# MDL Project Report

Team Name: Sentinels

Team Number : 25

Team Members:
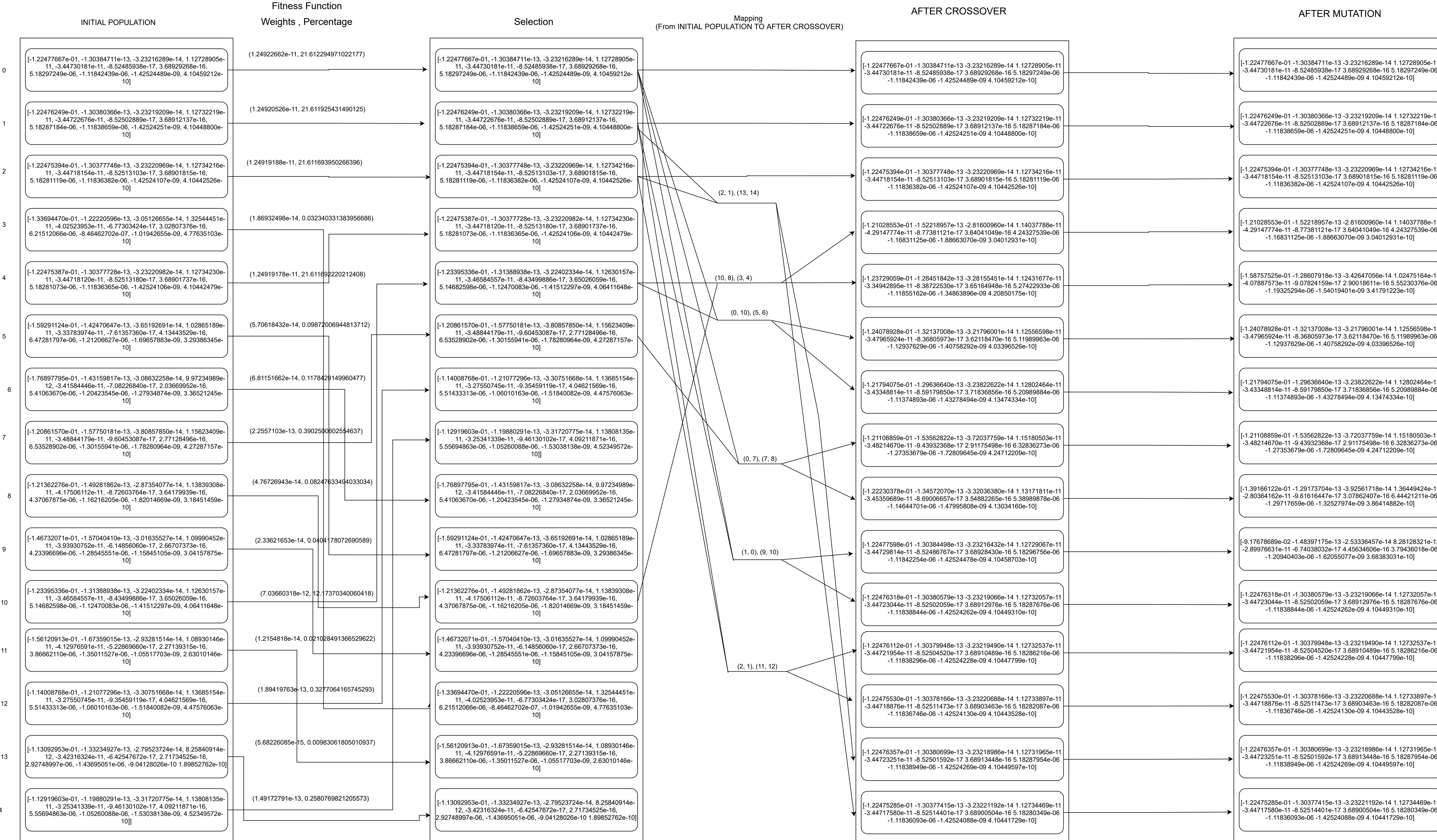
- Pratyanshu Pandey (2019101025)
- Bhaskar Joshi (2019111002)

## Genetic Algorithm

1. The algorithm begins with taking the overfit vector as the base vector and creating an initial population with it. The size of the initial population was 5 then was further reduced to 15.
2. Every feature in an individual in the initial population is nothing but the corresponding feature in the overfit vector multiplied by a random number between 0 and 2 both inclusive. This is done in order to reduce the search space from [-10,10] to a smaller range of values with the overfit vector as a guideline.
3. Thus we have an initial population for a generation. Now we query the server for the train and validation errors and calculate their fitness from the fitness function.
4. Now the vectors are sorted according to fitness in decreasing order. The higher the fitness the better the vector. The top 5 vectors are transferred directly to the population for the next generation. Let's call it the next population.
5. Other vectors for the next population are created by choosing 2 parents from the population and then do crossover on them to create 2 new children. This is done in a loop till the next population has the same size as the initial population.
6. By some probability mutation is applied on individuals of the next population using mutation function.
7. Now the next population becomes the initial population for the next generation.
8. Multiple generations are run till the solution converges and no substantial difference is to be found in new populations.
9. There are multiple heuristics and tricks used in combination with the genetic algorithm to get the result. This was just an overview. See the next sections for more details.

## Sample Iteration of Genetic Algorithm

Note : The generations displayed are Generation 250-252 of file generations_1.txt

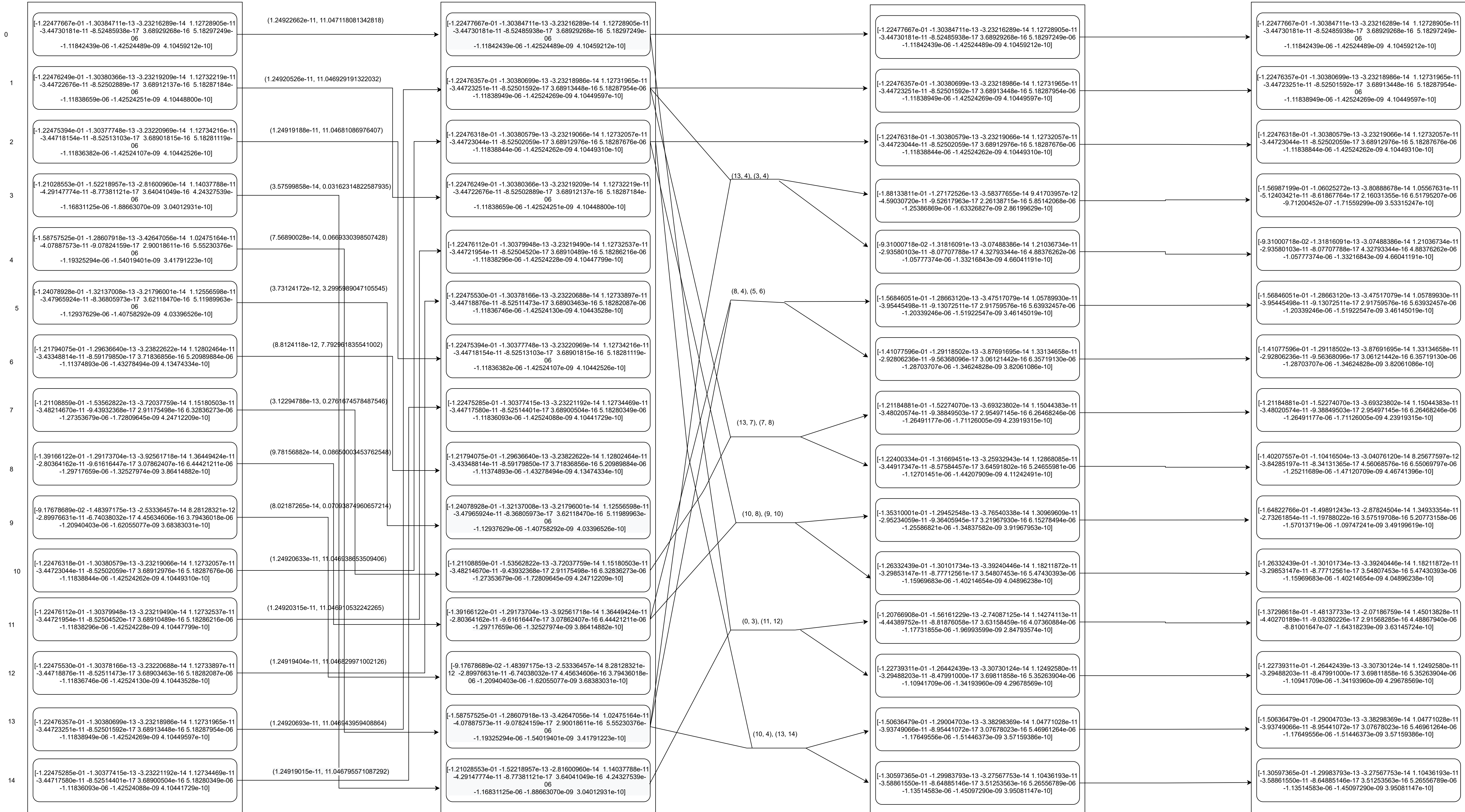Genetic algorithm flowchart showing the progression from INITIAL POPULATION through Fitness Function (Weights, Percentage), Selection, Mapping (From INITIAL POPULATION TO AFTER CROSSOVER), AFTER CROSSOVER, and AFTER MUTATION.
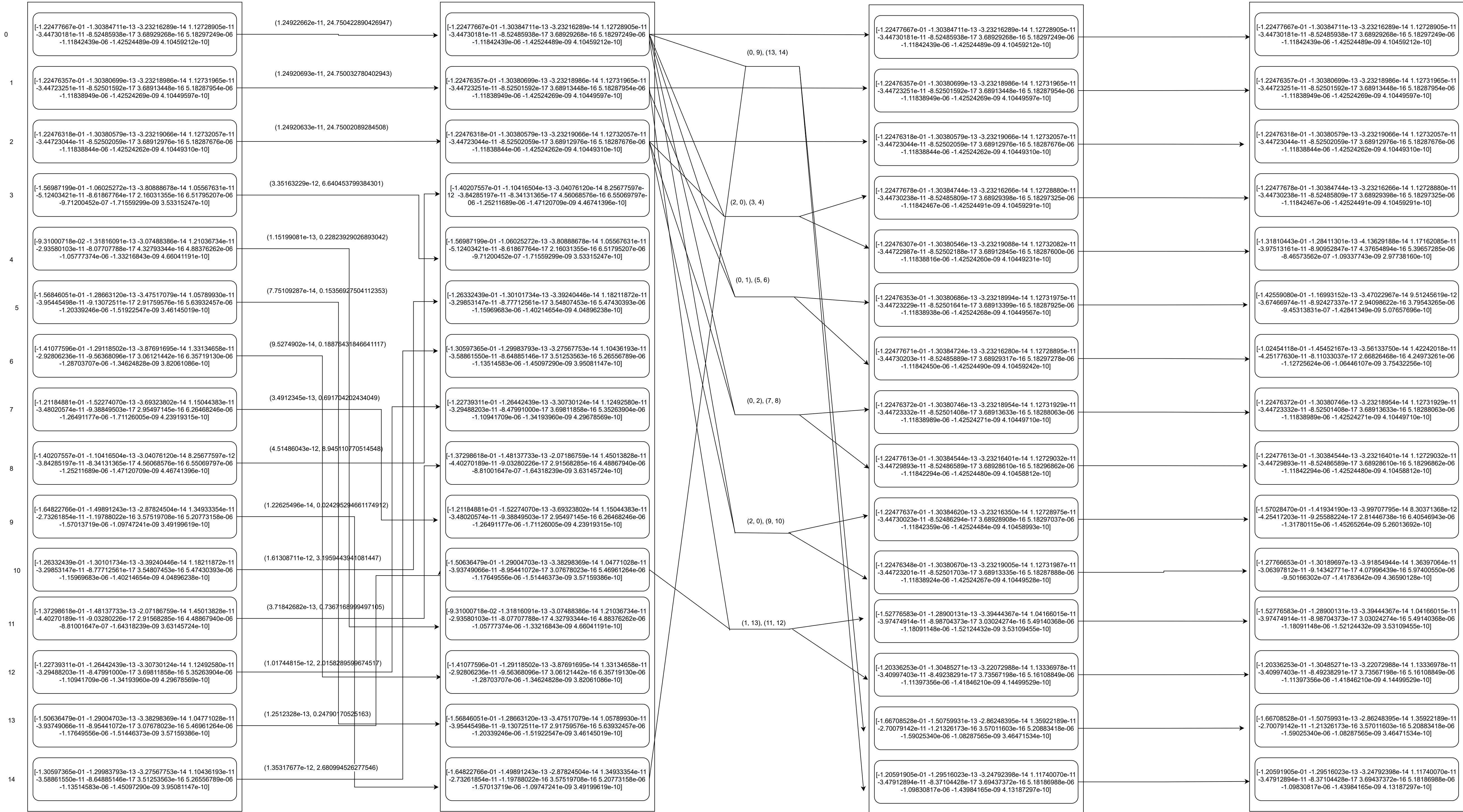
Fitness Function

Fitness Function

| INITIAL POPULATION | Weights , Percentage | Selection | Mapping (From INITIAL POPULATION TO AFTER CROSSOVER) | AFTER CROSSOVER | AFTER MUTATION |
|---|---|---|---|---|---|

**0**

INITIAL POPULATION 0:
[-1.22477667e-01 -1.30384711e-13 -3.23216289e-14 1.12728905e-11 -3.44730181e-11 -8.52485938e-17 3.68929268e-16 5.18297249e-06 -1.11842439e-06 -1.42524489e-09 4.10459212e-10]

Weights: (1.24922662e-11, 24.750422890426947)

Selection 0:
[-1.22477667e-01 -1.30384711e-13 -3.23216289e-14 1.12728905e-11 -3.44730181e-11 -8.52485938e-17 3.68929268e-16 5.18297249e-06 -1.11842439e-06 -1.42524489e-09 4.10459212e-10]

AFTER CROSSOVER 0:
[-1.22477667e-01 -1.30384711e-13 -3.23216289e-14 1.12728905e-11 -3.44730181e-11 -8.52485938e-17 3.68929268e-16 5.18297249e-06 -1.11842439e-06 -1.42524489e-09 4.10459212e-10]

AFTER MUTATION 0:
[-1.22477667e-01 -1.30384711e-13 -3.23216289e-14 1.12728905e-11 -3.44730181e-11 -8.52485938e-17 3.68929268e-16 5.18297249e-06 -1.11842439e-06 -1.42524489e-09 4.10459212e-10]

**1**

INITIAL POPULATION 1:
[-1.22476357e-01 -1.30380699e-13 -3.23218986e-14 1.12731965e-11 -3.44723251e-11 -8.52501592e-17 3.68913448e-16 5.18287954e-06 -1.11838949e-06 -1.42524269e-09 4.10449597e-10]

Weights: (1.24920693e-11, 24.750032780402943)

Selection 1:
[-1.22476357e-01 -1.30380699e-13 -3.23218986e-14 1.12731965e-11 -3.44723251e-11 -8.52501592e-17 3.68913448e-16 5.18287954e-06 -1.11838949e-06 -1.42524269e-09 4.10449597e-10]

AFTER CROSSOVER 1:
[-1.22476357e-01 -1.30380699e-13 -3.23218986e-14 1.12731965e-11 -3.44723251e-11 -8.52501592e-17 3.68913448e-16 5.18287954e-06 -1.11838949e-06 -1.42524269e-09 4.10449597e-10]

AFTER MUTATION 1:
[-1.22476357e-01 -1.30380699e-13 -3.23218986e-14 1.12731965e-11 -3.44723251e-11 -8.52501592e-17 3.68913448e-16 5.18287954e-06 -1.11838949e-06 -1.42524269e-09 4.10449597e-10]

**2**

INITIAL POPULATION 2:
[-1.22476318e-01 -1.30380579e-13 -3.23219066e-14 1.12732057e-11 -3.44723044e-11 -8.52502059e-17 3.68912976e-16 5.18287676e-06 -1.11838844e-06 -1.42524262e-09 4.10449310e-10]

Weights: (1.24920633e-11, 24.75002089284508)

Selection 2:
[-1.22476318e-01 -1.30380579e-13 -3.23219066e-14 1.12732057e-11 -3.44723044e-11 -8.52502059e-17 3.68912976e-16 5.18287676e-06 -1.11838844e-06 -1.42524262e-09 4.10449310e-10]

AFTER CROSSOVER 2:
[-1.22476318e-01 -1.30380579e-13 -3.23219066e-14 1.12732057e-11 -3.44723044e-11 -8.52502059e-17 3.68912976e-16 5.18287676e-06 -1.11838844e-06 -1.42524262e-09 4.10449310e-10]

AFTER MUTATION 2:
[-1.22476318e-01 -1.30380579e-13 -3.23219066e-14 1.12732057e-11 -3.44723044e-11 -8.52502059e-17 3.68912976e-16 5.18287676e-06 -1.11838844e-06 -1.42524262e-09 4.10449310e-10]

**3**

INITIAL POPULATION 3:
[-1.56987199e-01 -1.06025272e-13 -3.80888678e-14 1.05567631e-11 -5.12403421e-11 -8.61867764e-17 2.16031355e-16 6.51795207e-06 -9.71200452e-07 -1.71559299e-09 3.53315247e-10]

Weights: (3.35163229e-12, 6.640453799384301)

Selection 3:
[-1.40207557e-01 -1.10416504e-13 -3.04076120e-14 8.25677597e-12 -3.84285197e-11 -8.34131365e-17 4.56068576e-16 6.55069797e-06 -1.25211689e-06 -1.47120709e-09 4.46741396e-10]

AFTER CROSSOVER 3:
[-1.22477678e-01 -1.30384744e-13 -3.23216266e-14 1.12728880e-11 -3.44730238e-11 -8.52485809e-17 3.68929398e-16 5.18297325e-06 -1.11842467e-06 -1.42524491e-09 4.10459291e-10]

AFTER MUTATION 3:
[-1.22477678e-01 -1.30384744e-13 -3.23216266e-14 1.12728880e-11 -3.44730238e-11 -8.52485809e-17 3.68929398e-16 5.18297325e-06 -1.11842467e-06 -1.42524491e-09 4.10459291e-10]

**4**

INITIAL POPULATION 4:
[-9.31000718e-02 -1.31816091e-13 -3.07488386e-14 1.21036734e-11 -2.93580103e-11 -8.07707788e-17 4.32793344e-16 4.88376262e-06 -1.05777374e-06 -1.33216843e-09 4.66041191e-10]

Weights: (1.15199081e-13, 0.22823929026893042)

Selection 4:
[-1.56987199e-01 -1.06025272e-13 -3.80888678e-14 1.05567631e-11 -5.12403421e-11 -8.61867764e-17 2.16031355e-16 6.51795207e-06 -9.71200452e-07 -1.71559299e-09 3.53315247e-10]

AFTER CROSSOVER 4:
[-1.22476307e-01 -1.30380546e-13 -3.23219088e-14 1.12732082e-11 -3.44722987e-11 -8.52502188e-17 3.68912845e-16 5.18287600e-06 -1.11838816e-06 -1.42524260e-09 4.10449231e-10]

AFTER MUTATION 4:
[-1.31810443e-01 -1.28411301e-13 -4.13629188e-14 1.17162085e-11 -3.97513161e-11 -8.90952847e-17 4.37654894e-16 5.39657285e-06 -8.46573562e-07 -1.09337743e-09 2.97738160e-10]

**5**

INITIAL POPULATION 5:
[-1.56846051e-01 -1.28663120e-13 -3.47517079e-14 1.05789930e-11 -3.95445498e-11 -9.13072511e-17 2.91759576e-16 5.63932457e-06 -1.20339246e-06 -1.51922547e-09 3.46145019e-10]

Weights: (7.75109287e-14, 0.15356927504112353)

Selection 5:
[-1.26332439e-01 -1.30101734e-13 -3.39240446e-14 1.18211872e-11 -3.29853147e-11 -8.77712561e-17 3.54807453e-16 5.47430393e-06 -1.15969683e-06 -1.40214654e-09 4.04896238e-10]

AFTER CROSSOVER 5:
[-1.22476353e-01 -1.30380686e-13 -3.23218994e-14 1.12731975e-11 -3.44723229e-11 -8.52501641e-17 3.68913399e-16 5.18287925e-06 -1.11838938e-06 -1.42524268e-09 4.10449567e-10]

AFTER MUTATION 5:
[-1.42559080e-01 -1.16993152e-13 -3.47022967e-14 9.51245619e-12 -3.67466974e-11 -8.92427337e-17 2.94098622e-16 3.79543265e-06 -9.45313831e-07 -1.42841349e-09 5.07657696e-10]

**6**

INITIAL POPULATION 6:
[-1.41077596e-01 -1.29118502e-13 -3.87691695e-14 1.33134658e-11 -2.92806236e-11 -9.56368096e-17 3.06121442e-16 6.35719130e-06 -1.28703707e-06 -1.34624828e-09 3.82061086e-10]

Weights: (9.5274902e-14, 0.1887643184641117)

Selection 6:
[-1.30597365e-01 -1.29983793e-13 -3.27567753e-14 1.10436193e-11 -3.58861550e-11 -8.64885146e-17 3.51253563e-16 5.26556789e-06 -1.13514583e-06 -1.45097290e-09 3.95081147e-10]

AFTER CROSSOVER 6:
[-1.22477671e-01 -1.30384724e-13 -3.23216280e-14 1.12728895e-11 -3.44730203e-11 -8.52485889e-17 3.68929317e-16 5.18297278e-06 -1.11842450e-06 -1.42524490e-09 4.10459242e-10]

AFTER MUTATION 6:
[-1.02454118e-01 -1.45452167e-13 -3.56133750e-14 1.42242018e-11 -4.25177630e-11 -8.11033037e-17 2.66826468e-16 4.24973261e-06 -1.12725624e-06 -1.06446107e-09 3.75432256e-10]

**7**

INITIAL POPULATION 7:
[-1.21184881e-01 -1.52274070e-13 -3.69323802e-14 1.15044383e-11 -3.48020574e-11 -9.38849503e-17 2.95497145e-16 6.26468246e-06 -1.26491177e-06 -1.71126005e-09 4.23919315e-10]

Weights: (3.4912345e-13, 0.691704202434049)

Selection 7:
[-1.22739311e-01 -1.26442439e-13 -3.30730124e-14 1.12492580e-11 -3.29488203e-11 -8.47991000e-17 3.69811858e-16 5.35263904e-06 -1.10941709e-06 -1.34193960e-09 4.29678569e-10]

AFTER CROSSOVER 7:
[-1.22476372e-01 -1.30380746e-13 -3.23218954e-14 1.12731929e-11 -3.44723332e-11 -8.52501408e-17 3.68913633e-16 5.18288063e-06 -1.11838989e-06 -1.42524271e-09 4.10449710e-10]

AFTER MUTATION 7:
[-1.22476372e-01 -1.30380746e-13 -3.23216954e-14 1.12731929e-11 -3.44723332e-11 -8.52501408e-17 3.68913633e-16 5.18288063e-06 -1.11838989e-06 -1.42524271e-09 4.10449710e-10]

**8**

INITIAL POPULATION 8:
[-1.40207557e-01 -1.10416504e-13 -3.04076120e-14 8.256775997e-12 -3.84285197e-11 -8.34131365e-17 4.56068576e-16 6.55069797e-06 -1.25211689e-06 -1.47120709e-09 4.46741396e-10]

Weights: (4.51486043e-12, 8.945110770514548)

Selection 8:
[-1.37298618e-01 -1.48137733e-13 -2.07186759e-14 1.45013828e-11 -4.40270189e-11 -9.03280226e-17 2.91568285e-16 4.48867940e-06 -8.81001647e-07 -1.64318239e-09 3.63145724e-10]

AFTER CROSSOVER 8:
[-1.22477613e-01 -1.30384544e-13 -3.23216401e-14 1.12729032e-11 -3.44729893e-11 -8.52486589e-17 3.68928610e-16 5.18296862e-06 -1.11842294e-06 -1.42524480e-09 4.10458812e-10]

AFTER MUTATION 8:
[-1.22477613e-01 -1.30384544e-13 -3.23216401e-14 1.12729032e-11 -3.44729893e-11 -8.52486589e-17 3.68928610e-16 5.18296862e-06 -1.11842294e-06 -1.42524480e-09 4.10458812e-10]

**9**

INITIAL POPULATION 9:
[-1.64822766e-01 -1.49891243e-13 -2.87824504e-14 1.34933354e-11 -2.73261854e-11 -1.19788022e-16 3.57519708e-16 5.20773158e-06 -1.57013719e-06 -1.09747241e-09 3.49199619e-10]

Weights: (1.22625496e-14, 0.024295294661174912)

Selection 9:
[-1.21184881e-01 -1.52274070e-13 -3.69323802e-14 1.15044383e-11 -3.48020574e-11 -9.38849503e-17 2.95497145e-16 6.26468246e-06 -1.26491177e-06 -1.71126005e-09 4.23919315e-10]

AFTER CROSSOVER 9:
[-1.22477637e-01 -1.30384620e-13 -3.23216350e-14 1.12728975e-11 -3.44730023e-11 -8.52486294e-17 3.68928908e-16 5.18297037e-06 -1.11842359e-06 -1.42524484e-09 4.10458993e-10]

AFTER MUTATION 9:
[-1.57028470e-01 -1.41934190e-13 -3.99707795e-14 8.30371368e-12 -4.25417203e-11 -9.25588224e-17 2.81446738e-16 6.40546943e-06 -1.31780115e-06 -1.45265264e-09 5.26013692e-10]

**10**

INITIAL POPULATION 10:
[-1.26332439e-01 -1.30101734e-13 -3.39240446e-14 1.18211872e-11 -3.29853147e-11 -8.77712561e-17 3.54807453e-16 5.47430393e-06 -1.15969683e-06 -1.40214654e-09 4.04896238e-10]

Weights: (1.61308711e-12, 3.1959443941081447)

Selection 10:
[-1.50636479e-01 -1.29004703e-13 -3.38298369e-14 1.04771028e-11 -3.93749066e-11 -8.95441072e-17 3.07678023e-16 5.46961264e-06 -1.17649556e-06 -1.51446373e-09 3.57159386e-10]

AFTER CROSSOVER 10:
[-1.22476348e-01 -1.30380670e-13 -3.23219005e-14 1.12731987e-11 -3.44723201e-11 -8.52501703e-17 3.68913335e-16 5.18287888e-06 -1.11838924e-06 -1.42524267e-09 4.10449528e-10]

AFTER MUTATION 10:
[-1.27766653e-01 -1.30189697e-13 -3.91854944e-14 1.36397064e-11 -3.06397812e-11 -9.14342771e-17 4.07996439e-16 5.97400550e-06 -9.50166302e-07 -1.41783642e-09 4.36590128e-10]

**11**

INITIAL POPULATION 11:
[-1.37298618e-01 -1.48137733e-13 -2.07186759e-14 1.45013828e-11 -4.40270189e-11 -9.03280226e-17 2.91568285e-16 4.48867940e-06 -8.81001647e-07 -1.64318239e-09 3.63145724e-10]

Weights: (3.71842682e-13, 0.7367168999497105)

Selection 11:
[-9.31000718e-02 -1.31816091e-13 -3.07488386e-14 1.21036734e-11 -2.93580103e-11 -8.07707788e-17 4.32793344e-16 4.88376262e-06 -1.05777374e-06 -1.33216843e-09 4.66041191e-10]

AFTER CROSSOVER 11:
[-1.52776583e-01 -1.28900131e-13 -3.39444367e-14 1.04166015e-11 -3.97474914e-11 -8.98704373e-17 3.03024274e-16 5.49140368e-06 -1.18091148e-06 -1.52124432e-09 3.53109455e-10]

AFTER MUTATION 11:
[-1.52776583e-01 -1.28900131e-13 -3.39444367e-14 1.04166015e-11 -3.97474914e-11 -8.98704373e-17 3.03024274e-16 5.49140368e-06 -1.18091148e-06 -1.52124432e-09 3.53109455e-10]

**12**

INITIAL POPULATION 12:
[-1.22739311e-01 -1.26442439e-13 -3.30730124e-14 1.12492580e-11 -3.29488203e-11 -8.47991000e-17 3.69811858e-16 5.35263904e-06 -1.10941709e-06 -1.34193960e-09 4.29678569e-10]

Weights: (1.01744815e-12, 2.0158289599674517)

Selection 12:
[-1.41077596e-01 -1.29118502e-13 -3.87691695e-14 1.33134658e-11 -2.92806236e-11 -9.56368096e-17 3.06121442e-16 6.35719130e-06 -1.28703707e-06 -1.34624828e-09 3.82061086e-10]

AFTER CROSSOVER 12:
[-1.20336253e-01 -1.30485271e-13 -3.22072988e-14 1.13336978e-11 -3.40997403e-11 -8.49238291e-17 3.73567198e-16 5.16108849e-06 -1.11397356e-06 -1.41846210e-09 4.14499529e-10]

AFTER MUTATION 12:
[-1.20336253e-01 -1.30485271e-13 -3.22072988e-14 1.13336978e-11 -3.40997403e-11 -8.49238291e-17 3.73567198e-16 5.16108849e-06 -1.11397356e-06 -1.41846210e-09 4.14499529e-10]

**13**

INITIAL POPULATION 13:
[-1.50636479e-01 -1.29004703e-13 -3.38298369e-14 1.04771028e-11 -3.93749066e-11 -8.95441072e-17 3.07678023e-16 5.46961264e-06 -1.17649556e-06 -1.51446373e-09 3.57159386e-10]

Weights: (1.2512328e-13, 0.24790170525163)

Selection 13:
[-1.56846051e-01 -1.28663120e-13 -3.47517079e-14 1.05789930e-11 -3.95445498e-11 -9.13072511e-17 2.91759576e-16 5.63932457e-06 -1.20339246e-06 -1.51922547e-09 3.46145019e-10]

AFTER CROSSOVER 13:
[-1.66708528e-01 -1.50759931e-13 -2.86248395e-14 1.35922189e-11 -2.70079142e-11 -1.21326173e-16 3.57011603e-16 5.20883418e-06 -1.59025340e-06 -1.08287565e-09 3.46471534e-10]

AFTER MUTATION 13:
[-1.66708528e-01 -1.50759931e-13 -2.86248395e-14 1.35922189e-11 -2.70079142e-11 -1.21326173e-16 3.57011603e-16 5.20883418e-06 -1.59025340e-06 -1.08287565e-09 3.46471534e-10]

**14**

INITIAL POPULATION 14:
[-1.30597365e-01 -1.29983793e-13 -3.27567753e-14 1.10436193e-11 -3.58861550e-11 -8.64885146e-17 3.51253563e-16 5.26556789e-06 -1.13514583e-06 -1.45097290e-09 3.95081147e-10]

Weights: (1.35317677e-12, 2.680994526277546)

Selection 14:
[-1.64822766e-01 -1.49891243e-13 -2.87824504e-14 1.34933354e-11 -2.73261854e-11 -1.19788022e-16 3.57519708e-16 5.20773158e-06 -1.57013719e-06 -1.09747241e-09 3.49199619e-10]

AFTER CROSSOVER 14:
[-1.20591905e-01 -1.29516023e-13 -3.24792398e-14 1.11740070e-11 -3.47912894e-11 -8.37104428e-17 3.69437372e-16 5.18186988e-06 -1.09830817e-06 -1.43984165e-09 4.13187297e-10]

AFTER MUTATION 14:
[-1.20591905e-01 -1.29516023e-13 -3.24792398e-14 1.11740070e-11 -3.47912894e-11 -8.37104428e-17 3.69437372e-16 5.18186988e-06 -1.09830817e-06 -1.43984165e-09 4.13187297e-10]

Mapping labels:
(0, 9), (13, 14)
(2, 0), (3, 4)
(0, 1), (5, 6)
(0, 2), (7, 8)
(2, 0), (9, 10)
(1, 13), (11, 12)

# Fitness Function

```python
def calc_errors(individual, precalc):
    for item in precalc:
        if list(individual) == list(item[0]):
            return list(item[1][1:])
    else:
        return get_errors(ID, list(individual))


# population = (POP_SIZE * 11) matrix
def calc_weights(population, precalc):
    weights = np.zeros((len(population), 3), dtype=float)
    for i in range(len(population)):
        train_err, validation_err = calc_errors(population[i], precalc)
        weight = 1 / ((TRAIN_RATIO * train_err + VAL_RATIO *
validation_err) + DIFF_RATIO * abs(train_err - validation_err))
        weights[i] = weight, train_err, validation_err
    return weights
```

The fitness function is called calc_weights and it takes 2 arguments: the population for which we calculate the weights and precalc- the precalculated weights.
Since we are always storing 5 top individuals for the next population we need not waste a query on getting their errors as we already have them from the previous generation. Thus the top 5 individuals along with their respective train and validation errors are stored in precalc for saving queries.

The fitness value called weight is calculated using:
```
weight = 1 / ((TRAIN_RATIO * train_err + VAL_RATIO * validation_err) +
DIFF_RATIO * abs(train_err - validation_err))
```

$$weight = \frac{1}{TRAIN\_RATIO * train\_err + VAL\_RATIO * validation\_err + DIFF\_RATIO * abs(train\_err - validation\_err)}$$

Thus high values of train and validation errors and the difference between them will mean that the vector is not good and will have low fitness value while vectors having low train and validation errors and having less difference in errors will have higher fitness values.
This fitness function is a good measure of the model as it takes everything into account.

The values of TRAIN_RATIO, VAL_RATIO, DIFF_RATIO are changed from time to time depending upon the errors the current generation is producing. If train error is high, TRAIN_RATIO is increased. Similarly for validation error, difference.

## Selection Function

```python
def selection(population, weights):
    length = len(population)
    od_weights = weights[:, :1].flatten()
    sum_weights = sum(od_weights, start=0)
    normal_weights = [weight / sum_weights for weight in od_weights]
    indexes = np.random.choice(a=length, size=2, replace=False,
p=normal_weights)
    parents = [population[index] for index in indexes]
    return parents, np.array(indexes)
```

Selection function is much simpler. It takes all of the individuals in the population and their corresponding weights as arguments. Then we normalize the weights by dividing them by their sum so that they add up to 1. Now these normalized weights can act as probability for selection of a vector for being a parent in crossover. We thus choose and return 2 vectors that can act as parent for crossover.

## Crossover Function (Simulated Binary Crossover)

```python
# p1, p2 are parents each a vector of degree MAX_DEG
def simulated_binary_crossover(p1, p2):
    u = random.random()
    b = None
    if u <= 0.5:
        b = (2 * u) ** (1 / (DISTRIBUTION_INDEX + 1))
    else:
        b = (1 / (2 * (1 - u))) ** (1 / (DISTRIBUTION_INDEX + 1))
    c1 = 0.5 * ((1 - b) * p1 + (1 + b) * p2)
    c2 = 0.5 * ((1 + b) * p1 + (1 - b) * p2)
    return c1, c2
```

Simulated binary Crossover is used for the purpose of simulating single point binary crossover for real space. The idea is to calculate childrens c1 and c2 from parents p1 and p2 such that:

$$\frac{c_1 + c_2}{2} = \frac{p_1 + p_2}{2}$$

This is done by first choosing a random number u in [0,1). Then we calculate b using:

$$b = \begin{cases} (2u)^{\frac{1}{n_c+1}}, & if \ u \le 0.5 \\ \left(\frac{1}{2(1-u)}\right)^{\frac{1}{n_c+1}}, & if \ u > 0.5 \end{cases}$$

$$where \ n_c \ is \ the \ DISTRIBUTION\_INDEX$$

Then the children are calculated with the help of p1, and p2.

$$c1 = 0.5 * ((1 - b) * p1 + (1 + b) * p2)$$
$$c2 = 0.5 * ((1 + b) * p1 + (1 - b) * p2)$$

The graph below shows the probability density of offspring if the 2 parents are far enough.



The graph below shows the probability density of offspring if the 2 parents are close.



The graph below shows variation in probability density of offsprings for 2 values of distribution index 2 and 5.

## Mutation Function

```python
def mutation(population):
    for individual in population[ELITISM:]:
        if random.uniform(0, 1) < MUT_PROB:
            for index in range(11):
                individual[index] += individual[index] * random.uniform(-MUT_RANGE,
MUT_RANGE)

    for individual in population:
        for index in range(11):
            individual[index] = max(-10, individual[index])
            individual[index] = min(10, individual[index])
    return population
```

Mutation function takes population as input then it skips over the top elites that go directly to the next population. For the remaining individuals of the population, they get mutated with a probability that is defined by the hyperparameter MUT_PROB.
 If an individual is to be mutated then every feature of that individual is multiplied by a random number in range [1 - MUT_RANGE, 1 + MUT_RANGE]. This multiplication factor is calculated independently for each feature.

Once the mutation is done a check is applied to ensure that the features of vectors remain within the acceptable bounds of [-10,10]. If not they are set to the closest acceptable value.

# Hyperparameters

**POP_SIZE**
The population size of a generation i.e. the number of individuals in the generation. It was initially set to 25 to have more vectors and more variations. But after a point it was becoming difficult to manage 25 individuals with less queries so we shifted to 15.

**INITIAL_VECTOR**
The initial population is formed by taking an initial vector as the first individual of the population. Then a new individual is created by taking each feature of this vector and multiplying it with a randomly chosen number in the range [0,2]. This is done independently for each feature. This process is done in a loop till we reach POP_SIZE number of individuals.
The INITIAL_VECTOR is in a way a very rough guidance to the initial population. It also reduces the search space from [-10,10] for each feature to a very small amount. Thus the INITIAL_VECTOR was chosen to be the given OVERFIT VECTOR for the first time and was changed to other vectors when the situation demanded. (Explained in heuristics)

**ELITISM**
This is the number of top individuals transferred directly to the next generation. Elites have a not so subtle tendency to influence the vectors to converge in their direction. This speeds up convergence of solution and to ensure that it does not converge quickly to a local minima the MUT_RANGE is kept high.

**MUT_PROB**
The probability with which mutation occurs in a non-elite individual. This is kept quite stationary at 0.4 or maybe 0.5. These values ensure sufficient mutation and variation but also sufficient inheritance from parents.

**MUT_RANGE**
If an individual is to be mutated then every feature of that individual is multiplied by a random number in range [1 - MUT_RANGE, 1 + MUT_RANGE]. This multiplication factor is calculated independently for each feature. MUT_RANGE is kept high at 0.3 if large variations are needed for convergence. But once we have converged, this is systematically reduced to lower values like 0.1 and 0.05 to figure out finer variations in data.

**GENERATIONS**
Number of generations to run in every iteration of running the code. The lower the generation i.e. 10 to 15 the more tuning can be done at every point. But once vectors converge this can be increased to 20 to 30, as then the noticeable changes in vectors are hard to come by.

**VAL_RATIO**
VAL_RATIO is the coefficient of validation error in the fitness value. If VAL_RATIO is high then the algorithm prefers vectors with lower validation errors as this reduces the overall error and increases the fitness value.
Thus if validation error is very high compared to train error we increase VAL_RATIO and vice versa.

**TRAIN_RATIO**

TRAIN_RATIO is the coefficient of train error in the fitness value. If TRAIN_RATIO is high then the algorithm prefers vectors with lower train errors as this reduces the overall error and increases the fitness value.

Thus if train error is very high compared to validation error we increase TRAIN_RATIO and vice versa.

**DIFF_RATIO**

DIFF_RATIO is the coefficient of difference of train and validation errors in the fitness value. If DIFF_RATIO is high then the algorithm prefers vectors with closer train and validation errors as this reduces the overall error and increases the fitness value.

Thus if the difference between errors is very high, we increase DIFF_RATIO and vice versa.

**DISTRIBUTION_INDEX**

In Simulated Binary Crossover, this is a measure of how far the children are from the parents. Large values of DISTRIBUTION_INDEX allows children to be closer to parents. While small values mean children are farther from parents. 2 to 5 is the moderate range and the range we use. If the solution has not converged the DISTRIBUTION_INDEX is kept at 2, to get good variations. It is then systematically increased to 5 as the solution converges to have finer variations.

# Heuristics and Statistical Information

Initially we chose to use the following hyperparameters:
- INITIAL_VECTOR = OVERFIT VECTOR
- POP_SIZE = 25
- ELITISM = 5
- MUT_PROB = 0.4
- MUT_RANGE = 0.3
- GENERATIONS = 10
- VAL_RATIO = TRAIN_RATIO = DIFF_RATIO = 1
- DISTRIBUTION_INDEX = 2

This is done to have high variations in vectors but also to converge it. Initially the errors are of the order 1e14 but converge very quickly in 4 to 5 generations to order of 1e12. In another 3 to 4 generations they converge to the order of e11. The algorithm then got stuck in the range 1e11 to 3e11. In 70 generations we finally broke through the e10 barrier. Throughout this we would change TRAIN_RATIO, VAL_RATIO and DIFF_RATIO depending on the errors we received.

In a failed heuristic, we continued by reducing mutation range and increasing DISTRIBUTION_INDEX and running 100 to 200 generations more with 25 as POP_SIZE. This did not work too well and even when we reached 5e10 train and validation errors, the

leaderboard rank was not so good. We theorize that this was most probably due to overfitting.

So we changed the heuristic.
So we looked up test errors for different ranges of errors for different vectors. This was done in order to get a range of values where we are close to vectors with perfect fit but do not overfit.
We found such a vector and then took it as the INITIAL_VECTOR and restarted the genetic algorithm but this time with a reduced POP_SIZE of 15 as 25 vectors were taking a lot of queries without substantial results.
The MUT_RANGE value was 0.3 and we only changed it to 0.1 later but no further decrease. This was to ensure that overfitting does not take place and there is substantial variation in the data. The ratios were again altered dynamically according to latest data.
ELITISM took values between 1,3 and 5. The elite vectors do not alter themselves but do contribute towards convergence by participating in selection and crossover. Thus the more the elite vectors the more influence they have. So we sometimes reduced ELITISM if we wanted to keep more variations in the generations.
 Again the solutions converged very soon.

We repeated this process of restarting with an initial vector a few more times and then checked test errors for vectors in the range 1e11 to 3e11 and found very good vectors.

# Final Vectors and Errors

The top vector is :
[-0.11659813562929511, -1.2885110929990298e-13, -3.577791708504743e-14, 1.2646749212560554e-11, -3.3831773009699504e-11, -9.284248071512872e-17, 4.003342760912018e-16, 6.398320041396442e-06, -1.1245508241603073e-06, -1.4072152000100355e-09, 4.1834156007562354e-10]
With  train error = 103090039506.4589 and validation error = 312500882005.9999

While we were able to achieve train errors and validation errors as low as 4e10 and 2.8e10, these vectors did not give good results. This is because these vectors overfit both the train and validation datasets and thus performed badly on test dataset.

The vector we selected performs good on the test dataset because it does not overfit anything and all the 3 errors test, train and validation are close to each other. None of the errors are too high compared to one another. This is a good sign that the vector is a good fit.

# Tricks

1. Since we are using ELITISM, it is bound to be the case that some vectors get directly transferred to the next population. But we already have the errors and weights for

these vectors and there is no need to calculate them again by wasting queries. So these errors along with the corresponding vectors are stored in precalc and are reused for the next generation.

2. One trick used was to initially map out the sensitivity of each feature of the data and how reliable the OVERFIT VECTOR can be for the initial population.