

Genetic Algorithm Project

Team 38 - GENTOOmen

Roll No.	Name
2019114001	Akshett Rai Jindal
2019111031	Zishan Kazi

Task

The task was to get the weights of the features for a model, using [Genetic Algorithms](#) provided that the weights of the features corresponding to an overfit model are given to us.

Running the Program

```
$ python3 main.py
```

The output of the program is `output.json` which contains the weight vectors of the last generation.

About G.A.

A **genetic algorithm** is a search heuristic that is inspired by Charles Darwin's theory of natural evolution.

We have implemented it as follows:

- Initialize a population
- Determine fitness of population
- Until convergence *repeat* the following:
 - **Select parents** from the population
 - **Crossover** to generate children vectors
 - **Mutate** the new population
 - Calculate **fitness** for new population

Each population contains multiple individuals where each individual represents a point in search space and possible solution.

Every individual is a vector of size 11 with 11 floating point values in the range of $[-10, 10]$. These 11 values represent the weights of the different features of the 11 features of the data.

Selection

Selection is basically taking the individuals with best fitness and pass their genes/bits onto next generations to maintain good fitness scores.

Crossover

Here two vectors generate an offspring. These two vectors are taken from `selection` function. The two vectors are then manipulated to generate an offspring.

Mutations

We change each bit of the vector with a probability to diversify the population and to prevent similar vectors which can led to convergence.

Code

The first population is created using an initial given overfit vector. Copies of this vector are made on which we mutate to generate a population of size `POPULATION_SIZE` .

For each vector, at every index mutation is performed with a probability of $6/10$. Then the value at that index is replaced with the value at the overfit vector at that index multiplied by some factor chosen uniformly between (0.7, 1.3).

The **fitness** of the population is an linear combination of the train error and the validation error.

After the population is initialized, a mating pool is made containing the top `MATING_POOL_SIZE` parents sorted according to their fitness.

```
population_fitness = population_fitness[np.argsort(population_fitness[:, -1])]
mating_pool = population_fitness[:MATING_POOL_SIZE]
```

Then parents are uniformly chosen from the mating pool as follows:

```
index1 = random.randint(0, MATING_POOL_SIZE-1)
index2 = random.randint(0, MATING_POOL_SIZE-1)

parent1 = mating_pool[index1]
parent2 = mating_pool[index2]
```

Then offsprings are produced from the parents by applying **Simulated Binary Crossover** on them. This is followed by mutation of chromosomes, details of which are given [here](#).

The new population is created by choosing `POPULATION_SIZE` - `PASSED_FROM_PARENTS` top children generated and `PASSED_FROM_PARENTS` top parents.

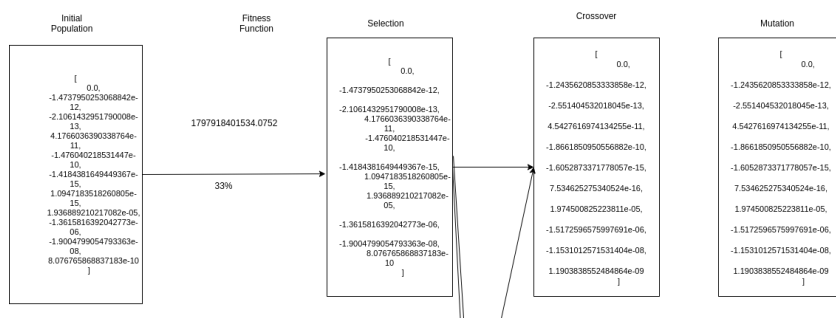
```
def create_next_gen(parents_fitness, children):
    child_fitness = get_population_fitness(children)[: (POPULATION_SIZE-PASSED_FROM_PARENTS)]
    parents_fitness = parents_fitness[:PASSED_FROM_PARENTS]
    generation = np.concatenate((parents_fitness, child_fitness))
    generation = generation[np.argsort(generation[:, -1])]
    return generation
```

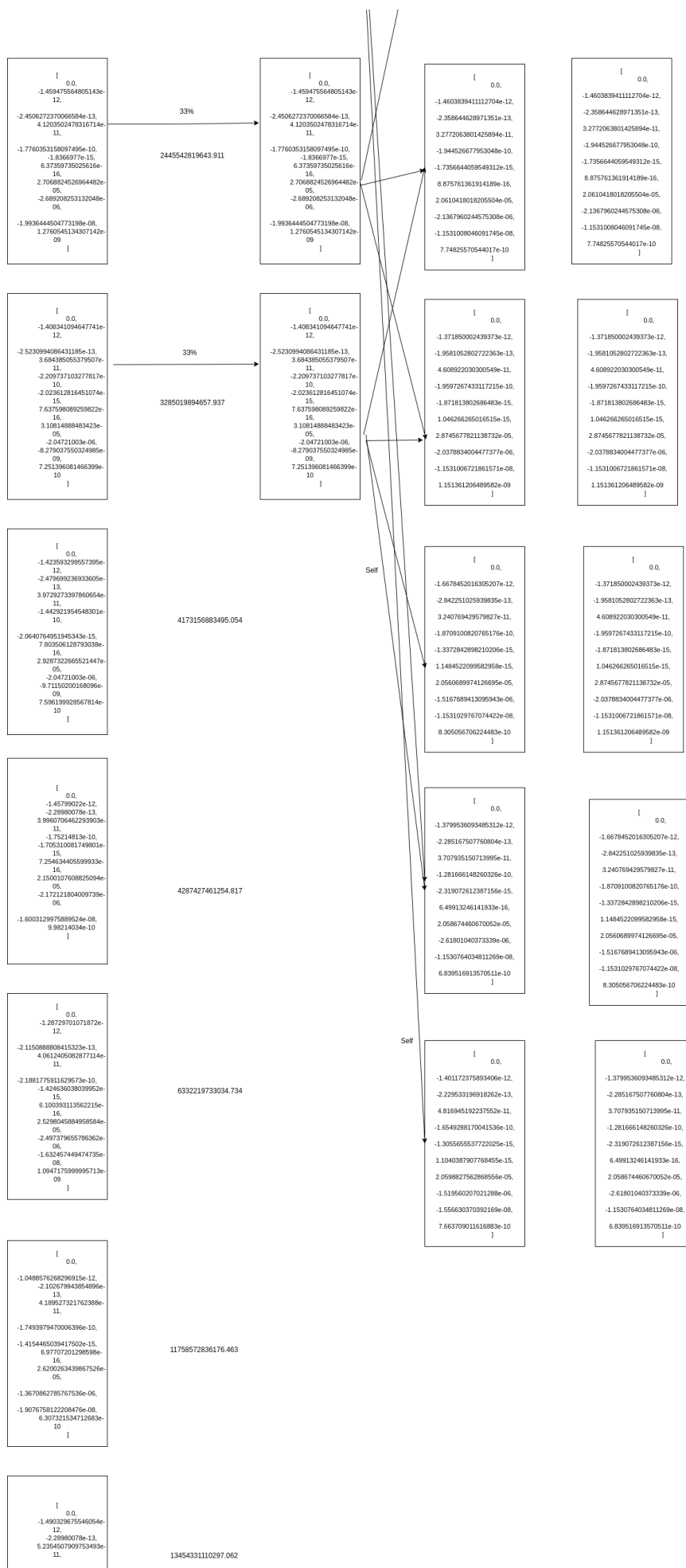
This process is repeated and the values are stored in a JSON object file from which it is read for next run of the program.

As you can see the code is **vectorized** and **completely modular** as separate function have been written for each significant step.

Iteration Diagrams

First Iteration





```
[
  0.0,
  -1.8722194120180694e-11,
  -2.411262591141495e-15,
  8.5294406e-16,
  2.143690960243266e-05,
  -2.04721003e-06,
  -1.59792834e-08,
  9.104897150956366e-10
]
```

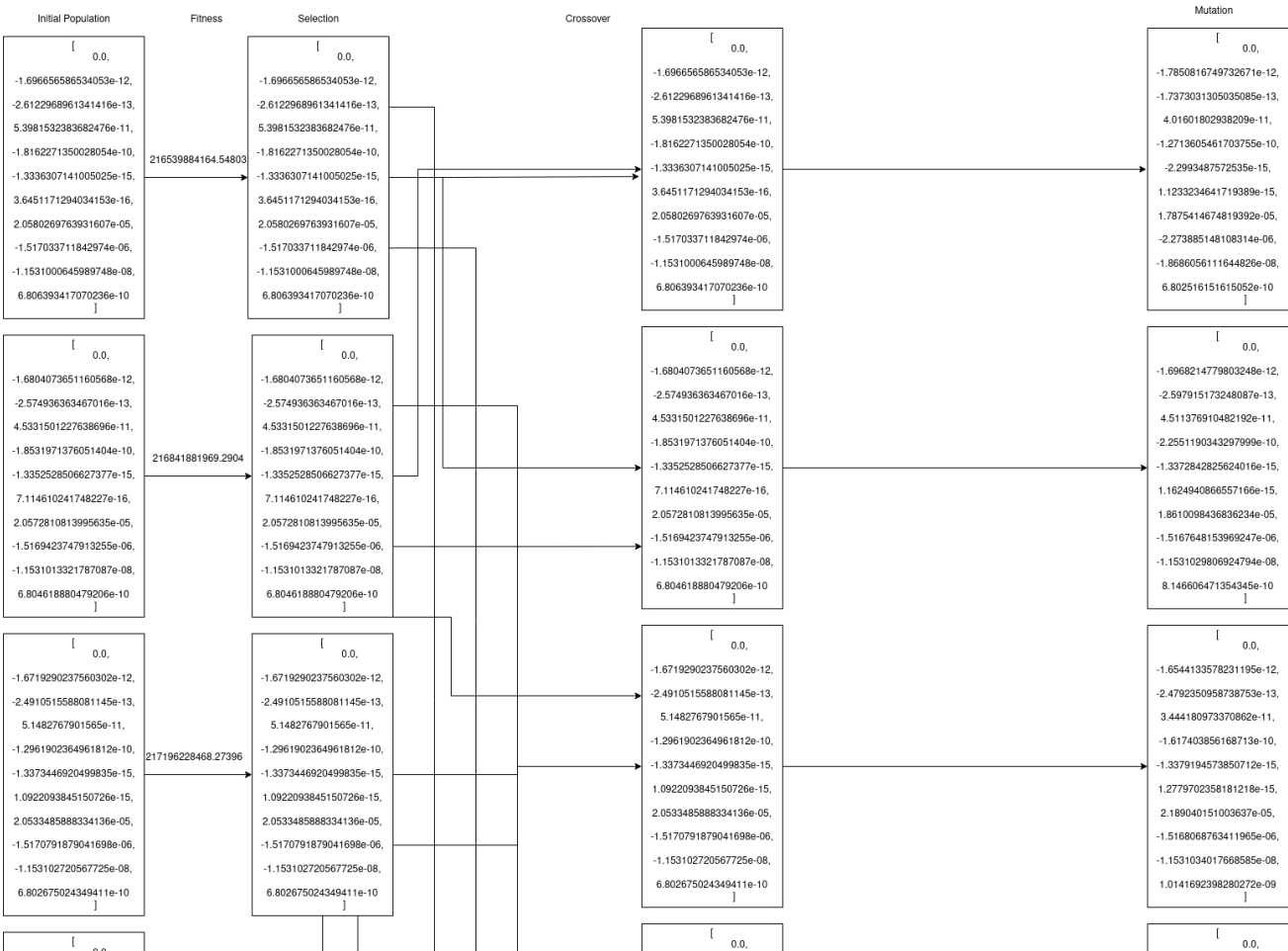
```
[
  0.0,
  -1.4201706752360025e-12,
  -2.1710402404103709e-13,
  5.115051527877852e-11,
  -1.75214813e-10,
  -1.794035914204154e-15,
  8.993648447926866e-16,
  3.019308074907925e-05,
  -2.5274155443288656e-06,
  -1.613702423903899e-08,
  1.2923362109872806e-09
]
```

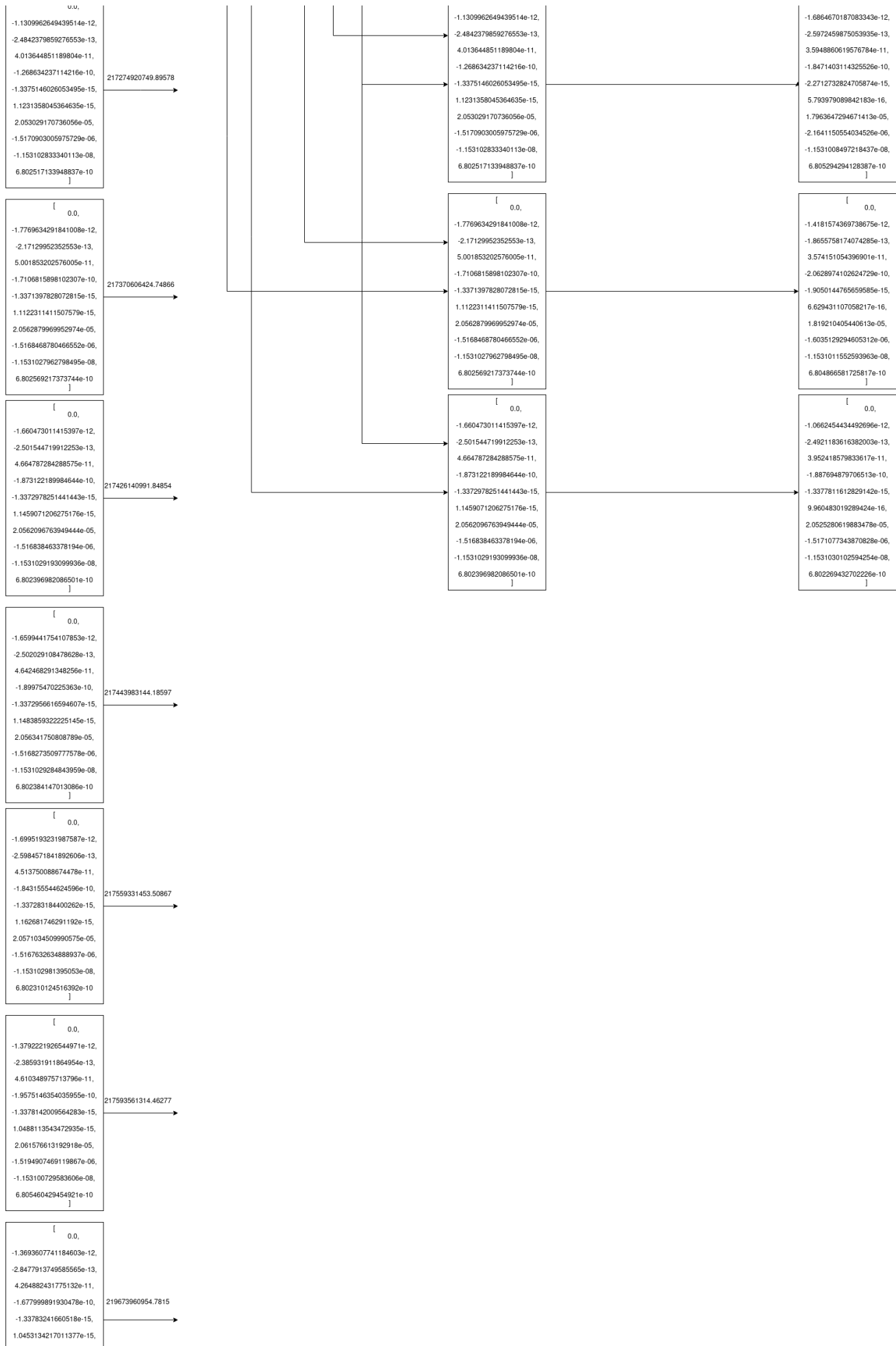
```
[
  0.0,
  -1.4743657583252537e-12,
  -2.1018964150770835e-13,
  4.178919126787864e-11,
  -1.75214813e-10,
  -1.4110980395430349e-15,
  6.96395851364423e-16,
  1.5458390338717505e-05,
  -1.3537868720888707e-06,
  -1.9110452623441447e-08,
  1.060946375407254e-09
]
```

34111041175538.746

87245049445526.34

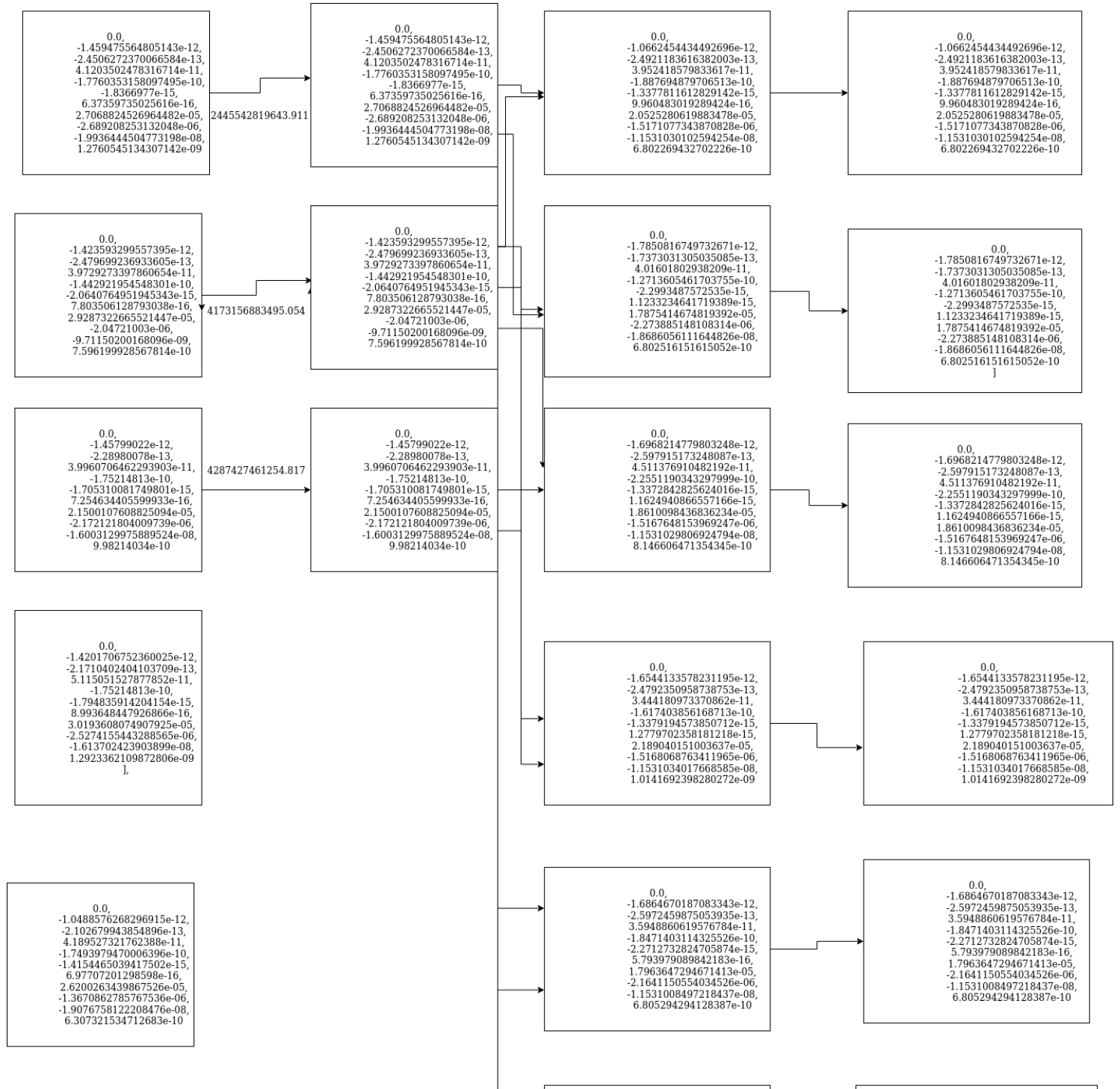
Second Iteration

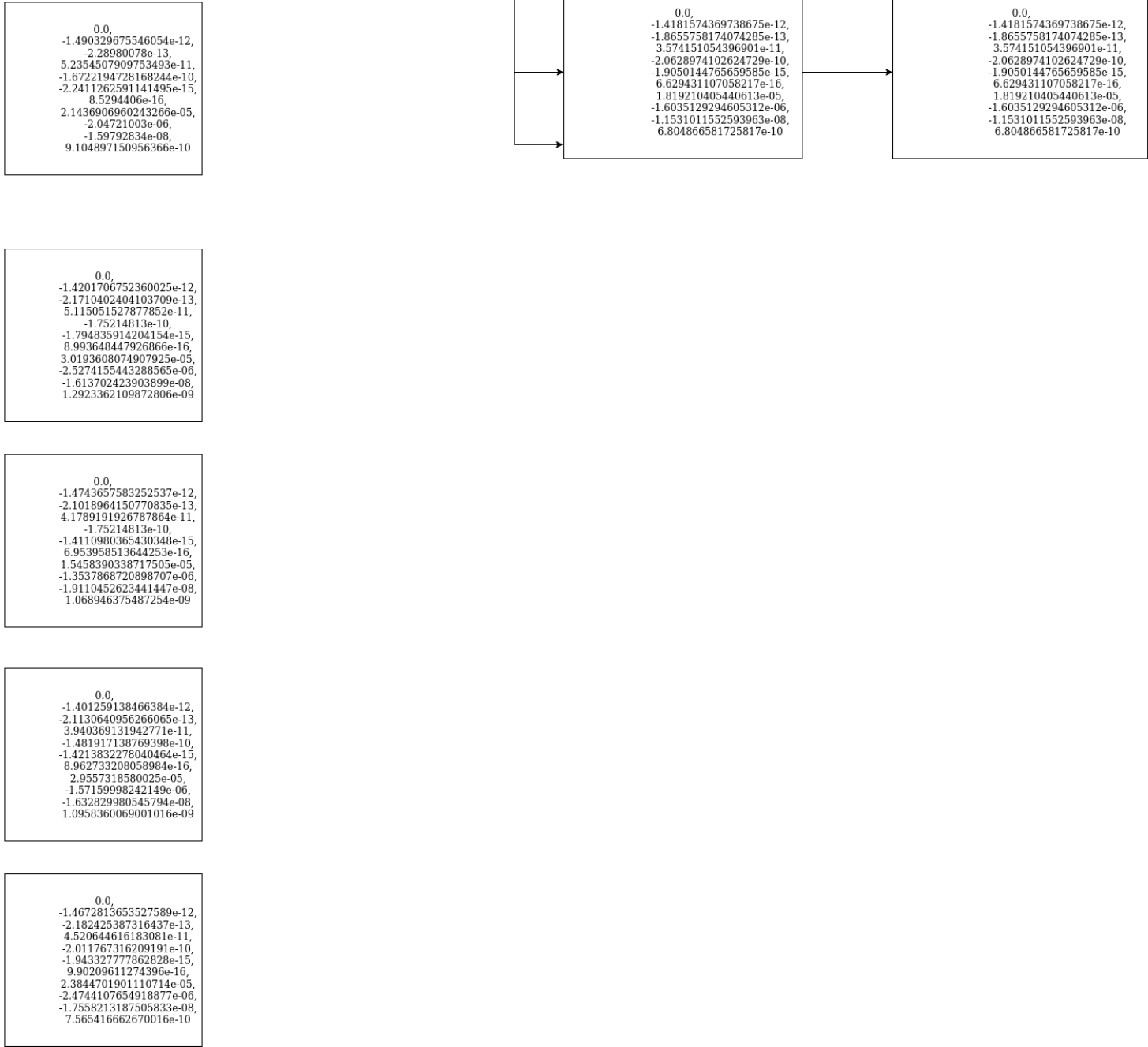




2.0617605074802182e-05,
-1.5161777650526904e-06,
-1.1531006523389213e-08,
6.805568495479519e-10
]

Third Iteration





Initial-Population

```
def initial_population():  
    first_population = []  
    for i in range(POPULATION_SIZE):  
        first_population.append(np.copy(first_parent))  
  
    for i in range(POPULATION_SIZE):  
        for j in range(VECTOR_SIZE):  
            vary = 0  
            mutation_prob = random.randint(0, 100)  
            if mutation_prob < 60:  
                vary = 1 + random.uniform(-0.05*j, 0.05*j)  
                rem = overfit_vector[j]*vary  
  
            if abs(rem) < 10:  
                first_population[i][j] = rem  
            elif abs(first_population[i][j]) >= 10:
```

```

        first_population[i][j] = random.uniform(-1,1)
    return first_population

```

We used overfit vectors provided to us as the initial population. We applied variation on the vectors (mutation). We changed every bit of every vector with a probability.

Fitness

```

def get_population_fitness(population):
    fit = np.empty((POPULATION_SIZE, 3))

    for i in range(POPULATION_SIZE):
        error = get_errors(SECRET_KEY, list(population[i]))
        fit[i][0] = error[0]
        fit[i][1] = error[1]
        fit[i][2] = abs(error[0]*TRAIN_RATIO + error[1])

    ret = np.column_stack((population, fit))
    ret = ret[np.argsort(ret[:, -1])]
    return ret

```

`get_errors` requests are sent to the server which returns the train error and validation error for every vector in the population.

$\text{Fitness} = \text{Train Ratio} * \text{Train Error} + \text{Validation Error}$

Train ratio was tested and selected so as to obtain less Train error, Validation error and better rank on leaderboard. It was tested on following values:

Train factor = 0.7

Idea behind it was to give more weightage to validation error than train error. But on testing, it gave more error. We also tested on lower values but the result wasn't positive. So we tested on greater values than 0.7

Train factor = 1

On testing, train factor=1 gave the best output. Giving equal weightage to both train and validation error.

Simulated-Binary-Crossover

Simulated binary crossover generates tow children from the two parents which follow the following equation while giving us control over the variation by manipulating the distribution index value.

$$\frac{x_1^{\text{new}} + x_2^{\text{new}}}{2} = \frac{x_1 + x_2}{2}$$

A random number is choosed u between $[0, 1)$. The distribution index is assigned desired value and then The crossover is done by choosing a random number in the range $[0, 1)$. The distribution index is assigned its value and then β is calculated:

$$\beta = \begin{cases} \frac{1}{(2u)\eta_c + 1} & u \leq 0.5 \\ \left(\frac{1}{2(1-u)}\right) \frac{1}{\eta_c + 1} & \text{Otherwise} \end{cases}$$

Distribution index that determines how far children go from parents. The greater its value the closer the children are to parents.

The distribution index is a value between $[2, 5]$ and the offsprings are calculated as follows:

$$x_1^{\text{new}} = 0.5[(1 + \beta)x_1 + (1 - \beta)x_2]$$

$$x_2^{\text{new}} = 0.5[(1 - \beta)x_1 + (1 + \beta)x_2]$$

The code is as shown:

```
def crossover(parent1, parent2):
    child1 = np.array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0])
    child2 = np.array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0])

    r1 = random.random()
    n = 3

    beta = (2 * r1)**((n + 1)**-1)
    if (r1 >= 0.5):
        beta = ((2*(1-r1))**(-1))**((n + 1)**-1)

    p1 = np.array(parent1)
    p2 = np.array(parent2)
    child1 = 0.5*((1 + beta) * p1 + (1 - beta) * p2)
    child2 = 0.5*((1 - beta) * p1 + (1 + beta) * p2)

    return [child1, child2]
```

We varied the **Distributed Index** value depending on the population.

Mutation

```
def mutation(vector):
    for i in range(VECTOR_SIZE):
        mutation_prob = random.randint(0, 100)
        if mutation_prob < 50:
            var = 1 + random.uniform(-0.3, 0.3)
            rem = overfit_vector[i]*var # multiply to create variations
            if abs(rem) <= 10:
                vector[i] = rem
    return vector
```

- Our mutations are probabilistic in nature. For the vector, at every index a mutation is decided to be performed with a probability of $1/2$.
- This was even tested for smaller values and larger, but due to limit of calls we stick to $1/2$ probability which gave us decent results.

Hyperparameters

Population size

The `POPULATION_SIZE` parameter is set to **10**. Due to limited api calls it is set so low. Increasing this size will definitely improve the working the code.

Mating pool size

The `MATING_POOL_SIZE` variable is **8**. We sort the parents by the fitness value and choose the top 8 that are selected for the mating pool.

Number of parents passed down to new generation

We set this variable to **8**.

- We kept the value small when we were just starting out and were reliant on more variations in the children to get out of the overfit. We did not want to forcefully bring down more parents as that would waste a considerable size of the new population.
- When we were unsure of how our mutations and crossover were performing or when we would change their parameters, we would increase this variable. We did this so that even if things go wrong, our good vectors are still retained in the new generations.

Distribution index (Crossover point)

This parameter was applied in the `Simulated Binary Crossover`. It determines how far children go from parents. The greater its value the closer the children are to parents. The Distribution Index is changed according to need. If the vectors felt too alike, the value was lowered, otherwise if they seemed too different, it was changed to 5.

Mutation Range

We varied our mutation range drastically throughout the assignment.

- We made the variation as little as between factors of (0.09, 1.01) for when we had to fine tune our vectors. We did this when we were confident we had a good vector and excessive mutations were not helping. So we tried small variations, to get its best features.
- When our vectors would reach a local minima - we would mutate extensively to get out of the minima. The factor for multiplication could vary anywhere from (0.3, 1.7).
- We would even assign random numbers at times to make drastic changes when no improvement was shown by the vectors. We did this initially when we first ran the 0 vector and it helped us achieve good results.

Heuristics

- **Initial vector:** Before we used overfit vector as initial population. After changing this to probabilistic mutation of every bit of overfit vector, we saw some better results.
- **Probabilistic mutation:** Earlier we were mutating on one index only. But we changed our code to mutate each index with a probability of $1/2$, this brought more variation in the genes and worked well for our populations.
- **Varying mutations for different indices :** After some days into assignment, we decreased the probability of mutation of each bit. This was because the older generation were coming out to be better than newer. So tried to make newer generation not a lot varying from older ones.

Generation_i.json

The trace for the final submitted vectors has all the generations in which they occurred.

The format is as follows,

```
[
  {
    "Generation": $genNumber,
    "Population": [
      $populationVector1,
      ...
    ],
    "Details": [
      {
        "Child Number": $childNumber,
        "Parent 2": $secondParentVector,
        "Parent 1": $firstParentVector,
        "After Crossover": $crossoverOfParents,
        "After Mutation": $mutationOfCrossoverVector
      }
    ]
  }
]
```

```
}  
]
```

As 8 parents are brought down to new generation, the last 8 children (when sorted by fitness) are not included in the new population for the next generation. s