# Genetic Algorithm for Variable Selection in Regression Problems

November 28, 2017

## 1  Object

Your task is to make an R package to implement a genetic algorithm for variable selection in regression problems, including both linear regression and GLMs.

## 2  Genetic Algorithm

### 2.1  Overall

Analogy:

- One Chromosome $\rightarrow$ One candidate model.

- number of alleles $C \rightarrow$ number of variables.

- $i$-th allele $\rightarrow$ $i$-th variable.

- $i$-th allele is 0 $\rightarrow$ not include $i$-th variable.

- $i$-th allele is 1 $\rightarrow$ include $i$-th variable.

Say $C = 4$. For a single chromosome, 1010, the suggested linear regression model is

$$y = \beta_0 + \beta_1 x_1 + \beta_3 x_3 + \varepsilon \tag{2.1}$$

while $\beta_2$, $\beta_4$ are zero.

### 2.2  GARS Genetic Algorithm for Regressors Selection [2, Section 3]

1. Each GA individual consists of a string of $C$ binary alleles: if the $i$-th cell $(i = 1, ..., C)$ has value 1, then $X_i$ is included in the model, otherwise not.

2. Every candidate solution is then evaluated with respect to a fitness function. The *AIC*, *BIC* and *SIC* criteria have been considered as possible fitness functions.

3. Randomly initializing the population and evaluating the population with respect to the chosen fitness function.

4. Generation evolution using (stochastic uniform sampling selection scheme?), single point crossover with $p_c = 0.8$, uniform mutation with $pm = 1/NBITS$ and (direct reinsertion of the best recorded candidate solution?).

1. Initialization:

   - randomized with all zero and all one.

2. Iteration:

(a) Evaluate all the initial points:
    evaluate AIC or other fitness function.

(b) Selection:
    continue on the next generations and use them to create the subsequent generation.

(c) Crossover:
    Crossover at all the points.

(d) Mutation:
    mutation.

(e) Next Generation and Redo the iteration

3. Until convergence

## 2.3 GARST Genetic Algorithm for Regressors Selection and Transformation

...

## 2.4 Primary Issues

- crossover in model selection? how many crossover point? Provide an option for the function, $crossover = TRUE/FALSE$, $crossover\_point =$ integer input, randomly distribute the points given the number or input a vector, customize the position of crossover. Another issue: should the number and position of crossover points be random?

## 2.5 Initialization and Parameter values [1, Chapter 3]

- Equal sizes for subsequent generations are not required.

- large generation size P for early generations to discourage premature convergence and promote search diversity.

- P can be decreased progressively as iterations continue

- variable mutation rate that is inversely proportional to the population diversity. Purpose: provides a stimulus to promote search diversity as generations become less diverse.

## 2.6 Evaluation

Evaluate the fitness function/Criteria:

- Akaike information criterion.
  From $\{stats\}$ package, the $AIC()$ function.

- Bayesian information criterion.

- Cross-validation.

- Deviance information criterion.

- False discovery rate.

- Focused information criterion.

For other fitness functions, user can supply the function defined themselves.

## 2.7 Selection

Select a few parents from the $i-th$ generation to continue on the next generations and use them to create the subsequent generation. The selection criteria:

- Based on the rank of the fitness function (AIC of the models).

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

where $P$ is the size of the generation, $r_i$ is the rank of $f(\theta_i^{(t)})$ among generation t.

# 3  Requirement

## 3.1  Generality

Allow reasonable inputs, in terms of specifying a dataset and regression model formula, as well as the type of regressions. Much of this is information you should just be able to pass along to lm() or glm().

## 3.2  Coding Style

- modular code

- functions or OOP methods that implement discrete tasks.

- overall design and style that is consistent across the components, in terms of functions vs. OOP methods, naming of objects, etc.

## 3.3  Efficiency

- Try to vectorize as much as possible to speed up.

- allow for shared memory parallel processing when working with the population in a given generation, in particular the evaluation of the fitness function,

## 3.4  Validation and Testing

- implementation on one or more real examples.

- Formal testing is required, with a set of tests where results are compared to some known truth.

## 3.5  Criterion

By default AIC as your objective criterion/fitness function. Allow specific fitness functions.

## 3.6  Operator

use the genetic operators described in Givens and Hoeting for variable selection. Allow additional operators.

# References

[1]  G. H. Givens and J. A. Hoeting. *Computational Statistics*. John Wiley & Sons, Inc., 2 edition, 2013.

[2]  S. Paterlini and T. MinervaI. Regression model selection using genetic algorithms. *Recent Advances in Neural Networks, Fuzzy Systems and Evolutionary Computing*, pages 19–27, 1 2010.