

Simple Evolutionary Exploration into Classification Algorithms

Davin Lin¹; Dirk Colbry, PhD²
¹Grinnell College, ²Michigan State University

Background

- Genetic algorithms (GAs) find good solutions in large solution spaces to search problems through a process inspired by evolution.
- There have been several studies that use GAs to search over hyperparameters of machine learning algorithms to learn values that work well for specific problems.
- The SEE Toolkit which uses and implements a Simple Genetic Algorithm (**Figure 1**) to support research workflows does not support searching over supervised-learning machine learning algorithms.

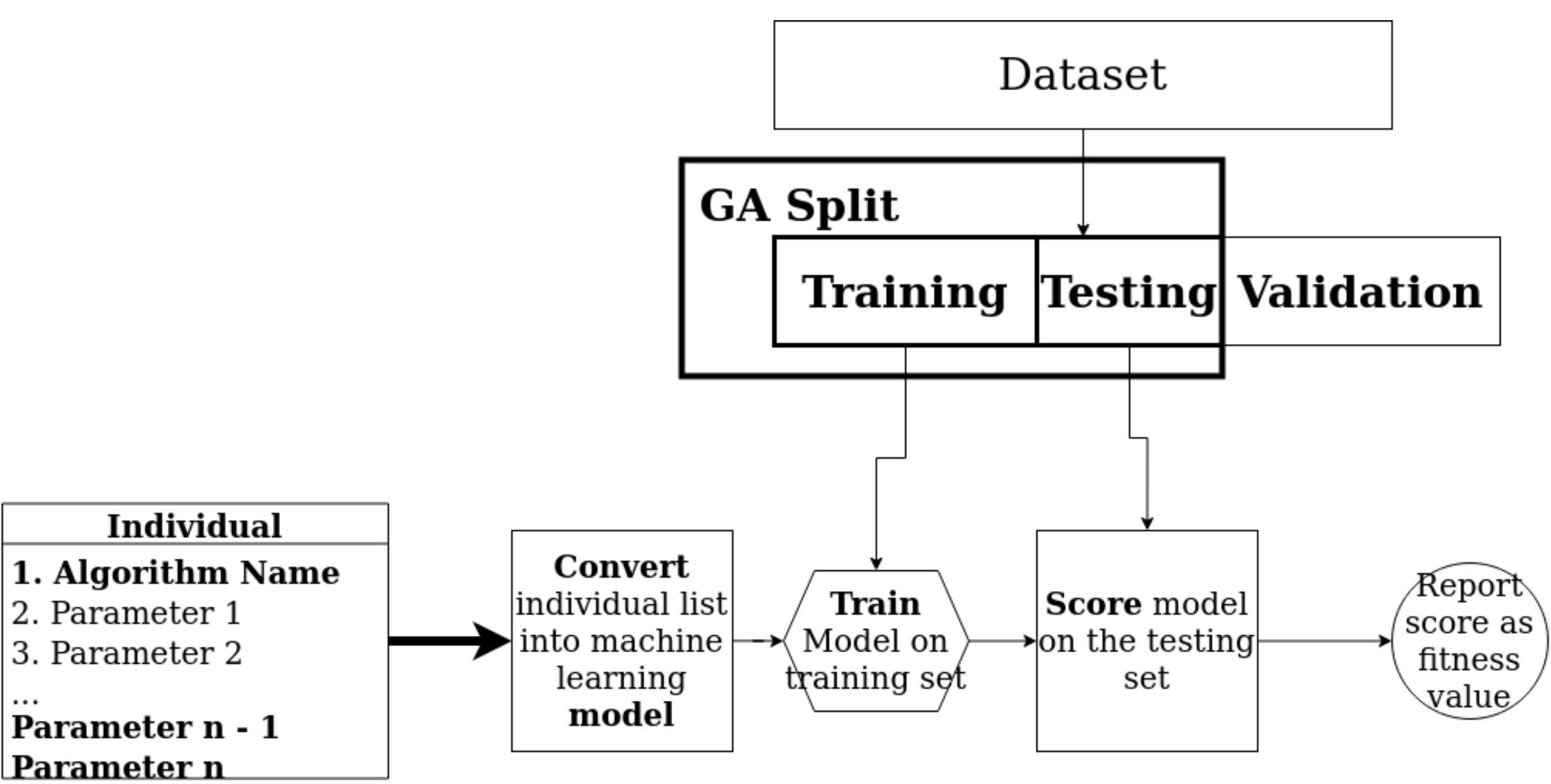
Goals

- To **extend** the SEE toolkit and framework to **support** supervised-machine learning classification algorithms
- To **benchmark** its performance via a comparison with a **previous work** on hyperparameter tuning via genetic search

Methods and Materials

- Methods:
- 30 trials of GA for 100 iterations and population of 100
 - Fitness Function: (**Figure 2**)
$$\text{Error Rate} = \frac{\text{number of incorrect}}{\text{number of classifications}}$$
- Materials:
- Breast Cancer Wisconsin (Diagnostic) Dataset
 - Dhahri et al. 2019 (study on hyperparameter tuning)
 - High Performance Computing Center at MSU ICER
 - Scikit-learn package for machine learning

Figure 2.
Fitness Function



Results*

- Our framework is able to find solutions similar to reported results (**Figure 3**)
- While the best found solutions converge during GA, when evaluated independently on the validation set, their fitness scores fall far from the range of performance during GA (**Figure 4**).

Figure 1. Simple Genetic Algorithm

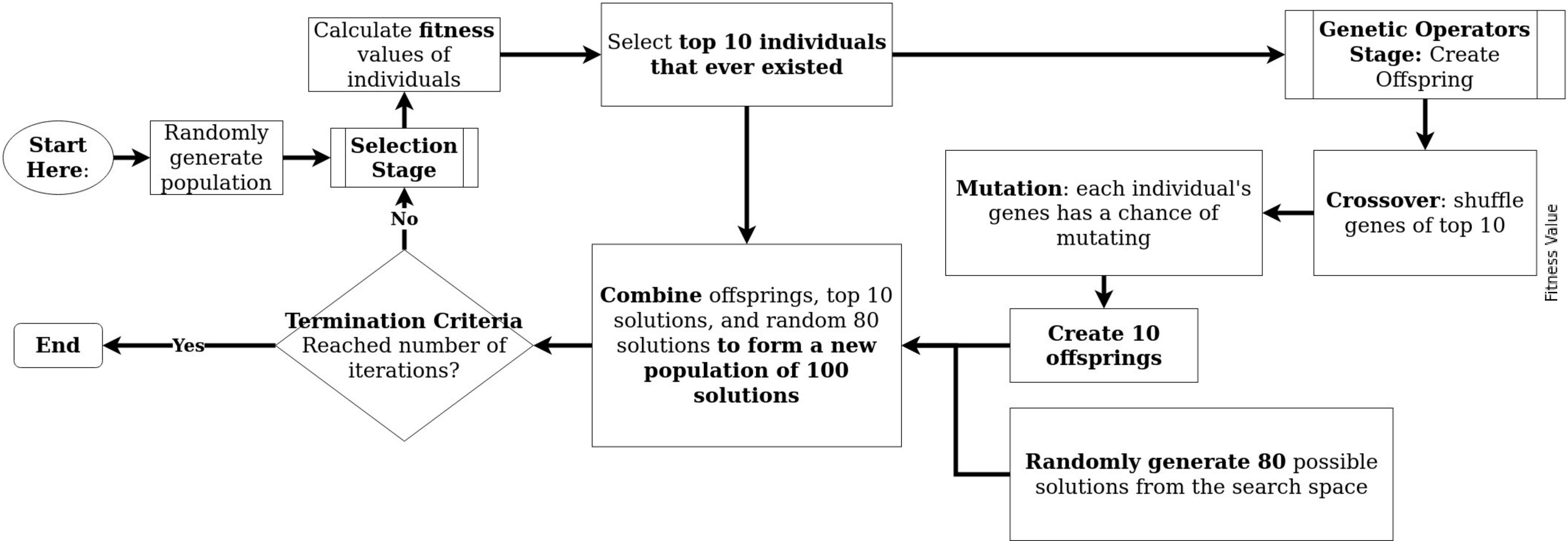


Figure 3.

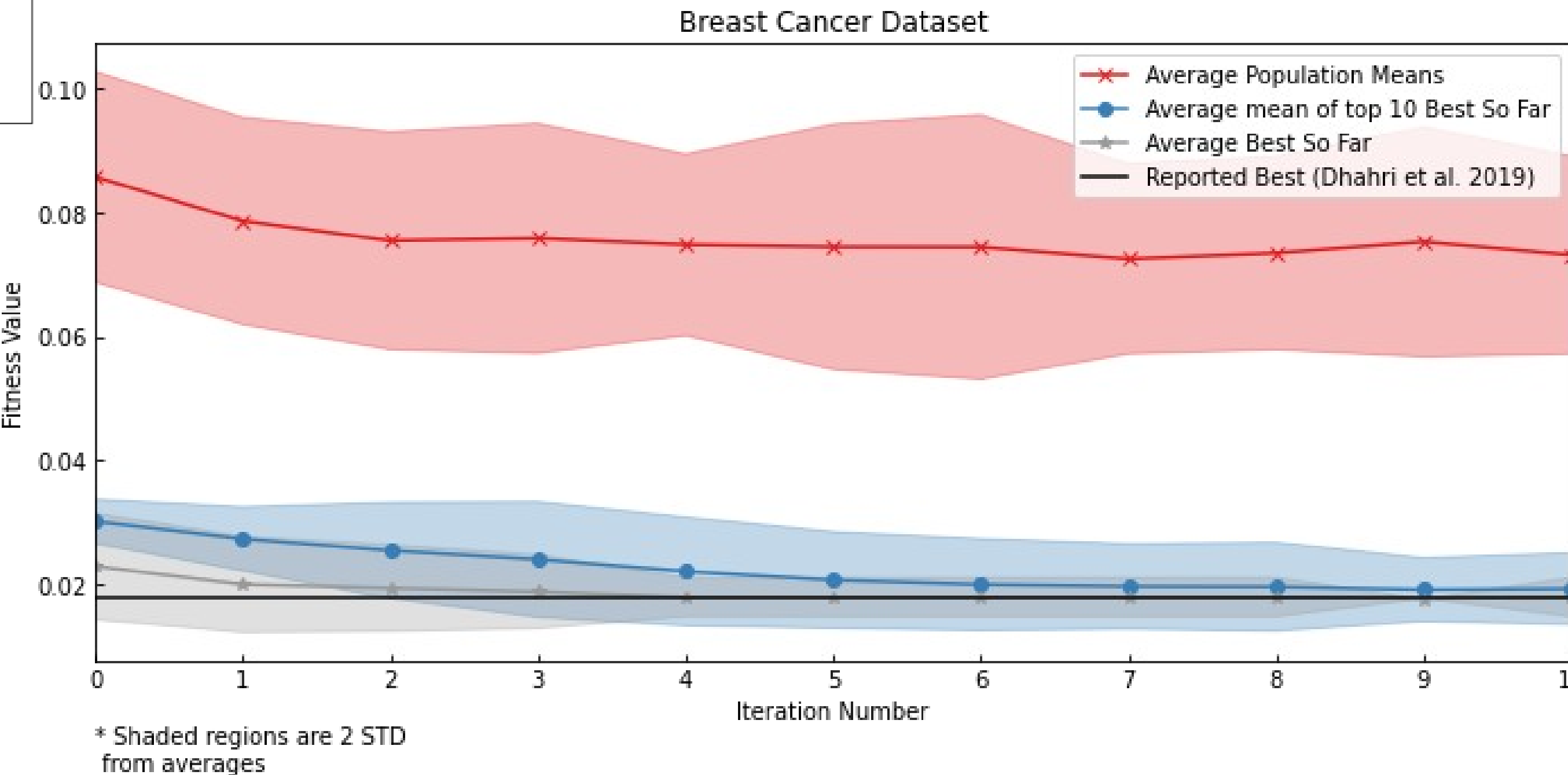
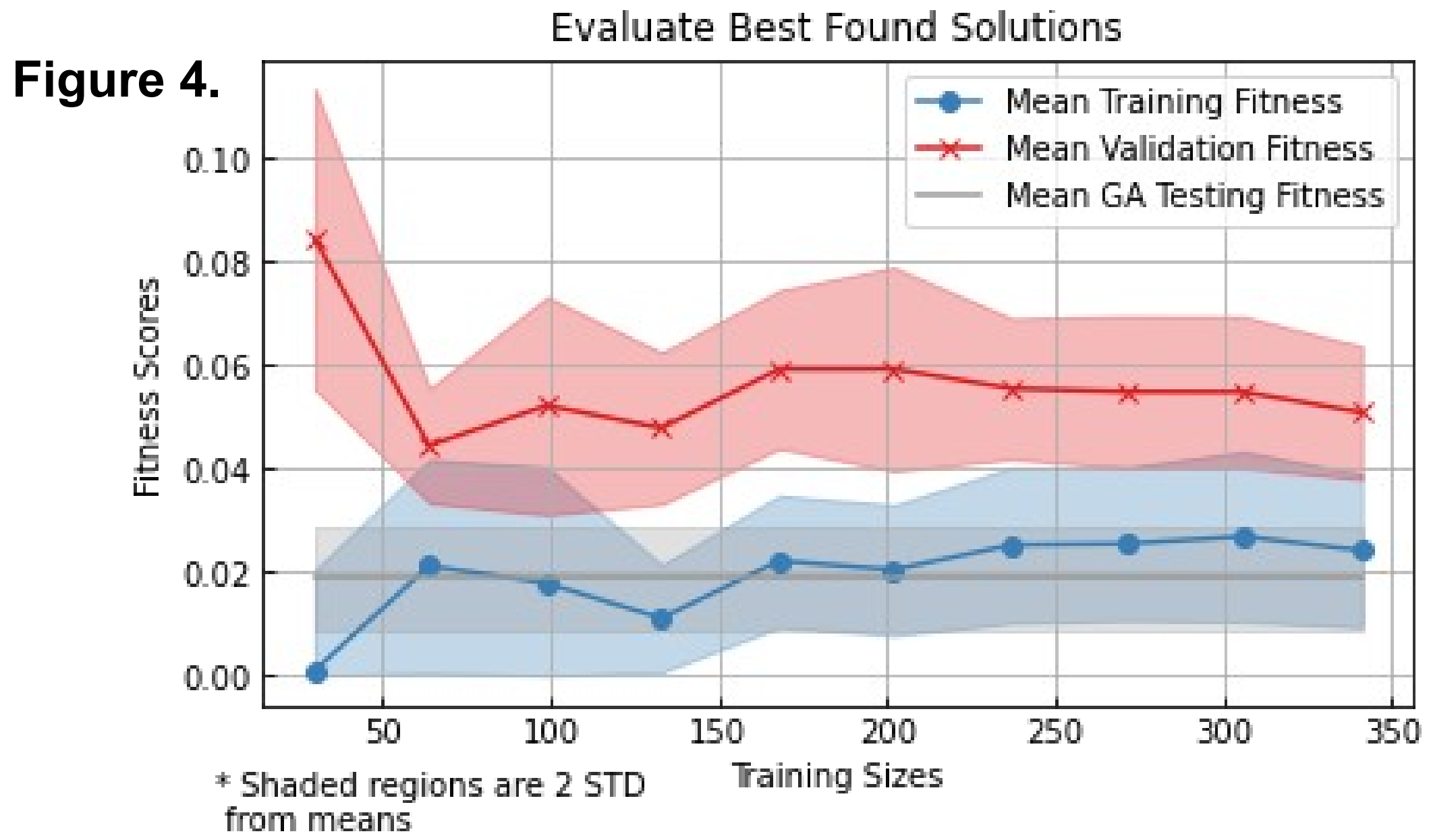


Figure 4.



Conclusions and Future Work

We extended the SEE framework to support classification algorithms and benchmarked the performance of its simple evolutionary approach.

Future Work:

- The implemented simple GA is a **naive** approach and can be greatly improved.
- Benchmark GA on alternate real datasets and using more complex metrics for the fitness function