

# GA Documentation

Mary Combs, Heejung Kim, Mohammad Soheilypour, Linqing(Waverly) Wei

December 15, 2017

## 1 Github

User Name: WaverlyWei

Link: <https://github.com/WaverlyWei/GA>

## 2 Contributions

Group Memebers: Mary Combs, Heejung Kim, Mohammad Soheilypour, Linqing(Waverly) Wei

Mary: Implemented Selection function; Implemented Helper functions(GHFitness, calculateAIC); Wrote helper page for selection function

Heejung:Implemented initiation function; Wrote test cases

Mohammad:Implemented mutation function; Wrote helper information for initiation, crossover and mutation functions; Organized the strucutre of overall Select function

Waverly:Implemented crossover function; Wrote pdf documentation; Created Git Repo; Created initial R package

## 3 Approach

The solution is composed of three levels: Main function (Select), Modular functions(initiation, selection, crossover, mutation) and Helper functions(calculateAIC, GHFitness).

## 4 Main function:Select

Input: User defined dataset(matrix), Model(lm/glm), Fitness function

Output: Selected variabels(list)

Stopping Criteria: AIC value converges  $\Delta AIC - > 0$

Select function first calls Initiation function to get starting data. The starting data are separated into two components: parents and intercept. Since converge happens asymptotically in this solution, the limiting steps of converge is preset within Select function. Then it continues to the next step: breeding the next generation. To generate the next generation, according to Genetic Algorithm, four modular functions are involved. Parents are passed into Selection function, ranked based on objective function(fitness) and parent matrix is updated for the next steps. Children are then passed into Crossover and Mutation function to be randomized. The outcome is then set to be the next generation of parents and will go into another round.

## 5 Modular Functions

### 5.1 Initiation

Input:vars(number of possible variables), pSize(population size for each generation), minC(min subset size in percent ), maxC(max subset size in percent)  
Output:Initial Matrix with 0s and 1s

Initiation function creates a subset of original dataset as the starting generation for Select function. The starting matrix passed onto Select function contains both the initial matrix and intercepts.

### 5.2 Selection

Input:dataset, outcome, parents(matrix), intercept, fitness function(Use GHFitness for now, can also be user defined)  
Output: parents selection, intercept selection, AIC ( A list of three components)

Selection function assigns initial fitness probability to each parent and then uses AIC as the objective function to rank parents. After sorting, fitness probability is updated for each parent. Meanwhile, updated intercept and AIC value is also recorded together with selected parents.

### 5.3 Crossover

Input:P1, P2(parent strings)  
Output:P3,P4(crossover strings, a list of two components)

Crossover function takes two parents strings and randomly choose a crossover site. Parents strings are cut at the site and ligated to produce children strings.

## 5.4 Mutation

Input: Parent, MutationProb

Output: Mutated Parent

Mutation functions takes a parent and choose random site to mutate under the constraint of mutation probability.

## 6 Helper Functions

### 6.1 GHFitness

Assign fitness probability to each parent based on the formula:

$$\phi(v_i^{(t)}) = \frac{2r_i}{P(P+1)}$$

This assignment gives the best individual a probability of  $2/(P+1)$

### 6.2 calculateAIC

AIC is calculated by applying linear model on response vector Y and covariates X. This function is used both for ranking individuals and for calculating the convergent AIC value (stopping criteria).

## 7 Testing

Testing is performed on four modular functions in two aspects: (1) Valid input  
(2) Expected return results

### 7.1 Example 1: Simulated Data

```
initData <- matrix(rnorm(2500, sd = 1:5), ncol = 5, byrow = TRUE)
```

```
initOutcome <- -1 + -1*initData[,1] + 2*initData[,3] + 1.1*initData[,5]
```

expected output: c(1,2,4,6)

actual output: c(1,6)

(1: intercept, 2: X1, 3:X2, etc.)

### 7.2 Example 2: Whitewine Quality Data

Input: White Wine Quality Dataset

Output: c(1,2,3,4,5,7,8,9,11,12)

(Selected variables indexed by column number)

```
##      fixed.acidity volatile.acidity citric.acid residual.sugar chlorides
## 1          7.0          0.27          0.36          20.7          0.045
## 2          6.3          0.30          0.34          1.6          0.049
## 3          8.1          0.28          0.40          6.9          0.050
## 4          7.2          0.23          0.32          8.5          0.058
## 5          7.2          0.23          0.32          8.5          0.058
## 6          8.1          0.28          0.40          6.9          0.050
##      free.sulfur.dioxide total.sulfur.dioxide density    pH sulphates alcohol
## 1                   45                   170 1.0010 3.00          0.45          8.8
## 2                   14                   132 0.9940 3.30          0.49          9.5
## 3                   30                   97 0.9951 3.26          0.44         10.1
## 4                   47                   186 0.9956 3.19          0.40          9.9
## 5                   47                   186 0.9956 3.19          0.40          9.9
## 6                   30                   97 0.9951 3.26          0.44         10.1
##      quality
## 1          6
## 2          6
## 3          6
## 4          6
## 5          6
## 6          6
## [1] 4898 12
```