

An Introduction to Evolutionary Computation

Lecture notes for Contemporary Intelligent Information Techniques
in 2012

Akira Imada
Brest State Technical University, Belarus

(last modified on)

December 9, 2012

1 What does GA look like?

We now assume to solve a problem which includes n variables. That is, our task is to determine an optimal set of n variables. Then we design GA as follows.

1.1 Individual is represented by a chromosome

Represent a series of x_i as a *population* of strings. Each of these strings is referred to as *chromosome* or sometimes called *individual*.

x_1	x_2	x_3	x_4	x_5	x_n
-------	-------	-------	-------	-------	-------	-------

Then we start an evolution as follows, expecting better solutions from generation to generation.

1. **(Initialization)** Generate an initial *population* of p individuals at random.
2. **(Fitness Evaluation)** Evaluate fitness of each chromosome and sort the chromosomes according to its fitness from the best to the worst.
3. **(Selection)** Select two chromosomes
 - Here, from the best half of the population at random, which is called a *Truncate Selection*.
4. **(Reproduction)** Produce a child by the following two operations:
 - *Uniform Crossover*, for example
 - *Mutation*
5. Create the next generation by repeating the steps from 3 to 4 n times.
6. Repeat the steps from 2 to 5 until (near) optimal solution is obtained.

1.2 How we select parents?

Hopefully, *the better the fitness the more likely to be selected.*

1.2.1 Three different versions of Selection

- *Truncation Selection (Simplest of the three)*: Select parents from the best some-percentage of the population.

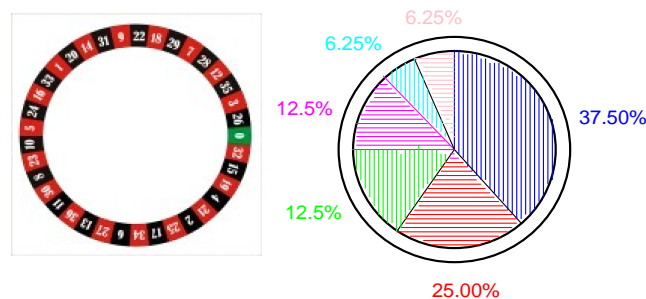
- *Roulette Wheel Selection (or Fitness Proportionate Selection)*:

With truncation selection, low fitness chromosomes have no chance to reproduce. In nature,



however, good children sometimes emerge from not good parents. So let's give a chance to a bad chromosomes to create children, though not with equal probability but proportionate to their fitness. This is called *fitness proportionate selection*, or *roulette wheel selection*.

Select such that the probability to be selected is proportional to the fitness value.



To be more specific, sort the individual from low to high and calculate cumulative value of fitness as follows: Then create a random number r from 0 to 1, and if $r < 0.3750$ then select #1, else if $r < 0.6250$ then select #2, else if $r < 0.7500$ then select #3 and so on.

- *Tournament Selection* Assume we have the original μ parents and their μ children. The fitness value of each of the 2μ individuals are compared to those of q individuals which are chosen randomly from the whole 2μ points at every time of the comparison. Then the 2μ individuals are ranked according to the number of wins, and the best μ individuals survive (q -tournament selection).

individual after sort	fitness	cumulative value of fitness
#1	0.3750	0.3750
#2	0.2500	0.6250
#3	0.1250	0.7500
#4	0.1250	0.8750
#5	0.0625	0.9375
#6	0.0625	1.0000

Table 1: Individuals are selected using random numbers from 0.0000 to 1.0000 according to the cumulative value of fitness after sorting the population.

Note that even the worst fitness individual have a chance to be selected under Roulette Selection however few it might be, while under Truncate Selection worse individual have no chance to be selected. Tournament Selection could also select a worse individual except for the worst q individuals. We can control the probability of selecting worse individual by changing q .

1.3 How parents produce a child?

1.3.1 Cross-over

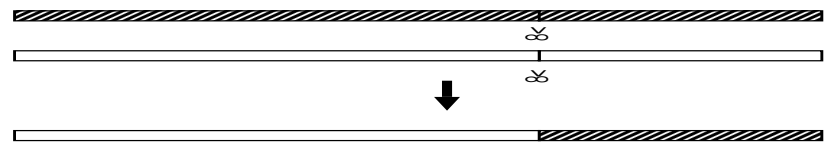
So-called *crossover* is exploited for the reproduce children. Here, we see three different versions of crossover bellow.

1.3.2 Mutation

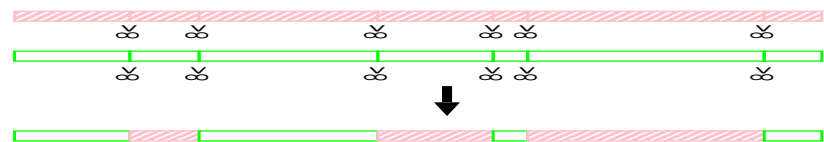
We should give a mutation to introduce new genes. This is to avoid for individuals in the population to be trapped into a *local minimum*. The probability for mutation to occur is small — typically $1/\text{number-of-genes}$. To be more specific,

If a randomly generated number from 0.0 to 1.0 is smaller than the mutation rate then mutate otherwise do nothing.

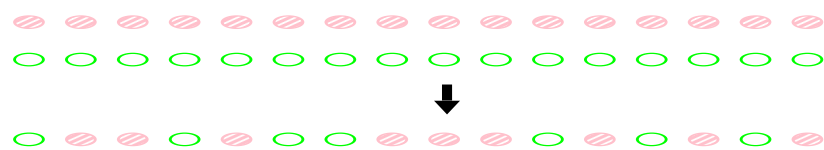
One-point Crossover



Multi-point Crossover



Uniform Crossover



2 All One Problem

To feel how the evolutionary algorithm works, let's try a toy program. Assume we have 20 genes in one chromosome. Each of the genes corresponds to a trait of us such as blue or brown eyes. All genes are binary. Also assume 1 is a good gene and 0 is a bad one.

Then starting with a population of random chromosomes, observe the evolution. Fitness can be estimated the number of 1 in the chromosom.

An example of the chromosome is

1 0 0 1 0 1 1 1 0 0 1 0 1 0 0 0 0 1 1 0.

The number of '1' is the fitness of this chromosome.

If your program works properly, you will get the best chromosom

1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1

after a several generations.

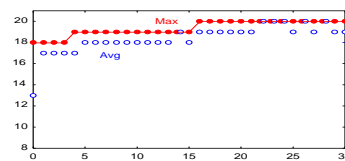


Figure 2: The generation vs fitness curve in the case of 20 genes.

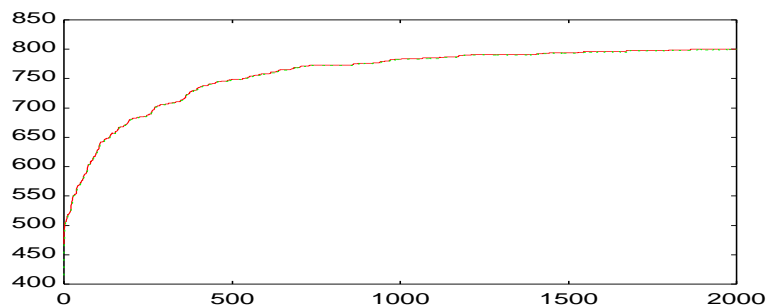


Figure 3: The generation vs fitness curve in the case of 800 genes.

3 Commonly Used Test Functions

Below we study two multi-dimensional test-functions whose location of the global minimum was already known. It is, therefore, useful to explore this test-function to learn how evolution works inside the PC.

3.1 Sphere model

3.1.1 Defined on 20-Dimensional space

Probably the simplest one is

$$y = x_1^2 + x_2^2 + x_3^2 + \cdots + x_{20}^2.$$

Or, equivalently,

$$y = \sum_{i=1}^{20} x_i^2 \quad (1)$$

which is defined, for example, on each $x_i \in [-5.12, 5.12]$ ($i = 1, 2, \dots, 20$). This function is an extension of well known $y = x^2$ to its 20-dimensional version. The unique global minimum is located on the Origin and no local minimum. Hence, this might be a good start to try a GA.

An example of chromosome is

3.2 4.4 -2.1 0.5 -3.8 -2.5 1.7 5.1 -0.3 -2.1 -3.8 5.0 0.4 4.2 -5.0 -1.3 3.3 4.0 -1.4 -3.9

By translating the i -th gene into x_i we will get all the value of x_1, \dots, x_{20} , and we can calculate the value of y . This is the *fitness* of this chromosome.

3.1.2 Defined on 1-dimensional space

To be able to imagine let's think its 2-dimensional version.

$$y = x^2$$

Then the graph is a well-know *parabola*. See Figure. 4.

3.2 Rastrigin Function

3.2.1 Defined on 20-Dimensional space

Now we observe a so-called the Rastrigin Function, which resembles a parabola but has many local minima wherever.

$$y = \{x_1^2 - \cos(2\pi x_1)\} + \{x_2^2 - \cos(2\pi x_2)\} + \cdots + \{x_{20}^2 - \cos(2\pi x_{20})\}.$$

Or, equivalently,

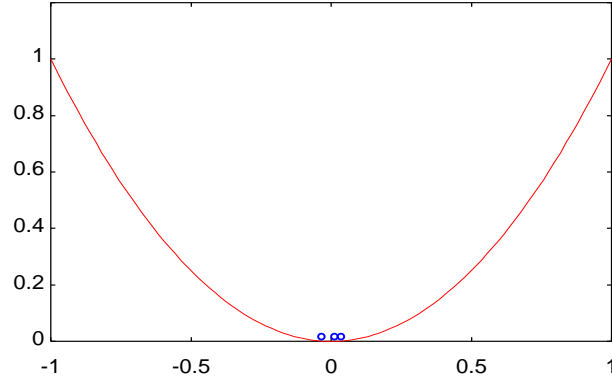


Figure 4: A two-dimensional version of the sphere model.

$$y = nA + \sum_{i=1}^n (x_i^2 - A \cos(2\pi x_i)), \quad x_i \in [-5.12 - 5.12].$$

Ruggedness might be controlled by modifying the value of A .

3.2.2 Defined on 1-dimensional space

A two dimensional example ($n = 1$):

$$y = 3 + x^2 - 3 \cos(2\pi x).$$

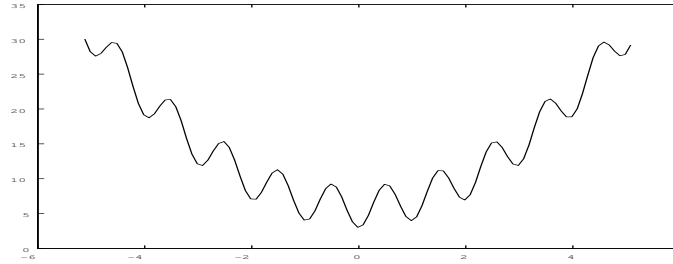


Figure 5: A 2-D version of Rastrigin function

4 Neural Network to solve even-parity problem

Here we determine a configuration of weight values of a neural network. Now think of a neural network to solve the *even-N-parity* Boolean function, that is, iff the number of "1" is even the output should be "1" otherwise "0". Assuming N-N-1 structure of feed forward neural network, we should determine $N^2 + N$ weight values so that any binary input results in the proper output.

Our chromosome in this case has $N^2 + N$ genes. Fitness can be evaluated by giving each of the neural net work corresponding to the chromosome all the possible 2^N query and counting the number of its correct output.

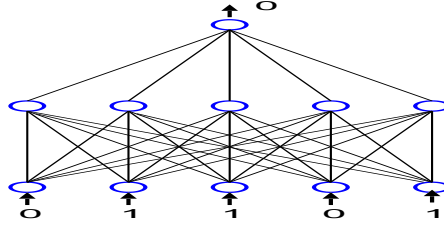


Figure 6: A 5-5-1 structure of Feedforward neural network

Now we see the set of input and output of the even-5-parity Boolean function.

input	output
00000	1
00001	0
00010	0
00011	1
00100	0
00101	1
00110	1
00111	0
01000	0
...	...
11110	1
11111	0

Table 7: Input and output of even-5-parity Boolean function.

5 Traveling Salse-person Problem

Assuming N cities all of whose coordinate are given, the Traveling Salse-person Problem (TSP) is a problem in which a sales-person should visit all of these cities once but only once and objective is to look for the shortest tour.

Assume we start from a city and return to the city after visiting all the other cities only once. When the number of cities is N , the number of all the possible routes is $(N - 1)!/2$ if we don't count the route with exactly opposite order of the cities of the route already count.

Let's try a TSP with 4 cities created at random in x - y coordinate where $0 < x, y < 10$. Then we calculate the distances between all the possible combinations of two cities, which can be shown by a distnace matrix. As the number of possible routes is just $3!/2 = 3$ we can directly calculate the distance every such routes. See an example in Figure 9.

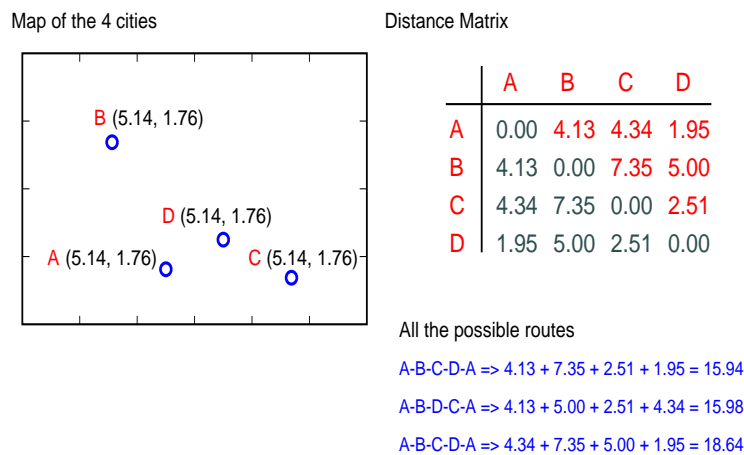


Figure 8: An example of four cities. Map, distance matrix, and all the possible routes with total distance of each route .

Problem is, the number of possible routes will explode as the number of cities N increases. How many routes do you guess when $N = 1000$ for example? Now we apply a genetic algorithm to this problem¹.

¹Here we explore the TSP by an GA. But there seems to be more direct way of computation called Ant Colony Optimization (ACO). ACO is an optimization technique we borrowed the metaphore of an intelligence of ant society as their collective behaviour. Ants are good at seeking a shortest path from their nest to a food when one of their colleagues finds it he communicates with others by using a chemical called pheromon. If we have a time we will study ACO later.

5.1 To solve TSP by a Genetic Algorithm

5.1.1 How to encode for a crossover to be valid?

First, we have to design our chromosome representing a tour? If we represent a candidate solution with a list of cities to be visited in order, such as the chromosome of the tour A-C-F-D-G-E-B is

$$(ACFDGEB) \quad (2)$$

Are the the results of crossover and mutation feasible? The answer is No! For example, the possible child of two parent $(ACFDGEB)$ and $(AGBFCDE)$ by *one-point-crossover* could be $(ACFFCDE)$ and this is not a feasible tour because C and F are visited twice and B is not visited. Or, if we give a standard mutation to $(ACFDGEB)$, for example, by replacing 4th gene with other randomly chosen city, such as $(ACFAGEB)$, which is not a feasible tour either.

Then, in order for the result of crossover and mutation to be feasible what representations are possible?

One idea is:

Algorithm 1 (Traveling Salesperson Problem) *Encoding a legitimate route to a chromosome made up of any integer.*

Step-1. Set $i = 1$.

Step-2. If i -th gene is n then n -th city in the list is the city to be currently visited.

Step-3. Remove the city from the list.

Step-4. Set $i = i + 1$ and repeat Step-2 to Step-4 while $i \leq n$.

For example, when the list of cities is

$$\{A, B, C, D, E, F, G, H, I\}$$

chromosome: (112141311) is the tour:

$$A-B-D-C-H-E-I-F-G.$$

Try some one-point crossover on two parents (112141311) and (515553321).

5.1.2 Mutation

Then next, how we design mutation? How about, for example, specifying two points at random and reverse the gene order?

5.1.3 When to stop a run?

Finally, when, on earth, we stop a run? We don't know the optimum. The answer is let's observe the evolution of fitness. We can expect the run converges to the optimum or near-optimum.

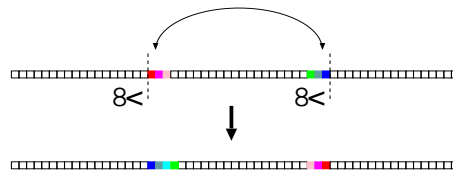


Figure 9: A possible mutation in Traveling Salesperson Problem.

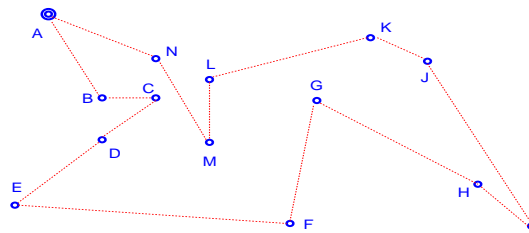


Figure 10: An example of fourteen cities and its one option of tour.

5.1.4 TSP with a huge number of cities

I found that 13,509 real cities in USA are given with their coordinates in the web-page <http://www.iwr.uni-heidelberg.de/groups/comopt/software/TSPLIB95>. Why don't we try this very challenging problem by ourselves. The plotted these cities are shown in Figure 3.

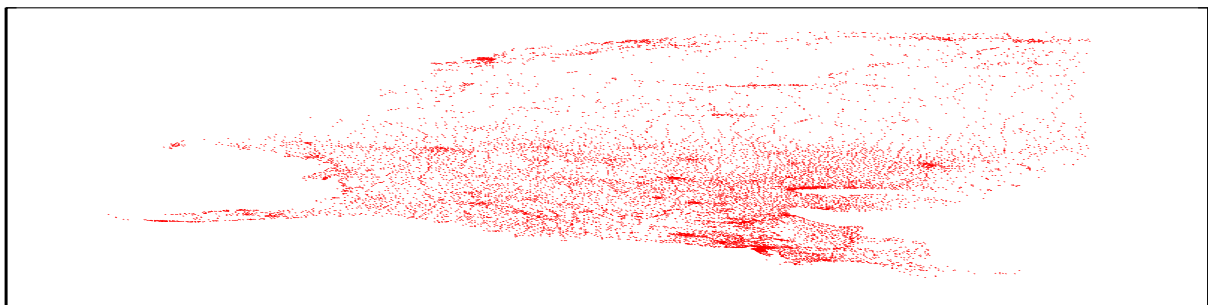


Figure 11: An example of 13509 real cities' location in US. Plotted with the data taken from <http://www.iwr.uni-heidelberg.de/groups/comopt/software/TSPLIB95>.

Assume now that we want to make an agent, or a robot, in a gridworld, a possible chromosome can be made up of integer gene from 1 to 4 where 1, 2, 3, and 4 correspond to one cell movement of the agent to north, south, east and west. Take a look at the below as an example.

(31111332332333131442411141)



6.1 Exploration of a gridworld with a limited energy

.

It would be not so difficult for us human to find such a route. But how about a computer? The question now would be, “How we designed fitness function?”

6.2 Jeep Problem

Assume we have a Jeep at the base which locates at the edge of a desert. The Jeep has a tank which can be filled with a maximum of one-unit of gasoline. With one unit

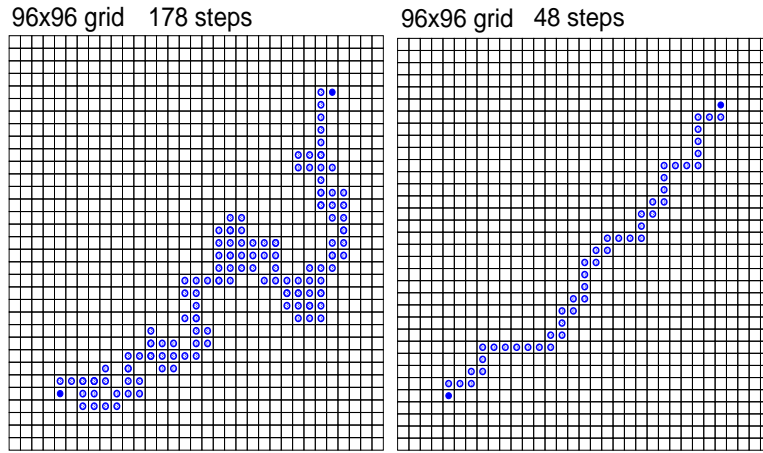


Figure 13: In the grid-world of 96 starting from (24,24) a robot walks aiming the goal at (72,72) of which the robot had no *a-priori* information. Left: The path of minimum length among 100 trials by random walk. Right: Minimal path the robot found after an evolutionary learning as shown in Fig. 3. (Marginal area is omitted.)

of gasoline, the jeep can move one unit of distance. The Jeep can only fill gasoline at the base. The jeep can carry containers to put its gasoline on the desert for a future use, Assume the tour should be on the strait line in the desert. ²

The question is “How far the jeep can penetrate in to the desert on the straight road when n unit of gasoline is available at the base.

For example when $n = 2$, the best strategy is to start the base with one unit of gasoline in its tank and go $1/3$ unit distance (it has spent $1/3$ unit of gasoline to reach the point, then put $1/3$ unit of gasoline in the container there (now jeep has $1/3$ unit of gasoline in the tank) and go back to the base. Exactly when the jeep arrive to base all gasoline filled at the start was spent. Then jeep fill the 2nd unit of gasoline given, go $1/3$ unit fill the gasoline he had put before and the tank is again full, then go forward until all the gasoline in the tank will be spent. Therefore the maximum distance the jeep can go is $4/3$ unit of distance.

Guess the maximum distance when $n = 2$. We already know the maximum distance with n unit of gasoline is D_n is expressed as the recursive relation $D_n = D_{n-1} + 1/(2n - 1)$.

Our interest is on whether an evolutionary computation can find a almost best strategy, say, $n = 5$ in which maximum distance is $1323/945=1.4$. (If my calcuration is correct. Try it by yourself).

²The problem first appeared as “a camel carrying grain in a desert” as the 52nd problem in the “*Propositions ad acuendos inventes*” (in Latin) attributed to Alcuin of York (around in B.C. 732–804). And now a jeep in a dessert, further, landrver in the Mars.

6.3 Finite State Automaton

Now let's make an artificial ant explore an gridworld. A strange assumption but there is



one grain of sugar, say 1mg, in the dark cell while nothing in the white cells. Now an ant explores the grid starting from the entrance with the aim of eating all those grains by following the shortest path to the last grain of sugar, that is to say, with 98 steps. Try to specify the transition table for such an FSA.

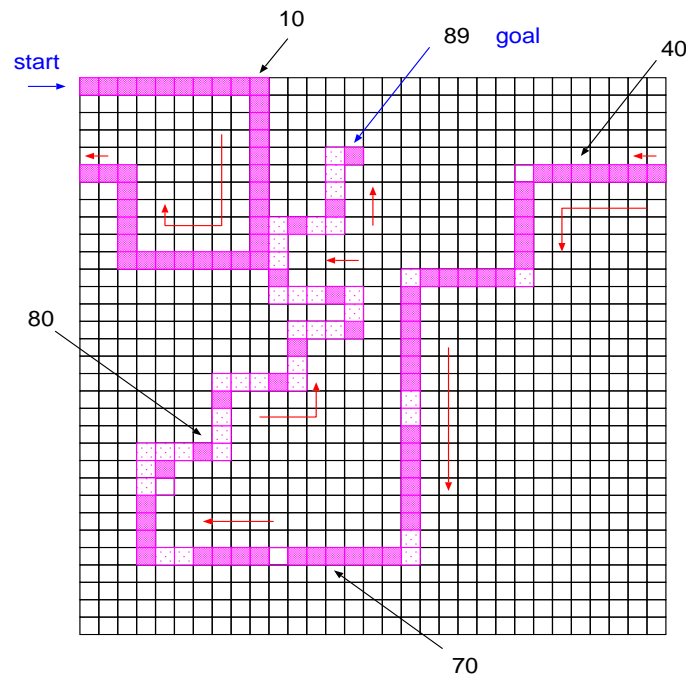


Figure 14: John-muir-trail. An agent should follow the path shown with shaded cells as effectively as possible. Note that the grid has a toroidal structure (the left neighbor cell of the leftmost cell is the rightmost cell).

Generally speaking, in order for an FSA to make a meaningful behavior, it will not be easy to specify a transition table. In this example, it would be interesting to make it using Genetic Algorithm. For the purpose, let's represent the transition table by all binary numbers. For the above example, the binary version of the transition table mentioned above is shown in Table 2.

Now we express this specific FSA with a chromosome

```
00001101101101001110100110000011110000001110101000010111010000110011001101010100001
```

current state	input	action	next state
000	0	00	001
	1	10	110
001	0	11	010
	1	01	110
010	0	10	011
	1	00	001
011	0	11	100
	1	00	011
100	0	10	101
	1	00	010
101	0	11	101
	1	00	011
110	0	00	110
	1	01	101
111	0	10	101
	1	00	001

Table 15: A binary version of the above mentioned transition table.

Now you may notice the number of states and the number of action is better to be in the form of 2^n .

We are now ready to apply a Genetic Algorithm to evolve random FSA to be clever enough to solve this John-muir-trail by the following procedure.

1. Create a population of, say 100, random binary chromosomes whose length is 80.
2. Make it try the trail one by one and evaluate how many sugar it can collect (fitness).
3. Select two chromosomes as parents (by, e.g., Roulette wheel selection) and create on child by chrossover and mutation.
4. repeat 3. (say, 100 times) untill we create a new generation.
5. repeat 3. and 4. untill we find a perfect chromosome.

7 Evolving Strategy — Iterated Prisoner's Dilemma

So, *To be or not to be?* — *That is the question.* Not only Shakespeare but in many literature works *dilemma* is their theme. The Opera “Tosca” by Puccini is one of those typical examples.³

7.1 When dilemma happens?

Assume n person are in the following game. Each of these n person is in the separate booth, where communication would not be available and not visible with each other. In each of those booths button are facilitated. You all are in the booth for one minute. If no one does not push the button, all of you will be given 100\$ each. If, on the other hand, someone push the button, or more people do so, the first one who push the button will be given 10\$ and other will not be given any money. What would you do, if you were one of these n people?

7.1.1 Condition to be a dilemma

What if the money given in case all do not touch button is 10 dollar, and otherwise the first person who push the button will be given 100 dollar? In this case no dilemma will be arisen. Push the button immediately without hesitation.

7.2 Prisoner's Dilemma

In the community of *Game Theory* we have the problem called Prisoner's Dilemma⁴ which is formulated as follows.

Problem (Prisoner's Dilemma) *Two newly arrested prisoners A and B are offered a deal:*

³Matt Ridley once wrote in his book, “The Origin of Virtue – Human Instincts and the Evolution of Cooperation.” Penguin Books (1996) about this opera. It reads: *In Puccini's opera Tosca, the heroine is faced with a terrible dilemma. Her lover Cavaradossi has been condemned to death by Scarpia, the police chief, but Scarpia has offered her a deal. If Tosca will sleep with him, he will save her lover's life by telling the firing squad to use blanks. Tosca decides to deceive Scarpia by agreeing to his request, but then stabbing him dead after he has given the order to use blanks. She does so, but too late discovers that Scarpia chose to deceive her too. The firing squad does not use blanks: Cavaradossi dies. Tosca commits suicide, and all three end up dead.* The book is regarding a Game Theory. The author continues: *Though they did not put it this way, Tosca and Scarpia were playing a game, indeed the most famous game in all of game theory, an esoteric branch of mathematics that provides a strange bridge between biology and economics. The game has been central to one of the most exciting scientific discoveries of recent years: nothing less than an understanding of why people are nice to each other. Moreover, Tosca and Scarpia each played the game in the way that game theory predicts they should, despite the disastrous outcome for each. How can this be?*

⁴Proposed by Merrill Flood and Melvin Dresher in the 1950's

- *If A confesses and B does not, A will be released and B will get 5 years in jail, and vice versa.*
- *If both confess, then both will get 4 years in jail.*
- *If both do not they will each get 2 years.*

Clearly, “0 year in prison” would be the best reward from the personal point of view. If we think of the both prisoners’ benefit, then the case “both in jail for 2 years” would be better than the case “both in jail for 4 years” unless you expected the prison as a “free hotel.”

7.2.1 Condition to be a dilemma

Now let me quote descriptions from a website⁵ where the parameter’s is as follows: R is a Reward for mutual cooperation. Therefore, if both players cooperate then both receive a reward R. (e.g., 3 points). If one player defects and the other cooperates then one player who defects receives T - the Temptation to defect - (e.g. 5 points), and the other player who cooperates receives M - the Minimum payoff (e.g., zero). If both players defect then they both receive P - the Punishment for mutual defection (e.g., 1).

If one player defects and the other cooperates then one player receives the Temptation to defect payoff (5 in this case) and the other player (the cooperator) receives the Sucker payoff (zero in this case). If both players defect then they both receive the Punishment for mutual defection payoff (1 in this case).

The question arises: what should you do in such a game?

Suppose you think the other player will cooperate. If you cooperate then you will receive a payoff of 3 for mutual cooperation. If you defect then you will receive a payoff of 5 for the Temptation to Defect payoff. Therefore, if you think the other player will cooperate then you should defect, to give you a payoff of 5.

But what if you think the other player will defect? If you cooperate, then you get the Sucker payoff of zero. If you defect then you would both receive the Punishment for Mutual Defection of 1 point. Therefore, if you think the other player will defect, you should defect as well.

So, you should defect, no matter what option your opponent chooses.

Of course, the same logic holds for your opponent. And, if you both defect you receive a payoff of 1 each, whereas, the better outcome would have been mutual cooperation with a payoff of 3. The choices of an individual is less than that could have been achieved by two cooperating players, thus the dilemma and the research challenge of finding strategies that promote mutual cooperation.

⁵<http://www.prisoners-dilemma.com/whatisit.html>. Author’s name is not visible in the page.

In defining a prisoners dilemma, certain conditions have to hold. The values we used above, to demonstrate the game, are not the only values that could have been used, but they do adhere to the conditions listed below.

Firstly, the order of the payoffs is important. The best a player can do is T (temptation to defect). The worst a player can do is to get the sucker payoff, S. If the two players cooperate then the reward for that mutual cooperation, R, should be better than the punishment for mutual defection, P. Therefore, the following must hold.

$$T > R > P > S$$

Secondly, players should not be allowed to get out of the dilemma by taking it in turns to exploit each other. Or, to be a little more pedantic, the players should not play the game so that they end up with half the time being exploited and the other half of the time exploiting their opponent. In other words, an even chance of being exploited or doing the exploiting is not as good an outcome as both players mutually cooperating. Therefore, the reward for mutual cooperation should be greater than the average of the payoff for the temptation and the sucker. That is, the following must hold.

$$R > (S + T)/2$$

However dilemma is not so serious like Hamlet's "*To be or not to be? That is a problem.*" It's better to always confess.

If the game is to be iterated, on the other hand, we have to see the game differently. Like any negotiation it's likely to have a dilemma – cooperate or betray?

7.2.2 Iterated Prisoner's Dilemma

The next question then is how about if the game is repeated? This is called the *Iterated Prisoner's Dilemma (IPD)*. In this case *strategy* to get a higher award as a result arises: What would be the optimal next action? For example, *Always Defect* strategy, or *Tit-for-Tat* strategy where the player cooperates on the first play, and afterward the same action as the opponent in the previous game.

Here, strategy determines the next action based on three previous moves of the two players in a row. Number of all possible previous three games is $2^6 = 64$ — 64 combination of Cooperate and Defect. Namely, all those possible combinations of 6 previous moves can be represented with a 64 bit binary chromosome. For example, if a history of 6 previous actions of opponent and the player is C-d-D-d-C-d we express the history by the binary number 100010 where "C" and "c" are expressed by 1 and "D" and "d" are expressed by 0 and uppercase "C" and "D" are opponent's and lowercase "c" and "d" is player's Cooperation and Defection, respectively.

Then the next action when the history was (000000) is put on the 1st bit with 0 being defection and 1 being cooperation. The same is repeated, that is, the next action when

Then each individual (chromosome) competes with each of other randomly chose p individuals, and the number of it wins against others is counted. This number of wins is the fitness value of the individual, which is called a p -tournament selection.

Assume now, for example, the game is

The history, in this case is

0, 1, 1, 1, 0, 0

Now if A's chromosome is

001010111100101010101001101010010101110100001111011100001111101

Chromosome

A: 1001000101011100100101010000110101010101101

B: 1001000101011100100101010000110101010101101

Iterated game:



Figure 16: A two-dimensional version of the Schwefel's Function's graph.

Now assume the number of population is N (e.g., 40). Each individual choose T (e.g., 5) other individuals to play a game from the one to the next. So fitness of the individual is the number of win, or the total amount of reward points during these T games.

8 Sorting Network

– How many minimum comparisons are necessary?

Which is clever? – Human or Computer?

When we write an algorithm, we often come across a necessity to sort a set of items in the order of some criteria. How, for example, do we create codes for sorting 16 integer inputs in ascending order? We pick up 2 items from one item to the next, compare, and swap them if the order is not the preferable one

Algorithm 2 (A Sorting Algorithm) *Now we assume to sort N numerical items from the smallest to the largest.*

- For $i = 1$ to $N-1$
 - ★ For $j = i+1$ to N
 - If $item(i) < item(j)$ Then swap $item(i)$ and $item(j)$

Now let us represent the sorting above by a graphics in the following way.

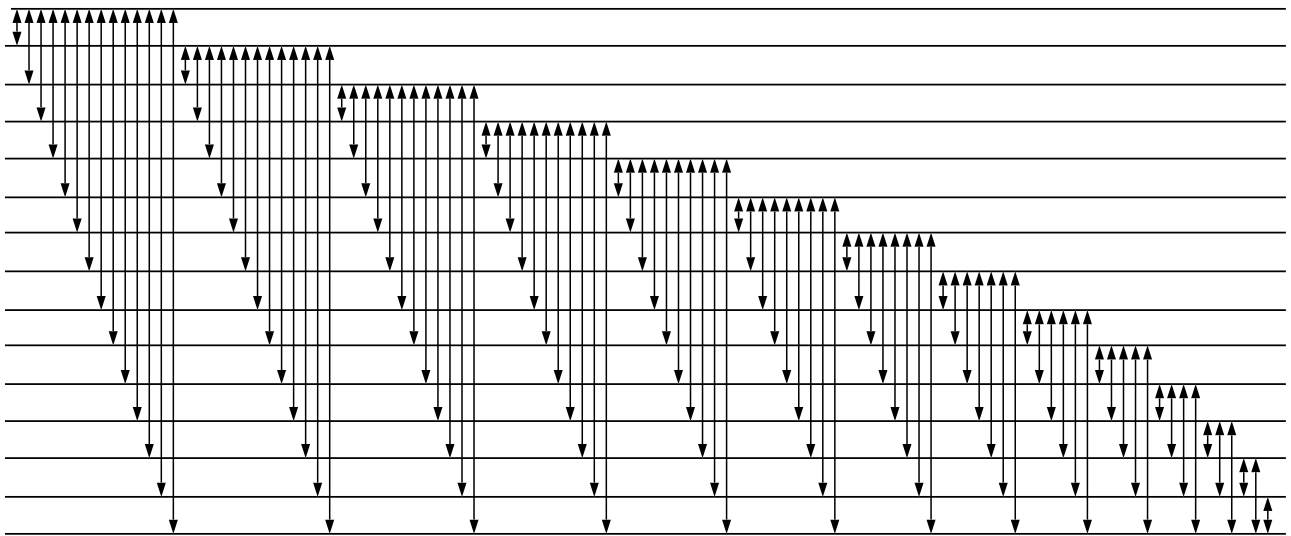


Figure 17: A typical network for sorting, not very efficient though.

The total number of comparison in this case is 120, though this is not a very efficient way. Then problem is what is the minimum number of comparisons with which any arbitrary set of 16 inputs are correctly sorted. That is,

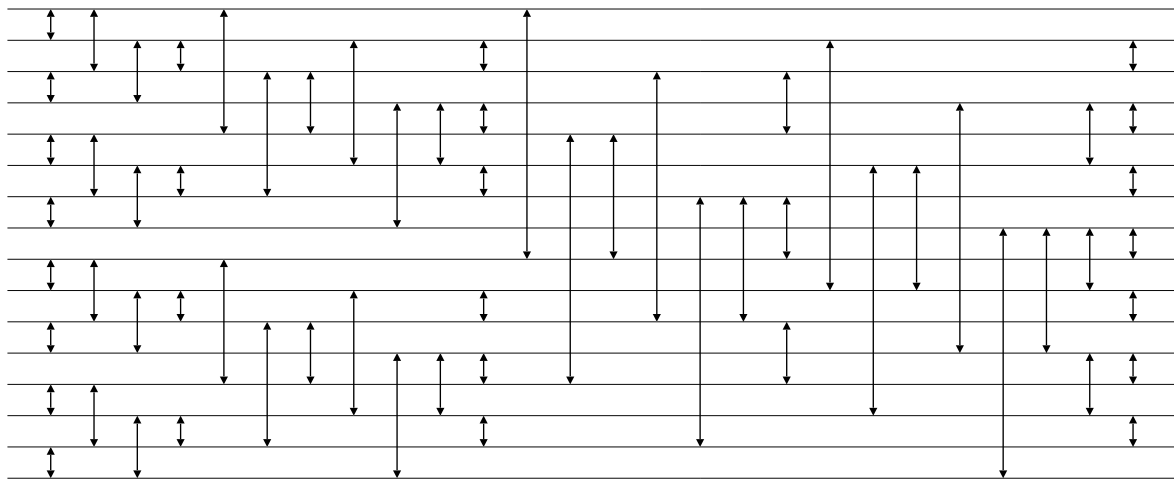
Problem (Sorting Network) *The task is to sort n items. For the purpose, the i -th and j -th element are compared and swap if necessary. and the goal is to find an algorithm which correctly sorts all n items with the minimum number of comparisons.*

It might be interesting to overview a little history of this topic. In 1960's, in a community of computer algorithms, there was a competition of what is the minimum number of comparisons when, say, ($n = 16$)? The result was

- ★ 65 comparisons Bose and Nelson (1962).
- ★ 63 by Batcher and by Floyd and Knuth (1964).
- ★ 62 by Shapiro (1969)
- ★ 60 by Green (1969)

See the Figure below.

Batcher Sort: 63 comparisons by Knuth (1973)



Comparisons in the same column can be made in parallel.

Figure 18: Batcher's proposition of sorting network with 63 comparisons (1964)

Up to now, however, we have had no proof for this optimality. Then let's make an Evolutionary Computation search for this minimum number. Would it work better than by human? Hillis (1992) challenged this. Hillis's innovation was that he employed *Diploidy Chromosome* as follows. We will show this more in detail later. Here, we show a simple version of GA implementation.

Now assume one chromosome corresponding to one sorting network is made up of 140 genes each of which takes an integer value from 01 to 16 permitting overlaps, such as

(12 01 05 04 16 12 04 14 01 02 06 07 15 08 10)

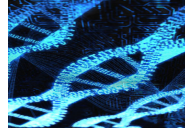
where an odd number gene and the next righthand side even number genes represent a comparison. In the example above

12 \leq 01; 05 \leq 04; 16 \leq 12;; 08 \leq 10)

Then one chromosome include 70 comparisons at most. Why at most? Because it could include a same comparison multiple time. Hence, the minimum number of comparison is one, which is very unlikely though.

9 Sorting Network Revisited

9.1 Let's be more biological - Exploitation of Diploidy Chromosomes



(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)	(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)	(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)
(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)	(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)	(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)
(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)	(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)	(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)
(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)	(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)	(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)
(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)	(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)	(1001 1000 1010 1101 1110 0100 1110 0011) (1001 1000 1010 1101 1110 0100 1110 0011)

Figure 19: An example of Hillis's set of diploidy chromosomes.

- Each individual consists of 15 pairs of 32-bit chromosomes.
- Each chromosome consists of eight 4-bit strings (called *codons*).

(0001 0010 0101 1000 0000 0100 1111 1001)
(0011 0100 0101 1000 1101 1100 1111 1001)

- Each codon represents an integer between 0 and 15 indicating which item is to be compared out of 16 items. That is, the above example is interpreted as.

(01 02 05 08 00 04 15 09)
(03 04 05 08 13 12 15 09)

- Each adjacent pair of codons in a chromosome specifies a comparison between two elements. Thus each chromosome encodes four comparisons, e.g.,

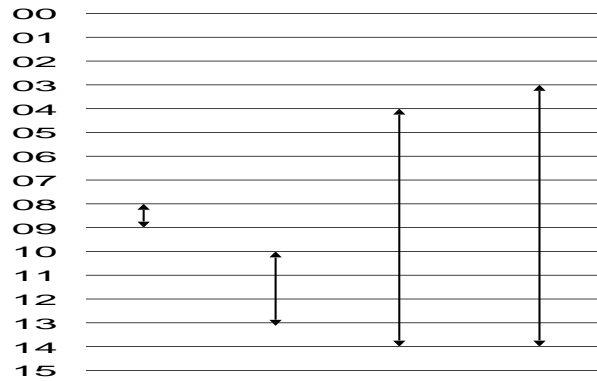


Figure 20: Four comparisons specified by the chromosome (09 08 10 13 14 04 14 03).

(09 08 10 13 14 04 14 03)

indicates the four comparisons below.

- The chromosome pairs is read off from left to right.
- If these adjacent two codons are identical at the same corresponding two positions of the chromosome pair – this is called *homozygous* – then only one pair of numbers is inserted in the phenotype. If it encodes different pairs of numbers – *heterozygous* – then both pairs are inserted in the phenotype. So in the previous example:

(01 02 05 08 00 04 15 09)
(03 04 05 08 13 12 15 09)

means the following six comparisons:

$01 \Leftrightarrow 02$, $03 \Leftrightarrow 04$, $05 \Leftrightarrow 08$, $00 \Leftrightarrow 04$, $13 \Leftrightarrow 12$, $15 \Leftrightarrow 09$

- Thus 15 pairs of chromosomes produce a phenotype with 60-120 comparisons. The more homozygous positions, the fewer comparisons.
- When two individuals are selected, one-point-crossover takes place within each chromosome pair inside each individual.
- For each of the 15 chromosome pairs, a crossover point is chosen at random and a single chromosome (called *gamete*) is formed.
- Thus 15 gametes from each parent are created.
- Each of the 15 gametes from the first parent is then paired with the gamete of the corresponding position from the second parent to form one offspring.

9.2 Pressure to homozygosity

Homozygous pair is more likely to survive than heterozygous pair, that is, two genes at the same location in a pair of chromosome will be more likely to be the same one after evolution.

For example, the probability of (1, 1) pair to be (1, 1) is $1/2$, while the probability of (1, 0) pair to be (1, 0) is $1/4$. The former is calculated as

$$1 \times ((1/4) \times (1/2) + (1/4) \times (1/2) + (1/4) \times 0),$$

while the latter as

$$(1/2) \times ((1/4) \times (1/2) + (1/4) \times (1/2) + (1/4) \times 0).$$

Hence, we can expect a more homozygous gene pair after a longer evolution. If, in our context, all the pair become to have the same chromosome, it implies the number of comparison is 60.

10 Visualization of high-dimensional space

Visualization of data in a high-dimensional space is important. Maybe you have already learned such methods like *Kohonen's Self Organizing Map (SOM)* or *Principal Components Analysis (ICA)*.

10.1 Why we need to reduce the dimension?

We, human, couldn't imagine the world of more than 3-dimensional space. In many scientific field, however, it is crucially important to grasp an image in high dimensional space. This is not only in scientific fields but also in real world around us.

Let me show an example. We now assume to assign newly employed soldires to appropriate mission according to their examination, say, of Mathematics and English.

	A	B	C	D	E	F	G	H	I	J	K	L	M
Mathematics	95	32	89	52	12	20	3	99	42	91	26	95	60
English	92	90	21	48	14	5	11	97	50	92	89	13	55

Table 21: A fictitious result of 2 examinations given to newly employed soldires.

The task of this classifying soldiers will be easier if you visualize the data.

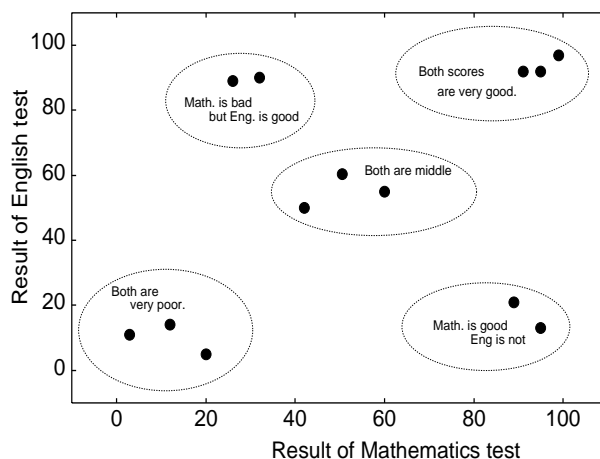


Figure 22: A visualization. It's easy to classify soldires into 5 groups.

What if, then, we have one additional score for each soldier, say, physical examination. In order to make a task like an espionage, it would be better to have a strong physical capability. In this situation we have to classify them with 3-dimensional data, or on 3-dimensional space if we want make it like the above mentioned 2-dimansuional case. Moreover, to be more practical, assume we have a set of scores of 10 different examinations. In this case, we cannot visualize any more in a usual sence.

So, visualization of high-dimensional space, or dimension reduction technique is very important topic, and so far many such techniques has been proposed, among which Kohonen's Self Organizing Map is very popular above all.

10.2 Sammon Mapping by GA

Here we learn about Sammon Mapping. Sammon Mapping is a mapping a set of points a in high-dimensional space to the 2-dimensional space with the distance relation being preserved as much as possible, or equivalently, the distances in the n -dimensional space are approximated by distances in the 2-dimensional distance with a minimal error.

This method was proposed in 1980's as an optimization problem to which they approached by Operations Research technique such as *Steepest Descend*, which is not so simple. Here, on the other hand, we employ Evolutionary Computatins which is quite simple. Let's see now what is the original Sammon Mapping look like.

Algorithm (Sammon Mapping)

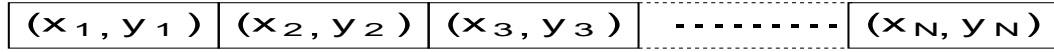
1. Assume N points are given in the n -D space.
2. Calculate distance matrix R ($N \times N$) whose i - j element is the Euclidean distance between the i -th and j -th point.
3. Also think of a tentative N points in the 2-D space that are located at random at the beginning.
4. The distance matrix Q is calculated in the same way as R .
5. Then the error matrix $P = R - Q$ is defined.
6. Search for the locations of N points in the 2-D space that minimizes the sum of element P .

This is an optimization problem which we now can solve quite simply by using EC. That is, by creating N points in 2-D space each of which corresponding N points in the n -D space with the distance relation being preserved as much as possible, or equivalently, such that the n -D distances are approximated by 2-D distances with a minimal error.

In an actual GA implementation of Sammon Mapping, chromosomes might be made up of n genes each of which corrisonds to $x - y$ coordinate of a candidate solution of n optimally distributed points in 2-dimensional space. Uniform crossover is employed and from time to time mutation is given by replacing one gene with other random $x - y$ coordinate. See the Figure 2. See also the Figure bellow.

Examples in $49^2 = 2401$ dimensional space:

Chromosome:



Recombination with Uniform Crossover:

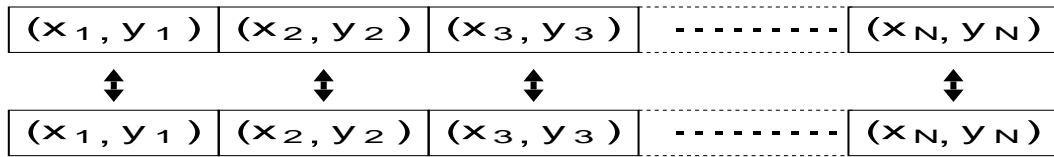


Figure 23: A chromosome representation and uniform crossover

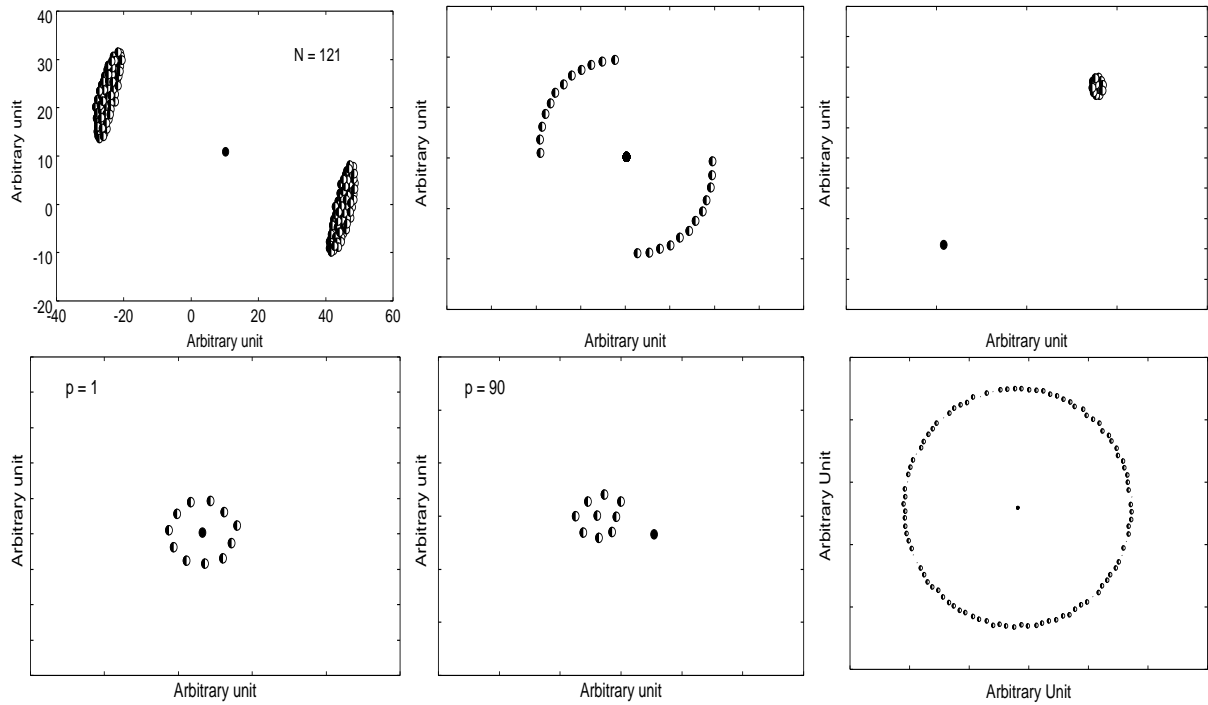


Figure 24: Six Examples of Mapping from 2401-dimensional space to the 2-dimensional space. Further explanations are shown in the text.

11 Multi Modal Problem

What if we have multiple meaningful different solutions?

In Traveling Salesperson Problem, we are interested in the result of minimum path, even if we have multiple possible paths. However, we sometimes are interested in all of possible solutions at a time in a run. For example, when we want a set of fuzzy rules for designing a fuzzy controller. The topic of this section is regarding this problem. Let's start with a simple mathematical functions.

11.1 Yet another Testfunction

A 2-D function but multi peaks

The question is how we design our chromosomes. In the multi-dimensional function $y = f(x_1, x_2, \dots, x_n)$ our genes might be continuous value each of which corresponds to the independent variable x_i ($i = 1, 2, \dots, n$), that is, our chromosomes are made up of n genes. On the other hand how should we design our chromosomes when the number of independent variable is only one. A chromosome with only one gene? How we crossover two parents?

O.K. we usually use binary chromosome in this situation. Any (decimal) real-valued variable $x_i \in [a, b]$ could be encoded by n -bit binary strings where a and b is represented by $(00 \dots 0)$ and $(11 \dots 1)$, respectively, and therefore accuracy (or granularity) is $(b - a)/(2^n - 1)$. For example, if our concern is the above $x \in [0, 1]$ then 10 bit of binary strings from 0000000000 to 1111111111 express decimals with the precision of $1/1024$. Or you might use and compare *Gray Code* where gray-code $a_1 a_2 \dots a_n$ is translated from binary number $b_1 b_2 \dots b_n$ as

$$a_i = \begin{cases} b_i & \text{if } i = 1 \\ b_{i-1} \oplus b_i & \text{otherwise} \end{cases} \quad (3)$$

where \oplus is addition modulo 2. In gray code a pair of adjacent decimals are different only with Hamming distance 1, while in the standard binary encode this does not hold.

The test-functions we studied in Section. 3 are what evolutionary computations are especially good at, since we can treat high dimensional function whatever large it may be, simply by setting the number of genes in a chromosome to the dimensionality.

Then what should we do, if we are interested in a 2-D function? For example,

$$y = \sin^6(5\pi x) \quad (4)$$

or

$$y = -2((x - 0.2)/0.8)^2 \sin^6(5\pi x) \quad (5)$$

are interesting functions in order for us to observe how randomly created chromosomes in the 1st generation evolve to find peaks. See Figure 2.

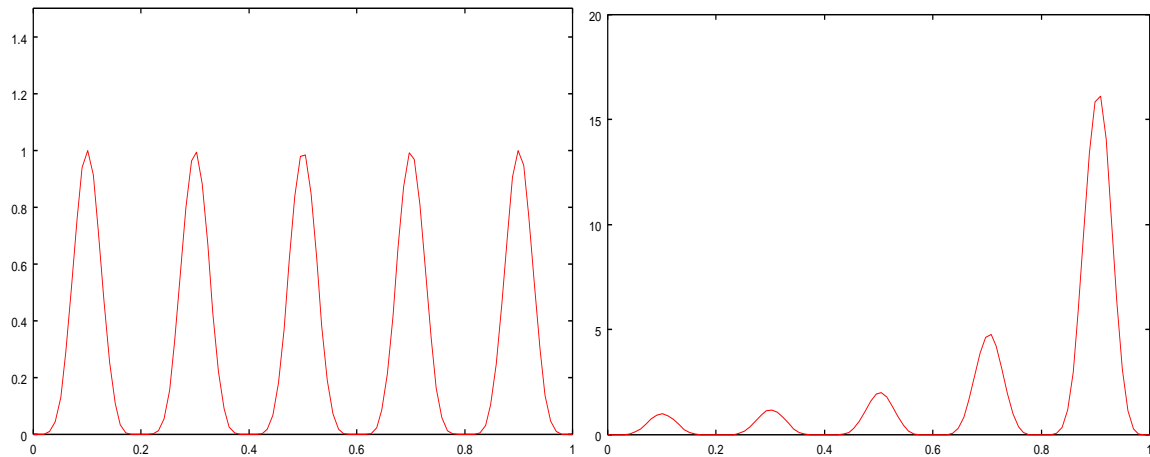


Figure 25: A multi-peak 2-D function and its variation

11.2 Multi-modal Optimization.

Sometimes we have multiple solutions. But EC usually converges only one solution out of them. Hence, to get those multiple solution we run the algorithm multiple time. Here we learn the technique in which individual construct niches and each species found a different solution at a run.

For the purpose, multiple *species* will be created and maintained in a population. These species independently search for a peak (hopefully an optimum solution), construct their *niche* and stay around the peak during a run.

In comparatively early days, essentially the following three methods were proposed. So-far proposed methods are

- Fitness sharing (Goldberg & Richardson, 1987)
 - Similar individuals share fitness with each other.
- Crowding (De Jong, 1975)
 - Similar individuals are replaced with random individuals
- Species Method.
 - Mating is restricted to among similar individuals.

These days, IMHO, the following two are popular among others.

- Deterministic Crowding (Mahfoud, 1992)
- Sequential Nicheing

Let's see some of the aspects more in detail.

Fitness Sharing Fitness of each individual is derated by an amount related to the number of similar individuals in the population. That is, shared fitness $F_s(i)$ of the individual i is

$$F_s(i) = \frac{F(i)}{\sum_{j=1}^{\mu} s(d_{ij})}$$

where $F(i)$ is fitness of individual i ; d_{ij} is distance between individual i and j ; Typically d_{ij} is *Hamming distance* if in *genotypic space* *Euclidean distance* if in *phenotypic space* and $s(\cdot)$ is called *sharing function* and defined as:

$$s(d_{ij}) = \begin{cases} 1 - (d_{ij}/\sigma_{\text{share}})^\alpha & \text{if } d_{ij} < \sigma_{\text{share}} \\ 0 & \text{otherwise} \end{cases}$$

where σ_{share} is interpreted as size of niche, and α determines the shape of the function. The denominator is called *niche count*. You see shape dependency of $s(d_{ij})$ on α in Figure 11.2.

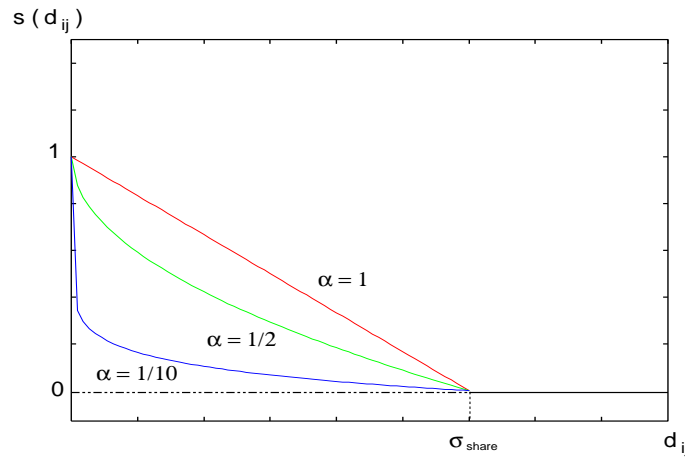


Figure 26: A shape dependency of $s(d_{ij})$ on α .

To be short (not so short though): Similar individual should share fitness. The number of individuals that can stay around any one of peaks (niche) is limited. The number of individuals stay near any peak will theoretically be proportional to the *height* of the peak

Deterministic Crowding: If the parents will be replaced or not with their children will be determined under a criteria of the distance between parents and children,

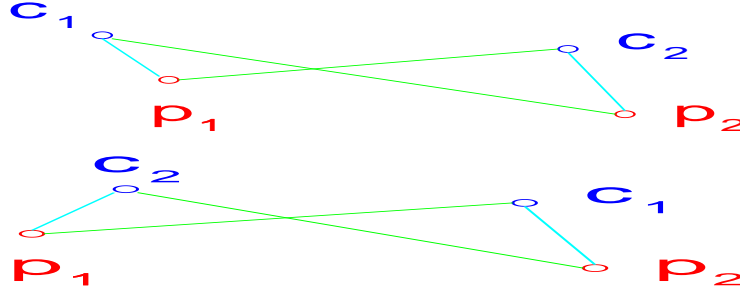


Figure 27: A typical two cases of distance between parents and children.

Algorithm Assuming crossover, mutation and fitness function are already defined

1. Choose two parents, p_1 and p_2 , at random, with no parent being chosen more than once.
2. Produce two children, c'_1 and c'_2 .
3. Mutate the children yielding c_1 and c_2 , with a crossover.
4. Replace parent with child as follows:
 - IF $d(p_1, c_1) + d(p_2, c_2) > d(p_1, c_2) + d(p_2, c_1)$
 - * IF $f(c_1) > f(p_1)$ THEN replace p_1 with c_1
 - * IF $f(c_2) > f(p_2)$ THEN replace p_2 with c_2
 - ELSE
 - * IF $f(c_2) > f(p_1)$ THEN replace p_1 with c_2
 - * IF $f(c_1) > f(p_2)$ THEN replace p_2 with c_1

where $d(\zeta_1, \zeta_2)$ is the Hamming distance between two points (ζ_1, ζ_2) in pattern configuration space. The process of producing child is repeated until all the population have taken part in the process. Then the cycle of reconstructing a new population and restarting the search is repeated until all the global optima are found or a set maximum number of generation has been reached.

Sequential Niche: Single run is repeated sequentially, keeping the best individual at each run.

Algorithm

1. Define niche radius r .
2. Define modified-fitness function $m(\mathbf{x})$ by equating it to the original fitness function $f(\mathbf{x})$ here at the start.

3. Run the GA and pick up the best individual at the end of the run.
4. Update $m(\mathbf{x})$ as ⁶

$$m_{n+1}(\mathbf{x}) = m_n(\mathbf{x}) \cdot g(\mathbf{x}, \mathbf{s}_n) \quad (7)$$

where n is the number of so-far run, \mathbf{s}_n is the best individual in the n -th run and

$$g(\mathbf{x}, \mathbf{s}_n) = \begin{cases} (d_{xs}/r)^\alpha & \text{if } d_{xs} < r \\ 1 & \text{otherwise} \end{cases} \quad (8)$$

is called derating function where d_{xs} is a distance between \mathbf{x} and \mathbf{s}_n while m_0 is the original raw fitness function of each individual.

5. Run the GA using the modified fitness function and keep the best individual found in the run,
6. Update the modified fitness function
7. If the raw fitness of the best individual exceeds the solution threshold, (See also below) then display this as a solution.
8. If all solutions have not been found, then return to step 5.

- *Solution Threshold is*

- Lower fitness limit for maxima of interest, assuming we know how many peaks. If it's not of the case, set the threshold to zero.

Excercise Like Figure 11.2, draw a graph of $y = (x/r)^\alpha$ with $r = 1$ and $\alpha = 0.5, 1, 2, 4, 8$ to know what $g(\mathbf{x}, \mathbf{s}_n)$ looks like.

It would be interesting to try a multi-modal EC to the following two test functions.

- (1) $y = \sin^6(5\pi x)$
- (2) $y = -2((x - 0.2)/0.8)^2 \sin^6(5\pi x)$

⁶This is called a *Power Derating Function* when we think of another alternative called *Exponential Derating Function*:

$$g_e(x, s) = \begin{cases} \exp((\log m(x, s)) \cdot (r - d_{xs}/r)) & \text{if } d_{xs} < r \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

12 Multi Objective Genetic Algorithm (MOGA)

So far we have learned how to get the possible solution(s) which fulfills one objective function for the problem, that is, the goal is maximize the fitness function. In real world problem, however, we have usually multiple objectives or criteria to be fulfilled simultaneously.

Those objectives sometimes conflict with each other. Like “time” and “money”: The more we want to earn money, the less time to spent the money; or “reliability” of the product and “cost” to produce it in a manufactural factory. Or, suppose an Opera Company trys to employ one Soprano singer. The criteria is voice, beauty-or-not), slim-or-not, language-capability (Italian, German, etc). However God tend not to give us two talents at a time, alas.

Then, first of all, when we have multiple objective function, we must define an important concept of parate optimal or equivalently non-dominated solution.

Definition (Parate Optimal or Non-dominated Solution) *A candidate solution is called a non-dominated iff there is no ohter better solution w.r.t. all the objectives.*

To be more specific, assume we have n objective functions;

$$f_1(\mathbf{x}), f_2(\mathbf{x}), f_3(\mathbf{x}), \dots f_n(\mathbf{x})$$

where \mathbf{x} is a candidate solution. Now if a new candidate solution \mathbf{y} improves all the objetives for \mathbf{x} , i.e.,

$$f_i(\mathbf{y}) > f_i(\mathbf{x}) \text{ for } \forall i$$

we say

“ \mathbf{y} dominate \mathbf{x} .”

When no such \mathbf{y} exists, we say

“ \mathbf{x} is non-dominated” or “Parete Optimum.”

A toy example: We now assume the two objective functions as follows.

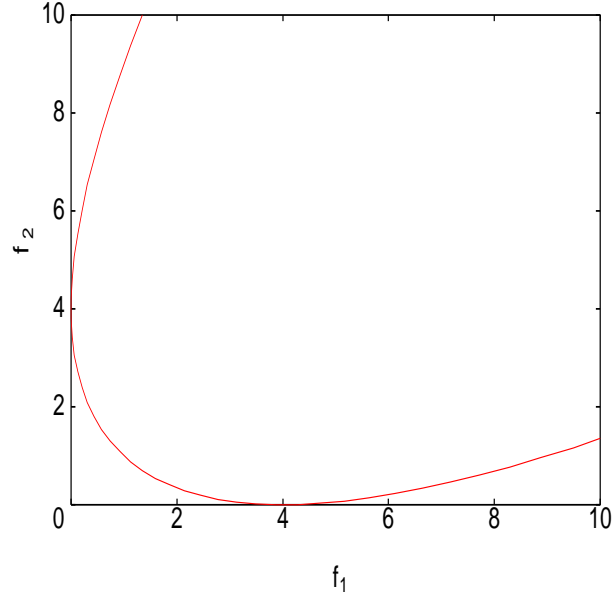
$$\begin{cases} f_1(x) = x^2 \\ f_2(x) = (x - 2)^2 \end{cases}$$

- $x=0$ is optimum w.r.t. f_1 but not so good w.r.t. f_2 .
- $x=2$ is optimum w.r.t. f_2 but not so good w.r.t. f_1 .

	test-1	test-1	test-1	test-1	test-1	dominated by	dominates	rank
A	9	9	9	8	7	0	3	1
B	5	4	5	3	6	1	2	2
C	3	3	4	2	3	2	1	3
D	2	2	3	1	2	3	0	4
E	1	1	1	1	9	0	0	4

Table 28: An example of who dominates whom and how rank is counted.

- Any other point in between is a compromise or trade-off and is a Pareto-optimum.
- But the solution $x=3$, e.g., is not a Pareto-optimum since this point is not better than the solution $x = 2$ w.r.t. either objective.
- If we plot in the f_1 - f_2 space, an increase in f_1 in some region means a decrease in f_2 , or vice versa which implies that the solutions in the region are Parete optimum, while in other region an increase in f_1 make f_2 increas (decrease). See Figure ??.
This f_1 - f_2 space is called a *Trade-off Space*.

Figure 29: Trade-off space for $f_1(x) = x^2$ and $f_2(x) = (x - 2)^2$.

Thought Experiment: What if we plot all individuals of generation 0 and, say, generation 100?

How an implementation goes? Evolution is rather similar more or less to a GA with single fitness function. The main difference is we have multiple objective function. Hence we merge these multiple objective function into one fitness function. So far many ideas have been proposed. Among all:

- Weighted sum approach.

- The fitness function is calculated as

$$f(\mathbf{x}) = \sum_{i=1}^N w_i f_i(\mathbf{x}) \quad (9)$$

where w_i expresses an importance of the i -th objectives.

- Note that for any set of weight > 0 , the optimum is always a non-dominated solution but opposite is not always true.

- The minimax approach

- The fitness function is calculated by minimizing the maximum of n objective functions.

- Target vector approach

- The fitness function is calculated by minimizing the vector

$$(f_1, f_2, f_3, \dots, f_n)$$

from a pre-defined goal.

- Median/Average ranking approach

- The rank $r(\mathbf{x}_i)$ of each individual x_i in the population w.r.t. i -th objective function is calculated. Then the fitness is defined as median/average of these r_i ($i = 1, \dots, n$).

- Parete ranking approach

- Ranking is according to “*how many individuals in the population they are dominated by.*”

We now take a look at a typical implemetation of MOGA.

Algorithm (A Multi Objective GA)

1. *Initialize the population.*
2. *Select individuals uniformly from population.*
3. *Perform crossover and mutation to create a child.*

4. Calculate the rank of the new child.
5. Find the individual in the entire population that is most similar to the child. Replace that individual with the new child if the child's ranking is better, or if the child dominates it.⁷
6. Update the ranking of the population if the child has been inserted.
7. Perform steps 2-6 according to the population size.
8. If the stop criterion is not met go to step 2 and start a new generation.

⁷Step 5 implies that the new child is only inserted into the population if it dominates the most similar individual, or if it has a lower ranking, i.e. a lower degree of dominance.

The restricted replacement strategy also constitutes an extreme form of elitism, as the only way of replacing a non-dominated individual is to create a child that dominates it.

The similarity of two individuals is measured using a distance function.

13 Data Mining

Data mining is a method to extract information or knowledge from a mountain of data.

See for example a headline from the New York Times on 13 November 2012. It reads:

Secret of the Obama Victory? Rerun Watchers, for One Thing ... By JIM RUTENBERG The Obama campaign found supporters by culling never-before-used data about viewing habits and combining it with more personal information.

13.1 Data classification

Specifically, here, we take a look at data classification. One method to classify data set into a set of classes, we exploit a rule set, such as

IF <condition> THEN <class>.

As an example, let's classify Iris flowers. Iris flower dataset⁸ is made up of 150 samples consists of three species of iris flower, that is, *setosa*, *versicolor* and *virginica*. Each of these three families includes 50 samples. Each sample is a four-dimensional vector representing four attributes of the iris flower, that is, *sepal-length*, *sepal-width*, *petal-length*, and *petal-width*. All data are given as crisp as below.

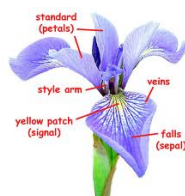


Figure 30: A part of the iris flower data-base.

x_1	x_2	x_3	x_4	class
0.65	0.80	0.20	0.08	1
0.62	0.68	0.20	0.08	1
0.89	0.73	0.68	0.56	2
0.87	0.70	0.71	0.60	2
0.80	0.75	0.87	1.00	3
0.73	0.61	0.74	0.76	3

⁸University of California Irvine Machine Learning Repository.
ics.uci.edu: pub/machine-learning-databases.

Now, for the time being we use only x_1 and x_2 just for the sake of simplicity to describe rules, unless otherwise notified. For example

IF $0.81 \leq x_1 \leq 0.89$ AND $0.70 \leq x_2 \leq 0.73$ THEN this is Class 2

classifies $x_1 = 0.89$, $x_2 = 0.73$ properly to Class 2 while $x_1 = 0.65$, $x_2 = 0.80$ and $x_1 = 0.80$, $x_2 = 0.75$ are not.

So far so good, but what if the region overlaps with each other between more than two pieces? Or, what if unknown somewhat irregular data are given?

13.1.1 Evaluation of how good is a rule

The rule for class 1 should accept all data of class 1, but at the same time this rule should reject all data of class 2 and class 3. So count (i) how many data from class 1 are successfully accepted, and then (ii) how many data from class 2 are successfully rejected and (iii) how many data from class 2 are successfully rejected.

One method to make it will be to evaluate the rule by

$$\frac{TP}{TP + FN} \times \frac{TN}{FP + TN} \quad (10)$$

where TP stands for true-positive, FN stands for false-negative, TN stands for true-negative, and FP stands for false-positive, and true-positive is the number of cases covered by the rule that have the class predicted by the rule; false-positives is the number of cases covered by the rule that have a class different from the class predicted by the rule; false-negatives is the number of cases that are not covered by the rule but that have the class predicted by the rule; true-negatives is the number of cases that are not covered by the rule and do not have the class predicted by the rule.

Our mission is to train our classifiers with known data shown in Appendix as "*data for training*" and then evaluate how good is the classifiers with "*data for checking*" also shown in appendix.

Evaluation is how the system appropriately classifies 69 data. Ask system which family by giving 69 data one by one. Score is incremented if the result is correct. Hence the maximum score is 69 and minimum is 0. Please note that even random guessing would score one out of three. So score 23 might be the most stupid classifier.

What you should show me are (1) run the algorithm and result of input x_1 , x_2 , x_3 , and x_4 (2) a rule set, and (3) success rate (true-positive, true-negative, false-positive, false-negative).

13.1.2 Genetic Algorithm (GA) Approach

- **Specify region of attributes by chromosom**

Let's specify <condition> like

IF $a_1 < x_1 < a_1 + \delta_1$ AND $a_2 < x_2 < a_2 + \delta_2$ THEN class i

Then our chromosome might be

$(a_1 \ \delta_1 \ a_2 \ \delta_2)$

13.1.3 Fuzzy Logic Approach

- *Still with an evolution of chromosomes*

With our chromosome being like

$(\text{large} \ \text{medium})$

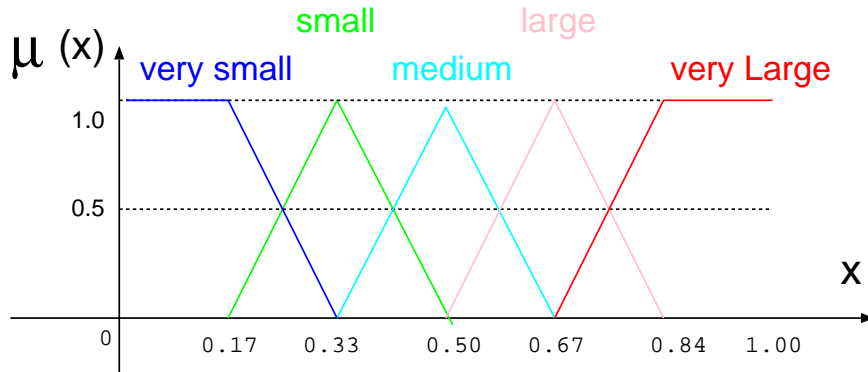
we can evolve a population of such chromosomes again with fitness being the Eq. (10).

- *More fuzzy approach - membership fuction*

If we can calculate a membershipfunction of, for example

IF x_1 is LARGE AND x_2 is MEDIUM THEN class 1

We can know the likelinass of this rule. While membership function of each linguistic term is assumed as



and use formulae of Fuzzy set theory

$$\mu_{A \cap B}(x) = \min\{\mu_A(x), \mu_B(x)\} \quad (11)$$

$$\mu_{A \cup B}(x) = \max\{\mu_A(x), \mu_B(x)\}, \quad (12)$$

$$\mu_{A \rightarrow B}(x) = \mu_A(x) \times \mu_B(x) \quad (13)$$

13.1.4 Artificial Immune System (AIS) Approach

• Clonal Selection in general

Assume an antigen invaded in our body. Then B cells with high affinity to the antigen are activated. These B cells are stimulated to proliferate producing large number of clones (Note that these clones are an exact copy of their parent) and these clones mutate and turn into plasma cells.

We now define a measure of how an antibody matches to the invading antigen. Let's call it *affinity*.

Although the maturation process produces a small number of high affinity matches it must also produce a large number of low affinity B cells.

This algorithm is proposed by ... for the purpose of pattern classification being based on the algorithm called CLONALG proposed by ... for the purpose of pattern recognition.

Algorithm 3 Artificial Clonal Selection algorithm

Assume we have a set of Antigens Ag^9 .

1. Generate an initial population of N antibodies Ab
2. Select an antigen Ag_i
3. For every member of Ab calculate its affinity to the Ag_i
4. Select the n highest affinity antibodies and generate a number of clones for each antibody in proportion to its affinity and construct a new population C^i (a population specifically for Ag_i)
5. Mutate every member of C^i with a probability being inversely proportional to its affinity. Let's call this population a mature population C^{i*}
6. Measure again affinity of the member C^{i*} to the Ag_i . (If its affinity is greater than the current memory cell Abm_i it become the new memory cell.
7. (a) Replace p highest affinity antibodies in Ab with p highest affinity matured antibodies in C^{i*}
(b) Replace $(N - p)$ low affinity antibodies in Ab with a new set of $(N - p)$ random antibodies.
8. Repeat 2-7 for all antigens Ag_i ($i=1, 2, \dots, M$)
9. Repeat 1-8 for all different antigens¹⁰.

Thus we can create N mature antibody population C^{i*} for each of antigen Ag_i . And one memory cell for each antigen.

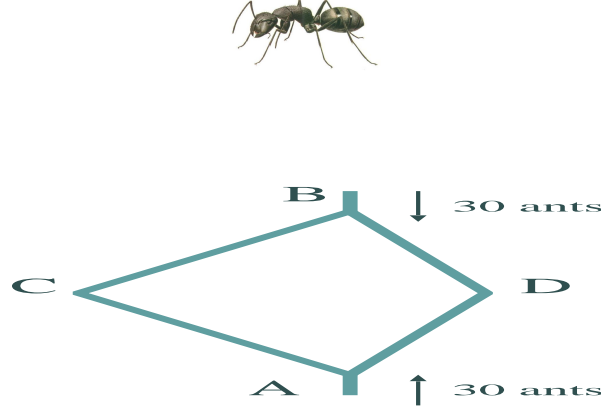
⁹Here we take each of three different iris flowers as antigen, that is, we have three Antigens called Setosa, Versicolor, and Virsinica.

¹⁰That is, Setosa, Versicolor, and Virsinica.

13.1.5 Ant Colony Optimization (ACO) Approach

- ACO in general

This subsection is edited from G. Sarigiannidis¹¹.



Which path the ants should choose

Each ant who is at the city s choose the next city j from $J(s)$ – its not-yet-visited-city from city s – with a probability:

$$p_{sj} \propto \tau_{sj}(t) \{\eta(s, j)\}^\beta.$$

Or, to normalize it:

$$p_{sj} = \frac{\tau_{sj}(t) \{\eta(s, j)\}^\beta}{\sum_{j \in J(s)} \tau_{ij}(t) \{\eta(s, j)\}^\beta}$$

where $\tau_{ij}(t)$ is the amount of pheromone on the path from city s to city j , $\eta(s, j) = 1/d_{sj}$ while d_{sj} is the distance between city s and city j , and β is a parameter which determines the importance of balance between pheromone and distance.

Local pheromon update

When all the ants reach the next city, the amount of pheromon is updated as:

$$\tau_{rs} = (1 - \rho) \cdot \tau_{rs} + \rho \cdot \tau_0$$

Global pheromon update

After all ants have completete their tours, they remain to the position for the next iteration, and the global pheromone updaite is applied to all the path (r, s) one by one as follows:

$$\tau_{rs} = (1 - \alpha) \cdot \tau_{rs} + \alpha \cdot \Delta\tau_{rs}$$

where $\Delta\tau_{rs} = 1/(L_{BS})$ if the path (r, s) belongs to the best tour so far, otherwise 0 where L_{BS} is the minimum total tour length found so far. An interpretation of $0 < \alpha < 1$ can

¹¹G. Sarigiannidis "Discussion on the implementation of the Ant Colony System for the Traveling Salesman Problem." Delft University of Technology

be the pheromone decay parameter, that is, evaporation rate of pheromone.

Parameters

Typically, the number of ants is 10, and the other parameters are $\beta = 2$, $\alpha = \rho = 0.1$, $\tau_0 = 1/(n \cdot L_{NN})$ where L_{NN} is the total tour length by Nearest Neighbor heuristic and n is a number of cities.

The ant who is in the node i chooses the next node to move from a set of possible nodes directly reachable with a probability p_{ij} which are pre-specified in advance a run constant during a run, by considering (i) η_{ij} – how the path to the node j looks good (called heuristic term) and (ii) τ_{ij} – how the path accumulate pheromone those previous ants left.

• **TSP by ACO** In order to understand how Ant Colony Optimization works, let's take a brief look at its application to the Traveling Salesperson Problem. We now assume N cities that a salesperson must visit once and only once. Now imagine N ants are assigned to each of those N cities – one ant at one city. The following parameters are to be used.

★ probability for one ant at node i to choose node j

$$p_{ij} \propto \eta_{ij} \tau_{ij}^{\beta}$$

where $\beta(> 0)$ determines the relative importance of pheromone versus heuristics, and η is called *heuristic parameter*. The larger value, the more preferable. In the TSP η_{ij} simply takes a value $1/d_{ij}$ where d_{ij} is the distance between city i and city- j . τ is called *pheromone density* and updated during a run in the following way.

• Data Classification by ACO

We use Iris-flower data again and the rules to classify iris-flower data are described with fuzzy linguistic term. A general form of such rule is:

IF x_1 is small & x_2 is large & \dots THEN y_1 is strong & y_2 is weak & \dots

Here, we assume those linguistic terms such as *small*, *large*, etc for x_i , and *strong*, *weak*, etc. for y_j are all pre-defined by a membership function.

Remember here how we define, for example, a membership function of

IF A & B THEN C & D.

Then we can define a membership function of these rules.

$$\mu_{A \cup B \Rightarrow C \cup D}.$$

Our purpose here is to classify iris-flower data with four features into either of three families. So the form of our rule is specifically like:

IF x_1 is A_1 & x_2 is A_2 & x_3 is A_3 & x_4 is A_4 THEN y is B_j

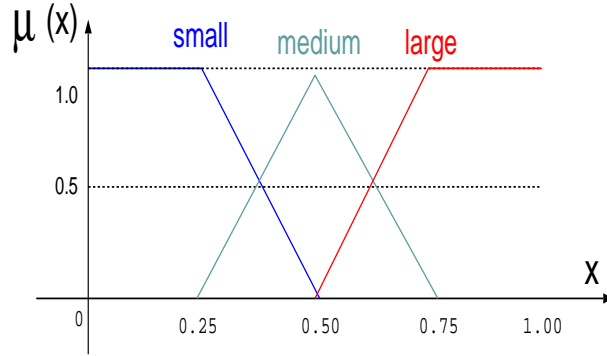


Figure 31: Three membership function of "small", "medium", "large" for all x_1, x_2, x_3 and x_4 in commune.

where A_i is either of *small*, *medium*, or *large* and B_j is either of 1, 2, or 3.

Then the number of all possible such rules is $3^4 \times 3 = 243$. Now let's denote IF-part of the i -th rule as R_i . Now we assign heuristic term for each of those links from A_i to B_j as

$$\eta_{ij} = \mu_{R_i \Rightarrow B_j}(x_1, x_2, x_3, x_4, y_1) \quad (14)$$

For example

$$\eta_{ij} = \max \min(\mu_{A_1}(x_1), \mu_{A_2}(x_2), \mu_{A_3}(x_3), \mu_{A_4}(x_4))$$

See Figure 33.

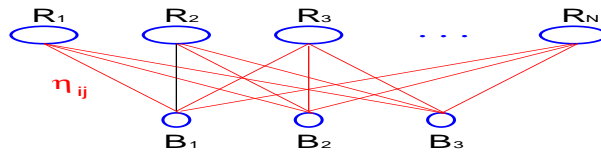


Figure 32: A schematic diagram of all possible link between IF-part and THEN-part of rule set. Heuristic term η_{ij} are assigned to those links.

With initial pheromone being arbitrary τ_0 , pheromone on the link from i to j is updated every time when one ant walks this link (local update) as

$$\tau_{ij} = (1 - \rho) \cdot \tau_{ij} + \rho \cdot \tau_0. \quad (15)$$

When one ant finish all the R_i pheromone on all the links are updated (global update) as

$$\tau_{ij} = (1 - \rho) \cdot \tau_{ij} + \rho \cdot \Delta_{ij} \quad (16)$$

where Δ_{ij} is fitness of the *best rule so far found* if the link from i to j belongs to the link of the best rule so far found, and 0 otherwise. Use Eq. (10) for fitness calculation of one

	x_1	x_1	x_1	x_1	B
R_1 :	small	small	small	small	1
R_2 :	small	small	small	small	2
R_3 :	small	small	small	small	3
R_4 :	small	small	small	medium	1
R_5 :	small	small	small	medium	2
R_6 :	small	small	small	medium	3
R_7 :	small	small	small	large	1
	
R_{243} :	large	large	large	large	3

Figure 33: All possible rules are something like this.

rule. The calculation is made by using the *data for training*. Note that it is this process that the system learn how to classify well.

The probability of one ant at the node i to choose the link to the node k is as already described:

$$p_{ij} = \frac{\tau_{ij}(t) \{\eta(i, j)\}^\beta}{\sum_{j \in J(i)} \tau_{ij}(t) \{\eta(i, j)\}^\beta}$$

with β being, say, 2.

Now let's denote input-output data set from our 120 training data of iris-flowers as

$$\mathbf{e}^k = (x_1^k, x_2^k, x_3^k, x_4^k, y^k) \quad k = 1, 2, \dots, 120.$$

Then our rules are like

R^s : IF x_1 is large & x_2 is small & x_3 is medium & x_4 is small THEN 2.

To be more general

R^s : IF x_1 is A_1^s & x_2 is A_2^s & x_3 is A_3^s & x_4 is A_4^s THEN B^s .

Then

$$\eta_{ij} = \max_{e^k} (\min(\mu_{A^i}(\mathbf{x}^k), \mu_{B^j}(y^k))) \quad (17)$$

where

$$\mu_{A^i}(\mathbf{x}^k) = \min(\mu_{A_1^i}(x_1^k), \mu_{A_2^i}(x_2^k), \mu_{A_3^i}(x_3^k), \mu_{A_4^i}(x_4^k)) \quad (18)$$

Let me take an example. Assume $s = 1$ and

R^1 : IF x_1 is small & x_2 is small & x_3 is small & x_4 is small THEN 1.

That is,

A_1^1 is small, A_2^1 is small, A_3^1 is small and A_4^1 is small.

Then from the iris-flower data for training for $k = 1$

$$x_1 = 0.65, x_2 = 0.80, x_3 = 0.20, x_4 = 0.08 \text{ and } y = 1.$$

That is,

$$\mathbf{e}^1 = (0.65, 0.80, 0.20, 0.08, 1).$$

Corresponding values of membership function

$$\mu_{small}(0.65), \mu_{small}(0.80), \mu_{small}(0.20), \mu_{small}(0.08)$$

are obtained from Fig. 31.

By repeating this procedure from $k = 1$ to $k = 120$, and then take the maximum in Eq.(17) we can obtain each value of η_{ij} .

As for B^s , they are generally also a fuzzy linguistic value such as *strong* and its membership function is defined in the same way as Fig. 31, but here B^s is integer value 1, 2, or 3, that is *singleton*. Therefore, for example $\mu_{setosa}(y) = 1$ if and only if $y=1$ otherwise $\mu_{setosa}(y) = 0$. Similarly, $\mu_{versicolor}(y) = 1$ if and only if $y=2$ otherwise $\mu_{versicolor}(y) = 0$, and $\mu_{virginica}(y) = 1$ if and only if $y=3$ otherwise $\mu_{virginica}(y) = 0$.

Now we are ready to run an ACO algorithm for Iris-flower classification

Algorithm 4 (ACO for iris-flower classification) Initialize τ_0 , calculate all η_{ij}

1. For $k = 1$ to K .
2. For $i = 1$ to 243 .
3. Put ant k to the node R_i .
4. Choose link according to p_{ij} updating t_{ij} locally.
5. Set THEN-part of the " R_i as B_j .
6. Repeat 2-5.
7. Update all of the τ_{ij} globally according to the fitness of the rule obtained..
8. Repeat 1-7.

14 Genetic Programming (GP)

The title of this section suggests *Genetic Programming* suggests *Evolution of Program*. What we want here is when we have specific task and need a program to solve the task, we start with a population of random programs and then evolve them. We expect from generation to generation a better programs than previous ones and perfect programeventually emerges.

As previous evolution, (i) we create a population of random chromosomes; (ii) evaluate fitness of each chromosome; (iii) select chromosomes so that better two are more likely to be selected; (iv) produce children by crossover and mutation (v) repeat (iii) to (iV) until the next generation includes the same number of chromosomes as previous generation; (vi) repeat (ii) to (v) until fitness saturates or pre-determined generations are repeated.

But here, chromosome is not a string like so far, but *tree*.

14.1 How we create random tree, and how we evolve trees.

We, first of all, prepare for function set and terminal set. Then we follow the below.

- (1) Choose one function from the *Function Set* at random, and assign it to the root node.
- (2) Assign each of arcs a function or a terminal chosen at random from *Function Set* or *Terminal Set*.
- (3) If the node is a terminal, then the node will not grow any more, but instead become a *leave*. Oterwise repeat (2).
- (4) These are repeated until all the end nodes are terminal set.

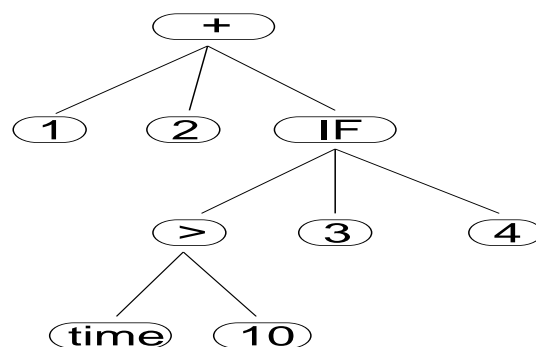
14.2 Evolution of program.

Like the program language *LISP*, some programming languages have a tree-structure.

14.2.1 An Example of Tree representation of a program

Program in LISP can be represent by a tree. The following is a simple example of one instruction from LISP program and its tree structure.

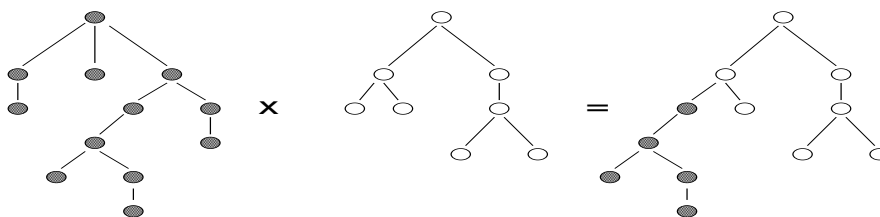
(+ 1 2 (IF (> time 10) 3 4))



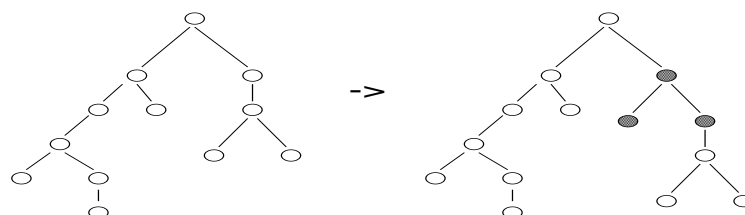
14.2.2 Crossover and Mutation

We now look at how we crossover and mutate two trees. See the figure below.

crossover



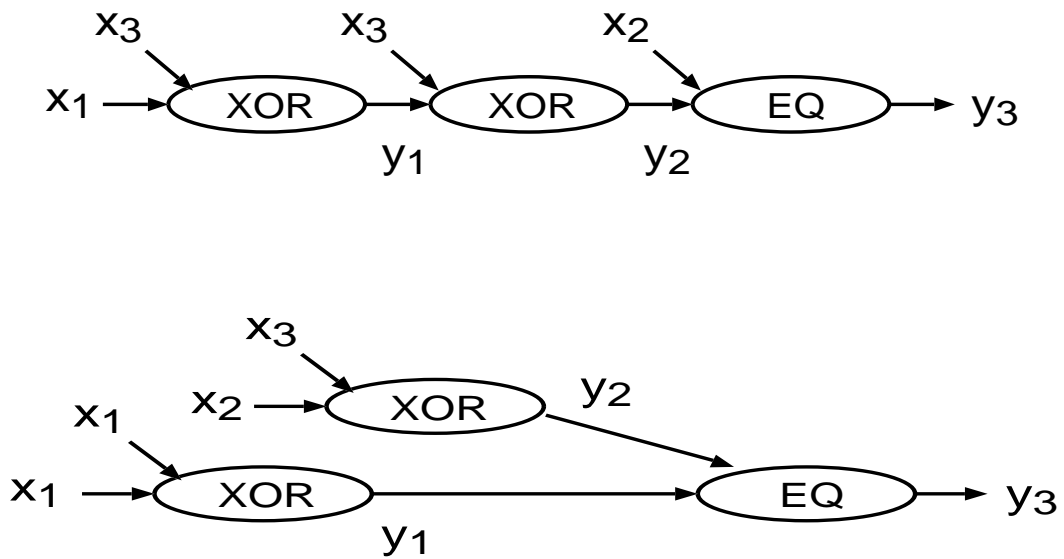
mutation



14.3 Evolution of hardware

14.3.1 Even-n-parity

Sometimes, we need automatic error-detection coding. Assume we need $(n - 1)$ binary bits for our encoding. Then in order to detect an accidental assignment of a code we add one additional bit called parity-check-bit. 0 or 1 will be assigned so that total number of bit 1 becomes even. Hence we can know the incorrectly coded line, though some lucky mistake can avoid this detection.



APPENDIX

I. The other commonly used test-function

Schwefel function

The function below is called Schwefel function and enormous amount of local minimum and the only unique global minimum at the Origin. Search for the global minimum when the function is defined on, say, $x_i \in [-500, 500]$.

$$y = \sum_{i=1}^n (x_i \sin(|x_i|)) \quad (19)$$

You might try to explore this hyper-surface by setting n to 20, for example. Then the surface is defined on 20-dimensional Euclidean space. If, however, you want to know the image of the hyper-surface, see a two dimensional version of this function. The graph of

$$y = x \sin(|x|) \quad (20)$$

is shown in Figure 1.

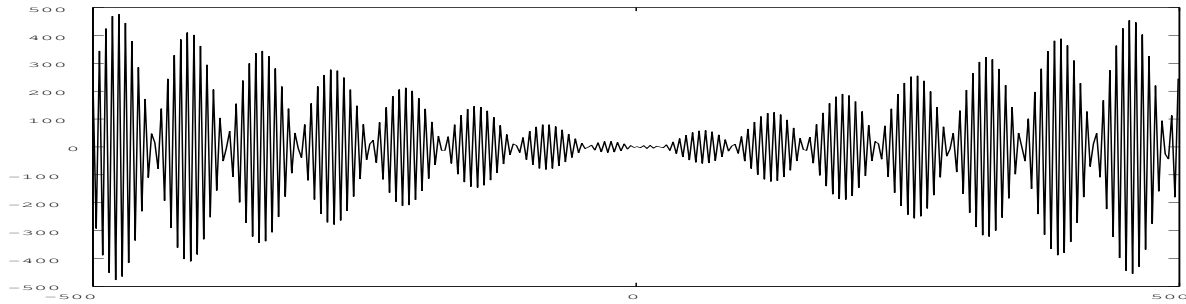


Figure 34: A two-dimensional version of the Schwefel's Function's graph.

14.4 Rastrigin's Function

★ Rastrigin's Function

$$y = nA + \sum_{i=1}^n (x_i^2 - A \cos(2\pi x_i)), \quad x_i \in [-5.12, 5.12].$$

· Ruggedness might be controlled by modifying the value of A .

- A two dimensional example ($n = 1$):

$$y = 3 + x^2 - 3 \cos(2\pi x).$$

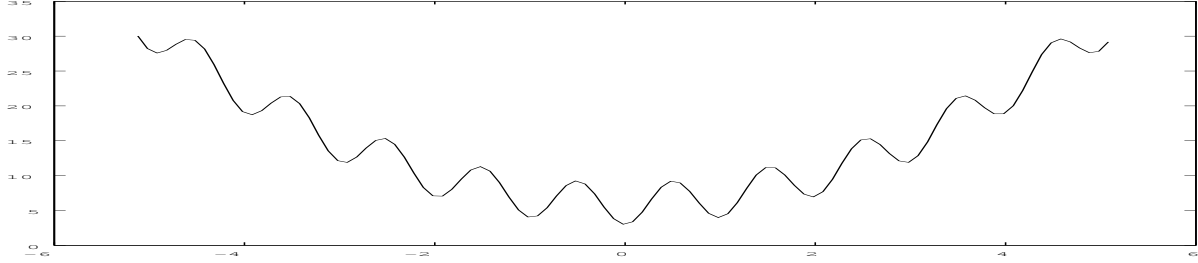


Figure 35: A 2-D version of Rastrigin function

- ★ Griewangk's Function F_8 :

$$y = \sum_{i=1}^n x_i^2/4000 - \prod_{i=1}^n \cos(x_i/\sqrt{i}) + 1, \quad x_i \in [-600, 600].$$

- A two dimensional example:

$$y = x^2/4000 - \cos x + 1.$$

- ★ Ackley's Function F_9 :

$$y = -20 \sum_{i=1}^n \exp(-0.2\sqrt{x_i^2/n}) - \exp((\sum_{i=1}^n \cos 2\pi x_i)/n) + 20 + e,$$

$$x_i \in [-30, 30].$$

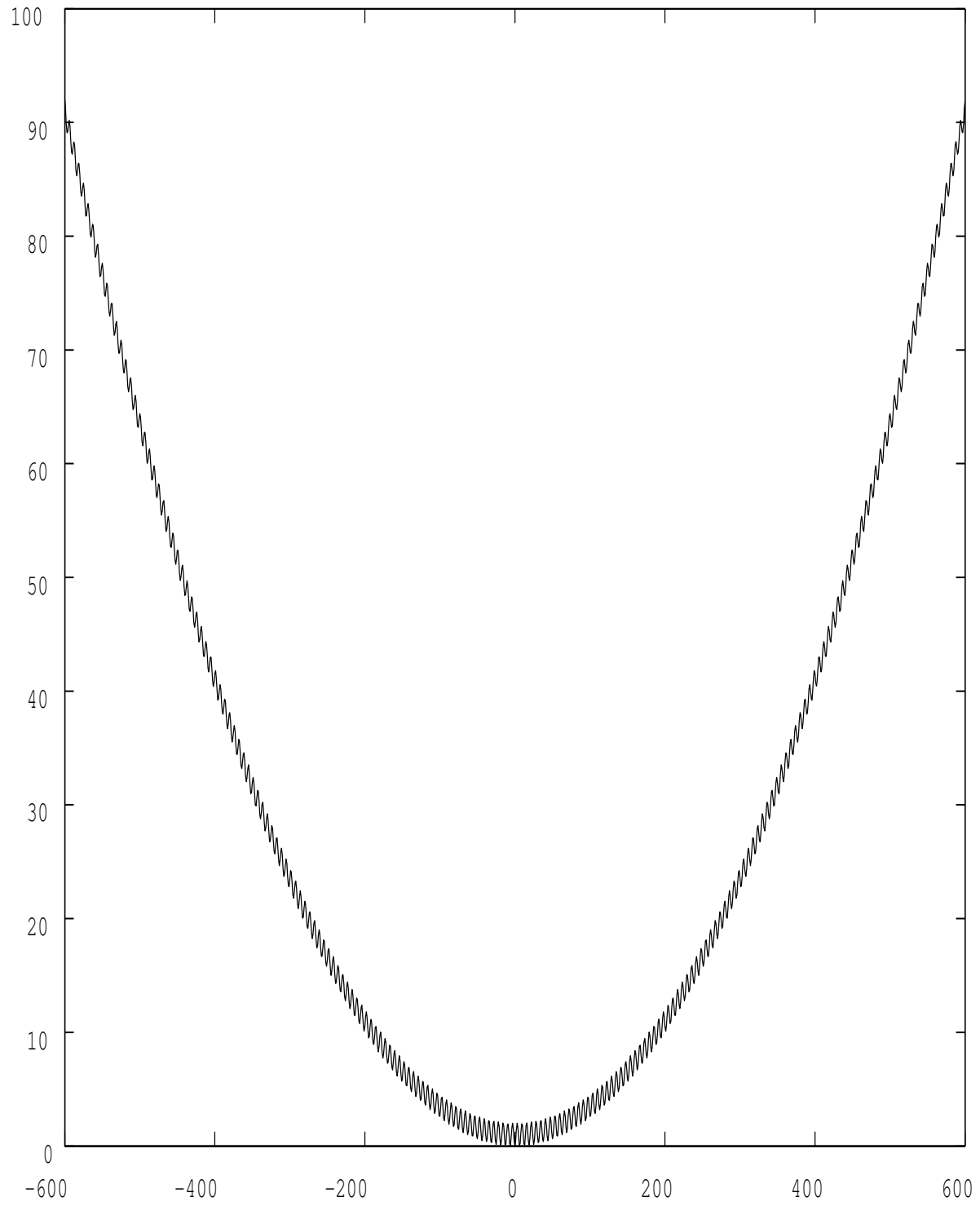


Figure 36: A 2-D version of Griewangk function

- A two dimensional example:

$$y = -20 \exp(-0.2\sqrt{x^2}) - \exp(\cos 2\pi x) + 20 + e.$$

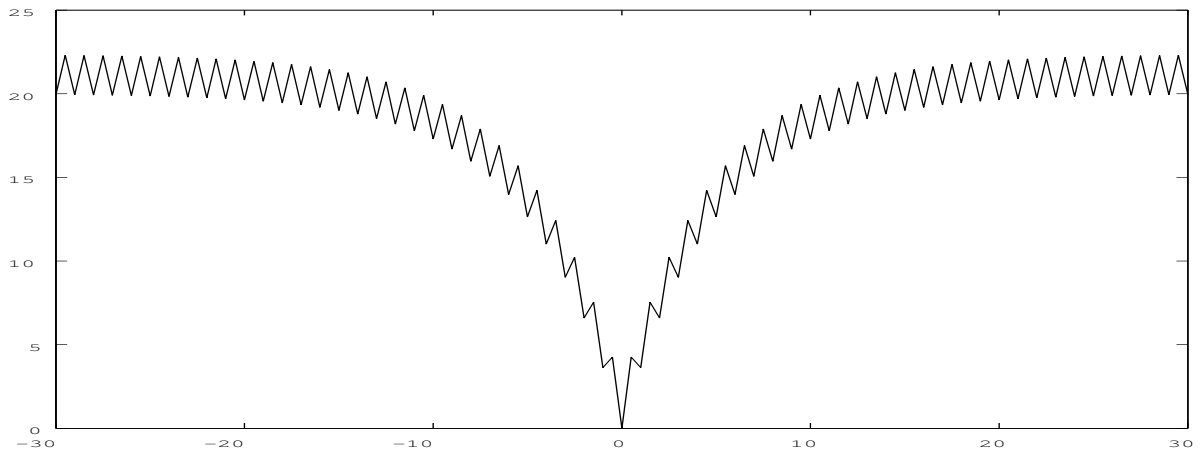


Figure 37: A 2-D version of Ackley function

II. Iris Flower Database

Data for training - $3 \times 40 = 120$ data



Setosa				Versicolor				Virginica			
x_1	x_2	x_3	x_4	x_1	x_2	x_3	x_4	x_1	x_2	x_3	x_4
0.65	0.80	0.20	0.08	0.89	0.73	0.68	0.56	0.80	0.75	0.87	1.00
0.62	0.68	0.20	0.08	0.81	0.73	0.65	0.60	0.73	0.61	0.74	0.76
0.59	0.73	0.19	0.08	0.87	0.70	0.71	0.60	0.90	0.68	0.86	0.84
0.58	0.70	0.22	0.08	0.70	0.52	0.58	0.52	0.80	0.66	0.81	0.72
0.63	0.82	0.20	0.08	0.82	0.64	0.67	0.60	0.82	0.68	0.84	0.88
0.68	0.89	0.25	0.16	0.72	0.64	0.65	0.52	0.96	0.68	0.96	0.84
0.58	0.77	0.20	0.12	0.80	0.75	0.68	0.64	0.62	0.57	0.65	0.68
0.63	0.77	0.22	0.08	0.62	0.55	0.48	0.40	0.92	0.66	0.91	0.72
0.56	0.66	0.20	0.08	0.84	0.66	0.67	0.52	0.85	0.57	0.84	0.72
0.62	0.70	0.22	0.04	0.66	0.61	0.57	0.56	0.91	0.82	0.88	1.00
0.68	0.84	0.22	0.08	0.63	0.45	0.51	0.40	0.82	0.73	0.74	0.80
0.61	0.77	0.23	0.08	0.75	0.68	0.61	0.60	0.81	0.61	0.77	0.76
0.61	0.68	0.20	0.04	0.76	0.50	0.58	0.40	0.86	0.68	0.80	0.84
0.54	0.68	0.16	0.04	0.77	0.66	0.68	0.56	0.72	0.57	0.72	0.80
0.73	0.91	0.17	0.08	0.71	0.66	0.52	0.52	0.73	0.64	0.74	0.96
0.72	1.00	0.22	0.16	0.85	0.70	0.64	0.56	0.81	0.73	0.77	0.92
0.68	0.89	0.19	0.16	0.71	0.68	0.65	0.60	0.82	0.68	0.80	0.72
0.65	0.80	0.20	0.12	0.73	0.61	0.59	0.40	0.97	0.86	0.97	0.88
0.72	0.86	0.25	0.12	0.78	0.50	0.65	0.60	0.97	0.59	1.00	0.92
0.65	0.86	0.22	0.12	0.71	0.57	0.57	0.44	0.76	0.50	0.72	0.60
0.68	0.77	0.25	0.08	0.75	0.73	0.70	0.72	0.87	0.73	0.83	0.92
0.65	0.84	0.22	0.16	0.77	0.64	0.58	0.52	0.71	0.64	0.71	0.80
0.58	0.82	0.14	0.08	0.80	0.57	0.71	0.60	0.97	0.64	0.97	0.80
0.65	0.75	0.25	0.20	0.77	0.64	0.68	0.48	0.80	0.61	0.71	0.72

(to be cont'd to the next page)

(cont'd)

Setosa				Versicolor				Virginica			
x_1	x_2	x_3	x_4	x_1	x_2	x_3	x_4	x_1	x_2	x_3	x_4
0.61	0.77	0.28	0.08	0.81	0.66	0.62	0.52	0.85	0.75	0.83	0.84
0.63	0.68	0.23	0.08	0.84	0.68	0.64	0.56	0.91	0.73	0.87	0.72
0.63	0.77	0.23	0.16	0.86	0.64	0.70	0.56	0.78	0.64	0.70	0.72
0.66	0.80	0.22	0.08	0.85	0.68	0.72	0.68	0.77	0.68	0.71	0.72
0.66	0.77	0.20	0.08	0.76	0.66	0.65	0.60	0.81	0.64	0.81	0.84
0.59	0.73	0.23	0.08	0.72	0.59	0.51	0.40	0.91	0.68	0.84	0.64
0.61	0.70	0.23	0.08	0.70	0.55	0.55	0.44	0.94	0.64	0.88	0.76
0.68	0.77	0.22	0.16	0.70	0.55	0.54	0.40	1.00	0.86	0.93	0.80
0.66	0.93	0.22	0.04	0.73	0.61	0.57	0.48	0.81	0.64	0.81	0.88
0.70	0.95	0.20	0.08	0.76	0.61	0.74	0.64	0.80	0.64	0.74	0.60
0.62	0.70	0.22	0.04	0.68	0.68	0.65	0.60	0.77	0.59	0.81	0.56
0.63	0.73	0.17	0.08	0.76	0.77	0.65	0.64	0.97	0.68	0.88	0.92
0.70	0.80	0.19	0.08	0.85	0.70	0.68	0.60	0.80	0.77	0.81	0.96
0.62	0.70	0.22	0.04	0.80	0.52	0.64	0.52	0.81	0.70	0.80	0.72
0.56	0.68	0.19	0.08	0.71	0.68	0.59	0.52	0.76	0.68	0.70	0.72
0.65	0.77	0.22	0.08	0.70	0.57	0.58	0.52	0.87	0.70	0.78	0.84

Data for evaluating the system after training - $3 \times 23 = 69$ data



Setosa				Versicolor				Virginica			
x_1	x_2	x_3	x_4	x_1	x_2	x_3	x_4	x_1	x_2	x_3	x_4
0.63	0.80	0.19	0.12	0.70	0.59	0.64	0.48	0.85	0.70	0.81	0.96
0.57	0.52	0.19	0.12	0.77	0.68	0.67	0.56	0.87	0.70	0.74	0.92
0.56	0.73	0.19	0.08	0.73	0.59	0.58	0.48	0.73	0.61	0.74	0.76
0.63	0.80	0.23	0.24	0.63	0.52	0.48	0.40	0.86	0.73	0.86	0.92
0.65	0.86	0.28	0.16	0.71	0.61	0.61	0.52	0.85	0.75	0.83	1.00
0.61	0.68	0.20	0.12	0.72	0.68	0.61	0.48	0.85	0.68	0.75	0.92
0.65	0.86	0.23	0.08	0.72	0.66	0.61	0.52	0.80	0.57	0.72	0.76
0.58	0.73	0.20	0.08	0.78	0.66	0.62	0.52	0.82	0.68	0.75	0.80
0.67	0.84	0.22	0.08	0.65	0.57	0.43	0.44	0.78	0.77	0.78	0.92
0.63	0.75	0.20	0.08	0.72	0.64	0.59	0.52	0.75	0.68	0.74	0.72
0.62	0.70	0.22	0.04	0.80	0.52	0.64	0.52	0.81	0.70	0.80	0.72
0.56	0.68	0.19	0.08	0.71	0.68	0.59	0.52	0.76	0.68	0.70	0.72
0.65	0.77	0.22	0.08	0.70	0.57	0.58	0.52	0.87	0.70	0.78	0.84
0.63	0.80	0.19	0.12	0.70	0.59	0.64	0.48	0.85	0.70	0.81	0.96
0.57	0.52	0.19	0.12	0.77	0.68	0.67	0.56	0.87	0.70	0.74	0.92
0.56	0.73	0.19	0.08	0.73	0.59	0.58	0.48	0.73	0.61	0.74	0.76
0.63	0.80	0.23	0.24	0.63	0.52	0.48	0.40	0.86	0.73	0.86	0.92
0.65	0.86	0.28	0.16	0.71	0.61	0.61	0.52	0.85	0.75	0.83	1.00
0.61	0.68	0.20	0.12	0.72	0.68	0.61	0.48	0.85	0.68	0.75	0.92
0.65	0.86	0.23	0.08	0.72	0.66	0.61	0.52	0.80	0.57	0.72	0.76
0.58	0.73	0.20	0.08	0.78	0.66	0.62	0.52	0.82	0.68	0.75	0.80
0.67	0.84	0.22	0.08	0.65	0.57	0.43	0.44	0.78	0.77	0.78	0.92
0.63	0.75	0.20	0.08	0.72	0.64	0.59	0.52	0.75	0.68	0.74	0.72