

# MDL Project

The Team (87)

March 17, 2021

## 1 Genetic Algorithm

### 1.1 Summary

Genetic algorithms are inspired by the process of natural selection, and our genetic algorithm uses the same principle. The first generation is initialized by randomly slightly increasing or decreasing the values in the given overfit vector. The future generations are generated by combining parents picked from this generation, with their probability of being picked depending on their evaluated fitness. These children are also randomly mutated by a mutation function to increase diversity. This process is repeated multiple times till satisfactory vectors are obtained.

### 1.2 Pseudocode

1. Initialize first generation of population size  $n$  by perturbing overfit vector,
2. Repeat till convergence
  - (a) Calculate fitness based on errors using fitness function
  - (b) Pick parents from pool size  $s$
  - (c) Generate  $n$  children with a crossover function
  - (d) Mutate children with mutation function
  - (e) Set generated children as next population

### 1.3 Algorithm and code snippets

#### 1.3.1 Initializing first generation

The first generation is initialized by slightly perturbing the values in the given overfit vector. The values are randomly increased or decreased by a percentage. We also manually experimented with different values of changes at different indices.

```
1 population = np.zeros(shape = (n, 11))
2     for i in range(n):
3         rng = np.random.uniform(low = low_limit, high = high_limit, size=(1, 11))
4         population[i, :] = overfit_vector + rng*overfit_vector
```

### 1.3.2 Generating next generations

The next generations are generated by first pruning the parent generation to a pool size based on their fitness, and then picking parents from this pool. The children are generated with a crossover function which takes both the parents as input and then randomly mutating them. We also retained some of the best parents between generations, which is explained in the heuristics section. The following code snippet shows the procedure for generating one generation of children.

```
1 def get_next_pop(initial_pop, fitness, errors, n_parents = n_parents):
2     log = {"initial_pop": initial_pop, "mutations": False, "parents": [], "next_pop": []}
3
4     triple = sorted(list(zip(initial_pop, fitness, errors)), key = lambda x: x[1])
5     fitness = np.array(sorted(fitness))
6     p = [(1/i)/np.sum(1/fitness[:pool_size]) for i in fitness[:pool_size]]
7
8     next_fitness = []
9     next_errors = []
10    next_pop = []
11    parents = []
12    mutations = []
13
14    for _ in range(population-n_parents):
15        parent1 = np.array(triple[np.random.choice(range(pool_size), p = p)][0])
16        parent2 = np.array(triple[np.random.choice(range(pool_size), p = p)][0])
17        parents.append((parent1, parent2))
18
19        child = crossover(parent1, parent2)
20        if(np.random.uniform() < .5 or (list(parent1) == list(parent2))):
21            mutations.append(child)
22            child = mutation(child, prob = 1)
23        else:
24            mutations.append(False)
25
26        c_fitness, c_error = fitness_function(child)
27        next_fitness.append(c_fitness)
28        next_errors.append(c_error)
29        next_pop.append(child)
30
31    for i in range(n_parents):
32        # print(i, initial_pop[i], errors[i])
33        next_pop.append(np.array(triple[i][0]))
34        next_fitness.append(triple[i][1])
35        next_errors.append(triple[i][2])
36
37    log["errors"] = next_errors
38    log["fitness"] = next_fitness
39    log["parents"] = parents
40    log["next_pop"] = next_pop
41    log["mutations"] = mutations
42
43    f = open(f'{path}/{time.time()}.txt', 'wb')
44    pickle.dump(log, f)
45    return log
```

## 2 Functions

### 2.1 Fitness function

The fitness function used is a weighted sum of the train and validation error. The weight of train error was modified between generations according to our need, which is explained in the heuristics section. Another type of fitness function we used is given in the comment code. We hypothesized that this would allow our algorithm to pick errors closer to a target error, and reduced the target error as the generations progressed further. A combination of the two error functions was used.

```
1 def fitness_function(vector, weight):
2     errors = get_errors(key, list(vector))
3     fitness = abs(weight*errors[0]+errors[1])
4     #fitness = (abs(errors[1]-5e10) +abs(errors[0]-5e10))
5     return fitness, errors
```

### 2.2 Crossover function

The crossover function used is **blend crossover**. Given parents  $p_1$  and  $p_2$  ( $p_2 > p_1$ ), it randomly picks children in the range  $[p_1 - \alpha(p_2 - p_1), p_2 + \alpha(p_2 - p_1)]$  with a uniform distribution.

```
1 def crossover(parent1, parent2):
2     alpha = 0.5
3     child = np.zeros(11)
4     for i in range(11):
5         max_p = max(parent2[i], parent1[i])
6         min_p = min(parent2[i], parent1[i])
7         low = min_p - alpha*(max_p - min_p)
8         high = max_p + alpha*(max_p - min_p)
9
10        child[i] = np.random.uniform(low = low, high = high)
11    return child
```

### 2.3 Mutation function

The mutation function randomly increases or decrease the value at an index by a certain small percentage. This introduces more diversity in the generations. Depending on the requirement, this range of modifications was from as little as .005 to as large as .5

```
1 def mutation(vector, prob = index_prob):
2     rng = (np.random.uniform(size = 11)) < prob
3     delta = np.random.uniform(low = low_limit, high = high_limit, size = 11)
4     vector = vector + delta*vector
5     return vector
```

## 3 Hyperparameters

- Population size - This was initially chosen to be 20. However, it was gradually reduced to 10 as we obtained better generations and to preserve requests.
- Pool size - This ranged from 5-10 based on the population size. Half the population was pruned. This was done because it results in faster convergence.

- Alpha - Based on recommended choices of alpha for blend crossover in our research, we initially chose alpha as 0.5. This was modified to reduce or increase the range of values generated depending on performance.
- Parents retained - Between 1-4 parents were retained between generations. We did not want to make this value too high so as to ensure that diversity is maintained.
- Mutation probability - The mutation probability ranged from .3 to .7 depending on how fast the error was decreasing between generations. We increased this if the error was dropping too slowly or if we got stuck in a local minima.

## 4 Heuristics

- Weighted sum as fitness function - As explained above, the fitness function we used was a weighted sum of the errors. This was done to prioritize train or validation error at different errors. Assigning a higher weight to validation error allowed us to reduce it quicker, while using equal weights allowed us to bring the values together. By assigning the weight as -1, we also used the absolute difference between the two errors as a fitness metric. This allowed our genetic algorithm to generate vectors that minimize the difference between the train and validation error, minimizing the chance of them overfitting to any particular dataset.
- Parents retained - The parents retained between generations allowed us to ensure that at the very least the performance does not degrade between generations. However, we took care to keep this value small, to ensure that diversity isn't compromised. On occasions, the parents retained continued achieving the best fitness for multiple generations. This was undesirable, as it indicated that the newer generations were not performing well. These iterations were scrapped.
- Mutation probability - The mutation probability was varied between various generations. When the performance wasn't improving, we increased the probability to introduce more diversity in the generations. We also changed the range of variation for certain indices manually according to our estimate of performance. A blanket range of change for all indices was tried earlier, but this resulted in erratic performance as certain indices required much finer modifications than other.
- Through observation and manual experimentation, we observed that even drastic change in values at certain indices did not affect the error in any appreciable way. We focused our efforts on fine-tuning the value of the other indices. Initializing the first generation on this basis was tried, but it did not result in any significant changes to performance.

## 5 Results

1. The best errors we obtained for a single vector are -
  - Train Error - 0.74e11
  - Validation - 0.83e11

The overfit vector we started with had much lower train error and higher validation error. That is a pattern observed when the algorithm overfits to the training set and does not generalize well to unseen data. By contrast, the vectors our genetic algorithm was able to generate had lower validation error and higher training error, indicating that these vectors no longer overfit the training set. These values are also close to each other, meaning that similar performance is

being attained in different datasets. This should generalize well to unseen data. We also obtained vectors which had much lower validation errors, but these had likely overfit to the validation data.

2. The errors converged in around 70 generations. Further generations allowed us to finetune the parameters.

## 6 Illustration







