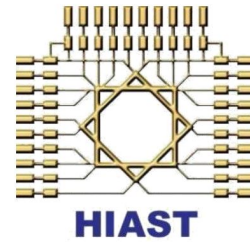


Syrian Arab Republic  
Higher Institute for Applied Sciences and Technology  
Informatics Department  
Fourth Year



# Using Genetic Algorithms to Find Approximations for the Minimum Vertex Cover Problem

Keywords: genetic algorithms, NP-complete, combinatorial optimization,  
non-deterministic algorithms, approximation algorithms, minimum vertex cover.

Author: *Farouk Hjabo*  
Academic Supervisor: *Dr. Said Desouki*  
General Supervisor: *Dr. Kadan Aljoumaa*  
Langauge Supervisor: *Mr. Fahmi Alammareen*

February 23, 2018

## Abstract

Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

*“It is not the strongest of the species that survives, nor the most intelligent, but the one most responsive to change.”*

— Charles Darwin

# Contents

Cover Page	i
Abstract	ii
Contents	iii
<b>1 Introduction</b>	<b>1</b>
<b>2 Genetic Algorithms</b>	<b>1</b>
2.1 The Intuition Behind GAs . . . . .	1
2.2 Typical GA Schema . . . . .	2
<b>3 The Minimum Vertex Cover Problem</b>	<b>4</b>
3.1 Related Algorithms . . . . .	5
3.2 GA for the mvcp . . . . .	7
<b>4 Experimental Results</b>	<b>8</b>
4.1 experimental observations . . . . .	11
<b>5 Conclusion</b>	<b>11</b>
References	13

# 1 Introduction

Once the NP-hardness of a combinatorial optimization problem is established, the search for an optimal solution is abandoned. The goal then becomes one of finding a good heuristic, i.e. a polynomial running time algorithm that can find solutions close to the optimal. In most cases, traditional heuristics are problem dependent; a heuristic is tailored to the specific problem it is trying to solve.

In this paper, we present an alternative approach that uses genetic algorithms as a generalized heuristic for solving NP-hard combinatorial optimization problems. The application of a genetic algorithm is demonstrated here for the *minimum vertex cover* problem. These algorithms have been successfully applied to a broad range of problems. This wide range can be tackled by genetic algorithms mainly due to the fact that they work with an encoding of the domain rather than with the problem domain itself.

## 2 Genetic Algorithms

### 2.1 The Intuition Behind GAs

Genetic Algorithms (GAs) [1] are population based search algorithms, where by repeated use of genetic operations, such as ***mutation***, ***selection***, ***crossover***, etc... Successive new generations of better populations in the direction of search objectives are created. They are inspired by Darwinian principles based on natural evolution. In other words, the main idea behind genetic algorithms is that only the ***fittest*** individuals will survive.

Main advantage of genetic algorithms is that they ideally do not make any assumption about the underlying problem, hence they are suitable to tackle a wide range of diverse problems in engineering, art, biology, economics, marketing, genetics, operations research, robotics, social sciences, physics, politics, chemistry, etc...

So what is a GA? A typical GA consists of the following:

1. a number, or ***population***, of candidate solutions to the problem.
2. a way of calculating how good or bad the individual solutions within the population are. i.e. how ***fit*** an individual is?

3. a method for mixing fragments of the better solutions to form new, on average even better solutions. Which permits the population to *evolve* naturally.
4. a *mutation* operator to avoid permanent loss of diversity within the solutions. This allows introducing new information in the population.

## 2.2 Typical GA Schema

The basic iteration cycle of a genetic algorithm proceeds on a population of individuals, each of which represents a search point in the space of potential solutions of a given optimization problem. In case of a canonical genetic algorithm, each individual is a binary vector  $\vec{x} = (x_1, \dots, x_n) \in \{0, 1\}^n$  of fixed length  $n$ . The fitness function  $f : \{0, 1\}^n \rightarrow \mathbb{R}$  provides a quality measure which is used by the selection procedure to direct the search towards regions of the search space where the average fitness of the population increases. The recombination operator allows for the exchange of information between different individuals, and mutation introduces innovation into the population.

This cycle is repeated until a termination criterion is fulfilled. In most cases, the algorithm is terminated after a certain number of iterations of the basic cycle or when a satisfactory fitness level has been reached. refer to figure 1 for a visual representation of this process.

**Initialization** The population size (denoted as  $u$ ) depends on the nature of the problem, but typically contains several hundreds or thousands of possible solutions. Often, the initial population is generated randomly, allowing the entire range of possible solutions (the search space). Occasionally, the solutions may be ‘seeded’ in areas where optimal solutions are likely to be found.

**Selection** Normally, selection in genetic algorithms is a probabilistic operator which uses the relative fitness  $p_s(\vec{x}_i) = f(\vec{x}_i) / \sum_{j=1}^u f(\vec{x}_j)$  to serve as selection probabilities ( $u$  denotes the population size). This selection operator is called proportional selection. If the problem under consideration is a minimization one, or if the fitness function can take negative values, then  $f(\vec{x}_i)$  has to be linearly transformed before calculating selection probabilities. This technique known as linear dynamic scaling is commonly used in genetic algorithms (see [13], pp. 123-124, or [14]).

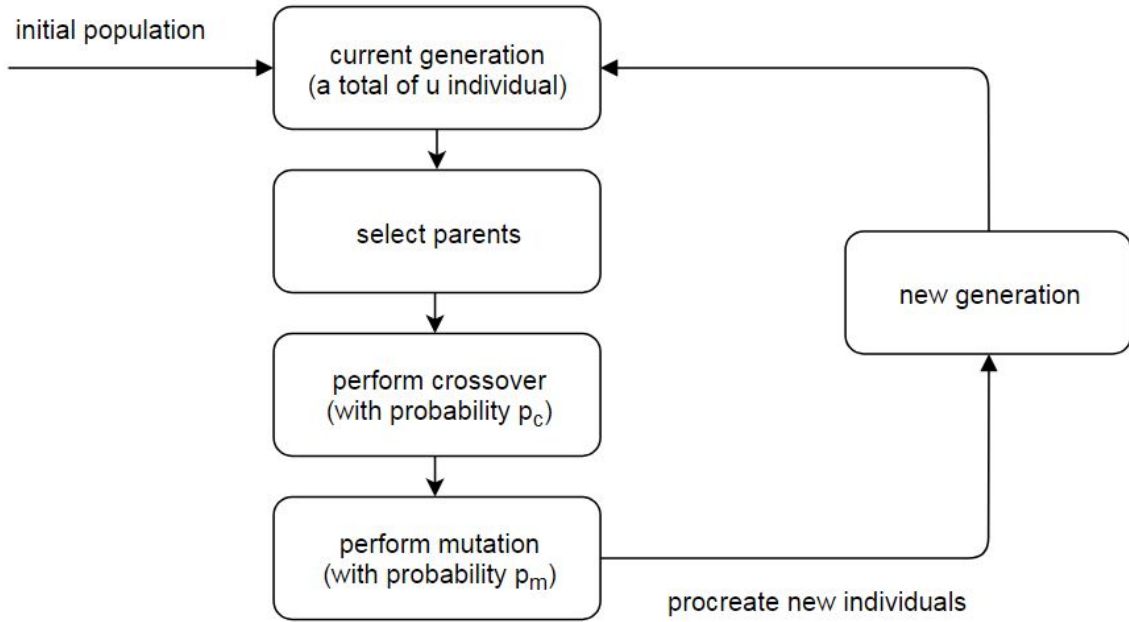


Figure 1: general schema for genetic algorithms

**Crossover** The recombination (crossover) operator allows for the exchange of information between different individuals. The original one-point crossover [9] works on two parent individuals (which are randomly chosen from the population, see the previous paragraph todo) by choosing a crossover point  $\chi \in \{1, \dots, n-1\}$  at random and exchanging all bits after the  $\chi^{\text{th}}$  one between both individuals. The crossover rate  $p_c$  (e.g.  $p_c \approx 0.6$ ) determines the probability to undergo crossover (see todo for more details). The crossover operator can be extended to a generalized multi-point crossover [10] or even to uniform crossover, where an it is randomly decided for each bit whether to exchange it or not [19]. The strong mixing effect introduced by uniform crossover is sometimes helpful to overcome local optima.

**Mutation** Innovation, i.e. new information, is introduced into the population by means of mutation, which works by inverting bits with a small probability  $p_m$  (e.g.  $p_m \approx 0.001$ ) Though mutation is often interpreted as a rather unimportant operator in genetic algorithms [9], recent theoretical work gives strong evidence for an appropriate choice of a mutation rate  $p_m = 1/n$  on many problems [2; 12].

For better understanding of figure 1

---

**Algorithm 1** Gentic Algorithm( $n, u, p_c, p_m, f$ )

---

```
1:  $P \leftarrow \{\}$  ▷  $P$  contains the population
2: loop  $u$  times ▷ initialization
3:   add a randomly generated binary vector of length  $n$  to  $P$ 
4: end loop
5: while termination criterion is not fulfilled do
6:    $Q \leftarrow \{\}$  ▷  $Q$  will hold the new generation
7:   loop  $u/2$  times
8:     select  $x, y$  from  $P$  using proportional selection ▷ selection
9:     with probability  $p_c$  perform crossover between  $x, y$  ▷ crossover
10:    for each bit in  $x$  flip it with probability  $p_m$  ▷ mutation
11:    for each bit in  $y$  flip it with probability  $p_m$  ▷ mutation
12:    add  $x, y$  to  $Q$ 
13:  end loop
14:   $P \leftarrow Q$ 
15: end while
```

---

### 3 The Minimum Vertex Cover Problem

The problem of finding a minimum vertex cover is a classical optimization problem in computer science and is a typical example of an NP-hard optimization problem that has an approximation algorithm. The mvcp of an undirected graph  $G = (V, E)$  where  $V$  is the set of vertices and  $E$  denotes the set of edges, consists in finding the smallest subset  $V' \subseteq V$  such that  $\forall \langle i, j \rangle \in E$ , we have  $i \in V'$  or  $j \in V'$  (or both).  $V'$  is said to be a vertex cover of  $G$ . The following is a formal definition of the mvcp in which we make use of Stinson's terminology for combinatorial optimization problems [14]:

**Problem instance:** A graph  $G = (V, E)$ , where  $V = \{1, 2, \dots, n\}$  is the set of vertices and  $E \subseteq V \times V$  the set of edges. An edge between vertices  $i, j$  is denoted by the pair  $\langle i, j \rangle \in E$ . We define the adjacency matrix  $e_{ij}$  according to:

$$e_{ij} = \begin{cases} 1 & \text{if } \langle i, j \rangle \in E \\ 0 & \text{otherwise} \end{cases}$$

especially, we have  $e_{ii} = 0$

**Feasible solution:** A set  $V' \subseteq V$  such that  $\forall \langle i, j \rangle \in E, i \in V' \vee j \in V'$

**Objective function:** The size  $|V'|$  of the vertex cover  $V'$ .

**Optimal solution:** a vertex cover  $V'$  that minimizes  $|V'|$ .

an interesting property of the mvcp is that by solving it, we also have a solution for two other graph problems: the *maximum independent set problem*<sup>\*</sup> and the *maximum clique problem*<sup>†</sup>. The close relationship between these problems is characterized by the following lemma [6].

### Lemma 1

For any graph  $G = (V, E)$  and  $V' \subseteq V$ , the following statements are equivalent:

- $V'$  is the minimum vertex cover of  $G$ .
- $V - V'$  is the maximum independent set in  $G$ .
- $V - V'$  is the maximum clique in  $\bar{G} = (V, \bar{E})$ , where  $\bar{E} = \{\langle i, j \rangle \in V \times V : \langle i, j \rangle \notin E\}$

Consequently, one can obtain a solution of the maximum independent set problem by taking the complement of the solution to the minimum vertex cover problem. A solution to the maximum clique problem is obtained by taking the complement of the minimum vertex cover of  $\bar{G}$ .

## 3.1 Related Algorithms

Since we are interested in covering all the edges of the given graph by using as few nodes as possible, one might be tempted to use a greedy-based heuristic to tackle the mvcp. The algorithm consists in repeatedly selecting a vertex of highest degree (the node that covers as many of the remaining edges as possible), and removing all of its incident edges. This is not a good strategy as was demonstrated by Papadimitriou and Steiglitz ([10], p. 407).

They considered regular graphs, each of which consists of three levels. The first two levels have the same number of nodes while the third level has two nodes less than the number of nodes found on the previous two levels. More precisely, each graph consists

---

<sup>\*</sup>misp: given  $G = (V, E)$ , find the largest subset  $V' \subseteq V$  such that  $\forall i, j \in V', \langle i, j \rangle \notin E$

<sup>†</sup>mcp: given  $G = (V, E)$ , find the largest subset  $V' \subseteq V$  such that  $\forall i, j \in V', \langle i, j \rangle \in E$



of  $n = 3k + 4$  ( $k \geq 1$ ) nodes;  $k + 2$  nodes on the first level, labeled  $1, 2, \dots, k + 2$ ; followed by  $k + 2$  nodes on the second level, labeled  $k + 3, \dots, 2k + 4$ ; while the  $k$  nodes of the third level are labeled  $2k + 5, \dots, 3k + 4$ . The regular graph for  $k = 3$  can be found in figure 2. A minimum vertex cover is obtained by choosing all the nodes of the second level, but the greedy strategy would start with the nodes of the third level, since these have highest degree. Consequently, the greedy strategy finds a solution of size  $2k + 2$  (since it has to select  $k + 2$  nodes in addition to the nodes of the third level), while the optimal solution has size  $k + 2$ . However this algorithm is not totally useless. It may be shown that it always achieves the ratio  $\Theta(\lg n)^\dagger$ , where  $n$  is the number of vertices. That is, if  $C^*(n)$  is the size of the optimal solution and  $C(n)$  is the size of the solution found by the aforementioned algorithm. Then we have  $\frac{C(n)}{C^*(n)} = \Theta(\lg n)$  [9], p-323.

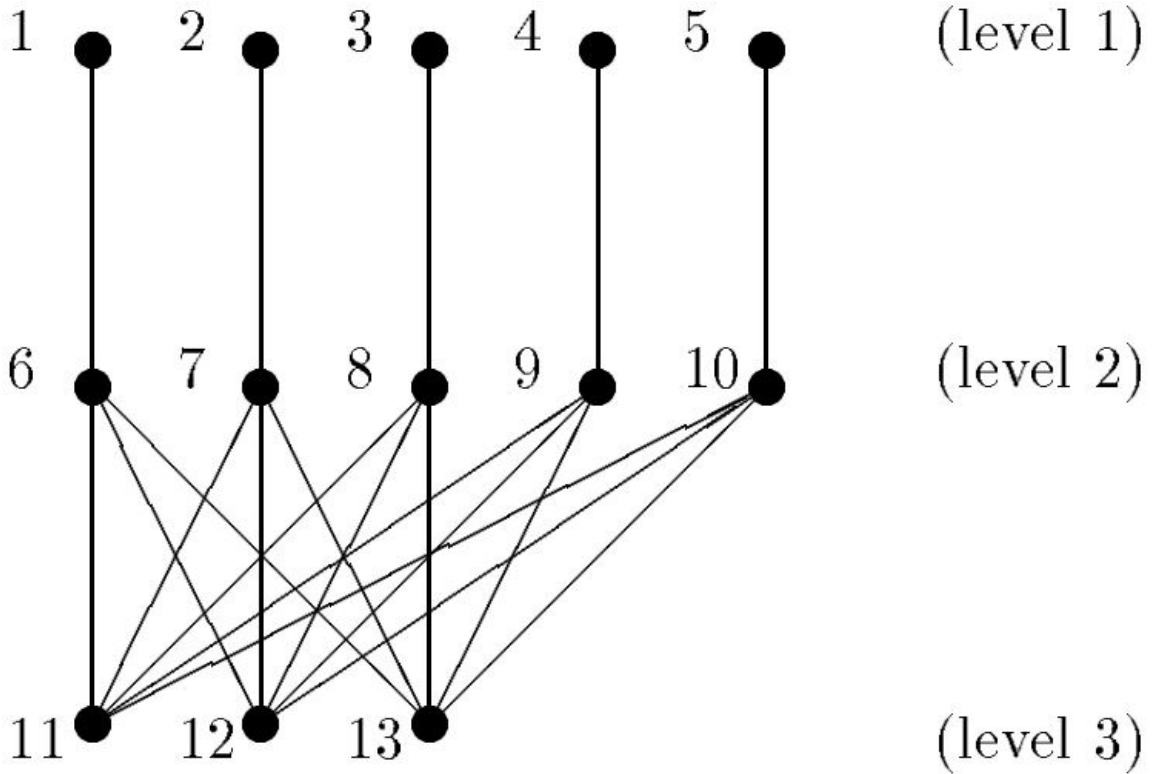


Figure 2: Construction of the regular graph after Papadimitriou and Steiglitz with  $k = 3$ .

Fortunately, a simpler and better algorithm exists. The following simple algorithm

---


$$^\ddagger f(n) = \Theta(g(n)) \iff \exists \alpha, \beta, n_0 : \forall n \geq n_0, \alpha \cdot g(n) \leq f(n) \leq \beta \cdot g(n)$$

is surprisingly the best approximation algorithm known for the mvcp ([9], p. 301).

---

**Algorithm 2**  $\text{vercov}(G = (V, E))$

---

```

1:  $C \leftarrow \{\}$   $\triangleright C$  contains the vertex cover being constructed
2:  $E' \leftarrow E$ 
3: while  $E' \neq \phi$  do
4:   randomly choose  $\langle u, v \rangle \in E'$ 
5:    $C \leftarrow C \cup \{u, v\}$ 
6:    $E' \leftarrow E' - \{\langle x, y \rangle \in E' : x = u \vee y = v \vee x = v \vee y = u\}$ 
7:    $\triangleright$  remove from  $E'$  every edge incident on either  $u$  or  $v$ 
8: end while
9: return  $C$ 

```

---

This is a *2-approximation algorithm*. Which means that the size of the vertex cover it returns is guaranteed to be no more than twice the size of an optimal vertex cover ([1], p. 968). More precisely, if  $C^*$  is the size of the optimal solution and  $C$  is the size of the solution returned by algorithm 2 then the inequality  $C \leq 2C^*$  must hold (notice that here the approximation doesn't depend on the number of vertices  $n$ ).

### 3.2 GA for the mvcp

In order to apply GA to solve the mvcp, a cover  $V'$  is represented by a binary vector  $\vec{x} = x_1 x_2 \dots x_n$  where  $x_i = 1$  if the  $i^{\text{th}}$  node is in  $V'$ , and  $x_i = 0$  if it is not. Using this representation, genetic algorithm is applied to minimize the following fitness function: [0]

$$\begin{aligned}
 f(\vec{x}) &= \sum_{i=1}^n \left( x_i + n \cdot (1 - x_i) \cdot \sum_{j=i}^n (1 - x_j) e_{ij} \right) \\
 &= \sum_{i=1}^n x_i + n \cdot \sum_{i=1}^n \sum_{j=i}^n (1 - x_i) \cdot (1 - x_j) \cdot e_{ij}
 \end{aligned} \tag{1}$$

The term  $\sum_{i=1}^n x_i$  of  $f(\vec{x})$  determines the size of the potential vertex cover represented by  $\vec{x}$ , while the term  $n \cdot \sum_{i=1}^n \sum_{j=i}^n (1 - x_i) \cdot (1 - x_j) \cdot e_{ij}$  penalizes sets  $V'$  that are not covers by adding a penalty of magnitude  $n$  for each edge  $e_{ij}$  for which  $i \notin V' \wedge j \notin V'$  (note that  $e_{ii}$  is previously defined to be zero). Consequently, any infeasible solution will have a fitness greater than or equal to  $n$ , and for any feasible solution the second term will drop to zero. This fitness function was developed according to the following

design principles that are important for a successful penalty function approach [11 || 13 || 7]:

- Infeasible binary vectors are guaranteed to yield fitness values which are inferior to fitness values of the worst feasible solutions.
- The penalty should be graded, i.e., fitness values should improve as solutions approach feasible regions of the search space. In other words, the more edges a solution covers, the fitter it is. notice that this is guaranteed by the second term in equation 1, where for each uncovered edge the fitness is penalized with the value of  $n$ .

The experiments reported in section 4 are performed by using a genetic algorithm with a population size of  $u = 50$ , a mutation rate  $p_m = 1/n$ , crossover rate  $p_c = 0.6$ , proportional selection, and two-point crossover. As reported in [2, 12], the latter is expected to perform better than the traditional one-point crossover. The termination criterion is specified in section 4.

## 4 Experimental Results

For the experimental tests, the genetic algorithm software package GENESYs is used [1]misp. This implementation is based on the widely used GENESIS software by Grefenstette (see [3], pp. 374-377 misp), but allows for more flexibility concerning genetic operators and data monitoring. The parameter settings for our experiments are given in section 3.2.

Multiple problem instances are used to demonstrate the applicability of GAs to tackle the mvcp. The results in table 3 and table 4 correspond to the first test set, which consists of five random graphs of size  $n = 100$  with different edge densities  $d$ ,  $d \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ , are used. An edge density of  $d = 0.1$  means that an edge is placed between two nodes with a probability of 0.1. Moreover we construct the graphs in a way that guarantees an optimum of quality 55 (we set  $k = 55$  in algorithm 3). In addition to the regular graphs introduced by Papadimitriou and Steiglitz and described in section 3.1. Recall that these graphs contain  $n = 3k + 4$  ( $k \geq 1$ ) nodes distributed on three levels. They can be scaled up by choosing high values for  $k$ . A graph instance of the regular graph of size  $n = 100$  ( $k = 32$ ) is chosen.

---

**Algorithm 3** Generate( $n, d, k$ )

---

```
1: randomly select  $V' = \{i_1, \dots, i_k\} \subseteq V = \{1, \dots, n\}$ 
2:    $\triangleright$  randomly select  $k$  different nodes, these nodes will hold a potential mvc
3: for  $i = 1$  to  $n - 1$  do
4:   for  $j = i + 1$  to  $n$  do
5:     if ( $\text{Random}(0, 1) \leq d$ )  $\wedge$  ( $i \in V' \vee j \in V'$ ) then
6:        $e_{ij} = 1$ 
7:     else
8:        $e_{ij} = 0$ 
```

---

For each of these graphs, a total of  $N = 100$  independent runs of the genetic algorithm is performed. The results are summarized in table 3 for the best fitness values that were encountered during the 100 runs. For each problem instance, the different fitness values that were obtained and their frequencies are recorded. Furthermore, the average fitness value  $\bar{f}$  over all 100 runs is indicated at the bottom of the table. For example, the first column in table 3 indicates that the best obtained fitness value in 1 of the 100 runs was  $f(\vec{x}) = 53$ , the best obtained fitness value in 1 of the 100 runs was  $f(\vec{x}) = 54$ , the best obtained fitness value in 3 of the 100 runs was  $f(\vec{x}) = 55$ , and so on... And the average of these values is indicated by  $\bar{f}$ , which is the average value of the best obtained fitness values over the 100 runs.

As a termination criterion for the genetic algorithm, a bound for the total number of function evaluations is set. For example The total number of function evaluations per single run is chosen to be  $2 \cdot 10^4$  in table 3. That is, the algorithm terminates when  $2 \cdot 10^4$  evaluations of the fitness function had occurred. This bound is indicated as an index  $t$  in the notation  $f_t(\vec{x})$ . Consequently, only a small fraction of the search space is explored by the genetic algorithm. This fraction is about  $2 \cdot 10^4 / 2^{100} \approx 1.6 \cdot 10^{-24}\%$  for the first test set. and it is about  $4 \cdot 10^4 / 2^{200} \approx 2.5 \cdot 10^{-54}\%$  for the second test set.

For each of the problems in the first test set, a total of  $N = 100$  independent runs of the vercov algorithm is performed. The results are summarized in table 4. Notice that for one problem instance, the results may differ in each run due to the non-determinism of the algorithm (see algorithm 2 line 4).

To test the behaviour of the genetic algorithm as well as the vercov algorithm for an even larger problem size, The same experiments were performed for larger graphs The results in table 5 and table 6 correspond to the second test set, which consists of five random graphs of size  $n = 200$  with different edge densities  $d$ ,  $d \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ , are used. Moreover we construct the graphs in a way that

mvcp100-01		mvcp100-02		mvcp100-03		mvcp100-04		mvcp100-05		PS100	
$f_{2 \cdot 10^4}(\vec{x})$	N	$f_{2 \cdot 10^4}(\vec{x})$	N	$f_{2 \cdot 10^4}(\vec{x})$	N	$f_{2 \cdot 10^4}(\vec{x})$	N	$f_{2 \cdot 10^4}(\vec{x})$	N	$f_{2 \cdot 10^4}(\vec{x})$	N
53	1	55	34	55	77	55	96	55	99	34	65
54	1	57	3	56	1	67	1	83	1	66	35
55	3	59	2	59	6	68	1				
56	6	60	3	63	3	75	1				
57	4	61	10	64	3	90	1				
58	4	62	1	66	1						
59	9	63	8	67	1						
60	4	64	1	68	1						
61	6	65	1	70	1						
62	12	66	6	71	1						
63	5	67	6	74	1						
> 63	45	> 67	25	> 74	4						
$\bar{f} = 62.61$		$\bar{f} = 62.75$		$\bar{f} = 57.62$		$\bar{f} = 55.80$		$\bar{f} = 55.28$		$\bar{f} = 45.20$	

Figure 3: general schema for genetic algorithms

mvcp100-01		mvcp100-02		mvcp100-03		mvcp100-04		mvcp100-05		PS100	
$f(\vec{x})$	N	$f(\vec{x})$	N	$f(\vec{x})$	N	$f(\vec{x})$	N	$f(\vec{x})$	N	$f(\vec{x})$	N
80	6	80		80		80	1	80		66	100
82	17	82	1	82		82		82	1		
84	23	84	2	84	1	84	1	84			
86	20	86	4	86	6	86		86	5		
88	17	88	18	88	12	88	5	88	7		
90	14	90	17	90	15	90	13	90	10		
92	1	92	31	92	18	92	24	92	21		
94	2	94	20	94	28	94	34	94	27		
96		96	4	96	16	96	19	96	17		
98		98	1	98	4	98	3	98	9		
100		100		100		100		100	3		
$\bar{f} = 86.46$		$\bar{f} = 89.22$		$\bar{f} = 92.22$		$\bar{f} = 92.96$		$\bar{f} = 93.12$			

Figure 4: general schema for genetic algorithms

guarantees an optimum of quality 110 (see algorithm 3). In addition to a graph instance of the regular graphs introduced by Papadimitriou and Steiglitz of size  $n = 220(k = 66)$ . In this case, the genetic algorithm was allowed to run for  $4 \cdot 10^4$  function evaluations. The corresponding results obtained by the genetic algorithm are shown in table 5, while the results from the vercov algorithm are presented in table 6. See the previous paragraphs for a better understanding of these results.

mvcp200-01		mvcp200-02		mvcp200-03		mvcp200-04		mvcp200-05		PS202	
$f_{4 \cdot 10^4}(\vec{x})$	N	$f_{4 \cdot 10^4}(\vec{x})$	N	$f_{4 \cdot 10^4}(\vec{x})$	N	$f_{4 \cdot 10^4}(\vec{x})$	N	$f_{4 \cdot 10^4}(\vec{x})$	N	$f_{4 \cdot 10^4}(\vec{x})$	N
110	2	110	55	110	95	110	99	110	100	68	60
113	1	120	10	128	1	140	1			134	40
116	4	121	7	129	6						
117	2	122	2	130	2						
119	3	123	7	134	1						
120	1	125	1								
121	2	126	1								
122	3	127	1								
124	3	129	1								
125	6	132	2								
126	2	134	1								
> 126	71	> 134	12								
$\bar{f} = 132.50$		$\bar{f} = 119.07$		$\bar{f} = 111.01$		$\bar{f} = 110.30$		$\bar{f} = 110.00$		$\bar{f} = 108.00$	

Figure 5: general schema for genetic algorithms

## 4.1 expermental observations

## 5 Conclusion

sit amet, consectetur adipisicing elit, sed do eiusmod

mvcp200-01		mvcp200-02		mvcp200-03		mvcp200-04		mvcp200-05		PS202	
$f(\vec{x})$	N	$f(\vec{x})$	N	$f(\vec{x})$	N	$f(\vec{x})$	N	$f(\vec{x})$	N	$f(\vec{x})$	N
172	2	172		172		172		172		134	100
174	3	174		174	1	174		174			
176	13	176	1	176	1	176	1	176	1		
178	14	178	1	178	3	178	4	178	1		
180	23	180	5	180	5	180	1	180	5		
182	14	182	11	182	5	182	9	182	5		
184	15	184	18	184	17	184	8	184	5		
186	7	186	16	186	19	186	16	186	13		
188	5	188	24	188	16	188	16	188	11		
190	3	190	15	190	16	190	16	190	21		
192	1	192	4	192	9	192	15	192	19		
> 192	0	> 192	5	> 192	8	> 192	14	> 192	19		
$\bar{f} = 180.98$		$\bar{f} = 186.46$		$\bar{f} = 186.92$		$\bar{f} = 188.06$		$\bar{f} = 189.16$			

Figure 6: general schema for genetic algorithms

## References

- [1] D. Goldberg. *Genetic Algorithms in Search, Optimization, and Machine Learning*. Artificial Intelligence. Addison-Wesley Publishing Company, 1989.