

C6.2 Speciation methods

Kalyanmoy Deb (C6.2.1–C6.2.4) and *William M Spears* (C6.2.5, C6.2.6)

Abstract

In nature, a species is defined as a collection of phenotypically similar individuals. Many biologists believe that individuals in a sexually reproductive species can be created and maintained by allowing restrictive mating only among individuals from the same species. The connection between the formation of multiple species in nature and in search and optimization problems lies in solving multimodal problems, where the objective is not only to find one optimal solution, but to find a number of optimal solutions. In those problems, each optimal solution may be assumed to constitute a species. Since evolutionary algorithms work with a population of solutions, the concept of natural speciation techniques can be implemented to allow formation of multiple subpopulations, each focusing its search for one optimal solution. This way, multiple optimal solutions can be discovered simultaneously. In this section, a number of speciation techniques are discussed.

C6.2.1 Introduction

Kalyanmoy Deb

Despite some controversy, most biologists agree that a species is a collection of individuals which resemble each other more closely than they resemble individuals of another species (Eldredge 1989). It is also clear that the reproductive process of the sexually reproducing organisms causes individuals to resemble their parents, thereby maintaining a phenotypic similarity among individuals of the community or the *species*. Thus, there is a strong correlation among the reproductively coherent individuals and a phenotypically similar cluster of individuals. Since in evolutionary algorithms a population of solutions is used, artificial species of phenotypically similar solutions can be formed and maintained in the population by restricting their mating to that with similar individuals. Before we outline how to form and maintain multiple species in a population, let us discuss why it could be necessary to form species in the applications of evolutionary algorithms.

In Section C6.1, we saw that multiple optimal solutions in a multimodal optimization problem can be found simultaneously by forming artificial niches (subpopulations) in the population. Each niche can be considered to represent a peak (in the spirit of maximization problems). To capture a number of peaks simultaneously and maintain them for many generations, a niching method is used. Niching helps to emphasize and maintain solutions around multiple optima. However, in niching, the main emphasis is devoted to distributing the population members across different peaks. Thus, the niching technique cannot quite focus its search on each peak and find the exact optimal solutions efficiently. This is because some of the search effort is wasted in the recombination of interpeak solutions, which, in turn, may produce some *lethal* solutions representing none of the peaks. A speciation method used in evolutionary computation (EC) studies, on the other hand, restricts mating to that among like solutions (likeness can be defined phenotypically or genotypically) and discourages mating among solutions of different peaks. If the likeness is defined properly, two parent solutions chosen for mating are likely to represent the same peak. Thus, when like individuals mate with each other, the created children solutions are also similar to the

parent solutions and are likely to be members of the same peak. This way, the restriction of mating to that among like solutions may reduce the creation of lethal solutions (which represent none of the peaks). This may allow the search to concentrate on each peak and help find the best or near-best optimum solution efficiently. However, in order to apply the speciation technique properly, solutions representing each peak must first be found. Thus, the speciation technique cannot be used independently. In the presence of both niching and speciation, niching finds and maintains subpopulation of solutions around multiple optima and the speciation technique allows us to make an inherent parallel search in each optimum to find multiple optimal solutions simultaneously.

Among the evolutionary algorithms, a number of speciation methods have been suggested and implemented in *genetic algorithms* (GAs). Of the earlier works related to mating restriction in GAs, B1.2 Hollstien's (1971) inbreeding scheme where mating was allowed between similar individuals in his simulation of animal husbandry problems, Booker's (1982) taxon–exemplar scheme for restrictive mating in his simulation of learning pattern classes, Holland's suggestion of a tag–template scheme (Goldberg 1989), Sannier and Goodman's (1987) restrictive mating in forming separate coherent groups in a population, Deb's (1989) phenotypic and genotypic mating restriction schemes, and Spears's (1994) and Perry's (1984) speciation using tag bits are a few studies. In the following, we discuss some of the above speciation methods in more detail.

C6.2.2 Booker's taxon–exemplar scheme

Kalyanmoy Deb

Booker (1982) used taxons and exemplars in his *learning algorithm* to reduce the formation of lethal B1.5.2 individuals. He defined a taxon as a string (constructed over the three-letter alphabet {0, 1, #}, with a # matching a 0 or a 1). The population is initialized with taxon strings. In his *restricted mating policy*, he wanted to restrict mating among similar taxon strings, which were identified by calculating a match score of the taxon strings with a given exemplar binary string. He allowed partial match scores depending on the matching of the taxon and the exemplar. For the following two taxon strings and the exemplar string, the first taxon matches the exemplar completely. The second taxon matches the exemplar partially (in first, third, and fourth positions):

Taxon	Exemplar
(1 # 0 0 #)	(1 0 0 0 0).
(# 1 # 0 1)	

If the taxon completely matches the exemplar, a score is assigned as the sum of the string length and the number of #s in the taxon. The partial credit is also assigned based on the number of correct matches and the number of #s in the taxon. In order to implement the restrictive mating policy, he chose parent taxon strings from a sample subpopulation determined based on the available matching taxon strings in the population. If a specified number of matching taxon strings are available in the population, parent strings are chosen uniformly at random from all the matching taxon strings. Otherwise, parent strings are chosen according to a probability distribution calculated based on the match score of the taxon strings. In a number of pattern discovery problems an improved performance is observed with the restricted mating policy.

After the patterns were discovered, Booker extended his above scheme to classify the discovered patterns using a modified string as follows:

Taxon	Tag
(1 # 0 0 #)	: (1 0 0 0 0).

In addition to the taxon string, a tag string is introduced to classify the discovered taxon strings (or patterns). The taxon strings matching a particular tag string were considered to be in the same class. A similar match score was used, except that this time the matching was performed with the taxon and tag strings. As discussed elsewhere (Goldberg 1989), there is one difficulty with the above tag–taxon scheme. The tag string must be of the same length as the taxon string. This increases the complexity of the classification problem, whereas the same concept can be implemented with shorter tag and template strings, as suggested by Holland; a brief description of this is given by Goldberg (1989).

C6.2.3 The tag–template method

Kalyanmoy Deb

In addition to the functional string (the taxon string in Booker’s pattern classification problem), a template and a tag string are introduced. The template string is constructed from the three-letter alphabet (1, 0, and #) as before, but the tag string is a binary string of the same length as the template string. A typical string with the tag and template strings would look like the following:

Template	Tag	Functional string
(#01)	: (100)	: (1011001101).

The size of tag and template strings depends on the number of desired solutions. A simple calculation shows that if q different optimal solutions (peaks) are to be found, the minimum string length for the tag and template is $\lceil \log_2 q \rceil$ (Deb 1989). The tag and template strings are created at random in the initial population along with the functional string. These two strings do not affect the fitness of the functional string. However, they are affected by the crossover and the mutation operators, as well. For the template string, the mutation operator must be modified to operate on a three-allele string. The purpose of these strings is to restrict mating. Before crossing a pair of individual strings, their tag and template strings are *matched*. If the match score exceeds a threshold value, the crossover is performed between the two strings as usual; otherwise some other string pair is tested for a possible mating. In this process, the tag and template strings corresponding to the good individuals in early populations are emphasized and an artificial tag is set for solutions in each peak. Later on, since crossing over is only performed between the matched strings, only similar strings (or strings from the same peak) tend to participate in crossover.

Although neither Holland nor Goldberg simulated this speciation method, Deb (1989) (with assistance from David Goldberg) implemented this scheme and applied this technique in solving multimodal test problems. In both cases, GAs with the tag–template scheme performed better than GAs without it.

C6.2.4 Phenotypic and genotypic mating restriction

Kalyanmoy Deb

Deb (1989) has developed two mating restriction schemes based on the phenotypic and genotypic distance between mating individuals. The mating restriction schemes are straightforward. In order to choose a mate for an individual, their distance (in phenotypic mating restriction the Euclidean distance and in genotypic mating restriction the Hamming distance) is computed. If the distance is closer than a parameter σ_{mating} , they participate in the crossover operation; otherwise another individual is chosen at random and their distance is computed. This process is continued until a suitable mate is found or all population members are exhausted, in which case a random individual is chosen as a mate. Deb has implemented both the above mating restriction schemes with a single-point crossover and applied them to solve a number of multimodal test problems. Although, in all his simulations, the parameter σ_{mating} was kept the same as the parameter σ_{share} used in the niching methods, other values of σ_{mating} may also be chosen. It is worthwhile to mention that niching with the σ_{share} parameter is implemented in the selection operator and the mating restriction with the σ_{mating} parameter is implemented in the crossover operator. GAs with niching and mating restriction were found to better distribute the population across the peaks than GAs with sharing alone. Here, we present simulation results for the phenotypic mating restriction scheme adopted in that study. In solving the single-variable, five-peaked function in the interval $0 \leq x \leq 1$

$$\text{maximize} \quad 2^{-2((x-0.1)/0.8)^2} \sin^6(5\pi x)$$

with $\sigma_{\text{share}} = \sigma_{\text{mating}} = 0.1$, 100 population members after 200 generations without and with phenotypic mating restriction are shown in figure C6.2.1. Stochastic remainder roulette wheel selection and single-point crossover operators are used. The crossover and mutation probabilities are kept as 0.9 and 0.0, respectively. The figures show that, with the mating restriction scheme (the right-hand panel), the number of lethal (nonpeak) individuals has been significantly decreased. This study also implemented a genotypic mating restriction scheme and similar results were obtained. Some guidelines in choosing the sharing and mating restriction parameters are outlined elsewhere (Deb 1989, Deb and Goldberg 1989).

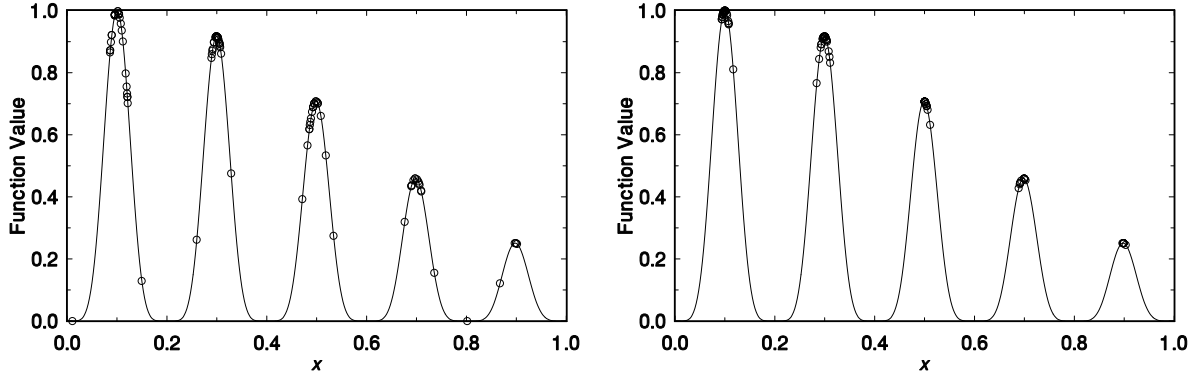


Figure C6.2.1. The distribution of 100 solutions without (left) and with (right) a mating restriction scheme.

C6.2.5 Speciation using tag bits

William M Spears

Another method for identifying species is via the use of tag bits, which are appended to every individual. Each species corresponds to a particular setting of these bits. Suppose there are k different sets of tag bit values at a particular generation of the evolutionary algorithm (EA). Denote these sets as $\{S_0, \dots, S_{k-1}\}$. The sets are numbered arbitrarily. Each individual belongs to one S_i and all individuals in a particular S_i have the same tag bit values. For example, suppose there is only one tag bit and that some individuals exist with a tag bit value zero and that the remainder exist with tag bit value one. Then (arbitrarily) assign the former set of individuals to S_0 and the latter set to S_1 . Let $\| \cdot \|$ denote the cardinality of the sets.

Spears (1994) uses the tag bits to restrict mating and to perform fitness sharing. With sharing, the perceived fitness, F_i , is a normalization of the objective fitness f_i :

$$F_i = \frac{f_i}{\|S_j\|} \quad i \in S_j$$

where $\|S_j\|$ is the size of the species that individual i is in.

The average fitness of the population, \bar{F} , becomes

$$\bar{F} = \frac{\sum_{i \in S_0} (f_i / \|S_0\|) + \dots + \sum_{i \in S_{k-1}} (f_i / \|S_{k-1}\|)}{\|S_0\| + \dots + \|S_{k-1}\|}$$

which is just

$$\bar{F} = \frac{\sum_{i \in S_0} (f_i / \|S_0\|) + \dots + \sum_{i \in S_{k-1}} (f_i / \|S_{k-1}\|)}{N}$$

since the species sizes have to total N (recall that no individual can lie in more than one species). The expected number of offspring for an individual is now F_i / \bar{F} .

Restricted mating is performed by only allowing recombination to occur between individuals with the same tag bit values. Mutation can flip all bits, including the tag bits, thus allowing individuals to change labels. Experimental results, as well as some modifications to the above mechanism can be found in the article by Spears (1994). Code for the algorithm can be found at <http://www.aic.nrl.navy.mil/~spears>.

Perry's thesis work (Perry 1984) with speciation is extremely similar to the above technique. Perry includes both species and environmental regions in an EA. Species are identified via tag bits and an environmental region is similar to an EA population. Recombination within an environment can occur only on individuals with the same tag bit values. Mutation is allowed to change tag bits, in order to introduce new species. The additional use of a 'migration' operator, which moves individuals from one environment to another, does not have an analog in the work of Spears (1994).

Perry gives an example of two species in an environment—fitness proportional selection is performed, and the average fitness of an environment is

$$\bar{f} = \frac{\sum_{i \in S_0} f_i + \sum_{i \in S_1} f_i}{\|S_0\| + \|S_1\|}$$

or

$$\bar{f} = \frac{\sum_{i \in S_0} f_i + \sum_{i \in S_1} f_i}{N}$$

where N is the population size of the environmental niche. The expected number of offspring is f_i / \bar{f} . One can see that the main difference between the two methods is the use of sharing in the computation of fitness in the work of Spears (1994). Thus it is not surprising that in many of Perry's experimental runs one particular species would eventually dominate an environmental niche (however, it should be noted that in the work of Perry (1984) the domination of an environment by a species was not undesirable behavior).

The use of tag bits makes restricted mating and fitness sharing more efficient because distance comparisons do not have to be computed. Interestingly, it is also possible to make Goldberg's implementation of sharing more efficient by sampling (Goldberg *et al* 1992). In other words the distance of each individual from the rest is estimated by using a subset of the remaining individuals.

C6.2.6 Relationship with parallel algorithms

William M Spears

Clearly this work has similarities to the EA research performed on *parallel architectures*. In a parallel EA, a topology is imposed on the EA population, resulting in species. However, there are some important differences between the parallel approaches and the sequential approach. For example, with the fitness sharing approaches the fitness of an individual and the species size are dynamic, based on the other individuals (and species). This concentrates effort on more promising peaks, while still maintaining individuals in other areas of the search space. This is typically not true for parallel EAs implemented on MIMD or SIMD architectures. When using a MIMD architecture, species are dedicated to particular processors and the species remain a constant size. In SIMD implementations, one or two individuals reside on a processor, and species are formed by defining overlapping neighborhoods. However, due to the overlap, one particular species will eventually take over the whole population. C6.3.1, C6.4.2

References

- Booker L B 1982 *Intelligent Behavior as an Adaptation to the Task Environment* Doctoral Dissertation, University of Michigan; *Dissertation Abstracts Int.* **43** 469B
- Deb K 1989 *Genetic Algorithms in Multimodal Function Optimization* Master's Thesis, University of Alabama; TCGA Report 89002
- Deb K and Goldberg D E 1989 An investigation of niche and species formation in genetic function optimization *Proc. 3rd Int. Conf. on Genetic Algorithms (Fairfax, VA, 1989)* ed J D Schaffer (San Mateo, CA: Morgan Kaufmann) pp 42–50
- Eldredge N 1989 *Macro-evolutionary Dynamics: Species, Niches and Adaptive Peaks* (New York: McGraw-Hill)
- Goldberg D E 1989 *Genetic Algorithms in Search, Optimization, and Machine Learning* (Reading, MA: Addison-Wesley)
- Goldberg D E and Richardson J 1987 Genetic algorithms with sharing for multimodal function optimization *Proc. 2nd Int. Conf. on Genetic Algorithms (Cambridge, MA, 1987)* ed J J Grefenstette (Hillsdale, NJ: Erlbaum) pp 41–9
- Goldberg D E, Deb K and Horn J 1992 Massive multimodality, deception, and genetic algorithms *Proc. Parallel Problem Solving from Nature Conf.* (Amsterdam: North-Holland) pp 37–46
- Hollstien R B 1971 *Artificial Genetic Adaptation in Computer Control Systems* Doctoral Dissertation, University of Michigan; *Dissertation Abstracts Int.* **32** 1510B
- Perry Z A 1984 *Experimental Study of Speciation in Ecological Niche Theory using Genetic Algorithms* Doctoral Dissertation, University of Michigan; *Dissertation Abstracts Int.* **45** 3870B
- Sannier A V and Goodman E D 1987 Genetic learning procedures in distributed environments *Proc. 2nd Int. Conf. on Genetic Algorithms (Cambridge, MA, 1987)* ed J J Grefenstette (Hillsdale, NJ: Erlbaum) pp 162–9
- Spears W M 1994 Simple subpopulation schemes *Proc. Conf. on Evolutionary Programming* (Singapore: World Scientific) pp 296–307