conf. on Data Engineering, 1995 pp. 25-33.
17. Kemenade van, C., Akker van den, M., Kok, J., Evolutionary Air Traffic Flow Management for large 3D-problems, in Proc. Int. Conf. Parallel Problem Solving from Nature, 1996, pp. 910-919.
18. Mannila, H., Toivonen, H., Verkamo, A.I., Efficient Algorithms for Discovering Association Rules, in Proc. AAI-94 Workshop on Knowledge Discovery, 1994, pp. 181-192.
19. Michalewicz, Z., Genetic Algorithms + Data Structures = Evolution Programs. Springer-Verlag, New York, USA.
20. Penning, H. de, Suurd, P., NLR Genetic Search Base, internal report, University of Twente, 1998.
21. Srikant, R., Agrawal, R., Mining Quantitative Association Rules in Large Relational Tables, in Proc. ACM SIGMOD '96 Int. Conf. on Management of Data, 1996, pp. 1-12.
22. Thierens, D., Goldberg, D., Elitist Recombination: an integrated selection recombination GA, 1st IEEE Conf. on Evolutionary Computing, 1994, pp. 508-512.