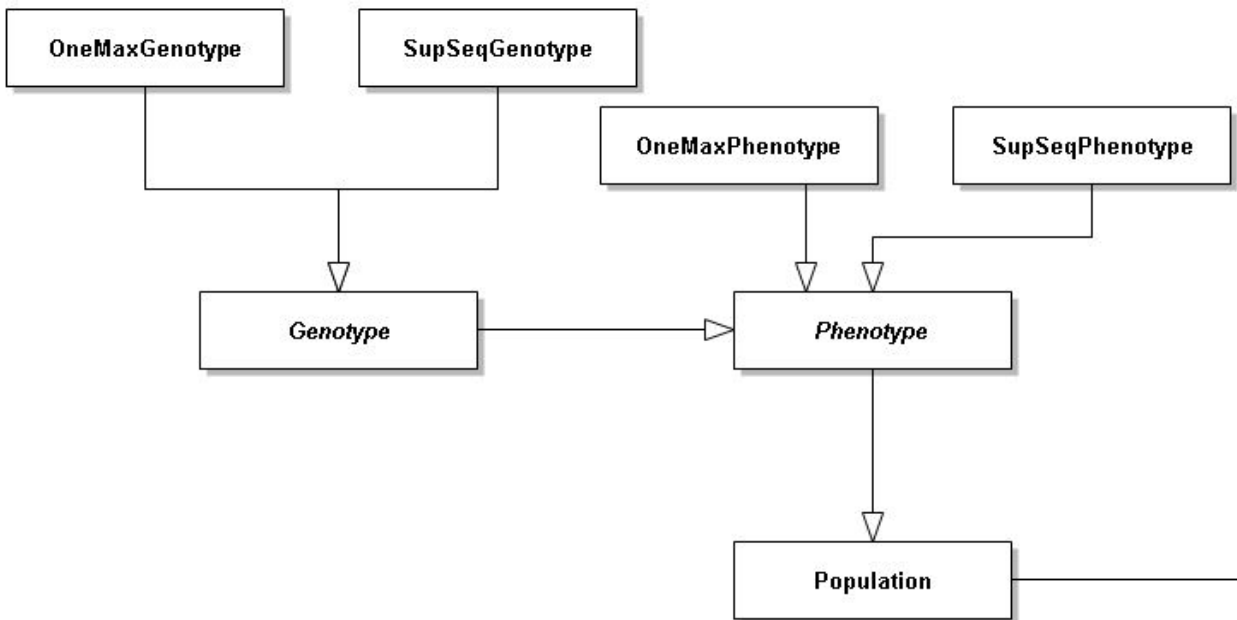


# Searching for Surprising Sequences with a Genetic Algorithm

By Kristoffer Hagen

## 1 The model

I used the framework that I built for the first part of the assignment here as well.



The SupSeqGenotype contains the basic genetic information such as the amount of symbols and the size of the string, together with an array of integers that makes up the genes. The phenotype contain the aforementioned genotype together with methods for mating.

## 2 The fitness function

The fitness function for this problem is basically the following;

I take a certain distance between two integers in the genes, then I go through the array of genes and add in pairs the integers that match this distance to a list. For the genes 1 2 3 4 3 2 1 and the distance 1, I would add the pairs 1 3, 2 4, 3 3, 4 2 and 3 1. Then I will go through the list I just made and look for duplicates, if any are found I

```
//@Override
public double findGlobalFitness() {
    ArrayList<IntPair> pairs;
    int conflict = 0;

    for (int i = 1; i < size; i++) { //distances'
        pairs = new ArrayList<IntPair>();
        for (int j = 0; j < size-i; j++) { //navigation
            int value1 = vector[j];
            int value2 = vector[j+i];
            pairs.add(new IntPair(value1, value2));
        }
        for (int k = 0; k < pairs.size(); k++) {
            IntPair ip = pairs.get(k);
            for (int l = k+1; l < pairs.size(); l++) {
                if(ip.isEqual(pairs.get(l)))
                    conflict++;
            }
        }
    }
}
```

will subtract one from the current fitness and move on.

Say I had list of pairs that looked like this: 2 1, 2 0, 0 1, 2 1. There is one duplicate here.

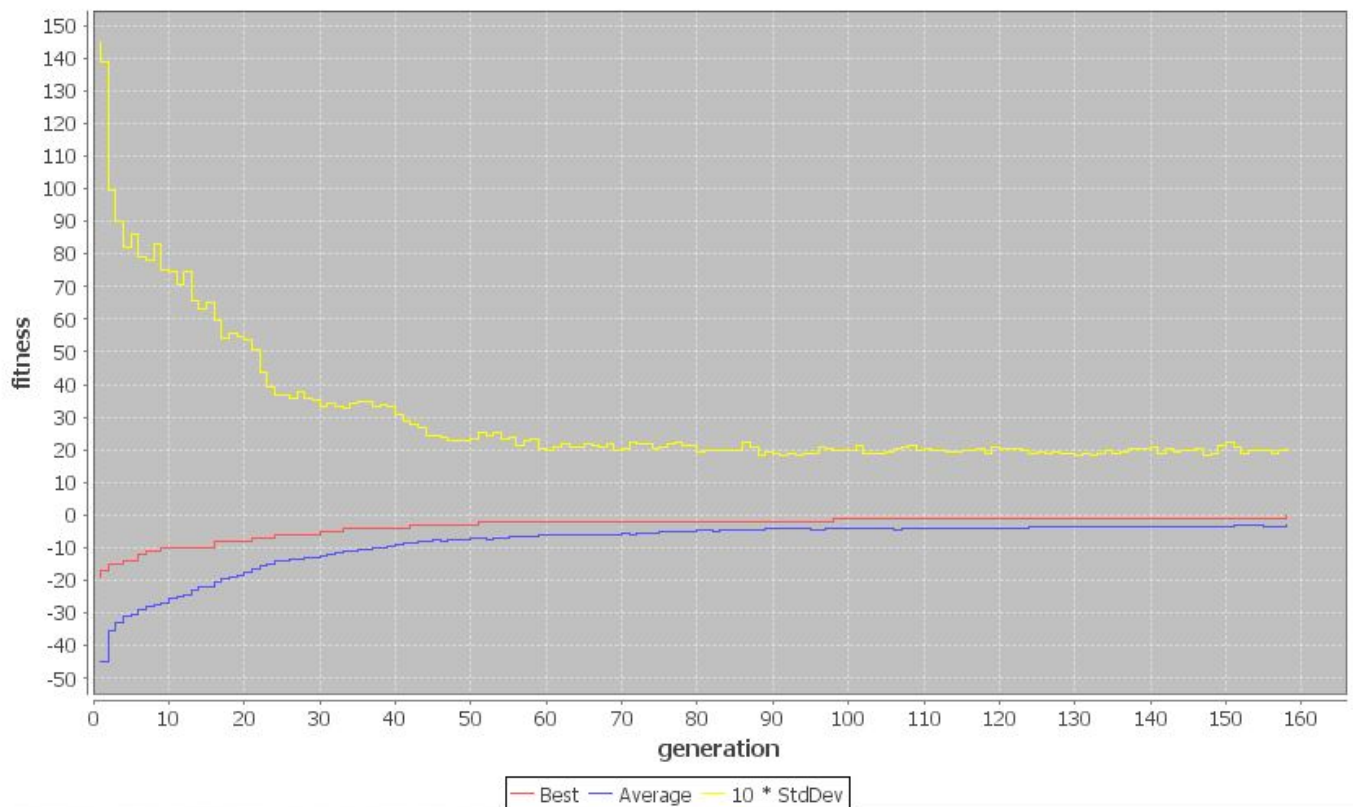
For locally surprising sequences I would just look at distance=0, but for globally I would look through all possible distances.

### 3 The table

For all these runs, in addition to the numbers shown I used the following: Generational adult selection, proportionate parent selection when population is 500, SigmaPieSelector when population is 50. I used a random split crossover (cut the genes at a random point, use one parent for one half, the other parent for the other). Mutation rate of 0.2 (Which I in retrospect wish was lower.)

Fitness plot and output from S=19 L=48 globally surprising:

**Fitness Plot**



```
Generation 153
Generation 154
Generation 155
Generation 156
Generation 157
Flawless string found!
Vector:
11 13 12 14 7 8 0 6 17 11 4 17 8 18 3 2 13 9 5 16 16 6 1 15 7 3 5 2 14 0 3 10 17 15 13 18 4 10 4 12 11 1 0 12 8 14 9
Flawless string end!
Generation 158
```

Locally Surprising				Globally Surprising			
Symbols	Population	Generations	Length	Sequence	Generations	Length	Sequence
3	50	17	10	1100201221	1	7	0121102
4	50	18	17	32122033110013023	160	10	2310120332
5	50	80	26	24032043013341221442311002	59	12	100432414230
6	50	288	37	0223240044541433051355201121034253150	177	16	513150430221453
7	50	54	49	5650603464204401523324135300251436312611662210554	281	17	16403055621531024
8	500	390	64	7422625400160530646110712032134414376736335045231517247027566557	115	20	24775106423631435072
9	500	484	80	20855363540160057452714317076410251221567323387818848377268662804424650611303475	186	23	08136546775341028582631
10	500	375	98	38934827317162205503002806083740754941976795132599856109018152691239644635772114786684533658704243	218	25	2967145304008625739824317
11	500	156	120		115	27	240469810102051710387493450126
12	500	380	136		201	30	27211561101038609910173411128950467
13	500	250	160		126	32	122101254811145703286617312909311510426
14	500	480	186		280	35	113475122819121711003106139841321312537481115
15	500	420	210		307	38	137108111311129401232198651414512010673120119528134
16	500	357	235		366	39	121068114145713030155121119131351598123710146102411610
17	500	272	260		238	42	13214117135812463314101901071511130201689161517414651231548
18	500	235	280		444	45	41215917951634125107126014781711691316131702111515414161251310873
19	500	218	310		158	48	1113121478061711417818321395161661157352140310171513184104121110128149
20	500	244	340		316	50	71161574181217261089131403955191412181161913319411968280151210171141627513

## 5 Theory

Given a 20-gene chromosome and onemax, locally and globally surprising sequence with 9 symbols. I would rank then globally, locally, onemax in order of difficulty. Onemax really isn't very problematic as every single improvement is always one step closer to a solution. The genes are all individual, or not connected. Locally surprising sequence is harder because when you modify one gene, that also affects the following one in the sequence. In the string 1 2 3, if the 2 is modified, that changes both the 1 2 pair and the 2 3 pair. This makes it harder to solve than onemax.

For globally surprising sequence it is taken to the next level again. Every gene is now connected to up to many other genes.

In the sequence 0 1 2 1 1 0 2, for example. The first 2 is a part 6 pairs and when it is changed it also changes the fitness of the sequence in a less predictable way, this makes globally surprising sequence the hardest one to solve.