



ELSEVIER

Information Sciences 131 (2001) 47–64

INFORMATION  
SCIENCES  
AN INTERNATIONAL JOURNAL

[www.elsevier.com/locate/ins](http://www.elsevier.com/locate/ins)

# Automatic test data generation for path testing using GAs

Jin-Cherng Lin, Pu-Lin Yeh \*

*Department of Computer Science and Engineering, Tatung University, Taipei 10451, Taiwan, ROC*

Received 1 March 1999; received in revised form 23 May 2000; accepted 30 July 2000

---

## Abstract

Genetic algorithms (GAs) are inspired by Darwin's *the survival of the fittest* theory. This paper discusses a genetic algorithm that can automatically generate test cases to test a selected path. This algorithm takes a selected path as a target and executes sequences of operators iteratively for test cases to evolve. The evolved test case will lead the program execution to achieve the target path. To determine which test cases should survive to produce the next generation of fitter test cases, a metric named normalized extended Hamming distance (NEHD, which is used to determine whether the final test case is found) is developed. Based on NEHD, a fitness function named *SIMILARITY* is defined to determine which test cases should survive if the final test case has not been found. Even when there are loops in the target path, *SIMILARITY* can help the algorithm to lead the execution to flow along the target path. © 2001 Elsevier Science Inc. All rights reserved.

**Keywords:** Software testing; Path testing; Genetic algorithms

---

## 1. Introduction

The basic concepts of genetic algorithms (GAs) were developed by Holland [7,8]. To date, genetic algorithms are well developed and widely used in the area

---

\* Corresponding author. Tel.: +886-2-25925252 ext. 3295.

E-mail address: [plyeh@ms1.tisnet.net.tw](mailto:plyeh@ms1.tisnet.net.tw) (P.-L. Yeh).

of artificial intelligence and have been successfully applied in many fields of optimization [17–21].

It is understood that all occurrences of pre-selected features in a tested program should be executed (e.g. *all-branch testing* requires that each branch in a program should be tested at least once). It is also understood that manually generating a large number of test cases to fulfill the testing criteria is too great an effort. As a result, certain degree of automation is necessary. As Ould [14] suggested, the automation of test case generation is the most important aspect of automatic testing. Automatic test case generation can keep testers from bias when preparing test cases.

Many automatic test case generators have been developed and many classical searching methods have been used to derive test cases [1–4,9]. However, these techniques work only for continuous functions (while most program domains are of discontinuous spaces). This is where GAs step in. GAs include a class of adaptive searching techniques which are able to search a discontinuous space.

### 1.1. The works that apply genetic algorithms to program testing

Genetic algorithms have been used to automatically find a program's longest or shortest execution times. In their paper about testing real-time systems using genetic algorithms [16], Wegener et al. investigated the effectiveness of GAs to validate the temporal correctness of real-time systems by establishing the longest and the shortest execution times. The authors declared that an appropriate fitness function for the genetic algorithms is found, and the fitness function is to measure the execution time in processor cycles. Their experiments using GAs on a number of programs have successfully identified new longer and shorter execution times than those that had been found using random testing or systematic testing. However, they did not show the fitness function in their work.

Genetic algorithms have also been used to search program domains for suitable test cases to satisfy the all-branch testing [6,10,11]. In their papers about automatic structural testing using genetic algorithms, Jones et al. showed that appropriate fitness functions are derived automatically for each branch predicate. All branches were covered with two orders of magnitude fewer test cases than random testing. The adequacy of the test cases is improved by the genetic algorithm's ability to generate test cases, which are at or close to the branch predicates.

Consider a predicate such as

If  $A = B$  then ...

The fitness function for this predicate is based on the Hamming distance\* between  $A$  and  $B$ . The authors declare that the reciprocal of Hamming distance leads to a suitable fitness function.

\*The Hamming distance is a count based on the number of different bits in the bit patterns for  $A$  and  $B$  (i.e.  $A \oplus B$ ). Hamming distance between two bit patterns of equal length is defined as the number of different bits in the bit patterns. For example, the distance between 000 and 111 is three, while the distance between 010 and 011 is one.

The results of the papers above confirm that genetic algorithms are more efficient and effective than random testing in all-branch testing and finding longest execution time. In this paper, we developed a new metric (which is a fitness function) to determine the distance between the exercised path and the target path. The genetic algorithm with the metric is used to generate test cases for testing the target path. Such path-testing algorithm is designed to:

1. select a target path from the program under test,
2. sample test cases from the whole domain of the program,
3. execute the program with the test cases to exercise paths, and calculate the fitness between the target path and each exercised path, and then
4. evolve into new generation of fitter test cases,
5. repeat from the step 3 to the step 4 until the fittest test case is found.

## 2. Genetic algorithms

Genetic algorithms begin with a set of initial individuals as the first generation, which are sampled at random from the problem domain. The algorithms are developed to perform a series of operations that transform the present generation into a new, fitter generation.

Each individual in each generation is evaluated with a fitness function. Based on the evaluation, the evolution of the individuals may approach the optimal solution.

The most common operations of a genetic algorithm are designed to produce efficient solutions for the target problem [15]. These primary operations include:

(a) *Reproduction*: This operation assigns the reproduction probability to each individual based on the output of the fitness function. The individual with a higher ranking is given a greater probability for reproduction. As a result, the fitter individuals are allowed a better survival chance from one generation to the next.

(b) *Crossover*: This operation is used to produce the descendants that make up the next generation. This operation involves the following crossbreeding procedures:

- (i) Randomly select two individuals as a couple from the parent generation.
- (ii) Randomly select a position of the genes, corresponding to this couple, as the crossover point. Thus, each gene is divided into two parts.
- (iii) Exchange the first parts of both genes corresponding to the couple.

(iv) Add the two resulted individuals to the next generation.

(c) *Mutation*: This operation picks a gene at random and changing its state according to the mutation probability. The purpose of the mutation operation is to maintain the diversity in a generation to prevent premature convergence to a local optimal solution. The mutation probability is given intuitively since there is no definite way to determine the mutation probability.

Upon completion of crossover processing and mutation operations, there will be an original parent population and a new offspring population. A fitness function should be devised to determine which of these parents and offsprings can be survived into the next generation. After performing the fitness function, these parents, and offsprings are filtered and a new generation is formed. These operations are iterated until the expected goal is achieved. Genetic algorithms guarantee high probability of improving the quality of the individuals over several generations according to the Schema Theorem [5].

### 3. Apply genetic algorithms to path testing

Path testing is one kind of software testing methodologies which searches the program domain for suitable test cases such that after executing the program with the test cases, a predefined path can be reached. Since a program may contain an infinite number of paths, it is only practical to select a specific subset path to perform path testing. Genetic algorithms could be applied to path testing if the target paths are clearly defined and an appropriate fitness function related to this goal is built.

In this paper, the selected target path is the goal for the algorithm to approach, and a test case is regarded as an individual of a generation. Test case generation for path testing consists of four basic steps:

(1) *Control flow graph construction*: In this step, the source program is transferred to a graph that represents the control flow of the program (e.g. Fig. 3 is the control flow graph of the program in Fig. 2). A program's control flow graph can be constructed manually or automatically with exiting tools. Each branch of the graph is denoted by a label and different branches correspond with different labels. Each branch of the tested program is instrumented with a procedure or statements that can generate a specific label as an output. The output label indicates which specific branch is executed. The example of original program and its instrumented version are shown in Figs. 2 and 4, respectively. In Fig. 4, the statement *strcat(path, "e")* is one of the instrumented statements. The instrumented program will generate two sets of output, one by its original function and the other the instrumented statements. The second output is a label string that represents an executed path. When loops exist in the path, the instrumentation reflects the loops with repeating sub-strings.

(2) *Target path selection*: In path testing, some paths might be very meaningful and need to be selected for testing. For example, the path “abc” in Fig. 2 is the most difficult path to be covered in random testing. So that it is selected as the target path in this example to show the ability of genetic algorithm.

(3) *Test case generation and execution*: This step is used to create test cases. After each unwanted path is executed, the algorithm automatically generates new test cases to execute new paths and leads the control flow to the target path. Finally, a suitable test case that executes along the target paths could be generated.

(4) *Test result evaluation*: This step is to execute the selected path and to determine whether the testing criteria or test oracles are satisfied.

A genetic algorithm to find a test case which generates a given target path is depicted below:

```

Initialize test cases from the domain of the program to be tested at random;
Do
    feed the program with the test cases;
    if any of the test case has reached to the target path
        output the successful message;
        exit;
    endif
    Determine which test case should survive with fitness function;
    Reproduce the survivors;
    Select parents randomly from the survivors for crossover;
    Select the crossover sites of the parents;
    Produce the next new generation of test cases;
    Mutate the new generation of test cases according to the mutation probability;
    if iteration limit exceeded
        output a failure message;
        exit;
    endif
loop;

```

The first generation of test cases is generated at random. The generated cases are fed to the program for execution. One test case will be exercised in one and only one correlated path. Survivors of test cases to the next generation are chosen according to the fitness function, which is a measurement function used to calculate the distances between the output paths and the target path. Such

distances are used to determine which test cases should survive. After all the test cases in the present generation are executed, the new generation of test cases is generated by *reproduction*, *crossover* and *mutation*. The system will automatically generate the next generation of test cases until one of the test cases covers the target path.

When applying genetic to the algorithms to the automation of software testing, variable values exist as patterns of bits. When encoding the individual, although the program inputs may be of different types and of complicated data structure, these inputs can be treated as a single, concatenated bit string denoted as  $\{b_1, b_2, \dots, b_n\}$ . For example, the parameters of the function *funl(float f, int m)* can be treated as a 48-bit (6 byte) string. In this string, *f* occupies 32 bits (4 bytes) and *m* occupies 16 bits (2 bytes). Each input bit string is referred to as an individual of a generation, an individual is subdivided into genes, and any bit or sub-string in the string can be regarded as a gene. Normally, genetic algorithms operate on the bit strings or an individual of fixed size. The mutation is an operation in which a bit within a bit string is chosen randomly and then flipped from 1 to 0 and from 0 to 1.

### 3.1. To develop the fitness function

In branch testing, Hamming distance has been used to measure the difference between the covered branches and the selected branches as the fitness [11]. This distance metric can only be used to measure the distance of two objects in which they have no specific sequences. For the tested elements in all-nodes or all-branches testing criterion, no tested sequence is needed. However for path testing, two different path may contain the same branches but in different sequences. The simple Hamming distance is no longer suitable.

One contribution of this paper is to extend the Hamming distance from the first-order to the  $n$ th order ( $n > 1$ ) to measure the distance between two path. Such extension is hereby named extended Hamming distance (EHD).

In the paragraph below, a fitness function is developed to measure the distance between two paths.

Hamming distance is derived from the *symmetric difference* in set theory. The *symmetric difference* of the set  $\alpha$  and the set  $\beta$  (denoted as  $\alpha \oplus \beta$ ) is a set containing the elements either in  $\alpha$  or  $\beta$  but not in both. In other words,  $\alpha \oplus \beta$  equals to  $(\alpha \cup \beta) - (\alpha \cap \beta)$ . In this paper, the cardinality of the *symmetric difference* is defined as the *distance* between  $\alpha$  and  $\beta$

$$D_{\alpha-\beta} = |\alpha \oplus \beta|.$$

The notation,  $|\alpha \oplus \beta|$ , denotes the cardinality of  $\alpha \oplus \beta$ .

The distance  $D_{\alpha-\beta}$  is normalized to become a real number

$$N_{\alpha-\beta} = \frac{D_{\alpha-\beta}}{|\alpha \cup \beta|},$$

where  $N_{\alpha-\beta}$  represents the *normalized distance* between sets  $\alpha$  and  $\beta$ , and

$$0 \leq N_{\alpha-\beta} \leq 1,$$

because  $D_{\alpha-\beta} = |\alpha \oplus \beta| = |\alpha \cup \beta| - |\alpha \cap \beta| \leq |\alpha \cup \beta|$ .

Replace  $D_{\alpha-\beta}$  with  $|\alpha \cup \beta| - |\alpha \cap \beta|$  to derive

$$N_{\alpha-\beta} = \frac{|\alpha \cup \beta| - |\alpha \cap \beta|}{|\alpha \cup \beta|} = 1 - \frac{|\alpha \cap \beta|}{|\alpha \cup \beta|}$$

$$\Rightarrow (1 - N_{\alpha-\beta}) = \frac{|\alpha \cap \beta|}{|\alpha \cup \beta|} = M_{\alpha-\beta},$$

where  $M_{\alpha-\beta} = (1 - \text{the normalized distance between } \alpha \text{ and } \beta)$ . The function  $M_{\alpha-\beta}$  is named *SIMILARITY*, and is used to measure the similarity between  $\alpha$  and  $\beta$ .

In the following section,  $M_{\alpha-\beta}$  is used to measure the similarity between two paths in a program control-flow diagram.

A *complete path* is a path whose initial node is the start node and whose final node is an exit node. If  $P$  is a control-flow diagram of a given program and  $Q$  is the set of all complete paths within  $P$ , then

$$\begin{aligned} Q &= \{\text{path}_i | \text{path}_i \text{ is a complete path within } P\} \\ &= \{\text{path}_1, \text{path}_2, \dots, \text{path}_z\}, \end{aligned}$$

where  $z$  = the number of complete paths in  $P$ ,

$$\text{path}_i = \text{the } i\text{th complete path in } P, \quad 1 \leq i \leq z.$$

Let  $S_1^1 = \{g | g \text{ is a branch of } \text{path}_i\}$ ,

$$S_1^2 = \{h | h \text{ is an ordered pair of cascaded branches of } \text{path}_1\},$$

.....

$$\begin{aligned} S_1^n &= \{k | k \text{ is an ordered } n\text{-tuple of cascaded} \\ &\quad \text{branches of } \text{path}_1, \quad n \leq |S_1^1|\}, \end{aligned}$$

.....

$$\begin{aligned} S_q^t &= \{r | r \text{ is an ordered } t\text{-tuple of cascaded branches} \\ &\quad \text{of } \text{path}_q, \quad t \leq |S_q^1| \text{ and } 1 \leq q \leq z\}, \end{aligned}$$

.....

The first-order *distance* between  $\mathbf{path}_i$  and  $\mathbf{path}_j$  is expressed as

$$D_{i-j}^1 = |S_i^1 \oplus S_j^1|.$$

The normalized first-order *distance* between  $\mathbf{path}_i$  and  $\mathbf{path}_j$  is expressed as

$$N_{i-j}^1 = \frac{D_{i-j}^1}{|S_i^1 \cup S_j^1|}.$$

The first-order *similarity* between  $\mathbf{path}_i$  and  $\mathbf{path}_j$  is defined as

$$M_{i-j}^1 = 1 - N_{i-j}^1.$$

The second-order *distance* between  $\mathbf{path}_i$  and  $\mathbf{path}_j$  is expressed as

$$D_{i-j}^2 = |S_i^2 \oplus S_j^2|.$$

The normalized second-order *distance* between  $\mathbf{path}_i$  and  $\mathbf{path}_j$  is expressed as

$$N_{i-j}^2 = \frac{D_{i-j}^2}{|S_i^2 \cup S_j^2|}.$$

The second-order *similarity* between  $\mathbf{path}_i$  and  $\mathbf{path}_j$  is defined as

$$M_{i-j}^2 = 1 - N_{i-j}^2.$$

The  $m$ th ( $m = 1..n$ ) order *distance* between  $\mathbf{path}_i$  and  $\mathbf{path}_j$  is expressed as

$$D_{i-j}^m = |S_i^m \oplus S_j^m|.$$

The normalized  $m$ th order *distance* between  $\mathbf{path}_i$  and  $\mathbf{path}_j$  is expressed as

$$N_{i-j}^m = \frac{D_{i-j}^m}{|S_i^m \cup S_j^m|}.$$

The  $m$ th order *similarity* between  $\mathbf{path}_i$  and  $\mathbf{path}_j$  is defined as

$$M_{i-j}^m = 1 - N_{i-j}^m.$$

The notation  $D_{i-j}^m$  is the  $m$ th order EHD between  $\mathbf{path}_i$  and  $\mathbf{path}_j$ . The notation  $N_{i-j}^m$  is the  $m$ th order normalized extended Hamming distance (NEHD) between  $\mathbf{path}_i$  and  $\mathbf{path}_j$ . The notation  $M_{i-j}^m$  is named the  $m$ th order *SIMILARITY* between  $\mathbf{path}_i$  and  $\mathbf{path}_j$  where  $1 \leq m \leq n$ . Note that  $M_{i-j}^m = 0$  (or  $N_{i-j}^m = 1$ ) if  $S_i^m \cap S_j^m = \phi$ . It means that  $\mathbf{path}_i$  and  $\mathbf{path}_j$  have no common  $m$ -tuple cascaded breaches. Larger  $N_{i-j}^m$  means greater *difference* between  $\mathbf{path}_j$  and  $\mathbf{path}_i$ . Contrarily, larger  $M_{i-j}^m$  means greater *similarity* between  $\mathbf{path}_j$  and  $\mathbf{path}_i$ . When  $\mathbf{path}_j$  and  $\mathbf{path}_i$  have no common branch, NEHD should take the form of  $\{N_{i-j}^1 = 1, N_{i-j}^2 = 1, \dots, N_{i-j}^n = 1\}$  or  $\{M_{i-j}^1 = 0, M_{i-j}^2 = 0, \dots, M_{i-j}^n = 0\}$ , which is resulted from a worst test case. When  $\mathbf{path}_j$  and  $\mathbf{path}_i$  are identical, NEHD should take the form of  $\{N_{i-j}^1 = 0, N_{i-j}^2 = 0, \dots, N_{i-j}^n = 0\}$  or  $\{M_{i-j}^1 = 1, M_{i-j}^2 = 1, \dots, M_{i-j}^n = 1\}$ , which is resulted from a perfect test case that forces the program to execute along the target path. Therefore, if NEHD is not in the



condition of  $\{N_{i-j}^1 = 0, N_{i-j}^2 = 0, \dots, N_{i-j}^n = 0\}$ , the fitness function *SIMILARITY* between  $\mathbf{path}_j$  and  $\mathbf{path}_i$  is defined as

$$\text{SIMILARITY}_{i-j} = M_{i-j}^1 \times W_1 + M_{i-j}^2 \times W_2 + M_{i-j}^3 \times W_3 + \dots + M_{i-j}^n \times W_n,$$

where  $W$ 's are the weighing factor of fitness and  $W_1 < W_2 < W_3 < \dots < W_n$ .

$n = |S_i^1|$  if  $\mathbf{path}_i$  is defined as the target path.

$n = |S_j^1|$  if  $\mathbf{path}_j$  is otherwise. Since, if the target path is constructed by  $n$  branched, the values of  $M_{i-j}^{n+1}, M_{i-j}^{n+2}, \dots$  should be zero and are insignificant in *similarity* comparison.

$\text{SIMILARITY}_{i-j}$  determines the fitness between current executed  $\mathbf{path}_j$  and the target  $\mathbf{path}_i$ . The greater  $\text{SIMILARITY}_{i-j}$  leads to the better fitness. The higher-order *similarity* is more significant than its lower-order counterpart. The highest-ordered *similarity* between  $\mathbf{path}_i$  and  $\mathbf{path}_j(M_{i-j}^n)$  is therefore the most significant one. The semantics of the tested program has much effect on the values of  $W$ 's. Determining the values of  $W$ 's is quite difficult and is usually done via experience.  $W_{k+1} = y * W_k$  means the  $(k+1)$ th order *similarity* is  $y$  times more significant than that of the  $k$ th order.

Let the least significant weight ( $W_1$ ) be 1. In experience, the  $W$ 's for  $\mathbf{path}_i$  may be assigned as

$$\begin{aligned} W_1 &= 1, \\ W_2 &= W_1 \times |S_i^1|, \\ W_3 &= W_2 \times |S_i^2|, \\ &\dots\dots \\ W_n &= W_{n-1} \times |S_i^{n-1}|. \end{aligned}$$

The distance between  $\mathbf{path}_i$  (the target path) and  $\mathbf{path}_j$  is larger than the distance between  $\mathbf{path}_i$  and  $\mathbf{path}_k$  if  $\text{SIMILARITY}_{i-k} > \text{SIMILARITY}_{i-j}$ . Therefore,  $\mathbf{path}_k$  is closer to the target than  $\mathbf{path}_j$  is. The *SIMILARITY* function can help the algorithm to search the program domain and to find fitter test cases. Even when there are loops in the target path, the function can help the algorithm to lead the execution to flow along the loops of the target path.

### 3.2. Calculating the distance similarity

An example of calculating the distance as well as the similarity between the target path and two executed paths is illustrated below.

In Fig. 1, each branch of the control flow graph is denoted with a label. A sequence of branches that construct a path can be denoted with a string of labels. Table 1 illustrates a minimal set of paths that satisfies the all-path

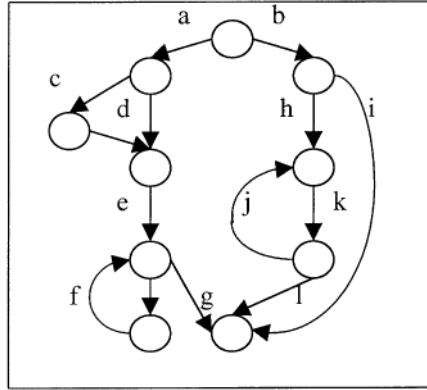


Fig. 1. An example of control flow diagram.

Table 1

The all-path list of the control flow diagram in Fig. 1

---

$\mathbf{path}_1 = \langle b, i \rangle$
$\mathbf{path}_2 = \langle b, h, k, l \rangle$
$\mathbf{path}_{3-x} = \langle b, h, \langle k, j \rangle^x, l \rangle$
$\mathbf{path}_4 = \langle a, d, e, g \rangle$
$\mathbf{path}_5 = \langle a, c, e, g \rangle$
$\mathbf{path}_{6-x} = \langle a, d, e, \langle f \rangle^x, g \rangle$
$\mathbf{path}_{7-x} = \langle a, c, e, \langle f \rangle^x, g \rangle$

---

strategy in Fig. 1. The number  $x$  in the superscript represents the number of loops. In the following illustrations,  $\mathbf{path}_{7-2}$  is selected as the target to test. The path  $\mathbf{path}_{7-2}$  has two loops in branch ' $f$ ' and is one of  $\mathbf{path}_{7-x}$ . To determine which path is more similar to the target path, two paths ( $\mathbf{path}_5, \mathbf{path}_4$ ) are selected for demonstration.

### 3.2.1. The distance and the similarity between $\mathbf{path}_{7-2}$ and $\mathbf{path}_5$

The following is a calculation of the distance and the similarity between  $\mathbf{path}_{7-2}$  and  $\mathbf{path}_5$ . If  $\mathbf{path}_{7-2}$  is the target path (where  $\mathbf{path}_{7-2} = \langle a, c, e, \langle f \rangle^2, g \rangle, |S_{7-2}^1| = 5$ ), then the distance between  $\mathbf{path}_{7-2}$  and  $\mathbf{path}_5 (= \langle a, c, e, g \rangle)$  is illustrated below.

(1) The calculation of  $D_{5-7-2}^1$ ,  $N_{5-7-2}^1$  and  $M_{5-7-2}^1$

Since

(a) the set of distinct branches in  $\mathbf{path}_5$  is  $S_5^1$  (where  $S_5^1 = \{a, c, e, g\}$ ),

(b) the set of distinct branches in both  $\mathbf{path}_{7-2}$  and  $\mathbf{path}_5$  is  $S_5^1 \cup S_{7-2}^1$  (where  $S_5^1 \cup S_{7-2}^1 = \{a, c, e, f, g\}$ , and

(c) the set of the first-order symmetric difference between  $\mathbf{path}_{7-2}$  and  $\mathbf{path}_5$  is  $S_5^1 \oplus S_{7-2}^1$  (where  $S_5^1 \oplus S_{7-2}^1 = \{f\}$ , thus they lead to

$$D_{5-7,2}^1 = |S_5^1 \oplus S_{7,2}^1| = 1,$$

$$N_{5-7,2}^1 = \frac{D_{5-7,2}^1}{|S_5^1 \cup S_{7,2}^1|} = 1/5, \text{ and}$$

$$M_{5-7,2}^1 = 1 - N_{5-7,2}^1 = 4/5.$$

(2) *The calculation of  $D_{5-7,2}^2$ ,  $N_{5-7,2}^2$  and  $M_{5-7,2}^2$*

Since

(a) the set of distinct ordered pairs of consecutive branches in **path**<sub>7,2</sub> and **path**<sub>5</sub> are  $S_{7,2}^2$  and  $S_5^2$  (where  $S_{7,2}^2 = \{\langle a, c \rangle, \langle c, e \rangle, \langle e, f \rangle, \langle f, f \rangle, \langle f, g \rangle\}$ ,  $|S_{7,2}^2| = 5$ ) and  $S_5^2 = \{\langle a, c \rangle, \langle c, e \rangle, \langle e, g \rangle\}$ ,

(b) the set of distinct ordered pairs of consecutive branches in both **path**<sub>7,2</sub> and **path**<sub>5</sub> is  $S_5^2 \cup S_{7,2}^2$  (where  $S_5^2 \cup S_{7,2}^2 = \{\langle a, c \rangle, \langle c, e \rangle, \langle e, f \rangle, \langle f, f \rangle, \langle f, g \rangle, \langle e, g \rangle\}$ ) and

(c) the set of the second-order *symmetric difference* between **path**<sub>7,2</sub> and **path**<sub>5</sub> is  $S_5^2 \oplus S_{7,2}^2$  (where  $S_5^2 \oplus S_{7,2}^2 = \{\langle e, f \rangle, \langle f, f \rangle, \langle f, g \rangle, \langle e, g \rangle\}$ ), thus they lead to

$$D_{5-7,2}^2 = |S_5^2 \oplus S_{7,2}^2| = 4,$$

$$N_{5-7,2}^2 = \frac{D_{5-7,2}^2}{|S_5^2 \cup S_{7,2}^2|} = 4/6, \text{ and}$$

$$M_{5-7,2}^2 = 1 - N_{5-7,2}^2 = 2/6.$$

(3) *The calculation of  $D_{5-7,2}^3$ ,  $N_{5-7,2}^3$  and  $M_{5-7,2}^3$*

Since

(a) the set of distinct ordered 3\_ tuple of consecutive branches in **path**<sub>7,2</sub> and **path**<sub>5</sub> are  $S_{7,2}^3$  and  $S_5^3$  (where  $S_{7,2}^3 = \{\langle a, c, e \rangle, \langle c, e, f \rangle, \langle e, f, f \rangle, \langle f, f, g \rangle\}$ ,  $|S_{7,2}^3| = 4$ , and  $S_5^3 = \{\langle a, c, e \rangle, \langle c, e, g \rangle\}$ ),

(b) the set of distinct ordered 3\_ tuple of consecutive branches in both **path**<sub>7,2</sub> and **path**<sub>5</sub> is  $S_5^3 \cup S_{7,2}^3$  (where  $S_5^3 \cup S_{7,2}^3 = \{\langle a, c, e \rangle, \langle c, e, f \rangle, \langle e, f, f \rangle, \langle f, f, g \rangle, \langle c, e, g \rangle\}$ ), and

(c) the set of the third-order *symmetric difference* between **path**<sub>7,2</sub> and **path**<sub>5</sub> is  $S_5^3 \oplus S_{7,2}^3$  (where  $S_5^3 \oplus S_{7,2}^3 = \{\langle c, e, f \rangle, \langle e, f, f \rangle, \langle f, f, g \rangle, \langle c, e, g \rangle\}$ ), thus they lead to,

$$D_{5-7,2}^3 = |S_5^3 \oplus S_{7,2}^3| = 4,$$

$$N_{5-7,2}^3 = \frac{D_{5-7,2}^3}{|S_5^3 \cup S_{7,2}^3|} = 4/5, \text{ and}$$

$$M_{5-7,2}^3 = 1 - N_{5-7,2}^3 = 1/5.$$

(4) The calculation of  $D_{5-7,2}^4$ ,  $N_{5-7,2}^4$  and  $M_{5-7,2}^4$

Since

(a) the set of distinctly ordered 4\_tuple of consecutive branches in  $\mathbf{path}_{7,2}$  and  $\mathbf{path}_5$  are  $S_{7,2}^4$  and  $S_5^4$  (where  $S_{7,2}^4 = \{\langle a, c, e, f \rangle, \langle c, e, f, f \rangle, \langle e, f, f, g \rangle\}$ , and  $S_5^4 = \{\langle a, d, f, i \rangle\}$ ),

(b) the set of distinct ordered 4\_tuple of consecutive branches in both  $\mathbf{path}_{7,2}$  and  $\mathbf{path}_5$  is  $S_5^4 \cup S_{7,2}^4$  (where  $S_5^4 \cup S_{7,2}^4 = \{\langle a, c, e, f \rangle, \langle c, e, f, f \rangle, \langle e, f, f, g \rangle, \langle a, d, f, i \rangle\}$ )

(c) the set of the fourth-order symmetric difference between  $\mathbf{path}_{7,2}$  and  $\mathbf{path}_5$  is  $S_5^4 \oplus S_{7,2}^4$  (where  $S_5^4 \oplus S_{7,2}^4 = \{\langle a, c, e, f \rangle, \langle c, e, f, f \rangle, \langle e, f, f, g \rangle, \langle a, d, f, i \rangle\}$ ), thus they lead to

$$D_{5-7,2}^4 = |S_5^4 \oplus S_{7,2}^4| = 4,$$

$$N_{5-7,2}^4 = \frac{D_{5-7,2}^4}{|S_5^4 \cup S_{7,2}^4|} = 1, \text{ and}$$

$$M_{5-7,2}^4 = 1 - N_{5-7,2}^4 = 0.$$

Since  $M_{5-7,2}^4 = 0$ , the  $D_{5-7,2}^5$ ,  $N_{5-7,2}^5$ ,  $M_{5-7,2}^5$  and the higher-order of distances need not to be figured out.

The NEHD between  $\mathbf{path}_5$  and  $\mathbf{path}_{7,2}$  is  $\text{NEHD}_{5-7,2}$ , i.e.  $\{1/5, 4/6, 4/5, 1\}$ . Since the fourth-order NEHD between  $\mathbf{path}_{7,2}$  and  $\mathbf{path}_5$  is 1, the  $w$ th order (where  $w \geq 4$ ). NEHD between  $\mathbf{path}_{7,2}$  and  $\mathbf{path}_5$  must also be 1. Therefore,  $\text{NEHD}_{5-7,2} = \{1/5, 4/6, 4/5, 1\} = \{1/5, 4/6, 4/5, 1, 1, 1, \dots\}$ .

### 3.2.2. The distance and similarity between $\mathbf{path}_{7,2}$ and $\mathbf{path}_4$

Using the same method as in Section 3.2.1, the values of  $D_{4-7,2}^1$ ,  $N_{4-7,2}^1$ ,  $M_{4-7,2}^1$ ,  $D_{4-7,2}^2$ ,  $N_{4-7,2}^2$  and  $M_{4-7,2}^2$  lead to

$$D_{4-7,2}^1 = 3,$$

$$N_{4-7,2}^1 = 3/6,$$

$$M_{4-7,2}^1 = 3/6,$$

$$D_{4-7,2}^2 = 8,$$

$$N_{4-7,2}^2 = 8/8 = 1, \text{ and}$$

$$M_{4-7,2}^2 = 0.$$

Since the second-order distance between  $\mathbf{path}_{7,2}$  and  $\mathbf{path}_4$  is 1, the  $w$ th order (where  $w > 2$ ) distance between  $\mathbf{path}_{7,2}$  and  $\mathbf{path}_4$  is also 1. The NEHD between  $\mathbf{path}_4$  and  $\mathbf{path}_{7,2}$  is  $\text{NEHD}_{4-7,2}$  (where  $\text{NEHD}_{4-7,2} = \{3/6, 1\} = \{3/6, 1, 1, 1, 1\} = \{3/6, 1, 1, 1, 1, \dots\}$ ).

To compare  $\text{NEHD}_{4-7,2} (= \{3/6, 1, 1, 1, 1\})$  with  $\text{NEHD}_{5-7,2} (= \{1/6, 4/6, 4/5, 1\})$ , it is clear that the distance between  $\mathbf{path}_4$  and  $\mathbf{path}_{7,2}$  is greater than the distance between  $\mathbf{path}_5$  and  $\mathbf{path}_{7,2}$ , because the higher the order is, the

more significant the weighting is (where  $N_{4-7,2}^4 = N_{5-7,2}^4, N_{4-7,2}^3 > N_{5-7,2}^3, N_{4-7,2}^2 > N_{5-7,2}^2$ , and  $N_{4-7,2}^1 > N_{5-7,2}^1$ ).

At the same time, the weighing factors are obtained below:

$$\begin{aligned} W_1 &= 1, \\ W_2 &= W_1 \times |S_{7,2}^1| = 1 \times 5 = 5, \\ W_3 &= W_2 \times |S_{7,2}^2| = 5 \times 5 = 25, \\ W_4 &= W_3 \times |S_{7,2}^3| = 25 \times 4 = 100. \end{aligned}$$

The fitness values of  $path_4$  and  $path_5$  can be derived below:

$$\begin{aligned} SIMILARITY_{4-7,2} &= M_{4-7,2}^1 \times W_1 + M_{4-7,2}^2 \times W_2 = (1 - 3/6) \times 1 \\ &\quad + (1 - 1) \times 5 = 0.5, \\ SIMILARITY_{5-7,2} &= M_{5-7,2}^1 \times W_1 + M_{5-7,2}^2 \times W_2 + M_{5-7,2}^3 \times W_3 + M_{5-7,2}^4 \times W_4 \\ &= (1 - 1/6) \times 1 + (1 - 4/6) \times 5 + (1 - 4/5) \times 25 \\ &\quad + (1 - 1) \times 100 = 7.5, \end{aligned}$$

which shows that  $path_5$  is closer to the target path ( $path_{7,2}$ ) than  $path_4$  is or that  $path_5$  is more similar to the target path ( $path_{7,2}$ ) than  $path_4$  is.

### 3.3. The experiment results

In the general path testing experiment, the four basic steps are processed below.

(1) *Control flow graph construction*: The tested program (Fig. 2 Triangle classifier) determines what kind of triangle can be formed by any three input lengths. The program's control flow diagram, which contains four paths, is shown in Fig. 3. After the instrumentation, the tested program is transformed

```

int TriangleA(int a, int b, int c)
{
    int Triangle = 0;                                /*b1*/
    int c1,c2,c3,c4,c5,c6,c7,c8;                     /*b1*/
    if (((a+b>c)&&(b+c>a)&&(c+a>b))&&((a>0)&&(b>0)&&(c>0))) /*b1*/
    {
        if ((a != b) && (b != c) && (c != a))           /*b2*/
            Triangle = 1;                                /* Scalene b3*/
        else
        {
            if (((a == b) && (b == c)) || ((b == c) && (c == a)) || ((c == a) && (a == b))) /*b4*/
                Triangle = 2;                            /* Isosceles b5*/
            else                                           /*b6*/
                Triangle = 3;                            /* Equilateral b6*/
        }
    }
    return ( Triangle );                                /*b7*/
}

```

Fig. 2. Triangle classifier [12,13].

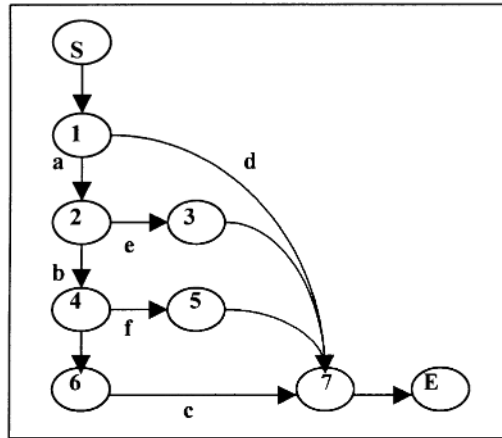


Fig. 3. The control flow diagram of the program shown in Fig. 2.

to the one shown in Fig. 4. In the implemented program, a set of labels is implemented to indicate when the path is executed. Therefore, executing the implemented program under test with each test case can produce a string of labels for the fitness function.

The comments  $/ * b_1 * / \dots / * b_7 * /$  denote the program blocks 1 to 7, respectively. In Fig. 3, the nodes 1 to 7 indicate the control flow locations of these blocks.

(2) *Target path selection*: The path “abc” in Fig. 2 is the most difficult path to be covered in random testing. The program in Fig. 2 has three integers as input parameters. While the three parameters are positive and equal, the path “abc” will be covered. The covered probability of this path is  $2^{-30} (= 1 * 2^{-15} * 2^{-15})$  for each positive integer is 15 bits). To show the ability of searching test cases for specific paths by using genetic algorithm is much greater than by using random testing, the path “abc” is selected as the target path in this example.

(3) *Test case generation and execution*: According to the genetic algorithms, an experimental tool for automatically generating test cases to test a specific path is developed. The tool divided total 48-bit input string into three genes using the conjunctions between two integers as the fixed crossover site. Pairs of test cases were combined using two-point crossover algorithm. Traditionally, the mutation probability is set to the reciprocal of the length of the bit string [11]. Hence, the mutation probability is set to 1/48.

In this experiment, the first generation of test cases was chosen by the tool from the tested program domain randomly. Then the tested program was executed with these test cases. The executed results are evaluated by fitness function to determine which test cases should survive to generate the next

```

#include <stdio.h>
#include <string.h>
char path[256];
int TriangleA(unsigned int a, unsigned int b, unsigned int c)
{
    int Triangle = 0;
    if ((a + b > c) && (b + c > a) && (c + a > b))
    {
        strcat(path, "a");
        if ((a != b) && (b != c) && (c != a))
        {
            strcat(path, "e");
            Triangle = 1;
        }
    }
    else
    {
        strcat(path, "b");
        if (((a == b) && (b != c)) || ((b == c) && (c != a)) || ((c == a) && (a != b)))
        {
            strcat(path, "f");
            Triangle = 2;
        }
    }
    else
    {
        strcat(path, "c");
        Triangle = 3;
    }
}
else
{
    strcat(path, "d");
}
return (Triangle);
}

```

Fig. 4. The example of the program in Fig. 2 after being instrumented.

generation. Again, new generations of test cases are generated by *reproduction*, *crossover* and *mutation*. The average fitness of each generation is steadily improved until the target path is achieved. In the initial generations, the execution of generated test cases was mostly group in the path  $\langle d \rangle$ . In the subsequent generations, the execution of generated test cases gathered in the path  $\langle a, e \rangle$ . According to the experiments, about 52.5% and 47.5% of the first generation test cases, which were generated at random, executed the path  $\langle d \rangle$  and the path  $\langle a, e \rangle$ , respectively, and none of the test cases executed the other paths. After the fifth generation, all the evolution of the test cases left the path  $\langle d \rangle$  and mostly gathered in the path  $\langle a, e \rangle$ . Afterward, the execution of generated test cases approached the path  $\langle abf \rangle$  gradually. Finally, at least a test case reached the path  $\langle abc \rangle$  and succeeded in generating the test case. After one hundred experiments, the results on Fig. 5 show that the target path was obtained within 10 100 test cases (i.e. 10 generations  $\times$  1000 test cases in each generation + 100 test cases in the 11th generation) by average. While, based on the theory of probability, it will take random testing  $2^{30}$  tests to reach the target.

Fig. 6 shows the evolution of test cases from the first generation to the 10th generation. These values were averaged from 100 experiments. Before the fifth

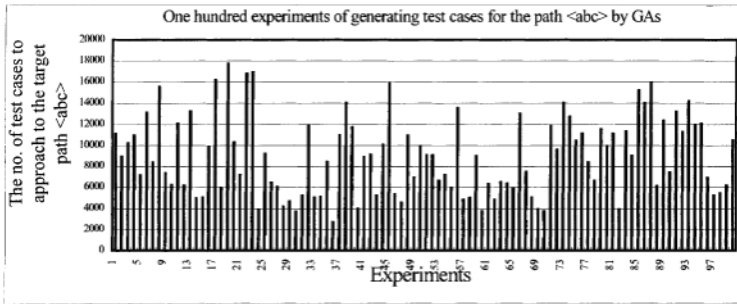


Fig. 5. On hundred experiments of generating test cases for executing the target path  $\langle abc \rangle$  in Fig. 4 by using GAs.

generation, the algorithm decreased the number of test cases on the path  $\langle d \rangle$  and increased the number of test cases on the path  $\langle ae \rangle$ . After the fifth generation, the number of test cases on the path  $\langle ae \rangle$  was decreased and the number of test cases on the path  $\langle abf \rangle$  increased speedily.

(4) *Test result evaluation*: This step is to execute the selected path with the test cases found in step (3) and to determine whether the outputs are correct or not.

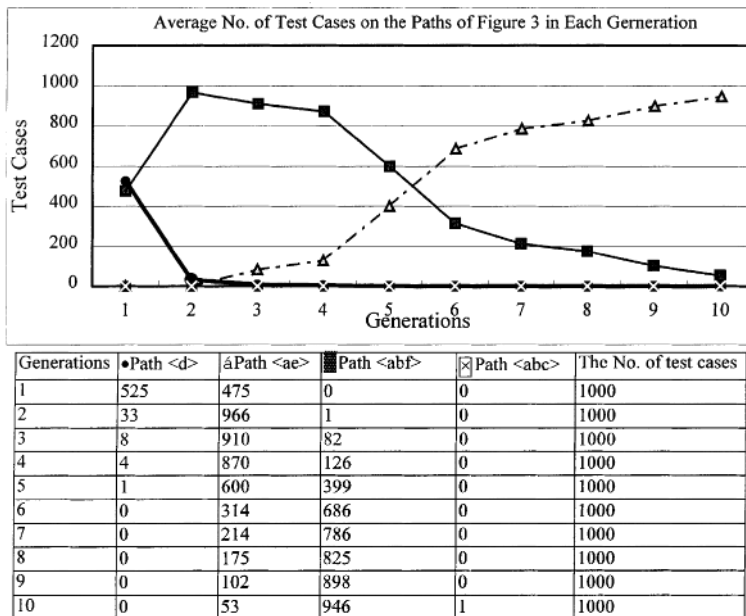


Fig. 6. The average number of test cases on the paths of Fig. 3 in each generation.



Since there is no guarantee that the algorithm can find the target within a definite number of runs, the execution of the algorithm was allowed to continue for 450 generations before it was stopped. Fortunately, the evolutions of the test cases to execute the target path were less than 18 generations in our experiments. We have also applied 100 test cases and 10 000 test cases in each generation in the same experiment. It is found that the best result is to apply 1000 test cases in one generation. As a summary, the number of individuals of one generation should be large enough to maintain diversity, yet small enough to avoid an excessive number of tests.

#### 4. Conclusion

In this paper, the genetic algorithms are used to automatically generate test cases for path testing. The greatest merit of using the genetic algorithm in program testing is its simplicity. Each iteration of the genetic algorithms generates a generation of individuals. In practice, the computation time cannot be infinite, so that the iterations in the algorithm should be limited. Within the limited generations, solutions derived by genetic algorithms may be trapped around a local optimum and, as a result, fail to locate the required global optimum. Although the tested cases generated by such algorithms may be trapped around unwanted paths and fail to locate the required paths, since the test cases of the first generation are normally distributed over the domain of the tested program, the probability of being trapped is very low.

The quality of test cases produced by genetic algorithms is higher than the quality of test cases produced by random way because the algorithm can direct the generation of test cases to the desirable range fast. Comparing with random testing, using extended Hamming distance to derive *SIMILARITY* (a fitness function) is a very useful approach for path testing. Although a program may contain infeasible paths where no method, including genetic algorithms, can generate suitable test cases, this paper shows that genetic algorithms are able to meaningfully reduce time required for lengthy testing by generating test cases for path testing in an automatic way.

#### References

- [1] J.F. Benders, Partitioning procedures for solving mixed variables programming problems, *Numer. Math.* 4 (1962) 238.
- [2] A. Bertolini, An overview of automated software testing, *J. System Software* 15 (1991) 133–138.
- [3] C.J. Burgess, The automated generation of test cases for compilers, *J. Software Test. Verif. Reliab.* 4 (2) (1994) 81–100.

- [4] R. Ferguson, B. Korel, The chaining approach for software test data generation, *ACM Trans. Software Engrg. and Methodology* 5 (1) (1996) 63–86.
- [5] D.E. Goldberg, *Genetic Algorithms in Search, Optimization and Machine Learning*, Addison-Wesley, Reading, MA, 1989.
- [6] B.F. Jones, D.E., H.-H. Sthamer, A strategy for using genetic algorithms to automate branch and fault-based testing, *The Computer Journal* 41 (1998) 98–107.
- [7] J.H. Holland, Genetic algorithms and the optimal allocation of trials, *SIAM J. Comput.* (1973) 89–104.
- [8] J.H. Holland, *Adaptation in Nature and Artificial Systems*, Addison-Wesley, Reading, MA, 1975.
- [9] D.C. Ince, The automatic generation of test data, *Computer J.* 30 (1) (1987) 63–69.
- [10] B.F. Jones, H.-H. Sthamer, X. Yang, D.E. Eyres, The automatic generation of software test data sets using adaptive search techniques, in: *Proceedings of the Third International Conference on Software Quality Management*, Seville, 1995, pp. 435–444 (BCS/CMP).
- [11] B.F. Jones, H.-H. Sthamer, D.E. Eyres, Automatic structural testing using genetic algorithms, *Software Engrg. J.* 9 (1996) 299–306.
- [12] P.N. Lee, Correspondent computing, in: *Proceeding of the ACM 1988 Computer Science Conference*, Atlanta, Georgia, 1988, pp. 12–19.
- [13] G. Myers, *The Art of Software Testing*, Wiley, New York, 1979.
- [14] M.A. Ould, Testing—a challenge to method and tool developers, *Software Engrg. J.* 6 (2) (1991) 59–66.
- [15] G. Syswerda, Uniform crossover in genetic algorithms, in: *Proceedings of the Third International Conference on Genetic Algorithms*, 1989, pp. 2–9.
- [16] J. Wegener, H. Sthamer, B.F. Jones, D.E. Eyres, Testint real-time system using genetic algorithms, *Software Quality J.* 6 (1997) 127–153.
- [17] E. Alba, J.F. Aldana, J.M. Troya, Genetic algorithms as heuristics for optimizing ANN design, in: *Proceedings of the International Conference on Artificial Neural Nets and Genetic Algorithms*, Innsbruck, Austria, 1993, pp. 682–690.
- [18] S.A. Harp, T. Samad, A. Guha, Towards the genetic synthesis of neural networks, in: *Proceedings of the Third International Conference on Genetic Algorithms*, 1989, pp. 360–369.
- [19] H. Kitano, Empirical studies on the speed of convergence of neural network training using genetic algorithms, in: *Proceedings of the Eighth National Conference on Artificial Intelligence*, 1990, pp. 789–795.
- [20] P. Robbins, A. Soper, K. Rennolls, Use of genetic algorithms for optimal topology determination in back propagation neural networks, in: *Proceedings of the International Conference on Artificial Neural Nets and Genetic Algorithms*, Innsbruck, Austria, 1993, pp. 726–729.
- [21] W. Schiffmann, M. Joost, R. Werner, Application of genetic algorithms to the construction of topologies for multi-layer perceptions, in: *Proceedings of the International Conference on Artificial Neural Nets and Genetic Algorithms*, Innsbruck, Austria, 1993, pp. 675–682.