

# Chapter 1

## Memetic Algorithms in Bioinformatics

Regina Berretta, Carlos Cotta, and Pablo Moscato

### 1.1 Introduction

Bioinformatics is an exciting research field for memetic algorithms (MAs). Its core activity is the integration of techniques from Computer Science, Mathematics and Statistics to address challenging computational problems related with the analysis of large volumes of data. Due to its huge relevance as a means to understand biology in the 21st Century, this field has attracted the attention of many pioneers in MAs, including the authors of this chapter.

During the past two decades, the field of molecular biology and the new high-throughput technologies associated with it has spawned a number of interesting problems. These problems can, in many cases, be posed as optimization problems which are combinatorial, non-linear, and often have aspects of both. Some examples arise in the analysis of large scale genetic datasets (e.g., gene expression using microarrays, massive datasets of single nucleotide polymorphisms derived from genome-wide association studies, etc.).

The field of bioinformatics is characterized by a constant evolution in computational methods for clustering and feature selection, analysis of phylogenetic trees (inference and reconstruction), image processing, protein analysis (structure predic-

---

Regina Berretta

Centre for Bioinformatics, Biomarker Discovery and Information-Based Medicine, School of Electrical Engineering and Computer Science, The University of Newcastle, University Drive, Callaghan, NSW, 2308, Australia. e-mail: [Regina.Berretta@newcastle.edu.au](mailto:Regina.Berretta@newcastle.edu.au)

Carlos Cotta

Escuela Técnica Superior de Ingeniería Informática, Universidad de Málaga, Campus de Teatinos, 29071 - Málaga, Spain e-mail: [ccottap@lcc.uma.es](mailto:ccottap@lcc.uma.es)

Pablo Moscato

Centre for Bioinformatics, Biomarker Discovery and Information-Based Medicine and Hunter Medical Research Institute, School of Electrical Engineering and Computer Science, The University of Newcastle, University Drive, Callaghan, NSW, 2308, Australia. e-mail: [Pablo.Moscato@newcastle.edu.au](mailto:Pablo.Moscato@newcastle.edu.au)

tion, sequence alignment), drug therapy design, among many others others. As we said before, many aspects of these problems are combinatorial in nature, involving the selection or the arrangement of discrete objects. Many of these combinatorial problems are NP-optimization problems, thus biologists are generally interested in finding the optimal solution of a given problem, but if that is impossible to obtain, they also rely for their investigations in high-quality solutions, provided by some metaheuristic technique. In this sense, MAs are a good strategy as they can provide solutions pretty fast, but then if they are coupled to an exact solver (thus forming a complete MA – check chapter ??), they can also prove the optimality of the final solution.

In general, researchers employ exact methods developed by themselves, and highly crafted for the problem at hand, or rely in Integer Programming reformulations of their problems. References in Mathematical Programming, Integer Programming for problems in computational biology can be found in works by Lancia [54] and Althaus *et al.* [2]. A hands-on approach to modeling using commercial packages can be found in [38] and [31]. Our experience with students, coming from different academic backgrounds, also suggest that the book by Williams [95], and the reviews of Greenberg, Hart and Lancia [37] and Festa [29], are not only useful but they have the added value of being very motivational for those interested in crossing fields and to jump into this new area. However, it is clear that since the size of the datasets associated to these challenges problems is in general is massive, in many cases it is necessary to develop efficient metaheuristics to deal with the large instances of these problems. As usual, research on metaheuristics is important as it can provide good upper bounding schemes to guide exact search procedures.

This chapter provides an review of MAs that have been developed to address some of the problems mentioned above. For an eagle's view of the contents, in Table 1.1 the reader can find a list of references grouped by application. For the sake of completeness we have also included in this table some applications in the wider area of biomedicine, where applications of memetic algorithms are also manifold. In particular, it is worth mentioning the deployment of MAs for optimizing cancer treatment, both in radiotherapy [40, 41] and chemotherapy [56, 57, 89]. Precisely related to this later issue of drug scheduling we can cite the work of Neri *et al.* for HIV multidrug therapy [65]. Imaging applications in tomography and imaging are also numerous [7, 10, 24, 25, 80] (please check [71] for a review of metaheuristic methods applied to microwave imaging). In the following sections we will focus on the purely bioinformatic tasks defined in the table though.

## 1.2 Microarray Data Analysis

With the introduction of DNA microarray technologies, it is now possible measure the expression of thousands of genes simultaneously. However, this obviously comes at a price as even a single microarray experiment leads to the need to deal with large datasets. This has posed a challenge primarily for statistics, as researchers

Table 1.1: An overview of MA applications in Bioinformatics

Area	Subarea	Reference
Microarray analysis	clustering	[46, 62, 70, 83, 84]
	gene ordering	[19, 59, 64]
	feature selection	[39, 45, 100, 103, 104, 105]
Phylogenetics	inference and reconstruction	[16, 17, 18, 33, 75, 96]
	consensus tree	[72]
Protein analysis	structure prediction	[3, 11, 13, 51, 52, 69, 81, 102]
	structure comparison	[8, 49]
Molecular design	ligand docking	[43, 63]
	PCR product primer design	[98]
Sequence analysis	DNA sequencing	[27]
	multiple sequence alignment	[87]
	supersequence problem	[14, 32]
Systems biology	cell models	[77]
	gene regulatory network	[47, 48, 68, 85, 86, 88]
Biomedicine	3D reconstruction of forensic objects	[80]
	Radiotherapy	[40, 41]
	Drug therapy design	[56, 57, 65, 89]
	Tomography	[7, 10, 24, 25]

now need to deal with the “large  $n$ , small  $m$ ” problem (where  $n$  denotes the number of measurements on a single sample and  $m$  is the total number of samples). Statisticians obviously prefer to deal with the reverse situation, with more samples than measurements. When multi-variate methods are required, researchers resort to obtaining “molecular signatures”, searching for a more coherent, reliable and robust set of molecular changes [66]. They count on Computer Science (allied of course with statistical methods) for the development of sophisticated algorithms to analyze such data.

The approaches for the analysis of microarray datasets can be primary classified as *unsupervised* and *supervised methods*. At this description level, we can understand that these microarray datasets are basically two-dimensional arrays of values (the measurements) and that a re-assignment of labels to the samples (and, analogously, to the measurements) helps to uncover some structure within the data.

Clustering algorithms are the most common example of unsupervised methods to find these structures. Another unsupervised method, which can be seen as a particular type of clustering algorithm is called *gene ordering*. In this case the overall objective is to find a permutation of either the rows or columns of this two-dimensional array such that those having the same patterns of global expression are relatively close in the permutation. An example of supervised method is feature selection, in which the aim is selecting a subset of features (genes in this case) such that a main goal is optimized, for example, classification accuracy.

We now give a brief description of some MAs that have been proposed to address the clustering and feature selection problems in microarrays.

### 1.2.1 Clustering

From the description we have given before, it is clear that clustering encompasses a wide number of different problems, as the word “scheduling” in Production Planning and Operations Research encompasses different specific problems. Merz and Zell’s proposal [62] for the clustering problem in microarray data analysis is based on a model in which the task is to define an assignment of objects into clusters, such that the sum of squared distances to the centroid of the cluster is minimized. They proposed a MA which uses the K-Means algorithm as a local search technique. They use uniform crossover and they also propose a new one denominated replacement recombination operator. They compare the MA with a multi-start  $k$ -means local search using five different microarray datasets.

Speer *et al.* used in [83, 84] a Minimum Spanning Tree (MST) to represent the data, where each node is a gene and each edge between nodes  $i$  and  $j$  represent the dissimilarity between genes  $i$  and  $j$ , thus modeling the clustering problem as tree partitioning problem, i.e., deleting a set of edges to find the clusters. They proposed a MA based on the framework presented by Merz and Zell in [62]. They use two fitness functions, the sum-of-squared-error criteria (the same used in [62]) and the Davies-Bouldin-Index [21], which minimizes the intra-cluster and maximizes the inter-cluster distances. Using four microarray datasets, they compared the MA with two other popular clustering algorithms, the average linkage algorithm [28] and the Best2Partition [99], which is also based on a MST-representation of the data.

Palacios *et al.* [70] present the results of different population based metaheuristics (genetic algorithms, MAs and estimation of distribution algorithms) to obtain biclusters from microarray datasets. According to the authors, the advantage of finding biclusters in microarray datasets (instead of traditional clusters) stems from the ability to find a group of genes that are similar in a specific subset of samples. To analyze the performance of each algorithm, they used a yeast expression dataset comprising 17 samples on 2,900 probes.

Gene Ordering is another unsupervised method that can be interpreted as a special type of clustering algorithm. The objective is, given a gene expression dataset, to rearrange the genes, such that genes with similar expression patterns stay close to each other. MAs to tackle this problem have been proposed in [19, 59, 64]. In [19], Cotta *et al.* represent a solution as a binary tree, using hierarchical clustering as a start point. The crossover operator is similar to the one used in [17], using subtrees from the parents to create a offspring. Flipping subtrees are used as the model for the mutation operator. Two local searches are applied, the first one work inverting branches of subtrees and the second one employs a pairwise interchange local search. They test the MA in instances with up to 500 genes. Mendes *et al.* [59] uses the same MA, but with the objective to evaluate the impact of parallel processing in

the performance of the MA and ability to apply it in larger instances (up to 1,000 genes). More recently, in [64] these MAs are improved significantly, with the inclusion of new local searches which employ Tabu Search. The MA is tested not only in microarray instances (containing more than 6,000 genes), but as well in images, where the objective is unscramble the rows of an image when the image has all its rows permuted at random. The images are excellent as benchmark instances and help to evaluate gene ordering and different clustering algorithms, making easier to understand the quality of the results. The MA proposed by Moscato *et al.* [64] has been successfully applied in different microarray studies [4, 20, 36, 44, 60, 76].

### 1.2.2 Feature Selection

Feature selection methods are used primarily in bioinformatics to reduce the dimensionality of a dataset to help to discriminate between classes of samples under study. We note that the definition of a feature is rather general, it can be gene expression (as in microarray datasets), a single nucleotide polymorphism (SNP) (as in genome-wide association studies), protein abundances (as in ELISA kit panels), among many others sources of biological information. Feature Selection methods can be classified as *filter* or *wrapper methods*. In filter methods, the features selected are evaluated based only on the characteristic of the data and in the wrapper methods, a classification algorithm is embedded in the method, giving constant feedback regarding the quality of the set of features selected.

Zhu *et al.* [104] present a MA for feature selection problem with the objective to improve classification performance. Each individual in the population is composed of a set of selected features ( $X$ ) and a set of excluded features ( $Y$ ). The local search procedure move features between sets  $X$  and  $Y$  based on some filter ranking methods, such as ReliefF, Gain Ratio and Chi-Square. They evaluated the performance of their approach using four UCI datasets (UC Irvine Machine Learning Repository<sup>1</sup>) and four microarray datasets, showing improvements in the classification accuracy.

In [100], Zhu and Ong present a similar MA, but now using a Markov blanket approach in the local search procedure. In [103], the same authors present a comparison study between the MAs presented in [104] and [100]. They evaluated the results on synthetic and real microarray datasets. Both MAs perform well in regards to classification accuracy, but the one that uses Markov blanket approach gives smaller feature sets. Finally, in [105], they present a memetic framework that combines the previous approaches with a hybridization of wrapper and filter feature selections methods. The computational tests were done in fourteen microarray data sets containing 1,000 to 24,481 genes. They have also tested their methods for hyperspectral imagery classification. The classification accuracy was good and the number of features selected varies depending on the local search used.

---

<sup>1</sup> <http://archive.ics.uci.edu/ml/>

Other MAs for feature selection problems were proposed in [39, 45]. However, as stated by Zhu *et al.* [105], due to the inefficient local search methods used a large amount of redundant computation is incurred on evaluating the fitness of feature subsets. This is an issue worth considering in detail when designing an MA as we rely on the power of local search, associated with good data structures, to speed-up the process. This is an area of great interest and we hope more sophisticated MAs will be developed during this decade.

### 1.3 Phylogenetics

The aim of phylogenetics is to study the evolutionary relationship between species, which can be represented by a phylogenetic tree. The inference of phylogenetic trees, known as Phylogeny Problem, is a very challenging task and is certainly important in molecular biology. It has connections with other problem domains in bioinformatics like *multiple sequence alignment*, *protein structure prediction*, among others [16]. The aim of the Phylogeny Problem is to find the tree (or in certain cases the network), that best represents the evolutionary history of a set of species. Several criteria have been defined in order to measure the quality of a certain tree given certain input data (typically, molecular data corresponding to a collection of different organisms or taxa); these can be broadly grouped into sequence-based methods (such as maximum parsimony and maximum likelihood) and distance-based methods (e.g., minimal ultrametric trees). Unfortunately, *NP*-hardness has been shown for phylogenetic inference under most of these models [22, 23, 30, 97]). Due to the complexity of the problem, the research focus in the development of powerful metaheuristics, like MAs [16, 17, 18, 33, 75, 96].

Cotta and Moscato proposed several MAs for hierarchical clustering from distance matrices under a minimum-weight ultrametric tree model (i.e., finding an ultrametric tree of minimal overall weight, such that its associated distance matrix bounds the observed distances from above). The first approaches [17] were based on the use of evolutionary algorithms endowed with heuristic decoders, which could be viewed as greedy hill-climbers for genotype-to-phenotype mapping. Although these provided much better results than other simpler decoder-based approaches and tree-based evolutionary algorithms, their computational cost was also large. Later [18] an orthodox memetic approach was presented based on the use of a tree representation and a local search operator based on tree rotations.

A scatter search method using path relinking was subsequently presented by Cotta [16]. Scatter Search (SS) [34, 35, 53] is a powerful metaheuristic which can be considered as a particular type of MA that often relies more on deterministic strategies rather than in randomization. In this work, the author used an ultrametric model and a minimum weight criterion as in previous works [17, 18]. The SS algorithm was evaluated using five real biological data sets from an online repository

–the TreeBase site<sup>2</sup>– and was shown to compare favorably to an evolutionary algorithm and a MA. Related to this, Gallardo *et al.* [33] propose an hybrid algorithm that combines Branch and Bound (BnB) and MA in an interleaved way. The idea is to have both techniques sharing information between them. They used the same five biological data sets from as [16] and showed improved results.

Williams and Smith [96] use maximum parsimony as the optimization criteria, which means that the tree with the less evolutionary events is the best. They propose a MA, which uses diverse and elitist populations (similar with the ones used in scatter search methods). More precisely, their approach is based on maintaining a collection of Rec-I-DCM3 trees (Recursive-Iterative DCM3, a powerful heuristic for designing maximum parsimony trees [78]) which cooperate within a selectore-combinative evolutionary algorithm. They evaluate their method using biological datasets with up to 4,114 sequences, obtaining better results than parsimony ratchet [67] and TNT (Tree Analysis using New Technology<sup>3</sup>). Richer *et al.* [75] also uses maximum parsimony as the optimization criteria. They propose a MA that uses progressive neighborhood as local search (similar with VNS - variable neighborhood search [42]). They used twelve instances from TreeBase, and obtained results that were generally equal or better than TNT.

A problem related to phylogenetic inference is that of finding consensus trees, namely finding a tree that summarizes the information comprised in a collection of trees (e.g., finding a unique tree that faithfully amalgamates the outcome of different phylogenetic inference methods). A seminal approach to this problem using evolutionary methods can be found in [15] on the basis of the TreeRank distance measure [92] between trees. Pirkwieser and Raidl [72] tackled this problem using VNS, evolutionary algorithms (EAs), MAs (using EAs endowed with local search on different tree-based neighborhood structures), and multi-level hybrids based on the intertwined execution of VNS and EA/MA which ultimately produced the best results.

## 1.4 Protein Structure Analysis and Molecular Design

Problems involving analysis of protein structure are fundamental in bioinformatics. We refer to Oakley *et al.* [69] who present a review of problems involving analysis of protein structure (including structure prediction, structure comparison, aggregation of structures, etc.).

The protein structure prediction (PSP) problem aims to find the 3D structure with minimum energy (based in a specific energy model) given the primary sequence of the protein (i.e., the linear sequence of amino acids composing the protein). Krasnogor *et al.* [51] analyzed three main factors affecting the efficacy of evolutionary algorithms for PSP: the encoding scheme, the way illegal shapes are considered

---

<sup>2</sup> [www.treebase.org](http://www.treebase.org)

<sup>3</sup> <http://www.zmuc.dk/public/phylogeny/tnt/>

by the search, and the energy (fitness) function used. In [11] the protein structure prediction problem on the hydrophobic-polar (HP) model was considered. The HP model [26] is based on classifying each amino acid into two classes: hydrophobic or non-polar (H), and hydrophilic or polar (P), according to their interaction with water molecules. In this case the binary sequence of H/P amino acids is embedded in a cubic lattice subject to non-overlapping constraints, with the aim of maximizing the number of H-H contacts, namely the number of H-H pairs that are adjacent in the lattice. The MA featured the inclusion of a backtracking operator in order to repair infeasible protein configurations. A similar approach was used in [13] in the context of the HPNX energy model, an extension of the HP model in which polar amino acids are split into three classes: positively charged (P), negatively charged (N), and neutral (X). Krasnogor *et al.* [52] presented a multimemetic algorithm for protein structure prediction using four different models (HP in square and triangle lattice, and functional model proteins in the square and diamond lattice). Bazzoli and Tettamanzi [3] also considered the the HP cubic lattice model. They presented a MA using a self-adaptive strategy, where the local search is applied with a probability guided by a function similar with the one used in simulated annealing, with the aim to either control exploitation or diversification. According with the authors, the MA was strongly based on the MA proposed by Krasnogor and Smith [50], where the authors compared self-adaptation against other local-search approaches for traveling salesman problem. Santos and Santos [81] presents a MA for the protein structure problem using 2D triangular HP lattice model, whose main feature was the use of caching in order to reuse computation and speed-up fitness evaluation. The study of Zhao [102] addressed HP models as well. They described several metaheuristics such as MAs, tabu search, ant colony optimization, self-organizing map-based computing approaches and chain growth algorithm PERM, highlighting their advantages and disadvantages.

Protein structure comparison or protein alignment is another important problem in the area of protein structure analysis problem. In this case the goal is to identify structural similarities between proteins. Some MAs developed to deal with this problem can be found in [8, 49, 58, 90]. Carr *et al.* [8] considered the maximum contact map overlap problem. They presented a multimemetic algorithm where a family of local searches is used: selection of the particular local search to be applied depends on the instance, stage of the search or which individual is using it. The MA proposed is a combination of the genetic algorithm proposed by Lancia *et al.* [55] and six different local searches. Their computational results have showed that that the results obtained by their method is compatible with the state of art in this problem. Related to this problem. Also, Krasnogor [49] proposed a self-generating MA to obtain structural alignment between pair of proteins using the Maximum Contact Map Overlap (MaxCMO) problem as model. MaxCMO is an alignment of two proteins that maximizes the structural similarity. They tested the approach in four different data sets, which one was composed of randomly generated proteins and the other three data sets with real world proteins.

A bioinformatics area closely related to protein structure analysis is that of molecular design, which actually can be regarded as a superset of the former. In-



deed, conformational analysis, namely determining the low-energy configurations a molecule can adopt is a natural generalization of the protein structure prediction problem (for example, Zacharias *et al.* [101] presented a MA based on a genetic algorithm endowed with simulated annealing to determine the ground state geometry of molecular systems). In general, molecular design is a very hard problem, and numerous evolutionary approaches have been proposed in the literature to deal with problem in this area, check, e.g., [9, 94].

Ligand docking, i.e., the identification of putative ligands based on the geometry of the latter and that of a receptor site, is a problem within the area of molecular design with paramount interest for structure-based drug discovery. MA approaches to this problem have been proposed by Hart *et al.* [43, 63] using an evolutionary algorithm endowed with the Solis-Wets method for local search (see Chapter ??), aimed to minimizing the free energy potential of the docking. This MA is used in the AutoDock<sup>4</sup> software package. MAs have been also used for PCR (Polymerase chain reaction) product primer design [98], taking into account constraints such as primer length, GC content, melting temperature, etc.

## 1.5 Sequence Analysis

Sequence analysis is arguably one of the lowest-level tasks in bioinformatics, albeit it remains a very important one due to its role in generating the input data for further biological problems. Within this general subarea we can cite problems such as DNA sequencing and the alignment of genomic/proteomic sequences.

DNA sequencing amounts to determining the correct order of nucleotides in a certain DNA sequencing. This order must be ascertained by assembling short fragments of DNA obtained from the fragmentation by chemical or mechanical means of a larger sequence. These fragments are typically randomly distributed across the sequence and partially overlap, thus leading to a permutational problem with strong similarities to that of finding a minimum weight Hamiltonian path. In [27] a spatially-structured evolutionary algorithm endowed with a so-called problem-aware local search (PALS) procedure is presented for this purpose.

Another important problem in sequence analysis is that of aligning sequences of nucleotides or amino acids. This problem actually bears some relationship with sequencing, since the determination of the best overlap among DNA fragments requires finding the best pairwise alignment. The applications of sequence alignment are not limited to this case though; thus, they are very important in phylogenetic studies to cite a relevant example. This alignment problem is easily solvable in polynomial time for two sequences using a dynamic programming approach, but its complexity quickly grows for when a multiple sequence alignment is sought. Not surprisingly, evolutionary methods have been commonly applied to this problem – see [82] for a survey. Some of these evolutionary approaches can be actually

---

<sup>4</sup> <http://autodock.scripps.edu/>

regarded as memetic. For example, the evolutionary Clustal/improver presented in [87] incorporates a seeding mechanism (using the outcome of the Clustal<sup>5</sup> software package) for creating a high quality initial population, and a improvement strategy based on the removal of matched gap columns which can be regarded as a simple form of local search.

Closely related to sequence alignment, the problem of finding the shortest common supersequence (SCS) for a collection of biological sequences stands as another important task. A supersequence of a given sequence is a possibly longer sequence in which all the symbols of the former can be found in the same order (although not necessarily consecutively). Finding the SCS for a given collection of sequences is a NP-hard problem that has been commonly dealt with metaheuristics [5, 6, 12] including MAs. Thus, Cotta [14] considered a MA defined on the basis of an evolutionary algorithm endowed with a repairing mechanism (based on a greedy heuristic) and a local search operator based on the iterative removal of symbols in the tentative supersequence. Later, Gallardo *et al.* [32] presented a multi-level MA that combined the previous algorithm with a beam search algorithm (see Chapter ??), executed in an intertwined way. This MA was shown to provide much better results than the combined algorithm as stand-alone techniques.

## 1.6 Systems Biology

Systems biology [1] is a prominent interdisciplinary area of bioscientific research focusing on the holistic study of cellular systems from the perspective of (and using tools from) complex systems and dynamical systems theory. This encompasses the analysis and modeling of cell systems, including the study of networks of genomic/proteomic/metabolomic interactions. The latter are very amenable to the use of network-theoretical results and graph-based algorithmic tools, among which MAs excel. Thus, Spieth *et al.* consider a memetic approach to gene regulatory network modeling using linear weight matrices [93] and S-systems [91]. They use a binary genetic algorithm to evolve the topology of the network, and an evolution strategy to do local search on the parameters of the model representing the network. They consider a so-called *feedback MA* in which the outcome of the local search is used to filter gene dependencies whose strength is below a certain threshold. This can be regarded as a Lamarckian learning procedure, as opposed to the Baldwinian learning of the simpler MA [85] without feedback. An analogous approach is followed by Norman and Iba [68]: they consider time series data of gene expression and use a differential evolution endowed with hill climbing to determine the structure of the network and the kinetic parameters; an information-based criterion is used for fitness evaluation. It is also worth mentioning the work of Kimura *et al.* [47] in which a genetic local search method is used to solve the inference problem in the context of S-systems. In a later work [48], they consider a cooperative approach based on mul-

---

<sup>5</sup> <http://www.clustal.org/>

multiple subpopulations and problem decomposition and use golden section search in order to do local improvement. Tsai and Wang [88] consider a differential evolution hybridized with local search for S-system inference too.

A wider perspective on cell models is provided by [77]. They consider the use of P-systems [73], a computing model included in the ampler paradigm of membrane computing [74]. These computational models are inspired in cellular processes, and can be roughly described as a system of so-called *membrane structures*, namely permeable (and potentially nested) containers that comprise collections of symbols and grammar-like rules for their evolution. By an appropriate definition of the rules and a wise arrangement of membranes it is possible to carry out an arbitrary computation. The biological inspiration of these systems make them specifically suited for cell modeling and simulation though. Romero-Campero *et al.* use a two-level genetic algorithm to evolve the structure of a P-system: the upper level is devoted to searching in the space of rules, and the lower level performs numerical adjustment of the kinetic parameters determining the probability of application of each rule.

**Acknowledgements** C. Cotta is supported by Spanish MICINN under project NEMESIS (TIN2008-05941) and Junta de Andalucía under project TIC-6083.



## References

1. Alon U (2006) An Introduction to Systems Biology – Design Principles of Biological Circuits, Mathematical and Computational Biology Series, vol 10. Chapman & Hall/Crc.
2. Althaus E, Klau G, Kohlbacher O, Lenhof HP, Reinert K (2009) Efficient algorithms. In: Albers S, Alt H, Näher S (eds) Integer Linear Programming in Computational Biology, Springer-Verlag, Berlin Heidelberg, Lecture Notes in Computer Science, vol 5760, pp 199–218
3. Bazzoli A, Tettamanzi A (2004) A memetic algorithm for protein structure prediction in a 3D-lattice HP model. In: Raidl GR, et al (eds) EvoWorkshops, Lecture Notes in Computer Science, vol 3005, Springer-Verlag, Coimbra, Portugal, pp 1–10
4. Berretta R, Costa W, Moscato P (2008) Combinatorial optimization models for finding genetic signatures from gene expression datasets. In: Keith JM (ed) Bioinformatics, Volume II: Structure, Function and Applications, Methods in Molecular Biology, Humana Press, chap 19, pp 363–378
5. Blum C, Cotta C, Fernández AJ, Gallardo JE (2007) A probabilistic beam search approach to the shortest common supersequence problem. In: Cotta C, van Hemert JI (eds) Evolutionary Computation in Combinatorial Optimization 2007, Springer-Verlag, Berlin Heidelberg, Lecture Notes in Computer Science, vol 4446, pp 36–47
6. Branke J, Middendorf M, Schneider F (1998) Improved heuristics and a genetic algorithm for finding short supersequences. *OR-Spektrum* 20:39–45
7. Cadieux S, Tanizaki N, Okamura T (1997) Time efficient and robust 3-D brain image centering and realignment using hybrid genetic algorithm. In: 36th SICE Annual Conference, IEEE Press, pp 1279–1284
8. Carr R, Hart W, Krasnogor N, Hirst J, Burke E (2002) Alignment of protein structures with a memetic evolutionary algorithm. In: Langdon WB, et al (eds) GECCO 2002, Morgan Kaufmann, New York NY, pp 1027–1034
9. Clark D, Westhead D (1996) Evolutionary algorithms in computer-aided molecular design. *Journal of Computer-aided Molecular Design* 10(4):337–358

10. Cordón O, Damas S, Santamaría J, Martí R (2008) Scatter search for the 3D point matching problem in image registration. *INFORMS Journal on Computing* 20(1):55–68
11. Cotta C (2003) Protein structure prediction using evolutionary algorithms hybridized with backtracking. In: Mira J, Álvarez J (eds) *Artificial Neural Nets Problem Solving Methods*, Springer-Verlag, Berlin Heidelberg, *Lecture Notes in Computer Science*, vol 2687, pp 321–328
12. Cotta C (2005) A comparison of evolutionary approaches to the shortest common supersequence problem. In: Cabestany J, Prieto A, Hernández FS (eds) *Computational Intelligence and Bioinspired Systems*, Springer-Verlag, Berlin Heidelberg, *Lecture Notes in Computer Science*, vol 3512, pp 50–58
13. Cotta C (2005) Hybrid evolutionary algorithms for protein structure prediction in the HPNX model. In: *Computational Intelligence, Theory and Applications, Advances in Soft Computing*, vol 2, Springer-Verlag, pp 525–534
14. Cotta C (2005) Memetic algorithms with partial lamarckism for the shortest common supersequence problem. In: Mira J, Álvarez J (eds) *IWINAC 2005*, Springer-Verlag, Berlin Heidelberg, *Lecture Notes in Computer Science*, vol 3562, pp 84–91
15. Cotta C (2005) On the application of evolutionary algorithms to the consensus tree problem. In: Gottlieb J, Raidl G (eds) *Evolutionary Computation in Combinatorial Optimization*, Springer-Verlag, Berlin, *Lecture Notes in Computer Science*, vol 3248, pp 58–67
16. Cotta C (2005) Scatter search with path relinking for phylogenetic inference. *European Journal of Operational Research* 169(2):520–532
17. Cotta C, Moscato P (2002) Inferring phylogenetic trees using evolutionary algorithms. In: [61], pp 720–729
18. Cotta C, Moscato P (2003) A memetic-aided approach to hierarchical clustering from distance matrices: Application to phylogeny and gene expression clustering. *Biosystems* 72(1-2):75–97
19. Cotta C, Mendes A, Garcia V, Franca P, Moscato P (2003) Applying memetic algorithms to the analysis of microarray data. In: *Applications of Evolutionary Computing 2003*, *Lecture Notes in Computer Science*, vol 2611, Springer-Verlag, pp 77–81
20. Cox M, Bowden N, Moscato P, Berretta R, Scott RI, Lechner-Scott JS (2007) Memetic algorithms as a new method to interpret gene expression profiles in multiple sclerosis. In: *Abstracts of the 23rd Congress of the European Committee for Treatment and Research in Multiple Sclerosis and the 12th Annual Conference of Rehabilitation in Multiple Sclerosis*, Prague, Czech Republic, vol 13, Suppl. 2, p S205
21. Davies DL, Bouldin DW (1979) A cluster separation measure. *IEEE Transactions on Pattern Analysis and Machine Intelligence* PAMI-1:224–227
22. Day W (1983) Computationally difficult problems in phylogeny systematics. *Journal of Theoretic Biology* 103:429–438
23. Day W (1987) Computational complexity of inferring phylogenies from dissimilarity matrices. *Bulletin of Mathematical Biology* 49(4):461–467

24. Di Gesù V, Bosco G, Lo G, Millonzi F, Valenti C (2008) A memetic algorithm for binary image reconstruction. In: Brimkov VE, Barneva RP, Hautman HA (eds) *Combinatorial Image Analysis, Lecture Notes in Computer Science*, vol 5958, Springer-Verlag, pp 384–395
25. Di Gesù V, Bosco GL, Millonzi F, Valenti C (2008) Discrete tomography reconstruction through a new memetic algorithm. In: Giacobini M et al (ed) *Applications of Evolutionary Computing, Lecture Notes in Computer Science*, vol 4974, Springer, pp 347–352
26. Dill K (1990) Dominant forces in protein folding. *Biochemistry* 29:7133–7155
27. Dorronsoro B, Alba E, Luque G, Bouvry P (2008) A self-adaptive cellular memetic algorithm for the dna fragment assembly problem. In: CEC 2008, IEEE Press, Hong Kong, pp 2656–2663
28. Eisen M, Spellman P, Brown P, Botstein D (1998) Cluster analysis and display of genome-wide expression patterns. *National Academy of Sciences* 95(25):14,863–14,868
29. Festa P (2007) On some optimization problems in molecular biology. *Mathematical Biosciences* 207(2):219–234
30. Foulds L, Graham R (1982) The steiner problem in phylogeny is NP-complete. *Advances in Applied Mathematics* 3(1):43–49
31. Fourer R, Gay DM, Kernighan BW (2003) *AMPL: A Modeling Language for Mathematical Programming*, 2nd Edition. Duxbury Press Brooks Cole Publishing Co.
32. Gallardo JE, Cotta C, Fernández AJ (2007) On the hybridization of memetic algorithms with branch-and-bound techniques. *IEEE Transactions on Systems, Man and Cybernetics, part B* 37(1):77–83
33. Gallardo JE, Cotta C, Fernández AJ (2007) Reconstructing phylogenies with memetic algorithms and branch-and-bound. In: Bandyopadhyay S, Maulik U, Wang J (eds) *Analysis of Biological Data: A Soft Computing Approach*, World Scientific, pp 59–84
34. Glover F (1999) *Scatter search and path relinking*. McGraw-Hill, Maidenhead, Berkshire, England, UK, pp 291–316
35. Glover F, Laguna M, Martí R (2000) Fundamentals of scatter search and path relinking. *Control and Cybernetics* 39(3):653–684
36. Gómez Ravetti M, Rosso OA, Berretta R, Moscato P (2010) Uncovering molecular biomarkers that correlate cognitive decline with the changes of hippocampus' gene expression profiles in alzheimer's disease. *PLoS ONE* 5(4):e10,153
37. Greenberg HJ, Hart WE, Lancia G (2004) Opportunities for combinatorial optimization in computational biology. *INFORMS Journal on Computing* 16(3):211–231
38. Guéret C, Prins C, Sevaux M (2002) *Applications of optimisation with Xpress-MP*. Dash Optimization
39. Guerra-Salcedo C, Chen S, Whitley D, Smith S (1999) Fast and accurate feature selection using hybrid genetic strategies. In: CEC 1999, IEEE Press, Washington DC, pp 177–184

40. Haas O, Burnham K, Mills J, Reeves C, Fisher M (1996) Hybrid genetic algorithms applied to beam orientation in radiotherapy. In: Proceedings of 4th European Conference on Intelligent Techniques and Soft Computing, Verlag Mainz, vol 3, pp 2050–2055
41. Haas O, Burnham K, Mills J (1998) Optimization of beam orientation in radiotherapy using planar geometry. *Physics in Medicine and Biology* 43(8):2179–2193
42. Hansen P, Mladenović N (2002) Variable neighborhood search. In: Glover F, Kochenberger G (eds) *Handbook of Metaheuristics*, Kluwer Academic Publishers, Boston MA, pp 145–184
43. Hart W, Rosin C, Belew R, Morris G (2000) Improved evolutionary hybrids for flexible ligand docking in autodock. In: Floudas CA, Pardalos PM (eds) *Optimization in Computational Chemistry and Molecular Biology, Nonconvex Optimization and Its Applications*, vol 40, Springer-Verlag, pp 209–230
44. Hourani M, Berretta R, Mendes A, Moscato P (2008) Genetic signatures for a rodent model of parkinson’s disease using combinatorial optimization methods. In: Keith JM (ed) *Bioinformatics, Volume II: Structure, Function and Applications, Methods in Molecura Biology*, Humana Press, pp 379–392
45. Il-Seok O, Jin-Seon L, Byung-Ro M (2004) Hybrid genetic algorithms for feature selection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 26:1424–1437
46. Inostroza-Ponta M, Berretta R, Moscato P (2011) Qapgrid: A two level qap-based approach for large-scale data analysis and visualization. *PLoS ONE* 6(1):e14,468
47. Kimura S, Hatakeyama M, Konagaya A (2004) Inference of s-system models of genetic networks from noisy time-series data. *Chem-Bio Informatics Journal* 4(1):1–14
48. Kimura S, et al (2005) Inference of S-system models of genetic networks using a cooperative coevolutionary algorithm. *Bioinformatics* 21(7):1154–1163
49. Krasnogor N (2004) Self-generating metaheuristics in bioinformatics: The protein structure comparison case. *Genetic Programming and Evolvable Machines* 5(2):181–201
50. Krasnogor N, Smith J (2000) A memetic algorithm with self-adaptive local search: TSP as a case study. In: Whitley LD, et al (eds) *GECCO 2000*, Morgan Kaufmann, Las Vegas NV, pp 987–994
51. Krasnogor N, Hart W, Smith J, Pelta D (1999) Protein structure prediction with evolutionary algorithms. In: Banzhaf W, et al (eds) *GECCO 1999*, Morgan Kaufmann, Orlando FL, pp 1569–1601
52. Krasnogor N, Blackburne B, Burke E, Hirst J (2002) Multimeme algorithms for proteine structure prediction. In: [61], pp 769–778
53. Laguna M, Martí R (2003) *Scatter Search. Methodology and Implementations in C*. Kluwer Academic Publishers, Boston MA
54. Lancia G (2008) *Mathematical programming in computational biology: an annotated bibliography*. *Algorithms* 1(2):100–129



55. Lancia G, Carr R, Walenz B, Istrail S (2001) 101 optimal pdb structure alignments: a branch-and-cut algorithm for the maximum contact map overlap problem. In: Fifth Annual International Conference on Computational Molecular Biology, RECOMB, ACM Press, pp 193–202
56. Liang Y, Leung KS, Mok TSK (2004) Evolutionary drug scheduling model for cancer chemotherapy. In: Deb K, et al (eds) GECCO 2004, Springer-Verlag, Seattle WA, Lecture Notes in Computer Science, vol 3102, pp 1126–1137
57. Liang Y, Leung KS, Mok TSK (2008) Evolutionary drug scheduling models with different toxicity metabolism in cancer chemotherapy. *Applied Soft Computing* 8(1):140–149
58. May A, Johnson M (1994) Protein-structure comparisons using a combination of a genetic algorithm, dynamic-programming and least-squares minimization. *Protein Engineering* 7(4):475–485
59. Mendes A, Cotta C, Garcia V, França P, Moscato P (2005) Gene ordering in microarray data using parallel memetic algorithms. In: Skie T, Yang CS (eds) 2005 International Conference on Parallel Processing Workshops, IEEE Press, Oslo, Norway, pp 604–611
60. Mendes A, Scott R, Moscato P (2007) Microarrays - identifying molecular portraits for prostate tumors with different gleason patterns. In: Trent R (ed) *Clinical Bioinformatics - Methods in Molecular Medicine*, Methods in Molecular Medicine, vol 141, Humana Press, pp 131–151
61. Merelo Guervós JJ, et al (eds) (2002) *Parallel Problem Solving from Nature VII*, Lecture Notes in Computer Science, vol 2439, Springer-Verlag, Granada, Spain
62. Merz P, Zell A (2002) Clustering gene expression profiles with memetic algorithms. In: [61], pp 811 – 820
63. Morris G, Goodsell D, Halliday R, Huey R, Hart W, Belew R, AJOlson (1998) Automated docking using a lamarckian genetic algorithm and an empirical binding free energy function. *Journal of Computational Chemistry* 19(14):1639–1662
64. Moscato P, Mendes A, Berretta R (2007) Benchmarking a memetic algorithm for ordering microarray data. *Biosystems* 88(1-2):56–75
65. Neri F, Toivanen J, Cascella GL, Ong YS (2007) An adaptive multimeme algorithm for designing HIV multidrug therapies. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 4(2):264–278
66. Nilsson R, Björkegren J, Tegnér J (2009) On reliable discovery of molecular signatures. *BMC Bioinformatics* 10(38)
67. Nixon KC (1999) The parsimony ratchet, a new method for rapid parsimony analysis. *Cladistics* 15:407–414
68. Noman N, Iba H (2007) Inferring gene regulatory networks using differential evolution with local search heuristics. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 4(4):634–647
69. Oakley M, Barthel D, Bykov Y, Garibaldi J, Burke E, Krasnogor N, Hirst J (2008) Search strategies in structural bioinformatics. *Current Protein and Peptide Science* 9(3):260–274

70. Palacios P, Pelta D, Blanco A (2006) Obtaining biclusters in microarrays with population-based heuristics. In: Rothlauf F, et al (eds) *Applications of Evolutionary Computing 2006*, Springer-Verlag, Budapest, Hungary, *Lecture Notes in Computer Science*, vol 3907, pp 115–126
71. Pastorino M (2007) Stochastic optimization methods applied to microwave imaging: A review. *IEEE Transactions on Antennas and Propagation* 55(3, Part 1):538–548
72. Pirkwieser S, Raidl GR (2008) Finding consensus trees by evolutionary, variable neighborhood search, and hybrid algorithms. In: [79], pp 323–330
73. Păun G (1998) Computing with membranes. Tech. Rep. TUCS Report 208, Turku Center for Computer Science
74. Păun G (2002) *Membrane Computing: An Introduction*. Springer-Verlag, Berlin
75. Richer JM, Goeffon A, Hao JK (2009) A memetic algorithm for phylogenetic reconstruction with maximum parsimony. In: *Evolutionary Computation, Machine Learning and Data Mining in Bioinformatics*, Springer-Verlag, Berlin Heidelberg, *Lecture Notes in Computer Science*, vol 5483, pp 164–175
76. Riveros C, et al (2010) A transcription factor map as revealed by a genome-wide gene expression analysis of whole-blood mRNA transcriptome in multiple sclerosis. *PLoS ONE* 5(12):e14,176
77. Romero-Campero F, Cao H, Camara M, Krasnogor N (2008) Structure and parameter estimation for cell systems biology models. In: [79], pp 331–338
78. Roshan U, Moret BME, Warnow T, Williams TL (2004) Rec-i-dcm3: A fast algorithmic technique for reconstructing large phylogenetic trees. In: *IEEE Computer Society Bioinformatics Conference 2004*, IEEE Press, pp 98–109
79. Ryan C, Keijzer M (eds) (2008) *Genetic and Evolutionary Computation Conference – GECCO 2008*, ACM Press, Atlanta GA
80. Santamaría J, Córdón O, Damas S, García-Torres JM, Quirin A (2009) Performance evaluation of memetic approaches in 3D reconstruction of forensic objects. *Soft Computing* 13(8-9):883–904
81. Santos E, Santos J (2006) Effective computational reuse for energy evaluations in protein folding. *International Journal of Artificial Intelligence Tools* 15(5):725–739
82. Shyu C, Sheneman L, Foster JA (2004) Multiple sequence alignment with evolutionary computation. *Genetic Programming and Evolvable Machines* 5:121–144
83. Speer N, Merz P, Spieth C, Zell A (2003) Clustering gene expression data with memetic algorithms based on minimum spanning trees. In: *CEC 2003*, IEEE Press, Canberra, Australia, pp 1848–1855
84. Speer N, Spieth C, Zell A (2004) A memetic clustering algorithm for the functional partition of genes based on the gene ontology. In: *IEEE Symposium on Computational Intelligence in Bioinformatics and Computational Biology*, IEEE Press, pp 252–259

85. Spieth C, Streichert F, Speer N, Zell A (2004) A memetic inference method for gene regulatory networks based on S-systems. In: CEC 2004, IEEE Press, Portland OR, pp 152–157
86. Spieth C, Streichert F, Supper J, Speer N, Zell A (2005) Feedback memetic algorithms for modeling gene regulatory networks. In: 2005 IEEE Symposium on Computational Intelligence in Bioinformatics and Computational Biology., IEEE Press, pp 61–67
87. Thomsen R, Munkegade N, Fogel GB, Krink T, Group E, Group E (2002) A clustal alignment improver using evolutionary algorithms. In: CEC 2002, IEEE Press, Honolulu HI, pp 121–126
88. Tsai KY, Wang FS (2005) Evolutionary optimization with data collocation for reverse engineering of biological networks. *Bioinformatics* 21(7):1180–1188
89. Tse SM, Liang Y, Leung KS, Lee KH, Mok T (2007) A memetic algorithm for multiple-drug cancer chemotherapy schedule optimization. *IEEE Transactions on Systems, Man, and Cybernetics, Part B* 37(1):84–91
90. Volk J, Herrmann T, Wuthrich K (2008) Automated sequence-specific protein nmr assignment using the memetic algorithm match. *Journal of Biomolecular Nmr* 41(3):127–138
91. Wang H, Qian L, Dougherty E (2010) Inference of gene regulatory networks using S-system: a unified approach. *IET Systems Biology* 4(2):145–156
92. Wang J, Shan H, Shasha D, Piel W (2003) Treerank: A similarity measure for nearest neighbor searching in phylogenetic databases. In: 15th International Conference on Scientific and Statistical Database Management, IEEE Press, Cambridge MA, pp 171–180
93. Weaver D, Workman C, Stormo G (1999) Modeling regulatory networks with weight matrices. *Pacific Symposium on Biocomputing* 4:112–123
94. Willett P (1995) Genetic algorithms in molecular recognition and design. *Trends in Biotechnology* 13(12):516–521
95. Williams HP (2000) *Model Building in Mathematical Programming*, 4th Edition. Wiley
96. Williams T, Smith M (2006) The role of diverse populations in phylogenetic analysis. In: Cattolico M (ed) *GECCO 2006*, ACM Press, Seattle WA, pp 287–294
97. Wu B, Chao KM, Tang C (1999) Approximation and exact algorithms for constructing minimum ultrametric trees from distance matrices. *Journal of Combinatorial Optimization* 3(2):199–211
98. Yang CH, Cheng YH, Chuang LY, Chang HW (2009) Specific PCR product primer design using memetic algorithm. *Biotechnol Prog* 25(3):745–753
99. Ying Xu DX V Olman (2002) Clustering gene expression data using a graph-theoretic approach: An application of minimum spanning tree. *Bioinformatics* 18(4):526 – 535
100. Z Zhu MD Y S Ong (2007) Markov blanket-embedded genetic algorithm for gene selection. *Pattern Recognition archive* 40(11):3236–3248

101. Zacharias CR, Lemes MR, Pino AD (1998) Combining genetic algorithm and simulated annealing: a molecular geometry optimization study. *Journal of Molecular Structure: THEOCHEM* 430:29 – 39
102. Zhao XC (2008) Advances on protein folding simulations based on the lattice hp models with natural computing. *Applied Soft Computing* 8(2):1029–1040
103. Zhu Z, Ong YS (2007) Memetic algorithms for feature selection on microarray data. In: *Advances in Neural Networks – ISNN 2007, Lecture Notes in Computer Science*, vol 4491, Springer-Verlag, pp 1327–1335
104. Zhu Z, Ong YS, Dash M (2007) Wrapper-filter feature selection algorithm using a memetic framework. *IEEE Transactions on Systems, Man, and Cybernetics, Part B* 37(1):70–76
105. Zhu Z, Jia S, Ji Z (2010) Towards a memetic feature selection paradigm. *Computational Intelligence Magazine* 5:41–53

# Index

- ant colony optimization, 8
- Baldwinian learning, 10
- bioinformatics, 1–11
  - cell models, 11
  - clustering, 4–5
  - conformational analysis, 9
  - consensus tree, 7
  - DNA sequencing, 9
  - feature selection, 5–6
    - filter vs wrapper methods, 5
  - gene ordering, 3
  - gene regulatory networks, 10
  - ligand docking, 9
  - microarray analysis, 2–6
  - molecular design, 8–9
  - molecular signature, 3
  - phylogeny, 6–7
    - maximum parsimony, 7
    - ultrametric tree, 6
  - polymerase chain reaction, 9
  - protein alignment, 8
  - protein structure analysis, 7–8
  - protein structure prediction, 7
    - HP model, 8
  - sequence alignment, 9
  - sequence analysis, 9–10
  - shortest common supersequence, 10
  - systems biology, 10–11
- biomedicine, 2
  - drug therapy scheduling, 2
  - radiotherapy, 2
  - tomography, 2
- branch and bound, 7
- differential evolution, 10
- hybridization
  - with beam search, 10
  - with evolution strategy, 10
  - with hill climbing, 10
  - with simulated annealing, 9
- local search
  - golden section search, 11
  - Solis-Wets, 9
- Markov blanket, 5
- membrane computing, 11
- memetic algorithm
  - multimeme, 8
  - self-adaptive, 8
- microarray analysis, 2–6
- minimum spanning tree, 4
- P-systems, 11
- path relinking, 6
- Rec-I-DCM3, 7
- S-systems, 10
- scatter search, 6
- tabu search, 8
- variable neighborhood search, 7