

TUTORIAL GENETIC ALGORITHMS

Alexandre P. Alves da Silva

Universidade Federal de Itajubá, Instituto de Engenharia Elétrica
Av. BPS, 1303 - Itajubá - MG - CEP 37500-903 - Brazil
alex@iee.efei.br

Abstract - Research on genetic algorithms (GAs) has shown that the initial proposals are incapable of solving hard problems in a robust and efficient way. Usually, for large-scale optimization problems, the execution time of first generation GAs dramatically increases while solution quality decreases. The focus of this tutorial is pointing out the main design issues in tailoring GAs to large-scale optimization problems. Important topics such as encoding schemes, selection procedures, self-adaptive and knowledge-based operators are discussed.

Index Terms - Genetic algorithms, large-scale optimization.

1. MODERN HEURISTIC SEARCH TECHNIQUES

Optimization is the basic concept behind the application of genetic algorithms (GAs), or any other evolutionary algorithm [1]-[3], to any field of interest. Over and above the problems in which optimization itself is the final goal, it is also a way for (or the main idea behind) modeling, forecasting, control, simulation, etc.. Traditional optimization techniques begin with a single candidate and iteratively search for the optimal solution by applying static heuristics. On the other hand, the GA approach uses a population of candidates to search several areas of a solution space, simultaneously and adaptively.

Evolutionary computation allows precise modeling of the optimization problem, although not usually providing mathematically optimal solutions. Another advantage of using evolutionary computation techniques is that there is no need for having an explicit objective function. Moreover, when the objective function is available, it does not have to be differentiable.

Genetic algorithms have been most commonly applied to solve combinatorial optimization problems. Combinatorial optimization usually involves a huge number of possible solutions, which makes the use of enumeration techniques (e.g., cutting plane, branch and bound, or dynamic programming) hopeless.

In large-scale combinatorial optimization problems, the number of possible solutions grow exponentially with the problem size. Therefore, the application of optimization methods to find the optimal solution is computationally impracticable. Heuristic search techniques are frequently employed in this case for achieving high quality solutions within reasonable run time.

Among the heuristic search methods there are the ones that apply local search (e.g., hill climbing) and the ones that use a non-convex optimization approach, in which cost-deteriorating neighbors are accepted also. The most popular methods which go beyond simple local search are GAs [4]-[7] (and other evolutionary techniques, like evolutionary programming, evolutionary strategies, etc.), simulated annealing (SA) [8], and tabu search (TS) [9]. Particle swarm [10] is another optimization technique that has shown great potential, lately. However, more experience is still necessary to prove its efficiency and robustness.

Simulated annealing uses a probability function that allows a move to a worse solution with a decreasing probability, as the search progresses. With GAs, a pool of solutions is used and the neighborhood function is extended to act on pairs of solutions. Tabu search uses a deterministic rather than stochastic search. Tabu search is based on neighborhood search with local optima avoidance. In order to avoid cycling, a short-term adaptive memory is used in TS. Genetic algorithms have a basic distinction when compared with other methods based on stochastic search. They can use coding (genotypic space) for representing the problem. The other methods solve the optimization problem in the original representation space (phenotypic).

The most rigorous global search methods have asymptotic convergence proof (also known as convergence in probability), i.e., the optimal solution is guaranteed to be found if infinite time is available. Among SA, GA and TS algorithms, simulated annealing and genetic algorithms are the only ones with proof of convergence. However, there is no such proof for the canonical GA [11], i.e., the one with proportional selection (Section 5.1.6) and crossover/mutation with constant probabilities (Section 5.4).

Although all the mentioned algorithms have been successfully applied to real world problems, several of their crucial parameters have been selected empirically. Theoretical knowledge of the impact of these parameters on convergence is still an open problem. In fact, there is no theoretical result for tabu and particle swarm searches.

The choice of representation for a GA is fundamental to achieving good results. Encoding allows a kind of *tunneling* in the original search space. That means, a particle has a non-zero probability of passing a potential barrier even when it does not have enough energy to jump over the barrier. The tunneling idea is that rather than escaping from local minima by random uphill moves, escape can be achieved with the quantum tunnel effect. It is not the height of the barrier that determines the rate of escape from a local optimum, but its width relative to current population variance.

The main shortcoming of the standard SA procedure is the slow asymptotic convergence with respect to the *temperature* parameter T . In the standard SA algorithm, the cooling schedule for asymptotic global convergence is inversely proportional to the logarithm of the number of iterations, i.e., $T(k) = c/(1 + \log k)$. The constant c is the largest depth of any local minimum that is not the global minimum. Convergence in probability cannot be guaranteed for faster cooling rates, e.g., lower values for c .

Tabu search owes its efficiency to an experience-based fine-tuning of a large collection of parameters. Tabu search is a general search scheme that must be tailored to the details of the problem at hand. Unfortunately, as mentioned before, there is little theoretical knowledge for guiding this tailoring process. Heuristic search methods utilize different mechanisms in order to explore the state space. These mechanisms are based on three basic features:

- the use of memoryless search (e.g., standard SA and GA) or adaptive memory (e.g., TS);
- the kind of neighborhood exploration used, i.e., random (e.g., SA and GAs) or systematic (e.g., TS); and
- the number of current solutions taken from one iteration to the next (GAs, as opposed to SA and TS, take multiples solutions to the next iteration).

The combination of these mechanisms for exploring the state space determines the search diversification (global exploration) and intensification (local exploitation) capabilities. The standard SA algorithm is notoriously deficient with respect to the diversification aspect. On the other hand, the standard GA is poor in intensification.

When the objective function has very many equally good local minima, wherever the starting point is, a small random disturbance can avoid the small local minima and reach one of the good ones, making this an appropriate problem for SA. However, SA is less suitable for a problem in which there is one global minimum that is much better than all the other local ones. In this case, it is very important to find that valley. Therefore, it is better to spend less time improving any set of parameters and more time working with an ensemble to examine different regions of the space. This is what GAs do best. Hybrid methods have been proposed in order to improve the robustness of the search.

2. INTRODUCTION TO GAS

Genetic algorithms operate on a population of individuals. Each individual is a potential solution to a given problem and is typically encoded as a fixed-length binary string (other representations have also been used, including character-based and real-valued encodings, etc.), which is an analogy with an actual chromosome. After an initial population is randomly or heuristically generated, the algorithm evolves the population through sequential and iterative application of three operators: selection, crossover and mutation. A new generation is formed at the end of each iteration.

For large-scale optimization problems, the initial population can incorporate prior knowledge about solutions. This procedure should not drastically restrict the population diversity, otherwise premature convergence could occur. Typical population sizes vary between 30 and 200. The population size is usually set as a function of the chromosome length. The execution of a GA iteration is basically a two stage process. It starts with the current population. Selection is applied to create an intermediate population (mating pool). Then, crossover and mutation are applied to the intermediate population to create the next generation of potential solutions. Although a lot of emphasis has been placed on the three above mentioned operators, the coding scheme and the fitness function are the most important aspects of any GA, because they are problem dependent.

The most popular explanation about how GAs can result in robust search relies on the argument of hyperplane sampling. In order to understand this concept, assume a problem encoded with 3 bits. The search space is represented by a cube with one of its vertices at the origin 000. For example, the upper surface of the cube contains all the points of the form *1*, where * could be either 0 or 1. A string that contains the symbol * is referred to as a schema. It can be viewed as a (hyper)plane representing a set of solutions with common properties.

The order of a schema is the number of fixed positions present in the string. The defining length is the distance between the first and last fixed positions of a particular schema. Building blocks are highly fit strings of low defining

length and low order. It can be shown that about $O(n^3)$ hyperplanes are simultaneously sampled when the number of strings contained in the population is n . Therefore, even though a GA never explicitly evaluates any particular hyperplane of the search space, it changes the distribution of strings as if it had.

Genetic algorithms process many hyperplanes implicitly in parallel when selection acts on the population. The true fitness of a hyperplane partition corresponds to the average fitness of all strings that lie in that hyperplane. Genetic algorithms use the population as a sample for estimating the fitness of that hyperplane partition. After the initial generation, the pool of new strings are biased toward regions that have previously contained strings that were above average with respect to previous populations. In order to further explore the search space, crossover and mutation generate new sample points, while partially preserving the distribution of strings that is observed after selection.

In the following sections, several important design stages of a GA are presented. Section 3 shows different possibilities for encoding. It emphasizes the importance of the encoding scheme on GA convergence. Section 4 treats the formulation of the fitness function. Section 5 presents different propositions for the selection, crossover and mutation operators. Parameter control in GAs is addressed in Section 5, too. This chapter is concluded with a short presentation of niching methods, which serve for multi-objective optimization via GAs.

3. ENCODING

In order to apply a GA to a given problem the first decision one has to make is the kind of genotype the problem needs. That means, a decision must be taken on how the parameters of the problem will be mapped into a finite string of symbols, known as genes (with constant or dynamic length), encoding a possible solution in a given problem space. The issue of selecting an appropriate representation is crucial for the search. The symbol alphabet used is often binary, though other representations have also been used, including character-based and real-valued encodings.

In the majority of GA applications, the strings use a binary alphabet and their length is constant during all the evolutionary process. Also, all the parameters decode to the same range of values and are allocated the same number of bits for the genes in the string. A problem occurs when a gene may only have a finite number of discrete valid values if a binary representation is used. If the number of values is not a power of two, then some of the binary codes are redundant, i.e., they will not correspond to any valid gene value. The most popular compromise is to map the invalid code to a valid one.

Another shortcoming of binary encoding is the so called Hamming cliffs (e.g., in the Appendix, although being neighbors in decimal representation, the Hamming distance between the binary strings for -0.6 and -0.5 is three (different bits)). It is worthwhile to mention that Gray coding, although frequently recommended as a solution to Hamming cliffs, because adjacent numbers differ by a single bit, has an analogous drawback for numbers at the opposite extremes of the decimal scale (e.g., the minimum and maximum gene values differ by only one bit, too). Binary encoding can also introduce an additional nonlinearity, thus making the combined objective function (the one in the genotypespace) more multimodal than the original one (in the phenotype space).

At the beginning of GA research, the binary representation was recommended because it was supposed to give the largest number of schemata (plural of schema), therefore providing the highest degree of implicit parallelism. However, new interpretations have shown that high-cardinality alphabets (e.g., real numbers) are more effective due to the higher expression power and low effective cardinality [12]-[14]. Complex applications need non-binary alphabets. Integer or continuous-valued genes are typically used in large-scale function optimization problems. Another advantage of non-binary representations, particularly the real-valued one, is the easy definition of problem specific operators.

When using binary coding, the positions of the genes in the chromosome is extremely important for a successful GA design, unless uniform crossover is applied (see Section 5.2). A bad choice can make the problem harder than necessary. Therefore, correlated binary genes should be coded together in order to form building blocks, thus diminishing the disruptive effects of crossover. However, this information is usually unavailable beforehand.

Epistasis is a measure of problem difficulty for GAs. It represents the interaction among different genes in a chromosome. This depends on the extent to which the change in chromosome fitness resulting from a small change in one gene varies according to the values of other genes. The higher the epistasis level, the harder the problem is. This is obviously also true when applying uniform crossover or real-valued encoding. As mentioned earlier, a possibility for making the gene ordering irrelevant is to apply uniform crossover, because the result of this operation is not affected by the positions of the genes. The same goal can be achieved with real-valued encoding and recombination operators that also turn the genes positions irrelevant (Section 5.2). However, making the gene ordering irrelevant does not necessarily mean an easier way to a good solution.

One possible answer for the binary gene position problem is to use an operator called inversion. This is implemented by extending every gene by adding the position it occupies in the string. Inversion is interesting because it

can freely mix the genes of the same string in order to put together the building blocks, automatically, during evolution (e.g., [(2 1) (3 0) (1 0) (4 1)], where the first number is a bit tag which indexes the bit and the second one represents the bit value, i.e., (3 0) means that the third bit is equal to zero). At first sight, the inversion operator looks very useful when the correlated parameters are not known a priori. With the association of a position to every gene, the string can be correctly reordered before evaluation. However, for large-scale problems, inversion is useless. Reordering greatly expands the search space, making the problem much more difficult to solve.

Therefore, the very hard encoding problem still remains in the hands of the designer. In order to achieve good performance for large tasks, GAs must be matched to the search problem at hand. The only way to succeed is by using domain-specific knowledge to select an appropriate representation.

4. FITNESS FUNCTION

Each string is evaluated and assigned a fitness value after the creation of an initial population. It is useful to distinguish between the objective function and the fitness function used by a GA. The objective function provides a measure of performance with respect to a particular set of gene values, independently of any other string. The fitness function transforms that measure of performance into an allocation of reproductive opportunities, i.e., the fitness of a string is defined with respect to other members of the current population. After decoding the chromosomes, i.e., applying the genotype to phenotype transformation, each string is assigned a fitness value. The phenotype is used as input to the fitness function. Then, the fitness values are employed to relatively ponder the strings in the population.

The specification of an appropriate fitness function is crucial for the correct operation of a GA [15]. As an optimization tool, GAs face the task of dealing with problem constraints [16]. Crossover and mutation, i.e., the perturbation (variation) mechanism of GAs, are general operators that do not take into account the feasibility region. Therefore, infeasible offspring appear quite frequently. There are four basic techniques for handling constraints when using GAs.

The simplest alternative is the rejecting technique in which infeasible chromosomes are discarded all over the generations. A different strategy is the repairing procedure, which uses a converter to transform an infeasible chromosome into a feasible one. Another possible technique is the creation of problem specific genetic operators to preserve feasibility of chromosomes.

The previous procedures do not generate infeasible solutions. This is not usually an advantage. In fact, for large-scale, highly constrained optimization problems, this is certainly a great drawback. Particularly for real world problems, where the optimal solutions usually are on the boundaries of feasible regions, the above mentioned techniques for handling constraints often lead to poor solutions. One possible way for overcoming this drawback is to apply the repairing procedure only to a fraction (10%, for instance) of the infeasible population.

It has been suggested that constraint handling for such type of optimization problem should be performed allowing search through infeasible regions. Penalty functions allow the exploration of infeasible subspaces [17]. An infeasible point close to the optimum solution generally contains much more information about it than a feasible point far from the optimum. On the other hand, the design of penalty functions is difficult and problem dependent. Usually, there is no a priori information about the distance to optimal points. Therefore, penalty methods consider only the distance from the feasible region. Penalties based on the number of violated constraints do not work well.

There are two possible forms to build a fitness function with penalty term: the addition and multiplication forms. The former is represented as $g(\underline{x}) = f(\underline{x}) + p(\underline{x})$; where for maximization problems $p(\underline{x}) = 0$ for feasible points, and $p(\underline{x}) < 0$ otherwise. The maximum absolute $p(\underline{x})$ value cannot be greater than the minimum absolute $f(\underline{x})$ value for any generation, in order to avoid negative fitness values. The multiplication form is represented as $g(\underline{x}) = f(\underline{x})p(\underline{x})$; where for maximization problems $p(\underline{x}) = 1$ for feasible points, and $0 \leq p(\underline{x}) < 1$ otherwise.

The penalty term should vary not only with respect to the degree of constraints violations, but also with respect to the GA iteration count. Therefore, besides the amount of violation, the penalty term usually contains variable penalty factors, too (one per violated constraint). The key for a successful penalty technique is the proper setting of these penalty factors. Small penalty factors can lead to infeasible solutions, while very large ones totally neglect infeasible subspaces. In average, the absolute values of the objective and penalty functions should be similar. At least in theory, the parameters of the penalty functions can, also, be encoded as GA parameters. This procedure creates an adaptive method, which is optimized as the GA evolves toward the solution.

In summary, the main problems associated with the fitness function specification are the following:

- dependence on whether the problem is related to maximization or minimization;

- when the fitness function is noisy for a non-deterministic environment [18];
- the fitness function may change dynamically as the GA is executed;
- the fitness function evaluation can be so time consuming that only approximations to fitness values can be computed;
- the fitness function should allocate very different values to strings in order to make the selection operator work easier (Section 5.1.6);
- it must consider the constraints of the problem; and
- it could incorporate different sub-objectives.

The fitness function is a black box for the GA. Internally, this may be achieved by a mathematical function, a simulator program, or a human expert that decides the quality of a string. At the beginning of the iterative search, the fitness function values for the population members are usually randomly distributed and wide spread over the problem domain. As the search evolves, particular values for each gene begin to dominate. The fitness variance decreases as the population converges. This variation in fitness range during the evolutionary process often leads to the problems of premature convergence and slow finishing.

4.1. Premature Convergence

A frequent problem with GAs, known as deception, is that the genes from a few comparatively highly fit (but not optimal) individuals may rapidly come to dominate the population, causing it to converge on a local maximum. Once the population has converged, the ability of the GA to continue searching for better solutions is nearly eliminated. Crossover (Section 5.2) of almost identical chromosomes generally produces similar offspring. Only mutation (Section 5.3), with its random perturbation mechanism, remains to explore new regions of the search space.

The schema theorem says that reproductive opportunities should be given to individuals in proportion to their relative fitnesses. However, by doing that, premature convergence occurs because the population is not infinite (basic hypothesis of the theorem). This is due to genetic drift (see Section 6). In order to make GAs work effectively on finite populations, the (proportional) way individuals are selected for reproduction must be modified. Different ways of performing selection are described in Section 5.1. The basic idea is to control the number of reproductive opportunities each individual gets. The strategy is to compress the range of fitnesses, without losing selection pressure (Section 4.2), and avoid any super-fit individual from suddenly dominating the population.

4.2. Slow Finishing

This is the opposite problem of premature convergence. After many generations, the population has almost converged, but it is still possible that the global maximum (or a high quality local one) has not been found. The average fitness is high, and the difference between the best and the average individuals is small. Therefore, there is insufficient variance in the fitness function values to localize the maxima.

The same techniques used to tackle premature convergence are used also for fighting slow finishing. An expansion of the range of population fitnesses is produced, instead of a compression. Both procedures are prone to bad re-mapping (underexpansion or overcompression) due to super-poor or super-fit individuals.

5. BASIC OPERATORS

In this section, several important design issues for the selection, crossover and mutation operators are presented. Selection implements the survival of the fittest according to some predefined fitness function. Therefore, high-fitness individuals have a better chance of reproducing, while low-fitness ones are more likely to disappear. Selection alone cannot introduce any new individuals into the population, i.e., it cannot find new points in the search space. Crossover and mutation are used to explore the solution space.

Crossover, which represents mating (recombination) of two individuals is performed by exchanging parts of their strings to form two new individuals (offspring). In its simplest form, sub-strings are exchanged after a crossover point is randomly determined. The crossover operator is applied with a certain probability, usually in the range [0.5, 1.0]. This operator allows the evolutionary process to move toward promising regions of the search space. It is likely to create even better individuals by recombining portions of good individuals. The new offspring created from mating, after being subject to mutation, are put into the next generation.

The purpose of the mutation operator is to maintain diversity within the population and inhibit premature convergence to local optima by randomly sampling new points in the search space. The GA stopping criterion may be specified as a maximal number of generations or as the achievement of an appropriate level for the generation average fitness (stagnation).

5.1. Selection

Selection, more than crossover and mutation, is the operator responsible for determining the convergence characteristics of GAs [19], [20]. Selection pressure is the degree to which the best individuals are favored [21]. The higher the

selection pressure, the more the best individuals are favored. The selection intensity of GAs is the expected change of average fitness in a population after selection is performed. Analyses of selection schemes show that the change in mean fitness at each generation is a function of the population fitness variance.

The convergence rate of a GA is largely determined by the magnitude of the selection pressure. Higher selection pressures imply in higher convergence rates. If the selection pressure is too low, the convergence rate will be slow, and the GA will unnecessarily take longer to find a high quality solution. If the selection pressure is too high, it is very probable that the GA will prematurely converge to a bad solution. In fact, selection schemes should also preserve population diversity, in addition to providing selection pressure. One possibility to achieve this goal is to maximize the product of selection intensity and population fitness standard deviation. Therefore, if two selection methods have the same selection intensity, the method giving the higher standard deviation of the selected parents is the best choice.

Many selection schemes are currently in use. They can be classified in two groups: proportionate selection and ordinal-based selection. Proportionate-based procedures select individuals based on their fitness values relative to the fitness of the other individuals in the population. Ordinal-based procedures select individuals not upon their fitness, but based on their rank within the population.

An ordinal selection scheme has a fundamental advantage over a proportional selection one. The former is translation and scale invariant, i.e., the selection pressure does not change when every individual's fitness is multiplied and added by a constant. The selection intensity of proportionate selection is the only one that is sensitive to the current population distribution [22]. However, conclusive statements about the performance of rank-based selection schemes are difficult to make because, by suitable (but tricky!) adjustment, proportionate selection can give similar performance.

5.1.1. Tournament Selection

This selection scheme is implemented by choosing some number of individuals randomly from the population and copying the best individual from this group into the intermediate population, and by repeating it until the mating pool is complete. Tournaments are frequently held only between two individuals. Bigger tournaments are also used with arbitrary group sizes (not too big in comparison with the population size). Tournament selection can be implemented very efficiently because no sorting of the population is required.

The tournament procedure selects the mating pool without re-mapping the fitnesses. By adjusting the tournament size the selection pressure can be made arbitrarily large or small. Bigger tournaments have the effect of increasing the selection pressure, since below average individuals do not have good chances of winning a competition.

5.1.2. Truncation Selection

In truncation selection, only a subset of the best individuals are chosen to be possibly selected, with the same probability. This procedure is repeated until the mating pool is complete. As a sorting of the population is required, truncation selection has a greater time complexity than tournament selection. As in tournament selection, there is no fitness re-mapping in truncation selection.

5.1.3. Linear Ranking Selection

The individuals are sorted according to their fitness values and the last position is assigned to the best individual, while the first position is allocated to the worst one. The selection probability is linearly assigned to the individuals according to their ranks. All individuals get a different selection probability, even when equal fitness values occur.

5.1.4. Exponential Ranking Selection

Exponential ranking selection differs from linear ranking selection only in that the probabilities of the ranked individuals are exponentially weighted.

5.1.5. Elitist Selection

Preservation of the elite solutions from the preceding generation assures that the best solutions known so far will remain in the population and have more opportunities to produce offspring. Elitist selection is used in combination with other selection strategies.

5.1.6. Proportional Selection

This is the first selection method proposed for GAs. The probability that an individual will be selected is simply proportionate to its fitness value. The time complexity of the method is the same as in tournament selection. This mechanism works only if all fitness values are greater than zero. The selection probabilities strongly depend on the scaling of the fitness function. In fact, most of the scaling procedures described in the next sections have been proposed to keep proportional selection working. One big drawback of proportional selection is that the selection intensity is usually low, because a single individual, either the fittest or the worst, dictates the degree of compression of the range of

fitnesses. This is quite common even during the early stage of the search, when the population variance is high. Negative selection intensity is also possible.

Notice that in ordinal-based selection schemes the effect of extreme individuals is negligible, irrespective of how much greater or smaller their fitnesses are than the rest of the population. Therefore, despite its popularity among practitioners, proportional selection (i.e., roulette wheel) is usually an inferior scheme. There are different scaling operators that help in separating the fitness values in order to improve the work of the proportional selection mechanism. The most common ones are linear scaling, sigma truncation and power scaling.

5.1.6.1. Linear scaling

Linear scaling (i.e., $f' = a.f + b$) works well except when most populations members are highly fit, but a few very poor individuals are present. The coefficients a and b are usually chosen to enforce equality of the objective and fitness functions average values, and also cause maximum scaled fitness to be a specified multiple (usually two) of the average fitness. These two conditions ensure that average population members receive one offspring copy on average, and the best receives the specified multiple number of copies. Notice that proportional selection with linear scaling is not the same as linear ranking selection.

5.1.6.2. Sigma truncation

In order to overcome the presence of super-poor individuals, the use of population variance information has been suggested to preprocess objective function values before scaling. This procedure subtracts a constant from the objective function values; $f' = \max[0, f - (\bar{f} - d.\sigma)]$, where \bar{f} is the mean objective function value in the population. The constant d is chosen as a multiple (between 1 and 3) of the population standard deviation, and negative results are arbitrarily set to zero.

5.1.6.3. Power scaling

Another possibility is power scaling, i.e., $f' = f^k$. In general, the k value is problem dependent and may require adaptation during a run to expand or compress the range of fitness function values. The problem with all fitness scaling schemes is that the degree of compression can be determined by a single extreme individual, degrading the GA performance.

5.2. Crossover

Crossover is a very controversial operator due to its disruptive nature (i.e., it can split important information). In fact, besides GAs, other evolutionary algorithms do not rely on crossover (or similar type of recombination). However, no definite answer about the necessity of using crossover has been reached so far.

The traditional GA uses one-point crossover (Fig. 1), where the two parents are each cut once at specific points, and the segments located after the cuts exchanged. The positions of the bits in the schema determine the likelihood that these bits will remain together after crossover. Obviously, an order-1 schema is not affected by recombination, since the critical bit is always inherited by one of the offspring.

$$\begin{array}{r}
 [1 \ 1 | 0 \ 1 \ 0 | 1 \ 0 \ 1] \ [1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1] \\
 \Rightarrow \\
 [1 \ 0 | 0 \ 1 \ 1 | 0 \ 0 \ 0] \ [1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]
 \end{array}$$

Fig. 1. Example of one-point crossover.

The crossover operator presented above can be generalized in order to apply multiple-point crossover. However, more than two crossover points, although giving a better exploration capacity, can be too much disruptive. The crossover mechanism can be better visualized treating strings as rings. In Fig. 2, two-point crossover is applied to the example shown in Fig. 1. Each offspring takes one ring segment, in between adjacent cut points, from each parent. The contiguous ring segment(s) is(are) taken from the other parent. For more than two crossover points, this procedure is repeated until the last segment is filled. An extra cut is assumed at the beginning of the string, i.e., between genes g8 and g1, for an odd number of cut points.

From the linear string point of view, the elements in between the two crossover points are swapped between two parents to form two offspring (Fig. 2). One-point crossover can be represented by the ring geometry as a two-point crossover with the first cut point always between genes g8 and g1. For multiple-point crossover, the cut points can be anywhere, as long as they are not the same.

$$\begin{array}{ccc}
 g1 & & 1 \\
 & / & \\
 & & 1
 \end{array}$$

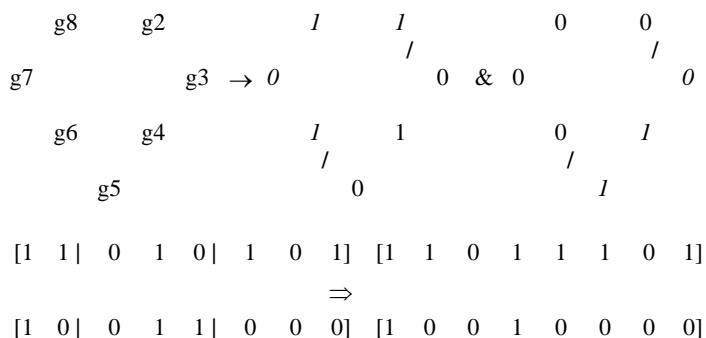


Fig. 2. Ring representation and two-point crossover.

Uniform crossover is another important recombination mechanism [23]. Offspring is created by randomly picking each bit from either of the two parent strings (Fig. 3). This means that each bit is inherited independently from any other bit. Uniform crossover has the advantage that the ordering of the genes is irrelevant in terms of splitting building blocks.

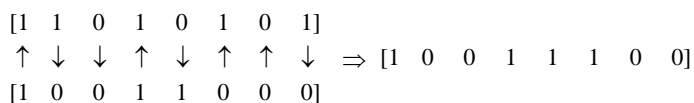


Fig. 3. Example of uniform crossover, where each arrow points to the randomly picked gene value.

Uniform crossover is more disruptive than two-point crossover. On the other hand, two-point crossover performs poorly when the population has largely converged because of the inability to promote diversity. For small populations, which is not usually the case for large-scale problems, more disruptive crossover operators such as uniform or m -point ($m > 2$) may perform better because they help overcome the limited amount of information.

Reduced surrogates can be used to improve two-point crossover exploration ability. It is highly recommended for large-scale problems. The idea is to ignore all bits that are equivalent in the two parent strings (Fig. 4). Afterwards, crossover is applied on the reduced surrogates, i.e., only one possible cut is considered between any pair of non-equivalent bits.

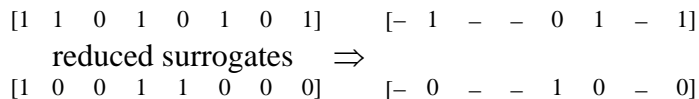


Fig. 4. Implementation of reduced surrogates to improve crossover exploration capability.

Notice that the reduced surrogate form implements the original crossover operation in an unbiased way. For example, the cut points between genes 2|3, 3|4 and 4|5 produce the same effect on offspring. Therefore, two-point reduced surrogate crossover considers these cut points as one single possible cross point.

The crossover operator can be redefined for real-valued encoding. Different combinations have been utilized (e.g., a convex combination such as $\lambda_1 x_1 + \lambda_2 x_2$). One possibility is to take the average of the two corresponding parent genes. The square-root of the product of the two values can also be used. Another possibility is to take the difference between the two values, and add it to the higher or subtract it from the lower.

5.3. Mutation

The GA literature has reflected a growing recognition of the importance of mutation in contrast with viewing it as responsible for re-introducing inadvertently lost gene values. The mutation operator is more important at the final generations when the majority of the individuals present similar quality. As is shown in Section 5.4, a variable mutation rate is very important for the search efficiency. Its setting is much more critical than that of crossover rate.

In the case of binary encoding, mutation is carried out by flipping bits at random, with some small probability (usually in the range [0.001; 0.05]). For real-valued encoding, the mutation operator can be implemented by random replacement, i.e., replace the value with a random one. Another possibility is to add/subtract (or multiply by) a random (e.g., uniformly or Gaussian distributed) amount. Mutation can also be used as a hill-climbing mechanism.

5.4. Control Parameters Estimation

Typical values for the population size, crossover and mutation rates have been selected in the intervals [30, 200], [0.5, 1.0] and [0.001, 0.05], respectively. Fixed crossover and mutation operators do not provide enough search power for tackling large-scale optimization problems. Parameter manual tuning is common practice in GA design. One parameter is tuned at a time in order to avoid the impossible task of simultaneous estimation. However, as they strongly interact in complex forms, this tuning procedure is prone to sub-optimality.

In fact, any static set of parameters is inappropriate, regardless of how they are tuned. The GA search technique is an adaptive process, which requires continuous tracking of the search dynamics. Therefore, the use of constant parameters leads to inferior performance. For example, it is obvious that large mutations can be helpful during early generations to improve the GA exploration capability. This is not the case for the end of the search, when small mutation steps are needed to fine tune sub-optimal solutions.

The proper way for dealing with this problem is by using parameters that are functions of the number of generations. Deterministic rules are frequently applied for implementing this idea. However, besides being very difficult to define, they fail to take into account the actual progress of the population performance. Adaptive rules based on population variance, or even the search for optimal parameters as part of the GA processing (i.e., including parameters as part of the chromosomes) seem to be more promising [24], [25].

6. NICHING METHODS

Two agents cause the reduction of population fitness variance at each generation. The first, selection pressure, multiplies copies of the fitter individuals. The other agent is independent of fitness. It is called genetic drift [26] and is due to the stochastic nature of the selection operator (i.e., bias on the random sampling of the population). When there is lack of selection pressure, genetic drift is responsible for premature convergence. The GA still ends up on a single peak, even when there are several ones of equal fitness.

Therefore, even when multi-objective optimization is not the main goal, the identification of multiple optima is beneficial for the GA performance. Niching methods extend standard GAs by creating stable subpopulations around global and local optimal solutions. Niching methods maintain population diversity and allow GAs to explore many peaks simultaneously. They are based on either fitness sharing or crowding schemes [27].

Fitness sharing decreases each element's fitness proportionally to the number of similar individuals in the population, i.e., in the same niche. The similarity measure is based on either the genotype (e.g., Hamming distance) or the phenotype (e.g., Euclidian distance). On the contrary, crowding schemes do not require the setting of a similarity threshold (niche radius). Crowding implicitly defines neighborhood by the application of tournament rules. It can be implemented as follows. When an offspring is created, one individual is chosen, from a random subset of the population, to disappear. The chosen one is the element which most closely resembles the new offspring.

Another idea used by niching methods is restricted mating. This mechanism avoids the recombination of individuals which do not belong to the same niche. Highly fit, but not similar, parents can produce highly unfit offspring. Restricted mating is based on the assumption that if similar parents (i.e., from the same niche) are mated, then offspring will be similar to them.

It is important to notice that similarity of genotypes does not necessarily imply similarity of the corresponding phenotypes. The hypothesis that highly fit parents generate highly fit offspring is valid only under the occurrence of building blocks and low epistasis. When the genes strongly interact, there is no guarantee that these offspring will not be lethals.

7. FINAL COMMENTS

This tutorial on GAs has pointed out the main topics on their design. The focus on the essential topics helps to not miss the forest for the trees. The first generation of GAs, based on the canonical algorithm, considering proportional selection and crossover/mutation with constant probabilities, was not originally proposed for solving static optimization problems [28]. Almost three decades of research has adapted the original proposal [29] to deal with this type of problem.

One important issue that has been avoided in this chapter is parallel GAs. They introduce new parameters such as the number of populations and their sizes, the topology of communications (e.g., each population is connected to all the others), and the migration rate. Although many implementations of parallel GAs have been described in the literature, the effect of these new parameters on the quality of the search is still under analysis [30]. Recently, another interesting idea, based on the theory of immunity in biology, has been proposed [31].

Since the first applications of GAs, they have been applied not only to pure optimization problems, but also to model identification, control, and neural network training. After a necessary period of maturing, GAs are being used now, frequently in combination with conventional optimization techniques, for solving large-scale problems.

APPENDIX

TABLE I
 HAMMING DISTANCE AND GRAY CODE

Binary	Gray	Real
[0000]	[0000]	-0.9
[0001]	[0001]	-0.8
[0010]	[0011]	-0.7
[0011]	[0010]	-0.6
[0100]	[0110]	-0.5
[0101]	[0111]	-0.4
[0110]	[0101]	-0.3
[0111]	[0100]	-0.2
[1000]	[1100]	-0.1
[1001]	[1101]	0.0
[1010]	[1111]	+0.1
[1011]	[1110]	+0.2
[1100]	[1010]	+0.3
[1101]	[1011]	+0.4
[1110]	[1001]	+0.5
[1111]	[1000]	+0.6

ACKNOWLEDGMENTS

This work was supported by the Brazilian Research Council (CNPq) under grant No. 300054/91-2. Alves da Silva would also like to thank PRONEX for the financial support and Professor Djalma M. Falcão, from the Federal University of Rio de Janeiro, for his time and effort reviewing this chapter.

REFERENCES

- [1] D.B. Fogel: "An introduction to simulated evolutionary optimization", IEEE Trans. on Neural Networks, Vol. 5, No. 1, 1994, pp. 3-14.
- [2] D.B. Fogel: *Evolutionary Computation - Toward a New Philosophy of Machine Intelligence*, IEEE Press, 1995.
- [3] T. Bäck, U. Hammel, and H.-P. Schwefel: "Evolutionary computation: comments on the history and current state", IEEE Trans. on Evolutionary Computation, Vol. 1, No. 1, 1997, pp. 3-17.
- [4] D.E. Goldberg: *Genetic Algorithms in Search, Optimization and Machine Learning*, Addison-Wesley, 1989.
- [5] D. Beasley, D.R. Bull, and R.R. Martin: "An overview of genetic algorithms: Part 1, fundamentals", University Computing, Vol. 15, No. 2, 1993, pp. 58-69.
- [6] D. Beasley, D.R. Bull, and R.R. Martin: "An overview of genetic algorithms: Part 2, research topics", University Computing, Vol. 15, No. 4, 1993, pp. 170-181.
- [7] D. Whitley: "A genetic algorithm tutorial", Statistics and Computing, Vol. 4, 1994, pp. 65-85.
- [8] E. Aarts and J. Korst: *Simulated Annealing and Boltzmann*, John Wiley, 1989.
- [9] F. Glover and M. Laguna: *Tabu Search*, Kluwer Academic, 1997.
- [10] J. Kennedy and R.C. Eberhart: *Swarm Intelligence*, Morgan Kaufmann, 2001.
- [11] G. Rudolph: "Convergence analysis of canonical genetic algorithms", IEEE Trans. on Neural Networks, Vol. 5, No. 1, 1994, pp. 96-101.
- [12] H.J. Antonisse: "A new interpretation of schema notation that overturns the binary encoding constraint", Proc. 3rd Int. Conf. on Genetic Algorithms, Morgan Kaufmann, 1989, pp. 86-91.
- [13] D.E. Goldberg: "The theory of virtual alphabets", in Parallel Problem Solving from Nature 1, Lecture Notes in Computer Science, Vol. 496, Springer, 1991, pp. 13-22.
- [14] L.J. Eshelman and J.D. Schaffer: "Real-coded genetic algorithms and interval-schemata", in Foundations of Genetic Algorithms 2, Morgan Kaufmann, 1993, pp. 187-202.

- [15] N.J. Radcliffe and P. D. Surry: “Fitness variance of formae and performance prediction”, in Foundations of Genetic Algorithms 3, Morgan Kaufmann, 1995, pp. 51-72.
- [16] Z. Michalewicz and M. Schoenauer: “Evolutionary algorithms for constrained parameter optimization problems”, Evolutionary Computation, Vol. 4, No. 1, 1996, pp. 1-32.
- [17] M. Gen and R. Cheng: “A survey of penalty techniques in genetic algorithms”, Proc. 3rd IEEE Conf. on Evolutionary Computation, IEEE Press, 1996, pp. 804-809.
- [18] B.L. Miller and D.E. Goldberg: “Genetic algorithms, selection schemes, and the varying effects of noise”, University of Illinois at Urbana-Champaign, IlliGAL Report, No. 95009, 1995.
- [19] D.E. Goldberg and K. Deb: “A comparative analysis of selection schemes used in genetic algorithms”, in Foundations of Genetic Algorithms, Morgan Kaufmann, 1991, pp. 69-93.
- [20] T. Blickle and L. Thiele: “A comparison of selection schemes used in genetic algorithms”, Swiss Federal Institute of Technology, TIK-Report, Nr.11, 2nd Version, 1995.
- [21] T. Bäck: “Selective pressure in evolutionary algorithms: A characterization of selection mechanisms”, Proc. 1st IEEE Conf. on Evolutionary Computation, IEEE Press, 1994, pp. 57-62.
- [22] L.D. Whitley: “The GENITOR algorithm and selection pressure: Why rank-based allocation of reproductive trials is best”, Proc. 3rd Int. Conf. on Genetic Algorithms, Morgan Kaufmann, 1989, pp. 116-121.
- [23] G. Syswerda: “Uniform crossover in genetic algorithms”, Proc. 3rd Int. Conf. on Genetic Algorithms, Morgan Kaufmann, 1989, pp. 2-9.
- [24] N. Saravanan, D.B. Fogel, and K.M. Nelson: “A comparison of methods for self-adaptation in evolutionary algorithms”, BioSystems, Vol. 36, 1995, pp. 157-166.
- [25] A.E. Eiben, R. Hinterding, and Z. Michalewicz: “Parameter control in evolutionary algorithms”, IEEE Trans. on Evolutionary Computation, Vol. 3, No. 2, 1999, pp. 124-141.
- [26] A. Rogers and A. Prügel-Bennet: “Genetic drift in genetic algorithm selection schemes”, IEEE Trans. on Evolutionary Computation, Vol. 3, No. 4, 1999, pp. 298-303.
- [27] B. Sareni and L. Krähenbühl: “Fitness sharing and niching methods revisited”, IEEE Trans. on Evolutionary Computation, Vol. 2, No. 3, 1998, pp. 97-106.
- [28] K.A. De Jong: “Genetic algorithms are NOT function optimizers,” in Foundations of Genetic Algorithms 2, Morgan Kaufmann, 1993, pp. 5-17.
- [29] J.H. Holland: *Adaptation in Natural and Artificial Systems*, University of Michigan Press, 1975.
- [30] E. Cantu-Paz: “Markov chain models of parallel genetic algorithms”, IEEE Trans. on Evolutionary Computation, Vol. 4, No. 3, 2000, pp. 216-226.
- [31] L. Jiao and L. Wang: “A novel genetic algorithm based on immunity”, IEEE Trans. on Systems, Man, and Cybernetics – Part A, Vol. 30, No. 5, 2000, pp. 552-561.